

IMPERIAL COLLEGE LONDON

DEPARTMENT OF COMPUTING

---

# Chameleon Text: Exploring ways to increase variety in artificial data

---

*Author:*  
Thien P. Nguyen

*Supervisor:*  
Lucia Specia

Submitted in partial fulfillment of the requirements for the MSc degree in  
Computing (Artificial Intelligence) of Imperial College London

May 2019

# Chapter 1

## Literature Survey

### 1.1 Language Modelling

Language modelling is the task of predicting a word  $w_i$  in a text  $w$  given some sequence of previous words  $(w_1, w_2, \dots, w_{i-1})$ . More formally, an unconditional language model assigns a probability to a sequence of words,  $w = (w_1, w_2, \dots, w_{i-1})$ . This probability can be decomposed using the chain rule:

$$p(w) = p(w_1) \times p(w_2|w_1) \times p(w_3|w_1, w_2) \times \dots \times p(w_i|w_1, w_2, \dots, w_{i-1}) \quad (1.1)$$

$$p(w) = \prod_{t=1}^{|w|} p(w_t|w_1, \dots, w_{t-1}) \quad (1.2)$$

Traditionally, assigning words to probabilities may conflate syntactically dubious sentences but it remains to be a useful method for representing texts. It is one of the core natural language processing problems, and it used in a variety of contemporary applications, ranging from machine translation, to email response generation, to document summarisation.

In particular, we are more interested in conditional language modelling. This slightly differs from the definition described above - A conditional language model assigns probabilities to sequences of words,  $w = (w_1, w_2, \dots, w_{i-1})$ , given a conditioning variable,  $x$ .

$$p(w|x) = \prod_{t=1}^{|w|} p(w_t|x, w_1, \dots, w_{t-1}) \quad (1.3)$$

Modern language models revolve around the use of neural networks, which was started by bengio in 2003, with a simple MLP that encoded words. Neural Networks are non-linear statistical models that generate complex relationships between input and output vectors. Note that the input and output vectors are of a fixed dimension, which becomes a problem for our task at hand. Consequently, neural language models took off and has become the go to approach, giving rise to the introduction to RNNs being used to solve language model problems.

### 1.1.1 Recurrent Neural Networks

Recurrent neural networks (RNNs) are a particular class of neural networks such that the outputs are not necessarily restricted and discrete (as opposed to the MLP). The architecture of RNNs are especially important in NLP related problems as words in sentences are typically conditioned on the previous words. RNNs operate over a sequence of variable-length vectors, making them popular in contemporary NLP problems. Given a sequence of inputs  $(x_1, \dots, x_T)$ , a standard RNN computes a sequence of outputs  $(y_1, \dots, y_T)$  by iterating the following equation:

#### Problems

Theoretically, RNNs showed promise but there existed a multitude of problems. Firstly, It became apparent that it was very difficult for RNNs to leverage relationships between potentially relevant inputs and outputs - there isn't necessarily a clear indicator in the architecture that could facilitate this feature. This is described as a long range dependency problem. Furthermore, they were especially impractical due to the vanishing and exploding gradient problems having a larger effect on training. This is by design, as the architecture of the network.

### 1.1.2 LSTM

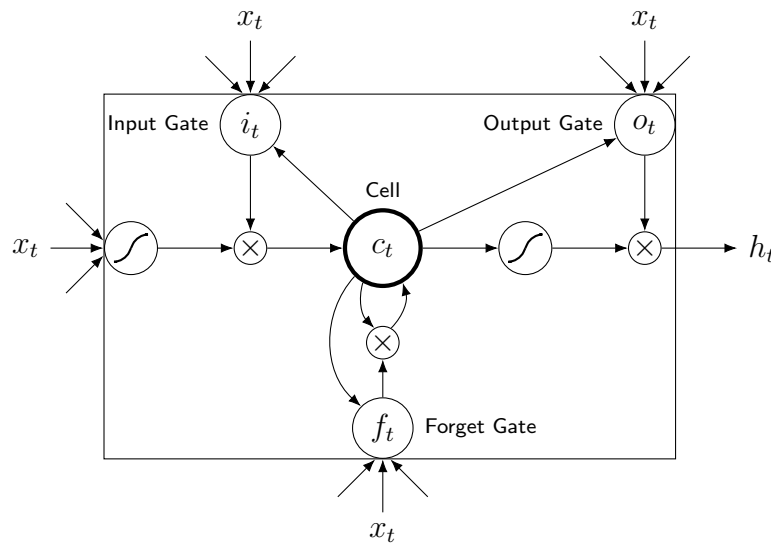


Figure 1.1: A diagram of the LSTM model.

