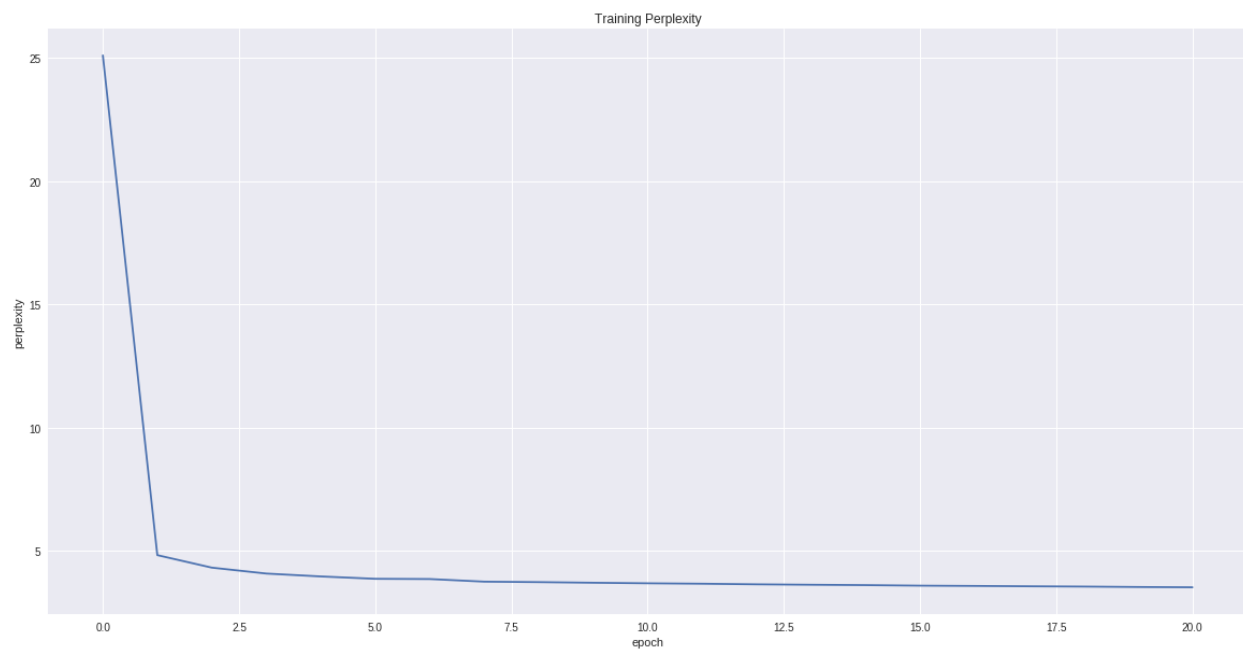
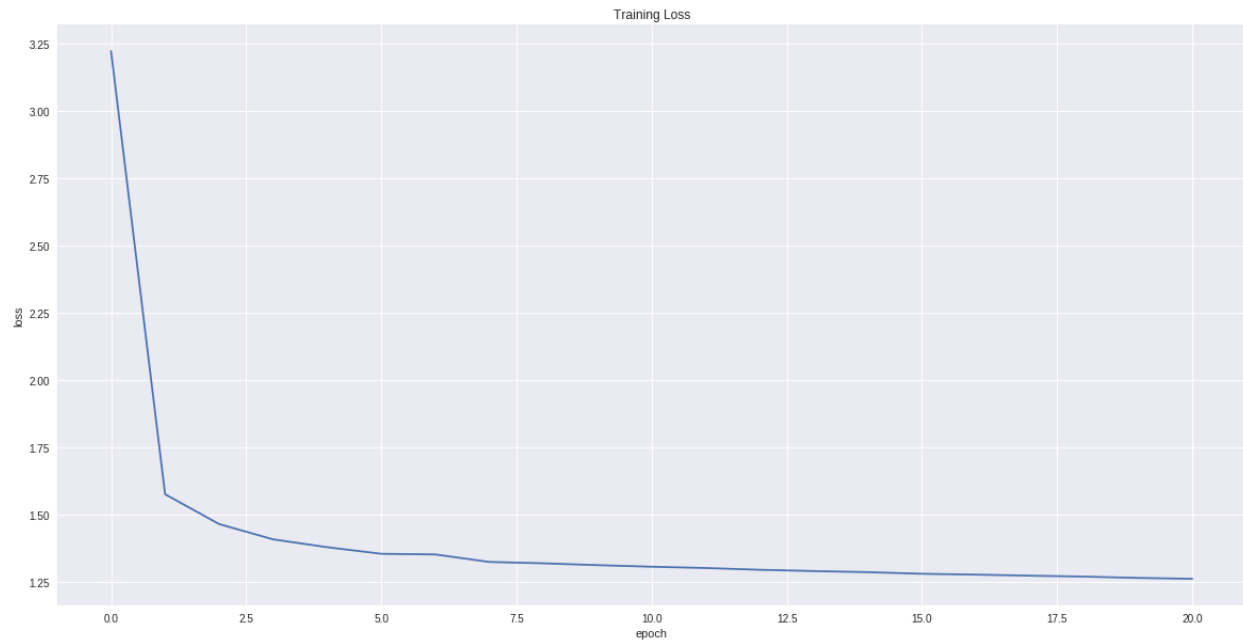
