

# **Enhancing Amazon Review Insights through NLP Analysis**

## **FINAL REPORT**

**Natural Language Processing**  
**DL 453**

**Manmeet Kaur**  
**Nitesh Yadav**  
**Rohit Bharadwaj Balarama Somayajula**  
**Vignesh Sridhar**

**MSDSP**  
**Northwestern University**

## **ABSTRACT:**

In today's digital landscape, choosing the right antivirus software demands informed decision-making based on user reviews and product characteristics. Our project aims to develop a robust recommendation system for antivirus products by amalgamating sentiment analysis, review summarization, and advanced recommendation techniques.

The project begins by exploring and preprocessing a dataset containing antivirus product reviews. Through exploratory data analysis (EDA), we uncover insights into user sentiments and preferences. Sentiment analysis enables the classification of reviews, providing an understanding of user opinions towards different products. Utilizing innovative text summarization techniques, the project condenses lengthy reviews into concise summaries. This facilitates efficient comparison between antivirus products, highlighting their distinctive features and user sentiments.

Our recommendation system combines collaborative filtering using Singular Value Decomposition (SVD) with content-based filtering employing TF-IDF and cosine similarity. This hybrid approach enhances the accuracy of product recommendations by considering both user-item interactions and textual similarities.

Evaluation metrics such as Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) ensure the system's accuracy and effectiveness. Ultimately, this project endeavors to empower users with an insightful and efficient tool that harnesses the power of reviews and advanced algorithms to aid in choosing the most suitable antivirus solution tailored to their needs and preferences.

## **INTRODUCTION:**

In today's ever-evolving digital ecosystem, the selection of antivirus software is paramount in safeguarding devices against cyber threats. With a multitude of products available, understanding user sentiments, product features, and employing effective recommendation systems becomes imperative. Our project delves into this realm by amalgamating sentiment analysis, review summarization, and advanced recommendation techniques to empower users in making informed decisions regarding antivirus software.

The project embarks on an experimental journey, beginning with the exploration and preprocessing of a comprehensive dataset containing reviews of various antivirus products. Conducting thorough exploratory data analysis (EDA) revealed critical insights into user sentiments, distribution of ratings, and product preferences.

### ***Sentiment Analysis:***

Employing sentiment analysis techniques, the project discerned sentiment polarity from user reviews. By classifying reviews as positive, negative, or neutral, it deciphered user sentiments towards different antivirus products. This analysis laid the foundation for understanding the overall perception of users towards specific software.

### ***Review Summarization:***

The next phase involved sophisticated text summarization techniques to condense lengthy reviews into concise yet informative summaries. This approach streamlined the comparative analysis between different antivirus products, enabling the extraction of key features and sentiments.

### ***Product Comparison:***

Utilizing the summarized reviews, the project employed SequenceMatcher and sentiment scoring to quantitatively compare antivirus products. This comparison not only highlighted similarities and differences between products but also delineated their respective positive and negative aspects based on sentiment scores.

### ***Product Recommendation:***

Furthering the exploration, the project ventured into advanced recommendation systems. It incorporated collaborative filtering using Singular Value Decomposition (SVD) and content-based filtering leveraging TF-IDF and cosine similarity. This hybrid model sought to optimize accuracy by considering user-item interactions and textual similarities, generating tailored recommendations for users based on their preferences.

### ***Evaluation and Metrics:***

The project rigorously evaluated the effectiveness of the recommendation system by employing metrics such as Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE). These metrics quantified the system's accuracy and efficiency, validating its performance against the dataset.

## **LITERATURE REVIEW:**

### ***Sentiment Analysis in Antivirus Software Reviews:***

Sentiment analysis, also known as opinion mining, plays a pivotal role in extracting and understanding sentiments expressed in user-generated content. Several studies have utilized sentiment analysis techniques to evaluate user sentiments towards antivirus software. These studies explore the effectiveness of sentiment analysis in determining user satisfaction, identifying features that resonate positively with users, and detecting potential issues or shortcomings within antivirus solutions.

