

LUDWIG-MAXIMILIANS-UNIVERSITY

DEPARTMENT OF STATISTICS



Sequence Modeling:
Recurrent and Recursive Nets

Seminar

Markus Dumke

Munich, 27th December 2016

Abstract

In the last years deep learning has received lots of attention and many applications and models were designed. One especially useful architecture developed were Recurrent Neural Networks (RNNs). RNNs are neural networks suited for sequential data. In particular Recurrent Neural Networks have been applied to language tasks such as speech recognition and generation (cite some papers?), machine translation, image captioning, as part of chatbots and all applications which involve processing text input and generating new text as an output. This paper will review different recurrent neural network architectures and present the challenges when optimizing such models. In the last part we will have a look at Google's Neural Machine Translation, recently introduced by (cite) and now deployed by Google Translate service.

Contents

1	Introduction	3
2	Recurrent Neural Network	4
3	Optimization and the Vanishing Gradient Problem	6
3.1	Vanishing Gradient Problem	6
3.2	Solutions to Optimization Difficulties	7
4	Gated RNNs	8
4.1	Long-Short-Term-Memory (LSTM)	8
4.2	Gated Recurrent Unit (GRU)	9
5	Application 1: LSTM	9
5.1	Sampling from an RNN	9
6	Extensions of Recurrent Neural Network	10
6.1	Recurrent Connections from output to hidden units	10
6.2	One-output RNN	11
6.3	RNN based on context	11
6.4	Deep RNN	11
6.5	Bidirectional RNN	12
6.6	Recursive Neural Network	12
6.7	Encoder-Decoder Network	13
6.8	Attention	13
6.9	Residual Network	15
7	Application 2: Google's Neural Machine Translation	16

1 Introduction

In recent years Recurrent Neural Networks (abbreviated RNNs) have been applied to various tasks involving sequential data, especially in the field of natural language processing. Applications include handwriting recognition and generation, machine translation, sentiment analysis, speech recognition and generation, image captioning and many more.

In comparison to standard feed-forward neural networks, RNNs share parameters across time and are therefore much more efficient for sequential tasks as they can learn the structure of language, which does not change over time, while a feed-forward network must learn all rules of language for each position of time anew. For example a sentence might contain some time information (e.g. a date) and the information content is independent whether the date appears in the beginning or the end of the sentence.

RNNs can therefore overcome the two major drawbacks of feed-forward neural networks: the independence assumption and the fixed length assumption. They can process (in theory) arbitrary length sequences, one input at each time step, whilst feed-forward neural networks are dependent on a fixed-size input.

The hidden states can be seen as some kind of memory about the past thus making RNNs very powerful learners and some RNNs are able to capture long-term dependencies.

In this paper we will explore basic RNN architectures and their applications, starting with Vanilla RNNs, going through the optimization and how to overcome the associated problems (especially the Vanishing Gradient Problem). Thereafter extensions of the basic RNN architecture like the LSTM, GRU, Deep RNNs, Bidirectional RNNs, Recursive Neural Networks, Encoder-Decoder Networks, the attention mechanism and Residual Networks will be covered. In the final part Google's Neural Machine Translation as a practical use-case will be explained.

2 Recurrent Neural Network

There exists a vast range of different Recurrent Neural Network architectures designed for different applications. Usually they have a hidden state h_t , which is shared over time, and inputs and outputs at one or more time steps. Recurrent Neural Networks can be seen as a dynamical system with the hidden state as a physical component like a biological cell, which is "updated" at each point in time. The same architecture can be unrolled over time showing the computational graph for each time step next to each other, compare Figure 1. In the following we will use the unrolled graph representation as this allows to look at which computations are performed in each time step.

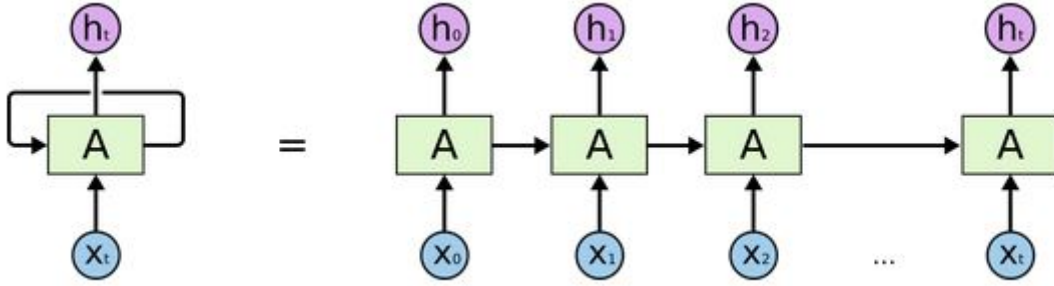


Figure 1: An RNN can be represented as a dynamical system or unrolled over time.
Source: Olah (2015b)

The probably most basic Recurrent Neural Network (known as Vanilla RNN) consists of the following structure: At each time step an input $x^{(t)}$ is processed and an output, e.g. a probability distribution over words, is computed. The input could be a word or symbol represented as a vector, the output could be a probability distribution over all words/symbols, which can be used to sample the next word or symbol from the model. Such a model can be applied to language modelling, e.g. as a part of a chatbot. In the middle of the network are hidden states, which are fixed-sized vectors capturing relevant information from all previous inputs by learning a useful representation of the data. Recurrent connections between hidden states allow information from the past to be passed forward in time. Figure 2 shows the unrolled computational graph. The RNN can be written as

$$\begin{aligned}
 a_t &= b + W h_{t-1} + U x_t \\
 h_t &= \tanh(a_t) \\
 o_t &= c + V h_t \\
 \hat{y}_t &= \text{softmax}(o_t).
 \end{aligned} \tag{1}$$

Note that parameters are shared over time, so at each time step the same activation function, weights W , U , V and biases b , c are used. This corresponds to a stationarity assumption. The most commonly used activation function for Recurrent Neural Networks is the tanh function,

which is symmetric around 0 with all values in the interval of $[-1, 1]$.

