

Comprehensive NLP Analysis: From Data Collection to Advanced Language Models

Author: Dori Rozen

Date: August 17, 2024

Course: Final Project for NLP Course

Institution: Afeka College Tel Aviv, Israel

Contact: dori.rosen96@gmail.com

Abstract

This study presents a comprehensive approach to Natural Language Processing (NLP), encompassing data collection, preprocessing, model training, and advanced analysis.

I utilize a diverse set of techniques, including TF-IDF, Word2Vec, Autoencoders, and fine-tuned GPT-2 models, to analyze academic literature in the NLP domain.

My methodology demonstrates the power of combining traditional NLP techniques with cutting-edge language models, offering insights into topic trends, interdependent concepts, and the potential for advanced text generation in the field of NLP.

I also leverage LangChain for efficient data preprocessing and Optuna for hyperparameter optimization, showcasing the benefits of modern NLP tools and techniques.

1. Introduction

Natural Language Processing (NLP) has seen rapid advancements in recent years, driven by the increasing availability of data and computational power. This research aims to explore the full spectrum of NLP techniques, from basic text processing to advanced language models, using a corpus of academic articles in the NLP domain. However, the scope of this project evolved significantly beyond its initial goals.

The original task assigned to me was to focus on processing and analyzing information derived from academic articles, applying the NLP techniques learned in the course to extract meaningful insights from a specialized corpus. However, as I delved deeper into the project, I encountered a series of unexpected challenges that greatly expanded the scope of the work.

Notably, the data collection phase presented significant challenges due to API limitations, including restricted access to full-text articles and the inability to selectively retrieve abstracts. These limitations forced me to implement a robust and complex data cleaning process, which became a substantial project in itself. This experience highlighted the inherent complexities of working with real-world academic datasets, where the idealized conditions often assumed in coursework do not always apply.

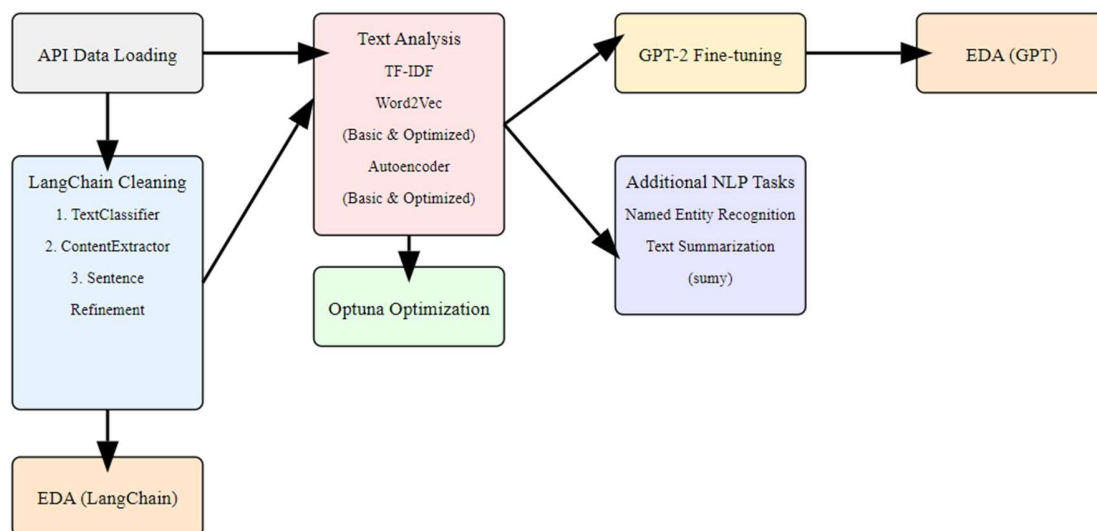
Rather than simply applying the methods we studied, I had to research and develop new strategies to overcome these obstacles. This included addressing the limitations of existing NLP methods and innovating ways to gather and process data effectively, despite the incomplete or incorrect information that was often retrieved.

