

# Machine Learning and Differential Equations

University of Bath



Darrah Kirkpatrick

April 20, 2021

## **Abstract**

In this report we aim to teach the theory that is required to understand Generative Latent Time-series Models and its applications to make time-series predictions. This model was first introduced in the Neural Ordinary Differential Equation paper [1], which developed a link between deep neural networks and (discretised) ordinary differential equations. The Neural Ordinary Differential Equation approach showed how theory from optimal control can be used to compute learning gradients of our model by solving a single integral. This approach demonstrated the ability to compute learning gradients quickly, whilst allowing the problem to scale linearly with size, representing a substantial advance in Deep Learning literature.

# Contents

1	Introduction . . . . .	2
2	Statistics . . . . .	2
2.1	Marginal Likelihood . . . . .	2
2.2	Kullback–Leibler Divergence . . . . .	4
2.2.1	Important Example . . . . .	4
3	Machine Learning Basics . . . . .	5
3.1	Parametric models . . . . .	5
3.2	Loss Function . . . . .	5
3.3	Gradient Descent . . . . .	6
3.4	Training . . . . .	7
3.5	Theory Example . . . . .	9
3.6	Neural Networks Architecture . . . . .	11
3.7	Backpropagation Algorithm . . . . .	11
3.8	Generative Modelling . . . . .	12
4	Variational Autoencoders . . . . .	13
4.1	Autoencoder . . . . .	13
4.2	Variational Autoencoders . . . . .	14
4.2.1	Motivation . . . . .	14
4.2.2	Probabilistic Framework . . . . .	15
4.2.3	The Problem With This Approach . . . . .	15
4.2.4	Overcoming This Problem . . . . .	16
4.2.5	Reparameterization Trick . . . . .	17
4.3	Gaussian VAE Example . . . . .	19
4.4	Application of VAE to MNIST Dataset . . . . .	20
5	Neural Ordinary Differential Equations . . . . .	22
5.1	Motivation . . . . .	23
5.1.1	Toy Circular Dynamics Example . . . . .	23
5.2	Formalisation and Improvements . . . . .	25
6	Generative Latent Time-series Models . . . . .	27
6.1	Overview . . . . .	27
6.2	Training . . . . .	28
6.3	Bi-directional Spiral Example . . . . .	31
<b>A Little Results and Derivations</b>		<b>33</b>
1	Gibbs’ Inequality [2] . . . . .	33
2	Adjoint and Augmented Dynamics Derivations . . . . .	33
2.1	Adjoint Derivation . . . . .	33
2.2	Augmented Dynamics . . . . .	34

# 1 Introduction

In this report, we will explore the theory required to understand generative latent time-series modelling and its applications to make time-series predictions. To achieve this, we will split the report into four main parts:

1. Background Theory
2. Variational Autoencoders
3. Neural Ordinary Differential Equations
4. Generative Latent Time-series Models

Where the first three parts will introduce the theory required to understand generative latent time-series models. With the final part, putting it all together and showing its application to a make time-series predictions.

## 2 Statistics

Whilst this paper will not focus heavily on the statistical side of machine learning, it is important that we recap some key statistical theory as it will provide motivation and deepen our understanding to what we discuss later.

### 2.1 Marginal Likelihood

In this section we will introduce some important machinery that will help us quantify how "good" a model is at modelling our data. This can extremely useful as fitting a model to our data can help us understand the underlying mechanisms at play which will help aid in future predictions.

We start off by framing the problem that we would like to solve. Suppose that we have the dataset  $\mathbf{X} = (\mathbf{x}^{(i)})_{i=1}^N$  with  $\mathbf{x}^{(i)} \in \mathbb{R}^n$ , and we want to find a parametric model  $p_\theta(\cdot)$ , that closely fits our data  $\mathbf{X}$ . As our parametric model depends on the value of  $\theta$ , we want to find the value of  $\theta$  that bests fits  $p_\theta(\cdot)$  to our data.

We will demonstrate this idea with the following small example. Suppose we knew our data  $\hat{X} = \{\hat{x}^{(i)}\}_{i=1}^N$  with  $\hat{x}^{(i)} \in \mathbb{R}$ , was a normally distributed with variance 0.1 and an unknown mean. Hence by using the PDF of the normal distribution  $\mathcal{N}(\theta, 0.1)$ , as our parametric model  $p_\theta(\cdot)$ . We aim to find the value of best value for  $\theta$  (which we will refer to as  $\theta^*$ ) such that  $p_\theta(\cdot)$  models our data. We will do this by minimising the difference between our model and the data (i.e,  $\theta^* = \operatorname{argmin}_\theta \sum_{i=1}^N ||p_\theta(x) - \hat{x}^{(i)}||^2$ , for  $x \in \mathbb{R}$ ).

We will now formally define how "good" our parametric model is at fitting our data, by how likely it is to see our data given our model,

**Definition 2.1** (Likelihood function [4]). Given a sample  $\mathbf{x}^{(i)}$  from our data set  $\mathbf{X}$ , and a parametric model  $p_\theta(\cdot)$ , the likelihood function  $L$  is defined as,

$$L(\theta, \mathbf{x}^{(i)}) = p_\theta(\mathbf{x}^{(i)}) \tag{1}$$

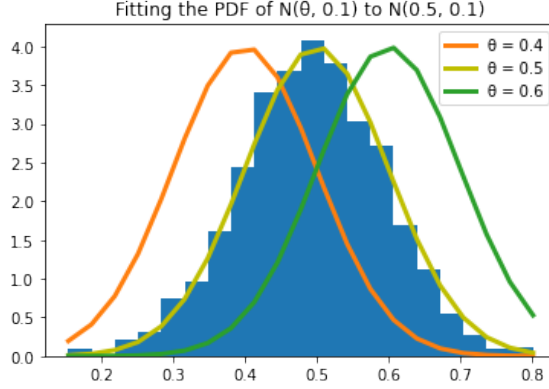


Figure 1: Fitting the PDF of  $\mathcal{N}(\theta, 0.1)$  to 1000 samples of  $\mathcal{N}(0.5, 0.1)$  with different values of  $\theta$  (i.e  $\theta = 0.4, 0.5, 0.6$ ). [3]

Therefore, given the sample  $\mathbf{x}^{(i)}$ , the likelihood function,  $L(\theta, \mathbf{x}^{(i)})$ , tells us how likely we would be to obtain that sample given the model parameter  $\theta$ .

Moreover, by assuming that our dataset  $\mathbf{X}$  is made up of  $N$  independent and identically distributed random variables. We define the joint likelihood function as the following,

$$L(\theta, \mathbf{X}) = \prod_{i=1}^N p_{\theta}(\mathbf{x}^{(i)}) \quad (2)$$

Now that we have the machinery to measure how likely it is to observe our data given the model parameter  $\theta$ . We want to find the value of  $\theta$  that maximises the likelihood of observing our data.

**Definition 2.2.** (Maximum likelihood estimator [4]) Given the data set  $\mathbf{X}$ , and parametric model  $p_{\theta}(\cdot)$ . The maximum likelihood estimate (MLE) is a value of  $\theta^*$  which maximises the likelihood function  $L(\theta, \mathbf{X})$ ,

$$\theta^* = \underset{\theta}{\operatorname{argmax}} L(\theta, \mathbf{X}) = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p_{\theta}(\mathbf{x}^{(i)})$$

This definitions produces the same outcome as our small example above, as we are maximising the probability of the data occurring given the parametric model (i.e minimising the difference between the parametric model and our data).

Up to now we are have only dealt with models that we are fully-observed as we are able to compute the maximum likelihood estimator directly from our data. Now suppose our model is made up with variables which we can not observe, and which are therefore not included in our dataset. These variables are what we refer to as latent variables, and which we will denote by  $\mathbf{z}$ .

Hence, to compute the likelihood as a function of  $\theta$  over our observed data, we define the marginal likelihood for the model as the following,

$$p_{\theta}(\mathbf{x}) = \int_{\Omega_{\mathbf{z}}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\Omega_{\mathbf{z}}} p_{\theta}(\mathbf{z}) p_{\theta}(\mathbf{x}|\mathbf{z}) d\mathbf{z} \quad (3)$$

Where  $p_\theta(\mathbf{x}, \mathbf{z})$  is the joint distribution over the observed and latent variables, and  $\Omega_{\mathbf{z}}$  is the domain of  $\mathbf{z}$ .

We can think of the marginal likelihood as "getting rid of what is unknown". Thus, given  $p_\theta(\mathbf{x}, \mathbf{z})$ , in principle we can work out maximum likelihood estimator  $\theta^*$ , as  $p_\theta(\mathbf{x})$  only contains variables which are observed. We will see later in Section 4.2.3 that the integral in (3) is sometimes intractable. Hence, we can not find the maximum likelihood parameter  $\theta^*$  directly. Therefore, we require some clever tricks to overcome this limitation.

## 2.2 Kullback–Leibler Divergence

The Kullback–Leibler Divergence  $D_{\text{KL}}$ , measures how much two probability distributions differ from each-other and is defined as the following,

**Definition 2.3** (Kullback–Leibler Divergence). Let  $p$  and  $q$  be two probability distributions densities, then the Kullback–Leibler Divergence  $D_{\text{KL}}$  is defined by,

$$D_{\text{KL}}(p||q) = \int_{\Omega_{\mathbf{x}}} p(x) \log \left( \frac{p(x)}{q(x)} \right) dx \quad (4)$$

An important remark is that by Gibbs' inequality (Appendix 1), we have  $D_{\text{KL}} \geq 0$  and equality iff  $p \equiv q$ .

For our purposes we will think of the quantity  $D_{\text{KL}}(p||q)$ , as a measure of how well the probability distribution  $p$  approximates  $q$ . This will be an important tool we will use in section 4.2.4, as we are unable to observe a distribution directly, so we require an approximation and hence we want to minimise difference between them.

### 2.2.1 Important Example

To illustrate this definition, we will compute the KL divergence of two Gaussians. This calculations will be important in Section 4.2.4, as we need to know the negative Kullback–Leibler Divergence between the two multivariate Gaussian distributions,  $q_\theta(\cdot) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$  and,  $p(\cdot) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  with dimension  $m$ .

To achieve this we firstly compute the two integrals,

$$\begin{aligned} \int q_\theta(\mathbf{z}) \log p(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}) \log \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) d\mathbf{z} \\ &= -\frac{m}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^m (\boldsymbol{\mu}_j^2 + \boldsymbol{\sigma}_j^2) \end{aligned}$$

and,

$$\begin{aligned} \int q_\theta(\mathbf{z}) \log q_\theta(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}) \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}) d\mathbf{z} \\ &= -\frac{m}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^m (1 + \log \boldsymbol{\sigma}_j^2) \end{aligned}$$

Therefore, we are able to combine these results and get the following,

$$\begin{aligned}
-D_{\text{KL}}(q_\phi(\cdot)||p_\theta(\cdot)) &= - \int q_\phi(\mathbf{z}) \log \left( \frac{q_\phi(\mathbf{z})}{p_\theta(\mathbf{z})} \right) d\mathbf{z} \\
&= \int q_\phi(\mathbf{z}) (\log p_\theta(\mathbf{z}) - \log q_\phi(\mathbf{z})) d\mathbf{z} \\
&= \frac{1}{2} \sum_{j=1}^m (1 + \log(\sigma_j^2) - \mu_j^2 - \sigma_j^2)
\end{aligned}$$

### 3 Machine Learning Basics

Throughout this section we will introduce the basic tools for solving modelling and prediction problems. However, we first must state some terminology that we will use throughout this section. Suppose that we were given the sets  $\mathbf{X}$  and  $\mathbf{Y}$ , and we want to predict  $\mathbf{Y}$  given  $\mathbf{X}$ , then we term  $\mathbf{X}$  is our data and  $\mathbf{Y}$  our observations respectively. The sets  $\mathbf{X}$  and  $\mathbf{Y}$  can take many of forms ranging from  $\mathbf{X}$  being the set of x-rays scans of the arm and  $\mathbf{Y}$  corresponding whether there is a fractured bone. To more simpler sets,  $\mathbf{X}$  being a set of real numbers and  $\mathbf{Y}$  corresponding to evaluations of an unknown linear function at each point of  $\mathbf{X}$ .

#### 3.1 Parametric models

As we described in the previous chapter, a parametric model is any model that is parameterised. These parameters are quantity that we can control and are used to describe what "features we can tune" in our model. This ability allows us to produce different model dynamics that we can use to find the parameters that best describes our data.

It is important to note that we have left the definition of a parametric model very open and loosely defined. This is on purpose as throughout this paper a parametric model may take a form of probability distributions, neural networks or functions.

We demonstrate the definition of a parametric model, by providing the following example. Suppose that our parametric model takes the form of,

$$f_\theta : \mathbb{R} \rightarrow \mathbb{R}, \quad f_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2, \quad \text{with } \theta \in \mathbb{R}^3 \quad (5)$$

We are able to see from Figure. (2), that the choice of the values of  $\theta_0, \theta_1, \theta_2$ , we are able to drastically change the dynamics of the model  $f_\theta$ .

#### 3.2 Loss Function

A loss function is a way of describing how well our parametric model preforms on given data point and is defined as the following,

$$L(\cdot, \cdot) : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$$

Where the inputs of the loss function is our parametric model,  $p_\theta$  evaluated at the data point  $\mathbf{x}^{(i)} \in \mathbf{X}$  and the corresponding observation,  $\mathbf{y}^{(i)} \in \mathbf{Y}$ . We are able

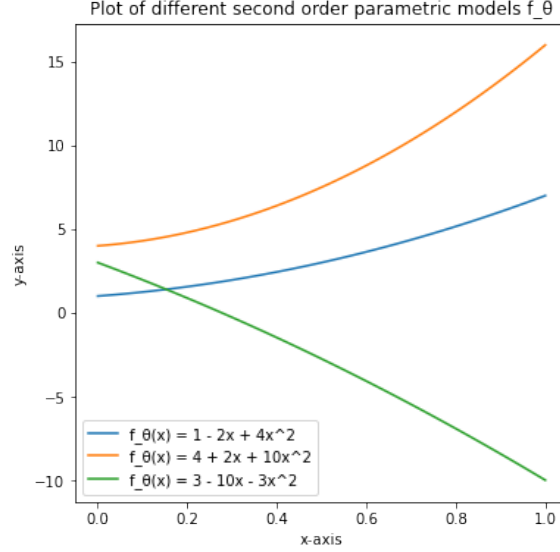


Figure 2: A plot of 3 different dynamics of the parametric models  $f_{\theta}$ , demonstrating how changing the value of  $\theta$  can produce different dynamics. (blue)  $f_{\theta}(x) = 1 - 2x + 4x^2$ . (orange)  $f_{\theta}(x) = 4 + 2x + 10x^2$ . (green)  $f_{\theta}(x) = 3 - 10x + 3x^2$ .

to think of the output of the loss function as a measure of how much information is lost by our parametric model given the data point  $\mathbf{x}^{(i)}$ .