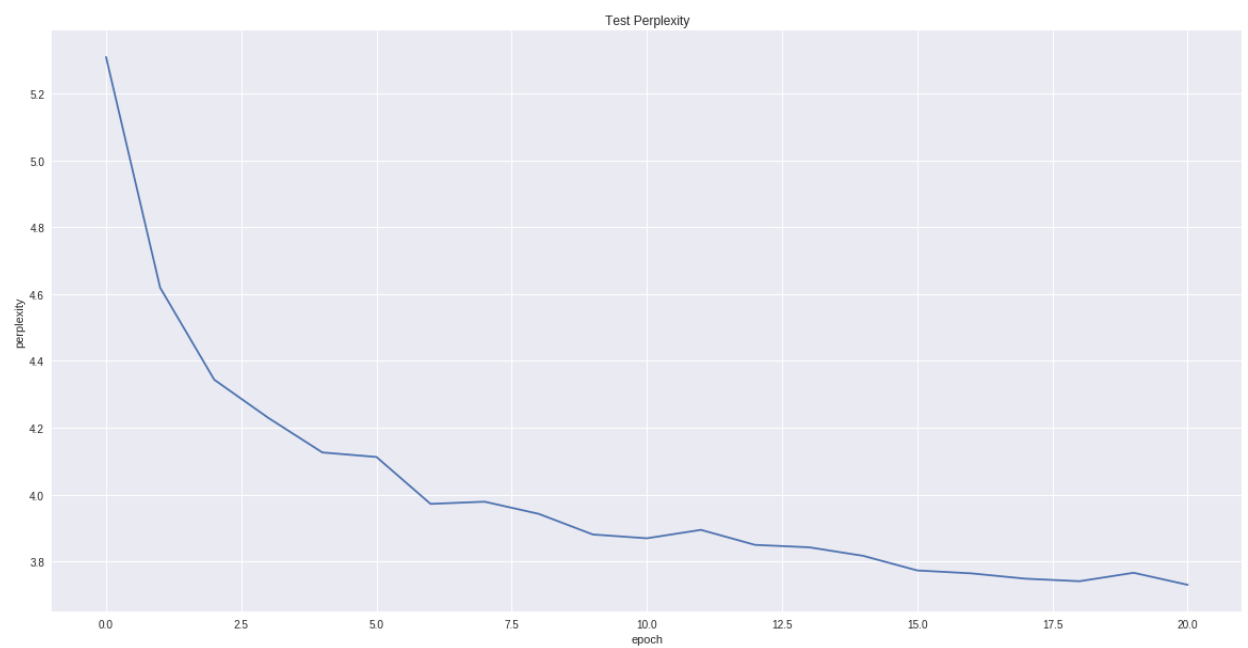