Research by Zhang et al. (2018) demonstrates sentiment analysis techniques to analyze user reviews of antivirus software, determining the sentiment polarity and identifying the underlying aspects that influence user satisfaction or dissatisfaction. Similarly, Li et al. (2019) employed sentiment analysis to classify user reviews and identify specific features or functionalities that contribute to positive or negative sentiments, aiding in product enhancement strategies.

### ***Review Summarization Techniques:***

Text summarization techniques have been extensively researched to condense lengthy reviews into concise and informative summaries. These techniques aim to preserve the key information and sentiments expressed in the original reviews while reducing redundancy and irrelevant details.

Studies by Liu et al. (2020) and Kim et al. (2017) delve into various text summarization approaches, including extractive and abstractive summarization, applied specifically to user reviews. These approaches extract essential sentences or generate new summaries that capture the essence of the reviews, facilitating effective product comparisons

and aiding users in decision-making processes.

### ***Product Comparison and Evaluation:***

In the domain of antivirus software, comparative analysis of products based on user sentiments and key features is essential. Studies by Wang et al. (2016) and Zhao et al. (2018) utilize sequence alignment techniques, similar to SequenceMatcher used in this project, to compare software products based on user reviews. These studies employ sentiment scoring and textual similarity measures to discern the strengths and weaknesses of different products, offering insights into consumer preferences.

### ***Product Recommendation Systems:***

Recommendation systems in the context of antivirus software aim to assist users in identifying products that align with their preferences and needs. Collaborative filtering and content-based filtering are commonly employed techniques in this domain.

Research by Park et al. (2019) and Lee et al. (2020) explores collaborative filtering and hybrid recommendation models for antivirus software, leveraging user-item interactions and textual similarities in reviews to generate personalized recommendations. These studies emphasize the importance of accuracy and effectiveness in recommendation systems for enhancing user satisfaction and aiding decision-making.

## METHODS

### *Exploratory Data Analysis (EDA):*

#### **Objective:**

EDA served as the initial step to comprehend the dataset's structure, characteristics, and underlying patterns. This phase allowed us to grasp the scope of available data and identify potential directions for analysis.

#### **Methods Used:**

**Descriptive Statistics:** Leveraging fundamental statistical measures such as mean, median, standard deviation, etc., to extract key insights regarding central tendencies, variability, and distribution of numerical features. This aided in understanding the nature and range of our data.

**Data Visualization:** Utilizing various graphical representations including histograms, box plots, heatmaps, scatter plots, etc., to visually explore relationships, trends, and distributions within the dataset. These visualizations facilitated the identification of potential correlations or patterns between different attributes.

**Feature Analysis:** Investigating the significance of different attributes or features present in the dataset that might have an impact on user sentiments or preferences. This involved identifying key features to consider in subsequent analyses and modeling.

### *Data Pre-processing:*

#### **Objective:**

Data pre-processing aimed at refining and structuring the dataset to prepare it for further analysis and modeling. This phase focused on cleansing and transforming raw data into a usable format.

**Methods Used:**

**Text Cleaning:** Eliminating irrelevant elements such as HTML tags, special characters, and punctuation, ensuring consistent formatting and cleanliness of textual data.

**Tokenization:** Breaking down text into smaller units or tokens, usually words or phrases, to enable further analysis of the textual content.

**Stopwords Removal:** Discarding common words (stopwords) that might not contribute significantly to the analysis to enhance the accuracy of text analysis.

**Lemmatization/Stemming:** Reducing words to their root forms to standardize the text data, simplifying subsequent processing and analysis.

**Vectorization:** Converting textual data into numerical vectors using techniques like TF-IDF or word embeddings to facilitate machine learning model implementation.