To demonstrate this point, suppose that our data came from the function,  $f(x) = 2x + 1$ , and our parametric model (5), takes the form  $f_{\theta}(x) = 2x + 4$ . By choosing the loss function to be the squared error between the true observation and our parametric model prediction (i.e  $L(f(x), f_{\theta}) = ||f(x) - f_{\theta}(x)||^2$ ), we can see that our model loses is 9 units of information given a data point,  $x \in \mathbb{R}$ .

As we alluded to above, suppose that our dataset consists of multiply data points,  $\mathbf{X} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  and  $\mathbf{Y} = \{\mathbf{y}^{(i)}\}_{i=1}^N$ , respectively. We define the total loss of our model given  $\mathbf{X}$  and  $\mathbf{Y}$ , as the following,

$$E(\theta) = \sum_{i=1}^N L(\mathbf{y}^{(i)}, p_{\theta}(\mathbf{x}^{(i)})) \quad (6)$$

Where  $p_{\theta}$  is a parametric model, parameterized by  $\theta$ . This quantity allows us to understand how our model performance on the total data set, as our model might model well certain observations and poorly on others.

### 3.3 Gradient Descent

Now that we have the total loss function  $E(\theta)$ , for the parametric model  $p_{\theta}$ , we want to find the value of  $\theta$  that minimises  $E(\theta)$  (i.e  $\theta^* = \operatorname{argmin}_{\theta} E(\theta)$ ). We do this by introducing an interactive approach to find the local minima of  $E(\theta)$ . We first start by initialising the parameter  $\theta_0$  to a random values and follow the procedure,

$$\begin{aligned} \mathbf{v}_n &= \eta_n \nabla_{\theta} E(\theta) \\ \theta_{n+1} &= \theta_n - \mathbf{v}_n \end{aligned}$$



Where  $\nabla_{\theta}E(\theta)$  is the derivative of  $E(\theta)$  with respect to  $\theta$  and  $\eta_n$  is the learning rate at step  $n$ . The learning rate  $\eta_n$ , describes how big of a step in the direction  $-\nabla_{\theta}E(\theta)$  that we will take. We take steps in the direction of  $-\nabla_{\theta}E(\theta)$  as  $\nabla_{\theta}E(\theta)$  give us the direction of maximal increase of  $E(\theta)$ ; Therefore by going in inverse direction, we will go towards the minimum of  $E(\theta)$ . It is important that we get the value of  $\eta_n$  correct, as if it is too small the algorithm will take too long and if we take too big of a step it will over shoot the minimum. This effect is demonstrated in Figure. 3, as choosing different values of the learning rate can cause various effects when finding the minima. This procedure is normally terminated when the values  $\theta_{n+1}$  and  $\theta_n$  differ within a tolerance that we choose (i.e  $\|\theta_{n+1} - \theta_n\|^2 < \varepsilon$  with  $\varepsilon \ll 1$ ).

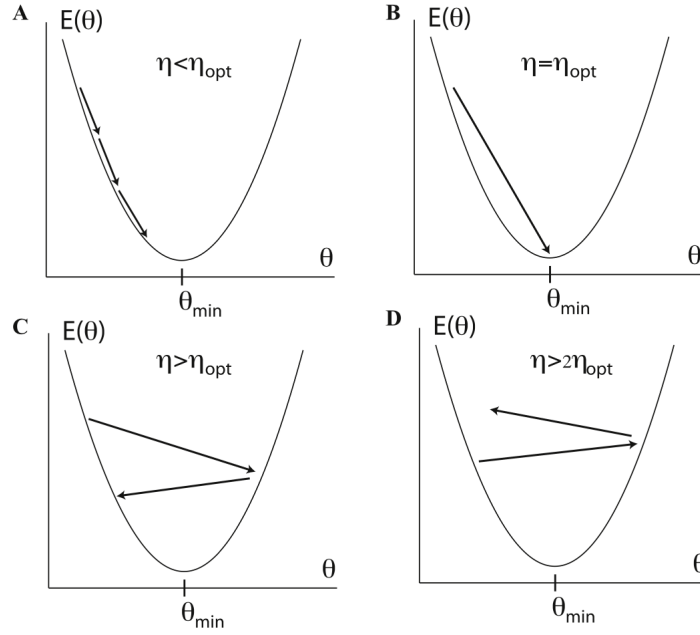


Figure 3: Effect of learning rate on convergence. For a one dimensional quadratic potential, one can show that there exists four different qualitative behaviours for gradient descent (GD) as a function of the learning rate  $\eta$  depending on the relationship between  $\eta$  and  $\eta_{\text{opt}} = [\partial_{\theta}^2 E(\theta)]^{-1}$ . (a) For  $\eta < \eta_{\text{opt}}$ , GD converges to the minimum. (b) For  $\eta = \eta_{\text{opt}}$ , GD converges in a single step. (c) For  $\eta_{\text{opt}} < \eta < 2\eta_{\text{opt}}$ , GD oscillates around the minima and eventually converges. (d) For  $\eta > 2\eta_{\text{opt}}$ , GD moves away from the minima. ([5], p.15)

### 3.4 Training

Up to now we have only dealt with fitting a model to a given data set. However one of the main advantages of modelling data comes from its ability to model data that we have not seen before. We formalise this notion by firstly splitting our data into two sets,  $\mathbf{X}_{\text{train}}$  and  $\mathbf{X}_{\text{test}}$ , with corresponding  $\mathbf{Y}_{\text{train}}$  and  $\mathbf{Y}_{\text{test}}$  respectively. This splitting allows us to train our model using  $\mathbf{X}_{\text{train}}$  and  $\mathbf{Y}_{\text{train}}$  then test the performance of our model on data that it has not seen before (i.e  $\mathbf{X}_{\text{test}}$  and  $\mathbf{Y}_{\text{test}}$ ). We measure this performance by defining the in-sample error  $E_{\text{in}} = E(\theta|\mathbf{X}_{\text{train}}, \mathbf{Y}_{\text{train}})$  and out-of-sample error  $E_{\text{out}} = E(\theta|\mathbf{X}_{\text{test}}, \mathbf{Y}_{\text{test}})$ .

The out-of-sample error  $E_{\text{out}}$ , will naturally decrease as the number of data points increase. This due to there being less sampling noise that appears in the data, and our model tending to its true distribution. As the model tends to its true distribution from where the data is drawn, the in-sample and out-of-sample error must therefore converge to the same value; This is termed the "bias" of our model. The bias tells us how good of a model we have if we had infinite data point. However, since in piratical application we only have a finite amount of data points, it is better to minimise the out-of-sample error  $E_{\text{out}}$  rather than the bias; As this will lead to a model which has better predicting abilities. This is due to the fact that the difference between  $E_{\text{in}}$  and  $E_{\text{out}}$  is the same as fitting and predicting new data. If the difference between the in-sample error  $E_{\text{in}}$ , is much smaller than out-of-sample error  $E_{\text{out}}$ , we know that we over-fitted our model, as performs poorly on new data.

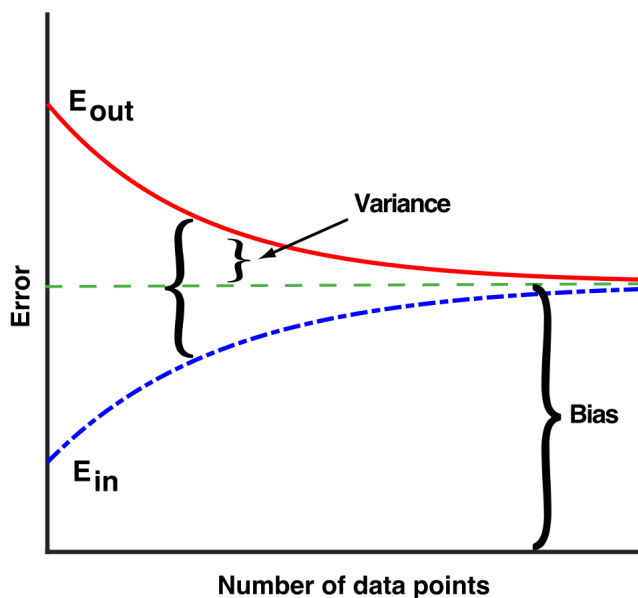


Figure 4: Schematic of typical in-sample and out-of-sample error as a function of training set size. The typical in-sample or training error,  $E_{\text{in}}$ , out-of-sample or generalisation error,  $E_{\text{out}}$ , bias, variance, and difference of errors as a function of the number of training data points. ([5], p.11)

Figure. 5, shows that the out-of-sample  $E_{\text{out}}$  is a function of "model complexity", which is how many parameters the model has. This supports the idea of over-fitting that we discussed above. As we expect increasing the number of parameters in our model we over-fit our model to the training data and hence, perform poorly on the at predicting new data. This is due to there being a finite amount of data points, even though increasing the model complexity will lead to a smaller bias due to a smaller in-sample error. Therefore it is a more favourable to choose a model with smaller variance than a less-biased model with large variance.

Finding the minima of a function is not inherently unique to this problem. Therefore, gradient decent with not the only algorithm that we can use to find the minima, with each algorithm offering their own take and limitations. To demonstrate this

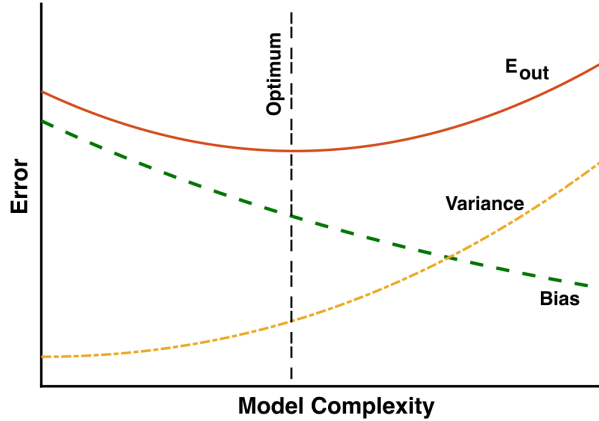


Figure 5: Bias–Variance tradeoff and model complexity. This schematic shows the typical out-of-sample error  $E_{\text{out}}$  as function of the model complexity for a training dataset of fixed size. ([5], p.12)

we will introduce RMSprop iteration, that is a form of adaptive learning method,

$$\begin{aligned} \mathbf{g}_t &= \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) \\ \mathbf{s}_t &= \beta \mathbf{s}_{t-1} + (1 - \beta) \mathbf{g}_t^2 \\ \boldsymbol{\theta}_{t+1} &= \boldsymbol{\theta}_t - \eta_t \frac{\mathbf{g}_t}{\sqrt{\mathbf{s}_t + \varepsilon}} \end{aligned}$$

Where  $\beta$  controls the averaging time of the second moment,  $\eta_t$  is the learning rate, and  $\varepsilon$  is a small regularisation constant.

Another way that we might control over-fitting is the use of epochs and batches. Epochs are the number of times in the we use the whole training data  $\mathbf{X}_{\text{train}}$  and  $\mathbf{Y}_{\text{train}}$  in training. Whilst batches split up the training data  $\mathbf{X}_{\text{train}}$  into smaller training data sets that we minimise the total cost over.

To demonstrate this idea suppose we we have a training data set consisting of 100 data points, and apply five epochs of batch sizes 25. Therefore we use the total data set five times, and each round of epoch contrasting of four evaluations of the training algorithm on the data set,

$$E_{\text{in}}^i = E(\boldsymbol{\theta} | \mathbf{X}_{\text{train}}^i, \mathbf{Y}_{\text{train}}^i), \quad \mathbf{X}_{\text{train}}^i = \sum_{j=25i+1}^{25(i+1)} \mathbf{x}^{(j)}, \quad \mathbf{Y}_{\text{train}}^i = \sum_{j=25i+1}^{25(i+1)} \mathbf{y}^{(j)}, \quad i \in \{0, 1, 2, 3\}$$

### 3.5 Theory Example

We will now combine the theory that we have introduced into a small example to demonstrate how each part plays a key role of solving a modelling problem.

Suppose that our data looked like it came from a quadratic distribution. We test this hypothesis by using the following quadratic parametric model,  $f_{\boldsymbol{\theta}}(x) = \theta_0 + \theta_1 x + \theta_2 x^2$ , and then trying to learn the value of  $\boldsymbol{\theta}$  that best fits  $f_{\boldsymbol{\theta}}$  to our

data. Finally we choose the squared error loss function to measure this difference and minimise over

$$\begin{aligned}
\theta^* &= \underset{\theta}{\operatorname{argmin}} E(\theta) \\
&= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N L(y^{(i)}, f_{\theta}(x^{(i)})) \\
&= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \|y^{(i)} - f_{\theta}(x^{(i)})\|^2 \\
&= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \|y^{(i)} - (\theta_0 + \theta_1 x^{(i)} + \theta_2 (x^{(i)})^2)\|^2
\end{aligned}$$

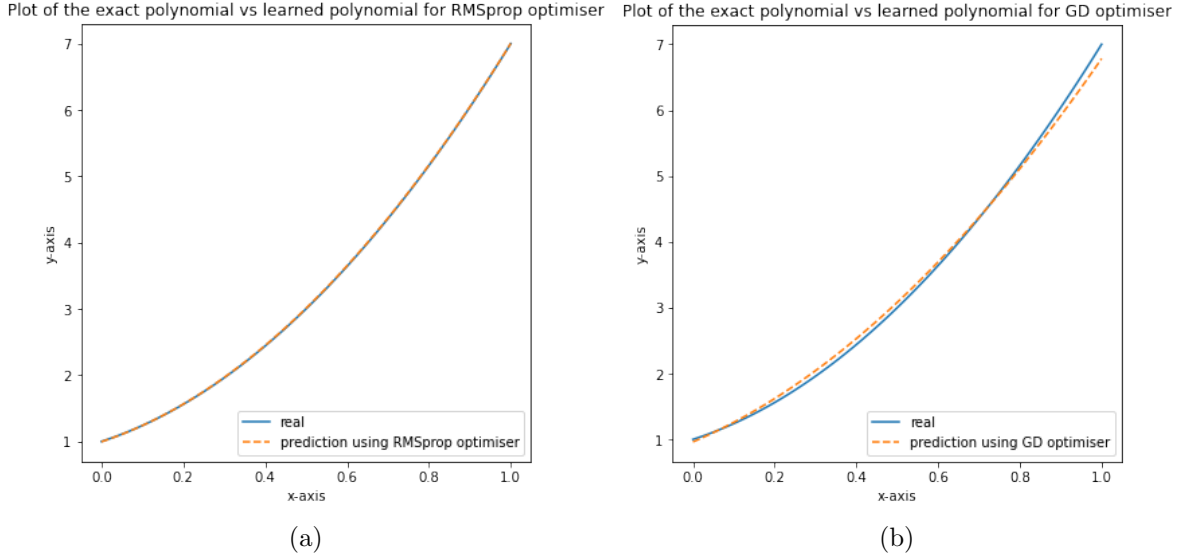


Figure 6: Plot of the exact polynomial  $f(x) = 1 + 2x + 4x^2$  vs learned polynomial for GD and RMSprop optimiser over the interval  $[0, 1]$  trained using 19200 data points, learning rates of 0.01, 4 epochs, and batch-sizes of 64. (blue) exact polynomial (orange) learned polynomial (a) Plot of the exact vs learned polynomial  $f_{\theta}(x) = 0.996 + 2.045x + 3.963x^2$  for the RMSprop optimiser. (b) Plot of the exact vs learned polynomial  $f_{\theta}(x) = 0.962 + 2.656x + 3.163x^2$  for the GD optimiser. [3]

Figure. 6, shows us that RMSprop performs well at approximating the true values of  $\theta$  whilst GD struggles slightly. This further emphasise how challenging modelling data can be, as choosing the parameters such as learning algorithm, number of data points, learning rate, epochs and batch-sizes, has to be done on a case to case basis.