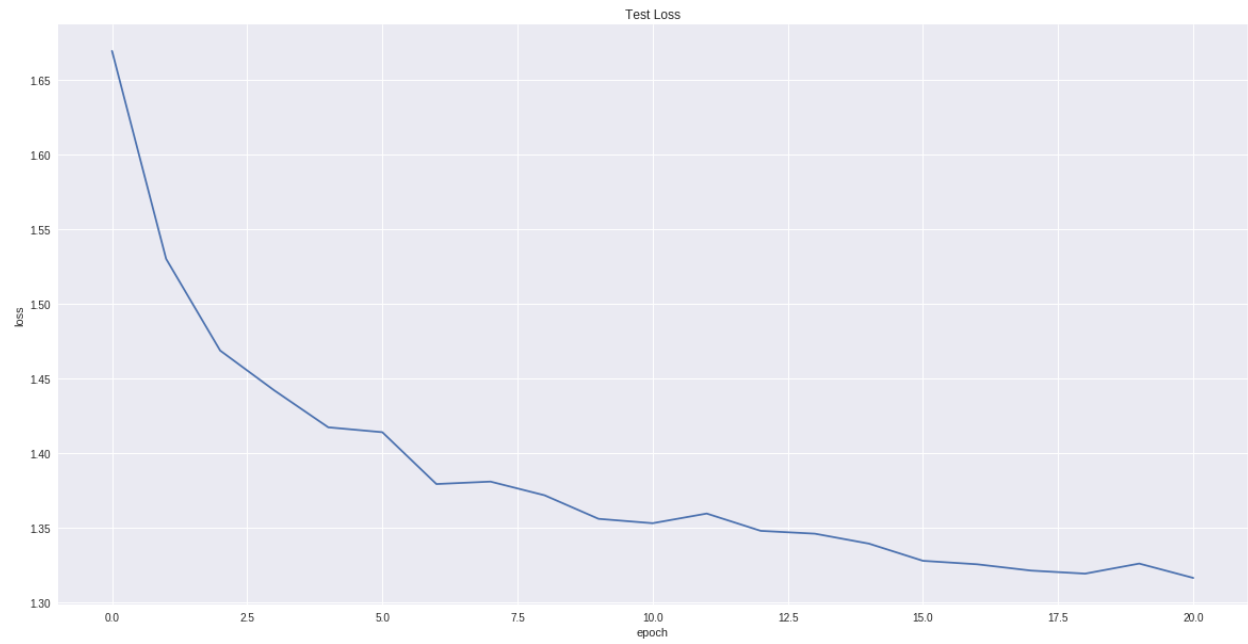
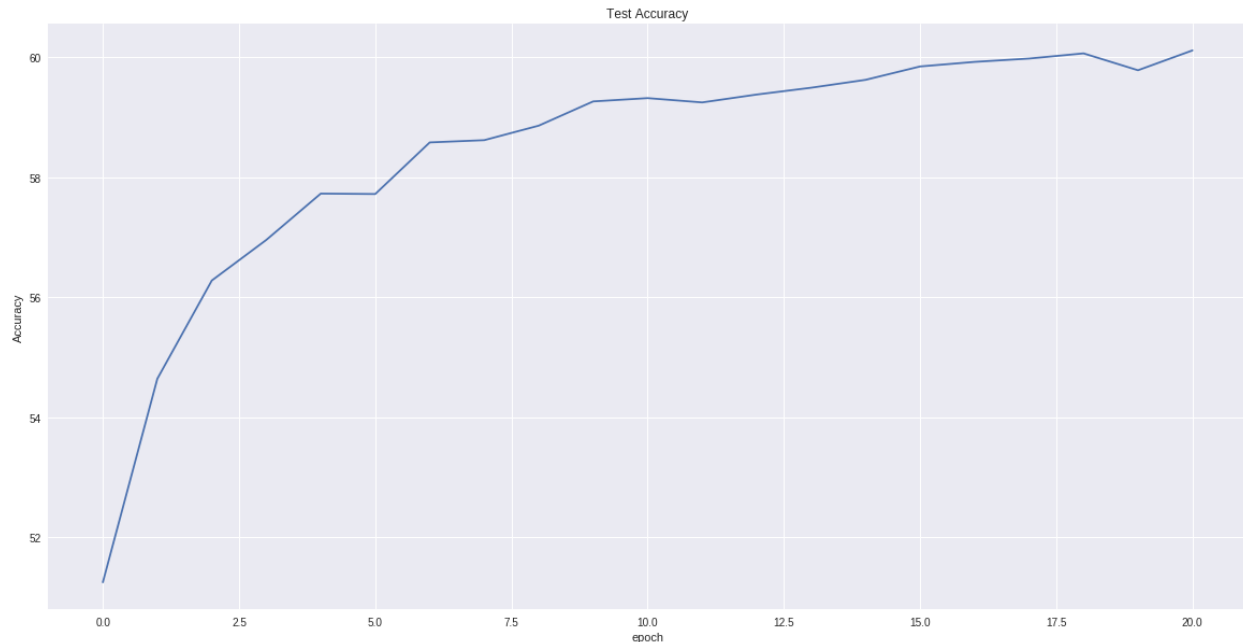


