

# Latent Variable Models in Neural Machine Translation

*John Isak Texas Falk*

A dissertation submitted in partial fulfillment  
of the requirements for the degree of  
**Master of Science**  
of  
**University College London.**

The Centre for Computational Statistics and Machine Learning  
University College London

August 23, 2017

I, John Isak Texas Falk, confirm that the work presented in this thesis is my own. Where information has been derived from other sources, I confirm that this has been indicated in the work.

# Abstract

Neural Machine translation is a direction in automated translation which learns the mapping from the source to the target language directly in an end-to-end fashion using the framework of neural networks. As NMT hinges on advances in neural network methods and architectures, it is able to directly use improvements there in its own domain.

We use recent improvements in latent variables models in order to train a completely probabilistic generative model with a latent variable representing a language-agnostic representation of a sentence mapping through deep neural networks to two different output languages. Following recent advances in training deep generative latent variable models we approximate the posterior of the latents given output with a recognition model, mimicking a VAE. Following the SGVB method to find a stochastic lower bound to the true log-likelihood of the observed data, we train the parameters of the generative model and the variational recognition model jointly to optimise this bound.

Since the recognition model acts as a pseudo-posterior (it approximates it given constraints on the distributional form of the recognition model) we can use this to translate from one language to another by finding the posterior  $q$ -distribution over  $z$  and then from sampled  $z$  find the most likely output of the languages.

# Acknowledgements

I would like to thank the help from my supervisor Harshil Shah for helping me make this thesis possible, and my parents, for always being there for me.

# Contents

<b>Nomenclature</b>	<b>9</b>
<b>1 Introduction</b>	<b>13</b>
<b>2 Background Knowledge</b>	<b>15</b>
2.1 Probability Theory and Statistics . . . . .	15
2.1.1 Rules and Theorems . . . . .	15
2.1.2 The Gaussian Distribution . . . . .	17
2.1.3 Maximum Likelihood Estimation . . . . .	18
2.1.4 Graphical models . . . . .	19
2.1.5 Approximate Inference . . . . .	21
2.2 Deep Learning . . . . .	22
2.2.1 Multilayer Perceptron . . . . .	23
2.2.2 Recurrent Neural Networks . . . . .	23
2.2.3 Convolutional Neural Networks . . . . .	24
2.3 Natural Language Processing . . . . .	25
2.3.1 Language model . . . . .	25
2.3.2 Word embeddings . . . . .	26
2.3.3 Neural machine translation . . . . .	27
2.4 Optimization . . . . .	28
2.4.1 Problem formulation . . . . .	28
2.4.2 Stochastic Gradient Descent . . . . .	28
2.4.3 ADAM . . . . .	29

<b>3</b>	<b>Methods and Theory</b>	<b>31</b>
3.1	ELBO . . . . .	31
3.2	Models . . . . .	32
3.2.1	Reconstruction . . . . .	33
3.2.2	Translation . . . . .	33
3.2.3	Training . . . . .	34
<b>4</b>	<b>Experiments</b>	<b>35</b>
4.1	Data . . . . .	35
4.1.1	Dataset . . . . .	35
4.1.2	Preprocessing . . . . .	35
4.2	Scores . . . . .	37
4.3	Reconstruction . . . . .	38
4.4	Translation . . . . .	38
<b>5</b>	<b>Conclusions</b>	<b>39</b>
	<b>Appendices</b>	<b>40</b>
<b>A</b>	<b>Proofs</b>	<b>40</b>
<b>B</b>	<b>Another Appendix About Things</b>	<b>41</b>
<b>C</b>	<b>Colophon</b>	<b>42</b>
	<b>Bibliography</b>	<b>43</b>

# List of Figures

2.1	RNN schematics . . . . .	24
2.2	Encoder decoder schematic . . . . .	27
4.1	Runs with various different recognition model . . . . .	38

# List of Tables

4.1	A randomly sampled sentence from the Europarl corpus . . . . .	36
-----	--	----



# Nomenclature

GD    Gradient Descent

MCMC   Markov Chain Monte Carlo

MLE   Maximum Likelihood Estimation

NLP   Natural Language Processing

NMT   Neural Machine Translation

SGD   Stochastic Gradient Descent

VI    Variational Inference

$\theta \in \Theta$    Parameters  $\theta$  belonging to the parameter space  $\Theta$

$\mathbf{v}_w$     Word embedding of  $w$

$\mathbf{w}$     One-hot encoding of  $w$

$\mathbf{x}_i \in \mathcal{Y}$    Sentence  $\mathbf{x}_i$  belonging to language  $X$

$\mathbf{y}_i \in \mathcal{Y}$    Sentence  $\mathbf{y}_i$  belonging to language  $Y$

$\mathcal{M}$     A model  $\mathcal{M}$

$\mathcal{X}$     Set of sentences belonging to language  $X$

$\mathcal{Y}$	Set of sentences belonging to language $Y$
$w$	The word as a symbol
$\sigma(x)$	An activation function with respect to $x$
$N$	Number of data points
$V$	Dictionary; set of all words in vocabulary
$\mathbf{0}_d$	The $d$ -dimensional zero vector
$\mathbf{0}_{d \times d}$	The $d \times d$ -dimensional zero matrix
$A \in \mathbb{R}^{n \times m}$	Real matrix $A$ of dimensions $n \times m$
$A \odot B$	Elementwise multiplication of matrices $A$ and $B$
$I_d$	The $d \times d$ unit matrix
$\mathbf{x}$	A real column vector
$\mathbb{R}$	Set of real numbers
$\mathbb{R}^d$	Real vector space of dimension $d$
$\nabla_{\mathbf{x}} Q(\mathbf{x})$	Gradient of function $Q$ with respect to vector $\mathbf{x}$
$\text{pa}(x)$	Set of parent nodes for node $x$ in a graphical model
$\mathbf{0}_d$	The $d$ -dimensional zero vector
$\mathbf{0}_{d \times d}$	The $d \times d$ -dimensional zero matrix
$A \in \mathbb{R}^{n \times m}$	Real matrix $A$ of dimensions $n \times m$
$A \odot B$	Elementwise multiplication of matrices $A$ and $B$
$I_d$	The $d \times d$ unit matrix

$\mathbf{x}$	A real column vector
$\mathbb{R}$	Set of real numbers
$\mathbb{R}^d$	Real vector space of dimension $d$
$\nabla_{\mathbf{x}}Q(\mathbf{x})$	Gradient of function $Q$ with respect to vector $\mathbf{x}$
$\text{pa}(x)$	Set of parent nodes for node $x$ in a graphical model
$\boldsymbol{\theta} \in \Theta$	Parameters $\boldsymbol{\theta}$ belonging to the parameter space $\Theta$
$\mathbf{v}_w$	Word embedding of $w$
$\mathbf{w}$	One-hot encoding of $w$
$\mathbf{x}_i \in \mathcal{Y}$	Sentence $\mathbf{x}_i$ belonging to language $X$
$\mathbf{y}_i \in \mathcal{Y}$	Sentence $\mathbf{y}_i$ belonging to language $Y$
$\mathcal{M}$	A model $\mathcal{M}$
$\mathcal{X}$	Set of sentences belonging to language $X$
$\mathcal{Y}$	Set of sentences belonging to language $Y$
$w$	The word as a symbol
$\sigma(x)$	An activation function with respect to $x$
$N$	Number of data points
$V$	Dictionary; set of all words in vocabulary
GD	Gradient Descent
MCMC	Markov Chain Monte Carlo
MLE	Maximum Likelihood Estimation

NLP Natural Language Processing

NMT Neural Machine Translation

SGD Stochastic Gradient Descent

VI Variational Inference

$\text{Cov}(\mathbf{x}, \mathbf{y})$  Multivariate covariance with respect to  $\mathbf{x}$  and  $\mathbf{y}$

$\text{Cov}(x, y)$  Covariance with respect to  $x$  and  $y$

$\mathbb{E}_x[f(x, y)]$  Expectation of a function  $f(x, y)$  with respect to a random variable  $x$ .  
If no ambiguity to which variable we are taking expectation with respect to exists the subscript may be left out