### 3.6 Neural Networks Architecture

The deep neural network architecture is defined as the following,

$$\mathbf{a}^l = \sigma(\mathbf{W}^l \mathbf{a}^{l-1} + \mathbf{b}^l) =: \sigma(\mathbf{z}^l), \quad l \in \{1, \dots, L\}, \quad L \in \mathbb{N} \quad (7)$$

$$n_l \in \mathbb{N}, \quad \mathbf{a}^l \in \mathbb{R}^{n_l}, \quad l \in \{0, \dots, L\} \quad (8)$$

$$\mathbf{W}^l \in \mathbb{R}^{n_l \times n_{l-1}}, \quad \mathbf{b}^l \in \mathbb{R}^{n_l}, \quad l \in \{1, \dots, L\} \quad (9)$$

We refer to  $\mathbf{a}^l$  as an activation layer. This is because this quantity represents how much "activation" each neuron/node gets when it is fed into the next layer.  $\mathbf{W}^l$  is the weight matrix, as it represents how much of weight each neuron of the previous activation layer has on the next layer (i.e by the matrix-multiplication  $\mathbf{W}^l \mathbf{a}^{l-1}$ ).  $\mathbf{b}^l$  is the bias (or bias vector) as it shifts the each value of  $\mathbf{W}^l \mathbf{a}^{l-1}$  either to the left (upwards) or right (downwards). Finally,  $\sigma(\cdot)$  is refereed to as the activation function, as it output how much "activation" is carried on to the next layer.

Since,  $\mathbf{W}^l$  and  $\mathbf{b}^l$  are not explicated defined, a neural network can be though of as a parametric model. Therefore we are apply all the theory we have just introduced to them. This is not the only way we can define neural network, as they can come in many of forms such as, Recurrent Neural Networks [6] or Convolutional Neural Networks [7].

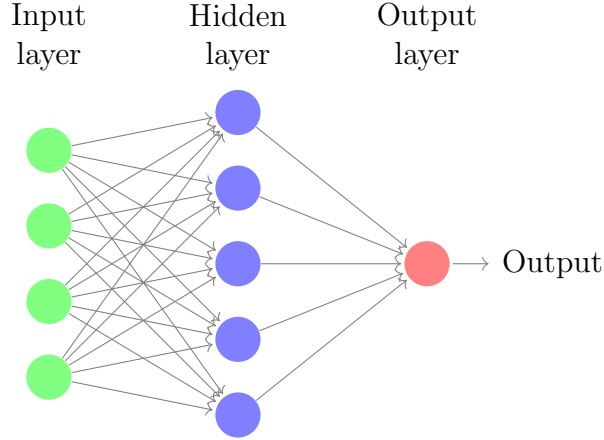


Figure 7: A deep neural network with an input dimension of 4 ( $n_0 = 4$ ), a hidden state of dimension of 5 ( $n_1 = 5$ ), and an output dimension of 1 ( $n_2 = 1$ ).

### 3.7 Backpropagation Algorithm

The backpropagation algorithm for a deep neural network allows us take gradients with respect to the weights and biases to find the minima of our loss function,  $C$ . Since the loss function depends on the output layer of our neural network, we define the error at the  $j$ -th neuron of the output layer, as the following,

$$\Delta_j^L = \frac{\partial C}{\partial \mathbf{z}_j^L} \quad (10)$$

Where  $L$  denotes the output layer of our neural network, defined in (7). Moreover as the error of our output indirectly depends on the previous layers of our network. We can generalise (10), to get the error at the  $j$ -th neuron on the  $l$ -th layer,  $\Delta_j^l$ , as the change in the loss function with respect to  $\mathbf{z}_j^l$ .

$$\Delta_j^l = \frac{\partial C}{\partial \mathbf{z}_j^l} = \frac{\partial C}{\partial \mathbf{a}_j^l} \frac{\partial \mathbf{a}_j^l}{\partial \mathbf{z}_j^l} = \frac{\partial C}{\partial \mathbf{a}_j^l} \sigma'(\mathbf{z}_j^l) \quad (11)$$

Where  $\sigma'(\mathbf{z}_j^l)$  is the derivative of the activation function  $\sigma(\cdot)$ , evaluated at  $\mathbf{z}_j^l$ . It is important to note that the dimension of the activation function's derivative may change per layer.

We can also express (11) in terms of the partial derivative of the loss function with respect to the bias, by noting that  $\partial \mathbf{b}_j^l / \partial \mathbf{z}_j^l = \mathbf{1}$ ,

$$\Delta_j^l = \frac{\partial C}{\partial \mathbf{z}_j^l} = \frac{\partial C}{\partial \mathbf{b}_j^l} \frac{\partial \mathbf{b}_j^l}{\partial \mathbf{z}_j^l} = \frac{\partial C}{\partial \mathbf{b}_j^l} \quad (12)$$

Hence by using the chain rule again, we are able to express the error at the  $l$ -th layer in terms of the  $l+1$ -th layer. This is because error can be through of as being carried forward through the neural network.

$$\Delta_j^l = \frac{\partial C}{\partial \mathbf{z}_j^l} = \sum_k \frac{\partial C}{\partial \mathbf{z}_k^{l+1}} \frac{\partial \mathbf{z}_k^{l+1}}{\partial \mathbf{z}_j^l} \quad (13)$$

$$= \sum_k \Delta_k^{l+1} \frac{\partial \mathbf{z}_k^{l+1}}{\partial \mathbf{z}_j^l} = \left( \sum_k \Delta_k^{l+1} \mathbf{W}_{kj}^{l+1} \right) \sigma(\mathbf{z}_j^l) \quad (14)$$

Finally, we derive the derivative of the loss function with respect to the weight  $\mathbf{W}_{jk}^l$ ,

$$\frac{\partial C}{\partial \mathbf{W}_{jk}^l} = \frac{\partial C}{\partial \mathbf{z}_j^l} \frac{\partial \mathbf{z}_j^l}{\partial \mathbf{W}_{jk}^l} = \Delta_j^l \mathbf{a}_k^{l-1} \quad (15)$$

---

**Algorithm 1** Backpropagation algorithm ([5] p.55)

---

- 1: **Feedforward:** Start with the first layer, exploit the feed-forward architecture through (7) to compute  $\mathbf{z}^l$  and  $\mathbf{a}^l$  for each subsequent layer.
  - 2: **Error at output layer:** Calculate the error of the output layer using (11).
  - 3: **"Backpropagate" the error:** We now use (14) to propagate the error backwards and calculate  $\Delta_j^l$  for all the layers.
  - 4: **Calculate gradient:** Now use (12) and (15) to calculate  $\frac{\partial C}{\partial \mathbf{b}_j^l}$  and  $\frac{\partial C}{\partial \mathbf{W}_{jk}^l}$ .
- 

### 3.8 Generative Modelling

Generative modelling is a way that we can produce new data. It does this by learning the joint probability  $p(\mathbf{X}, \mathbf{Y})$  or just  $p(\mathbf{X})$  if there are no observations  $\mathbf{Y}$  given. Therefore by sampling from its learned distribution, we are also produce samples would of likely came from our test data  $\mathbf{X}$ .



Figure 8: Samples from a generative model trained on Frey's face. [8]

## 4 Variational Autoencoders

### 4.1 Autoencoder

Before we understand variational autoencoders (VAEs), we firstly must understand their primitive cousin, autoencoders. Autoencoders encapsulate the idea that data can be expressed by a small number of key features, whilst the rest is noise. For example, if you were to communicate the key ideas of a face to someone, the main features would be the position and shape of the eyes, mouth, and nose, whilst the rest of the details like facial hair and piercings are noise.

This feature selection process is expressed formally by performing dimensional reduction to the input data. This can be done in various ways and we will refer to this action as the Encoder. We will also term the output of the Encoder the feature/latent space, as it expresses the key features that make up our input data. The autoencoder then projects the output of the Encoder back to the space containing the data, and we will refer to this operation as the Decoder. As we have assumed that our data can be represented by a small number of features, the aim is to minimise the difference between our input and reconstructed data. We will use the squared error loss function to quantify how accurate our autoencoder is at reconstructing the input data. To summarise, we have the following model,

$$(\text{Encoder}) \ \phi_E : \mathbb{R}^n \rightarrow \mathbb{R}^m, \quad (\text{Decoder}) \ \phi_D : \mathbb{R}^m \rightarrow \mathbb{R}^n, \quad m \ll n, \ m, n \in \mathbb{N}$$

$$(\text{L2 loss function}) \ L = \|\mathbf{x} - \hat{\mathbf{x}}\|^2 = \|\mathbf{x} - \phi_D(\phi_E(\mathbf{x}))\|^2$$

As we alluded to earlier, both the Encoder and Decoder can take many forms with the simplest case being linear maps, however in this paper we will only look at them being neural networks. Therefore, the autoencoder can be thought of as a single neural network, as it is just a combination of two smaller neural networks sequentially. By ensuring that  $m \ll n$ , our architecture creates a bottleneck for the data, allowing the key features to be used for the reconstruction. Hence, by

performing gradient descent, we are able to find the parameters that minimise the loss function for our data.

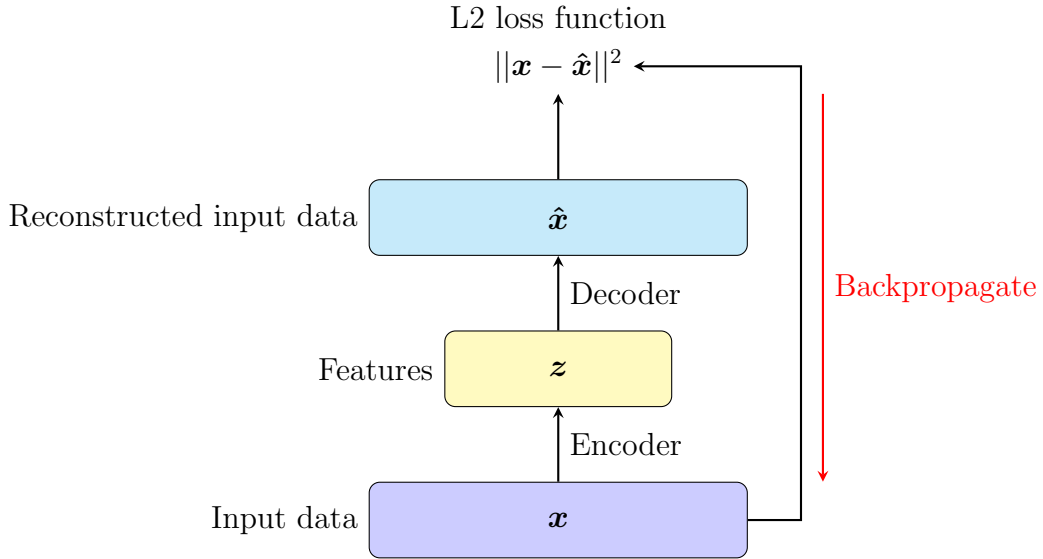


Figure 9: Training model for an autoencoder

However, this approach leads to the limitation that two points close in the feature space, may not necessary corresponds to two similarly points in data space. This renders the autoencoder model useless for generative modelling, as we want points that are close in feature space to be mapped to points close in data space because we expect them to have "similar qualities".

## 4.2 Variational Autoencoders

### 4.2.1 Motivation

Variational autoencoders are a probabilistic approach of autoencoders that aims to overcome the irregularity that appear between the feature and data space in autoencoders. This is achieved by mapping our input data to a distribution over the feature space rather than to a single point, we then sample from a decoder distribution to get our reconstructed input data. This procedure allows us to introduce regularisation at the time of training to prevent overfitting and ensures that our feature spaces has good properties for generative modelling.

Putting together the model described above, the variational autoencoder evauled at the point  $\mathbf{x}$  takes the form:

1. Our input data  $\mathbf{x}$  is encoded as a distribution  $p_{\theta}(\cdot|\mathbf{x})$  in feature space.
2. We then sample the point  $\mathbf{z}$  from the distribution  $p_{\theta}(\cdot|\mathbf{x})$ .
3. Then  $\mathbf{z}$  is decoded into the distribution  $p_{\theta}(\cdot|\mathbf{z})$ .
4. Finally,  $\hat{\mathbf{x}}$  is then drawn from the decoder distribution  $p_{\theta}(\cdot|\mathbf{z})$ .