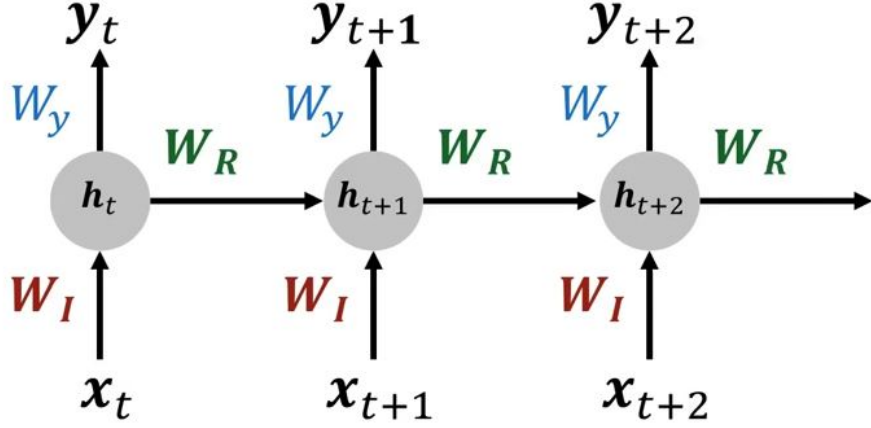


Figure 2: A Vanilla Recurrent Neural Network unrolled over time.
Source: Nervana (2016)

A simple character-level language model could look like this:

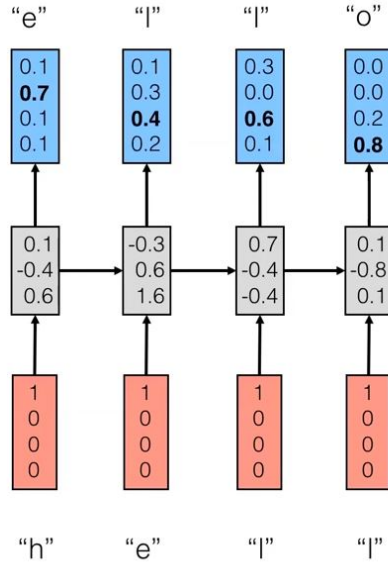


Figure 3: A simple application of a vanilla RNN for language modeling.
Source: Nervana (2016)

At each time point an input vector x_t representing a single character in one-hot-coding is processed, a single hidden state, a vector h_t , is computed and an output \hat{y}_t , a probability distribution over all characters is omitted. Based on

$$P(y_1, y_2, \dots, y_n) = \prod_{i=1}^n P(y_i | y_1, \dots, y_{i-1}) \quad (2)$$

the next character can be sampled based on the history of all previous characters.

3 Optimization and the Vanishing Gradient Problem

Recurrent Neural Networks can be trained with Backpropagation-through-time (BPTT), which is similar to the usual back-propagation algorithm used for training feed-forward neural networks, but contains connections over time. Usually minibatches of the training data are used for training. In the forward pass one minibatch (e.g. 32 symbols represented as vectors) is processed and moving forward in time each hidden state and output is computed according to Equation (1). For each time step the loss L_t between the predicted output \hat{y}_t and the training target y_t is computed. The choice of loss function depends on the distribution of the targets, e.g. the cross-entropy loss for a categorical output distribution or more general the negative log-likelihood of the target distribution. The overall loss is then the sum of losses for each time step,

$$L(y, \hat{y}) = \sum_t L(y_t, \hat{y}_t) = - \sum_t \log p(y_t | x_t) \quad (3)$$

In the backward pass the derivations of the final loss regarding all parameters are computed. Starting with the final time step we work our way backwards through time to obtain the gradient corresponding to the parameters using the chain rule of calculus. Then a Stochastic Gradient Descent (SGD) update can be applied to all parameters:

$$W = W - \epsilon \nabla_W L \quad (4)$$

and similar updates to all other parameters U , V , b , c , with learning rate ϵ .

As we are training on minibatches of training cases the computed gradient is a stochastic version of the true gradient obtained using all training cases. Each parameter is updated accordingly to its "share" of the final loss. SGD minimises the loss towards a local minimum.

3.1 Vanishing Gradient Problem

The backpropagation of the gradient through a large number of time steps results in many multiplications of the same weights similar to the power method. If weights are smaller than 1, this leads to Vanishing Gradients else to exploding gradients, if weights are greater than 1. Both causes the RNN to fail learning long-term dependencies, so that only the very last time steps influence the current update. This is known as the Vanishing Gradient or Exploding Gradient problem. Whilst exploding gradients can be dealt with by Gradient Clipping, vanishing gradients are a much more sincere problem and special architectures are needed to overcome the problem. The most famous architecture is the Long-Short-Term-Memory (LSTM) which

replaces the simple computation of the hidden state by a much more complex computation involving several equations.

