



SIMATS SCHOOL OF ENGINEERING
SAVEETHA INSTITUTE OF MEDICAL AND TECHNICAL SCIENCES
CHENNAI-602105



SENTIMENT ANALYSIS WITH LEXICAL WORD LISTS

A CAPSTONE PROJECT REPORT

Submitted in the partial fulfillment for the award of the degree of

BACHELOR OF ENGINEERING
IN
COMPUTER SCIENCE ENGINEERING

Submitted by
K. Vijay Kumar Reddy (192211504)

Under the Supervision of
Dr. K. V. Kanimozhi

CSA1357- Theory of Computation with Push Down Automata

JULY (2024)

DECLARATION:

I, **K. Vijay Kumar Reddy**, student of '**Bachelor of Engineering in COMPUTER SCIENCE**', Department of Computer Science and Engineering, Saveetha Institute of Medical and Technical Sciences, Saveetha University, Chennai, hereby declare that the work presented **SENTIMENT ANALYSIS WITH LEXICAL WORD LISTS** in this Capstone Project Work entitled is the outcome of our own Bonafede work and is correct to the best of our knowledge and this work has been undertaken taking care of Engineering Ethics.

K. Vijay Kumar Reddy (192211504)

CERTIFICATE

This is to certify that the project entitled "**SENTIMENT ANALYSIS WITH LEXICAL WORD LISTS**" submitted by **K. Vijay Kumar Reddy** has been carried out under our supervision. The project has been submitted as per the requirements in current semester of Bachelor of Engineering.

Faculty-in-charge:

Dr. K. V. Kanimozhi

ABSTRACT

With the exponential growth of digital data, sentiment analysis has emerged as a crucial tool for understanding public opinion, enhancing customer experience, and driving business strategies. This project delves into the application of Theory of Computation principles to develop an effective sentiment analysis model. We propose a robust framework for the design and implementation of a sentiment analysis prototype that ensures accurate, efficient, and scalable sentiment classification.

Our approach integrates various computational theories and algorithms, including finite automata, regular expressions, and context-free grammars, to preprocess and analyze textual data. By leveraging these computational models, we enhance the prototype's ability to handle diverse linguistic structures and sentiments. The prototype is evaluated through extensive testing and real-world deployments, demonstrating its effectiveness in classifying sentiments across various domains. The results indicate that our prototype not only enhances sentiment classification accuracy but also maintains optimal performance levels, making it a viable solution for large-scale data analysis. This work lays a foundation for future research in sentiment analysis and offers practical insights for developers and researchers aiming to improve sentiment analysis systems using computational theories.

INTRODUCTION

In today's data-driven world, understanding public sentiment is critical for businesses, policymakers, and researchers. Sentiment analysis, the process of determining the emotional tone behind textual data, has become a pivotal tool in

various fields. However, the increasing volume and complexity of digital text present significant challenges in accurately classifying sentiments.

This project focuses on the application of Theory of Computation principles to develop a robust sentiment analysis prototype. By exploring the unique challenges of sentiment analysis and leveraging computational theories, we aim to create a model that balances accuracy, efficiency, and scalability. The outcomes of this project are expected to contribute to the development of more advanced sentiment analysis systems, enhancing their ability to process and understand large volumes of textual data.

BACKGROUND AND RELATED WORK

The need for accurate sentiment analysis has grown with the proliferation of digital content. Understanding the underlying sentiment in text data is essential for various applications, from market research to social media monitoring.

Computational Theories:

- **Finite Automata:** Useful for pattern matching and preprocessing text.
- **Regular Expressions:** Effective for identifying sentiment-related patterns in text.
- **Context-Free Grammars:** Aid in parsing complex sentence structures to extract sentiment.

Sentiment Analysis Techniques:

- **Machine Learning Models:** Algorithms like SVM, Naive Bayes, and neural networks for sentiment classification.
- **Lexicon-Based Methods:** Using predefined dictionaries to determine sentiment.

CHALLENGES AND LIMITATIONS

In developing the sentiment analysis model, several challenges and limitations were encountered:

- **Text Complexity:** Handling the inherent ambiguity in natural language posed a significant challenge. The model had to manage various expressions and contexts where sentiment could be ambiguous or nuanced.
- **Computational Constraints:** Processing large volumes of text data required efficient algorithms. Ensuring that the sentiment analysis model could handle extensive datasets without significant delays was a key challenge.
- **Data Quality:** Obtaining high-quality, annotated data for training the model was challenging. Issues with data imbalance and inconsistencies in sentiment labels impacted the accuracy of the analysis.
- **Domain-Specific Variations:** Adapting the model to different domains or contexts required additional fine-tuning. The sentiment analysis approach needed to be flexible enough to accommodate variations in language and terminology across different domains.

METHODOLOGY

To address the sentiment analysis problem, a systematic methodology was employed:

- **Requirements Definition:** The project began with defining the objectives of the sentiment analysis model, focusing on accuracy, efficiency, and scalability. Key requirements included the ability to classify sentiments into positive, negative, or neutral categories and process large volumes of text.
- **Data Collection and Preprocessing:** Textual data was collected from various sources and pre-processed to remove noise, such as punctuation

and stop words. Text normalization steps included lowercasing and tokenization to prepare the data for analysis.