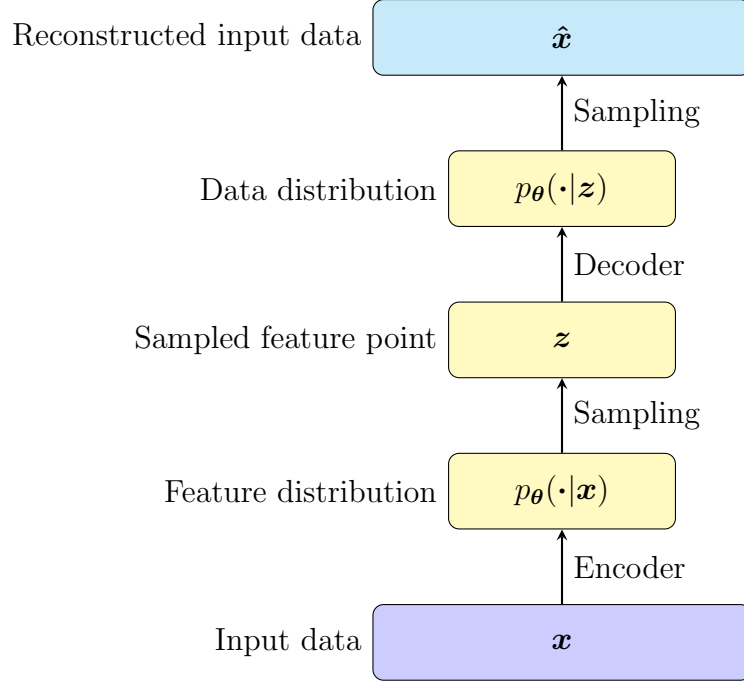


Figure 10: Graphical representation of a variational autoencoder evaluated at the point  $\mathbf{x}$ .

#### 4.2.2 Probabilistic Framework

To support the above claims, we will create a probabilistic framework. Variational autoencoders build off the assumption that our dataset  $\mathcal{D} = (\mathbf{x}^{(i)})_{i=1}^N$ , was generated by  $N$  independent samples of an unchanging underlying distribution. Where each sample  $\mathbf{x}^{(i)}$ , was generated from some random process that involves the continuous unobserved latent random variable  $\mathbf{z}^{(i)}$ . This process happens in two distinct steps:

1. The value  $\mathbf{z}^{(i)}$  is generated from some prior distribution  $p_{\theta^*}(\mathbf{z})$
2. Then  $\mathbf{x}^{(i)}$  is generated from some conditional distribution  $p_{\theta^*}(\mathbf{x}|\mathbf{z}^{(i)})$

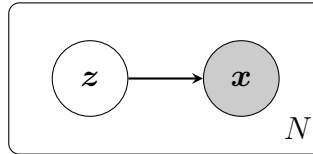


Figure 11: Probabilistic graphical model for our data generation process

Where  $\theta^*$ , are the true parameters of the parametric families of differentiable distributions,  $p_{\theta}(\mathbf{z})$  and  $p_{\theta}(\mathbf{x}|\mathbf{z})$ . This in turn allows us to define the joint parametric distribution  $p_{\theta}(\mathbf{x}, \mathbf{z})$ , as  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})$ .

#### 4.2.3 The Problem With This Approach

Therefore by using the the joint parametric distribution above  $p_{\theta}(\mathbf{x}, \mathbf{z}) = p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z})$ , the marginal likelihood is defined as the following,

$$p_{\theta}(\mathbf{x}) = \int_{\Omega_{\mathbf{z}}} p_{\theta}(\mathbf{x}, \mathbf{z}) d\mathbf{z} = \int_{\Omega_{\mathbf{z}}} p_{\theta}(\mathbf{x}|\mathbf{z})p_{\theta}(\mathbf{z}) d\mathbf{z}$$

So in principle we are able find the optimal parameter  $\boldsymbol{\theta}^*$  by applying GD to the likelihood given data set. However, due to integral being intractable as it does not have an analytical solution or an efficient estimator. We are unable take derivative with respect to the parameter  $\boldsymbol{\theta}$  and find  $\boldsymbol{\theta}^*$  that maximises  $p_{\boldsymbol{\theta}}(\mathbf{x})$ .

The intractability of  $p_{\boldsymbol{\theta}}(\mathbf{x})$ , is related to the intractability of the unobserved posterior distribution  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ , by the relationship,

$$p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}) = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{x})} = \frac{p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})}{\int_{\Omega_{\mathbf{z}}} p_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z}) d\mathbf{z}}$$

Therefore, if we have a tractable  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$  (or  $p_{\boldsymbol{\theta}}(\mathbf{x})$  respectably) this will lead to a tractable  $p_{\boldsymbol{\theta}}(\mathbf{x})$  (or  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$  respectably), as  $p_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{z})$  is tractable due to it being modelled as the product of two known parametric distribution.

#### 4.2.4 Overcoming This Problem

To overcome this problem we try to approximate the true unobserved posterior distribution  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$  by using the encoder model  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$ , parameterised by  $\boldsymbol{\phi}$ . We aim to finding the value for  $\boldsymbol{\phi}$  such that the distribution  $q_{\boldsymbol{\phi}^*}(\mathbf{z}|\mathbf{x})$ , closely approximates the true intractable posterior  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ . This also equivalent to solving the minimisation problem,

$$\boldsymbol{\phi}^* = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} D_{\text{KL}}(q_{\boldsymbol{\phi}}(\cdot|\mathbf{x})||p_{\boldsymbol{\theta}}(\cdot|\mathbf{x})). \quad (16)$$

Since we have assumed that each data point  $\mathbf{x}^{(i)}$  are independent, the total marginal likelihood is the sum of the individual likelihoods. Hence, by using our encoder model, we will compute the following marginal likelihood,

$$\begin{aligned} \log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) &= \mathbb{E}_{\mathbf{z} \sim q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} [\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})] \quad (\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \text{ does not depend on } \mathbf{z}) \\ &= \mathbb{E}_{\mathbf{z}} \left[ \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})} \right] \quad (\text{Bayes' Rule twice}) \\ &= \mathbb{E}_{\mathbf{z}} \left[ \log \frac{p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})} \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})}{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})} \right] \quad (\text{Multiplication by 1}) \\ &= \mathbb{E}_{\mathbf{z}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})] - \mathbb{E}_{\mathbf{z}} \left[ \log \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})}{p_{\boldsymbol{\theta}}(\mathbf{z})} \right] + \mathbb{E}_{\mathbf{z}} \left[ \log \frac{q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})}{p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})} \right] \quad (\text{Logs}) \\ &= \underbrace{\mathbb{E}_{\mathbf{z}} [\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})] - D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z}))}_{\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \boldsymbol{\phi})} + \underbrace{D_{\text{KL}}(q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)}))}_{\geq 0} \end{aligned}$$

The true intractable posterior  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$  still appears in the last term of the marginal likelihood. However, we are able to use the fact that the Kullback–Leibler divergence is non-negative and zero if and only if  $q_{\boldsymbol{\phi}}(\mathbf{z}|\mathbf{x})$  equals the true posterior, and obtain an evidence lower bound (ELBO),  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \boldsymbol{\phi})$ , for our marginal likelihood.

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

Hence, we aim to maximise  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \boldsymbol{\phi})$  through gradient descent and obtain the optimal parameters  $\boldsymbol{\theta}^*$  and  $\boldsymbol{\phi}^*$  that maximises the evidence lower bound.

$$\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) \geq \mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}^*, \boldsymbol{\phi}^*) \geq \mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \boldsymbol{\phi})$$

We will now look at the terms in  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$ , as it will allow us to understand what the ELBO is maximising for:

$\mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})]$ , can be thought of as how well the model done at reconstructing the input data. Therefore maximising for  $\boldsymbol{\theta}$  will give use a good reconstruction of our input data.

$D_{KL}(q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\mathbf{z}))$ , can be though as regularising  $\phi$  such that  $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$  approximate the distribution  $p_{\boldsymbol{\theta}}(\mathbf{z})$ . It is not entirely obvious that this will yield the correct  $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$  that we set out to obtain; However, if we fix  $\boldsymbol{\theta}$  and think about finding the value of  $\phi$  that minimises the KL-divergence between  $q_{\phi}(\cdot|\mathbf{x})$  and  $p_{\boldsymbol{\theta}}(\cdot|\mathbf{x})$ , we get the following,

$$\begin{aligned}
\phi^* &= \underset{\phi}{\operatorname{argmin}} (D_{KL}(q_{\phi}(\cdot|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\cdot|\mathbf{x}^{(i)}))) \quad (\text{Using the notation } \mathbb{E}_{\mathbf{z}}(\cdot) \equiv \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})}(\cdot)) \\
&= \underset{\phi}{\operatorname{argmin}} (\mathbb{E}_{\mathbf{z}}[\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})] - \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x}^{(i)})]) \quad (\text{Logs}) \\
&= \underset{\phi}{\operatorname{argmin}} \left( \mathbb{E}_{\mathbf{z}}[\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})] - \mathbb{E}_{\mathbf{z}} \left[ \frac{p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})p_{\boldsymbol{\theta}}(\mathbf{z})}{p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})} \right] \right) \quad (\text{Bayes' Rule twice}) \\
&= \underset{\phi}{\operatorname{argmin}} (\mathbb{E}_{\mathbf{z}}[\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})] - \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{z})] - \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})] + \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})]) \quad (\text{Logs}) \\
&\text{Since, } \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})] \text{ is independent of } \phi, \text{ we can remove it without effecting the argmin.} \\
&= \underset{\phi}{\operatorname{argmin}} (\mathbb{E}_{\mathbf{z}}[\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})] - \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{z})] - \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})]) \\
&= \underset{\phi}{\operatorname{argmax}} (-\mathbb{E}_{\mathbf{z}}[\log q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})] + \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{z})] + \mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})]) \quad (\text{argmin to argmax}) \\
&= \underset{\phi}{\operatorname{argmax}} (\mathbb{E}_{\mathbf{z}}[\log p_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}|\mathbf{z})] - D_{KL}(q_{\phi}(\cdot|\mathbf{x}^{(i)})||p_{\boldsymbol{\theta}}(\cdot))) \quad (D_{KL} \text{ definition})
\end{aligned}$$

Hence, we arrive at the same expression as  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$ . Therefore, maximising  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$  for  $\phi$  will produce the correct  $\phi^*$ .

Therefore, by maximising  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$ , we have the maximal trade-off between good reconstruction of our input data and  $q_{\phi}(\mathbf{z}|\mathbf{x}^{(i)})$  approximating the true intractable posterior,  $p_{\boldsymbol{\theta}}(\mathbf{z}|\mathbf{x})$ .

#### 4.2.5 Reparameterization Trick

Now that we have a computable lower bound for the marginal likelihood, we want to differentiate  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$  with respect to  $\boldsymbol{\theta}$  and  $\phi$  to find their optimal parameters  $\boldsymbol{\theta}^*$  and  $\phi^*$ , such that  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}^*, \phi^*) \geq \mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$ . However, we are unable to take the able to take the derivative of  $p_{\boldsymbol{\theta}}(\cdot|\mathbf{z})$  as depends on a random and non-deterministic sample  $\mathbf{z} \sim q_{\phi}(\cdot|\mathbf{x}^{(i)})$ . We overcome this by choosing our encoder model  $q_{\phi}(\cdot|\mathbf{x}^{(i)})$  under certain conditions (explained below) such that we are able to reparameterise the random variable  $\mathbf{z}$ , as a differentiable transformation function  $g_{\phi}(\boldsymbol{\varepsilon}, \mathbf{x}^{(i)})$ , where  $\boldsymbol{\varepsilon}$  is some auxiliary noise.

$$\mathbf{z} = g_{\phi}(\boldsymbol{\varepsilon}, \mathbf{x}^{(i)}), \quad \boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon}) \quad (17)$$

We will now state the conditions where we can choose a differentiable transformation  $g_\phi(\boldsymbol{\varepsilon}, \mathbf{x}^{(i)})$  and auxiliary variable  $\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})$  to replace  $\mathbf{z} \sim q_\phi(\cdot|\mathbf{x})$  for  $\mathbf{z} = g_\phi(\boldsymbol{\varepsilon}, \mathbf{x}^{(i)})$ . [8]

1.  $q_\phi(\cdot|\mathbf{x})$  has a tractable inverse CDF. In this case, let  $\boldsymbol{\varepsilon} \sim \mathcal{U}(\mathbf{0}, \mathbf{I})$ , and let  $g_\phi(\boldsymbol{\varepsilon}, \mathbf{x})$  be the inverse CDF of  $q_\phi(\cdot|\mathbf{x})$ . Examples: Exponential, Cauchy, Logistic, Rayleigh, Pareto, Weibull, Reciprocal, Gompertz, Gumbel and Erlang distributions.
2.  $q_\phi(\cdot|\mathbf{x})$  is from any "location-scale" family of distributions. In this we can choose the standard distribution (with location = 0, scale = 1) as the auxiliary variable  $\boldsymbol{\varepsilon}$ , and then let  $g_\phi(\boldsymbol{\varepsilon}, \mathbf{x}) = \text{location} + \text{scale} \times \boldsymbol{\varepsilon}$ . Examples: Laplace, Elliptical, Student's t, Logistic, Uniform, Triangular and Gaussian distributions.
3.  $q_\phi(\cdot|\mathbf{x})$  is a composition of random variables as different transformations of auxiliary variables. Then we take  $g_\phi(\boldsymbol{\varepsilon}, \mathbf{x})$  their composition of the above reparameterizations. Examples: Log-Normal (exponentiation of normally distributed variable), Gamma (a sum over exponentially distributed variables), Dirichlet (weighted sum of Gamma variates), Beta, Chi-Squared, and F distributions.

Hence, by using the reparameterization (17), we are able to form a Monte Carlo estimate of expectation for the function  $\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})$  with respect to our encoder model  $q_\phi(\cdot|\mathbf{x}^{(i)})$ , as the following,

$$\mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot|\mathbf{x}^{(i)})}[\log p_\theta(\mathbf{x}^{(i)}|\mathbf{z})] = \mathbb{E}_{\boldsymbol{\varepsilon} \sim p(\boldsymbol{\varepsilon})}[\log p_\theta(\mathbf{x}^{(i)}|g_\phi(\boldsymbol{\varepsilon}, \mathbf{x}^{(i)}))] \quad (18)$$

$$\simeq \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}|g_\phi(\boldsymbol{\varepsilon}^{(l)}, \mathbf{x}^{(i)})), \quad \boldsymbol{\varepsilon}^{(l)} \sim p(\boldsymbol{\varepsilon}) \quad (19)$$

Now applying the method above to our ELBO, we obtain a differentiable estimator  $\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi) \simeq \tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$ .

$$\tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi) = \frac{1}{L} \sum_{l=1}^L \log p_\theta(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}^{(i)})||p_\theta(\mathbf{z})) \quad (20)$$

$$\text{where } \mathbf{z}^{(i,l)} = g_\phi(\boldsymbol{\varepsilon}^{(i,l)}, \mathbf{x}^{(i)}), \quad \text{and } \boldsymbol{\varepsilon}^{(l)} \sim p(\boldsymbol{\varepsilon})$$

Now suppose our dataset  $\mathbf{X}$  consisting of  $N$  data points. We are able to preformed minibatch training to the above result, to get the following approximation of  $\mathcal{L}$ ,

$$\mathcal{L}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi) \simeq \tilde{\mathcal{L}}^M(\mathbf{X}^M; \boldsymbol{\theta}, \phi) = \frac{N}{M} \sum_{i=1}^M \tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$$

Where,  $\mathbf{X}^M = (\mathbf{x}^{(i)})_{i=1}^M$  are  $M$  randomly chosen data points from our training data  $\mathbf{X}$ . If we perform minibatch training with  $M$  is large enough (i.e.  $M = 100$ ), we are able to set  $L = 1$  in (20). As experiments show that there is no significant performance decrease by taking a single sample [8]. This allows us to save time training as we don't have to preform multiple samples for each data point, decreasing the complexity of training.

