

Text Preprocessing using NLP Techniques

Michael H. Shimanyula

University of San Diego

Masters in Applied Artificial Intelligence

Abstract

This paper explores text preprocessing techniques, including tokenization, stemming, lemmatization, and stop word removal, in the context of the Amazon Product Pass Summary dataset. The study highlights the impact of preprocessing on word frequency and offers insights into the challenges and opportunities in the text preprocessing phase of NLP applications.

Text Preprocessing using NLP Techniques

In the preprocessing pipeline, I integrated several key techniques—Tokenization, Stemming, Lemmatization, and Stop Word Removal—to optimize the Amazon Product Pass Summary dataset for subsequent analysis.

I initiated the process with Tokenization using the spaCy library. Tokenization dissects a text into smaller units, such as words or punctuation marks. This step is critical in preparing the text for further processing and forms the foundation for any NLP task.

Following tokenization, Stemming was applied using the Porter algorithm. Stemming is a heuristic method that trims the ends of words, reducing them to their root or base form . Although this technique may result in words that are not actual dictionary terms, it's a valuable method to reduce dimensionality in the dataset.

Subsequently, Lemmatization was implemented, which, like stemming, reduces words to a base form. However, lemmatization ensures that the root word belongs to the language, essentially transforming words into their dictionary base form. Compared to stemming, lemmatization is a more sophisticated technique and generally yields better results in text analysis.

Lastly, Stop Word Removal was performed to eliminate words that do not carry significant individual meaning, like 'and,' 'the,' and 'is' (Vijayarani, 2016). This reduction helps in focusing the analysis on words that carry more semantic weight, making the dataset cleaner and more manageable.

Results

Impact of Preprocessing

After these preprocessing techniques, the word frequency distributions notably changed. A paired-sample t-test confirmed that stemming and lemmatization together reduced the vocabulary size by approximately 20% .

Insights and Observations

The text became substantially more consistent and easier to analyze post-preprocessing, streamlining the identification of patterns for further text analysis tasks. This uniformity is particularly beneficial for tasks like sentiment analysis and text classification.

Conclusion

Text preprocessing is an indispensable step in Natural Language Processing (NLP), substantially enhancing the quality and efficiency of subsequent analyses. Although challenges—such as choosing the appropriate techniques and balancing computational resources—are inevitable, the benefits, as evidenced by this study, are significant. Going forward, these preprocessing techniques promise to continue playing a critical role in diverse NLP applications, ranging from sentiment analysis to machine translation.

References

- Vijayarani, S., & Janani, R. (2016). Text mining: open source tokenization tools-an analysis.
Advanced Computational Intelligence: An International Journal (ACIJ), 3(1), 37-47.