# Word/Character Level Language Modeling With RNN

Lefteris Soulas
Computer Science Department
New York University
es3431@nyu.edu

## 1. Question Answers

Q1. The lua script that executes the code is included separately.

Q2. The $i, prev\_h$ variables are used in a function inside the LSTM cell and returned as a linear combination. This function maps to $T_{2n,4n}\begin{pmatrix} h_t^{l-1} \\ h_{t-1}^l \end{pmatrix}$, in the paper. The $prev\_c$ maps to $c_{t-1}^l$, which is the memory cells from the previous time step.

Q3. The create_network() function, returns a single instance of the network which is later cloned many times in order to form the various time steps. The network at this point is folded.

Q4. $model.s$, is a vector which contains the hidden state of the model. $model.ds$, is a vector which contains the gradients from every hidden state and is updated in the backprop. $model.start\_s$, is a helper variable that initializes the very first hidden state of the model ($s[0]$). $model.start\_s$, is reset to 0 in the beginning of the program and in the very end, if the last batch contains less inputs than $batch\_size$, in which case we discard the last mini-batch and retrain on the beginning of the sequence as it was the very end.

Q5. In order to deal with the vanishing/exploding gradients, we perform gradient clipping above a certain threshold ($max\_grad\_norm$). If the norm of the gradients is above this threshold, the gradients are normalized so they have maximum value equal to $max\_grad\_norm$.

Q6. The optimization method used, is mini batch gradient descent. The batch length is determined by $batch\_size$. This algorithm doesn't use momentum, but uses a decay factor which decreases the learning rate after a number of epochs.

Q7. We have to zero out the extra output in the back propagation step in order not to affect the gradients.

## 2. Word level language model

For this part of the assignment we created the function *query_sentences.lua* which takes as input a number and some words and tries to complete the sentence in the best possible way. We didn't train a very elaborate model for this part, but we rather tested the functionality on a model similar to the baseline. In order to allow variations in our sequences, we didn't simply select the word with the highest probability, but used a multinomial sample generator to pick the predicted word.

## 3. Character level language model

For this part of the assignment we experimented with the LSTM RNNs in order to model a character based language model. In particular, our implementation is based on Wojciech Zaremba implementation [3].
As a first step we changed the code so it can predict characters instead of words and added an extra output, which is the log probabilities of each word given a certain input. In the beginning we trained the baseline model with no change at all and noted validation perplexity of 430. Following the assignment's instructions we increased the sequence length to 50 and retrained the model. There new validation perplexity was 350.

In order to improve the model even more we tried to tune the optimization parameters. An interesting observation we noticed is that the model's perplexity doesn't decrease if we increase the RNN's size. On the contrary, we tried three different models, with rnn_size =$\{1000, 1500, 2000\}$ and all of them were very slow and very difficult to converge. Eventually we gave up the effort as the perplexity was extremely large even after the $10^{th}$ epoch. On the other hand, we noticed that relatively smaller RNN sizes (around 400-500) perform much better. That's the range we kept for the rnn_size in the rest of our experiments.

During the hyper parameter tuning, we noticed that the initial learning rate of 1 can be too harsh in most cases and the mode's perplexity starts to decrease only after reaching the *max_epoch*, when the learning rate is divided by *decay*.

Thus, we kept an initial learning of 0.5 during the rest of our experiments.

The network has as a first layer a LookupTable module, which is responsible for learning the word/character embeddings. As we know very well, the initialization of the weight parameters is really important for the convergence of the model. Instead of assigning them arbitrarily to some value, we thought of using the vector representations from GloVe[1] or Word2Vec[2], in order to initialize each 200-D parameter vector for each of the 10000 words that exist in our vocabulary. While this could be a meaningful warm start for the word language modeling, there is little that we can do when it comes to character model, as our class space is limited (only 50 characters). Therefore we intuitively initialize the weights to 0.02, which is the uniform probability of each of the 50 words.

In our continuous effort to optimize the hyper parameters we run a few more models, trying to mimic the simple modification that the increase of the sequence length brought to the perplexity of the model. We tried to test the hypothesis that by increasing the sequence length, we are lowering the perplexity (something which intuitively makes sense as the RNN will use more "state" in order to distinguish different sequences of characters).

The best model we found till now achieves training perplexity 163.713 and validation perplexity of 262. We used a sequence length of 100, with batch size of 50 and rnn size of 500. The model was trained for 20 epochs with a decay factor of 2. As you can see in Figure 1.1, the model already started to overfit after 20 epochs. That's logical as we used a dropout rate of 0.

In order to mitigate this phenomenon and be able to scale, we decided to commit and train another model that will take into account all the above lessons we have learned so far. More specifically it will have a relatively small rnn size (we picked 400), a large sequence length (we picked 400) and a dropout rate of 0.6 in order to under-perform on the training set but generalize better. The use of a dropout layer is proven to increase generalization [4].

We trained the above network for 200 epochs, but because we wanted to have absolute control on the learning process, we controlled manually the learning rate after the 100th epoch. The model achieved the perplexity of 187.84 in the validation set and 444.5 in the training set 2. The discrepancy between the two, shows that we were able to learn a well generalized model and achieve a very good perplexity comparing to our other efforts. The total time of tuning the model was 2068 minutes.

The most bizarre/disappointing realization, is that when we run the a4_grading.py script, the perplexity gets to be close to 300. When we use an interactive session to load the model and train with the run_valid() function from the initial code, the perplexity is 187. We couldn't figure out why this bug appears. You can find the core_network for this model inside the folder *forThe187PerplexityModel*, as well as a script that loads it into memory and sets up the model so you can test the validation function yourself.

Moreover, we tried changing the lstm cell functions as the ones in the code weren't exactly the same as the ones presented in [5], but the accuracy of the model didn't improve.

## References

[1] Pennington et al, *GloVe: Global Vectors for Word Representation*

[2] Mikolov et al, *Efficient Estimation of Word Representations in Vector Space*

[3] Wojciech Zaremba, *Long Short Term Memory Units Code*

[4] Wojciech et al, *Recurrent Neural Network Regularization*

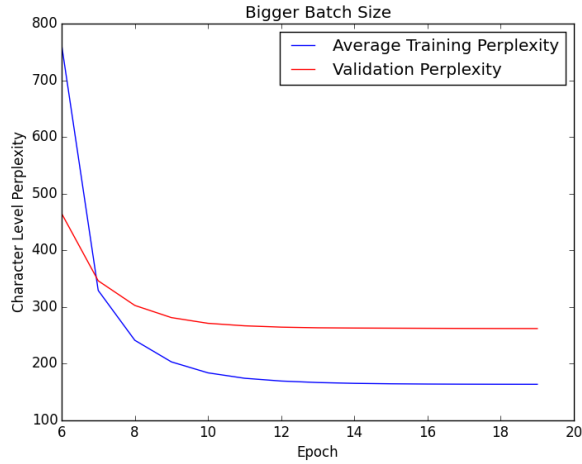[5] Alex Graves, *Generating Sequences With Recurrent Neural Networks*

Figure 1. Perplexity of qualifying model



Figure 2. Output from the 187 Perplexity model