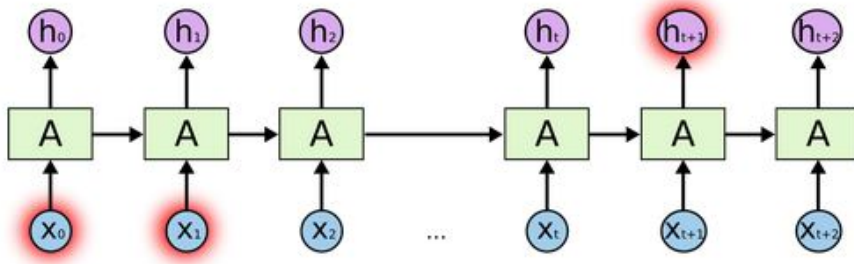


Figure 4: If the distance in time gets large a Vanilla RNN fails to learn the dependency on inputs that are many time steps away due to the Vanishing Gradient Problem.

Source: Olah (2015b)

3.2 Solutions to Optimization Difficulties

As already mentioned Gradient Clipping provides a simple and often used trick to prevent the Explosion of Gradients. Therefore the norm of the Gradient is clipped if it exceeds a certain threshold.

Another solution is to use a special regularization term, which ensures that information can flow over long time ranges (Goodfellow et al. (2016)). The gradient of the loss with respect to a hidden state h_t should be of similar size to the product of this gradient with the derivation of the hidden state to the previous hidden state h_{t-1} :

$$\nabla_{h_t} L \approx (\nabla_{h_t} L) \frac{\partial h_t}{\partial h_{t-1}} \quad (5)$$

Therefore the following regularization term can be added to the loss function:

$$\Omega = \sum_t \left(\frac{\|(\nabla_{h_t} L) \frac{\partial h_t}{\partial h_{t-1}}\|}{\|\nabla_{h_t} L\|} - 1 \right)^2 \quad (6)$$

Another idea is to add skip connections through time, which enable the gradient to flow back longer time ranges and learn dependencies from further back in time.

4 Gated RNNs

4.1 Long-Short-Term-Memory (LSTM)

The LSTM provides a different way of dealing with vanishing gradients. The simple architecture described in Section 2 is now modified. A cell state s_t is introduced, which can be manipulated by different gates to forget old information, add new information and read information out of it. Each gate is a vector of the same size as the cell state and each element of the vector is a number between 0 and 1, with 0 meaning "let nothing through" and 1 "let everything through". The gates are computed as a function of the previous hidden state h_{t-1} and the input at the current time step x_t multiplied by gate-specific weights and squashed through a Sigmoid function into the range of $[0, 1]$. The cell state allows the recurrent neural network to keep information over long time ranges and therefore overcome the vanishing gradient problem.

The equations for the three gates are

$$\begin{aligned} f_t &= \sigma(b^f + U^f x_t + W^f h_{t-1}) \\ i_t &= \sigma(b^i + U^i x_t + W^i h_{t-1}) \\ q_t &= \sigma(b^q + U^q x_t + W^q h_{t-1}). \end{aligned} \quad (7)$$

The cell state s_t is computed as a function of the previous cell state s_{t-1} multiplied by the forget gate f_t , which tells us how much of the old information is discarded. New information is added as a function of the input x_t and the previous hidden state h_{t-1} scaled by the input gate i_t . The hidden state of the network is then a function of the LSTM cell state multiplied by the output gate.

$$\begin{aligned} s_t &= f_t s_{t-1} + i_t \sigma(b + U x_t + W h_{t-1}) \\ h_t &= \tanh(s_t) q_t \end{aligned} \quad (8)$$

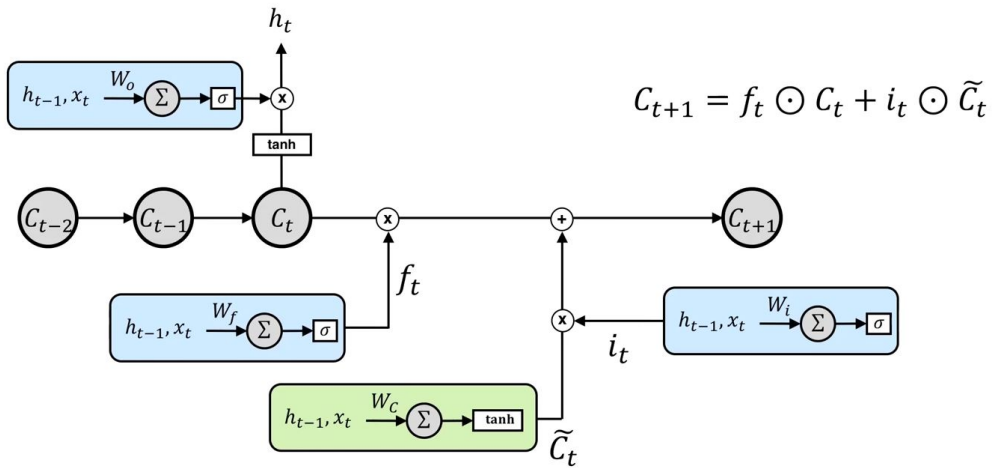


Figure 5: The LSTM architecture is used to model long-term-dependencies.

