

LSTM neural network for textual ngrams

Shaun C. D'Souza

Accenture

shaun.c.dsouza@accenture.com

Abstract

Cognitive neuroscience is the study of how the human brain functions on tasks like decision making, language, perception and reasoning. Deep learning is a class of machine learning algorithms that use neural networks. They are designed to model the responses of neurons in the human brain. Learning can be supervised or unsupervised. Ngram token models are used extensively in language prediction. Ngrams are probabilistic models that are used in predicting the next word or token. They are a statistical model of word sequences or tokens and are called Language Models or Lms. Ngrams are essential in creating language prediction models. We are exploring a broader sandbox ecosystems enabling for AI. Specifically, around Deep learning applications on unstructured content form on the web.

Keywords— Deep learning, Cognitive, LSTM, Neural network, Ngrams

1 Introduction

Ngram models work on the basis that we can predict the next token given the previous n-1 tokens. We use the following notation to compute the probability of a word sequence. In order to represent a random variable X taking on the value y we use $P(X = "y")$ or the simplification $P(y)$. Now, to compute the joint probability of a sequence of words $w_1...w_n$ we use $P(w_1, w_2, ..., w_n)$.

We compute the probability of an entire sequence of word by decomposing this probability using the chain rule of probability

$$P(w_1, w_2, ..., w_n) = P(w_1|w_2, ..., w_n) * P(w_2|w_3, ..., w_n) * P(w_3|w_4, ..., w_n) * P(w_n) \quad (1)$$

The chain rule shows the relation between computing the joint probability of a sequence of words given the conditional probability of the previous sequence of words. Using a Ngram model allows us to further simplify this equation as we estimate the probability of a word given its history by approximating the last N words. A Bigram model for example would use the conditional probability of a word given the word before it.

$$P(w_1|w_2, ..., w_n) \approx P(w_1|w_2) = \frac{C(w_1, w_2)}{C(w_2)} \quad (2)$$

A trigram model allows us to improve our predictability by using the preceding 2 word tokens.

$$P(w_1|w_2, ..., w_n) \approx P(w_1|w_2, w_3) * P(w_2|w_3) \quad (3)$$

The easiest way to estimate these probabilities is using the count value of the token sequences in the training data.

$$P(w_1|w_2, ..., w_n) = \frac{C(w_1, w_2, ..., w_n)}{C(w_2, ..., w_n)} \quad (4)$$

To obtain the count values for our tokens we use the ngram utility [10] to obtain a set of token sequence counts for all the data in the training set. For the purposes of our investigation we used the textual book data from the Gutenberg project and Brown data set. The Google Books Ngram Corpus [8] is available at <http://books.google.com/ngrams>.

Total 9624 unique ngrams in 890415 1-grams	
the	50869
N	32481
of	24406
to	23662
a	21639

Table 1: 1-grams

Total 273906 unique ngrams in 890414 2-grams	
N_N	7204
of_the	5293
in_the	4507
N_million	4493
to_N	2865

Table 2: 2-grams

2 Research Method

We process the data to obtain the 5-gram tokens using the ngram utility which gives us a count value of all consequent tokens in the training data. We are using textual data for the purposes of our investigation on computing the ngram probabilities.

We use a variety of smoothing techniques to normalize the data and since a large number of token sequences are not in the training data. We use a ngram log probability (NGLP) to estimate the probability of our language model. This allows us to maintain a sum value of the log probability for the training data as it is difficult to compute accumulated product value on decimal values. We use these values to compute the cross entropy of the data.

$$Crossentropy = H = -(\log_2(P(w_1|w_2, w_3) + \log_2(P(w_2|w_3, w_4)) + \dots \log_2(P(w_n - 2|w_n - 1, w_n))) \quad (5)$$

This allows us to compute the cross entropy of the data on the test set which gives us a measure of how well our language model is able to predict the tokens in the code. We calculate a perplexity value equivalent to 2^H .

A low cross entropy means that we are able to accurately predict the next token. If the model predicts every token correctly with a probability of 1, then our cross entropy is 0. We use this data to study multiple types of identifiers in the code including variable and class field definitions, method names and function calls in the code.

We measured the zipf distribution - Fig. 1 of the data in the text [3]. As per Zipf's law we see that the frequency of the tokens in the data set is inversely proportional to its rank in the number count of tokens. The training data contained 9624 unique tokens. The slope is -0.0479.

We plan to extend this code to deep learning applications on unstructured content form on the web along the lines of the Google Brain project [2] and TensorFlow [1]. This will allow us to build a knowledge base [12] using existing projects and reuse code as per the application.

