

7. Assignment in “Machine Learning for Natural Language Processing”

Summer Term 2021

1 General Questions

1. What are the differences between a recurrent neural network that processes inputs of length $l = 3$ and a multilayer perceptron (fully-connected network) that has three layers?
2. Describe two tasks that are better suited for RNN models than for MLP models and provide a short explanation.

? Something to think about

3. Two problems in Machine Translation that still remain with biRNN Seq2Seq models are a) unknown words (e.g., names) and b) forgetting/duplicating parts of the input. Think about possible solutions to these problems!

? Something to think about

4. Propose a neural network architecture that would be suited for generating image descriptions, that is, given an image as input, creates a sentence describing the content of the image! Provide reasoning why your architecture is a good choice for this task.

2 Beam Search

In the lecture, you learned that Beam Search can sometimes lead to better translations than simple Greedy Search, where the decoder always chooses the locally most likely token.

To show that this can indeed happen, look at the following situation: Given an Encoder-Decoder network trained to translate from English to German. The encoder has read

the sentence "I won't tell you nothing!" and produced an internal representation h that encodes this sentence. This representation is now fed to the decoder and a translation is generated. The decoder will return the following probabilities for output tokens:

$$\begin{aligned}
 P(\text{Ich}|\langle s \rangle) &= 0.9 \\
 P(\text{Er}|\langle s \rangle) &= 0.01 \\
 P(\text{verrate}|\langle s \rangle, \text{Ich}) &= 0.3 \\
 P(\text{verrate}|\langle s \rangle, \text{Er}) &= 0.05 \\
 P(\text{sage}|\langle s \rangle, \text{Ich}) &= 0.5 \\
 P(\text{sage}|\langle s \rangle, \text{Er}) &= 0.08 \\
 P(\text{euch}|\dots, \text{verrate}) &= P(\text{euch}|\dots, \text{sage}) = 0.3 \\
 P(\text{dir}|\dots, \text{verrate}) &= P(\text{dir}|\dots, \text{sage}) = 0.45 \\
 P(\text{nichts}|\dots, \text{euch}) &= P(\text{nichts}|\dots, \text{dir}) = 0.35 \\
 P(\text{nicht}|\dots, \text{euch}) &= P(\text{nicht}|\dots, \text{dir}) = 0.4 \\
 P(!|\dots, \text{nichts}) &= 0.8 \\
 P(!|\dots, \text{nicht}) &= 0.1 \\
 P(\text{nichts}|\dots, \text{nichts}) &= 0.01 \\
 P(\text{nichts}|\dots, \text{nicht}) &= 0.2
 \end{aligned}$$

1. Perform a greedy search to get the output of the decoder network!
2. Perform a beam search with $B = 2$ to get the output of the decoder network!
3. Compare the two outputs.

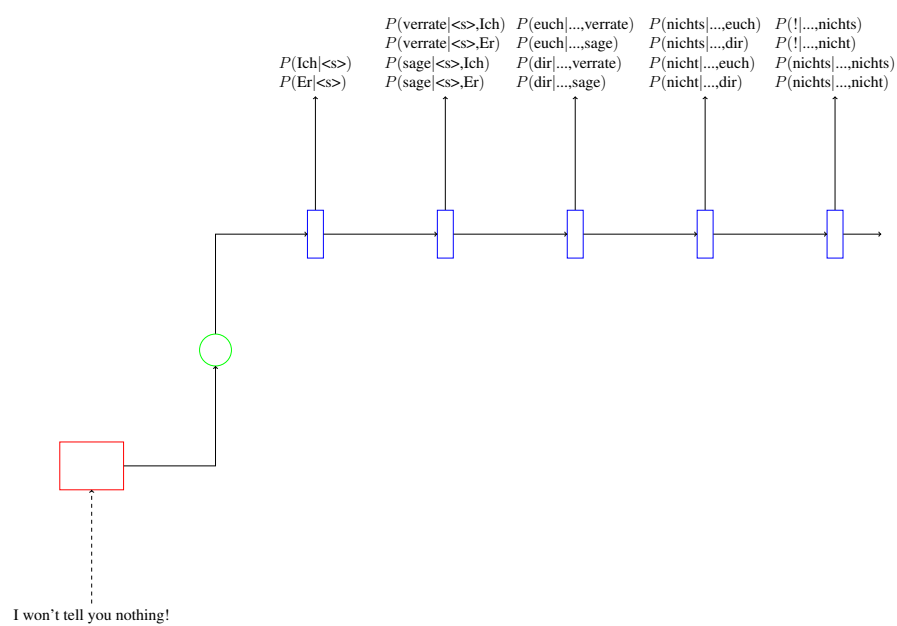


Figure 1: Visualisation of the Encoder-Decoder network used in this task