Source: Nervana (2016)

There are lots of variants of the described LSTM structure, including connections from the cell state to the gates or combining different gates into one, e.g. the forget gate and input gate into one update gate.

4.2 Gated Recurrent Unit (GRU)

The Gated Recurrent Unit (GRU) is an often used variant of the LSTM with a simpler structure. Instead of a forget, input and output gate there are only an update gate and a reset gate. There is no cell state s_t , instead the hidden state h_t is directly manipulated.

$$h_t = u_t h_{t-1} + (1 - u_t) \sigma(b + U x_{t-1} + W r_t h_{t-1}) \quad (9)$$

LSTM and GRU architectures achieve similar performance on many tasks, see Goodfellow et al. (2016).

5 Application 1: LSTM

5.1 Sampling from an RNN

As a generative model an RNN allows us to sample from the conditional distributions at each time step. This is especially useful, if we want to generate sequences ourselves, e.g. texts, programming code, handwriting or speech.

A specialty of sequential data is to determine the length of the sequence. A network without an input parameter determining the length of the sequence would perform poorly as the sequence could stop in the middle of a sentence. There are several ways to control the sequence length:

- special end symbol
- Bernoulli random variable
- integer value giving the length

If we want to generate sequences from a dictionary we could create a special end token, which ends the sequence when its sampled. Another more general way would be to sample an additional Bernoulli output at each time step, determining whether to continue or stop the sequence. A third way would be to sample an integer value determining the length of the sequence. Therefore we need to give the network an additional input at each time step regarding of how close it is to the end of episode. This way we can avoid sequences that stop in the middle of a sentence.

6 Extensions of Recurrent Neural Network

6.1 Recurrent Connections from output to hidden units

Recurrent Neural Networks that have only recurrent connections between the output at one time step to the hidden state at the next time step are a special kind of RNN. They are less powerful because the only information carried on from the past is the actual output of the previous time step, which is trained to match the training labels and therefore cannot carry arbitrary information about the whole past sequence as is possible with recurrence between hidden units. The upside of this is that they are much easier to train as each time step can be decoupled and therefore parallelization speeding up the training is possible.

The training algorithm for such output-to-hidden RNNs is called **Teacher forcing**. Whilst training the true outputs y from the training data are used as an input in the next time step instead of the model-generated outputs, because those are trained to be close to the true outputs. Therefore each time step is independent of another and can be trained in parallel.

At test time the true outputs y are usually not available so the self-generated outputs are fed back into the system. A problem that can arise from this is that inputs at train and test time can vary greatly in size. Possible solutions include training on both teacher-forced and free-running inputs and increase as kind of a curriculum learning strategy the ratio of generated values to real data.

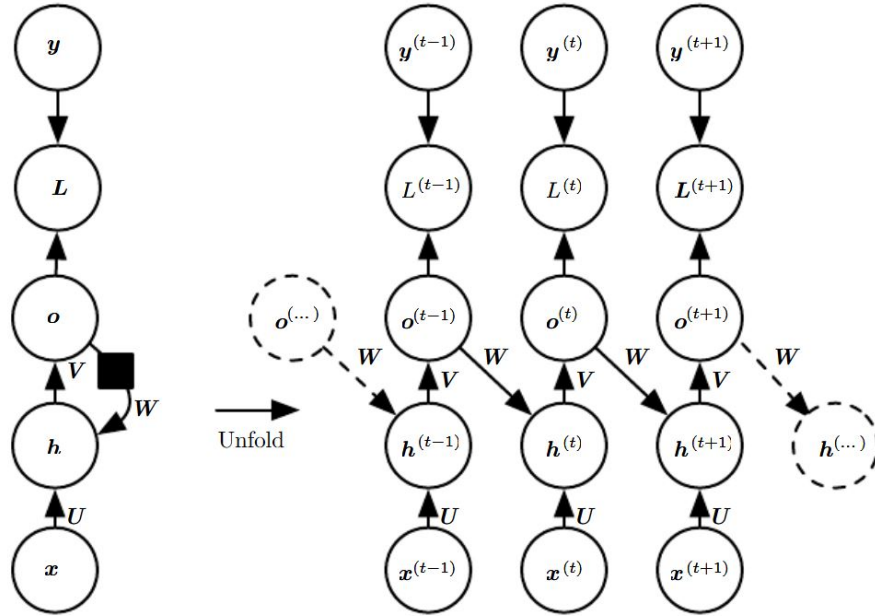


Figure 6: A Recurrent Neural Network with recurrent connections from the output to the hidden units.

Source:

6.2 One-output RNN

Recurrent Neural Networks do not need to have an output at each time step, instead they can only have outputs at a few time steps. A common variant is an RNN with only one output at the end of the sequence. Information from the whole input sequence is incorporated into the final hidden state, which is then used to create an output, e.g. a sentiment ("positive", "neutral" or "negative") for a movie review, see Figure 7. Other applications of such an architecture are sequence labeling, e.g. classify an article into different categories ("sports", "politics" etc.).