$D_{KL}(p, q||)$  Kullback-Leibler divergence from distribution  $q$  to  $p$

$\mathcal{D}$  Set of data points

$\mathcal{N}(\boldsymbol{\mu}, \Sigma)$  Normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\Sigma$

$p(x)$  Probability distribution with respect to random variable  $x$

$\text{Cov}(\mathbf{x}, \mathbf{y})$  Multivariate covariance with respect to  $\mathbf{x}$  and  $\mathbf{y}$

$\text{Cov}(x, y)$  Covariance with respect to  $x$  and  $y$

$\mathbb{E}_x[f(x, y)]$  Expectation of a function  $f(x, y)$  with respect to a random variable  $x$ .  
If no ambiguity to which variable we are taking expectation with respect to exists the subscript may be left out

$D_{KL}(p, q||)$  Kullback-Leibler divergence from distribution  $q$  to  $p$

$\mathcal{D}$  Set of data points

$\mathcal{N}(\boldsymbol{\mu}, \Sigma)$  Normal distribution with mean  $\boldsymbol{\mu}$  and covariance matrix  $\Sigma$

$p(x)$  Probability distribution with respect to random variable  $x$

## Chapter 1

# Introduction

Natural language processing (NLP) is an old field with roots in many different disciplines including but not exclusive to electrical engineering, artificial intelligence and linguistics [1, p. 10-15]. Machine learning often view language as a statistical problem. A language model is a statistical model that trained on a training corpus assign probabilities to new sentences. Using latent variable models enables these language models to generate novel sentences from an underlying stochastic variable. While using probabilistic models enables us to do inference in a principled manner, it is often non-trivial to find the optimal parameters of latent variable models. Variational Inference approximates the log-likelihood through the Evidence Lower Bound (ELBO) which bound it from below, by introducing an approximate distribution of the conditional probability of the latent variable,  $p(\mathbf{z}|\mathbf{x})$  [2].

Although Variational Inference specifies ways of approximating the log-likelihood by a variational distribution  $q$  it does not specify the form of this  $q$ . Variational Autoencoder (VAE) [3], also under the name of stochastic backpropagation [4], is a framework for training deep generative models with latent variables by introducing a recognition model  $q_\phi$ . Except for the most basic of models the ELBO function is not tractable, however, using the reparametrisation trick and the law of large numbers, it is possibly to approximate the ELBO by sampling the latent variable instead of integrating it out, giving us the Stochastic Gradient Variational Bayes (SGVB) algorithm. The recognition model may then be optimized jointly with

the generative model  $p_\theta$  with respect to the SGVB objective [3]. As the recognition model maps the input sentences to the latent space, it effectively tries to compress information about the sentences through  $\mathbf{z}$ , resulting in distributed latent representations of sentences [5]. After learning the recognition mapping from  $\mathbf{x}, \mathbf{y}$  to  $\mathbf{z}$  we may use this form to do translation by omitting one of the input languages, effectively letting a source language decide the distribution over the latent space, which can be mapped back through the generative model to both  $\mathcal{X}$  and  $\mathcal{Y}$  performing both reconstruction and translation concurrently.

## Chapter 2

# Background Knowledge

This chapter will introduce the necessary background knowledge to follow the theoretical derivations later.

## 2.1 Probability Theory and Statistics

We go over the rules of probability and statistics.

### 2.1.1 Rules and Theorems

Rules of probability that will be useful to us are the following, if we let  $X \in \mathcal{X}$  be a random variable and  $p(X)$  the probability distribution with respect to this variable, then

**Unit volume**

$$\int_{\mathcal{X}} p(X) \, dX = 1 \quad (2.1)$$

**Non-negativity**

$$p(X) \geq 0 \quad (2.2)$$

where the integral is interpreted as the Lebesgue integral if  $X$  is continuous and as a sum over the possible values of  $X$  if it is discrete.

Most manipulations of random variables reduces to the following two rules of probability: given two random variables  $X, Y$  defined on the domains  $\mathcal{X}, \mathcal{Y}$  respectively, then

**Sum rule**

$$p(X) = \int_{\mathcal{Y}} p(X, Y) dY \quad (2.3)$$

**Product rule**

$$p(X, Y) = p(Y|X)p(X) \quad (2.4)$$

. The product rule leads directly to Bayes Rule

**Bayes Rule**

$$p(X|Y) = \frac{p(Y|X)p(X)}{p(Y)} \quad (2.5)$$

.

Two very important operations involving probabilities of random variables are those of *Expectation* and *Covariance*. Assume  $X \in \mathcal{X}$  is a random variable with probability distribution  $p(X)$  and  $f : \mathcal{X} \rightarrow \mathbb{R}$ :

**Expectation**

$$\mathbb{E}_X[f] = \int_{\mathcal{X}} f(X)p(X) dX \quad (2.6)$$

**Covariance**

$$\text{Cov}(X, Y) = \mathbb{E}_{XY}[(X - \mathbb{E}_X[X])(Y - \mathbb{E}_Y[Y])] \quad (2.7)$$

We then define the variance operator as:

**Variance**

$$\text{Var}(X) = \text{Cov}(X, X) \quad (2.8)$$

The generalisation from  $f : \mathcal{X} \rightarrow \mathbb{R}$  to  $f : \mathcal{X} \rightarrow \mathbb{R}^n$  is defined in the straightforward manner. If  $\mathbf{f} = f(X)$  then:

$$\mathbb{E}_X \begin{bmatrix} f_1 \\ \vdots \\ f_n \end{bmatrix} = \begin{bmatrix} \mathbb{E}_X f_1 \\ \vdots \\ \mathbb{E}_X f_n \end{bmatrix}$$

similarly  $\text{Cov}(\mathbf{f})$  is a  $D \times D$ -dimensional matrix where  $\text{Cov}(\mathbf{f})_{i,j} = \text{Cov}(f_i, f_j)$  [6].



### 2.1.2 The Gaussian Distribution

For a  $D$ -dimensional random vector  $X$ , the multivariate Gaussian distribution takes the form

$$\mathcal{N}(X|\boldsymbol{\mu}, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(X - \boldsymbol{\mu})^T \Sigma^{-1}(X - \boldsymbol{\mu})\right) \quad (2.9)$$

where  $\boldsymbol{\mu}$  is a  $D$ -dimensional mean vector,  $\Sigma$  is a  $D \times D$  dimensional positive definite covariance matrix, and  $|\Sigma|$  denotes the determinant of  $\Sigma$ . It is straightforward to show that these parameters correspond to the mean and covariance as defined in equations (2.6) and (2.7)[].

The Gaussian distribution can be seen as a unit  $D$ -dimensional cube which is translated, sheared and rotated, giving rise to the fact that we can write any Gaussianly distributed random variable  $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \Sigma)$  as a linear combination of a unit Gaussian random variable  $\mathbf{z} \sim \mathcal{N}(\mathbf{x}|\mathbf{0}_D, \mathbf{I}_{D \times D})$ . If we let  $\Lambda\Lambda^T = \Sigma$  be the Cholesky decomposition[7, p. 100-102] of  $\Sigma$ , then we also have that

$$\mathbf{x} = \boldsymbol{\mu} + \Lambda\mathbf{z} \quad (2.10)$$

, where the equality is in terms of distribution. If we further assume that  $\mathbf{x}$  is parametrised by  $\boldsymbol{\mu}$  and  $\Sigma$  such that  $\Sigma$  is diagonal positive definite with diagonal  $\boldsymbol{\sigma}$ , then  $\Sigma = \boldsymbol{\sigma} \odot \mathbf{I}_{D \times D}$ . Finally this means that if we want to sample a random variable  $\mathbf{x}$  with diagonal covariance structure, then we can do this by sampling a unit normal  $\mathbf{z}$  which we then transform, which we can express as

$$\mathbf{x} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma} \odot \mathbf{I}_{D \times D}) \quad (2.11)$$

As the Gaussian distribution is part of the exponential family, the density of joint distribution of iid Gaussian variables are themselves Gaussian distributed where the natural parameters of this joint distribution is the sum of the natural parameters of each random variable in the joint. In particular for the Gaussian distribution, this means that if we have a collection of iid gaussian random vari-

ables  $\{\mathbf{x}_i\}_i^n$ , such that  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{x}_i|\boldsymbol{\mu}_i, \Sigma_i)$ , then the joint can be found to be Gaussian distributed as

$$\mathcal{N}(\boldsymbol{\mu}, \Sigma)$$

, where

$$\Sigma = \left( \sum_i^n \Sigma_i^{-1} \right)^{-1} \quad (2.12)$$

$$\boldsymbol{\mu} = \Sigma \left( \sum_i^n \Sigma_i^{-1} \boldsymbol{\mu}_i \right) \quad (2.13)$$

[6, p. 78-84].

