

Sentiment Analysis Task Report

Overview

The objective of the assignment is to develop an approach to identify subthemes and their respective sentiments from given reviews. The approach needs to extract insights from the reviews, identifying specific aspects (subthemes) and their associated sentiments, and then visualize these results.

Approach

The approach involves multiple steps, from data preprocessing and sentiment analysis to visualization and topic modeling. Each part of the provided code contributes to achieving the objective.

1. Data Loading/Preprocessing - Start the code by loading the dataset with ``pd.read_csv``. A column is removed for cleaning. The data distribution is shown by counting each column's unique values.

2. Sentiment Analysis - Count positive and negative attitudes across all columns. Iterate through each row and check for 'positive' and 'negative' keywords.

3. Visualization - Data is visualized through plots:

- A pie chart reflecting positive evaluations.
- A heatmap showing good reviews in the top 10 categories.
- A line plot of positive review patterns across columns.
- A bar chart showing positive and negative mood.
- Pie and bar charts showing sentiment across rows.

4. Data Cleaning - Lowercase text data in each column, removing punctuation and digits.

5. Subtheme Identification using Pretrained Models VADAR, Roberta and Topic Modeling (LDA)

- Performing sentiment analysis using the VADER model and comparing it with the RoBERTa model helps in quantifying sentiments in the dataset. This step is crucial for determining the sentiment polarity towards each subtheme identified from the reviews.

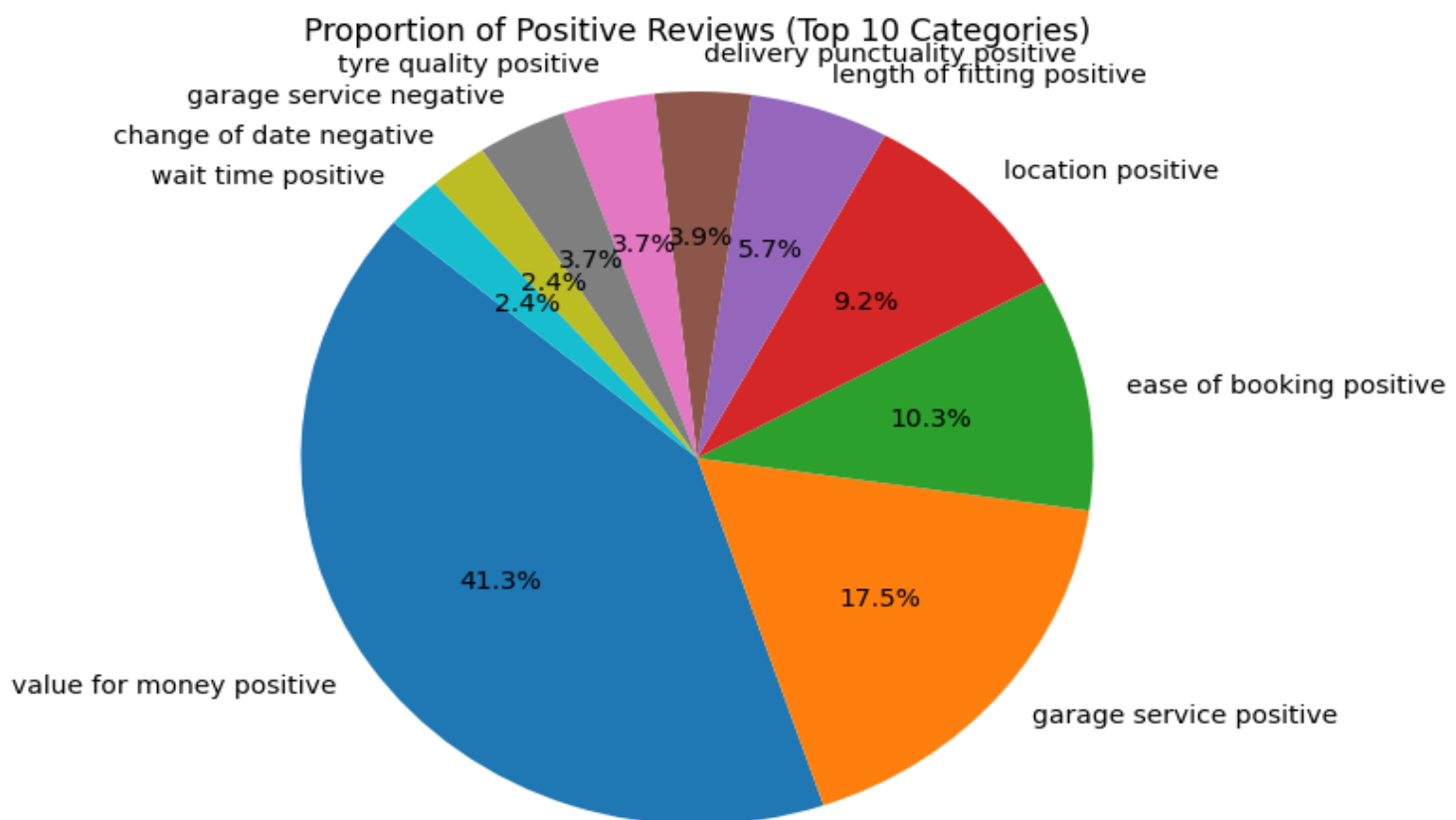
Topic Modelling with LDA:

