# Social Media Sentiment Analysis Report

**Objective**

The primary goal of this project was to analyze and classify social media text data to identify sentiment trends and determine the key features influencing sentiment classification. This analysis aimed to provide actionable insights to help understand public sentiment and support strategic decision-making.

---

**Key Features**

- **Data Preprocessing:**
    - The text data was cleaned by removing special characters, punctuation, and numbers to ensure that irrelevant information did not impact the analysis.
    - Text was converted to lowercase for consistency.
    - Tokenization was applied to break text into meaningful units, and stopwords were removed to focus on relevant terms.
    - Lemmatization was performed to reduce words to their base forms, ensuring uniformity and reducing redundancy in the dataset.
- **Machine Learning Pipeline:**
    - A pipeline was developed that integrated TF-IDF (Term Frequency-Inverse Document Frequency) vectorization for feature extraction and Random Forest for sentiment classification.
    - TF-IDF was utilized to transform the textual data into numerical features, capturing the importance of terms within the corpus while minimizing the effect of common but less meaningful words.
    - The Random Forest classifier was chosen for its robustness and ability to handle high-dimensional data effectively.
- **Feature Importance:**
    - The pipeline included a feature importance analysis to identify the top 10 terms that significantly influenced sentiment decisions.
    - Key features like "good," "love," and "hate" were visualized to provide insights into their impact on the classification process.

---

**Tools & Technologies**

- **Programming Language:** Python (pandas, scikit-learn, matplotlib, nltk)
- **Feature Extraction:** TF-IDF Vectorization
- **Modeling:** Random Forest Classifier
- **Visualization Tools:** Matplotlib for graphical representations
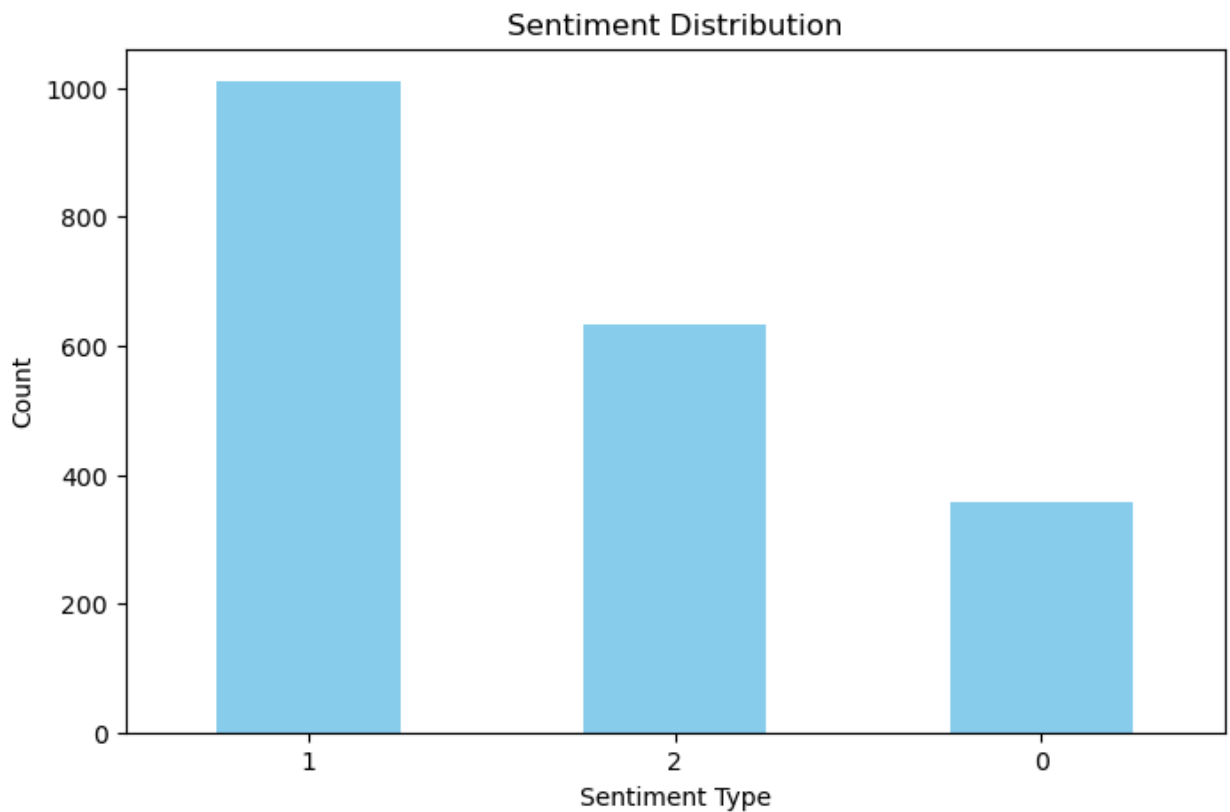
## Results