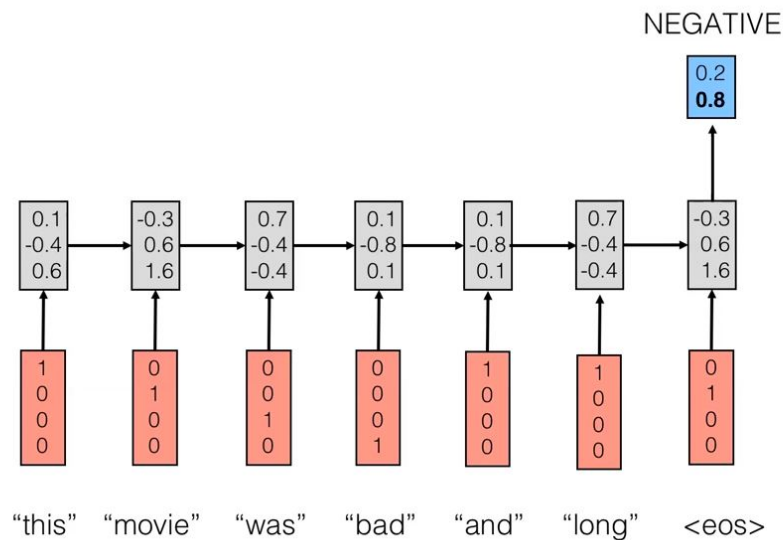


Figure 7: A Recurrent Neural Network with one output at the end of the sequence. Such a model can be used for sentiment analysis.

Source: Nervana (2016)

6.3 RNN based on context

There are multiple ways to include additional inputs that do not change over time in an RNN. The information could be incorporated into the hidden state h_0 or given as a bias term in each time step.

6.4 Deep RNN

RNNs can be made deep in all computations that are performed. One way to do this is to add multiple layers of hidden units in hierarchical groups, another way is to introduce deep computations between hidden to hidden units, input to hidden units or hidden to output units. Instead of a simple nonlinearity as the tanh function, a feed-forward neural network (MLP) is used to compute the next unit. This can be combined with skip connections, which skips the current MLP calculation and enables faster flow of information through the layers. Deep

RNNs are very powerful but have the disadvantage of a vast number of parameters and therefore difficult training.

6.5 Bidirectional RNN

Another generalization of the RNN structure are bidirectional RNNs. Those allow to process sequential data depending on both the past and future inputs, e.g. an application predicting missing words, which probably depend on both the preceding words and the following words. Bidirectional RNNs are also used for the encoding of a sequence in Machine Translation.

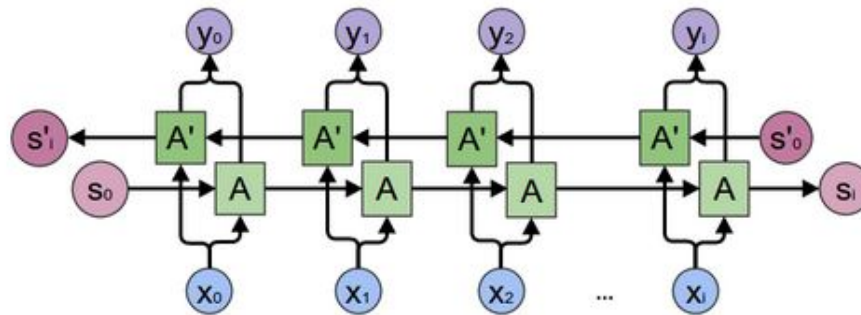


Figure 8: A bidirectional RNN

Source: Olah (2015a)

6.6 Recursive Neural Network

Recursive Neural Networks are a generalization of Recurrent Neural Networks. A tree structure instead of a chain structure is used for the computations of the RNN.

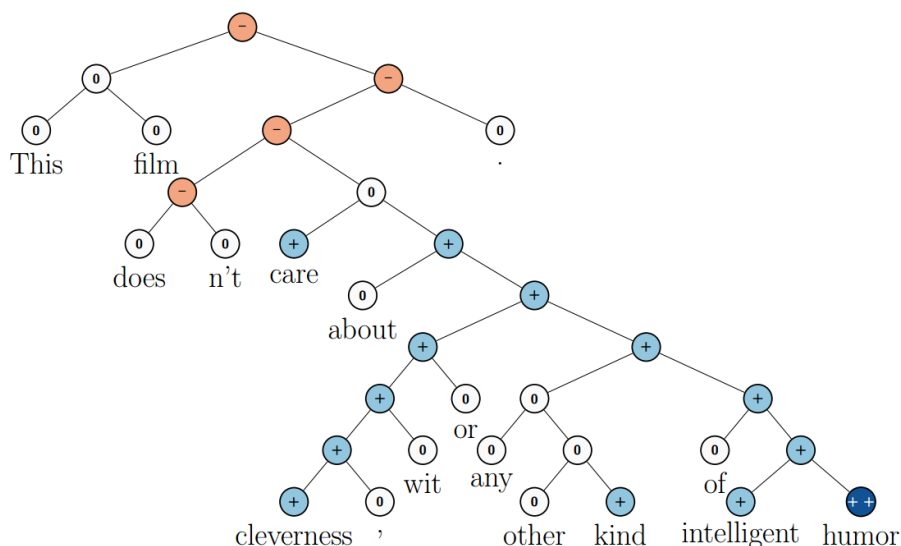


Figure 9: A recursive neural network

Source:

6.7 Encoder-Decoder Network

The encoder-decoder (also called sequence-sequence) framework allows us to generalize the idea of an RNN to different-length input sequences and output sequences. An input/encoder-RNN processes the input of length n_x and omits a typically fixed-length vector C , the context, usually a function of the last hidden state of the encoder-RNN. The output/decoder-RNN processes this context and outputs a sequence of length n_y , which could vary from n_x . Such an architecture can be used for machine translation or text summaries. Another application is image captioning, where the encoder RNN is typically a Convolutional Neural Network, which extracts relevant features from the image, on which the decoder RNN can focus and generate a description.