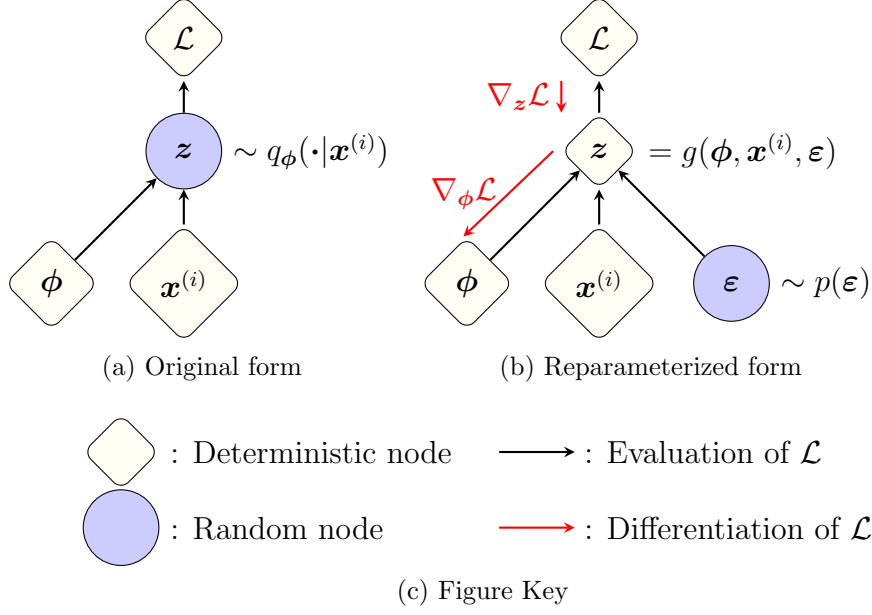


Figure 12: Illustration of the reparameterization trick used to get a differentiable  $\mathcal{L}$  w.r.t  $\phi$ . ([9], p.22)

### 4.3 Gaussian VAE Example

We will now demonstrate this reparameterization trick, by assuming that our prior over the latent variables, is given by an isotropic multivariate Gaussian  $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$  and our true posterior  $p_{\theta}(\mathbf{z}|\mathbf{x})$  is a multivariate Gaussian. Therefore, we will approximate the true posterior  $p_{\theta}(\mathbf{z}|\mathbf{x})$  with a multivariate parametric Gaussian model that has a diagonal covariance,  $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I})$ . Where  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}^2$  are function of  $\phi$ . We chose a diagonal covariance structure to promote the latent variables learning unique generalise features.

Since our parametric model  $q_{\phi}(\mathbf{z}|\mathbf{x})$  is a "location-scale" distribution, we will use the following reparameterisation,

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

Where  $\odot$  is element-wise multiplication. Therefore, by using the calculation that we carried out in Section 2.2, we obtain the following approximation for  $\tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi)$ ,

$$\tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi) \simeq \frac{1}{2} \sum_{j=1}^m (1 + \log((\boldsymbol{\sigma}_j^{(i)})^2) - (\boldsymbol{\mu}_j^{(i)})^2 - (\boldsymbol{\sigma}_j^{(i)})^2) + \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)})$$

$$\text{where } \mathbf{z}^{(i,l)} = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\varepsilon}^{(l)}, \quad \boldsymbol{\varepsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (21)$$

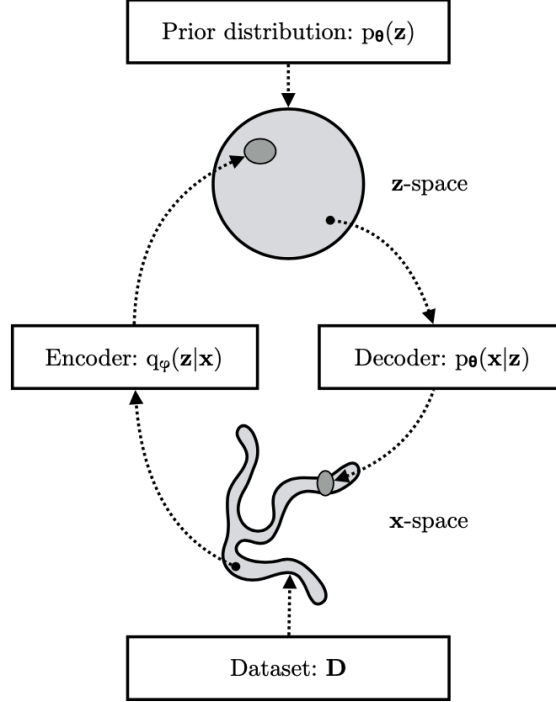


Figure 13: A VAE learns the mappings between our observed data space which its distribution which is typically complicated, and a feature space, whose distribution can be relatively simple (such as Gaussian, as in this figure). ([9], p.17)

#### 4.4 Application of VAE to MNIST Dataset

In this section we will show how variation autoencoder trained on the MNIST data set, can be used to generate images of handwritten digits. We demonstrate this by following the same set up as the previous section. Suppose, that the prior over the latent variable that are this given by an isotropic multivariate Gaussian  $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$  and our true posterior  $p_{\theta}(\mathbf{z}|\mathbf{x})$  is a multivariate Gaussian. Therefore, we will approximate the true posterior  $p_{\theta}(\mathbf{z}|\mathbf{x})$ , by a multivariate parametric Gaussian model that has a diagonal covariance,  $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I})$ . Where  $\boldsymbol{\mu}$  and  $\sigma^2$  are functions of  $\phi$ . Therefore by the previous section we have,

$$\tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \boldsymbol{\theta}, \phi) \simeq \frac{1}{2} \sum_{j=1}^m (1 + \log((\sigma_j^{(i)})^2) - (\boldsymbol{\mu}_j^{(i)})^2 - (\sigma_j^{(i)})^2) + \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}^{(i)} | \mathbf{z}^{(i,l)})$$

$$\text{where } \mathbf{z}^{(i,l)} = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\varepsilon}^{(l)}, \quad \boldsymbol{\varepsilon}^{(l)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

To produce the  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$  that will be fed into  $q_{\phi}(\mathbf{z}|\mathbf{x})$ , we define EncoderNeuralNet $_{\phi}$  that encodes our data  $\mathbf{x}^i$  into a  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$ . EncoderNeuralNet $_{\phi}$  is composed of a convolutional neural network (CNN) and dense layers, which inputs an image from the MNIST data sets (which are a  $28 \times 28$  matrix) and produces the encoded  $\boldsymbol{\mu}$  and  $\boldsymbol{\sigma}$ .

We will also make the additional assumption that  $p_{\theta}(\mathbf{x}|\mathbf{z})$  is a multivariate parametric Gaussian of the form,  $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\hat{\mathbf{x}}^{(i)}, \sigma_2 \mathbf{I})$ , with  $\hat{\mathbf{x}}^{(i)} = \text{DecoderNeuralNet}_{\theta}(\mathbf{z}^{(i,l)})$ . Where DecoderNeuralNet $_{\theta}$  is composed of a convolutional neural network (CNN) and dense layers, that takes in our sampled point  $\mathbf{z}^{(i,l)}$  and produces our mean  $\hat{\mathbf{x}}^{(i)}$ .

parameterised by  $\theta$  and  $\sigma_2$  is variant in our data is. formally, we would learn the parameter  $\sigma_2$ , however, we will simplify our problem by setting it to 1 (i.e  $\sigma_2 = 1$ ). Therefore,  $\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)})$  is given by,

$$\log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i,l)}) = -\frac{1}{2}\|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|^2 + C$$

Where  $C$  is a constant. Also since we are planning on using minibatch training when in the training procedure, we can simplify our expression for  $\tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \theta, \phi)$  as we can set  $L = 1$ . Therefore  $\tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \theta, \phi)$  takes the form of,

$$\begin{aligned} \tilde{\mathcal{L}}(\mathbf{x}^{(i)}; \theta, \phi) &\simeq \frac{1}{2} \sum_{j=1}^m (1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2) + \log p_{\theta}(\mathbf{x}^{(i)}|\mathbf{z}^{(i)}) \\ &= \frac{1}{2} \sum_{j=1}^m (1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2) - \frac{1}{2}\|\mathbf{x}^{(i)} - \hat{\mathbf{x}}^{(i)}\|^2 + C \end{aligned}$$

$$\text{where } \mathbf{z}^{(i)} = \boldsymbol{\mu}^{(i)} + \boldsymbol{\sigma}^{(i)} \odot \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

After training our model (described in Figure. 15) with 70,000 samples, batch sizes of 128, and 30 epochs. We get the following plots that describes the behaviour in latent space,

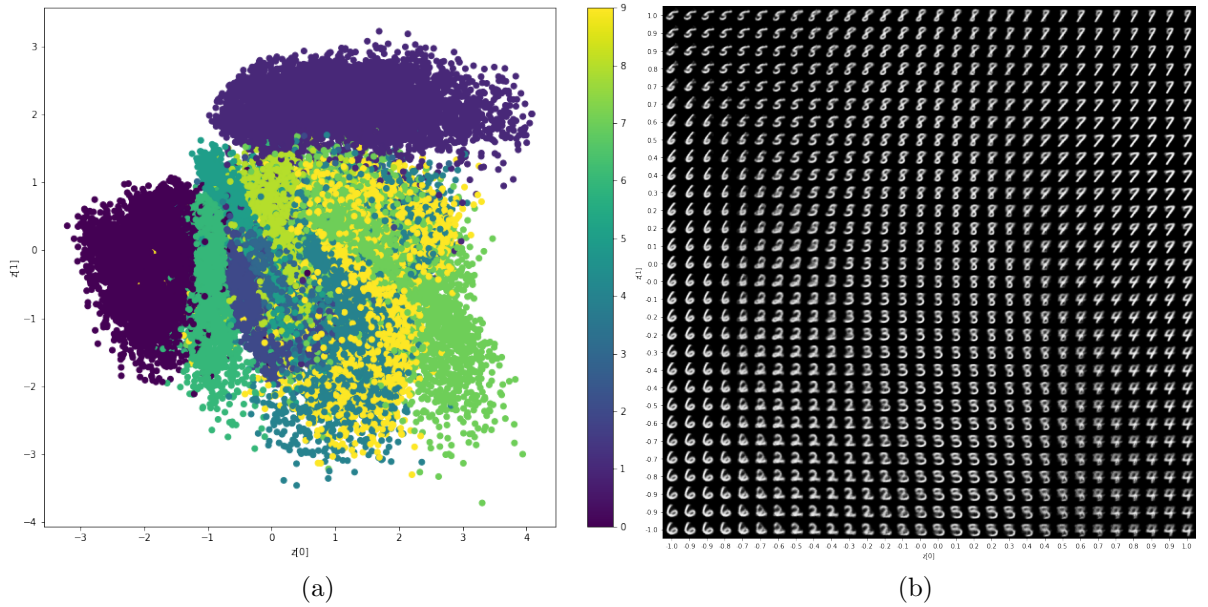


Figure 14: (a) Spatial distribution of latent space for MNIST digits (b) Learnt MNIST manifold in data-space.

We can see in both the figures that the model does a good job at regularising the distribution in latent space, since we chose a Gaussian centred at the origin. This allows for the KL divergence term in the ELBO to be smaller if the encoder is well-behaved, (i.e. is a relatively smooth distribution which is small far away from the origin). Thus, the KL term regularises the VAE. We can see this effect taking hold as most of the MNIST digits get regularised close to the origin.

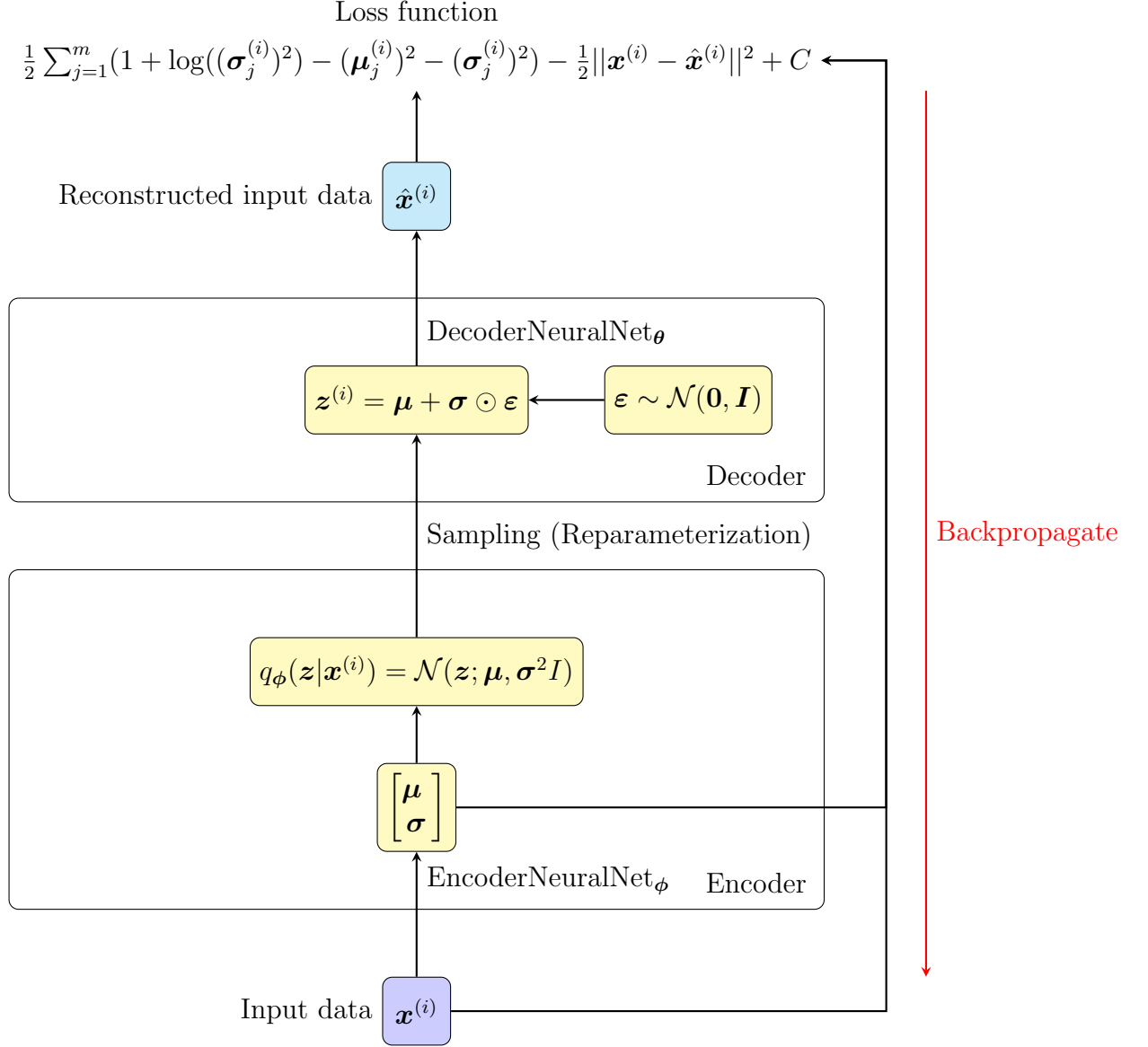


Figure 15: Trainable variational autoencoder

## 5 Neural Ordinary Differential Equations

We will now introduce the final piece of the background knowledge required before we tie it all together in the next chapter. As the title of the chapter may suggest, we will be looking at the combination of neural networks and ordinary differential equations (ODEs). The reason for studying this architecture is because in the next chapter we are going to be looking at modelling time series data with the use of VAE. However, we have only looked at VAEs without a time dependency and therefore the Neural ODE architecture will be used to help us close that gap.