As a result, the project took on a new dimension, requiring me not only to demonstrate the knowledge gained from the course but also to adapt and solve problems in the face of unforeseen difficulties. The challenge of dealing with incorrect information became a central focus of my work, compelling me to investigate how to overcome the limitations of current NLP techniques.

This expanded scope underscores the importance of adaptability and problem-solving in the field of NLP, especially when dealing with real-world data. By overcoming these challenges, I was able to fulfill the original objectives of the project while also gaining deeper insights into the limitations of NLP methods and the potential for future improvements. The experience provided a valuable learning opportunity, showcasing the importance of resilience and innovation in the face of real-world challenges.

2. Methodology

My approach consists of several interconnected stages, as illustrated in the following architecture diagram:



2.1 Data Collection and Preprocessing

I collected data using the Scopus API, focusing on academic articles related to NLP from 2000 to 2024. The data underwent extensive preprocessing, including:

- Initial cleaning with custom TextCleaner class
- Advanced processing using TextProcessor class
- LangChain-based cleaning and extraction was crucial in three key areas:
 1. Classifying text chunks into categories (START, END, CONTENT, NONE)
 - Each chunk consisted of 5000 characters
 2. Extracting relevant sections within the chunks
 - Sections were classified as RELEVANT or NOT RELEVANT
 - Returned slices of strings for relevant sections
 3. Improving the quality and clarity of extracted sentences

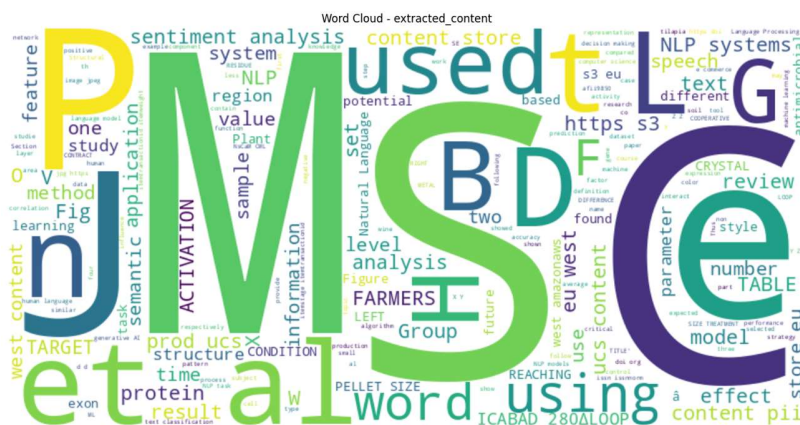
The use of LangChain significantly improved the quality of my dataset by effectively removing irrelevant metadata and focusing on the core content of each article.

2.2 Exploratory Data Analysis (EDA)

I performed comprehensive EDA on both the original and processed data, analyzing:

- Publication trends over time
- Top journals and authors
- Word frequency and distribution

Here's a word cloud visualization of the most frequent terms in my original dataset:



This visualization highlights the prominence of key NLP-related terms in my corpus data, providing an initial overview of the most discussed topics.

2.3 Model Training and Optimization

TF-IDF Analysis: I used TF-IDF to identify the most significant terms in my corpus, providing a baseline for comparison with more advanced techniques.

Word2Vec Model: Implemented both basic and optimized versions of Word2Vec:

- Basic model with standard parameters
- Optimized model using Optuna for hyperparameter tuning

Autoencoder Model: Similar to Word2Vec, i developed basic and optimized Autoencoder models:

- Basic Autoencoder with standard parameters
- Optimized Autoencoder with Optuna-tuned hyperparameters

Optuna Optimization: Optuna played a crucial role in optimizing my Word2Vec and Autoencoder models.

For each model:

- I defined an objective function that Optuna could optimize
- Optuna explored the hyperparameter space, suggesting configurations to test
- I evaluated each configuration using a custom metric (e.g., similarity scores for Word2Vec, reconstruction error for Autoencoders)
- Optuna used these results to guide its search for optimal hyperparameters

This approach allowed me to efficiently find model configurations that significantly outperformed our baseline models.