LSTM (Long Short Term Memory) is a type of RNN unit that attempts to retain information based on the previous inputs through the introduction of gated architectures.

### 1.1.3 Gated Recurrent Units

Introduced by Cho et al. (2014), GRUs are used to solve the common issue with LSTMs where the training time was relatively slow due to the number of derivatives

necessary to compute. Within the GRU architecture, a feature to retain the previous weights remain, but there exists an direct path to the input data from the output, allowing a reduction in training time.

They are unable to clearly distinguish between the performance of the two gated units they tested. However, research from Chung et al. (2014) that GRUs were found to perform better than LSTMs on smaller datasets.

## 1.2 Related Work

### 1.2.1 Autoencoders

Autoencoders are a specialised form of MLPs where the model attempts to recreate the inputs on the output. Autoencoders typically have a neural network layer in the model where its dimension is smaller than the input space, therefore representing a dimensionality reduction in the data. Autoencoders are composed of two different principal components, an encoder network  $\alpha$  and a decoder network  $\beta$ , such that  $\alpha : X \rightarrow F$  and  $\beta : F \rightarrow X$ . Measuring the success of the reconstruction is deduced by a reconstruction loss formula. This reconstruction loss compares the output of the decoder and compares it against the input of the encoder. The two networks are trained together in a manner that allows them to preserve the input as much as possible.

Autoencoders are popularised through their use in Machine Translation, Word Embeddings, and document clustering.

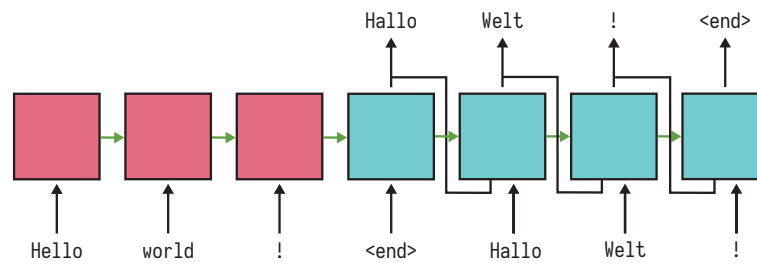
denoising autoencoders, where the input is a noisy representation of some input data, and it'll try and remove the noise (by comparing the reconstruction loss against a clean representation of the input).

### 1.2.2 Seq2Seq

Seq2Seq, introduced by Sutskever et al. (2014) is a modern interpretation of the autoencoder model,

Further work by XXX improved the original seq2seq model by providing an attention mechanism. The attention mechanism looks at all of the inputs from the hidden states of the encoders so far. This allows the decoder network to "attend" to different parts of the source input at each step of the output generation. This circumvents the need to encode the full source input into a fixed-length vector, helping it deal with long-range dependency problems.

Furthermore, this attention mechanism allows the model to learn what to attend to based on the input sequence and what it has so far, represented through a weighted combination of the two. Research from Bahdanau et al. (2014) has found this to work especially well in machine translation problems where languages which are relatively well aligned (such as English and German) - The decoder is most likely able to choose to attend to the response generation sequentially, such that the first output from the decoder can be generated based on the properties of the first input of the encoder and so on.

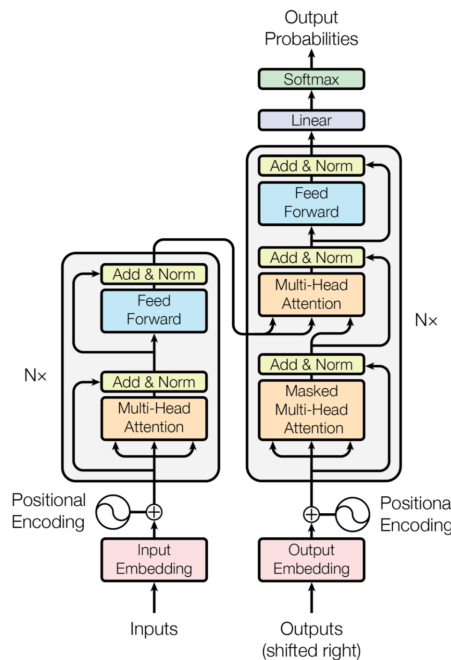


**Figure 1.2:** An abstracted model of the seq2seq architecture, where the encoder (pink) takes in the input sequence, and the decoder (blue) shows the output sequence.