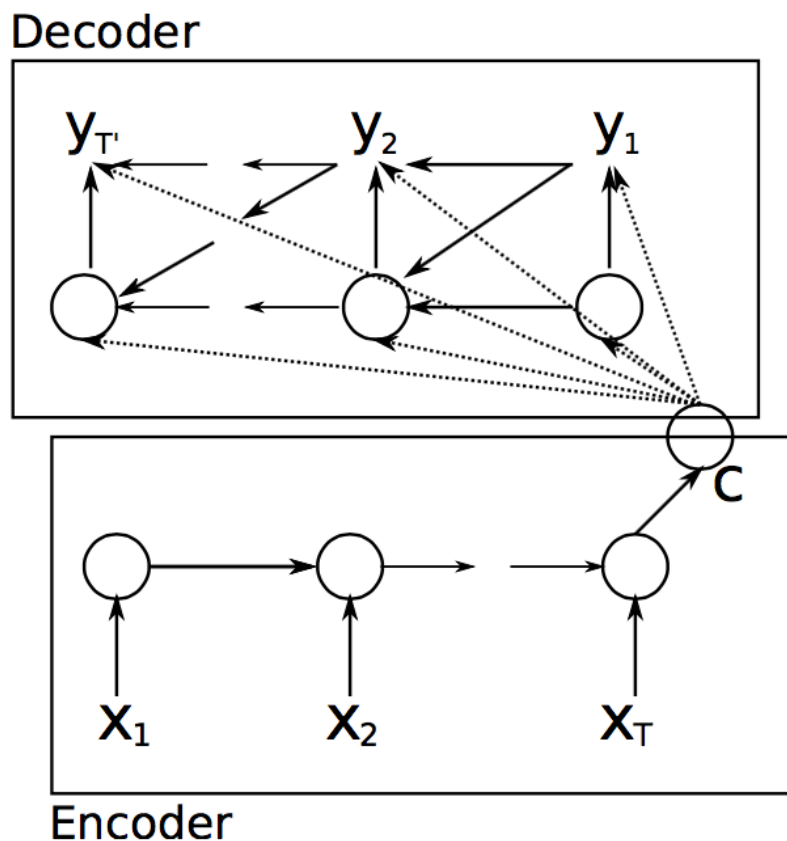


Figure 10: Encoder decoder
Source:

6.8 Attention

Bahdanau et al. (2015) introduced the concept of attention to Neural Machine Translation. In a classical decoder-encoder RNN all information about the input sequence must be incorporated into the final hidden state h_t , which is then passed as an input to the decoder network. With a long input sequence the fixed-sized hidden state is unlikely to capture all relevant information about the past and therefore learning is poor. An attention mechanism allows the decoder

network now to focus on different parts of the input sequence by adding connections from all hidden states to the output RNN. The attention mechanism can be seen as a set of weights, which allow the network to focus on relevant information from the input sequence

Luong et al. (2015) distinguish between a global and local attention mechanism. The global attention model exhibits connections to all hidden states h_1 to h_t of the input RNN, the attention distribution is therefore a variable-length vector of weights. The context vector c_t is then a weighted average of

$$c_t = \sum_t a_t h_t. \quad (10)$$

In comparison a local attention mechanism only focuses on a subset of the hidden states centered around an alignment position p_t to predict the current target y_t . The context c_t is a weighted average of the hidden states, which are inside the alignment window.