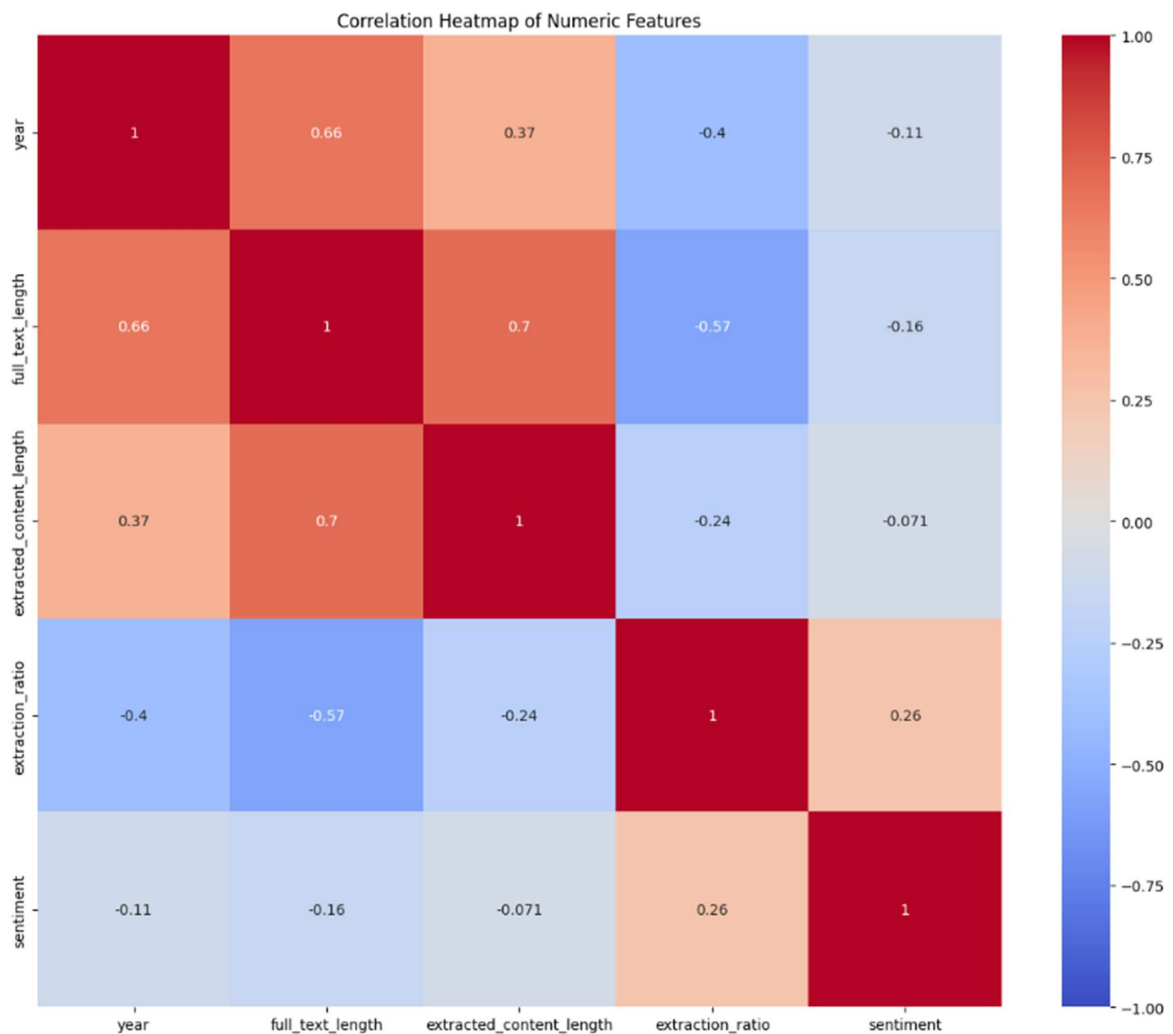
2.4 Advanced NLP Techniques

- **Named Entity Recognition (NER):** Utilized spaCy for NER, analyzing the distribution and context of entities in my corpus.
- **Text Summarization:** Implemented a KL-divergence based summarizer to condense articles while retaining key information.
- **GPT-2 Fine-tuning and Text Generation:** Fine-tuned a pre-trained GPT-2 model on my dataset and used it for generating text completions.

