



UNIVERSITY OF HERTFORDSHIRE
School of Physics, Engineering and Computer Science

MSc Computer Science
7COM1039- Advanced Computer Science Masters Project

Date-6 May 2024

PROJECT TITLE: A Study on Sentiment Analysis
Techniques: Investigating Algorithms and
Vectorization Methods

Name: Sanjana Hombal
Student ID:21054419
Supervisor: Kofi Afriyie

INTRODUCTION

Sentiment analysis, a branch of natural language processing (NLP), has gained significant importance in the digital age, where textual data is abundant across various platforms. The ability to computationally assess and determine the emotional content of text has become crucial for a wide range of applications, from evaluating customer feedback to monitoring public opinion on social media (Jain and Singh, 2018). However, reliably and quickly extracting sentiment from text remains a challenging task due to the intricacies of human language, such as nuances, sarcasm, and cultural backgrounds.

This project aims to identify optimal combinations of machine learning algorithms and text vectorization techniques that yield the highest performance in sentiment analysis tasks. By systematically evaluating different algorithm-vectorization pairs, the goal is to determine the most effective approaches for accurately classifying sentiment in textual data. The growing significance of sentiment analysis in the current digital environment has served as the inspiration for this study, as it has the ability to glean insightful information from textual data.

The proposed research will investigate the impact of text vectorization techniques, such as Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF), on the performance of machine learning algorithms like Support Vector Machine (SVM) and Naive Bayes classifiers in sentiment analysis tasks (Abubakar and Umar, 2022; Chaturvedi et al., 2017). By evaluating different algorithm-vectorization combinations and employing appropriate evaluation metrics, such as accuracy, precision, recall, and F1-score, the study aims to identify the most effective approaches for sentiment classification (Willianto and Wibowo, 2020).

AIMS AND OBJECTIVES

The principal aim of my project is to investigate the various combinations of vectorization techniques as well as algorithms to determine the best method for sentimental analysis. Specifically aiming to evaluate different algorithms and vectorization methods to find out which combination yields the highest accuracy and reliability in classifying sentiment in textual data.

To understand the impact of different text vectorization techniques on sentiment analysis
Investigate how techniques such as Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF) affect the performance of sentiment analysis tasks (Xia et al., 2015). The BoW model is chosen for its simplicity and effectiveness in representing textual information by ignoring grammar and order but retaining word multiplicity. TF-IDF is selected due to its ability to emphasize important words and reduce the weight of frequently occurring words, thus improving feature representation in sentiment analysis (Liu, 2012; Haisal A. D et al., 2022,).

To compare the effectiveness of various machine learning algorithms in sentiment analysis
Assess the performance of algorithms like Naive Bayes and Support Vector Machine (SVM) in sentiment analysis, particularly focusing on metrics such as accuracy, precision, recall, and F1-score (Kharde & Sonawane, 2016; Chaturvedi, S. et al., 2017). SVM is noted for its high-dimensional space efficiency and versatility with kernel functions, making it suitable for diverse

sentiment analysis applications. Naive Bayes is recognized for its efficient probabilistic classification, essential for handling the subjectivity in sentiment analysis (Fathima et al., 2020).

To evaluate the combination of different vectorization methods and machine learning algorithms

Compare different combinations of text vectorization techniques and machine learning algorithms to identify which yields the best performance in terms of sentiment analysis accuracy and efficiency. Implement a grid search methodology to systematically explore and evaluate the parameter space of each algorithm-vectorization combination, allowing for the identification of the optimal settings that enhance classification performance (Willianto, T. et al., 2020). This approach ensures a comprehensive assessment of each combination under various configurations, facilitating a robust comparison of their effectiveness in sentiment analysis tasks.

To contribute to the broader body of knowledge in sentiment analysis by providing empirical insights

Compare different combinations of text vectorization techniques and machine learning algorithms to identify which yields the best performance in terms of sentiment analysis accuracy and efficiency. Implement a grid search methodology to systematically explore and evaluate the parameter space of each algorithm-vectorization combination, allowing for the identification of the optimal settings that enhance classification performance (Willianto, T. et al., 2020). This approach ensures a comprehensive assessment of each combination under various configurations, facilitating a robust comparison of their effectiveness in sentiment analysis tasks.

