

# HW3: Language Translation

Jonah Philion  
jonahphilion@college.harvard.edu  
github

March 4, 2018

## 1 Introduction

In this problem set, we attempt to translate German to English. The two models take the form of

$$y_{i+1} \sim \sigma(W\phi(y_{:i+1}, \mathbf{x}) + b)$$

where  $\mathbf{x}$  is an input phrase in German and  $\mathbf{y}$  is the translated phrase in English.

- Encoder-Decoder with no attention
- Encoder-Decoder with attention

## 2 Problem Description

Given a sentence in German of up to 20 words, we would like to predict a sentence in English with the same meaning. For all models, a sentence  $\mathbf{x}_i$  is encoded as a sequence  $x_1, \dots, x_n$  where each  $x_j$  is a one-hot vector of length the vocabulary  $\mathcal{V}$ . The model outputs a categorical distribution over the vocabulary  $\mathcal{V}$ . Embeddings  $\mathcal{E}_d$  map a one hot vector  $x_j$  to a dense vector of size  $d$ . Each language gets its own word embedding.

## 3 Model and Algorithms

All models are trained on the IWLST. For models requiring gradient descent, training loss and validation loss are recorded in real time with visdom. Reported PPL is on the validation set with teaching.

### 3.1 Evaluation

1. All models assume a vocabulary of 10000 and have a predict function which given a batch of sentence fragments outputs an array containing the probabilities for the next word for all batches.
2. All models are tested with the same evaluation function. The function evaluates the model on batches of size 32 with the same metric used during training.

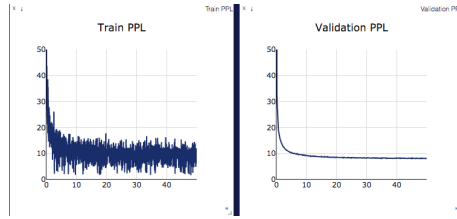


Figure 1: The network just barely reaches a PPL of 8.

## 4 Experiments

Both models are trained with Adam Optimizer .001 learning rate, weight\_decay .0001 for 10-20 epochs

### 4.1 No Attention

To encode, we reverse the input sentence, then embed the sentence in a space of dimension 200. An LSTM of one direction, 4 layers deep, with initialized state and hidden layers to 0, reads through the sentence and we store the final hidden layers and state, as well the history of the hidden layer throughout the sentence.

To decode, we embed the ground truth  $trg$ , run a LSTM with initial hidden state and state given by the output of the encoder. The hidden state  $h_t$  of the decoder LSTM at any  $t$  is projected to  $V$  by a linear layer  $D \rightarrow V_{en}$ . The graph during training is shown in Figure 1.

With beam search of  $B = 5$  the top scoring sentences of the first three sentences are

src: Als ich in meinen 20ern war , hatte ich meine erste Psychotherapie-Patientin .  
 model: And when I was in my <unk> , I 've got my first <unk> . ”  
 Google: When I was in my 20s, I had my first psychotherapy patient.

src: Ich war Doktorandin und studierte Klinische Psychologie in Berkeley .  
 model: So , I was <unk> and <unk> <unk> <unk> in <unk> . <unk> .  
 Google: I was a PhD student and studied Clinical Psychology in Berkeley.

src: Sie war eine 26-jährige Frau namens Alex .  
 model: And it was a <unk> man called <unk> <unk> .  
 Google: She was a 26-year-old woman named Alex.

### 4.2 Attention

The encoder is no different from that of the non-attention model.

For the decoder, the embedding of the source sentence is first padded with zeros to be a length of 20. Next, the most recent hidden state is concatenated to the most recent word vector and fed to a linear layer  $2D \rightarrow 20$ . The softmax of this vector gives the “attention” to the at most 20 words in the source sentence. The input to the decoder is a linear combination of the encoding with coefficients given by the “attention” vector concatenated with the input word vector. As in the no-attention model, a linear layer  $D \rightarrow V_{en}$  gives the output distribution over the vocabulary. I report the score for  $D = 300$ .

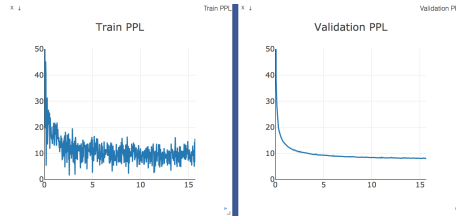


Figure 2: The network performs almost identically to the non attention model which means that the attention component of the decoder was just noise.

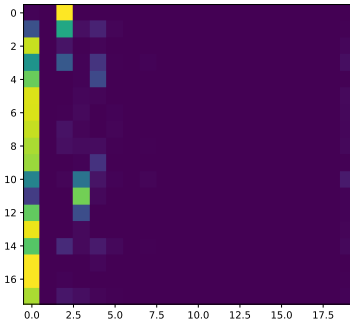


Figure 3: The attention does not look good. The horizontal axis is the maximum sentence length of 20, the vertical axis corresponds to each word generated in the output sequence. The sum of each row is 1.

PPL: 7.97

The training graph is shown in Figure 2. An example translation is below.

src: Als ich in meinen 20ern war , hatte ich meine erste Psychotherapie-Patientin .

model: And when I was <unk> <unk> , I was my first <unk> . I said .

Google: When I was in my 20s, I had my first psychotherapy patient.

src: Ich war Doktorandin und studierte Klinische Psychologie in Berkeley .

model: And it was <unk> <unk> , and I was <unk> in <unk> <unk> .

Google: I was a PhD student and studied Clinical Psychology in Berkeley.

src: Sie war eine 26-jährige Frau namens Alex .

model: She was a woman called <unk> <unk> . He said .

Google: She was a 26-year-old woman named Alex.