In the case of two random variables distributed according to the form as laid out in (2.10),  $\mathbf{x} \sim \mathcal{N}(\mathcal{N}(\boldsymbol{\mu}_x, \sigma_x \odot \mathbf{I}))$  and  $\mathbf{y} \sim \mathcal{N}(\mathcal{N}(\boldsymbol{\mu}_y, \sigma_y \odot \mathbf{I}))$  we have that the resulting distribution  $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x})p(\mathbf{y})$  is distributed such that

$$p(\mathbf{x}, \mathbf{y}) = \mathcal{N}(\boldsymbol{\mu}_{x,y}, \sigma_{x,y})$$

where

$$\sigma_{x,y} = \frac{1}{\sigma_x^{-1} + \sigma_y^{-1}} \quad (2.14)$$

$$\boldsymbol{\mu}_{x,y} = \frac{\sigma_x^{-1} \boldsymbol{\mu}_x + \sigma_y^{-1} \boldsymbol{\mu}_y}{\sigma_x^{-1} + \sigma_y^{-1}} \quad (2.15)$$

.

### 2.1.3 Maximum Likelihood Estimation

Assume we have a model  $\mathcal{M}$  parametrised by  $\boldsymbol{\theta}$  constrained to live in the parameter space  $\Theta$ . Given data  $\mathcal{D}$  we want to be able to fit the parameters  $\boldsymbol{\theta}$  such that these generalise to unseen data. The MLE of of the parameters of the model is defined to be

$$\hat{\boldsymbol{\theta}}_{ML} = \underset{\boldsymbol{\theta} \in \Theta}{\operatorname{argmax}} \mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) \quad (2.16)$$

.

While the original MLE is defined in terms of the likelihood function  $\mathcal{L}(\theta; \mathcal{D})$ , it's often more practical to work with the logarithm of this function, the log-likelihood function  $\ell(\theta; \mathcal{D})$ . Using the common assumption of i.i.d datapoints, the joint distribution becomes a product of individual probabilities for each datapoint,

$$\mathcal{L}(\theta|\mathcal{D}) = p(\mathbf{x}_1, \dots, \mathbf{x}_n|\theta) = \prod_i^n p(\mathbf{x}_i|\theta) \quad (2.17)$$

. Using the log-likelihood we transform this product into a form involving sums

$$\ell(\theta|\mathcal{D}) = \log \mathcal{L}(\theta|\mathcal{D}) = \sum_i^n \log p(\mathbf{x}_i) \quad (2.18)$$

. Besides from simplifying notation and calculation, it has the added benefit of reducing the risk of arithmetic underflow due to the small magnitude of individual probabilities<sup>1</sup>.

MLE estimators have a number of nice properties such as consistency and asymptotic normality, however, the optimization problem itself is often non-convex, making it hard to find the actual estimator[8].

### 2.1.4 Graphical models

tool for probabilistic modelling is the notion of using diagrams to specify the conditional relationships between random variables. A graphical model is a diagrammatic way of specifying this relationship by creating a Directed Acyclic Graph, a directed graph without any cycles. This representation is called a *Graphical Model* and provide a powerful way to visualize the structure of the probabilistic model and also how to use and abuse the structure of the model in order to infer variables in a computationally efficient way.

A graph in this setting consists of a set of *vertices* connected by *edges*, fol-

---

<sup>1</sup>Also, with the use of the log-sum-exp-trick,

$$\log \sum_{i=1}^n \exp(x_i) = \max_i x_i + \log \sum_{i=1}^n \exp(x_i - \max_i x_i)$$

this problem can be reduced even further

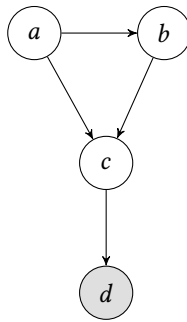
lowing the notation and nomenclature of graph theory as used in mathematics. While there are many different kinds of graphical models depending on the type of graph structure used (directed graphical models, undirected graphical models, factor graphs, etc.), we will only focus on the subset of graphical models called directed graphical models.

Directed graphical models specify how the joint distribution of a set of random variables  $\mathcal{X}$  factors in a conditional manner. In general, the relationship between a given directed graph and the corresponding distribution over the variables in  $\mathcal{X}$  is such that the joint distribution defined by the graph is given by the product, over all vertices of the graph, of a conditional distribution for each vertex conditioned on the variables corresponding to the parents of that vertex in the graph. So for a graph with  $K$  vertices, the joint distribution is given by

$$p(\mathcal{X}) = \prod_{k=1}^K p(x_k | \text{pa}(x_k)) \quad (2.19)$$

where  $\text{pa}(x_k)$  is defined as the set of random variables corresponding to the parent vertices of the random variable  $x_k$ .

Although a graphical model is completely defined in terms of its vertex and edge set, it is really the most powerful when visualized as a diagram. As an example I will repeat the above formulation in the context of the directed graphical model



There are two different vertices in this graphical model, the greyed out vertex indicates that the random variable is *observed* such that its value is fixed and known. The white vertices indicate latent random variables which we don't observe.

For this example we have the following factorisation of the joint distribution, following the rules laid out in equation (2.19),

$$p(a, b, c, d) = p(d|c)p(c|a, b)p(b|a)p(a)$$

### 2.1.5 Approximate Inference

While MLE is in many ways the optimal way that we can fit the model, it's only analytically and/or computationally feasible for very simple models which relies simple transformations and tractable distributional relationships. In cases where more powerful models are used it is very hard to find the MLE or even local maximum to the likelihood  $\mathcal{L}$ , or equivalently the log-likelihood  $\ell$ .

While in theory most of these problems can be resolved by MCMC sampling, which also practically have been implemented in the way of probabilistic programming with some success[9, Ch. 1][10, 11], most often this is too computationally intensive and has a large cost in terms of time. Approximate inference forms an alternative to MCMC for solving intractable densities by recasting this problem into an optimization problem instead of a sampling one as in MCMC.

The most common setting is in latent variable models, where a latent variable, often denoted by  $z$  is introduced to explain underlying causes to the observed variables, such as different clusters for Mixture of Gaussians or a lower-dimensional manifold in terms of the Factor Analysis model[6, page. 430-439, 583-586]. Latent variable models which rely on Gaussian distributions and linear relationships may be learned in an exact manner in the context of the EM algorithm or its many variations which guarantees parameters  $\hat{\theta}$  such that this point in the parameter space will yield a local maximum of the likelihood function[12, 13].

As theory and computing power has progressed, so has the advances in more powerful models. Compared to older models for which solutions often could be found analytically, these models were too non-linear, non-gaussian and complex to train in a straightforward manner, which can be seen directly from the numerous

heuristics that exist with regards how to train deep neural networks[14, 15].

The problem setting of variational inference is a latent variable model such that it can be split up into the observed variables  $\mathbf{x}$  and latent variables  $\mathbf{z}$  such that

$$p(\mathbf{z}, \mathbf{x}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) \quad (2.20)$$

. This also happens to cover the case of our generative model.