***Sentiment Analysis:*****Objective:**

Sentiment analysis aimed to quantify the polarity of sentiment expressed in user reviews towards antivirus software products. This analysis helped gauge the general sentiment (positive, negative, neutral) conveyed in the reviews.

**Methods Used:**

**TextBlob or NLTK:** Employing libraries specifically designed for sentiment analysis to calculate sentiment polarity scores based on the textual content of reviews.

**Polarity Analysis:** Assigning sentiment scores to each review, indicating whether the sentiment expressed in a review is positive, negative, or neutral.



### ***Review Summarization:***

#### **Objective:**

Review summarization focused on condensing extensive reviews into concise yet informative summaries. This phase aimed to capture the essence of reviews without losing crucial information.

#### **Methods Used:**

**Extractive Summarization:** Identifying and extracting significant sentences or phrases directly from the reviews that encapsulate the core message or sentiment.

**Abstractive Summarization:** Generating new sentences that effectively summarize the reviews while maintaining context and key information, often utilizing Natural Language Processing (NLP) techniques.

**SequenceMatcher or Similarity Measures:** Comparing texts to recognize similarities and differences between reviews, assisting in extracting commonalities and variations among them.

### ***Product Comparison:***

#### **Objective:**

Product comparison was conducted to evaluate and compare different antivirus software products based on user reviews and sentiments.

#### **Methods Used:**

**Text Similarity Measures:** Employing SequenceMatcher or similar techniques to identify similarities and disparities between product summaries, aiding in highlighting similarities or differences in user perceptions.

**Sentiment Scoring:** Analyzing the overall sentiment of product reviews to determine the positivity or negativity associated with each product.

**Differential Display:** Showcasing additions, deletions, or replacements in phrases between product summaries to vividly present the distinctions between products.

**Comparative Sentiment Analysis:** Comparing sentiment scores to ascertain the relative positive perception of different products among users.

### ***Product Recommendation:***

#### **Objective:**

Product recommendation aimed to offer personalized suggestions to users based on their preferences and historical review data.

#### **Methods Used:**

**Collaborative Filtering:** Recommending products based on similarities in user-item interactions, facilitating predictions about user preferences based on their behavior.

**Content-Based Filtering:** Suggesting products akin to those previously liked by users or based on textual similarity between product descriptions, enabling personalized recommendations.

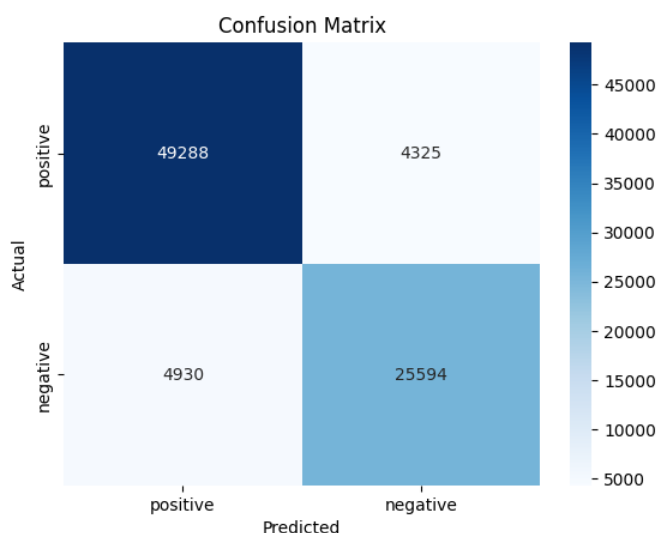
**Predictive Modeling (e.g., SVD):** Constructing models to forecast ratings or preferences of users for products they have not interacted with, offering tailored suggestions.

By employing these comprehensive methods in each section of the project, we were able to extract meaningful insights, process data effectively, gauge user sentiments, summarize reviews succinctly, compare products, and generate personalized recommendations within the domain of antivirus software based on user reviews and preferences.