## 5.1 Motivation

Many models in machine learning can be described as discrete transformations of hidden states, and can be expressed by the following relationship,

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \phi(\mathbf{z}_t, \boldsymbol{\theta}_t), \quad t \in \{0, \dots, T\}, \quad \mathbf{z}_t \in \mathbb{R}^m, \quad \phi(\cdot, \boldsymbol{\theta}_t) : \mathbb{R}^m \rightarrow \mathbb{R}^m \quad (22)$$

Where  $t$  is the depth of the network and  $\boldsymbol{\theta}_t$  parameterize the network at depth  $t$ . This relationship closely resembles the forward Euler method for solving ordinary differential equations for discrete time steps,  $t \in \{0, \dots, T\}$ . Therefore, analogously we can think of the network at depth  $t$ , representing the network at time  $t$ .

Hence, with the idea that (22) is the solution of the forward Euler method for solving an ODE at discrete layers  $\{0, \dots, T\}$ , we can hypothesise that  $\mathbf{z}(t)$  is continuous and governed by,

$$\frac{d\mathbf{z}(t)}{dt} = \mathbf{f}(\mathbf{z}(t), t, \boldsymbol{\theta}), \quad \mathbf{z}(0) = \mathbf{z}_0 \quad (23)$$

Where  $\mathbf{f} : \mathbb{R}^m \rightarrow \mathbb{R}^m$  is a neural network, parameterized by  $\boldsymbol{\theta}$ . This can always be achieved as we can ensure by our choice of  $\boldsymbol{\theta}$  that (22) holds (i.e  $\mathbf{f}(\mathbf{z}(t), t, \boldsymbol{\theta}) = \phi(\mathbf{z}_t, \boldsymbol{\theta}_t), \forall t \in \{0, \dots, T\}$ ).

Thus, we can think of the layers of our network in the continuous sense. Therefore, given an initial state of the network  $\mathbf{z}_0$ , we can solve (23) forwards to obtain the network's value at an arbitrary layer. This solution is unique and can always be achieved by treating the ODE solver as a black-box.

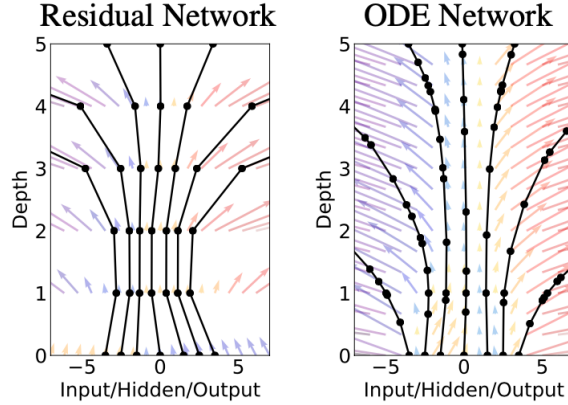


Figure 16: Left: A Residual network defines a discrete sequence of finite transformations. Right: An ODE network defines a vector field, of continuously state transforms. [1]

### 5.1.1 Toy Circular Dynamics Example

Since we are able to make the analogous link between layer depth of our network and time. We will use this idea to try and learn a continuous function of time. To achieve this, we will frame our problem using (23) and find the value of  $\boldsymbol{\theta}$  that

minimises the difference between  $\mathbf{z}(t)$  and the true solution.

We will make life easier for ourselves, we will try to learn a continuous function of time that we know its form as an initial value problem,

$$\dot{\mathbf{x}}(t) = \frac{d}{dt} \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \begin{bmatrix} -y(t) \\ x(t) \end{bmatrix} =: \mathbf{f}(\mathbf{x}(t)), \quad \mathbf{x}(0) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad t \in [0, 2\pi] \quad (24)$$

We notice the following relationship,

$$\ddot{x}(t) = -\dot{y} = -x(t), \quad x(0) = 1$$

Therefore, by the symmetry of our problem, we obtain,

$$\ddot{y}(t) = -y(t), \quad y(0) = 0$$

Hence, the unique solution to our initial value problem is given by,

$$\mathbf{x}(t) = \begin{bmatrix} \cos(t) \\ \sin(t) \end{bmatrix}, \quad t \in [0, 2\pi] \quad (25)$$

To learn the correct  $\mathbf{z}(t)$ , we will use the linear function,  $\mathbf{f}_\theta : \mathbb{R}^2 \rightarrow \mathbb{R}^2$  as our neural network.

$$\dot{\mathbf{z}}(t) = \frac{d}{dt} \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} = \mathbf{f}_\theta(\mathbf{z}(t)) := \begin{bmatrix} \theta_1 & \theta_2 \\ \theta_3 & \theta_4 \end{bmatrix} \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} + \begin{bmatrix} \theta_5 \\ \theta_6 \end{bmatrix}, \quad \mathbf{z}(0) = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad t \in [0, 2\pi] \quad (26)$$

Since we have the explicit form of the continuous function that we want to model (given by (25)), we will artificially simulate data points for our problem. This is done by randomly sampling a stopping time  $t_{\text{stop}}$  from  $\mathcal{U}(0, 2\pi)$  (i.e  $t_{\text{stop}} \sim \mathcal{U}(0, 2\pi)$ ) and using (25) to calculate the true value of the solution at that time. After doing this procedure 12,800 times, we have a large enough data set that we can use for training. To form a prediction from our network, we will use the fourth-order Runge–Kutta (RK4) method to integrate (26) forward in time to  $t_{\text{stop}}$ . Then we will use the Euclidean norm between our prediction and the true solution to quantify how well our prediction is and finally, we will backpropagate the error through the RK4 integration, as it has discretised the ODE.

$$\begin{aligned} \text{L2 loss function } L &= \|\mathbf{x}_{\text{stop}} - \mathbf{z}_{\text{stop}}\|^2 \\ &= \|\mathbf{x}_{\text{stop}} - \text{RK4}((1, 0)^T, \mathbf{f}_\theta, 0, t_{\text{stop}})\|^2 \\ \implies \frac{\partial L}{\partial \theta} &= \frac{\partial \|\mathbf{x}_{\text{stop}} - \text{RK4}((1, 0)^T, \mathbf{f}_\theta, 0, t_{\text{stop}})\|^2}{\partial \theta} \end{aligned}$$

We demonstrate this whole training process by the following diagram,

After training using the process described in Fig (17), with 12,800 sample points, batch sizes of 64, and 4 epochs, we get the following forcing function,

$$\mathbf{f}_\theta(\mathbf{z}(t)) = \begin{bmatrix} -0.00 & 1.01 \\ -1.01 & 0.00 \end{bmatrix} \begin{bmatrix} x(t) \\ y(t) \end{bmatrix} + \begin{bmatrix} -0.01 \\ -0.00 \end{bmatrix} \approx \begin{bmatrix} y(t) \\ -x(t) \end{bmatrix}$$

Due to the symmetrical nature of our ODE example we still get same set of equations to solve,

$$\ddot{x}(t) = \dot{y}(t) = -x(t), \quad x(0) = 1, \quad \text{and} \quad \ddot{y}(t) = -y(t), \quad y(0) = 0$$

As this example was done this a linear function, this reinforces that our theory should hold for neural network as they are a compositions of linear functions.

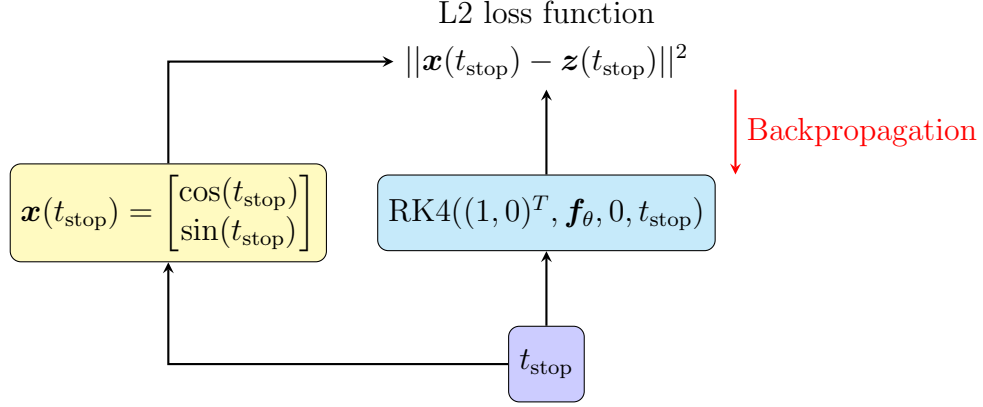


Figure 17: Training model for learning (25) using (26) and the RK4 method.

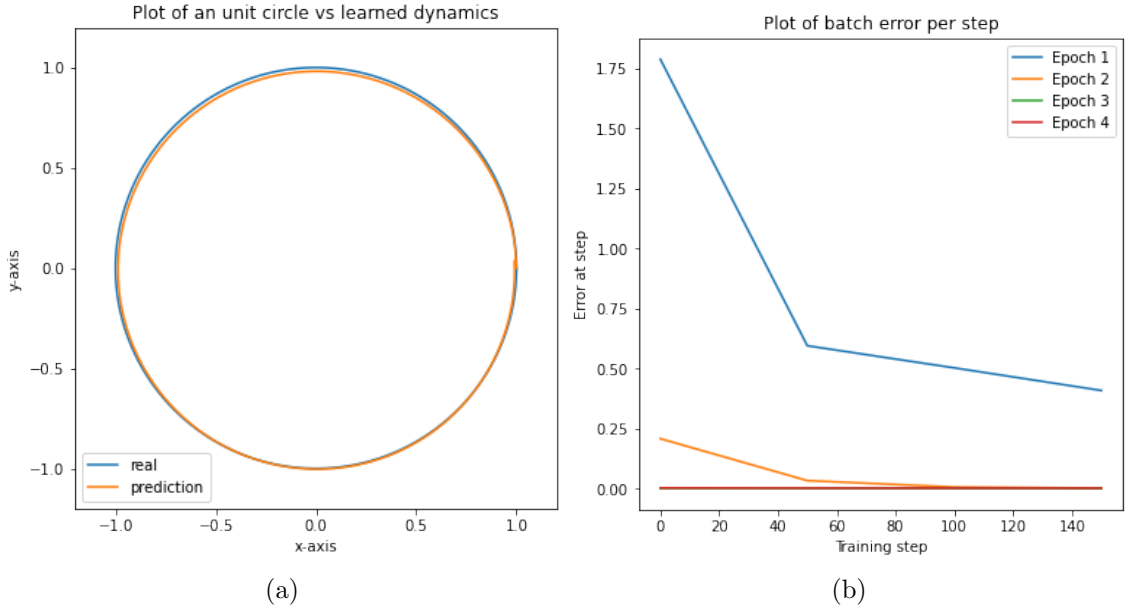


Figure 18: (a) Plot of the networks learnt trajectory  $\mathbf{z}(t)$  (orange) vs the true solution's trajectory  $\mathbf{x}(t)$  (blue) to the problem (24). (b) Plot of the every 50th batch error over four epochs (converge to an error of  $0.5 \times 10^{-3}$ ). [3]

## 5.2 Formalisation and Improvements

Like we did in our toy example (5.1.1), the native approach for training (23) is to apply backpropagation through the black-box solver. This works due to the block-box solver discretising the ODE, allowing us to take the derivative with respect to the parameters. However, this leads to high memory cost and introduces additional numerical error. This is because the black-box solver can be thought of as discrete transformation and compositions of the forcing function. As this procedure will be done a lot, this causes the derivative with respect to the parameters to be expensive, as they appear in the compositions. Hence, we seek to find a better solution to overcome this problem. We achieve this by deriving an expression for the derivative of the loss function with respect to the model parameters, that involves solving an ODE backward through the layers. This is more appealing approach as modern ODE solvers allow for the problem to scale linearly with size. We can also to control

the tolerance of our solution, allowing for the ability to dictate how accurate our training is; Thus, helping us prevent over-fitting.

To make understanding the theory more digestible, we will impose the following simplification to our problem. We will only consider the network at being evaluated at two discrete points,  $t_{\text{start}}$  and  $t_{\text{stop}}$ , and we will look at minimising the loss produced over that interval. This simplification does not cause any initiations to our theory, as we can extend this theory to multiple observations (i.e minimising the loss function over  $t_0, \dots, t_N$  observations ) and is shown in the Appendix 2.2.

