

Report - Homework #1

1 Getting Started with PyTorch

One correct prediction of my improved network would be ('unsympathetic', 'ADJ'). It is likely that the n-gram module added was able to recognize the suffix 'tic' as a characteristic of adjective. The pooling network got it wrong as noun, since the mean pooling didn't take the order of characters into consideration so wouldn't be able to identify the suffix 'tic'.

An incorrect prediction would be ('downwind', 'VERB') where 'downwind' should be an adjective. I guess in this example the model was confused by the prefix 'down' because many words with 'down' prefix, like 'downsample', are verbs. The pooling network also got it wrong.

2 Training Language Models

The extensions I tried alongwith my motivations include the following.

- Add a L2 regularization penalty on the activation of the LSTM output with $\lambda = 10^{-3}$, as there was overfitting, the decreasing training loss and the increasing validation loss in my original LSTM model, as shown in figure 1a.
- Add layer normalization to all the linear layers in the network, to enable smoother gradients, faster training, and better generalization accuracy.
- Decrease the learning rate while training: first train the model for 10 epochs with $lr = 10^{-3}$ and then 10 epochs with $lr = 10^{-4}$. This intuitively make sense from ml/optimization courses.

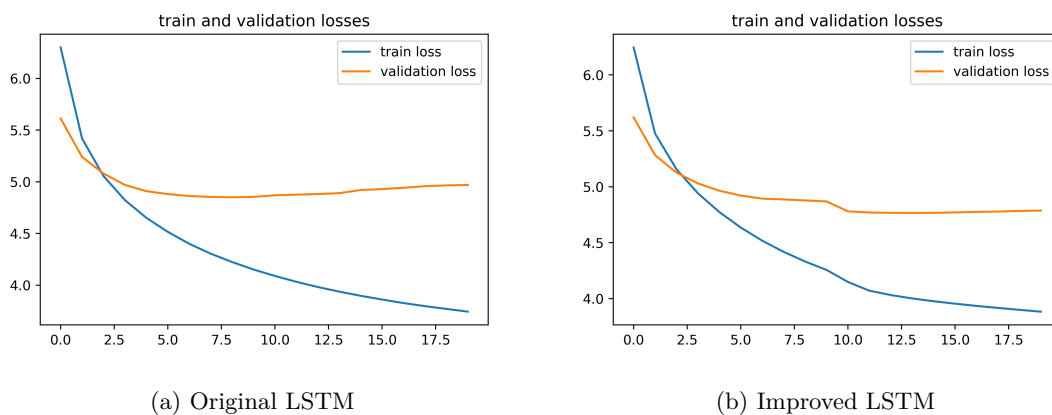


Figure 1: Learning Curves

The extensions were effective. As shown in the figure above, there was no longer overfitting in the improved LSTM after adding activation regularization. We can also see a sudden drop of validation loss after epoch 10 which is exactly when learning rate changed from 10^{-3} to 10^{-4} . The decrease of learning rate actually brought a big boost to our results.

Comparing my improved model to the original LSTM, the validation perplexity dropped from 127.9912383431524 to 117.62607399046803.