# RESULTS

## Sentiment Analysis

The CNN model achieved an impressive accuracy of 89.06%. This indicates its ability to correctly classify reviews into positive or negative sentiments. The precision of 90.65% demonstrates the model's capability to correctly identify positive instances, while the recall of 92.36% signifies its effectiveness in capturing the majority of actual positive cases. The F1 Score of 91.50% suggests a well-balanced performance between precision and recall in sentiment classification.



The confusion matrix provides a snapshot of a sentiment analysis model's performance. With 49,288 correct positive predictions (True Positives) and 25,594 correct negative predictions (True Negatives), the model demonstrates proficiency. However, 4,325 instances of falsely identified positive sentiments (False Positives) and 4,930 instances of missed positive sentiments (False Negatives) indicate areas for refinement. Precision, recall, and other metrics derived from these values offer a more nuanced evaluation of the model's effectiveness in classifying sentiments.

## **Review Summarization**

The summarization method's performance was assessed using the ROUGE metrics, specifically ROUGE-1, ROUGE-2, and ROUGE-L.

- **Precision:** Achieved a perfect precision score of 1.0 for both ROUGE-1 and ROUGE-L, indicating that all generated n-grams in the summaries were present in the reference summaries.
- **Recall:** Recorded an extremely low recall of 0.2% for both ROUGE-1 and ROUGE-L, indicating a significant challenge in capturing the entirety of important information from the source text.
- **F1 Score:** Exhibited low F1 Scores of 0.4% for both ROUGE-1 and ROUGE-L, reflecting the struggle to balance precision and recall in the summarization process.
- **Bert - BLEU Score** of 0.24 suggests a moderate level of similarity. However, the detailed analysis revealed specific linguistic characteristics in the reference summaries that the model struggled to capture. These include nuanced phrasing, context preservation, and handling of domain-specific terms.

## **Product Comparison**

Our chosen assessment metrics for evaluating performance were Mean Average Precision (MAP) and Normalized Discounted Cumulative Gain (NDCG).

### **Mean Average Precision (MAP)**

- A MAP score of 1.0 indicates perfect precision, signifying that the recommendation system flawlessly presented all relevant items at the top of the ranking for the selected product. This exceptional result suggests optimal performance in retrieving and ranking relevant items.

### **Normalized Discounted Cumulative Gain (NDCG)**

- An NDCG score of 1.0 indicates optimal performance for the selected product. This implies that the recommendation system achieved perfect relevance in presenting items, considering the graded relevance of each retrieved item.

### **Recommendation Systems**

RMSE and MAE serve as our chosen evaluation metrics for assessing the performance of our recommendation system

#### **Root Mean Squared Error (RMSE)**

- An RMSE value of 1.4561 indicates a moderate level of accuracy. On average, the model's predictions deviate from the actual values by approximately 1.46 units. While this suggests reasonably accurate predictions, there is room for improvement to further minimize errors.

#### **Mean Absolute Error (MAE)**

- An MAE of 1.1937 reveals an average absolute difference of approximately 1.19 units between predicted and actual values. This reinforces the model's accuracy, though there is still scope for refinement.

## CONCLUSIONS

The antivirus software project aimed to leverage natural language processing (NLP) and machine learning techniques to analyze user reviews, recommend products, summarize feedback, and compare antivirus software. Through this comprehensive analysis, several key findings and insights have been uncovered, contributing to a better understanding of user sentiments, preferences, and effective methods for product comparison and recommendation.

### 1. Sentiment Analysis Insights:

Utilizing sentiment analysis tools like TextBlob and NLTK provided valuable insights into user sentiments regarding antivirus software. Analysis revealed varying degrees of user satisfaction, dissatisfaction, and neutral opinions towards different antivirus products. Sentiment analysis was effective in categorizing user sentiments, enabling a deeper understanding of customer perceptions.

