

# Sequence to Sequence Modelling

Patrick Kahardipraja

September 2019

Mapping of a sequence to another sequence is an important paradigm because of the vast amount of problems that can be formulated in this manner. For instance, in automatic speech recognition (ASR), chunks of speech signals can be mapped to sequence of phonemes while in machine translation, a sequence of words in one language can be mapped to another language. Interestingly, many other tasks such as text summarization, question answering, and image caption generation can be phrased as a sequence to sequence problem. In this paper, I will attempt to distill how sequence to sequence learning works and the motivation behind it, with a particular focus on machine translation.

## Introduction

Prior to neural machine translation (NMT), phrase-based statistical machine translation (SMT) systems are widely used as it offers reliable performance. Despite its success, most of them are extremely complex and require a huge amount of effort, as it is often tailored to a specific language pair and do not generalize well to another languages. Furthermore, a lot of feature

engineering are required in order to capture a specific language phenomena, which prompt researchers to explore another approach.

The resurgence of deep neural networks (DNNs) in early 2010s, thanks to faster, parallel computation using GPUs and availability of large and high-quality datasets, bring a new wave of enthusiasm in deep models. With the capability to learn features automatically with multiple, hierarchical representation, DNNs achieve excellent performance on difficult tasks in computer vision [AlexNet] and speech recognition []. Albeit powerful, DNNs has its own limitation, as it requires input and output vectors with a fixed dimension and thus not suitable for sequence to sequence problem whose lengths are unknown beforehand. In addition, DNNs also do not generalize well across temporal patterns, because each neuron has its own specific connection and as a result, a single pattern may look totally different at different timesteps.

The natural remedy for this problem is to look onto recurrent neural networks (RNNs), as it allows operations over sequences of vectors. However, mapping using RNNs typically have one-to-one correspondence between the input vectors and the output vectors. It also has another problem, as the input and output sequences can have different lengths and non-monotonic alignments. Standard RNN architecture is also not reliable for learning long-range dependencies due to the vanishing gradient problem. This issue is addressed by Sutskever et al. [Seq2Seq], where they introduce a novel and straightforward method to solve general sequence to sequence mapping using Long-Short Term Memory (LSTM) architecture. With the success of sequence to sequence learning in machine translation tasks, research in neural machine translation continue to thrive, eventually resulting in many significant improvements such as attention mechanism [Bahdanau] and subword units to deal with rare words [WordPiece]. But, before delving in too deep, I will give some brief insight into the mechanism behind RNN and LSTM in the next section.

# Recurrent Neural Networks

Recurrent neural networks [Rumelhart] are type of neural network that is able to process arbitrary sequential input via combination of its internal state and input vector. At every timestep  $t$ , the hidden state vector  $h_t$  is overwritten as a function of the hidden state at the previous timestep  $h_{t-1}$  and the current input vector  $x_t$ . The input vector  $x_t$  itself could be a representation of  $t$ -th word in a sentence, which is usually obtained using pre-trained word embeddings [GloVe, Word2Vec, ElMo]. The hidden state of RNNs can be perceived as a memory with a fixed dimensionality that can be tuned, containing distributed representation of the processed input sequence up to time  $t$ .

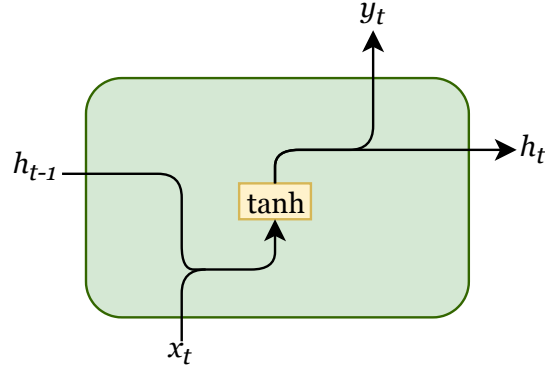
In a RNN, the forward step function consists of an affine transformation followed by a non-linear activation function. The hidden state then can be used to make predictions:

$$\begin{aligned}h_t &= a(W^{(x)}x_t + W^{(h)}h_{t-1} + b_h) \\ y_t &= g(W^{(y)}h_t + b_y)\end{aligned}$$

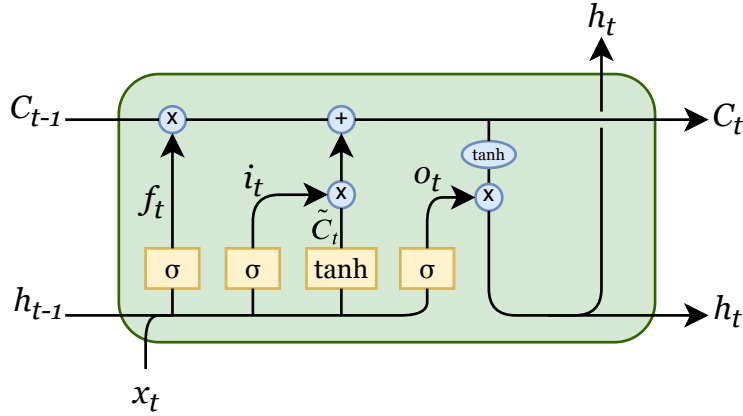
where in a typical application,  $a$  is the hyperbolic tangent function and  $g$  is the softmax function.

