

# Machine translation with Recurrent Neural Networks

Vera Neplenbroek

Kamiel Fokkink

Baran canl

veraneplenbroek@icloud.com kamielfokkink@gmail.com barantevitol@gmail.com

## Abstract

The goal of this project is to create a functioning English-Dutch translation model. The data set contains English and Dutch translations of the same movie subtitles, which will be connected in the translation model to build the translator. The model consists of two Recurrent Neural Networks, an Encoder and a Decoder, which will transform the English input sentence to a tensor, and then to a Dutch output sentence. It is based on a PyTorch implementation, but extended by reversing the input sentences, and using different data formats. Evaluation will be done using the BLEU measure, giving each translated Dutch sentence/text a score out of 1, based on how much it corresponds to the target sentence/text.

## 1 Introduction

Translating one word from one language to another can be done easily using traditional paper dictionaries. However, knowing the individual translation of each word is not enough to understand sentences or documents in the other language. The difficulty comes from language being compositional, thus the combining of words builds up meaning in the sentence. The final meaning of a sentence can not be reduced to individual meanings of the words, it also comes from their combinations. This is where traditional dictionaries fail, and machine translation comes in.

There are several approaches to translating a piece of text. In recent years, the statistical approach, based on observations of large

amounts of text and generating patterns from it, has seen great improvements, thanks to developments in the field of machine learning. For our project, the particular technique that we decided to use was Sequence to Sequence translation. This technique was pioneered in 2013, by Kalchbrenner and Blunsom<sup>1</sup>, and makes use of recurrent neural networks (RNNs). Our aim is to produce a translation model that is based on this architecture. It will be trained on English-Dutch sentence pairs, and take in a sentence in English, and return its Dutch translation.

For the training of our model we used sentence pairs from movie subtitles. The dataset was retrieved from Open Subtitles<sup>2</sup>. The data consists of translation pairs coupled together, so for each line uttered in the movie there is the English and the equivalent Dutch subtitle. An important advantage of using subtitles is that the text is well aligned, as compared to more freely translated datasets, such as books, articles or blogposts. For those sets the individual sentence pairs would fit less well to each other, making it much more difficult for the model to train on.

## 2 Preprocessing

The dataset we work with came in the TMX format. It was a dataset of 5.6 GB size. Thus, in order to work with it, we needed to divide the dataset into chunks or use samples we extract from it. Furthermore, we decided that while the TMX format might have its benefits, we did not need most of them. Therefore, we converted the dataset into a list of sentence

---

<sup>1</sup>N. Kalchbrenner, Phil Blunsom 2013. Recurrent Continuous Translation Models Association for Computational Linguistics

<sup>2</sup>Open Subtitles

pairs, and stored them within a pickle. Finally, we realized that our model required some normalization on each sentence via punctuation removal, lowercasing each word, and riddance from sentences beyond the maximum length.

## 2.1 Format Handling and Dividing the Dataset

As mentioned above, our dataset was too large to read all of it at once. Therefore, we decided to use basic file-streaming to read, process and save chunks of it. After deciding on an approximate size for each chunk, which was roughly 500 MB's, we created a pipeline. The pipeline took a chunk, cleared the irrelevant tags off, grabbed each English-Dutch sentence pair with a regex. It treated information left in two different chunks. After pairing them, it put the pair in a list, and wrote the list in a pickle file. This way, we acquired clean sentence pairs in different chunks, which we can read one by one.

## 2.2 Sampling

The developing process of models required small samples from the dataset to see the efficiency and quality of the model in a rough manner, before a proper evaluation process. However, just taking a chunk of the dataset would result with a high variation of low amounts of data. Therefore, we created a sampling function, which uses a preset of rules and extracts sentence pairs fitting that rule, then pickles them in another file for usage. This way, the accuracy of the model can be observed faster without an evaluation model on early stages. Also, we cleared the sentences off of some punctuation, transformed uppercase letters to lowercase and eliminated sentences beyond certain word counts. This was done to fit our dataset better for our model and remove noise.

## 3 The RNN Encoder-Decoder model

While regular neural networks are a powerful tool to train a model, they have one important limitation: they only take inputs of fixed length, and produce outputs of fixed lengths.