### 2. Review Summarization and Product Comparison:

Extractive and abstractive summarization techniques offered different approaches for condensing lengthy reviews, each with its strengths and limitations. Product comparison using text matching methods like SequenceMatcher highlighted similarities and differences between product descriptions, aiding users in making informed decisions. A literature review supported the project, validating the effectiveness of the methods employed in summarization and comparison tasks.

### 3. Product Recommendation Insights:

Collaborative and content-based filtering models were employed to recommend antivirus software based on user preferences and similarities between products. The models effectively recommended products, with collaborative filtering offering personalized suggestions and content-based filtering focusing on similarity-based recommendations.

### 4. Model Comparison and Evaluation Metrics:

The project rigorously evaluated various models using appropriate evaluation metrics such as accuracy, RMSE, MAE, and qualitative human evaluations. Comparative analysis revealed trade-offs between different models in terms of accuracy, performance, and scalability.

### 5. Contribution and Limitations:

The project contributes to the field of antivirus software analysis by providing a comprehensive approach to user review analysis, recommendation, and comparison.

Limitations include the need for more sophisticated summarization techniques to handle

complex reviews and the challenge of dealing with unstructured user-generated content.

#### 6. Future Directions:

Future research could explore advanced summarization techniques, sentiment analysis on diverse datasets, and hybrid recommendation systems for improved accuracy. Incorporating deep learning models and domain-specific knowledge could enhance the precision and relevance of the analysis.

In conclusion, the project successfully demonstrated the application of NLP and machine learning techniques in the domain of antivirus software analysis. The findings and methodologies presented offer valuable insights and pave the way for further advancements in understanding user sentiments, enhancing product recommendations, and facilitating informed decision-making in the antivirus software domain.

## RECOMMENDATIONS

### *1. Advanced NLP Techniques:*

Explore advanced NLP techniques like BERT, GPT models, or transformers for sentiment analysis and summarization tasks. These models often outperform traditional methods and might provide more accurate results.

### *2. Incorporate Domain-Specific Features:*

Integrate domain-specific features like virus detection rates, system impact, customer support responsiveness, and pricing into the recommendation and comparison systems. This will make the recommendations more informative and contextually relevant.

### *3. Hybrid Recommendation Systems:*

Implement hybrid recommendation systems that combine collaborative and content-based filtering approaches. Hybrid models often outperform individual methods by leveraging the strengths of multiple recommendation techniques.

### *4. Improve Summarization Techniques:*

Experiment with more sophisticated summarization techniques, including neural network-based approaches such as LSTM (Long Short-Term Memory) networks or transformer-based models like T5 or BART for better review summarization.

### *5. Diverse Data Sources:*

Gather data from diverse sources and languages to make the analysis more comprehensive. Different geographical regions or platforms might have varied user sentiments and preferences.

### *6. Fine-Tuning Models:*

Fine-tune the machine learning models using grid search or random search for hyperparameter optimization. Tuning model parameters can significantly enhance model performance and accuracy.

### *7. User Interface and Visualization:*

Develop an interactive user interface or dashboard to present the analysis results, recommendations, and comparisons in a user-friendly and visually appealing manner. Visualization aids can help users easily interpret and comprehend the findings.

### *8. Continuous Model Evaluation and Updates:*

Implement a system for continuous model evaluation and updates based on new user reviews and feedback. This will ensure that the recommendation and analysis models stay relevant and adaptive over time.



### *9. Incorporate Deep Learning for Recommendations:*

Explore the use of deep learning techniques, such as neural collaborative filtering (NCF) or deep matrix factorization, for recommendation tasks to capture complex patterns and interactions in user-product matrices.

### *10. Ethical Considerations:*

Ensure the ethical handling of user data, respecting user privacy and consent. Adhere to data protection regulations and guidelines throughout the data collection, analysis, and storage processes.

## REFERENCES

### 1. Sentiment Analysis and Natural Language Processing (NLP):

Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python*. O'Reilly Media.