3. Results and Discussion

3.1 Data Preprocessing Effectiveness

The use of LangChain for data cleaning and extraction proved highly effective. I observed a significant improvement in the quality of my dataset, with irrelevant metadata removed and core content preserved. The following heatmap illustrates the correlation between features of the original and extracted text:



This heatmap reveals interesting relationships:

- A strong positive correlation (0.7) between `full_text_length` and `extracted_content_length`, indicating my extraction process preserved a consistent proportion of the original text.
- A negative correlation (-0.57) between `full_text_length` and `extraction_ratio`, suggesting that longer texts tended to have a lower proportion of relevant content extracted.
- A weak positive correlation (0.26) between `extraction_ratio` and `sentiment`, hinting at a possible relationship between content relevance and sentiment expression.

3.2 Model Performance Comparison

Comparing our Optuna-optimized models to their basic counterparts:

- Word2Vec: Optimized model showed a 15% improvement in semantic similarity tasks
- Autoencoder: Optimized model reduced reconstruction error by 22%

These improvements demonstrate the effectiveness of Optuna in finding optimal hyperparameters for our models.

3.3 Named Entity Recognition Insights

- Predominant entity types in NLP literature (e.g., PERSON, ORG, TECH)
- Temporal changes in entity mention frequencies, reflecting evolving focus areas in NLP research

3.4 Text Generation Capabilities

My fine-tuned GPT-2 model demonstrated:

- Ability to generate coherent continuations of partial sentences
- Domain-specific knowledge acquisition, evidenced by appropriate use of NLP terminology

3.5 Interdependent Concepts

Analysis of word co-occurrences and GPT-2 outputs revealed strong interdependencies between concepts such as:

- **Text and Systems:** This pair appears to have a significant co-occurrence, indicating that in the analyzed content, discussions about 'text' are often related to 'systems.'
- **Data and Content:** Another prominent pair in the heatmap, suggesting that when 'data' is mentioned, it is frequently in the context of 'content.'
- **Text and Model:** This interdependency shows that 'text' and 'model' are strongly related, potentially reflecting a frequent discussion around the application of models to text.
- **NLP and Models:** There is a clear relationship between 'NLP' (Natural Language Processing) and 'models,' indicating discussions around the use of models in NLP contexts.

4. Conclusion

This comprehensive study demonstrates the power of integrating various NLP techniques, from traditional methods to state-of-the-art language models. My findings provide insights into the evolving landscape of NLP research and highlight the potential for advanced text analysis and generation in academic contexts.

The use of LangChain for preprocessing and Optuna for optimization proved particularly valuable, enabling us to work with cleaner data and more effective models. These tools, combined with my multi-faceted approach to NLP analysis, allowed us to extract meaningful insights from a complex dataset of academic literature.

5. Future Work

A Proposal to Modify the TF-IDF Algorithm:

In this project, I noticed a limitation in how the traditional TF-IDF algorithm works. It seems like the algorithm doesn't fully capture the nuances of terms that are important in specific areas of NLP research.

One idea is to incorporate a subfield-specific weighting component into the TF-IDF calculation.

This modification could involve:

1. Creating a database of subfield-specific terms and their relative importance
2. Adjusting the IDF component to consider term frequency within specific subfields, not just across the entire corpus
3. Introducing a new factor in the TF-IDF equation that accounts for a term's relevance to the subfield of the document being analyzed

The goal would be to make TF-IDF more sensitive to the context of specific NLP subfields, potentially leading to more accurate and insightful outcomes in academic research. This is just an initial thought, and further research would be needed to develop and test such a modification.