For the translation of a sequence of words from English to Dutch this is not very useful, because both the length of the input and output sequence will vary between different sentences. We want to train a model that can take in any sentence of arbitrary length, and produce a corresponding translation. The solution to train this kind of model is provided by Recurrent Neural Networks(RNN)<sup>3 4 5</sup>

### 3.1 Encoder-Decoder

The key element to take in a sequence of any length, and produce any desired output is the hidden state. To work with an input sequence  $\mathbf{x} = (x_1, x_2, \dots, x_t)$ , the RNN looks at the elements of the input step by step. In our case, the input is a sentence, and the elements are words. Beginning at the first word, the RNN creates a hidden state, and an output vector  $y$ . For each next word, the RNN looks at that word and the previous hidden state, and updates to a new hidden state by

$$\mathbf{h}_{(t)} = f(\mathbf{h}_{(t-1)}, x_t),$$

where  $f$  can be any non-linear activation function. After the final word  $x_t$  has been evaluated, we have the final output vector  $y$ , also called context vector. This is the first half of the model, or the Encoder. In this step, the meaning of the English input sentence gets encoded into a vector, thus we have a numerical representation of the input.

In the second step, we decode the context vector into a sentence in Dutch, again by using an RNN. The approach is very similar, just the other way round. The initial hidden state of the decoder is the context vector, and its first input is the `¡SOS!` (start of sentence) token. At this first step, the decoder also creates its own hidden state. From there, it generates each new word of the output sequence, by looking at its previous output, the context

<sup>3</sup>D. Bahdanau, K. Cho, Y. Bengio. 2014. Neural machine translation by jointly learning to align and translate. Retrieved from the arXiv database.

<sup>4</sup>I. Sutskever, O. Vinyals, Q. Le. 2014. Sequence to sequence learning with neural networks. Retrieved from the arXiv database.

<sup>5</sup>K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. Retrieved from the arXiv database.

vector, and its own hidden state. Thus the activation function for the decoder looks a bit different:

$$\mathbf{h}_t = f(\mathbf{h}_{(t-1)}, y_{t-1}, \mathbf{c})$$

Finally, it will output a generated sentence in Dutch.

### 3.2 Training

During the training of the model, both these RNNs are jointly trained via gradient-descent. The goal of the training is to maximize the probability that our model assigns to each Dutch sentence, given each English sentence.

$$\max \frac{1}{N} \sum_{n=1}^N \log p_{\theta}(y_n | x_n)$$

The sentences  $y_n$  and  $x_n$  are those that we have from the data. The probability depends on the parameters theta of our model, and is logarithmic for easier summation and to stay within computational bounds.

We used an implementation from a PyTorch tutorial <sup>6</sup> as the basis for our model. We looked into adapting the code to our needs, have it fit to our dataset, and interacted with the code in various places to implement some tweaks to the model, as discussed in the next section.

### 3.3 Changes to the model

When implementing the changes to the model we used a small sample from the dataset with sentences that are similar in grammar. We looked at the BLEU score computed after evaluating the trained model, which was compared to the one for the base model: 0.01850. After experimenting with multiple changes, we have decided to implement the attention mechanism, teacher forcing, and reversing the input sentences, due to their higher scores. The results for all changes are below.

#### Attention, Score: 0.02501

When using a regular decoder, the output vector from the encoder is all the decoder has to learn from and to base its translation on. This means the full input sentence and its meaning is now represented by one vector only.

<sup>6</sup>PyTorch, 2017

When using a decoder with attention mechanism, however, for each word that the decoder returns, it is provided with a set of attention weights. These weights tell the decoder which parts of the input vector are important to translate the word. This set of attention weights will be multiplied with the output vector from the encoder, to create a vector of weighted inputs for the decoder. For shorter sentences not all weights will be used <sup>7</sup>.

#### Teacher Forcing, Score: 0.04614

Teacher forcing is used in training the model, where instead of the last output from the decoder, the target tensor is fed to the decoder. Teacher forcing helps the model to learn coherent grammar, but if overexploited, the finalized model will have difficulty creating a grammatical structure by itself. Since the model is used to getting the first few words provided by the teacher, when it does not get these in the evaluation phase it does not know how to starting creating the output sentence<sup>8 9</sup>. We chose to implement teacher forcing with a ratio of 0.5, to prevent it's overuse.

#### Reversing the input, Score: 0.02149

Reversing the input (source) sentences when feeding them to the model to train has proven<sup>10</sup> to decrease the distance between words in the source and in the target language. In other words, long-term dependencies are being replaced by short-term dependencies. This results in a higher probability of two words, each from a different language, which context is depended on each other to be close in distance.

#### Activation functions, Score: 0.0

At first we tried to implement the log sigmoid and leaky relu activation functions, but in the original log softmax activation function a dimension was specified, which could

<sup>7</sup>Bahdanau, 2014

<sup>8</sup>PyTorch, 2017

<sup>9</sup> H. Jaeger. 2002. A tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the echo state network approach

<sup>10</sup>Sutskever, 2014

not be specified for the other activation functions. Since the model would not train using those functions, we tried using activation functions for which the dimension could be specified, namely the softmax and softmin activation functions. For these functions, however, regardless of the number of iterations, lengthy sentences filled with infrequent words were created. We believe this is due to the way the model was built, and that trying different activation functions is justified from a theoretical viewpoint <sup>11</sup>.