Although proven to be effective, RNN still has its own shortcoming. The main problem with RNN is that during training, magnitude of gradient can get weaker or stronger exponentially when backpropagating the error through time, especially with long sequences [Hochreiter, Bengio]. This phenomena is called vanishing or exploding gradient problem, which causes RNN model to experience difficulty when handling "long-term dependencies" that occur in a sequence.

Long Short Term Memory (LSTM) architecture [Hochreiter and Schmidhuber] addresses the problem of "long-term dependencies" by integrating a memory cell that is capable to memorize state that span over long sequences of time. The memory cell is controlled by gates, which have the ability to regulate how much information are added or removed in the memory cell. This means that while in a RNN a completely new hidden state is computed at every new timestep, in LSTM the hidden state is not completely overwritten, and updated according to the memory cell. The architecture of both RNNs and LSTMs are depicted in Figure 1.



(a) RNN unit



(b) LSTM unit, colah

Figure 1: Architecture of RNN (a) and LSTM (b)

A LSTM unit consists of 3 gates (input gate  $i_t$ , forget gate  $f_t$ , output gate  $o_t$ ), memory cell  $C_t$  and hidden state  $h_t$ . In a high-level sense, the input gate decides how much and which values will be updated, the forget gate controls the amount of information to be forgotten in the previous memory cell, and the output gate decides the hidden state by filtering the internal memory cell for each timestep. Each gate produces vectors, where their values are between 0 (completely closed) and 1 (completely open) using the sigmoid activation function.

The formula of LSTM is described with the following equations:

$$\begin{aligned} i_t &= \sigma(W^{(i)}x_t + U^{(i)}h_{t-1} + b_i) \\ f_t &= \sigma(W^{(f)}x_t + U^{(f)}h_{t-1} + b_f) \\ o_t &= \sigma(W^{(o)}x_t + U^{(o)}h_{t-1} + b_o) \\ \tilde{C}_t &= \tanh(W^{(c)}x_t + U^{(c)}h_{t-1} + b_c) \\ C_t &= f_t \odot C_{t-1} + i_t \odot \tilde{C}_t \\ h_t &= o_t \odot \tanh(C_t) \end{aligned}$$

where  $x_t$  is the input vector for timestep  $t$ ,  $\sigma$  is the sigmoid activation function and  $\odot$  denotes the Hadamard product of matrices.

## Variants

Beside the standard LSTM architecture in literature that is introduced above, there also exist several popular variants which are commonly used. One of the variations which is introduced by [Gers, Schmidhueber] use "peep-hole connections". These connections allow the gates to look into the memory cell state in order to learn precise and stable timings.

Other notable LSTM variants is the Gated Recurrent Unit (GRU), in-

roduced by [Cho et al.]. Unlike LSTM, GRU is less complex and requires less computation. GRU only have 2 gates, the update gate that decides how much information will be transferred from previous and candidate hidden state to the current one and reset gate that controls to what extent the previous hidden state will affect the candidate hidden state. In this manner, the update gate can be thought as a combination of input and forget gates of LSTM unit. This architecture also merges the internal memory cell and hidden state of LSTM into a single hidden state.

## Seq2Seq Model

The first model that is able to map a sentence into a vector and then to its translation is introduced by [Kalchbrenner et al.]. The model, which they called Recurrent Continuous Translation Models (RCTM), is composed of 2 separate parts: convolutional neural network for modelling the source sentence as the encoder and recurrent neural network for translation generation as a language modelling task, conditioned on the source sentence as the decoder. With this approach, the encoder can capture all the information contained in the source word representations and create a representation of the source sentences. The representation for the source sentences also restraint the generation of the target words in the language modelling phase.

The Recurrent Continuous Translation Models estimate the probability distribution over the sentence in the target language given sentences in the source language. Suppose that there exist a target sentence  $f = f_1, f_2, \dots, f_m$ , which is a translation of source sentence  $e = e_1, e_2, \dots, e_n$ . Then  $P(f|e)$  can be obtained with the formula:

$$P(f|e) = \prod_{i=1}^m P(f_i|f_{1:i-1}, e)$$

As can be seen in the formulation above, the model estimates  $P(f|e)$  by calculating the conditional probability  $P(f_i|f_{1:i-1}, e)$  for every translated word occurring at position  $i$ , given the preceding generated words  $f_{1:i-1}$  in the target sentence and the source sentence  $e$ . Conditioning the translation model to the preceding target words also ensure that it incorporates the target language model [Kalchbrenner et al.].

In RCTM, prediction of the target sentence use a language model based on a recurrent neural network [Mikolov et al.]. The recurrent language model predict the  $i$ -th word of the target sentence depending on all the previous generated words  $f_{1:i-1}$ , making no Markov assumption about the words dependencies in the target sentence. However, using the standard RNN architecture makes the prediction to be strongly affected by words close to  $f_i$  and weakly influenced by long-range dependencies that occur in the target language due to the nature of RNN.

Describe RLM. dreaming

## Recent Advances

## Conclusion