Short Answer Questions – Word Model

Train / Test Performance Figures







Final Test Performance

Accuracy = 60.1%, Perplexity = 3.73

Sentence Generation Results

Prompt: ‘Harry Potter and the’

Sample output (funny parts are underlined>):

- [illegible]

It is interesting that max and sample methods have got into repeating the same sequence of words (to be able to) a little sooner than beam search (of the side). It is likely because beam search explores the space more by following different beams for several steps.

Comparison of Sampling Methods

It seems that sample method generates the most reasonable sentences, which counters my expectation. I was expecting beam search to work better. It essentially obtains multiple samples in each step and explores each further for several steps. It thus explores the space more. With a larger BEAM_WIDTH, the advantage can become apparent but that comes at the cost of longer sampling; I could not test with any BEAM_WIDTH > 15 as colab kept disconnecting.

Comparison of Temperature Values

Sample output (interesting parts highlighted):

- generated with sample (temperature = 0.5) → Harry Potter and the room was a self sign of watch and the fire of the hall was wandering the expression of him and staring at the back of the boy had been purpled out of the portrait of the corridors and said, "Harry wobeam
- generated with beam (temperature = 1.0, beam width = 5) → Harry Potter and the match of the wand was looking at her wand and looked up at the end of the case of the case of the case of the wand was still beneath the back of the wand and looked up at the end of the case of the s

It seems that sample produces the most reasonable results. Beam results still show repetitive phrases.

For temperature in range (0, 1) we exponentially scale the values before normalizing over them which means the values we normalize over are either much smaller than their original (if the exponent was negative) or are much bigger than their original value (if the exponent was positive). That is, the probability distribution has more pronounced peaks and valleys. This makes sampling of the max more likely (under the asymptotic case of temperature → 0, sampling is equivalent to max). The opposite is true if temperature > 1, as we are smoothing out the probability distribution. This would encourage exploration for beam sampling that would allow us possibly produce better results. It will likely confuse sampling as the distribution we are working with now is smoother. If temperature = 1, we are working with the original probabilities.