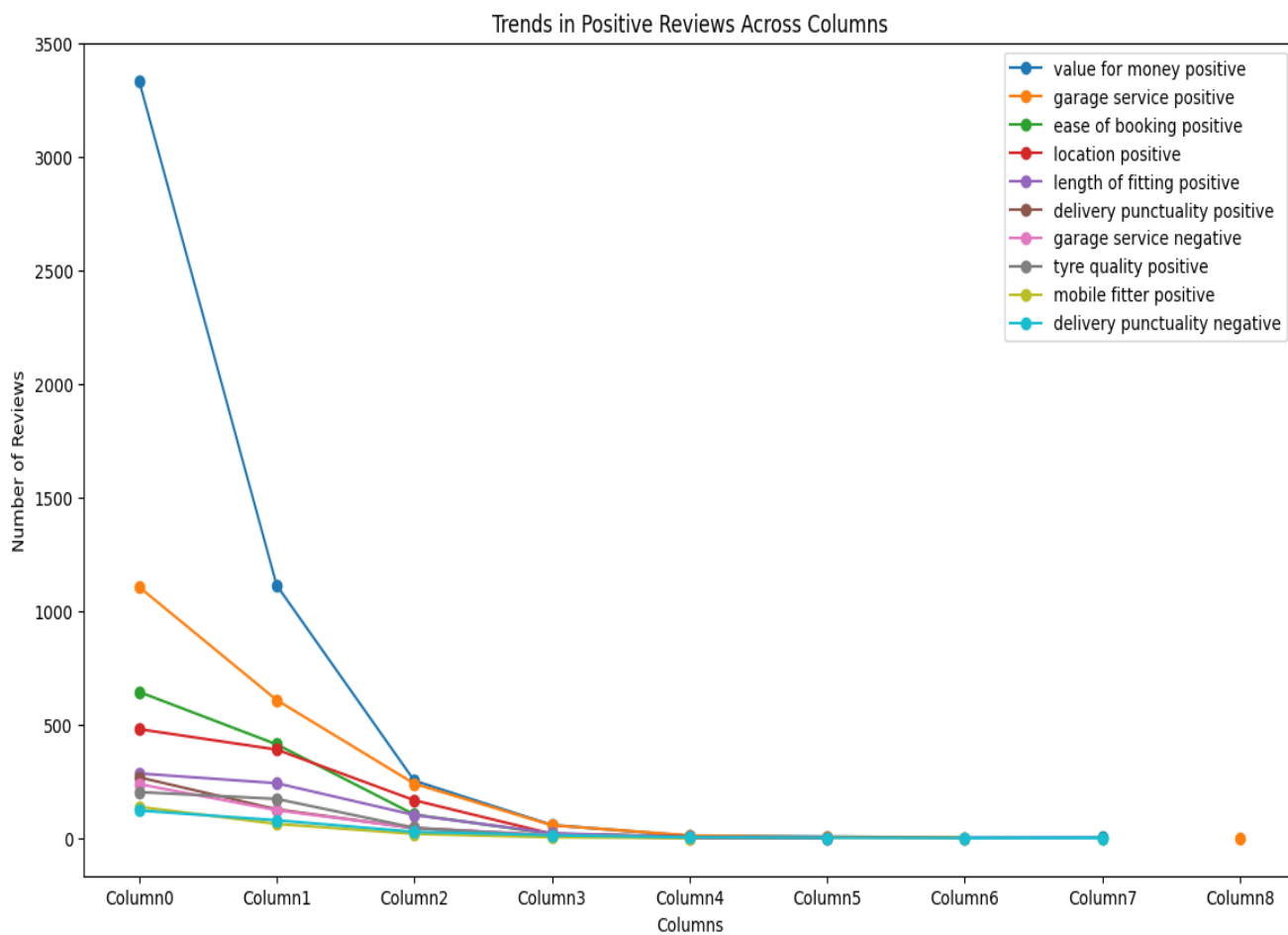
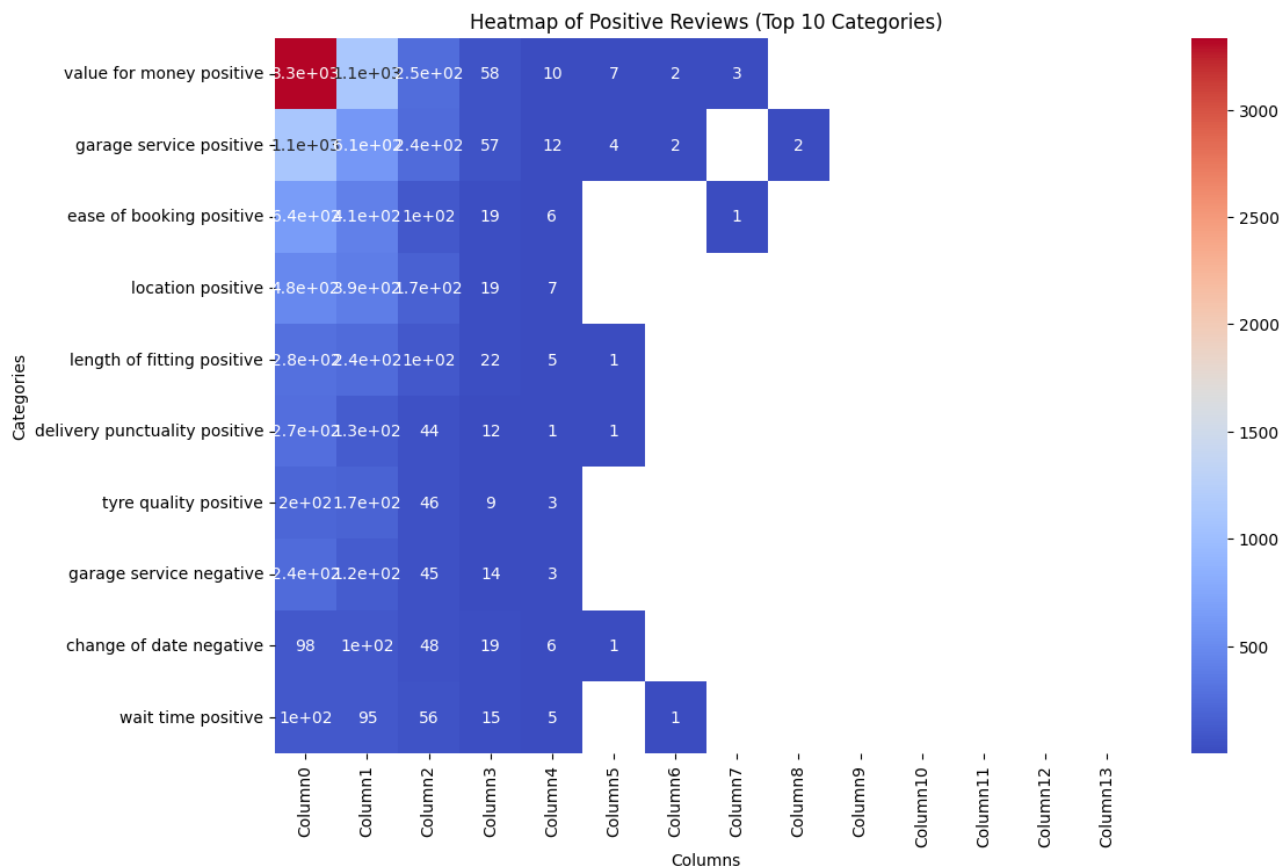
- Utilizing Latent Dirichlet Allocation (LDA) for topic modeling helps in identifying key topics or subthemes present in the text data. This is directly relevant to the task of identifying subthemes from reviews along with their respective sentiments.