- **Model Accuracy:** The sentiment classification model achieved an accuracy of **80.75%**, demonstrating a strong ability to correctly predict sentiments across the dataset.
- **Sentiment Insights:** The analysis revealed that positive sentiments were the most prevalent, followed by neutral and negative sentiments. This insight provides valuable information for understanding public opinion trends.
- **Feature Importance:** Terms like "good," "love," and "hate" emerged as the most influential features, reflecting their strong association with sentiment polarity.

## Visualizations

- **Sentiment**                                                          **Distribution:**
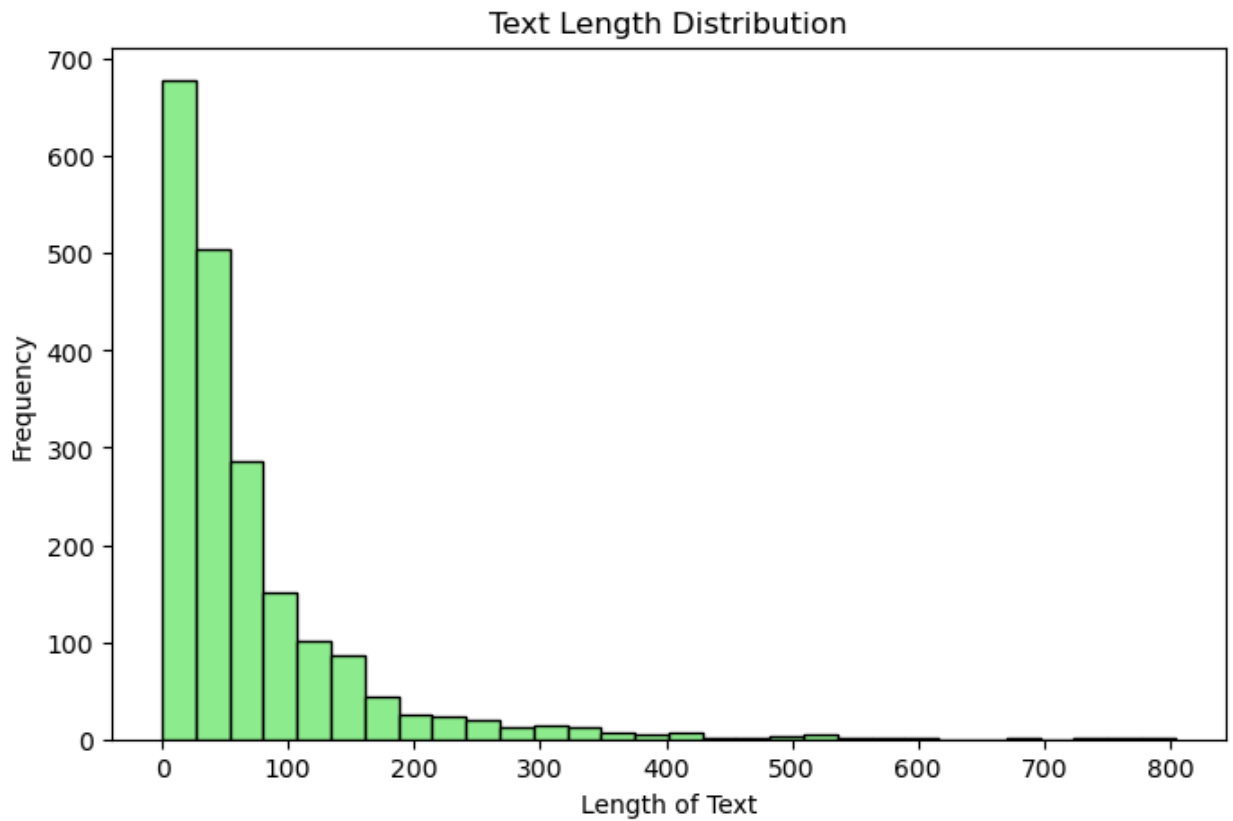


