

COLLOQUIUM REPORT-II

Ups and Downs of Word Embeddings

Submitted By: Agasthya Harekal(adh5677@psu.edu)

This Colloquium the speaker spoke about word embeddings which is one of the fundamental concepts of Natural Language Processing. Word Representation were described as one of the most primitive concepts in NLP. Word Embedding are described as vector representations, where prediction of words can be done based on other words in context. The subtraction of a word from another word and addition would result in a new word where the semantics of the underlying words are understood. Word Representations can be used to get accuracy for various NLP tasks. Many facets of the research were discussed that would aide towards getting better results in NLP. First one discussed was interpretability, where discussion was done how different people interpret the things in many ways. Here, discussion was done of the word representations as to the ways to derive the interpretability of them. Example was given of the wordnet which has rich semantic relations encoded where things are very easily interpreted. Word embedding/representations are often used as features for sentiment analysis, text classification and other tasks. To be meaningful, word embeddings need to pre-trained, and trained on all classes jointly. When a model is given for use in domain the interpretability plays a major role. Next aspect is the generalisability, where certain words are meant to be similar to each other than others. Also, models perform in the same domain much better than models that are trained in other domain. Third facet is the fairness, where unbiased, trustworthy, privacy preserving embeddings are expected. Fourth is the robustness, where stability of embeddings are expected where specifically stability across different training seeds, across different domains, across different algorithms are sorted after. In the research done by the speaker, stability is defined as percentage overlap between ten nearest neighbouring words in an embedding space. Here, different dataset, algorithms is considered and stability across different dataset/algorithms are considered. Careful research was done using dataset such as New York Times data, Europarl data and Algorithms such as word2vec, skip-gram model, glove, ppm model was considered. The conclusion was to fit a regression model using features such as word, algorithm, data to get stability of these models. After using this model, there were multiple findings, it was understood that frequency is not a major factor in stability. Curriculum learning was understood as an important aspect of stability which is the order of training data given to the algorithm. Also, it was found that part of speech Tagging(POS) is important aspect for stability and in terms of the models GLOVE model was the most stable embedding algorithm among others. The stability of the embedding were checked using other definitions such as using twenty, thirty neighbours and also using threshold distance to the neighbours were also considered, which gave similar results. The stability of various languages was studied where various insights of different languages were gathered.

The energy and cost of various NLP models were considered to find the most energy and cost efficient of the models. Also, efficiency of the models

were considered as opposed to accuracy of the model. The talk concluded that one should look towards solving the NLP problem using traditional NLP techniques. Also, there are always trade-offs to be considered when considering these facets of the model.

Word embeddings are used in field of Natural Language Processing since inception in areas such as distributional semantics. “ A word is characterised by the company it keeps” was said by John Rupert Firth. Word embedding is a numeric vector that represents a word. Words with similar meaning have similar representation. The word embedding can be of varying sizes such as 50, 100, 200 etc. They are basically fed into the machine learning models for training them for various tasks such as text classification, document classification, statistical analysis of documents, sentiment analysis. Based on usage of these words in sentences, word embeddings try to capture the semantics, contextuality and meaning of words. The three types of word embeddings used in general are 1) Glove 2)Word2Vec 3) ELMo

Glove Model-The Glove Model is trained on non-zero entries of a global word-word co-occurrence matrix which captures the frequency of a co-occurrence of a word with another in a given corpus.

Word2Vec Model- The Word2Vec model is an algorithm built on distributional hypothesis which states that words occurring in related linguistic contexts will have related meaning.

Elmo Model- The ELMO Model derives the word vectors using Bi-Directional LSTM, by pre-training on large text corpus.

The overall talk focused on the importance of word embeddings in context of Natural Language Processing and how the various facets of word embeddings have their individual contribution towards overall quality of the word embeddings which lead to the better performance of the models trained especially in terms of the efficiency of the models.