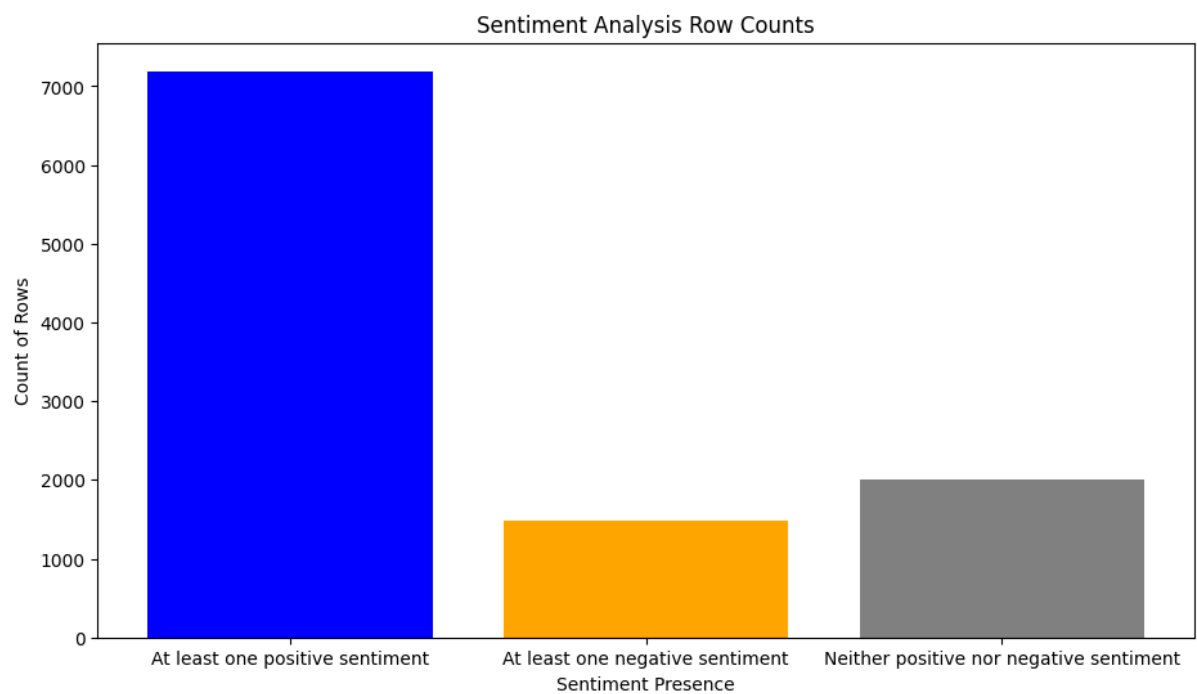
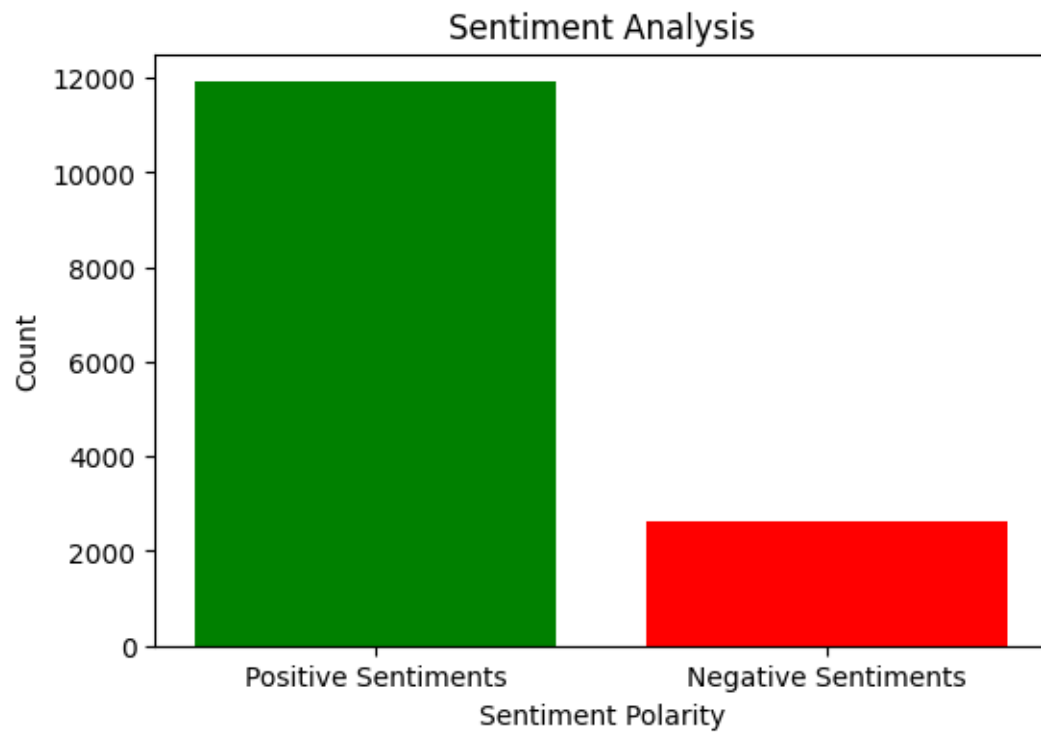
Conclusion:

Results include:

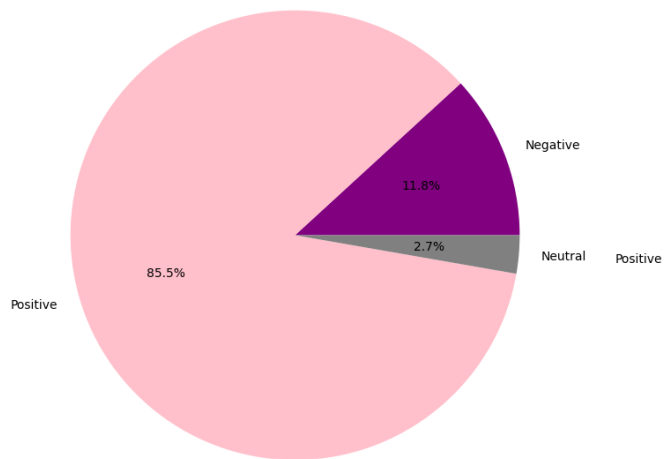
- Positive and negative sentiments are counted.
- Row counts by sentiment presence: The number of rows with one positive, one negative, or neither sentiment is calculated.
- Sentiment distributions and trends are visualized using multiple charts.
- Subtheme Classification: The RoBERTa model may classify subthemes using sentiment ratings, though not explicitly.
- LDA: Perform topic modeling to identify subthemes in the text data.
- Extracts key topics from the reviews, aiding in the identification of subthemes.