Technically, Approximate inference is a way of approximating a complicated distribution  $p(\mathbf{z}|\mathbf{x})$  by a distribution  $q(\mathbf{z})$  belonging to some constrained family of distributions  $\mathcal{Q}$ , for continuous distributions often such that  $\mathcal{Q} = \{\mathcal{N}(\boldsymbol{\mu}, \Sigma) | \boldsymbol{\mu} \in \mathbb{R}^d, \text{p.d } \Sigma \in \mathbb{R}^{d \times d}\}$ . The goal is then to find the element of  $\mathcal{Q}$  that minimizes some distance from  $p$  to  $q$ , most often the Kullback-Leibler divergence,

$$q^*(\mathbf{z}) = \underset{q(\mathbf{z}) \in \mathcal{Q}}{\operatorname{argmin}} KL(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}))) \quad (2.21)$$

. This optimal  $q^*$  may then be used as a pseudo-correct distribution in order to calculate other statistics and quantities.

## 2.2 Deep Learning

Until recently the field of NLP were dominated by older machine learning techniques utilising linear models trained over very high-dimensional and sparse feature vectors. Recently the field has switched over to neural networks over dense inputs instead using embeddings[16, p. 1 - 2].

What all neural networks have in common is that they are trying to find a functional relationship for the data, with the specific form of the function depending on the task. For NMT this reduces to finding the function  $f : \mathbf{x} \in \mathcal{X} \rightarrow \mathbf{y} \in \mathcal{Y}$  such that this  $f$  maximizes the likelihood  $P(\mathbf{x}|\mathbf{y})$ . Indeed many of the neural networks in existence has been shown to be universal approximators, theoretically being able to simulate a big set of nice functions[17].

### 2.2.1 Multilayer Perceptron

MLP's are neural networks represented by functional composition, where each function is interpreted as a layer of the network. The original MLP can be defined in terms of a recurrence relation such that if we have input vectors of the form  $\mathbf{x} \in \mathbb{R}^{d_{in}}$  and output vectors of the form  $\mathbf{y} \in \mathbb{R}^{d_{out}}$ , then an MLP with  $L$  layers have the functional form

$$f(\mathbf{x}|\boldsymbol{\theta}) = \sigma_L(\mathbf{W}_L \mathbf{z}_{L-1} + \mathbf{b}_L) \quad (2.22)$$

where for any  $l \in \{2, \dots, L-1\}$

$$\mathbf{z}_l = \sigma_l(\mathbf{W}_l \mathbf{z}_{l-1} + \mathbf{b}_l) \quad (2.23)$$

and with the base case

$$\mathbf{z}_1 = \sigma_1(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1) \quad (2.24)$$

.

$\mathbf{W}_l$  and  $\mathbf{b}_l$  may be of any dimension as long as it is dimensionally consistent with the input and output of the layer and conform to the original input and output dimensions. In this case we have that the parameters of the network are all of the biases and weights for the layers,  $\boldsymbol{\theta} = \{(\mathbf{W}_l, \mathbf{b}_l)_{l=1}^L\}$  (CITE PAPER ON MLP).

### 2.2.2 Recurrent Neural Networks

RNN's keeps information over long time-spans. It uses gates in order to let information flow forward in time. A recurrent neural network has the following functional form, given an input vector  $\mathbf{x}$  of length  $L$ :

$$\mathbf{s}_0 = \mathbf{0} \quad (2.25)$$

$$\mathbf{s}_t = f(\mathbf{U}\mathbf{x}_t + \mathbf{V}\mathbf{s}_{t-1}) \quad (2.26)$$

$$\mathbf{o}_t = \text{softmax}(\mathbf{V}\mathbf{s}_t) \quad (2.27)$$

for  $t \in \{0, \dots, L\}$  (CITE RNN PAPER).

In order to avoid vanishing or exploding gradients which means information does not propagate properly(CITE VANISHING GRADIENTS PROBLEM), we use the LSTM unit as the the function  $f$  which is defined implicitly as:

$$f_t = \sigma_g(U_f \mathbf{x}_t + V_f \mathbf{s}_{t-1}) \quad (2.28)$$

$$\mathbf{i}_t = \sigma_g(U_i \mathbf{x}_t + V_i \mathbf{s}_{t-1}) \quad (2.29)$$

$$\mathbf{o}_t = \sigma_g(U_o \mathbf{x}_t + V_o \mathbf{s}_{t-1}) \quad (2.30)$$

$$\mathbf{c}_t = f_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot \sigma_c(U_c \mathbf{x}_t + V_c \mathbf{s}_{t-1} + \mathbf{b}_c) \quad (2.31)$$

$$\mathbf{s}_t = \mathbf{o}_t \odot \sigma_s(\mathbf{c}_t) \quad (2.32)$$

where  $\mathbf{c}_0 = \mathbf{0}, \mathbf{s}_0 = \mathbf{0}$ , and  $f_t, \mathbf{i}_t, \mathbf{o}_t$  are the forget, input and output gates at time  $t$ . The activation functions are such that  $\sigma_g$  are sigmoid functions and  $\sigma_s, \sigma_c$  are hyperbolic tangents(CITE LSTM PAPER).

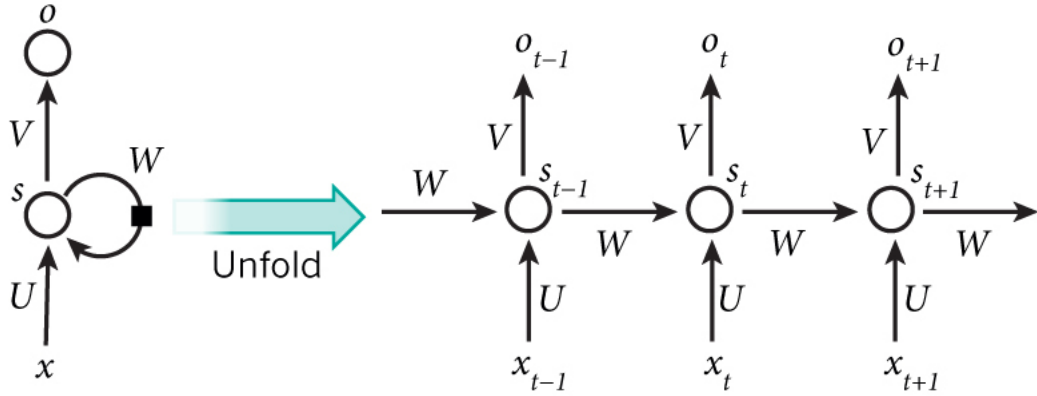


Figure 2.1: RNN schematics

### 2.2.3 Convolutional Neural Networks

The CNN we will use will be highly specialised for our case. We use the WaveNet as laid out in the paper by DeepMind.

WaveNet works on data of the form  $\mathbf{x}$ , where we describe the distribution as

$$p(\mathbf{x}) = \prod_{l=1}^L p(\mathbf{x}_l | \mathbf{x}_1, \dots, \mathbf{x}_{l-1})$$



which is the form common in NLP and NMT. The output of the model has the same time dimensionality as the input, meaning that the model has the ability to output a categorical distribution over the next value  $x_l$  with a softmax layer. In our case this is a parametrisation of the output probability at point  $l$ .

An important part is that it used dilated convolutions, convolution where the filter is applied to an area larger than its length by skipping input values with a certain step. Stacking the dilation layers with dilations being powers of 2 leads to output of the same dimensions as the input. This configuration of dilation factor results in exponential receptive field growth with depth (CITE WAVENET PAPER).

## 2.3 Natural Language Processing

Humans use natural language every day to convey concepts and abstractions to each other in an efficient manner. Compared to formal languages found in mathematics and programming, the natural languages we use are often ambiguous systems filled with rules and exceptions[18, 19].

Natural Language Processing is an old field that for a long time developed in parallel with the field of machine learning and computational statistics which deals with how to process information coming from human languages and is split up into several subfields such as Machine Translation, statistical parsing and sentiment analysis[19].