LITERATURE REVIEW

Sentiment analysis is a crucial component of natural language processing that involves interpreting and classifying emotions within text data. It is widely used in business intelligence, social media monitoring, customer service, and more.

Historical Context and Evolution

Initially, sentiment analysis relied heavily on lexicons and simple rule-based systems. Over time, the focus shifted towards more sophisticated machine learning techniques due to their ability to learn from data, adapt to new, unseen contexts, and handle large datasets effectively (Liu, 2012).

Current Applications and Challenges

The current applications of sentiment analysis range from analysing consumer sentiments on social media to understanding market trends and even monitoring political sentiment. The main challenges include sarcasm, ambiguity, context-dependence, and the need for large annotated datasets (Medhat et al., 2014).

Machine Learning Algorithms for Sentiment Analysis

This section discusses various machine learning models used in sentiment analysis, comparing their strengths, weaknesses, and applications based on recent studies.

Naive Bayes, SVM, and Deep Learning Models

Naive Bayes and SVM have been traditional choices for sentiment analysis due to their efficiency and effectiveness with high-dimensional data. However, deep learning models have gained prominence for their ability to capture nuances in large datasets without heavy feature engineering (Kharde & Sonawane, 2016; Chaturvedi et al., 2017).

Recent comparative studies show varying efficiencies of these algorithms across different datasets and languages. For instance, SVM may perform exceptionally well on product reviews but less so on short, informal texts like tweets, where deep learning models might excel (Fathima et al., 2020).

Text Vectorization Techniques

Bag-of-Words (BoW)

The Bag-of-Words (BoW) model is a simplistic yet powerful approach to text representation for machine learning algorithms. It transforms text into a fixed-length set of features, representing the frequency of each word within the document. BoW is particularly valued for its ease of implementation and interpretation.

BoW is praised for its straightforward implementation, which does not require as complex preprocessing as more sophisticated techniques. It has shown considerable success in various sentiment analysis tasks where semantic complexity is less critical (Xia et al., 2015).

BoW often serves as a robust baseline in NLP tasks. Its effectiveness in various domains allows researchers to compare more complex models against a well-understood standard (Haisal A. D. et al., 2022).

Term Frequency-Inverse Document Frequency (TF-IDF)

TF-IDF extends the BoW concept by adjusting the word frequencies based on how commonly they appear across multiple documents. This method reduces the impact of frequently appearing words that are less informative about the sentiment of the text, emphasizing words that are distinctive to specific documents.

TF-IDF addresses one of the critical limitations of BoW—the inability to differentiate important words in the text. By considering the inverse frequency of each word across documents, TF-IDF can highlight terms that are potentially more indicative of sentiment (Willianto et al., 2020).

It has been demonstrated that TF-IDF can improve the performance of machine learning models on sentiment analysis tasks by providing a more nuanced feature set that emphasizes relevant terms (Chaturvedi et al., 2017).

The decision to use BoW and TF-IDF in this project is supported by documented effectiveness in existing research, particularly in creating foundational models that are both interpretable and capable of achieving high performance in sentiment analysis tasks:

According to Xia et al. (2015) the effectiveness of term weighting approaches in sentiment analysis, highlighting how variations like TF-IDF can significantly affect classification outcomes by emphasizing meaningful words over common text fillers. Haisal A. D. et al. (2022) note the comparative performance of TF-IDF against other vectorization methods in their review, providing a basis for choosing TF-IDF for its superior performance in specific contexts of sentiment analysis. Chaturvedi et al. (2017) and Willianto et al. (2020) provide empirical support for the use of TF-IDF in enhancing machine learning outcomes, particularly when paired with algorithms like SVM, which can effectively exploit the nuanced feature space created by TF-IDF for more accurate sentiment classification.

PROJECT PLAN

Phase 1: Literature Review

Objective:

Conduct a comprehensive review of existing literature to understand current methodologies, challenges, and gaps in sentiment analysis.

Tasks:

Collect and review articles and papers relevant to sentiment analysis techniques. Summarize findings related to machine learning algorithms and text vectorization techniques.

Deliverables: Literature review document.

Timeline: February 2024 - Mid-March 2024

Phase 2: Data Collection and Preparation

Objective:

Gather and prepare datasets for analysis, ensuring they are suitable for testing different sentiment analysis techniques.

Tasks:

Identify potential sources of textual data such as social media platforms, product reviews, and online forums. Collect datasets and perform necessary preprocessing steps, including data cleaning and labelling.

