

Pre-Processing issues:

- Vocab size: $\sim 261k$ from just $\sim 17k$ entries

Learning custom embedding for the vocabulary of this magnitude provides more scope to overfit.

- Lack of punctuations \leftarrow default tokenizers)

Punctuations were understood to be important in certain cases of positive classes. (is definition) based on our manual analysis over the data.

- More ambiguity in samples (Example sentences)

The guidelines for definition were not clear in certain cases. Some examples from Notes.md

Feature solution

- Meaningful replacement for oov

Replace with pos tags instead of <unk>

- pos with Punctuations

Using "PUNCT" instead of all punctuation is ambiguous. So punctuations from the set indicated in Notes.md are used as is.

- Dependency feature:

Concatenation of modifier, head and label representations.

- Tokens + Pos + punct

Concatenation of $\begin{matrix} \text{Token} + \text{pos} \\ \text{rep} \quad \text{rep} \end{matrix}$ $\left[\begin{matrix} \text{Our Inclusion updated} \\ \text{to this paper work} \end{matrix} \right]$

Assumption:

Syntactic is all we need:

Initial assumption: Semantics might not have any effect for this classification task.

Turns out: Using Pretrained word vector has a positive impact

Todo: May need to run an Exclusion Exp to support this case.