



Automated Text Classification

Group: SOLO

Gayathri

Kodakandla

a1939114

The University of Adelaide

4533_COMP_SCI_7417_7717 Applied Natural Language

Processing

Lecturer: Dr. Orvila Sarker

Section	Title	
1	Abstract	2
2	Introduction	2
2.1	Data Collection Process.....	2
2.2	Summarisation Methods.....	2
3	Preprocessing	2
3.1	Text Cleaning.....	2
3.2	Feature Extraction.....	3
4	Data Visualisation	3
4.1	Word Cloud.....	3
4.2	Sentiment Distribution.....	3
4.3	Topic Distribution.....	4
4.4	ML vs VADER Confusion Matrix.....	4
5	Data Categorisation	5
5.1	Definition of Categories.....	5
5.2	ML Classification Results.....	6
6	Sentiment Analysis	6
6.1	VADER vs ML Agreement.....	6
6.2	Classification Report.....	6
7	Conclusion	7
8	References	7
	GITHUB REPOSITORY LINK	7

1. Abstract

This project aims to automate the classification of Stack Overflow posts related to Natural Language Processing (NLP) using a combination of rule-based keyword categorization and machine learning models. A dataset of 16,005 posts was collected using the Stack Exchange API, pre-processed with standard NLP techniques, and analyzed to assign categories and sentiment labels. The classification model achieved an accuracy of 73.26%, and sentiment analysis using VADER showed an ensemble agreement of 82.4%. This work demonstrates a pipeline for text preprocessing, categorization, and sentiment extraction from community-generated technical content.

2. Introduction

The increasing volume of developer discussions on platforms like Stack Overflow makes it challenging to filter and categorize posts effectively. Automating this process can help identify key NLP subtopics and sentiment trends within the community. This project explores a classification framework for NLP-tagged questions using traditional machine learning and lexicon-based sentiment analysis.

2.1 Data Collection Process

Posts tagged with [nlp] were fetched using the Stack Exchange API. A total of 16,005 unique questions were collected and saved in CSV format. The collected fields included question_id, title, body, and tags.

2.2 Summarisation Methods

Extractive Summarisation: This method selects and combines the most important sentences or phrases from the original text to form a summary.

Abstractive Summarisation: This method generates a summary by interpreting and paraphrasing the original content, possibly including new words or phrasing not present in the source.

3. Preprocessing

3.1 Text Cleaning

Each post's title and body were concatenated into a single text field. The following preprocessing steps were applied:

- Lowercasing
- Removal of non-alphanumeric characters
- Tokenization

- Stop word removal
- Lemmatization (using WordNetLemmatizer)

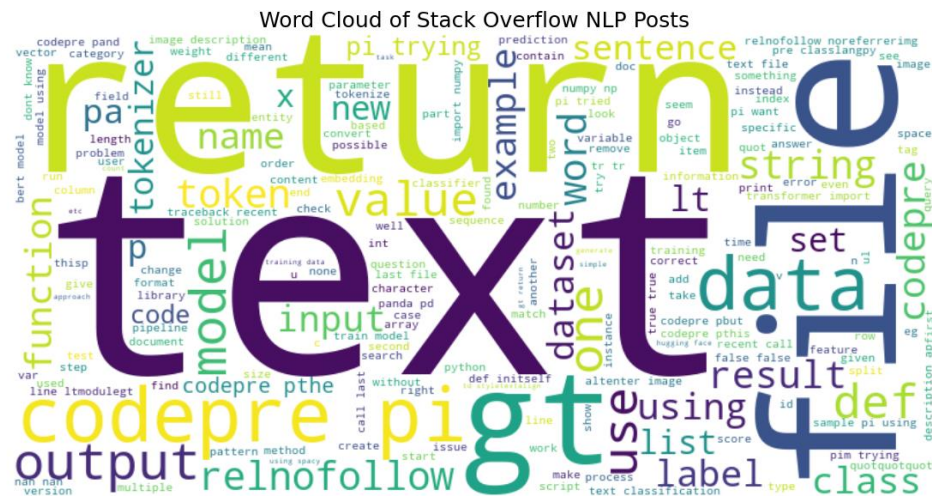
3.2 Feature Extraction

TF-IDF vectorization was used to transform the cleaned text into a numerical feature matrix suitable for model training.