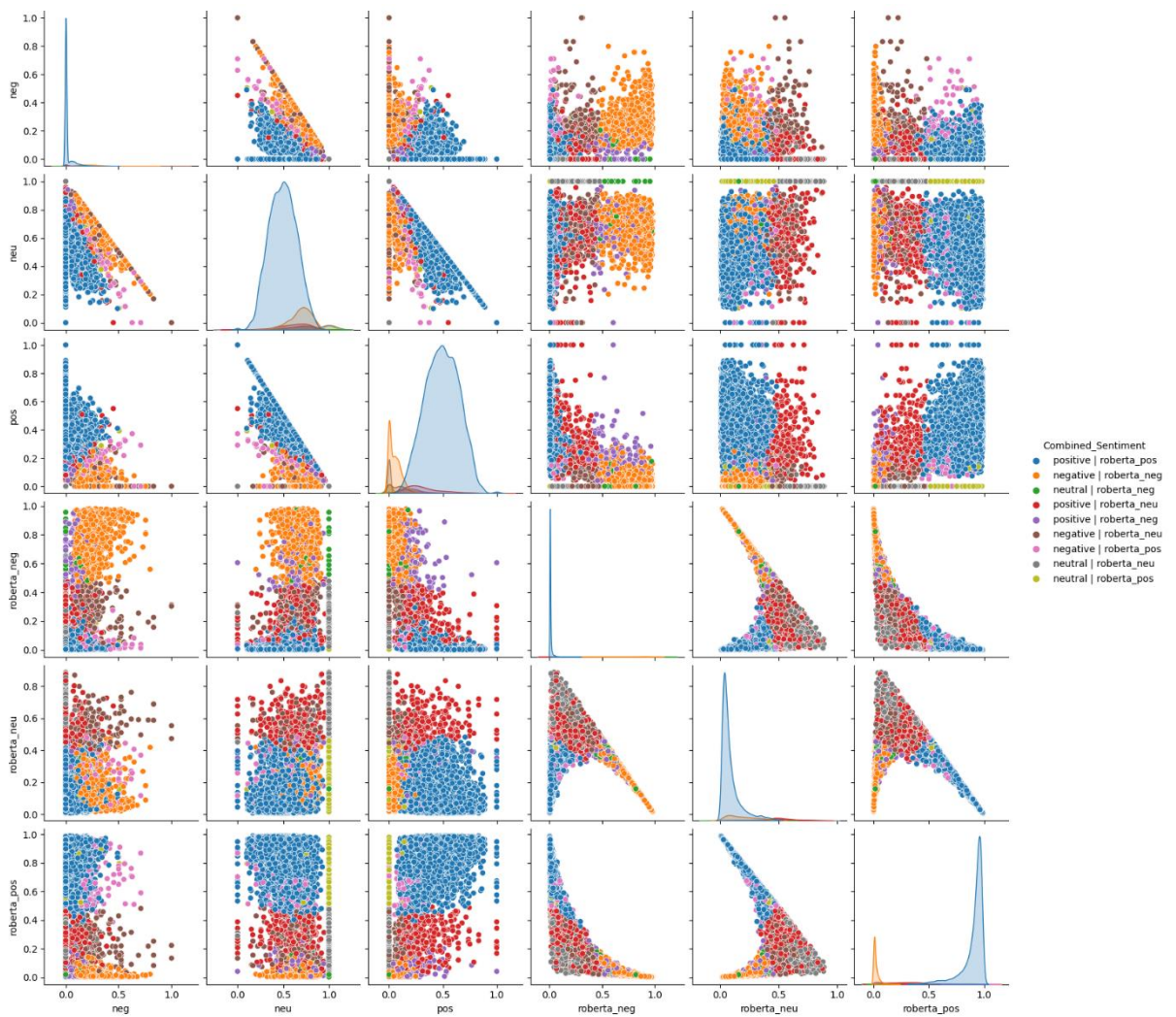
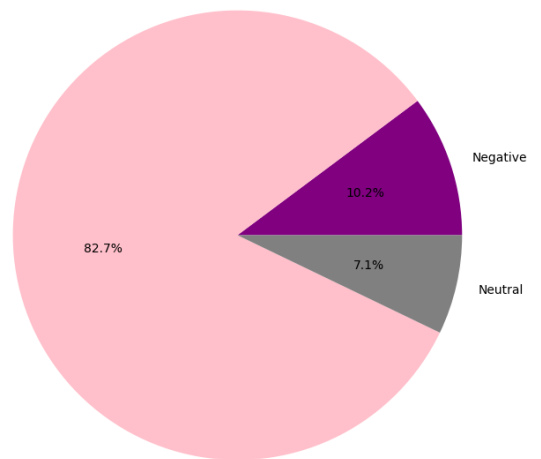