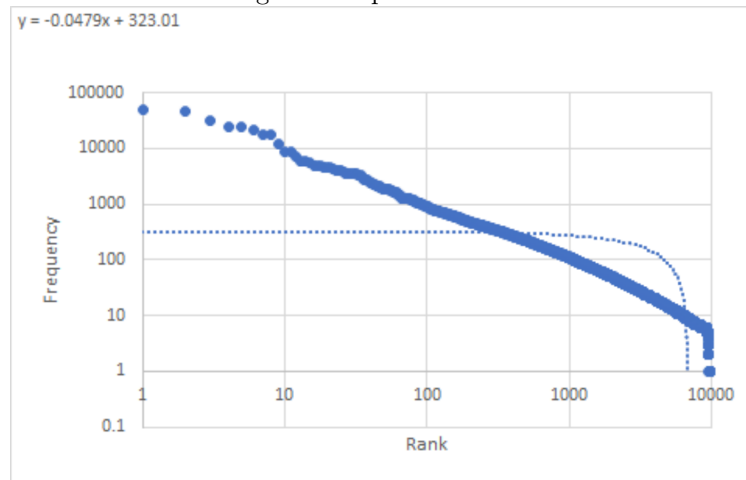
Total 613476 unique ngrams in 890413 3-grams	
the_u_s	920
N_million_or	687
of_N_million	665
N_a_share	631
million_or_N	621

Table 3: 3-grams

3-grams	entropy	7.95
	perplexity	247

Table 4: Cross entropy

Figure 1: Zipf distribution.



3 Results

TensorFlow is an open source library for deep learning developed by Google. It is a python library that is similar to numpy, scipy and uses data flow graphs and tensors for numerical computation. They support the development of neural networks using a set of libraries. A perceptron is a simple neural network designed to use a threshold activation function. It computes the activation of a neuron using the dot product of the input and weight vectors.

$$O(x) = \text{sgn}\left(\sum_{i=0}^n w_i * x_i\right) \quad (6)$$

$$\text{where } \text{sgn}(y) = \begin{cases} 1 & \text{if } y > 0 \\ -1 & \text{otherwise.} \end{cases}$$

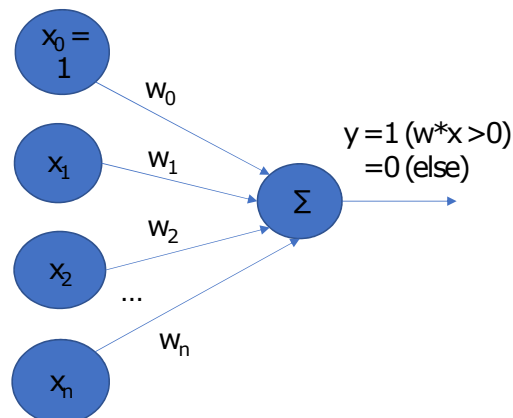


Figure 2: Single layer perceptron.

Sigmoid	Softmax	Hyperbolic tangent
$\sigma(y) = \frac{1}{1 + e^{-y}}$	$f(x) = \frac{e^x}{\sum_{j=1}^N e^{x_j}}$	$\tanh(x) = \frac{1 - e^{-2x}}{1 + e^{-2x}}$

Table 5: Activation functions

Fig. 2 shows a single layer perceptron. For a given dataset the perceptron is guaranteed to find a linear plane of separation described as the hyperplane decision surface in the n-dimension space. The perceptron training rule is used iteratively to update the network weights. The weight vectors are initialized randomly and updated using the rule

$$\begin{aligned} w_i &\leftarrow w_i + \Delta w_i \\ \Delta w_i &= \eta(t - o)x \end{aligned} \quad (7)$$

Additionally, multiple layers can be used in a multi-layer perceptron. This is effective on uni-dimensional data and finds application in a number of natural language processing tasks [4, 9] including part of speech (POS) tagging. The OpenNLP library uses a multi-layer perceptron in its trained model. These have been effectively used in applications in information extraction as shown in [5]. [6] demonstrated the use of a chunker model in detection of semantic triples.

A neural network uses a continuous activation function in each of the layers. Some of the activation functions are in Table 5.

Figure 3: Neural Network.

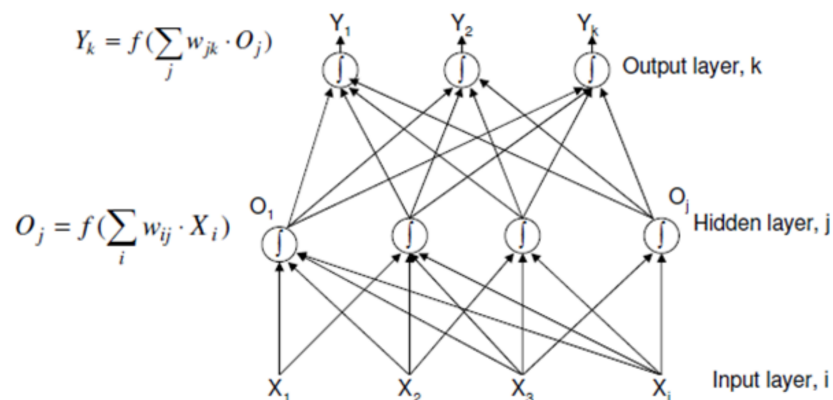


Fig. 3 shows an artificial neural network with Sigmoid activations. The weights in the hidden layer are used in determination of word vectors used in the Continuous Bag-of-words (CBOW) and Skip-gram models [11]. These are used in predicting semantic similarity. We use a TensorFlow Keras LSTM layer - Fig. 4 to output word sentences in the PTB corpus. Below is a generated text sequence corresponding to an input seed value. This can be used in applications like the Google Smart compose for autocompletion of input text given a suitable corpus of training data.