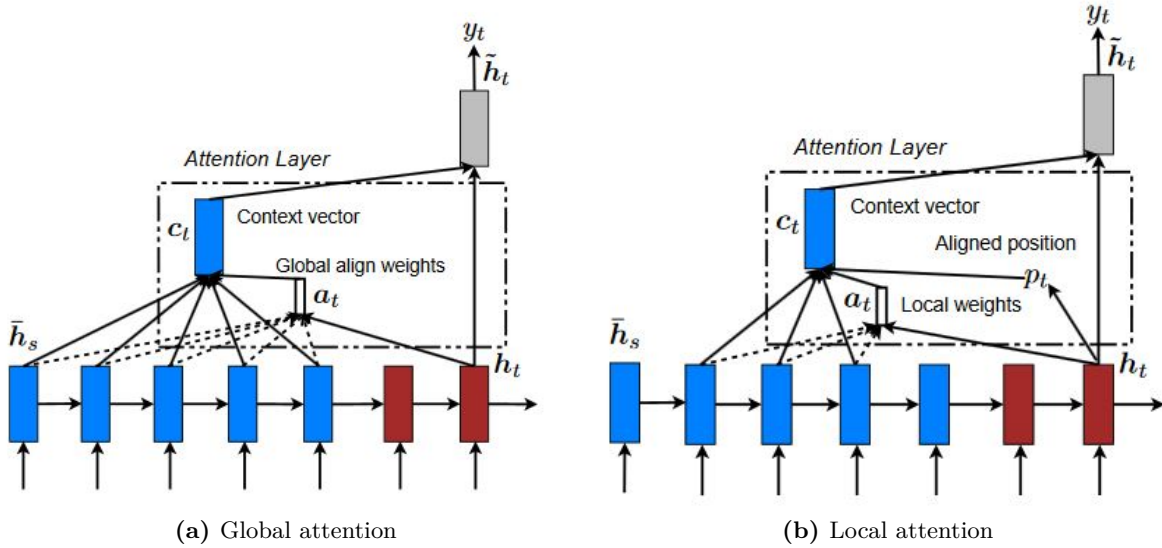


Figure 11: Attention mechanisms.

Source: Luong et al. (2015)

Attention works very well for different tasks such as Machine Translation or Image Captioning by allowing the model to look at relevant parts of the input, e.g. a specific part of the image, for which then a text description is generated.

6.9 Residual Network

Residual Networks (called Res-Nets) He et al. (2015) are useful when training deep neural network architectures. It is similar to dropout, which samples at each time step if a unit is included in the model and creating an ensemble of models. Instead of crossing out single units residual networks operate on whole layers. This can be seen as skip connections through time.

7 Application 2: Google's Neural Machine Translation

For decades Machine Translation of texts was build upon phrase-based concepts, which translated words and parts of sentences by looking its translation up in a dictionary. With the renaissance of neural networks and lots of promising results in diverse applications like image recognition and other classification tasks it was unsurprising that deep learning could be used to improve performance on translation tasks. Recurrent Neural Networks build the main tool to do so. They have been applied by ... with promising results. So it was only a matter of time till a large AI oriented company like Google would use neural network architectures to improve upon their translation service Google Translate. In September 2016 and November 2016 the Google Research team published two papers, in which a new Recurrent neural network architecture for machine translation was presented, which achieves state-of-the-art results on benchmark data sets. In the following we will have a in-depth look at these architecture.

Google's Neural Machine Translation (Wu et al. (2016)) is an application of a Recurrent Neural Network to the task of machine translation achieving state-of-the-art performance.

The architecture consists of an encoder-decoder architecture each build of several stacked LSTMs. In the encoder the first two hidden layers are processed by a bidirectional RNN, thereafter LSTMs with forward flow of information are stacked. For training this deep architecture residual connections between layers are used. An attention mechanism allows the decoder RNN to focus on different parts of the hidden states from the encoder RNN. In the decoder the translation of the input sequence is done with multiple stacked LSTMs (8 layers) and residual connections starting in the second layer. In the final layer a softmax is applied to get a probability distribution over words.

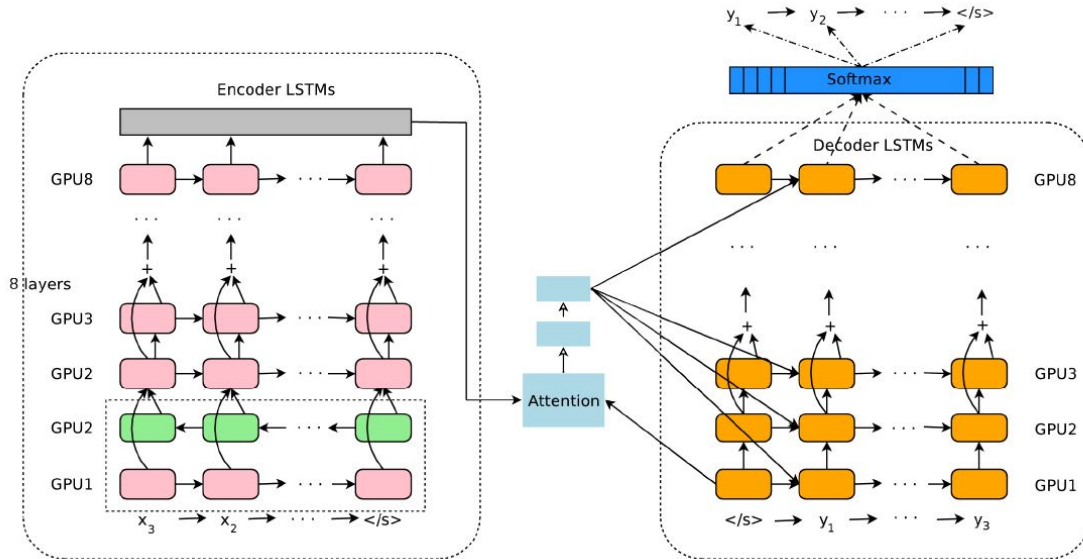


Figure 12: Neural Machine Translation
Source: Wu et al. (2016)

In November 2016 Google announced that the Neural Machine Translation is now capable of zero-shot machine translation, which is translating between language pairs which were not specifically trained. The model often finds a vector representation of sentences (called "interlingua"), where in the meaning similar sentences are clustered together across different languages.