We will now define the scalar loss function  $L : \mathbb{R}^m \rightarrow \mathbb{R}$ , that we alluded to above, where the input is the output of our black-box solver.

$$L(\mathbf{z}(t_1)) = L \left( \mathbf{z}(t_0) + \int_{t_0}^{t_1} \mathbf{f}(\mathbf{z}(t), t, \theta) dt \right) = L(\text{ODESolve}(\mathbf{z}(t_0), \mathbf{f}, t_0, t_1, \theta))$$

We aim to find the value of  $\theta$  that minimises  $L$ , as this will provide us with a model that captures the dynamics of our data closely. To obtain this we firstly define the adjoint, which is gradient of  $L$  with respect to the hidden state  $\mathbf{z}(t)$ ,

$$\mathbf{a}(t) = \frac{\partial L}{\partial \mathbf{z}(t)} \quad (27)$$

It is important to remember that  $t$  is the layer depth and it is now continuous, so this holds for all  $t \geq 0$ .

We show in Section 2.1, that the derivative of the adjoint with respect to the layer is given by,

$$\frac{d\mathbf{a}^T(t)}{dt} = -\mathbf{a}^T \frac{\partial \mathbf{f}(\mathbf{z}(t), \theta, t)}{\partial \mathbf{z}} \quad (28)$$

Since this does not tell use how  $L$  changes with respect to the parameters  $\theta$  and  $t$ . Therefore we define the augmented state  $\mathbf{z}_{\text{aug}}(t)$ , to allow us to do so.

$$\mathbf{z}_{\text{aug}}(t) = \begin{bmatrix} \mathbf{z}(t) \\ \theta \\ t \end{bmatrix}$$

We now derive the new dynamical system,

$$\frac{d\mathbf{z}_{\text{aug}}(t)}{dt} = \mathbf{f}_{\text{aug}}([\mathbf{z}, \theta, t]) := \begin{bmatrix} \mathbf{f}([\mathbf{z}, \theta, t]) \\ \mathbf{0} \\ 1 \end{bmatrix}$$

$$\mathbf{a}_{\text{aug}}(t) := \begin{bmatrix} \mathbf{a} \\ \mathbf{a}_\theta \\ \mathbf{a}_t \end{bmatrix} (t), \quad \mathbf{a}_\theta(t) := \frac{dL}{d\theta(t)}, \quad \mathbf{a}_t(t) := \frac{dL}{dt(t)}$$

Therefore, we are able to use (28) on our new dynamics to obtain,

$$\frac{d\mathbf{a}_{\text{aug}}^T(t)}{dt} = -\mathbf{a}_{\text{aug}}^T \frac{\partial \mathbf{f}_{\text{aug}}}{\partial [\mathbf{z}, \theta, t]} \quad (29)$$

Thus, by the calculating (29) (which is shown in Appendix 2.2), we obtain the derivative of  $\mathbf{a}_\theta(t)$  with respect to the layer,

$$\frac{d\mathbf{a}_\theta^T(t)}{dt} = \frac{d}{dt} \left( \frac{dL^T}{d\theta} \right) = -\mathbf{a}^T \frac{\partial \mathbf{f}(\mathbf{z}(t), \theta, t)}{\partial \theta}$$

Hence, by integrating this expression with backwards over the layers and assuming that  $\mathbf{a}_\theta(t_1) = \mathbf{0}$  (explained in the Appendix 2.2), we achieved what we set out to obtain,

$$\frac{dL^T}{d\theta} = - \int_{t_1}^{t_0} \mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{z}(t), \theta, t)}{\partial \theta} dt \quad (30)$$

To understand why we integrate backwards over the layers, we draw parallels from the discrete case. As the loss function depends of the output of the network (or in our case at  $t_1$ ) we backpropagate the error through the network updating the models as we go. This is the same idea in the continuous case by integrating from  $t_1$  to  $t_0$  (i.e from the output to the input of our network) allowing us to compute the overall sensitivity of the model parameters.

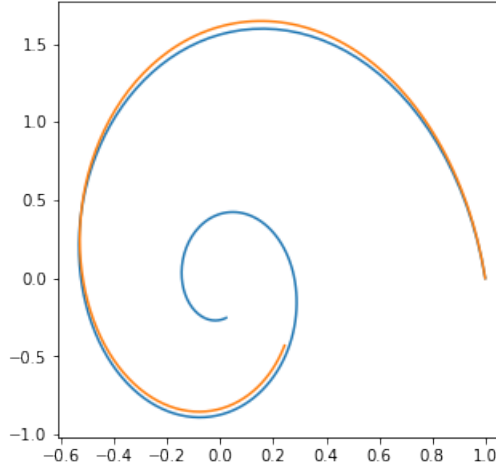


Figure 19: Plot of the networks learnt trajectory (orange) vs the true solution's trajectory (blue) using the adjoint sensitivity method to compute the gradients for a spiral dynamics. [3]

## 6 Generative Latent Time-series Models

### 6.1 Overview

We will now tie together the theory established throughout this paper by developing a generative approach for modelling continuous time series-data. This approach is a continuation of variational autoencoders that learns the time-trajectory of our sampled feature point  $z_{t_0} \sim p_\theta(\mathbf{z})$  in feature space.

We pose the question, we want to solve. Suppose that we have the data,  $\mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N}$  which are  $N + 1$  data observation at times  $t_0, \dots, t_N$ . Therefore

We do this by making the following adaptations to the variational autoencoder model:

1. The input data  $\mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N}$  is encoded as a distribution  $p_{\theta}(\cdot | \mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N})$  in feature space
2. We then sample the point  $\mathbf{z}_{t_0}$  from the distribution in feature space
3. Then we integrate the sampled point to the observation times  $t_0, \dots, t_N$ ,  $\mathbf{z}_{t_0}, \dots, \mathbf{z}_{t_N} = \text{ODESolve}(\mathbf{z}_{t_0}, \mathbf{f}, t_0, \dots, t_N, \theta)$ . Where  $\mathbf{f}$  is the forcing function of that describes the feature point through time,

$$\frac{d\mathbf{z}(t)}{dt} = \mathbf{f}(\mathbf{z}(t)), \quad \mathbf{z}(t_0) = \mathbf{z}_{t_0} \quad (31)$$

4. Finally, the points  $\mathbf{z}_{t_0}, \dots, \mathbf{z}_{t_N}$  are then decoded into data-space by  $\hat{\mathbf{x}}_{t_i} = p_{\theta}(\mathbf{x} | \mathbf{z}_{t_i})$  to get the reconstructed input data  $\hat{\mathbf{x}}_{t_0}, \dots, \hat{\mathbf{x}}_{t_N}$

Because the forcing function  $\mathbf{f}$  is unknown, we will use a neural ordinary differential equation to learn the correct trajectory in feature space. To ensure that our trajectory in feature space is uniquely defined, we make the assumption that the neural network,  $\mathbf{f}_{\theta} : \mathbb{R}^m \rightarrow \mathbb{R}^m$  parameterized by  $\theta$  is a time-invariant function.

We can express this whole process probabilistically by the following the sampling procedure,

$$\mathbf{z}_{t_0} \sim p_{\theta}(\mathbf{z} | \mathbf{x}) \quad (32)$$

$$\mathbf{z}_{t_0}, \dots, \mathbf{z}_{t_N} = \text{ODESolve}(\mathbf{z}_{t_0}, \mathbf{f}, t_0, \dots, t_N, \theta) \quad (33)$$

$$\text{draw each } \mathbf{x}_{t_i} \sim p_{\theta}(\cdot | \mathbf{z}_{t_i}) \quad (34)$$

## 6.2 Training

Since this model is a modification of variational autoencoders, we will build off its training model (Fig 15) to accommodate these changes. We start by modifying our  $\text{EncoderNeuralNet}_{\phi}$  to allow for multiple observations to be used in the encoding of  $\mu$  and  $\sigma$ . This is done by feeding our observation through a composition of a recurrent neural network (RNN) and a dense layer,

$$\text{RNN}(\mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N}) : \mathbb{R}^n \times \dots \times \mathbb{R}^n \rightarrow \mathbb{R}^n, \text{DenseLayer}(\cdot) = \begin{bmatrix} \mu \\ \sigma \end{bmatrix} : \mathbb{R}^n \rightarrow \mathbb{R}^{2m}$$

$$\text{EncoderNeuralNet}_{\phi}(\cdot) := \text{DenseLayer}(\text{RNN}(\cdot)) : \mathbb{R}^n \times \dots \times \mathbb{R}^n \rightarrow \mathbb{R}^{2m}$$

We choose our encoder in this way to ensure that the temporal nature of our data gets encapsulates into the encoding of our  $\mu$  and  $\sigma$ .

Like we did in our example 4.4, we assume that the our prior over the latent variables is given by an isotropic multivariate Gaussian  $p_{\theta}(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$  and our true posterior  $p_{\theta}(\mathbf{z} | \mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N})$  takes the form of a multivariate Gaussian; Therefore, we will approximate  $p_{\theta}(\mathbf{z} | \mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N})$  with a multivariate Gaussian with diagonal covariance,  $q_{\phi}(\mathbf{z} | \mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N}) = \mathcal{N}(\mathbf{z}; \mu, \sigma^2 \mathbf{I})$ .

Hence, our sampling procedure from our encoder now takes the form,

$$\mathbf{z}_{t_0} \sim q_{\phi}(\mathbf{z} | \mathbf{x}_{t_0}, \dots, \mathbf{x}_{t_N}) = \mathcal{N}(\mathbf{z}; \mu, \sigma^2 \mathbf{I})$$

We will use the same reparameterization trick described in section 4.3, to get,

$$\mathbf{z}_{t_0} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Now we will integrated  $\mathbf{f}_{\boldsymbol{\theta}}$  forwards in feature-space to our observations times,

$$\mathbf{z}_{t_0}, \dots, \mathbf{z}_{t_N} = \text{ODESolve}(\mathbf{z}_{t_0}, \mathbf{f}_{\boldsymbol{\theta}}, t_0, \dots, t_N, \boldsymbol{\theta})$$

We then decode these points into data-space by the following mapping,

$$\text{DecoderNeuralNet}_{\boldsymbol{\theta}}(\mathbf{x}|\cdot) : \mathbb{R}^m \rightarrow \mathbb{R}^n$$

$$\hat{\mathbf{x}}_{t_i} = \text{DecoderNeuralNet}_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}_{t_i}), \quad i \in \{0, \dots, N\}$$

This whole processes is described by the following figure,

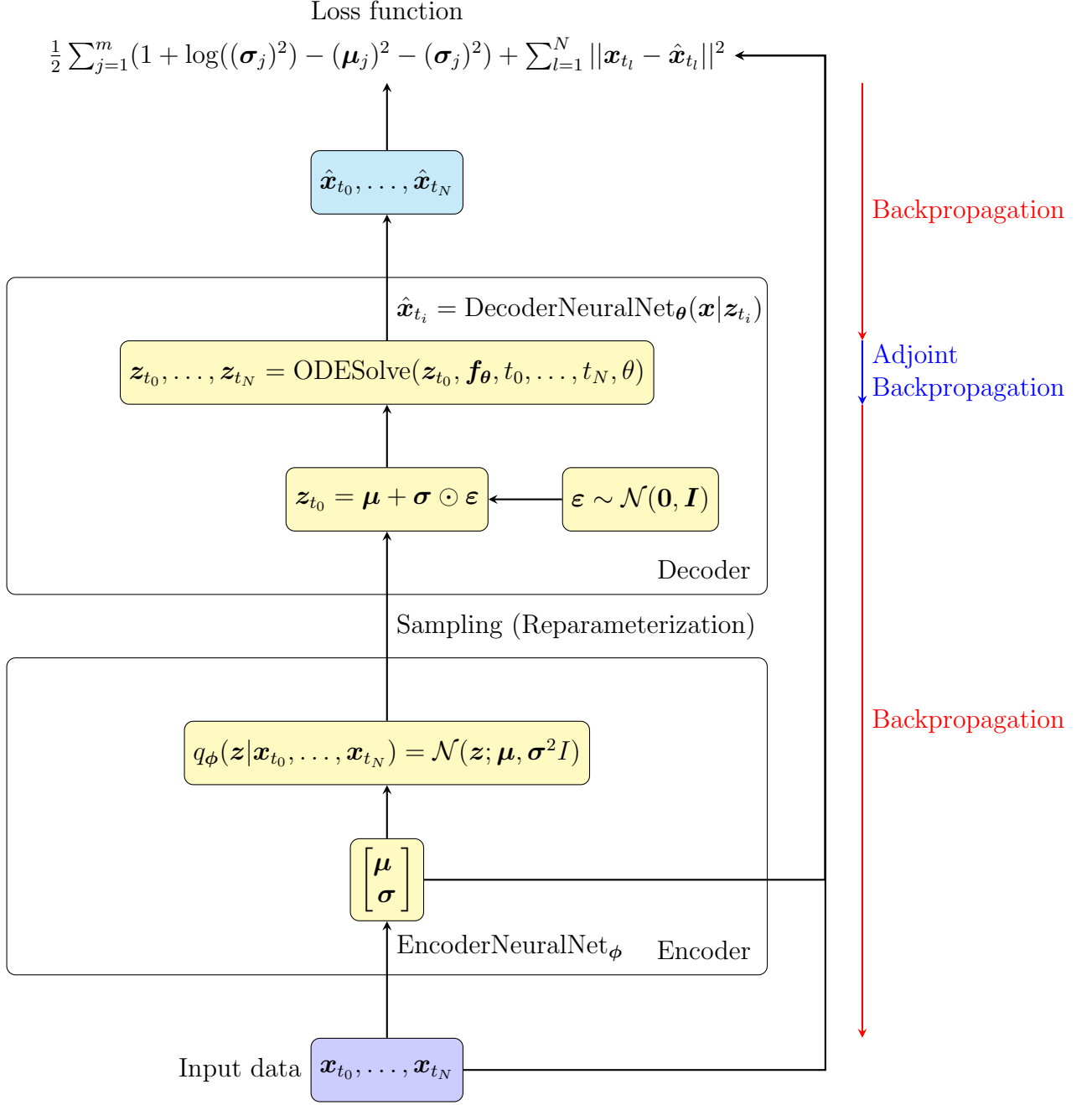


Figure 20: Trainable generative latent time-series model

It may not be apparently clear what form the adjoint takes at our observed



values which is used in alg ??.