### 2.3.1 Language model

We define a sentence to be a vector of words  $\mathbf{w}_{1:L} = (w_1, \dots, w_l)^\top$  such that each word is an atomic element  $w_i \in V$ , where  $V$  is the dictionary of words in our language. Then the joint distribution of a word with respect to the underlying probability measure can be rolled out using the probabilistic chain rule which is just repeated application of the original product rule in equation (2.4)

$$P(\mathbf{w}_{1:L}) = \prod_{l=1}^L P(w_l | w_1, \dots, w_{l-1}) \quad (2.33)$$

where  $w_l$  is the  $l$ 'th word of the sentence  $\mathbf{w}_{1:L}$ [20].

### 2.3.2 Word embeddings

Breaking down sentences at a word level and processing them into a form that encodes information efficiently is a problem which has gained notorious recognition, leading to algorithms such as word2vec and Glove[21, 22, 23]. However, these techniques work less well in a neural network setting where instead finding the best embedding jointly with the parameters of the model is preferred [16, p. 5-7].

A straightforward way to represent the various words of the dictionary is as one-hot-encoded vectors such that a word  $w \in V$ , such that size of  $V$  is  $|V|$ , with an index  $i$  given by its place in the dictionary sorted alphabetically in descending order will have the vector representation

$$\mathbf{w} = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad (2.34)$$

[16, p. 6] such that  $w_j = \delta_{ij}$ .

While this is a form conceptually easy to understand, it fails to account for the curse of dimensionality as the size of the vocabulary might grow to millions of entries and the fact that the cosine similarity of two words  $v_1, v_2 \in V$  is zero unless they are the same word

$$\cos_{\text{similarity}}(\mathbf{v}_1, \mathbf{v}_2) = \delta_{v_1 v_2} \quad (2.35)$$

. This means that no meaning is embedded in the vector space except for the location in the sorted dictionary. Instead we would like to associate each word in the vocabulary with a distributed *word feature vector*, a dense, real-valued vector in

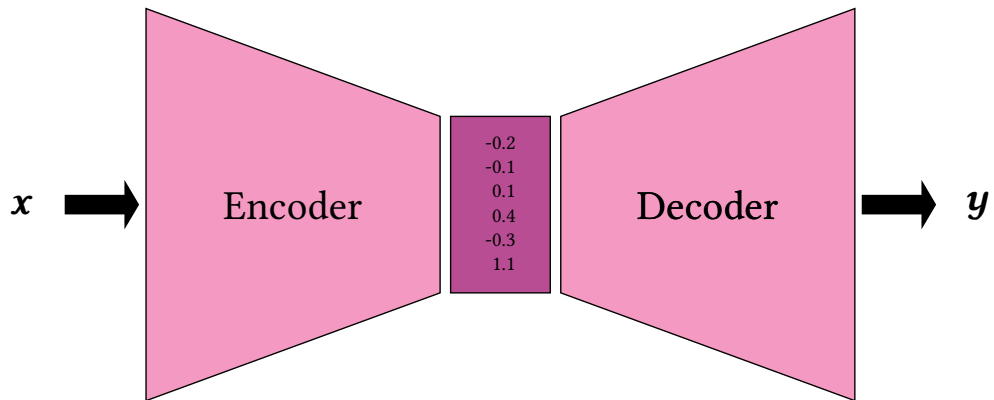
$\mathbb{R}^m$ . Express the joint probability function (2.33) of word sequences of a sentence in terms of the feature vectors of these words in the sequence and simultaneously learn the word feature vectors and the parameters of the model which dictates the form of the probability function,  $\theta$ . After this is done words which share similarities in some sense such as Dog, Puppy would have a higher similarity score than unrelated concepts such as Dog, Bulwark[20].

We may represent this in a mathematical form by trying to find a linear map  $C$  from any element  $w \in V$  such that  $C(w) \in \mathbb{R}^m$ . Using the usual canonical basis of an Euclidean space, we can express this linear map in terms of a matrix  $C \in \mathbb{R}^{m \times |V|}$ , thus the word feature vector of the learned embedding can be represented by the matrix multiplication  $Cw$ .

### 2.3.3 Neural machine translation

For a long time the dominant paradigm within machine translation was to use phrase based machine translation systems[24, 25], however since a couple of year back, modelling the word or character level directly with neural networks, so called NMT has become the best performing method[26, 27].

Most NMT models work in terms of an encoder-decoder architecture where the encoder extracts a fixed length representation  $c$ , often called a context vector, from a variable length input sentence  $x \in \mathcal{X}$ , and the decoder uses this representation to generate a correct translation  $y \in \mathcal{Y}$  from this representation[28].



**Figure 2.2:** Encoder decoder schematic

## 2.4 Optimization

### 2.4.1 Problem formulation

Numerical optimization has been applied to many parts of science and thus many different algorithms have been developed to tackle different problems. However, most of the theoretical results that exist deal with convex optimization, where the function we are trying to optimize is particularly well-behaved, leading to theoretical guarantees on the solution converged to. The surface of the function we are trying to optimize in machine learning in general and deep learning in particular is not well-behaved, being highly non-linear and non-convex, meaning most theoretical results that exist do not apply to this domain[29].

Nevertheless, while theoretical results are lacking, there has been substantial advances in various optimisation techniques fit to attack the highly non-linear, non-convex and high-dimensional optimization problems of learning in deep models, particularly from the Stochastic Gradient Descent. This has led to numerous gradient descent-like algorithms used machine learning[30].

### 2.4.2 Stochastic Gradient Descent

For a normal probabilistic machine learning problem, we have data  $\mathcal{D}$  that we try to model with a model  $\mathcal{M}$  parametrised by parameters  $\theta \in \Theta$ . The optimization problem can in the most general case be recast as an effort to find the parameters  $\theta_{ML}$  that maximizes the log-likelihood function (2.18).

As this can't be found analytically we have to resort to numerical schemes to find good candidates which hopefully will be close to the true solution  $\theta_{ML}$ . The first candidate is called Gradient Descent (GD). GD approximates the gradient as if it was linear and takes steps to maximize the probability through a hill-climbing scheme. If we let  $\mathcal{D} = \{z_i\}_{i=1}^n$  such that  $z_i$  is any general datapoint,  $Q$  an objective function that we try to maximize, such that  $Q : \Theta \times \mathcal{Z} \rightarrow \mathbb{R}$ , then the update using GD looks as follows

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} Q(z_i, \theta_t) \quad (2.36)$$

. While  $\gamma_t$  may vary with  $t$ , it is often fixed for practical reasons. As we can see

GD takes into account all of the datapoints in the set, in some sense updating the parameters with respect to the true linear approximation of the gradient.

However, with the huge size of data existing today, it is often not feasible to calculate the update for all datapoints in the dataset due to the time it would take and memory it would take to store the gradients.

Stochastic Gradient Descent is a similar algorithm to GD that instead of calculating the gradient with respect to the whole dataset calculates an approximate gradient, hence the word *stochastic*, as it picks a random subset of the dataset to update  $\theta$  with regards to. If we let  $I_t$  be a random subset of the indices of  $V$  of size  $m$  then SGD does the following update

$$\theta_{t+1} = \theta_t - \gamma_t \frac{1}{m} \sum_{i \in I_t} \nabla_{\theta} Q(z_i, \theta_t) \quad (2.37)$$

[31][6, p. 240].

The stochasticity has been shown to act as a regularizer and several analyses of this in terms of a Bayesian framework has been done, explaining the stationary behaviour of SGD with constant learning rate after convergence [32, 33].

### 2.4.3 ADAM

SGD has found widespread use within the machine learning community due to strong experimental results and ease of use, especially in deep learning. Lately though, a number of alternatives have sprung up that aims to improve the vanilla SGD, such as RMSProp[34] and AdaGrad[35].

Adam takes inspiration from RMSProp and AdaGrad. Technically, Adam keeps an exponential running average of the first and second order statistics of the gradient, using these to calculate an adaptive learning rate.