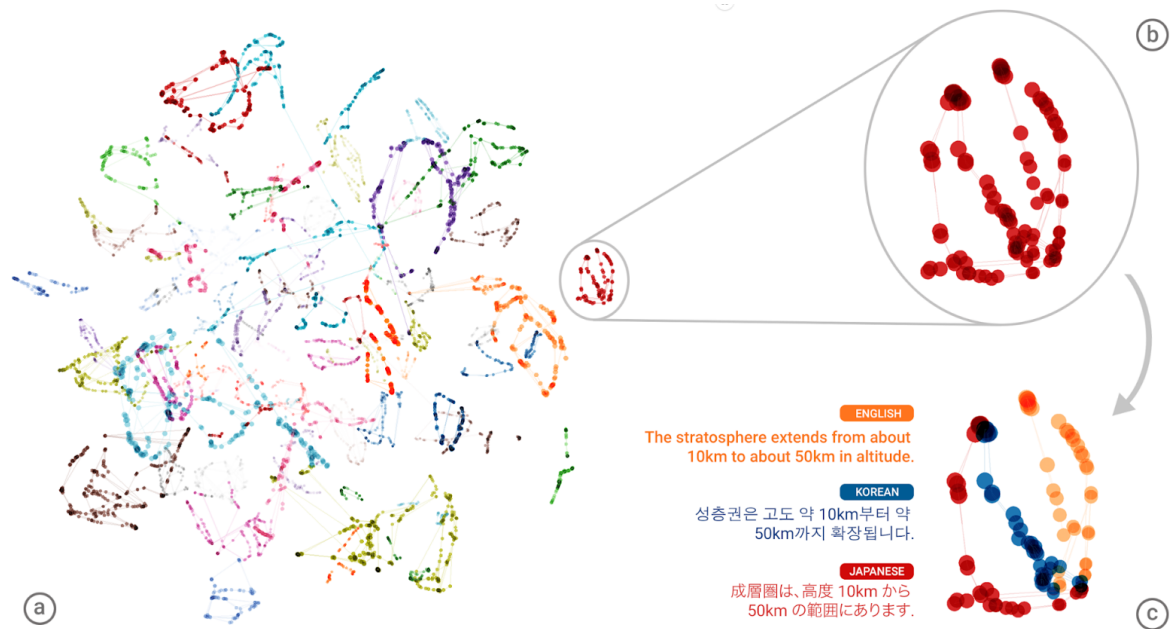


Figure 13: Neural Machine Translation language embeddings

Source: Johnson et al. (2016)

References

- Goodfellow, I., Bengio, Y. and Courville, A. (2016). Deep learning. Book in preparation for MIT Press.
URL: <http://www.deeplearningbook.org>
- He, K., Zhang, X., Ren, S. and Sun, J. (2015). Deep residual learning for image recognition, *CoRR* **abs/1512.03385**.
URL: <http://arxiv.org/abs/1512.03385>
- Johnson, M., Schuster, M., Le, Q. V., Krikun, M., Wu, Y., Chen, Z., Thorat, N., Viégas, F. B., Wattenberg, M., Corrado, G., Hughes, M. and Dean, J. (2016). Google’s multilingual neural machine translation system: Enabling zero-shot translation, *CoRR* **abs/1611.04558**.
URL: <http://arxiv.org/abs/1611.04558>
- Luong, M., Pham, H. and Manning, C. D. (2015). Effective approaches to attention-based neural machine translation, *CoRR* **abs/1508.04025**.
URL: <http://arxiv.org/abs/1508.04025>
- Nervana (2016). Understanding lstm networks, <https://www.nervanasys.com/recurrent-neural-networks>.
- Olah, C. (2015a). Neural networks, types, and functional programming, <http://colah.github.io/posts/2015-09-NN-Types-FP/>.
- Olah, C. (2015b). Understanding lstm networks, <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>.
- Wu, Y., Schuster, M., Chen, Z., Le, Q. V., Norouzi, M., Macherey, W., Krikun, M., Cao, Y., Gao, Q., Macherey, K., Klingner, J., Shah, A., Johnson, M., Liu, X., Kaiser, L., Gouws, S., Kato, Y., Kudo, T., Kazawa, H., Stevens, K., Kurian, G., Patil, N., Wang, W., Young, C., Smith, J., Riesa, J., Rudnick, A., Vinyals, O., Corrado, G., Hughes, M. and Dean, J. (2016). Google’s neural machine translation system: Bridging the gap between human and machine translation, *CoRR* **abs/1609.08144**.
URL: <http://arxiv.org/abs/1609.08144>

List of Figures

1	An RNN can be represented as a dynamical system or unrolled over time.	4
2	A Vanilla Recurrent Neural Network unrolled over time.	5
3	A simple application of a vanilla RNN for language modeling.	5
4	If the distance in time gets large a Vanilla RNN fails to learn the dependency on inputs that are many time steps away due to the Vanishing Gradient Problem. .	7
5	The LSTM architecture is used to model long-term-dependencies.	8
6	A Recurrent Neural Network with recurrent connections from the output to the hidden units.	10
7	A Recurrent Neural Network with one output at the end of the sequence. Such a model can be used for sentiment analysis.	11
8	A bidirectional RNN	12
9	A recursive neural network	12
10	Encoder decoder	13
11	Global attention mechanism	14
12	Local attention mechanism	14
13	Attention mechanisms.	14
14	Neural Machine Translation	16
15	Neural Machine Translation language embeddings	17