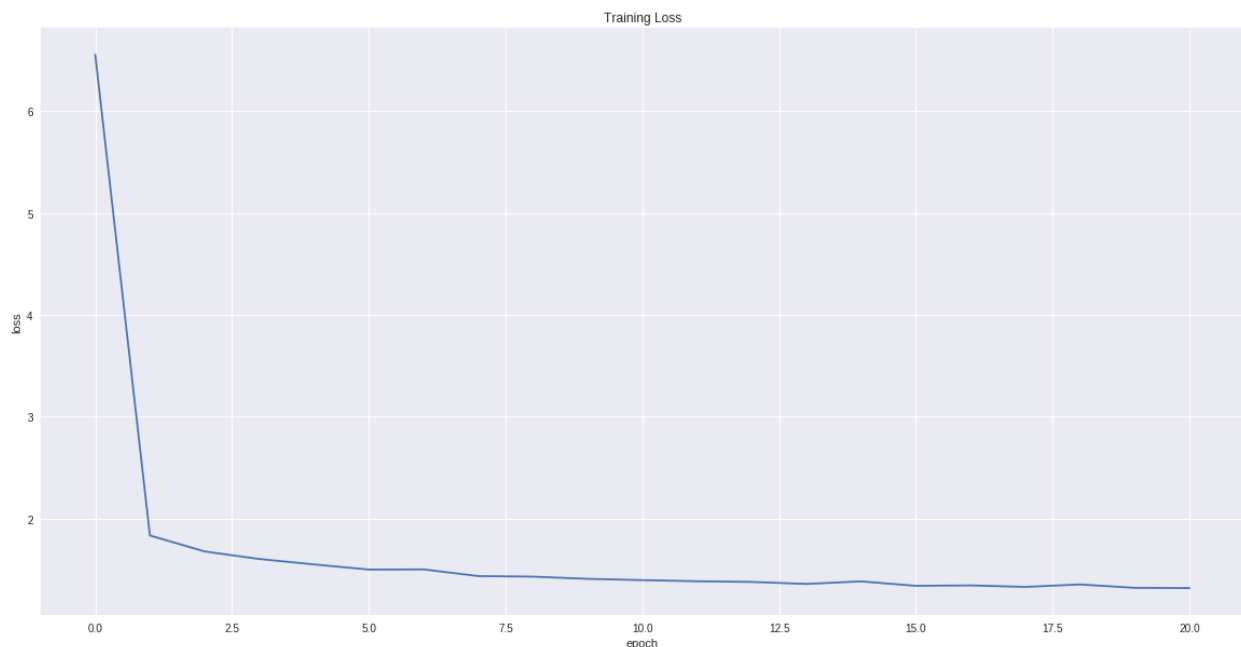
Short Answer Questions – New Corpus

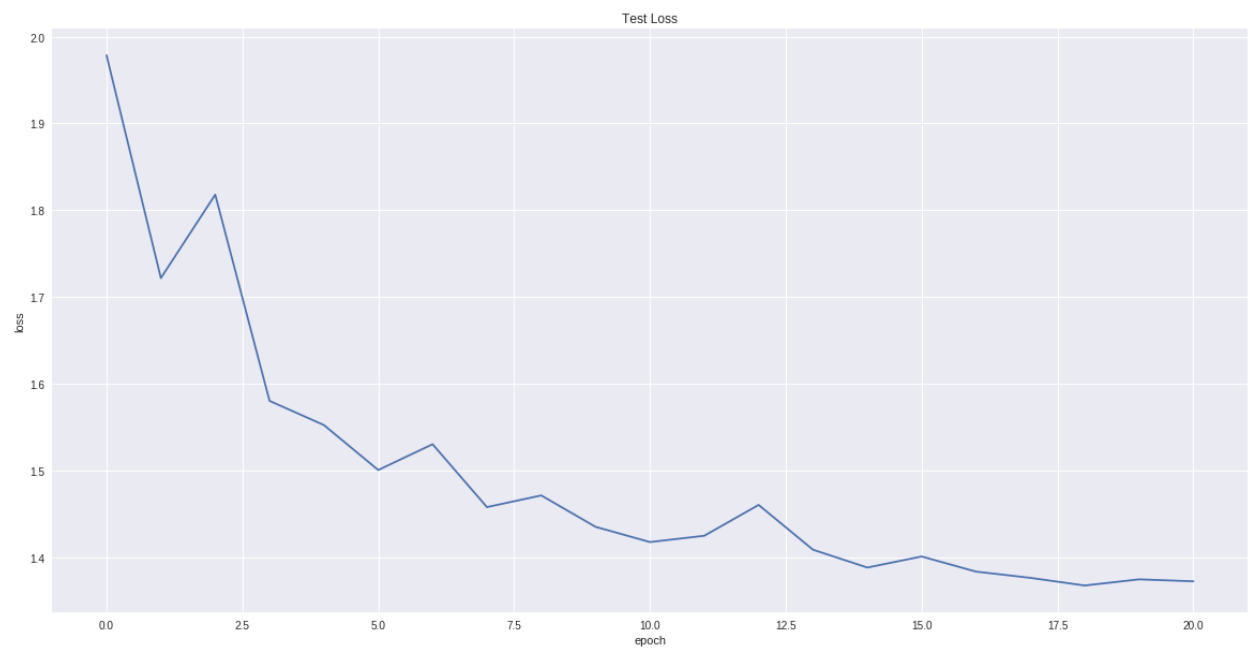
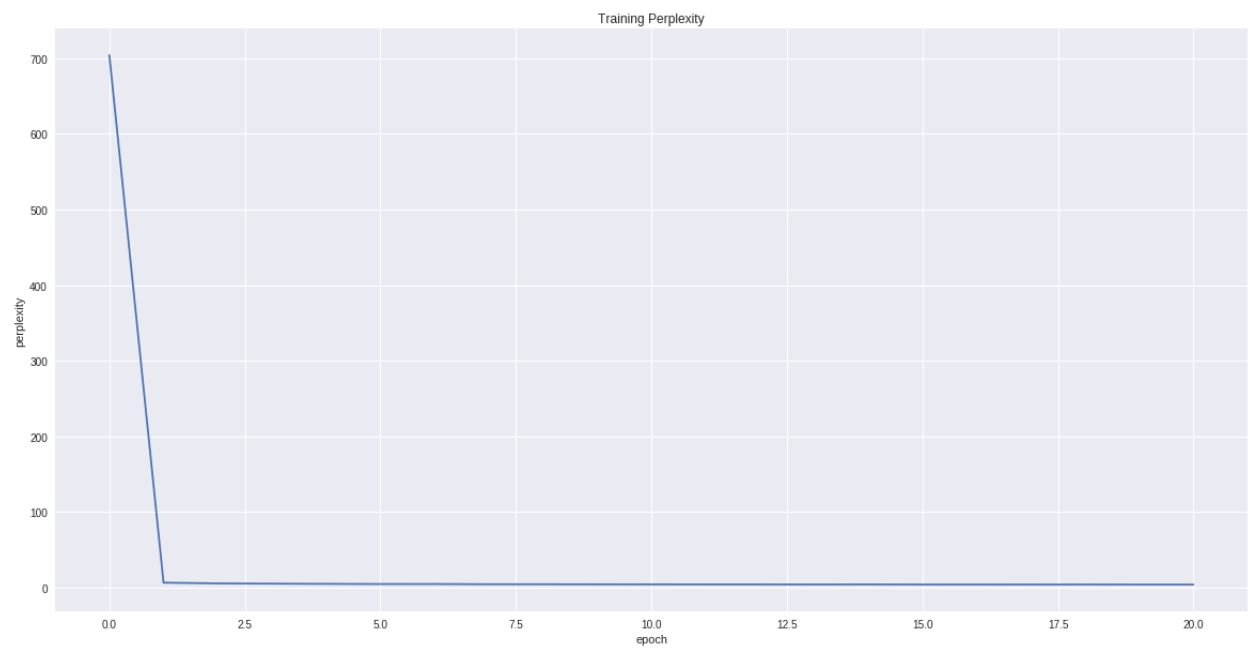
Corpus Information

Name: Lord of the Rings, Number of Characters = 2579193

Comparison of Sentence Generation on the Old and New Corpus

The sentences get into repeating a particular sentence a lot sooner with this new data. The test performance is lower too (see images below). This is likely because the new data set I have used is smaller (one third of the harry potter dataset in terms of file size).





- generated with beam (temperature = 1.5, beam width = 5) → found the hobbits of his face. There was a shadow of the hobbits were shall been silence. There was a dark shadow of the lands of the roads of the roads of the roads of the roads of the roads of the road. There was a

It's cool that 'Ring' but not 'ring' appears in the generated sentences. Also note that with a small temperature beam max has been caught up in repetition (probability distribution has more pronounced valleys and peaks). Beams performance improves as temperature goes above 1 (much less repetitions happen). However, sample generates more gibberish sentences as the temperature goes up.

Short Answer Questions – Student Forcing

Difficulties with Student Forcing

Training becomes slow and unstable.

Comparison with Teacher Forcing

The results are better with teacher forcing as by feeding in the correct value we prevent the error from being accumulated. By helping the network in this way we prevent it from falling in a local minima.