As laid out in the original paper, ADAM operates on a stochastic objective function  $f(\theta)$ , which in our case would be the probability averaged over our input batch.  $f_t(\theta)$  denotes this function over different timesteps, equally  $g_t = \nabla_{\theta} f_t(\theta)$ .  $m_t$  and  $v_t$  denotes the first and second order moments of the gradients, however, these are biased and are corrected in the algorithm.

The whole algorithm is as follows.

---

**Algorithm 1** ADAM
 

---

**Require:**  $\alpha$ : Stepsize

**Require:**  $\beta_1, \beta_2 \in [0, 1)$ : Exponential decay rates of the moment estimates

**Require:**  $f(\theta)$ : Stochastic objective function with parameters  $\theta$

**Require:**  $\theta_0$ : Initial parameter vector

```

1:  $m_0 \leftarrow 0$  (Initialize 1st moment vector)
2:  $v_0 \leftarrow 0$  (Initialize 2nd moment vector)
3:  $t \leftarrow 0$  (Initialize timestep)
4: while  $\theta_t$  not converged do
5:    $t \leftarrow t + 1$ 
6:    $g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$  (get gradients w.r.t stochastic objective at timestep  $t$ )
7:    $m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$  (Update biased first moment estimate)
8:    $v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$  (Update biased second raw moment estimate)
9:    $\hat{m}_1 \leftarrow m_t / (1 - \beta_1^t)$  (Compute bias-corrected first moment estimate)
10:   $\hat{v}_1 \leftarrow v_t / (1 - \beta_2^t)$  (Compute bias-corrected second raw moment estimate)
11:   $\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_1 / (\sqrt{\hat{v}_1} + \epsilon)$  (Update parameters)
return  $\theta_t$  (Resulting parameters)

```

---

Experimentally Adam has shown very good results on training various deep learning models such as MLP's, CNN's and RNN's so we will use it to train our models throughout this dissertation[36].

## Chapter 3

# Methods and Theory

### 3.1 ELBO

Ideally we would like to optimize the log-likelihood of the data,  $\ell(\theta; \mathcal{D})$ . For many models, it is not possible to evaluate this quantity directly due to it being computationally and/or analytically intractable.

Instead we make the following observation, letting  $\mathcal{D} = \{\mathbf{x}_i\}_{i=1}^N$

$$\log p_{\theta}(\mathcal{D}) = \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i) \quad (3.1)$$

and looking at each term in the sum we can see that we may rewrite this as

$$\begin{aligned} \log p_{\theta}(\mathbf{x}_i) &= \log \int_{\mathcal{Z}} p_{\theta}(z, \mathbf{x}_i) dz \\ &= \log \int_{\mathcal{Z}} q_{\varphi}(z|\mathbf{x}_i) \frac{p_{\theta}(\mathbf{x}_i, z)}{q_{\varphi}(z|\mathbf{x}_i)} dz \\ &= \log \mathbb{E}_{q_{\varphi}(z|\mathbf{x}_i)} \left[ \frac{p_{\theta}(\mathbf{x}_i, z)}{q_{\varphi}(z|\mathbf{x}_i)} \right] \\ &\geq \mathbb{E}_{q_{\varphi}(z|\mathbf{x}_i)} \left[ \log \frac{p_{\theta}(\mathbf{x}_i, z)}{q_{\varphi}(z|\mathbf{x}_i)} \right] \\ &= \mathbb{E}_{q_{\varphi}(z|\mathbf{x}_i)} \left[ \log \frac{p_{\theta}(\mathbf{x}_i|z)p_{\theta}(z)}{q_{\varphi}(z|\mathbf{x}_i)} \right] \\ &= \mathbb{E}_{q_{\varphi}(z|\mathbf{x}_i)} \left[ \log p_{\theta}(\mathbf{x}_i|z) \right] - D_{KL}(q_{\varphi}(z|\mathbf{x}_i) || p_{\theta}(z)) \end{aligned}$$

Which leads to the following observation, if we call the lower bound  $\mathcal{L}(\theta, \varphi; \mathbf{x})$ ,

that

$$\log p_{\theta}(\mathbf{x}) \geq \mathcal{L}(\theta, \varphi; \mathbf{x}) = \mathbb{E}_{q_{\varphi}(z|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|z)] - D_{KL}(q_{\varphi}(z|\mathbf{x})||p_{\theta}(z)) \quad (3.2)$$

In general this lower bound is still not tractable, however appealing to the Strong Law of Large Numbers (SLLN) we can sample  $z^l$  and use the reparametrisation trick (Include how the reparametrisation trick works) to approximate the lower bound by

$$\tilde{\mathcal{L}}^A(\theta, \varphi; \mathbf{x}) = \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}, z) - \log q_{\varphi}(z|\mathbf{x}) \quad (3.3)$$

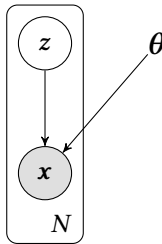
This can also be expressed in the closed form of the KL so

$$\tilde{\mathcal{L}}^A(\theta, \varphi; \mathbf{x}) = \left( \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(\mathbf{x}|z) \right) - D_{KL}(q_{\varphi}(z|\mathbf{x})||p_{\theta}(z)) \quad (3.4)$$

Since all quantities in this form have analytically known forms, it is simple to evaluate and sample from  $\tilde{\mathcal{L}}^A(\theta, \varphi; \mathbf{x})$ .

## 3.2 Models

The generative model follows the one laid out in [3]:



where  $z$  is the latent variable representing the latent sentence structure and the joint factorizes as  $p_{\theta}(\mathbf{x}, z) = p_{\theta}(\mathbf{x}|z)p(z)$ . In all models we let  $p(z) = \mathcal{N}(z|\mathbf{0}, I_{D \times D})$ .



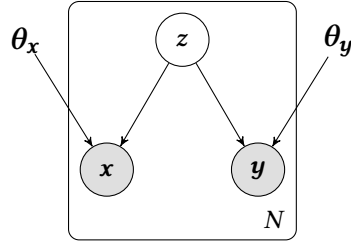
### 3.2.1 Reconstruction

The reconstruction model follows the original model. We assume a latent variable  $z$  that generates sentences  $\mathbf{x}$  belonging to a language  $\mathcal{X}$ .

The ELBO and related calculations follow the ones laid out in the ELBO section.

### 3.2.2 Translation

Translation adds an extra observed variable  $\mathbf{y}$ , a sentence belonging to  $\mathcal{Y}$ . The probabilistic model looks like:



From the graphical model we see that the joint distribution factorises such that

$$p(z, \mathbf{x}, \mathbf{y}) = p(z)p(\mathbf{x}|z)p(\mathbf{y}|z) = \mathcal{N}(0, I)\mathcal{N}(\mu_{\theta_x}(z), \sigma_{\theta_x}(z))\mathcal{N}(\mu_{\theta_y}(z), \sigma_{\theta_y}(z)) \quad (3.5)$$

. The notation  $\mathcal{N}(\mu_{\theta}(z), \sigma_{\theta}(z))$  means that the mean and variance diagonal of the normal are represented by a functional relationship  $f : z \mapsto (\mu, \sigma)$  where  $f$  is a function represented by a learnable neural network with parameters  $\theta$ .

The recognition model is parametrised by the probability distribution  $q_{\phi}(z|x, y)$ . We are free to choose this however we like, but the closer this is to the actual conditional distribution over the latent,  $p(z|x, y)$  the tighter the ELBO inequality  $\ell(x, y) \geq \mathcal{L}(x, y)$  will be. We let the distributions all be gaussians to have analytical tractability

$$\begin{aligned} q(z|x, y) &= q(z|x)q(z|y) \\ &= \mathcal{N}(\mu(x), \sigma(x)I)\mathcal{N}(\mu(y), \sigma(y)I) \end{aligned}$$

where both normals defined by the output of neural networks.

