

1 Question 1

The basic self-attention mechanism be improved by:

- 1- Using a 2-D matrix to represent the embedding, with each row of the matrix attending on a different part of the sentence. This allows us to interpret the sentence embedding in depth in our model
- 2- Adding a regularization term for the model, to discourage the redundancy in the embedding.

2 Question 2

The main motivations for replacing recurrent operations with self-attention:

Self-attention mechanism is more parallelizable and requiring significantly less time to train. In fact, the total computational complexity per layer is reduced. The third motivation is the path length between long-range dependencies in the network, in fact, the shorter these paths between any combination of positions in the input and output sequences, the easier it is to learn long-range dependencies. And the last benefit is that self-attention could yield more interpretable models.

3 Question 3

For this question I decided to plot the coefficients for the last review, let's start with coefficients for each sentence in the review:

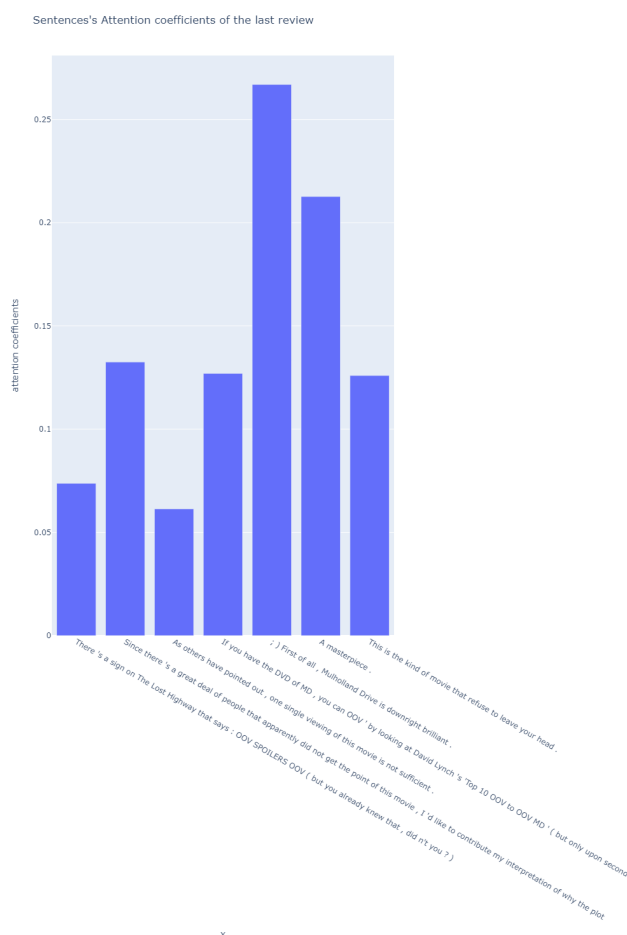


Figure 1: Attention coefficients for sentences.

We can notice that the most important sentence was: "First of all, Mulholland Drive is downright brilliant", which gives a positive impression about a scene of the movie. The second most important sentence is "A masterpiece" which is a sufficient sentence to tell that the review is positive. We can say in total the algorithm can give fair weights to each sentence depending on its contribution to building the review.

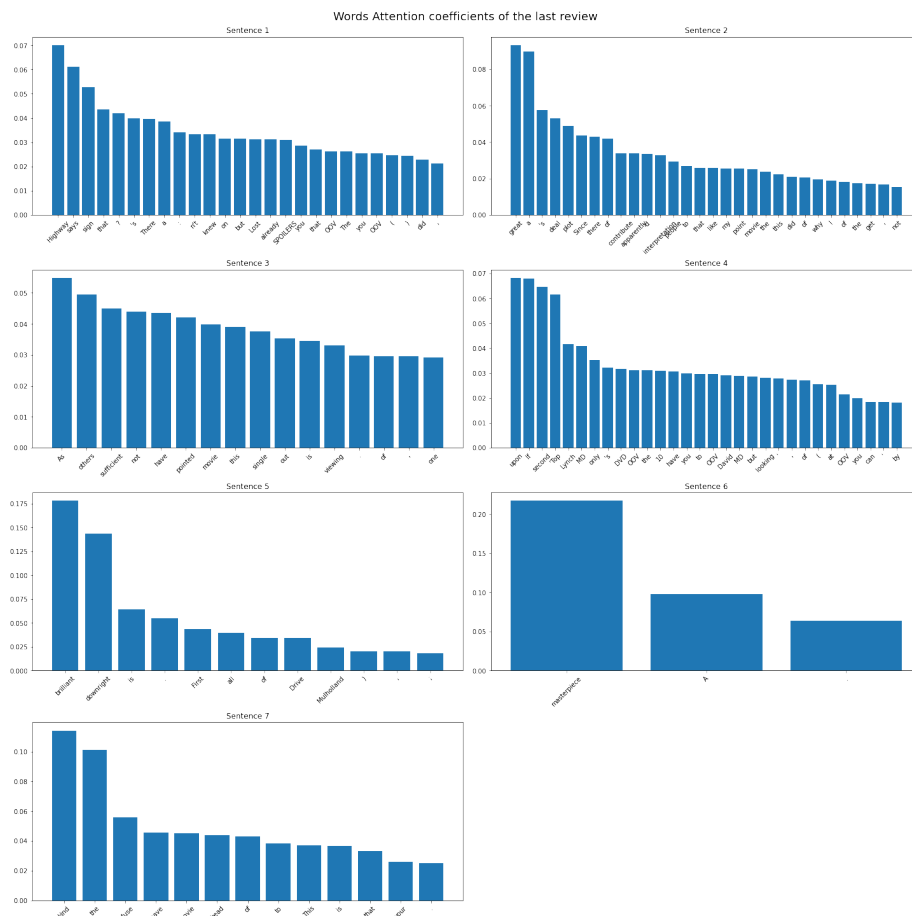


Figure 2: Attention coefficients for words.

We can see also that the algorithm perform well on attributing more weights to the key words of each sentence.

4 Question 4

The major limitation of HAN architecture is that each sentence is encoded in complete isolation. In other words, while encoding the representation of a given sentence in the document, HAN completely ignores the other sentences, thus HAN is a context-blind self-attention mechanism (even if at level 2 we attribute importance scores to sentences their encoding vectors have already been formed, and it is too late to modify them).