# Hotel Review Sentiment Analysis Report

## Introduction

This report presents the results of a sentiment analysis model developed to evaluate hotel reviews. The model is designed to classify reviews into three sentiment categories: Negative, Neutral, and Positive. It was trained on a dataset containing over 20,000 hotel reviews sourced from TripAdvisor.

## Data Overview

The dataset consists of hotel reviews along with their corresponding ratings. The ratings were used to create sentiment labels as follows:

- **Negative**: Ratings of 0 to 2
- **Neutral**: Ratings of 3
- **Positive**: Ratings of 4 to 5

The reviews were preprocessed to remove noise and enhance the quality of the text data. This included converting text to lowercase, removing special characters and numbers, tokenizing the text, and lemmatizing the words while removing stopwords.

## Model Development

### Methodology

The model employs a Support Vector Classifier (SVC) with a linear kernel for sentiment classification. The text data was vectorized using the Term Frequency-Inverse Document Frequency (TF-IDF) method, which transforms the text into a numerical format suitable for machine learning algorithms.

### Training and Validation

The dataset was split into training and testing sets using an 80-20 split. The model was trained on the training set, and its performance was evaluated using cross-validation and a separate test set. The following metrics were used to assess the model's performance:

- **Cross-validation scores**: The model was evaluated using 5-fold cross-validation, yielding the following scores:
  - [0.8533, 0.8597, 0.8548, 0.8520, 0.8588]

  The average cross-validation score was approximately **0.86**, indicating a strong performance across different subsets of the data.

## Model Performance

The model's performance was further evaluated on the test set, which consisted of 4,099 reviews. The classification report generated from the model's predictions is as follows:

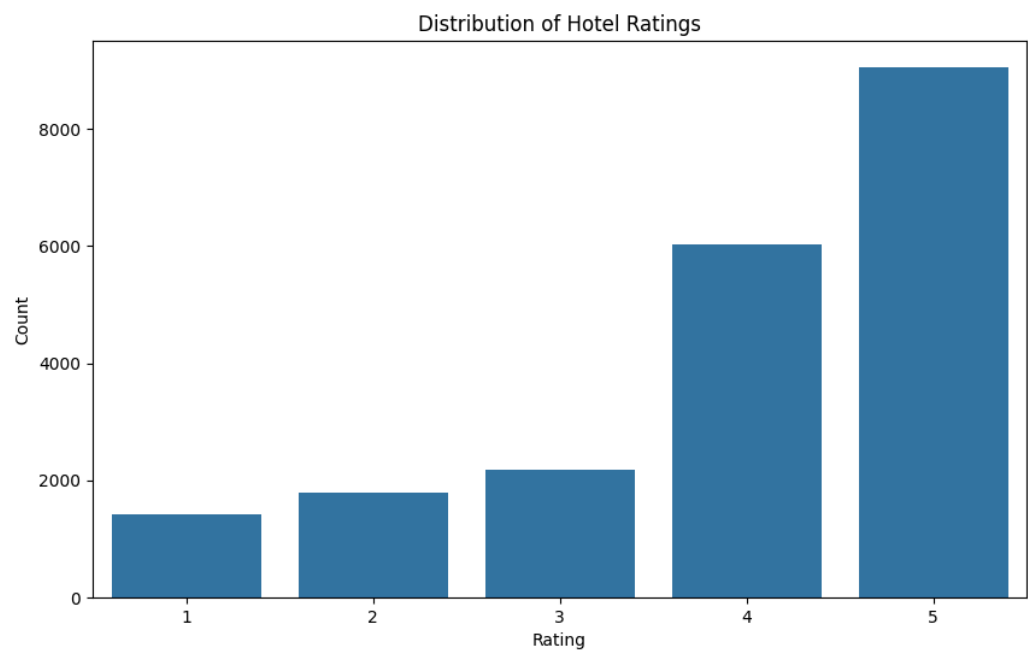| Sentiment | Precision | Recall | F1-Score | Support |
| --- | --- | --- | --- | --- |
| Negative | 0.77 | 0.79 | 0.78 | 625 |
| Neutral | 0.51 | 0.23 | 0.31 | 432 |
| Positive | 0.90 | 0.97 | 0.93 | 3042 |
| Accuracy | | | **0.86** | 4099 |
| Macro Avg | 0.73 | 0.66 | 0.68 | 4099 |
| Weighted Avg | 0.84 | 0.86 | 0.85 | 4099 |