## Different NN units

We considered changing the linear units that make up the RNN in the Pytorch implementation into RNN or LSTM units, but since Pytorch created a custom RNN from the linear unit specific for this implementation, it did not work and did not make sense to change the linear unit.

## 4 Evaluation

To evaluate the performances of our model we chose to write our own implementation of the widely used BLEU score. The BLEU score takes all the n-grams up until the specified n and calculates the amount of similarities between two sentences, where for each n-gram each word can only be used once to count towards a match with the other sentence. To calculate the BLEU score over a set of sentences, the geometric mean of the scores for each individual sentence is used <sup>12</sup>.

## 5 Discussion

One of the limitations of our model is that the sentence pairs we used to train the model on were not that well aligned. While for each pair in the data both sentences expressed the same line in the movie, there are various cases where they were phrased in a different manner. This could have affected the training, and eventually lower its performance. Out of 50 training

pairs, 9 were counted to be poorly-aligned in this manner. It is difficult to estimate a ratio for the whole data, but at least a significant portion of the pairs are sub-optimal. Filtering these pairs out would require human intervention on a major scale, which is beyond the scope of this project. Hence, we just trained on the data that we had.

There were a few changes that we tried to implement into our model, but that did not fit with the rest of the model. One of them was to try and write the attention mechanism ourselves. While during this course we have focused on the probabilistic approach towards NLP, there is also the symbolic approach, by writing top-down rules. While this sounded pleasant in theory, creating such a weight system which is better than the existing one proved to be too demanding, due to the necessary training trials needed to test the results. Hence, we decided to not implement it.

Finally, because of scale issues, we did not train over the whole dataset. Instead, we filtered sentence pairs on a maximum length of 15 words per sentence, starting with a pronoun, and followed by a form of to be or to have. We split this filtered dataset into training and test set in a 80/20 ratio. The total training time was 3 days. After training, our final model still had not learned a large vocabulary. However, learning this would just be a matter of time and computing power. The printout in Figure 1 shows three examples of the target sentence, and the translation that our model came up with.

## 6 Conclusion

When evaluating the translations that our model gave for the test set, we got the following BLEU scores between 0 and 1 for the various n-grams. 4-grams: 0.03, 3-grams: 0.09, 2-grams: 0.17, 1-grams: 0.29. So a considerable portion of the words gets correctly translated. We can conclude that the basis of our model works, albeit with a limited vocabulary, and is able to produce a somewhat meaningful translated sentence.

<sup>11</sup>PyTorch, 2017

<sup>12</sup>K. Papineni, S. Roukos, T. Ward, W. Zhu. 2002 BLEU: a Method for Automatic Evaluation of Machine Translation Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)

---

```

['we', 'moesten', 'immers', 'precies', 'weten', 'in', 'wel
ke', 'richting', 'we', 'onze', 'waarnemingen', 'deden']
['we', 'moesten', 'het', 'dat', 'we', 'in', 'de', 'het',
'<EOS>']
['zij', 'waren', 'als', 'vampier', 'vleermuizen'] ['zij',
'waren', 'als', 'de', '<EOS>']
['ik', 'heb', 'niets', 'meer', 'te', 'verbergen'] ['ik',
'heb', 'niets', 'niets', 'te', '<EOS>']

```

---

Figure 1: Sample from the output

## References

- [Bahdanau et al.2014] D. Bahdanau, K. Cho, Y. Bengio. 2014. *Neural machine translation by jointly learning to align and translate*. Retrieved from the arXiv database.
- [Cho et al.2014] K. Cho, B. van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio. 2014. *Learning phrase representations using RNN encoder-decoder for statistical machine translation*. Retrieved from the arXiv database.
- [Jaeger2002] H. Jaeger. 2002. *A tutorial on training recurrent neural networks, covering BPPT, RTRL, EKF and the "echo state network" approach*
- [Kalchbrenner et al.2013] N. Kalchbrenner, Phil Blunsom. 2013. *Recurrent Continuous Translation Models* Association for Computational Linguistics
- [Papineni et al.2002] K. Papineni, S. Roukos, T. Ward, W. Zhu. 2002. *BLEU: a Method for Automatic Evaluation of Machine Translation* Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)
- [PyTorch2017] *NLP From Scratch: Translation with a Sequence to Sequence Network and Attention*
- [Sutskever et al.2014] I. Sutskever, O. Vinyals, Q. Le. 2014. *Sequence to sequence learning with neural networks*. Retrieved from the arXiv database.