4. Data Visualisation

4.1 Word Cloud

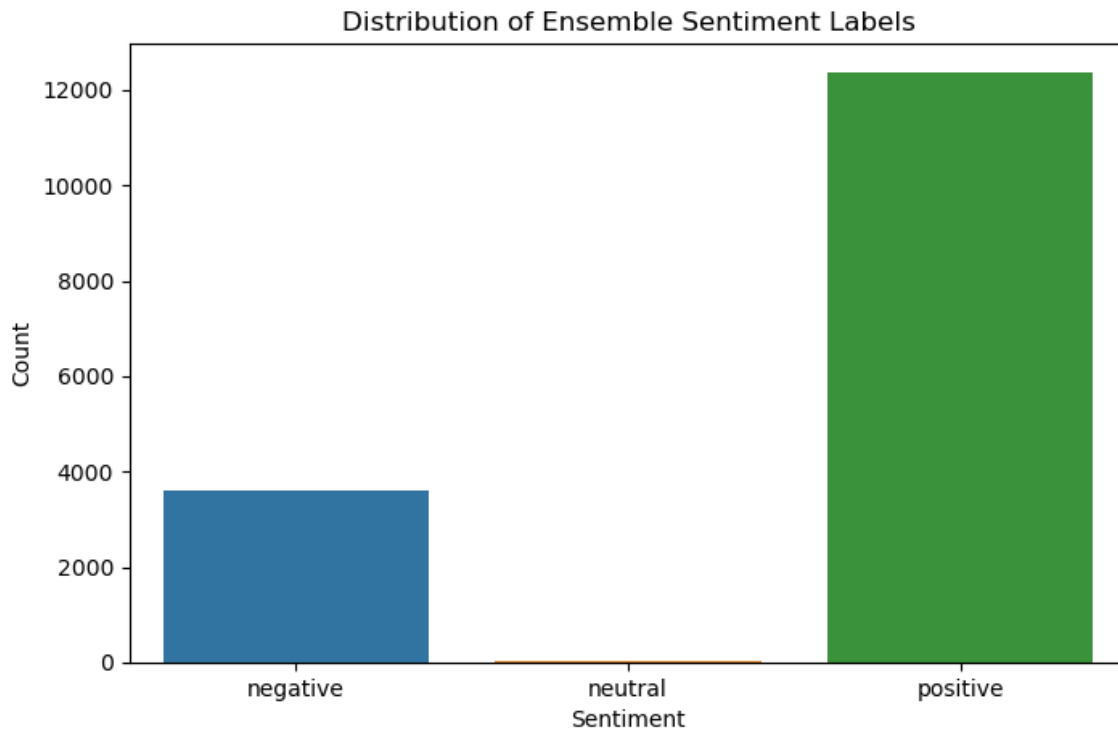
A word cloud was generated from the cleaned text, highlighting frequent terms such as "token", "model", "text", and "bert".



4.2 Sentiment Distribution

VADER sentiment scores showed that:

- Positive: 77%
- Negative: 22%
- Neutral: 1%

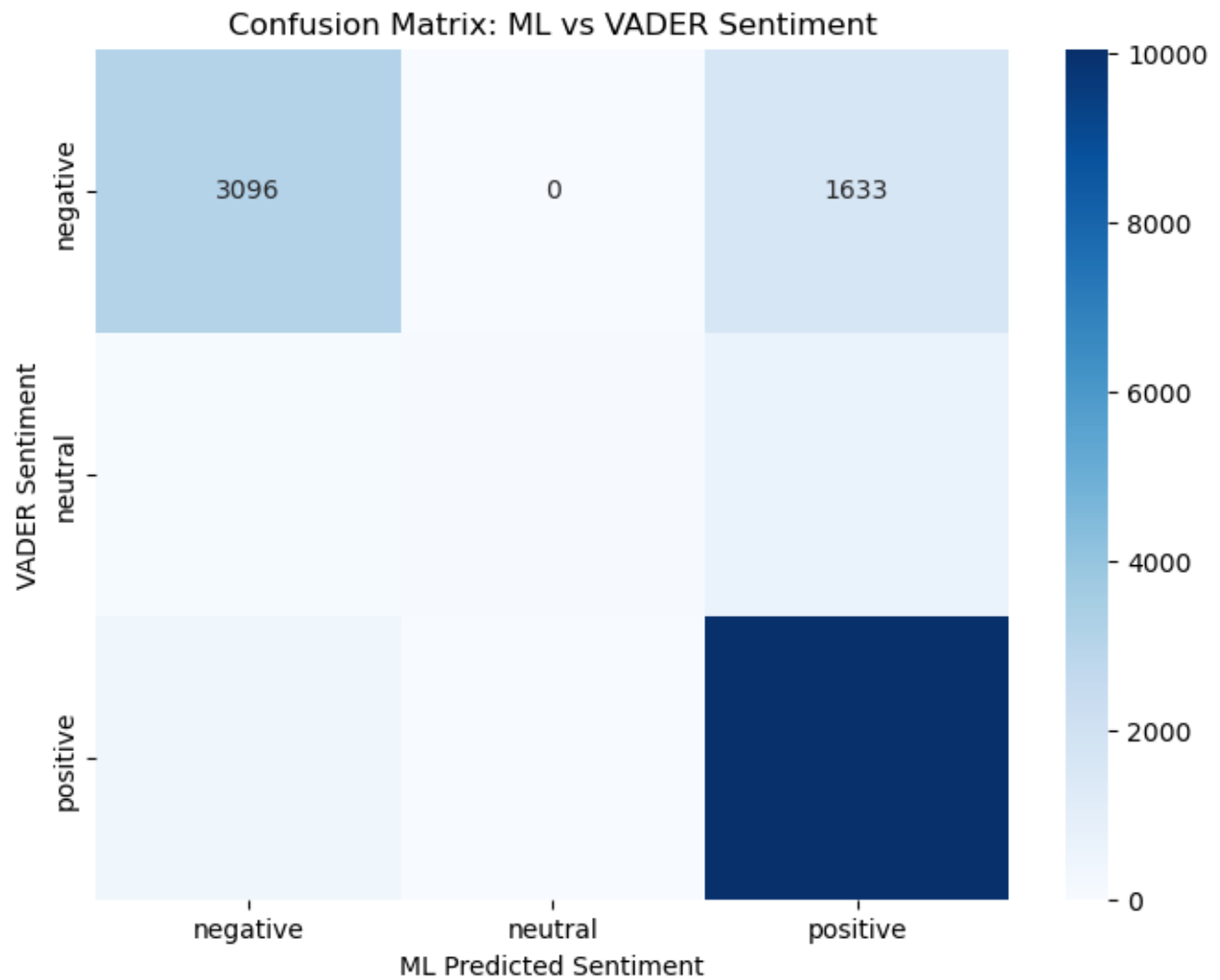


4.3 Topic Distribution

A bar chart was plotted to show the number of posts under each predefined NLP category.

4.4 ML vs VADER Confusion Matrix

A confusion matrix was plotted to compare the sentiment predictions made by the ML model and VADER.



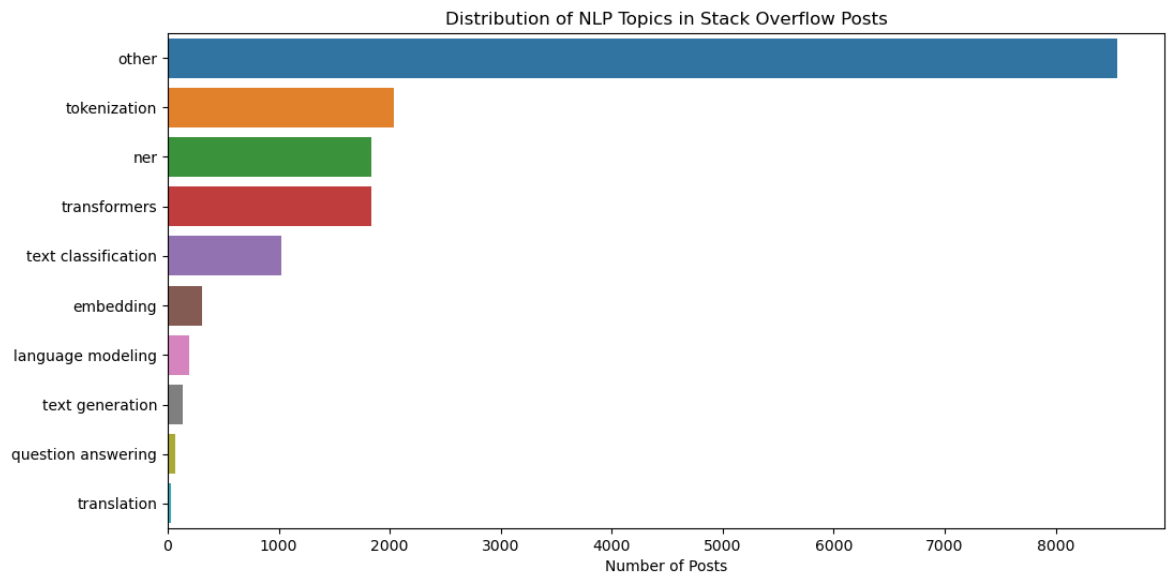
5. Data Categorization

5.1 Definition of Categories

Categories were manually defined based on keyword matches in the post titles:

- tokenization
- transformers
- ner
- language modeling
- text classification
- text generation
- translation
- question answering

- embedding
- other



-

5.2 ML Classification Results

- Accuracy: 73.26%
- Strongest classes: ner, tokenization, other
- Weakest classes (low support): translation, text generation

6. Sentiment Analysis

6.1 VADER vs ML Agreement

- Ensemble agreement rate: **82.4%**
- VADER Sentiment Distribution:
 - Positive: 10,448
 - Negative: 4,729
 - Neutral: 828

6.2 Classification Report (ML vs VADER)

- Positive: F1-score = 0.84
- Negative: F1-score = 0.63
- Neutral: F1-score = 0.02 (sparse)

7. Conclusion

This mini-project successfully implemented an end-to-end NLP pipeline to classify Stack Overflow posts by subfield and sentiment. While keyword-based labeling helped with bootstrapping categories, the classification model showed strong performance despite class imbalance. Sentiment analysis revealed a largely positive tone in the posts. Future improvements may include using BERT embeddings, class balancing techniques, and hierarchical classification.

8. References

- Jurafsky, D., & Martin, J. H. (2023). *Speech and Language Processing*. Stanford University.
- NLTK Documentation: <https://www.nltk.org/>
- Stack Exchange API: <https://api.stackexchange.com/>
- scikit-learn Documentation: <https://scikit-learn.org/>
- Hutto, C., & Gilbert, E. (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis.

GITHUB REPOSITORY LINK:

Please use the below link for redirecting to my github repo which has all the files

https://github.com/Gayathri224/NPL_STACKOVERFLOW.git

This

3. Conclusion

This

4. References

- [1] Afsharizadeh M., Ebrahimpour-Komleh H., Bagheri A 2020. Automatic text summarization of COVID-19 research articles using recurrent neural networks and coreference resolution. *Frontiers in Biomedical Technologies* 7, 236–248. doi:10.18502/fbt.v7i4.5321
- [2] Askdata. (2021). Train T5 for Text Summarization. Available at: <https://medium.com/askdata/train-t5-for-text-summarization-a1926f52d281>
- [3] Beltagy I, Lo K & Cohan A 2019, ‘SciBERT: A Pretrained Language Model for Scientific Text’, arXiv.org.
- [4] Cai X, Liu S, Yang L, Lu Y, Zhao J, She, D & Liu T 2022, ‘COVIDSum: A linguistically enriched SciBERT-based summarization model for COVID-19 scientific papers’, *Journal of Biomedical Informatics*, vol. 127, pp. 103999–103999.
- [5] Gasic, M. (2019). Text Summarization Part 2-State of the Art. Besedo Engineering. Available at: <https://medium.com/besedo-engineering/text-summarization-part-2-state-of-the-art-ae900e2ac55f>
- [6] Gu Y, Tinn R, Cheng H, Lucas M, Usuyama N, Liu X, ... Poon H 2022, ‘Domain-Specific Language Model Pretraining for Biomedical Natural Language Processing’, *ACM Transactions on Computing for Healthcare*, vol. 3, no. 1, pp. 1–23.
- [7] Gupta, A. (2021). Simple Abstractive Text Summarization with Pretrained T5 Text-to-Text Transfer Transformer. *Towards Data Science*. Available at: <https://towardsdatascience.com/simple-abstractive-text-summarization-with-pretrained-t5-text-to-text-transfer-transformer-10f6d602c426#:~:text=T5%20is%20an%20abstractive%20summarization,directly%20from%20the%20original%20text>
- [8] Jurafsky D, Martin J.H, 2023, ‘Speech and Language Processing’. Available at: https://web.stanford.edu/~jurafsky/slp3/ed3book_jan72023.pdf
- [9] T5 for Text Summarization in 7 Lines of Code. Artificialis. Available at: <https://medium.com/artificialis/t5-for-text-summarization-in-7-lines-of-code-b665c9e40771>