Since the product of Gaussian pdf's are gaussian itself, we have that may write that

$$q(z|\mathbf{x}, \mathbf{y}) = \mathcal{N}(z|\mu(\mathbf{x}, \mathbf{y}), \sigma(\mathbf{x}, \mathbf{y}))$$

which reduces to the expressions 2.13 for the mean and 2.12 for the covariance. Pulling all of these statements together, we get that the SGVB of the translation case reduces to the equation

$$\begin{aligned} \tilde{\mathcal{L}}^A(\theta, \varphi; \mathbf{x}) = & \left( \frac{1}{L} \sum_{l=1}^L \log(\mathcal{N}(\mathbf{x}|\mu_{\theta_y}(z), \sigma_{\theta_y}(z))) + \log(\mathcal{N}(\mathbf{y}|\mu_{\theta_y}(z), \sigma_{\theta_y}(z))) \right) + \\ & \frac{1}{2} \sum_{j=1}^J (1 + \log((\sigma_j^{(i)})^2) - (\mu_j^{(i)})^2 - (\sigma_j^{(i)})^2) \end{aligned} \quad (3.6)$$

where the  $\mu, \sigma$  is from the normal  $\mathcal{N}(z|\mu(\mathbf{x}, \mathbf{y}), \sigma(\mathbf{x}, \mathbf{y}))$ .

### 3.2.3 Training

Variational Autoencoders in NLP have a tendency to collapse the KL term in the SGVB objective to zero, effectively being reduced a RNNLM (CITE PAPER). As the KL can be seen as a measure of how much information is encoded in  $z$  this event makes the use of VAE's to model language-agnostic features in the latent space useless. One of the objectives is to make sure this KL-collapse doesn't happen.

Following the tips laid out in (CITE PAPER VAE FOR NLP) we use KL-annealing where we anneal the ELBO by a pseudo-objective

$$ELBO_{\beta_t} = \mathbb{E}_{q_{\varphi}(z|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|z)] - \beta_t * D_{KL}(q_{\varphi}(z|\mathbf{x})||p_{\theta}(z)) \quad (3.7)$$

where  $\beta_t$  is an annealing term going from 0 to 1 as we train. We do this by choosing a warm-up  $\beta$  and let  $\beta_t = \min(1, \frac{t}{\beta})$ , letting this warm up be done linearly until we recover the original ELBO objective, similar to normal simulated annealing (CITE ANNEALING PAPER).

## Chapter 4

# Experiments

### 4.1 Data

#### 4.1.1 Dataset

The dataset we have chosen to evaluate the models on is the Europarl dataset between languages English and French. Europarl is a dataset of the proceedings of the European Parliament, comprising in total of the 11 official languages of the European Union.

The dataset was chosen as the number of sentences for English and French is enough to be able to generalise (the uncompressed size of the full dataset is 619MB, 288MB for English, 311MB for French) to new sentences, and furthermore has established baseline for NMT in the form of BLEU scores for all the different language pairs in the full dataset, English-French in particular.[37]

#### 4.1.2 Preprocessing

The raw data is unfit for use directly with the model. For one thing the raw data is in the form of strings and in order to use neural networks we need to transform this into a form using vectors in some space  $\mathbb{R}^m$ .

From a linguistical point of view, a language contain a huge number of words at any point in time. Furthermore, words never leave the language completely, rather they go in and out of fashion and differ in how commonly they are use, also depending on what context they are used in.

**English:** I declare resumed the session of the European Parliament adjourned on Friday 17 December 1999, and I would like once again to wish you a happy new year in the hope that you enjoyed a pleasant festive period.

**French:** Je déclare reprise la session du Parlement européen qui avait été interrompue le vendredi 17 décembre dernier et je vous renouvelle tous mes vux en espérant que vous avez passé de bonnes vacances.

**Table 4.1:** A randomly sampled sentence from the Europarl corpus

- We specify the maximum and minimum length of the words. This means that when training and testing, we can pad the shorter sentences with Null tokens until they are of the same length as the longest sentence.
- We calculate the word frequencies in order to sort all of the words in the dataset in terms of how often it appear in absolute terms. This is then used to only retain the 30000 most common words. Words which are not part of this list gets replaced by an <UNK> token, specifying that it's an unknown word outside of the dictionary. This makes sure that only words which are prevalent enough such that the model can derive its relation to other words are part of the dictionary.
- Newline characters were removed and replaced by <EOS>, end-of-sentence tokens, signifying the end of a sentence.
- In parts where we have a dataset of 2 or more language sentences in parallel, we make sure that both of the the languages both satisfy the above criteria.

It is important to note that due to how languages differ, even though the dataset might consist of sentence pairs this will still not mean that in general both of the languages will have the same dictionary of words. Partially this is due to the different ways that languages are built up when expressing meaning, but on a more basic level, there are no bijection between languages as words have slightly different meaning and contexts, with some words only existing in one language but not the other.

## 4.2 Scores

We will evaluate our models on a variety of scores:

**ELBO** ELBO is the lower bound of the actual log-likelihood of the observed data

$$\sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

**Qualitative** Since natural language is not a formal in the sense that it is ambiguous, inconsistent and with exceptions to rules; any of these scores will be imperfect insofar as taking into account the feel of the generated sentences. Due to this we will inspect the sentences manually.

**BLEU** BLEU compares the generated sentences with sentences translated by professional translators, yielding a score telling us how well the generated translation does in relation to the translated benchmarks for each sentence.

We follow the thing as laid out in this link <http://www.statmt.org/wpt05/mt-shared-task/#TEST>.

**KL** Part of our investigation is about building models that take into account the latent space, enforcing the model to encode the information in the latent variable  $\mathbf{z}$  instead of the encoder/decoder part. Luckily, we have a quantitative measure of this, the KL divergence between the prior and the posterior  $q$ -distribution over  $\mathbf{z}$ ,

$$KL[q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\theta}(\mathbf{z})]$$

, where the KL is a measure of how much information is put into the  $q$ -distribution compared to just using the prior isotropic gaussian over  $\mathbf{z}$ ,  $p_{\theta}(\mathbf{z})$ .

**Perplexity** After training the model we want the sentences in the test set to be probable under the probability distribution parametrised by  $\theta$ . Perplexity is one measure of model fit. Given a sentence  $\mathbf{w}_{1:L}$  the perplexity, PP, is defined

as the per-word inverse probability,

$$\text{PP}(\mathbf{w}_{1:L}) = p(\mathbf{w}_{1:L})^{-\frac{1}{L}} \quad (4.1)$$

. In practice, this can be expressed in the form

$$\text{PP}(\mathbf{w}_{1:L}) = \exp\left(-\frac{1}{L} \prod_{l=1}^L \log p(\mathbf{w}_l | \mathbf{w}_1, \dots, \mathbf{w}_{l-1})\right) \quad (4.2)$$

### 4.3 Reconstruction

We ran a couple of experiments on reconstruction as a proof of concept to show that the model worked on data with only one language.

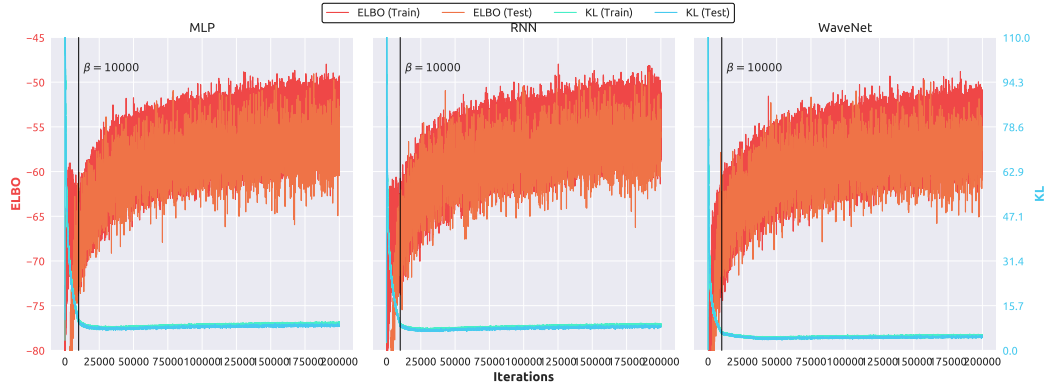


Figure 4.1: Runs with various different recognition model

### 4.4 Translation

## Chapter 5

# Conclusions

Possible extensions:

The Authors of WaveNet has introduced ByteNet, might be worthwhile looking into that as well as it is explicitly created to work well on character level generative models (<https://arxiv.org/pdf/1610.10099.pdf>).