- A bar chart was created to display the counts of positive, neutral, and negative sentiments in the dataset. This provided an overview of sentiment trends.
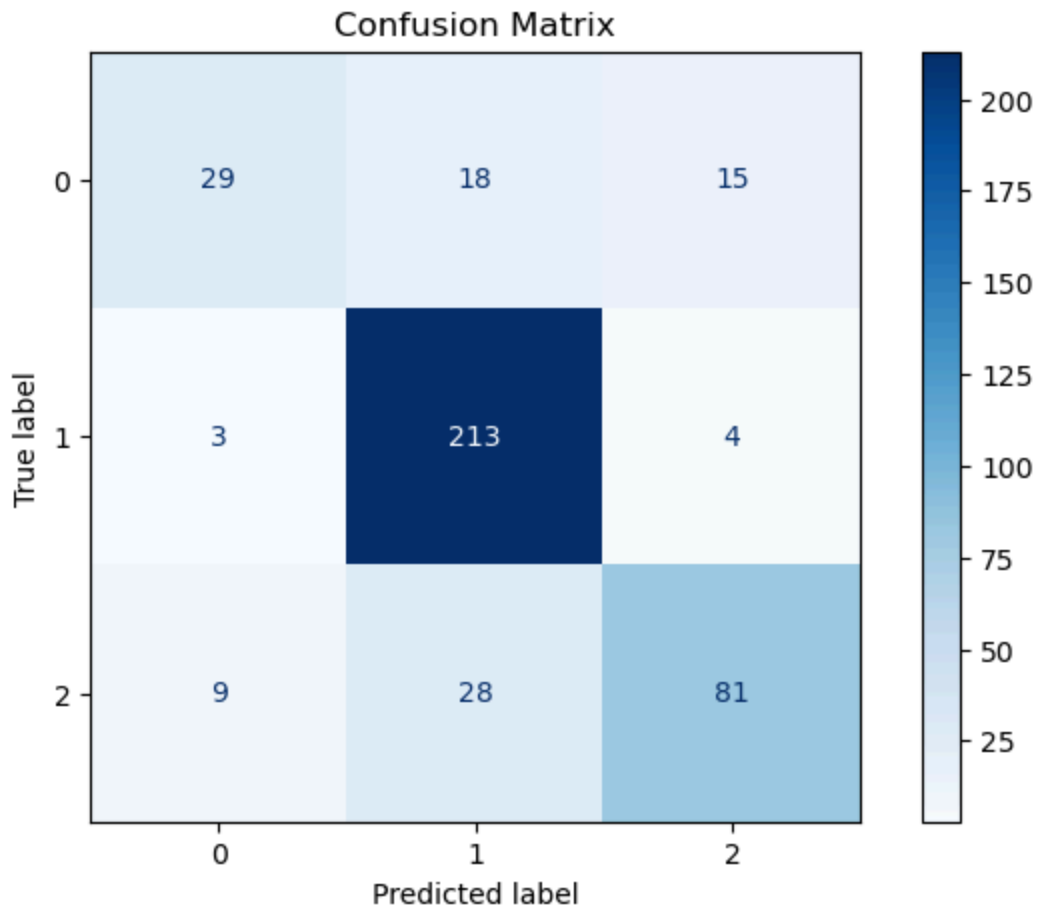
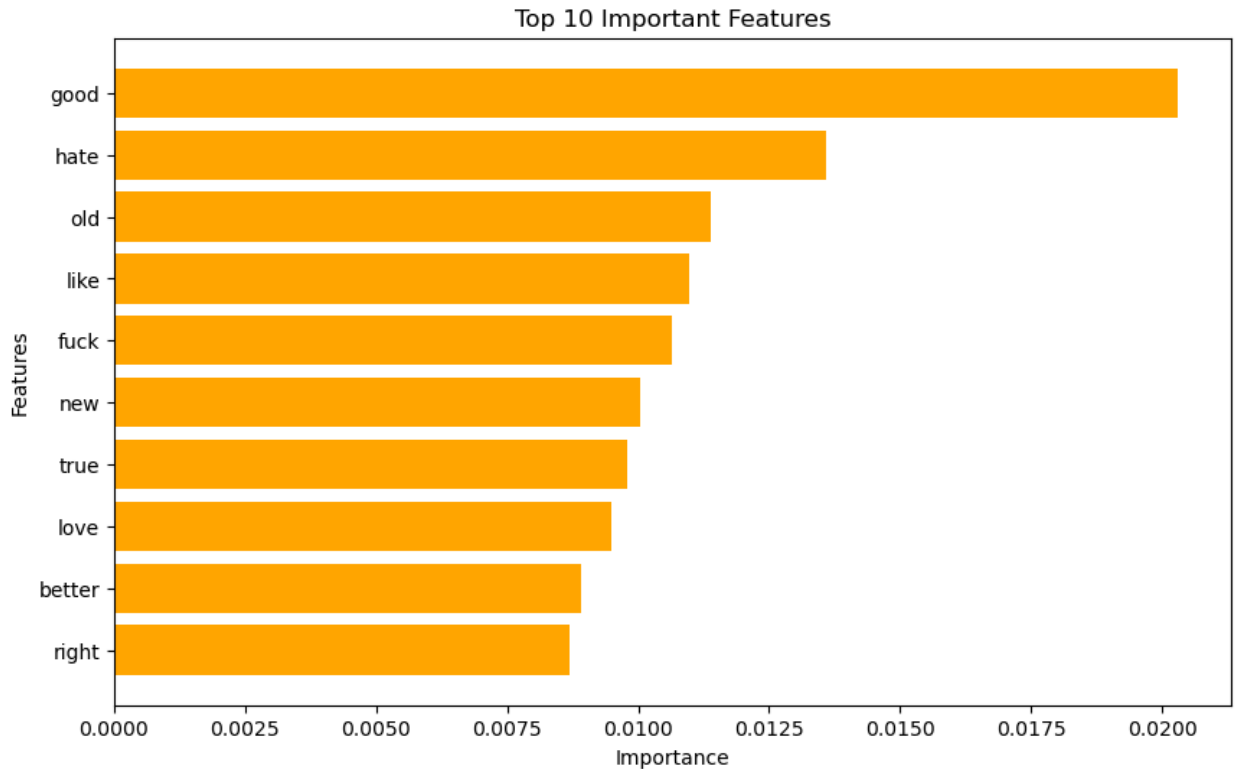● **Text** **Length** **Distribution:**

### Text Length Distribution



○ A histogram was used to analyze the distribution of text lengths, offering insights into the typical message size and its correlation with sentiment.

● **Confusion Matrix:**

Confusion Matrix

- ○ A confusion matrix visualized the performance of the model by showing the relationship between actual and predicted labels, highlighting areas of strength and opportunities for improvement.
- **Top 10 Features:**

Top 10 Important Features

- ○ A bar chart illustrated the top 10 influential terms in the sentiment classification model. This visualization highlighted terms that significantly impacted the model's decisions, aiding interpretability.

**Conclusion**

The Social Media Sentiment Analysis project successfully demonstrated the ability to classify sentiments accurately and identify key factors influencing sentiment trends. The insights provided by this analysis can be leveraged for strategic decision-making in areas such as marketing, customer feedback analysis, and public relations. Future work could explore deep learning approaches for further improvements in classification performance and expand the analysis to include multilingual datasets for broader applicability.