## Key Observations

- The model achieved an overall accuracy of **86%** on the test set.
- The **Positive** sentiment class exhibited the highest precision (0.90) and recall (0.97), indicating that the model is particularly effective at identifying positive reviews.
- The **Negative** sentiment class also performed reasonably well, with a precision of 0.77 and recall of 0.79.
- The **Neutral** sentiment class had lower performance metrics, with a precision of 0.51 and recall of 0.23, suggesting that the model struggles to accurately classify neutral reviews.

# Visualisations

Several visualisations were generated to provide insights into the dataset and model performance:

1. **Rating Distribution**: A count plot illustrating the distribution of hotel ratings.

2. **Word Clouds**: Word clouds for each sentiment category, highlighting the most frequently used words in positive, neutral, and negative reviews.
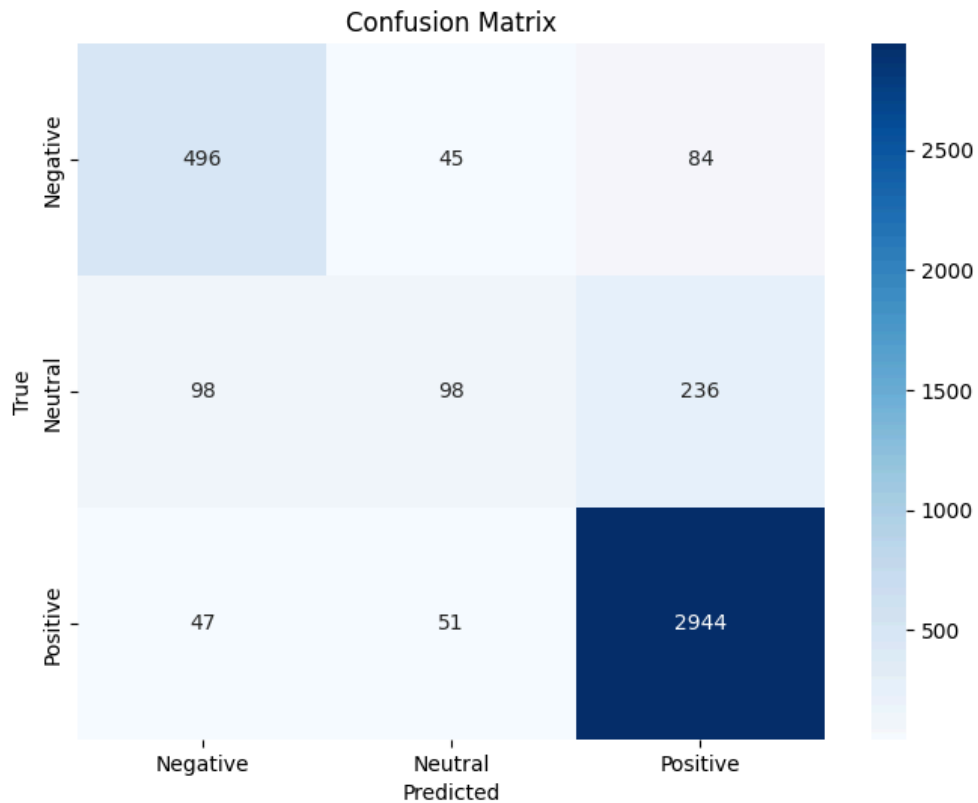
Word Cloud for Negative Reviews



Word Cloud for Neutral Reviews



Word Cloud for Positive Reviews

3. **Confusion Matrix**: A heatmap representing the confusion matrix, which visualises the model's performance across different sentiment classes.



## Conclusion

The sentiment analysis model developed for hotel reviews demonstrates strong performance, particularly in identifying positive sentiments. While the model performs adequately for negative sentiments, there is room for improvement in classifying neutral reviews. Future work may involve exploring additional features, or employing more advanced models to enhance classification accuracy across all sentiment categories.