

Q1. Implement an encoder-only **transformer**^[1] with the following specifications:

Non-Linearity	\tanh
Embedding Size	64
Attention heads	4
Encoder layers	4

1. Noun: NN
2. Verb: VB
3. Adjective
4. Others: A

To reduce the dimensionality of the learned word vectors from **word2vec** to **word2vec** word vectors, word vectors are also provided in csv

Initialization

Random Initialization in high dimensional spaces can lead to issues with convergence. This is why we will use **He-initialization**^[4] for initialization of weights. Biases are to be initialized to 0.

Training:

While Stochastic Gradient Descent works, it requires 4x epochs as compared to Adam Optimizer. The bonus goal is to implement Adam Optimizer.

Result:

Report the accuracy, precision and recall for each of the classes.

References:

[1] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. "Attention is all you need." Advances in neural information processing systems (2017).

[2] Erik F. Tjong Kim Sang and Fien De Meulder. "Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition." In Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003, pages 142– 147. (2003)

[3] Mikolov, D., S. Corrado, and K. Chen. "Efficient estimation of word representations in vector space." arXiv:1301.3781 (2013).

[4] He, Kaiming, and Z. Zhang. "Delving deep into rectifiers: Learning a layer-wise activation function for deep neural classification." In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 1026-1034. 2015.