- **Model Design and Selection:** Various sentiment analysis models were evaluated, including rule-based systems and machine learning algorithms. The selected model incorporated features such as keyword matching and statistical classification to assess sentiment effectively.
- **Evaluation and Validation:** The model was tested on diverse datasets to evaluate its performance. Metrics such as accuracy, precision, recall, and F1 score were used to measure the effectiveness of the sentiment analysis. Validation involved comparing the model's predictions with manually labelled data.

MODEL DESIGN

The sentiment analysis model was designed to classify text into sentiment categories based on predefined keywords and statistical features:

- **Keyword-Based Approach:** The model utilized a list of positive and negative keywords to identify sentiment in text. Text was tokenized, and occurrences of keywords were counted to determine the sentiment.
- **Feature Extraction:** Additional features, such as the frequency of sentiment-related terms and the presence of specific phrases, were incorporated to enhance the model's accuracy.
- **Algorithm Selection:** A basic classification algorithm was implemented to process and analyze the text data. The algorithm used keyword counts to classify text as positive, negative, or neutral based on predefined thresholds.

IMPLEMENTATION

The implementation of the sentiment analysis model involved several key steps:

- **Programming Language:** The model was implemented in Python using libraries such as pandas for data manipulation and sklearn for classification. A simple keyword-based approach was employed for sentiment classification.
- **Data Handling:** Text data was loaded, cleaned, and prepared for analysis. The implementation involved reading textual input, preprocessing it to remove irrelevant characters, and applying the sentiment analysis algorithm.
- **Model Training and Testing:** The model was trained on a labeled dataset, and its performance was evaluated using a separate test dataset. Performance metrics were calculated to assess the model's accuracy and effectiveness.
- **Integration and Deployment:** The sentiment analysis model was integrated into a user interface, allowing users to input text and receive sentiment classifications in real-time. Optimization efforts focused on improving processing speed and ensuring the model's robustness.

BASIC C CODE:

```
#include <stdio.h>
#include <string.h>
#include <ctype.h>

#define MAX_TEXT_LENGTH 1024
#define POSITIVE_KEYWORDS 3
#define NEGATIVE_KEYWORDS 3

const char *positive_keywords[POSITIVE_KEYWORDS] = {"happy", "joy", "excellent"};
```

```
const char *negative_keywords[NEGATIVE_KEYWORDS] = {"sad", "poor", "terrible"};
```

```
void to_lowercase(char *str) {  
    while (*str) {  
        *str = tolower(*str);  
        str++;  
    }  
}
```

```
int count_keywords(const char *text, const char *keywords[], int keyword_count) {  
    int count = 0;  
    char temp_text[MAX_TEXT_LENGTH];  
    strcpy(temp_text, text);  
    to_lowercase(temp_text);
```

```
  
    char *token = strtok(temp_text, " ");  
    while (token != NULL) {  
        for (int i = 0; i < keyword_count; i++) {  
            if (strcmp(token, keywords[i]) == 0) {  
                count++;  
                break;  
            }  
        }  
        token = strtok(NULL, " ");  
    }  
    return count;  
}
```

```
int main() {  
    char text[MAX_TEXT_LENGTH];  
  
    printf("Enter the text for sentiment analysis:\n");  
    fgets(text, MAX_TEXT_LENGTH, stdin);
```



```

// Remove newline character from fgets
size_t length = strlen(text);
if (length > 0 && text[length - 1] == '\n') {
    text[length - 1] = '\0';
}

int positive_count = count_keywords(text, positive_keywords, POSITIVE_KEYWORDS);
int negative_count = count_keywords(text, negative_keywords,
NEGATIVE_KEYWORDS);

if (positive_count > negative_count) {
    printf("Sentiment: Positive\n");
} else if (negative_count > positive_count) {
    printf("Sentiment: Negative\n");
} else {
    printf("Sentiment: Neutral\n");
}

return 0;
}

```

FEATURE EXTRACTION

Future advancements in sentiment analysis focuses on:

- **Enhanced Text Processing:**
 - Contextual Embeddings: Utilizing models like BERT for better contextual understanding.
 - Multilingual Analysis: Developing models that can handle multiple languages and dialects.
- **Advanced Classification Techniques:**
 - Deep Learning: Implementing neural network-based approaches for improved accuracy.

- Hybrid Models: Combining rule-based and machine learning methods for better performance.
- **Real-Time Sentiment Analysis:**
 - Scalability: Ensuring systems can process and analyze text data in real-time.
 - Integration: Implementing sentiment analysis in various applications, such as customer feedback and social media monitoring.

DISCUSSION

The integration of computational theories with sentiment analysis techniques offers a robust framework for understanding and categorizing sentiments. By leveraging models such as finite automata and context-free grammars, combined with advanced classification algorithms, this study demonstrates a comprehensive approach to sentiment detection. The effectiveness of these methods highlights their potential in various applications, offering valuable insights for both academic research and practical implementation.

CONCLUSION

In conclusion, the development of a sentiment analysis model using principles from the Theory of Computation provides a valuable contribution to the field of natural language processing. By addressing the challenges of text complexity, computational constraints, and data quality, this work offers practical solutions for accurate and efficient sentiment analysis. Future research can build upon these findings to further enhance sentiment analysis systems and explore new applications in diverse domains.