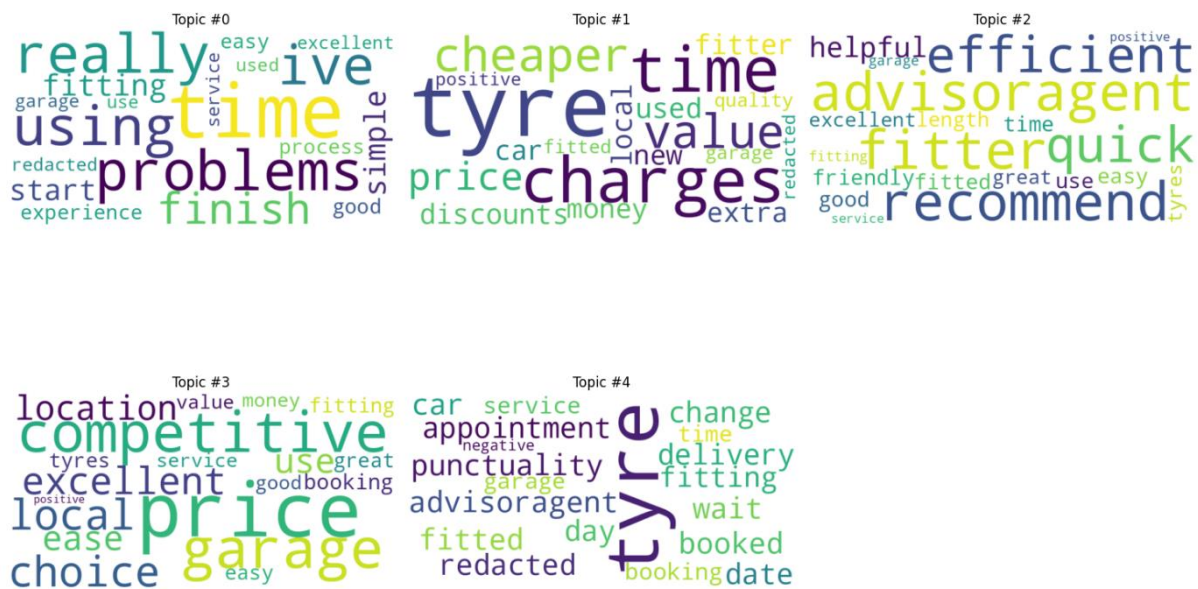
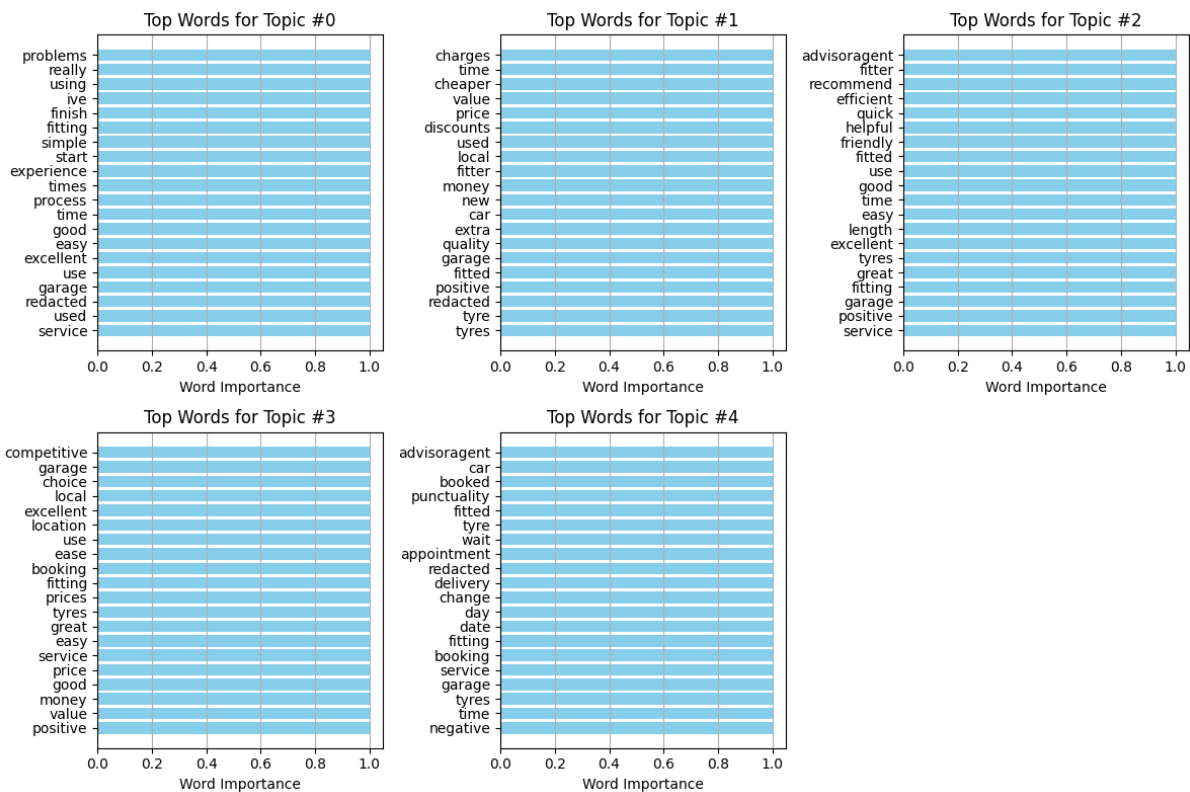


VADER Sentiment Distribution



RoBERTa Sentiment Distribution





Quality and completeness of code

1. **Functionality:** The algorithm counts dataset sentiments. It offers powerful sentiment analysis and visualizations using pretrained algorithms.
2. **Clarity and Readability:** The code is well-structured with step-by-step comments. This simplifies the analyzing procedure.
3. **Advanced Techniques:** Pretrained models like VADER and RoBERTa improve sentiment analysis by using advanced categorization methods.

Improvement Ideas

1. **Subtheme Identification:** The code should provide techniques to extract and classify subthemes to completely identify them and their attitudes. Possible methods include: - Text extraction for subtheme identification. Models can be used such as BERT, FLAN-T5
2. **Model Fine-tuning:** Pretrained models are reliable, but fine-tuning them on the dataset improves accuracy and relevance.
3. **Evaluation Metrics:** Precision, recall, and F1-score can quantify sentiment classification accuracy and performance.