

# Sentiment Analysis Report

## Introduction

The aim of this project is to analyse the sentiment of user reviews using the VADER (Valence Aware Dictionary and Entiment Reasoner) sentiment analysis tool. This involves loading data from a CSV file, cleaning and preprocessing the text data, analysing sentiment, and saving the results and summary to CSV files.

## Steps Involved

### 1. Loading Data

- a. The data is loaded from a CSV file using the **pandas** library.
- b. Error handling is implemented to manage cases where the file might not be found, is empty, or cannot be parsed.

### 2. Cleaning Data

- a. The data is cleaned by removing any null values.
- b. Only the **review** column is retained for analysis to focus on the relevant text data.

### 3. Preprocessing Text

- a. The text is converted to lowercase to ensure uniformity.
- b. Punctuation is removed from the text to standardize it and make it suitable for sentiment analysis.

### 4. Sentiment Analysis Using VADER

- a. The VADER sentiment analysis tool is used to analyze the sentiment of each review.
- b. Sentiments are categorized as **positive**, **negative**, or **neutral** based on the compound sentiment score provided by VADER.

### 5. Saving Results

- a. The processed reviews along with their corresponding sentiment labels are saved to a new CSV file.

- b. A summary report of the sentiment analysis, showing the count of each sentiment category, is generated and saved to another CSV file.

## Challenges Faced

### 1. Data Quality Issues

- a. Missing Values: The dataset contained null values which needed to be handled by removing the corresponding rows to ensure the analysis was accurate.
- b. Irrelevant Columns: The dataset might have included columns not relevant to sentiment analysis. We had to identify and retain only the necessary columns, specifically the 'review' column.

### 2. Text Preprocessing

- a. Inconsistent Text Formats: User reviews were in various formats, requiring standardization through conversion to lowercase and removal of punctuation.
- b. Handling Special Characters and Emojis: The presence of special characters and emojis could affect sentiment analysis. Ensuring these were appropriately handled or removed was crucial.

### 3. Sentiment Analysis Limitations

- a. Context Understanding: VADER, like other sentiment analysis tools, may not fully understand the context or subtleties of the text, such as sarcasm or idiomatic expressions.
- b. Score Thresholds: Determining appropriate thresholds for categorizing sentiments (e.g., positive, negative, neutral) required careful consideration to ensure accurate classification.

## Assumptions Made

### 1. Accuracy of VADER Sentiment Scores

- a. We assumed that the VADER sentiment analysis tool provides reasonably accurate sentiment scores for user reviews. While VADER is well-regarded for its accuracy, especially for social media text, it may still have limitations with certain types of reviews.

### 2. Importance of Text Preprocessing

- a. We assumed that converting text to lowercase and removing punctuation would sufficiently standardize the text for sentiment analysis. This preprocessing step is crucial to reduce variability and improve analysis accuracy.
- 3. Representative Sample
  - a. We assumed that the sample of reviews provided in the CSV file is representative of the broader user sentiment. This assumption is important for the generalizability of the results.

## Detailed Explanation

1. Loading Data
  - a. The data is loaded from a CSV file. This step ensures that the file exists, is not empty, and can be parsed correctly. If any of these conditions fail, an appropriate error message is displayed.
2. Cleaning Data
  - a. The cleaning process involves removing null values and retaining only the 'review' column, which contains the text data to be analyzed.
3. Preprocessing Text
  - a. Preprocessing involves converting the text to lowercase and removing punctuation. This standardizes the text, making it suitable for sentiment analysis.
4. Sentiment Analysis Using VADER
  - a. The VADER sentiment analysis tool is used to determine the sentiment of each review. Reviews are classified into three categories based on their sentiment scores:
    - i. Positive: Reviews with a compound score greater than or equal to 0.05.
    - ii. Negative: Reviews with a compound score less than or equal to -0.05.
    - iii. Neutral: Reviews with a compound score between -0.05 and 0.05.
5. Saving Results
  - a. The processed reviews and their sentiment labels are saved to a CSV file for future reference and analysis. Additionally, a summary report showing the count of each sentiment category is generated and saved to another CSV file.

## Conclusion

This project successfully demonstrates the process of sentiment analysis using the VADER tool. By cleaning and preprocessing the text data, we ensure it is suitable for analysis. The sentiment analysis provides valuable insights into user opinions, which are summarized and saved for further use. This structured approach to analyzing text data helps in understanding customer feedback and improving services or products accordingly. The challenges faced and assumptions made during this process highlight the importance of careful data handling and preprocessing to achieve accurate and meaningful results.