----- Generating with seed: "n plant near <unk> ill. was completed in n <eos> in a disputed n ruling the commerce commission said commonwealth edison could raise its electricity rates by \$ n million to pay for the plant <eos> but state courts upheld"

n plant near <unk> ill. was completed in n <eos> in a disputed n ruling the commerce commission said commonwealth edison could raise its electricity rates by \$ n million to pay for the plant <eos> but state courts upheld a challenge of last year 's that list of the low 's action action of the low are to debt higher prices the in financing <eos> the u.s. until that the latest will leave <eos> the and and japan is expected to slip <eos> the funds have money funds have received these rich since that the decline plot increase <unk> with the low and japan automobile said that the u.s. is that that u.s. is is already for a request <eos> the decline

and japan of their program managers the highest has authority dropped n n in september <eos>
first n n for export and \$ n million debt and year plant <eos>

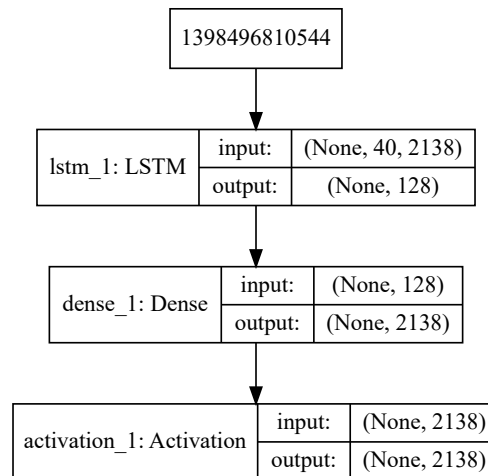


Figure 4: Keras LSTM model.

4 Conclusion

We explore the use of the TensorFlow library in creating a recurrent neural network (RNN). We train a LSTM neural network on textual data from the Penn Tree bank corpus. We see that the LSTM is able to accurately predict the word ngrams using a seed sentence. We plan to extend the work to use POS and chunker sequences [7] in phrase construction. We continue to explore the intersection of Technology and Business in the context of AI, Globalization, CSR and the Last mile with an emphasis on Deep learning applications in the broader web.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [2] M. Abadi, P. Barham, J. Chen, Z. Chen, A. Davis, J. Dean, M. Devin, S. Ghemawat, G. Irving, M. Isard, et al. Tensorflow: a system for large-scale machine learning. In *OSDI*, volume 16, pages 265–283, 2016.
- [3] M. Allamanis and C. Sutton. Mining source code repositories at massive scale using language modeling. In *Proceedings of the 10th Working Conference on Mining Software Repositories*, pages 207–216. IEEE Press, 2013.
- [4] X. Carreras and L. Màrquez. Introduction to the conll-2005 shared task: Semantic role labeling. In *Proceedings of the ninth conference on computational natural language learning*, pages 152–164. Association for Computational Linguistics, 2005.
- [5] S. D’Souza. Parser extraction of triples in unstructured text. *IAES International Journal of Artificial Intelligence (IJ-AI)*, 5(4):143–148, 2017.
- [6] S. C. D’souza. System and method for extracting information from unstructured text, June 19 2018. US Patent App. 15/474,194.

- [7] M. Honnibal and M. Johnson. An improved non-monotonic transition system for dependency parsing. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, pages 1373–1378, Lisbon, Portugal, September 2015. Association for Computational Linguistics.
- [8] Y. Lin, J.-B. Michel, E. L. Aiden, J. Orwant, W. Brockman, and S. Petrov. Syntactic annotations for the google books ngram corpus. In *Proceedings of the ACL 2012 system demonstrations*, pages 169–174. Association for Computational Linguistics, 2012.
- [9] D. Metzler and O. Kurland. Experimental methods for information retrieval. In *Proceedings of the 35th international ACM SIGIR conference on Research and development in information retrieval*, pages 1185–1186. ACM, 2012.
- [10] J.-B. Michel, Y. K. Shen, A. P. Aiden, A. Veres, M. K. Gray, J. P. Pickett, D. Hoiberg, D. Clancy, P. Norvig, J. Orwant, et al. Quantitative analysis of culture using millions of digitized books. *science*, page 1199644, 2010.
- [11] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013.
- [12] R. Rastogi. Building knowledge bases from the web. In *Proceedings of the 18th International Conference on Management of Data*, pages 5–5. Computer Society of India, 2012.