Kingma et. al has come with a new extension to VAE's which uses normalising flows <https://arxiv.org/pdf/1606.04934.pdf>.

## **Appendix A**

# **Proofs**

These are the proofs for some of the results.



## **Appendix B**

# **Another Appendix About Things**

## Appendix C

# Colophon

*This is a description of the tools you used to make your thesis. It helps people make future documents, reminds you, and looks good.*

*(example)* This document was set in the Times Roman typeface using  $\text{\LaTeX}$  and Bib $\text{\TeX}$ , composed with a text editor.

# Bibliography

- [1] Daniel Jurafsky and James H. Martin. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 1st edition, 2000.
- [2] David M. Blei, Alp Kucukelbir, and Jon D. McAuliffe. Variational Inference: A Review for Statisticians. *Journal of the American Statistical Association*, 112(518):859–877, April 2017. arXiv: 1601.00670.
- [3] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes. *arXiv:1312.6114 [cs, stat]*, December 2013. arXiv: 1312.6114.
- [4] D. Jimenez Rezende, S. Mohamed, and D. Wierstra. Stochastic Backpropagation and Approximate Inference in Deep Generative Models. *ArXiv e-prints*, January 2014.
- [5] Samuel R. Bowman, Luke Vilnis, Oriol Vinyals, Andrew M. Dai, Rafal Jozefowicz, and Samy Bengio. Generating Sentences from a Continuous Space. *arXiv:1511.06349 [cs]*, November 2015. arXiv: 1511.06349.
- [6] Christopher M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- [7] William H. Press, Saul A. Teukolsky, William T. Vetterling, and Brian P. Flannery. *Numerical Recipes 3rd Edition: The Art of Scientific Computing*. Cambridge University Press, New York, NY, USA, 3 edition, 2007.

- [8] George Casella and Roger Berger. *Statistical Inference*. Duxbury Resource Center, June 2001.
- [9] Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. *Handbook of Markov Chain Monte Carlo*. CRC press, 2011.
- [10] Bob Carpenter, Daniel Lee, Marcus A. Brubaker, Allen Riddell, Andrew Gelman, Ben Goodrich, Jiqiang Guo, Matt Hoffman, Michael Betancourt, and Peter Li. Stan: A probabilistic programming language.
- [11] John Salvatier, Thomas V. Wiecki, and Christopher Fonnesbeck. Probabilistic programming in python using pymc3. *PeerJ Computer Science*, 2:e55, 2016.
- [12] A. P. Dempster, N. M. Laird, and D. B. Rubin. Maximum likelihood from incomplete data via the em algorithm. *JOURNAL OF THE ROYAL STATISTICAL SOCIETY, SERIES B*, 39(1):1–38, 1977.
- [13] Radford Neal and Geoffrey E. Hinton. A view of the em algorithm that justifies incremental, sparse, and other variants. In *Learning in Graphical Models*, pages 355–368. Kluwer Academic Publishers, 1998.
- [14] Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures. *arXiv:1206.5533 [cs]*, June 2012. arXiv: 1206.5533.
- [15] Hugo Larochelle, Yoshua Bengio, Jérôme Louradour, and Pascal Lamblin. Exploring strategies for training deep neural networks. *J. Mach. Learn. Res.*, 10:1–40, June 2009.
- [16] Yoav Goldberg. A primer on neural network models for natural language processing, 2015. cite arxiv:1510.00726.
- [17] K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Netw.*, 2(5):359–366, July 1989.
- [18] Ronald Rosenfeld. Two decades of statistical language modeling: Where do we go from here. In *Proceedings of the IEEE*, page 2000, 2000.

- [19] Lenhart Schubert. Computational linguistics. In Edward N. Zalta, editor, *The Stanford Encyclopedia of Philosophy*. Metaphysics Research Lab, Stanford University, spring 2015 edition, 2015.
- [20] Yoshua Bengio, Réjean Ducharme, Pascal Vincent, and Christian Janvin. A neural probabilistic language model. *J. Mach. Learn. Res.*, 3:1137–1155, March 2003.
- [21] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
- [22] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. Glove: Global vectors for word representation. In *In EMNLP*, 2014.
- [23] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado, and Jeffrey Dean. Distributed representations of words and phrases and their compositionality. In *Proceedings of the 26th International Conference on Neural Information Processing Systems, NIPS’13*, pages 3111–3119, USA, 2013. Curran Associates Inc.
- [24] Philipp Koehn, Franz Josef Och, and Daniel Marcu. Statistical phrase-based translation. In *Proceedings of the 2003 Conference of the North American Chapter of the Association for Computational Linguistics on Human Language Technology - Volume 1, NAACL ’03*, pages 48–54, Stroudsburg, PA, USA, 2003. Association for Computational Linguistics.
- [25] Philipp Koehn, Hieu Hoang, Alexandra Birch, Chris Callison-Burch, Marcello Federico, Nicola Bertoldi, Brooke Cowan, Wade Shen, Christine Moran, Richard Zens, Chris Dyer, Ondřej Bojar, Alexandra Constantin, and Evan Herbst. Moses: Open source toolkit for statistical machine translation. In *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions, ACL ’07*, pages 177–180, Stroudsburg, PA, USA, 2007. Association for Computational Linguistics.

- [26] Krzysztof Wołk and Krzysztof Marasek. Neural-based machine translation for medical text domain. Based on European Medicines Agency leaflet texts. *Procedia Computer Science*, 64:2–9, 2015. arXiv: 1509.08644.
- [27] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Lukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. Google’s Neural Machine Translation System: Bridging the Gap between Human and Machine Translation. *arXiv:1609.08144 [cs]*, September 2016. arXiv: 1609.08144.
- [28] Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the Properties of Neural Machine Translation: Encoder-Decoder Approaches. *arXiv:1409.1259 [cs, stat]*, September 2014. arXiv: 1409.1259.
- [29] Anna Choromanska, Mikael Henaff, Michael Mathieu, Gérard Ben Arous, and Yann LeCun. The Loss Surfaces of Multilayer Networks. *arXiv:1412.0233 [cs]*, November 2014. arXiv: 1412.0233.
- [30] Sebastian Ruder. An overview of gradient descent optimization algorithms, 2016.
- [31] Léon Bottou. Stochastic gradient descent tricks. In Grégoire Montavon, Genevieve B. Orr, and Klaus-Robert Müller, editors, *Neural Networks: Tricks of the Trade (2nd ed.)*, volume 7700 of *Lecture Notes in Computer Science*, pages 421–436. Springer, 2012.
- [32] Stephan Mandt, Matthew D. Hoffman, and David M. Blei. Stochastic Gradient Descent as Approximate Bayesian Inference. *arXiv:1704.04289 [cs, stat]*, April 2017. arXiv: 1704.04289.

- [33] Stephan Mandt, Matthew D. Hoffman, and David M. Blei. A Variational Analysis of Stochastic Gradient Algorithms. *arXiv:1602.02666 [cs, stat]*, February 2016. arXiv: 1602.02666.
- [34] T. Tieleman and G. Hinton. Lecture 6.5—RmsProp: Divide the gradient by a running average of its recent magnitude. COURSERA: Neural Networks for Machine Learning, 2012.
- [35] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. Technical Report UCB/EECS-2010-24, EECS Department, University of California, Berkeley, Mar 2010.
- [36] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. *arXiv:1412.6980 [cs]*, December 2014. arXiv: 1412.6980.
- [37] Philipp Koehn. Europarl: A Parallel Corpus for Statistical Machine Translation. In *Conference Proceedings: the tenth Machine Translation Summit*, pages 79–86, Phuket, Thailand, 2005. AAMT, AAMT.