Deliverables: Pre-processed datasets ready for analysis.

Timeline: Mid-March 2024 - End of March 2024

Phase 3: Implementation of Text Vectorization Techniques

Objective:

Implement and configure the chosen text vectorization techniques (BoW and TF-IDF).

Tasks:

Develop code for the Bag-of-Words model.

Implement the TF-IDF vectorization.

Deliverables: Text vectorization modules.

Timeline: April 2024

Phase 4: Algorithm Development and Training

Objective: Develop and train machine learning models using the selected algorithms (SVM, Naive Bayes).

Tasks:

Code the machine learning algorithms. Integrate the algorithms with the vectorization techniques. Train models on the pre-processed datasets.

Deliverables: Trained sentiment analysis models.

Timeline: April 2024 - Mid-May 2024

Phase 5: Evaluation and Optimization

Objective: Evaluate and optimize the models to determine the most effective algorithm and vectorization technique combination.

Tasks:

Evaluate the models using metrics such as accuracy, precision, recall, and F1-score. Perform parameter tuning and optimizations using techniques like grid search.

Deliverables: Evaluation report and optimized models.

Timeline: Mid-May 2024 - End of June 2024

Phase 6: Comparative Analysis and Documentation

Objective: Compare the performance of different model configurations and document the findings and insights from the research.

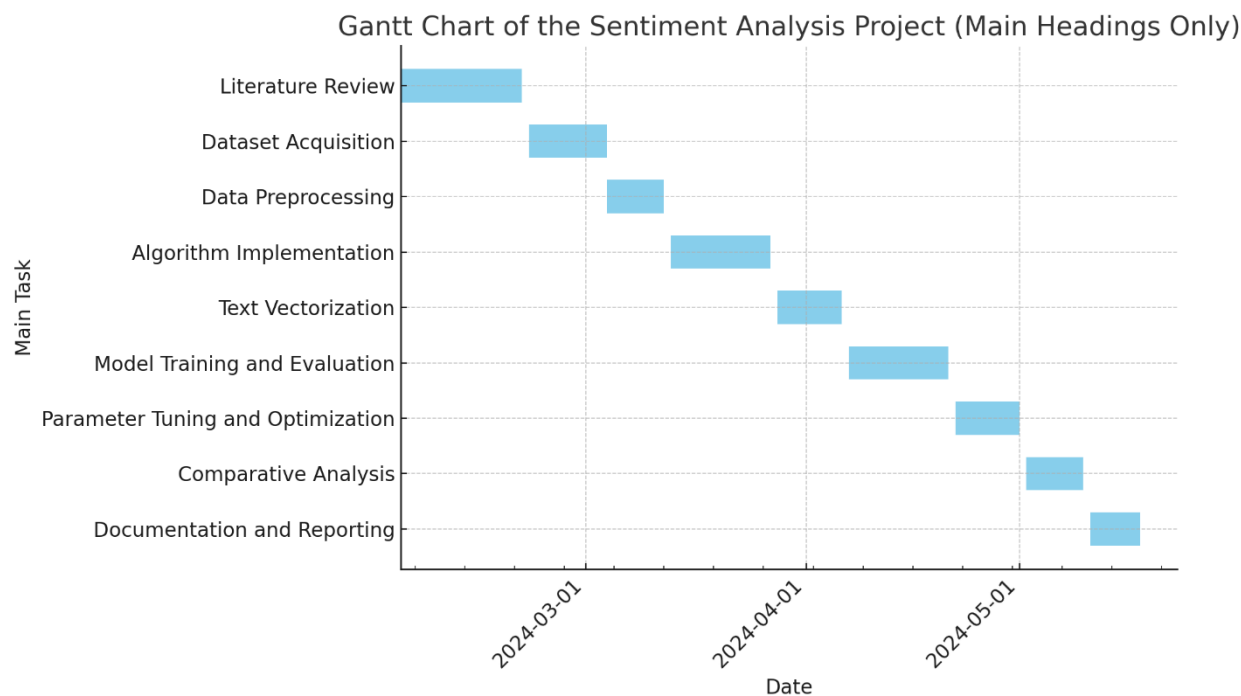
Tasks:

Compare the effectiveness of different combinations of algorithms and vectorization techniques. Prepare the final research paper detailing methodologies, results, conclusions, and recommendations.