$$\begin{aligned}
\mathbf{a}(t_i) &= \frac{\partial L}{\partial \mathbf{z}(t)} \Big|_{t=t_i} \\
&= \frac{\partial}{\partial \mathbf{z}(t)} \left( \frac{1}{2} \sum_{j=1}^m (1 + \log((\boldsymbol{\sigma}_j)^2) - (\boldsymbol{\mu}_j)^2 - (\boldsymbol{\sigma}_j)^2) - \frac{1}{2} \sum_{l=1}^N \|\mathbf{x}_{t_l} - \hat{\mathbf{x}}_{t_l}\|^2 \right) \Big|_{t=t_i} \\
&= \frac{\partial}{\partial \mathbf{z}(t)} \left( -\frac{1}{2} \sum_{l=1}^N \|\mathbf{x}_{t_l} - \text{DecoderNeuralNet}_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}_{t_i})\|^2 \right) \Big|_{t=t_i} \\
&= \|\mathbf{x}_{t_i} - \text{DecoderNeuralNet}_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}_{t_i})\| \frac{\partial \text{DecoderNeuralNet}_{\boldsymbol{\theta}}(\mathbf{x}|\mathbf{z}_{t_i})}{\partial \mathbf{z}(t)} \Big|_{t=t_i} \quad (\text{Chain Rule})
\end{aligned}$$

This applications has made the implicit assumption that our observation in data space is uniformly sampled; However, this is not normally the case. To overcome this to make the following adaptation to the above model,

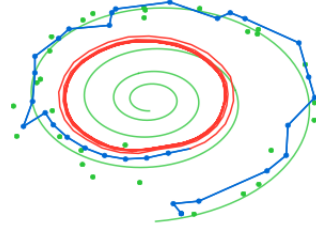
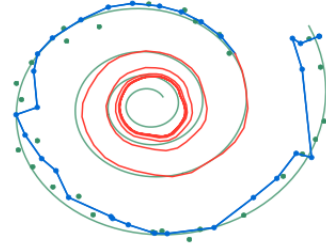
$$\text{RNN}(\{\mathbf{x}_{t_0}, t_0\}, \dots, \{\mathbf{x}_{t_N}, t_N\}) : \mathbb{R}^{(n+1)} \times \dots \times \mathbb{R}^{(n+1)} \rightarrow \mathbb{R}^{(n+1)}$$

$$\text{DenseLayer}(\cdot) = \begin{bmatrix} \boldsymbol{\mu} \\ \boldsymbol{\sigma} \end{bmatrix} : \mathbb{R}^{(n+1)} \rightarrow \mathbb{R}^{2m}$$

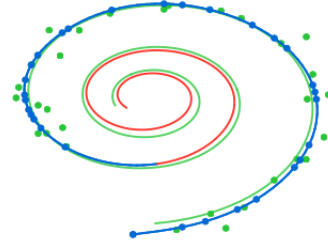
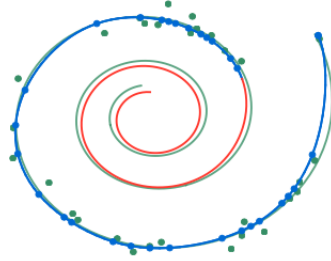
$$\text{EncoderNeuralNet}_{\phi}(\cdot) := \text{DenseLayer}(\text{RNN}(\cdot)) : \mathbb{R}^{(n+1)} \times \dots \times \mathbb{R}^{(n+1)} \rightarrow \mathbb{R}^{2m}$$

This concatenation of the observation and its time, allows for the learning of the distribution of the sampled observations times.

### 6.3 Bi-directional Spiral Example

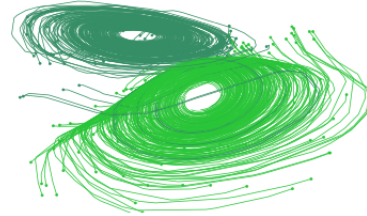


(a) Recurrent Neural Network



(b) Latent Neural Ordinary Differential Equation

— Ground Truth  
 • Observation  
 — Prediction  
 — Extrapolation



(c) Latent Trajectories

Figure 21: (a): Reconstruction and extrapolation of spirals with irregular time points by a recurrent neural network. (b): Reconstructions and extrapolations by a latent neural ODE. Blue curve shows model prediction. Red shows extrapolation. (c) A projection of inferred 4-dimensional latent ODE trajectories onto their first two dimensions. Color indicates the direction of the corresponding trajectory. The model has learned latent dynamics which distinguishes the two directions.[1]

# Appendix A

## Little Results and Derivations

### 1 Gibbs' Inequality [2]

Let  $p = \{p_1, \dots, p_k\}$  and  $q = \{q_1, \dots, q_k\}$  be two arbitrary discrete probability distributions, then Gibbs' inequality states that,

$$-\sum_{i=1}^k p_i \log p_i \leq -\sum_{i=1}^k p_i \log q_i$$

and equality if and only if,

$$p_i = q_i, \quad i \in \{1, \dots, k\}$$

### 2 Adjoint and Augmented Dynamics Derivations

We will now expand upon the derivations given in Neural Ordinary Differential Equations paper [1] to include details that were omitted from the original paper.

#### 2.1 Adjoint Derivation

In this section we will prove the statement claimed in Section 5.2, that the layer derivative of  $\mathbf{a}(t)$  is given by,

$$\frac{d\mathbf{a}^T(t)}{dt} = -\mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)}. \quad (\text{A.1})$$

It is important to remember that  $t$  is the layer depth and is now continuous.

We start by drawing parallels from a general (discrete) neural network. We use the fact that the gradient of the loss function with respect to the hidden layer  $\mathbf{z}_t$ , can be expressed in terms of the next layer  $\mathbf{z}_{t+1}$ ,

$$\frac{dL}{d\mathbf{z}_t} = \frac{dL}{d\mathbf{z}_{t+1}} \frac{d\mathbf{z}_{t+1}}{d\mathbf{z}_t}$$

We now generalise this approach to a neural network with continuous hidden states, by defining  $\mathbf{z}(t + \varepsilon)$ , as the following,

$$\mathbf{z}(t + \varepsilon) = \int_t^{t+\varepsilon} \mathbf{f}(\mathbf{z}(t), t, \theta) dt + \mathbf{z}(t) =: T_\varepsilon(\mathbf{z}(t), t), \quad \varepsilon > 0.$$

Hence, by applying the chain rule to  $\mathbf{z}(t + \varepsilon)$ , we are able to express the gradient of the loss function with respect to the continuous hidden state  $\mathbf{z}(t)$ , in terms of gradients that include,  $\mathbf{z}(t + \varepsilon)$ ,

$$\frac{dL^T}{d\mathbf{z}(t)} = \frac{dL^T}{d\mathbf{z}(t + \varepsilon)} \frac{d\mathbf{z}(t + \varepsilon)}{d\mathbf{z}(t)} \quad \text{or} \quad \mathbf{a}^T(t) = \mathbf{a}^T(t + \varepsilon) \frac{\partial T_\varepsilon(\mathbf{z}(t), t)}{\partial \mathbf{z}(t)} \quad (\text{A.2})$$

Thus, we have the required machinery to prove (A.1),

$$\begin{aligned} \frac{d\mathbf{a}^T(t)}{dt} &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}^T(t + \varepsilon) - \mathbf{a}^T(t)^T}{\varepsilon} \quad (\text{Derivative definition}) \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}^T(t + \varepsilon) - \mathbf{a}^T(t + \varepsilon) \frac{\partial}{\partial \mathbf{z}(t)} T_\varepsilon(\mathbf{z}(t))}{\varepsilon} \quad (\text{using (A.2)}) \\ \text{Now by expanding } T_\varepsilon(\mathbf{z}(t)) \text{ around } \mathbf{z}(t), \text{ we get,} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}^T(t + \varepsilon) - \mathbf{a}^T(t + \varepsilon) \frac{\partial}{\partial \mathbf{z}(t)} (\mathbf{z}(t) + \varepsilon \mathbf{f}(\mathbf{z}(t), t, \theta) + \mathcal{O}(\varepsilon^2))}{\varepsilon} \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{\mathbf{a}^T(t + \varepsilon) - \mathbf{a}^T(t + \varepsilon) \left( I + \varepsilon \frac{\partial \mathbf{f}(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} + \mathcal{O}(\varepsilon^2) \right)}{\varepsilon} \quad (\text{Applying the derivative}) \\ &= \lim_{\varepsilon \rightarrow 0^+} \frac{-\varepsilon \mathbf{a}^T(t + \varepsilon) \frac{\partial \mathbf{f}(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} + \mathcal{O}(\varepsilon^2)}{\varepsilon} \quad (\text{Expanding bracket}) \\ &= \lim_{\varepsilon \rightarrow 0^+} -\mathbf{a}^T(t + \varepsilon) \frac{\partial \mathbf{f}(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} + \mathcal{O}(\varepsilon) \quad (\text{Cancelling terms}) \\ &= -\mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)} \quad (\text{Taking limit}) \end{aligned}$$

Hence, we have arrived at what we set out to prove, that the layer derivative of  $\mathbf{a}(t)$  is given by,

$$\frac{d\mathbf{a}^T(t)}{dt} = -\mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{z}(t), t, \theta)}{\partial \mathbf{z}(t)}$$

## 2.2 Augmented Dynamics

By assuming that the model parameters  $\boldsymbol{\theta}$  is layer invariant, we are able to generalise (28) to obtain the gradient of  $L$  with respect to  $\boldsymbol{\theta}$  and  $t$  respectively.

To achieve this we introduce the augmented state dynamics for our system, defined as the following,

$$\frac{d\mathbf{z}_{\text{aug}}(t)}{dt} = \frac{d}{dt} \begin{bmatrix} \mathbf{z} \\ \boldsymbol{\theta} \\ t \end{bmatrix} (t) = \mathbf{f}_{\text{aug}}([\mathbf{z}, \boldsymbol{\theta}, t]) := \begin{bmatrix} \mathbf{f}([\mathbf{z}, \boldsymbol{\theta}, t]) \\ \mathbf{0} \\ 1 \end{bmatrix} \quad (\text{A.3})$$

$$\mathbf{a}_{\text{aug}}(t) := \begin{bmatrix} \mathbf{a} \\ \mathbf{a}_\theta \\ \mathbf{a}_t \end{bmatrix} (t), \quad \mathbf{a}_\theta(t) := \frac{dL}{d\boldsymbol{\theta}(t)}, \quad \mathbf{a}_t(t) := \frac{dL}{dt(t)} \quad (\text{A.4})$$

$$\frac{\partial \mathbf{f}_{\text{aug}}}{\partial [\mathbf{z}, \theta, t]} = \begin{bmatrix} \frac{\partial \mathbf{f}}{\partial \mathbf{z}} & \frac{\partial \mathbf{f}}{\partial \theta} & \frac{\partial \mathbf{f}}{\partial t} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix} (t) \quad (\text{A.5})$$

Where each  $\mathbf{0}$  is a matrix of zeroes with the appropriate dimensions.

Therefore, by applying (A.1) to our augmented system, we can understand how the augmented adjoint changes with respect to the layer,

$$\begin{aligned} \frac{d\mathbf{a}_{\text{aug}}^T(t)}{dt} &= - [\mathbf{a}^T(t) \quad \mathbf{a}_{\theta}^T(t) \quad \mathbf{a}_t^T(t)] \frac{\partial \mathbf{f}_{\text{aug}}}{\partial [\mathbf{z}, \theta, t]}(t) \\ &= - [\mathbf{a}^T \frac{\partial \mathbf{f}}{\partial \mathbf{z}} \quad \mathbf{a}^T \frac{\partial \mathbf{f}}{\partial \theta} \quad \mathbf{a}^T \frac{\partial \mathbf{f}}{\partial t}] (t) \end{aligned} \quad (\text{A.6})$$

Thus, by looking at the second element of equation (A.6), we get the following relationship,

$$\frac{d\mathbf{a}_{\theta}^T(t)}{dt} = \frac{d}{dt} \left( \frac{dL^T}{d\theta} \right) = \mathbf{a}^T(t) \frac{\partial \mathbf{f}}{\partial \theta}$$

Therefore, by integrating backwards over the layers and setting the adjoint of  $\theta$  at time  $t_1$  to zero as we have not "accumulated" any error yet. We therefore obtain the overall sensitivity of the model parameters,

$$\implies \frac{dL^T}{d\theta} = \mathbf{a}_{\theta}^T(t_0) = - \int_{t_1}^{t_0} \mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{z}(t), \theta, t)}{\partial \theta} dt$$

Finally, we can work out expressions for the derivative of  $L$  with respect to the times  $t_0$  and  $t_1$ , respectively. To derive the expression at  $t_1$ , we apply the chain rule to get,

$$\frac{\partial L^T}{\partial t_N} = \frac{\partial L^T}{\partial \mathbf{z}(t)} \frac{d\mathbf{z}(t)}{dt} \Big|_{t=t_N} = \mathbf{a}^T(t_N) \mathbf{f}(\mathbf{z}(t_N), \theta, t) \quad (\text{A.7})$$

Hence, by using the final term in equation (A.6), and integrating backwards over the layers we get,

$$\implies \frac{\partial L^T}{\partial t_0} = \mathbf{a}_t^T(t_0) = \mathbf{a}_t^T(t_1) - \int_{t_1}^{t_0} \mathbf{a}^T(t) \frac{\partial \mathbf{f}(\mathbf{z}(t), \theta, t)}{\partial t} dt \quad (\text{A.8})$$

Up to now we have only considered our loss function  $L$  depends between  $t_0$  and  $t_1$ . However, if we wanted to evaluate the ODE at multiple time observations  $t_0, \dots, t_N$ , we would adapt the above process by iterating the process over the intervals,  $[t_{N-1}, t_N], [t_{N-2}, t_{N-1}], \dots$ , accumulating the corrections for each interval.

# Bibliography

- [1] Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in neural information processing systems*, pages 6571–6583, 2018.
- [2] Pierre Bremaud. *An Introduction to Probabilistic Modeling*. Undergraduate Texts in Mathematics. Springer New York : Imprint: Springer, New York, NY, 1st ed. 1988. edition, 1988.
- [3] Kirkpatrick Darrah. MMath Code. <https://github.com/DarrahK/MMath-Code>, 2021.
- [4] Dr. Jonathan Bartlett. MA20226: Statistics 2A Lecture Notes. *University of Bath*, January 2019.
- [5] Pankaj Mehta, Marin Bukov, Ching-Hao Wang, Alexandre GR Day, Clint Richardson, Charles K Fisher, and David J Schwab. A high-bias, low-variance introduction to machine learning for physicists. *Physics reports*, 810:1–124, 2019.
- [6] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature (London)*, 323(6088):533–536, 1986.
- [7] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*, 1(4):541–551, 12 1989.
- [8] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2014.
- [9] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *CoRR*, abs/1906.02691, 2019.