Traditionally, Seq2seq would be the most common method of tackling this particular problem, but it also presents problems that make it suboptimal. Firstly, its discrete nature suggests that it is prone to noise in a similar fashion to how linear regression is not necessarily optimal as opposed to

### 1.2.3 Attention

### 1.2.4 Transformers

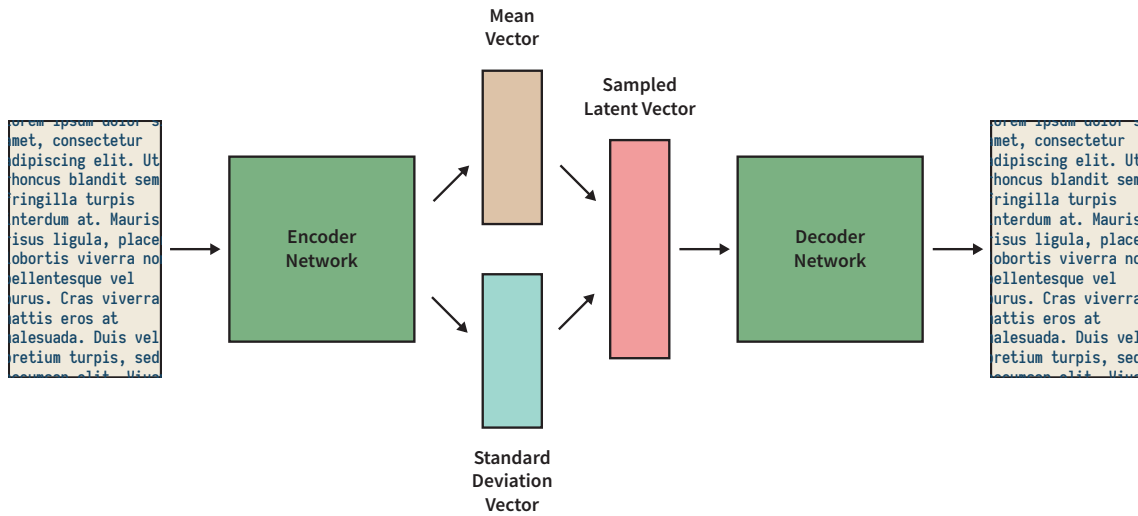


**Figure 1.3:** The transformer model architecture. Vaswani et al. (2017)

Vaswani et al. (2017) introduced the idea that it is possible to avoid the use of RNNs altogether and focus on leveraging the attention mechanism introduced in seq2seq. The resulting network architecture utilises stacked layers of residual networks

### 1.2.5 Variational Autoencoders

Our proposed solution revolves around the use of Variational Autoencoders. Variational Autoencoders (VAEs) introduce a constraint on the encoding network that forces the model to generate latent vectors that roughly follow a gaussian distribution, as opposed to creating a fixed latent vector. Consequently, VAEs introduce two extra vectors, a mean vector and a standard deviation vector, which is fed from the encoder. A sample is taken from the distribution and that is then fed into the decoder.



**Figure 1.4:** A simplified model architecture for a variational autoencoder, which takes as input some text, and its predicted output being the same text as the input.

Note that the decoder receives samples from a non-standard normal distribution produced by the encoder. The average of the samples of the different distributions should approximate to a standard normal.

Due to the stochastic nature of the network, we use a reconstruction loss that involves an expectation of the output; but we also use the KL divergence, which measures the relative difference of two probability distributions. In this particular case, we will be comparing the distribution of the decoder outputs against a standard gaussian  $\mathcal{N}(0, 1)$ .

$$\mathcal{L}(\theta, \phi, x, z) = \mathbb{E}_{q_{\phi}(z|x)}[\log p_{\theta}(x|z)] - D_{KL}(q_{\phi}(z|x) || p(z))$$

$$D_{KL}(P||Q) = \sum_{x \in X} P(x) \cdot \log\left(\frac{P(x)}{Q(x)}\right)$$

In other words, it is the expectation of the logarithmic difference between the probabilities  $P$  and  $Q$ , where the expectation of  $P$  is already known.