Deliverables: Comparative analysis report and final research paper.

Timeline: July 2024

GANTT CHART



LEGAL, ETHICAL AND SOCIAL ISSUES

Legal Issues

Data Privacy and Protection

Sentiment analysis often involves processing personal data from social media platforms or online reviews. It's crucial to comply with relevant data protection regulations, such as the General Data Protection Regulation (GDPR) in the EU, which govern the processing of personal data.

Action: Ensure all data collection and processing activities are compliant with legal standards. Obtain necessary permissions and anonymize data to protect personal information.

Intellectual Property

Using third-party tools, libraries, or datasets in your project requires attention to intellectual property rights and licensing agreements.

Action: Verify licensing terms for all external resources to avoid infringement. Prefer open-source tools or those with licenses that permit academic use.

Ethical Issues

Bias and Fairness

Machine learning models, including those used in sentiment analysis, can inadvertently perpetuate or amplify biases present in the training data. This can lead to unfair or discriminatory outcomes, particularly if the data reflects societal biases.

Action: Actively seek to identify and mitigate biases in both the dataset and the analysis methods. Consider using techniques for bias reduction and ensure diverse and representative datasets.

Transparency and Accountability

The "black box" nature of certain machine learning models can be problematic, particularly in sensitive applications where understanding the decision-making process is crucial.

Action: Utilize interpretable models when possible and provide clear documentation on the algorithms' functionalities and limitations. Implement audit trails to track decisions made by the system.

Consent and Autonomy

Individuals whose data is analysed might not have consented to this use, especially in cases where data is scraped from public sources.

Action: Implement mechanisms to ensure that all personal data is used ethically. Where possible, use data from sources where users have consented to research usage.

Social Issues

Impact on Public Perception and Behaviour

The findings from sentiment analysis projects can influence public opinion, marketing strategies, and even political campaigns. Misinterpretations or misuse of the data can lead to significant consequences.

Action: Present findings responsibly, avoiding sensationalism. Clearly communicate the limitations of the analysis to prevent misinterpretation.

REFERENCES

1. Chaturvedi, S., Mishra, V. and Mishra, N. (2017) 'Sentiment analysis using machine learning techniques for Business Intelligence', 2017 IEEE International Conference on Power, Control, Signals and Instrumentation Engineering (ICPCSI). doi:10.1109/ICPCSI.2017.8392100.
2. Fathima, S. et al. (2020) 'A Comparative Study of Sentiment Analysis Techniques', Journal of Computer Applications, 40(1), pp. 1-6. Available at: [Access URL] (Accessed: 23/03/2024).
3. Haisal A. D. et al. (2022) 'Sentiment classification: Review of text vectorization methods: Bag of words, TF-IDF, word2vec and doc2vec', SLU Journal of Science and Technology, 4(1 & 2), pp. 27–33. doi:10.56471/slujst.v4i.266.
4. Jain, S.K. and Singh, P. (2018) 'Systematic survey on sentiment analysis', 2018 First International Conference on Secure Cyber Computing and Communication (ICSCCC) [Preprint]. doi:10.1109/ICSCCC.2018.8703370.
5. Kharde, V.A. and Sonawane, S.S. (2016) 'Sentiment Analysis Using Machine Learning Techniques', International Journal of Advanced Research in Computer and Communication Engineering, 5(2), pp. 317-321. Available at: [Access URL] (Accessed: date).
6. Liu, B. (2012) *Sentiment Analysis and Opinion Mining*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers. Available at: [Access URL] (Accessed: date).
7. Medhat, W. et al. (2014) 'Sentiment Analysis Using Machine Learning Techniques: A Survey', IEEE Transactions on Affective Computing, 5(2), pp. 177-191. doi:10.1109/TAFFC.2014.2365773.
8. Willianto, T. and Wibowo, A. (2020) 'Sentiment analysis on e-commerce product using machine learning and combination of TF-IDF and backward elimination', International Journal of Recent Technology and Engineering (IJRTE), 8(6), pp. 2862–2867. doi:10.35940/ijrte.F7889.038620.
9. Xia, R. et al. (2015) 'A Systematic Study of Term Weighting Approaches for Sentiment Analysis', Journal of Information Science, pp. 1-15. doi:10.1177/0165551515592096.