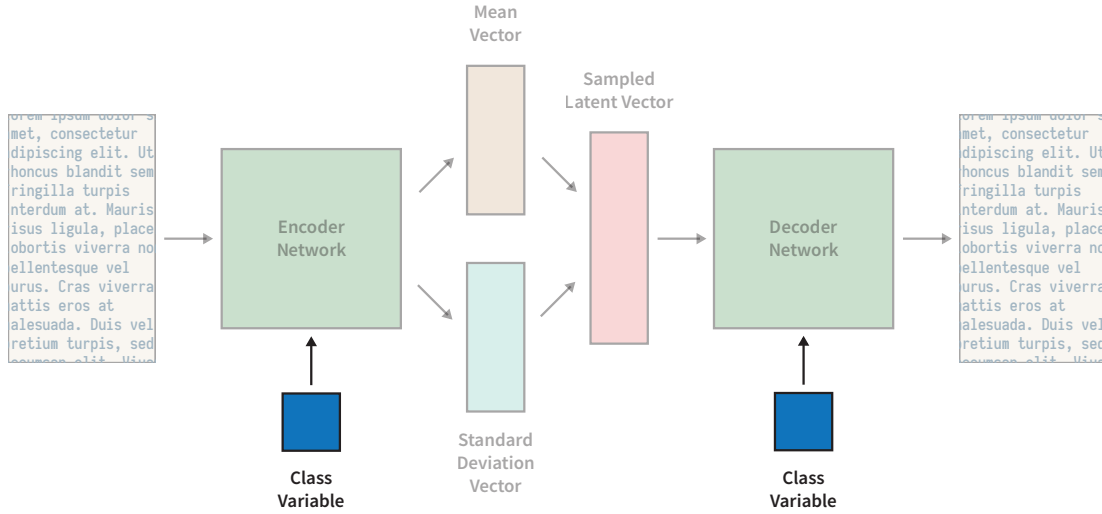
- You'll need to perform a reparameterisation trick (since you can't calculate derivatives of samples.) in order to perform backpropagation. (you can't push gradients through a sampling node.)

$$z = \mu + \sigma \cdot \epsilon$$

where  $\epsilon \sim \mathcal{N}(0, 1)$ . You want to learn  $\mu, \sigma$ .

### 1.2.6 Conditional Variational Autoencoders

Although VAEs are more robust when compared to their original autoencoder counterparts, the decoder class cannot produce outputs of a particular class on demand. CVAEs are an improved model of the original VAE architecture by conditioning on another description of the data, a class descriptor  $y$ .



**Figure 1.5:** A model architecture for a CVAE, which includes the label being fed into the encoder and decoder networks.

During training time, a class (represented by some arbitrary vector) is fed at the same time to the encoder and decoder. To generate an output that depends on  $y$  we feed that number to the decoder along with a random point in the latent space sampled from a standard normal distribution.

Samples can be generated from the conditional distribution  $p(x|y)$ . By changing the value of  $y$ , we can get corresponding samples  $x \sim p(x|y)$ . The system no longer relies on the latent space to encode what output is necessary; instead the latent space encodes other information that can distinguish itself based on the differing  $y$  values.

## 1.3 Variational Autoregressive Decoders

Introduced by Du et al. (2018), Variational Autoregressive Decoders (VADs) attempt to circumvent the sampling problem introduced from CVAEs by introducing multiple latent variables into the autoregressive Decoder. At different time-steps, this allows the decoder to produce a multimodal distribution of text sequences, allowing a variety of responses to be produced.

VADs use the seq2seq architecture as the base with variable-length queries  $x = \{x_1, x_2, \dots, x_n\}$ , and  $y = \{y_1, y_2, \dots, y_n\}$  representing the input and output responses respectively. The encoder network is a Bidirectional RNN with GRUs. The decoder network is an unidirectional RNN with GRUs. For each timestep  $t$ , each GRU in the decoder network is encoded with hidden state  $h_t^d$ .

### **Sequential Model**

Sequential

### **Prior Model**

### **Sequential Bag of Words**

### **Learning Mechanism**

VAD's propose



# Bibliography

- Bahdanau, D., Cho, K., and Bengio, Y. (2014). Neural Machine Translation by Jointly Learning to Align and Translate. *arXiv:1409.0473 [cs, stat]*. arXiv: 1409.0473. pages 3
- Cho, K., van Merriënboer, B., Bahdanau, D., and Bengio, Y. (2014). On the Properties of Neural Machine Translation: Encoder-Decoder Approaches. *arXiv:1409.1259 [cs, stat]*. arXiv: 1409.1259. pages 2
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. (2014). Empirical Evaluation of Gated Recurrent Neural Networks on Sequence Modeling. *arXiv:1412.3555 [cs]*. arXiv: 1412.3555. pages 3
- Du, J., Li, W., He, Y., Xu, R., Bing, L., and Wang, X. (2018). Variational Autoregressive Decoder for Neural Response Generation. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3154–3163, Brussels, Belgium. Association for Computational Linguistics. pages 6
- Sutskever, I., Vinyals, O., and Le, Q. V. (2014). Sequence to Sequence Learning with Neural Networks. *arXiv:1409.3215 [cs]*. arXiv: 1409.3215. pages 3
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., and Polosukhin, I. (2017). Attention Is All You Need. *arXiv:1706.03762 [cs]*. arXiv: 1706.03762. pages 4