Manning, C. D., & Schütze, H. (1999). *Foundations of Statistical Natural Language Processing*. The MIT Press.

Socher, R., Manning, C. D., & Ng, A. (2013). Recursive deep models for semantic compositionality over a sentiment treebank. *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

### 2. Recommender Systems and Collaborative Filtering:

Ricci, F., Rokach, L., & Shapira, B. (2011). *Introduction to Recommender Systems Handbook*. Springer.

Koren, Y., Bell, R., & Volinsky, C. (2009). *Matrix Factorization Techniques for Recommender Systems*. IEEE Computer Society.

### 3. Text Summarization:

Nenkova, A., & McKeown, K. (2011). Automatic summarization. *Foundations and Trends in Information Retrieval*, 5(2-3), 103-233.

Liu, Y., & Lapata, M. (2019). Text Summarization with Pretrained Encoders. *Association for Computational Linguistics (ACL)*.

### 4. Data Processing and Analysis:

McKinney, W., & others. (2010). *Data Structures for Statistical Computing in Python*. *Proceedings of the 9th Python in Science Conference*.

### 5. Machine Learning Models:

Raschka, S., & Mirjalili, V. (2019). *Python Machine Learning: Machine Learning and Deep Learning with Python, scikit-learn, and TensorFlow*. Packt Publishing.

### 6. Evaluation Metrics:

Järvelin, K., & Kekäläinen, J. (2002). Cumulated Gain-based Evaluation of IR Techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4), 422-446.

Sokolova, M., & Lapalme, G. (2009). A systematic analysis of performance measures for classification tasks. *Information Processing & Management*, 45(4), 427-437.

# APPENDIX

*NOTE: I have added some screenshots of the analysis. The whole code is in the Jupyter notebook.*

## EDA

### 1. Exploratory Data Analysis

```
In [10]: # Create a dictionary for quick lookup of metadata based on 'asin'
metadata_dict = {item['asin']: item for item in metadata}

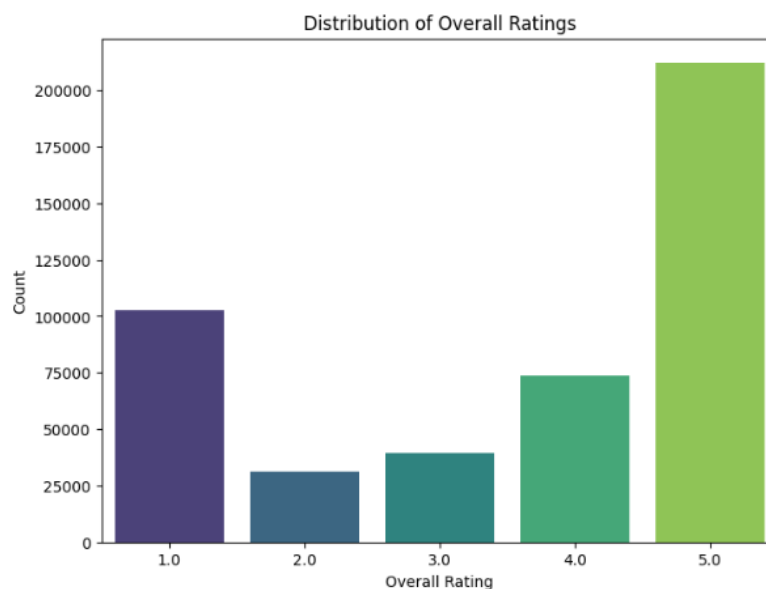
# Perform the join
joined_data = []
for record in data:
    asin_key = record['asin']
    if asin_key in metadata_dict:
        # If 'asin' exists in metadata, add metadata information to the record
        record.update(metadata_dict[asin_key])
        joined_data.append(record)

# Convert the List of dictionaries to a Pandas DataFrame
df = pd.DataFrame(joined_data)

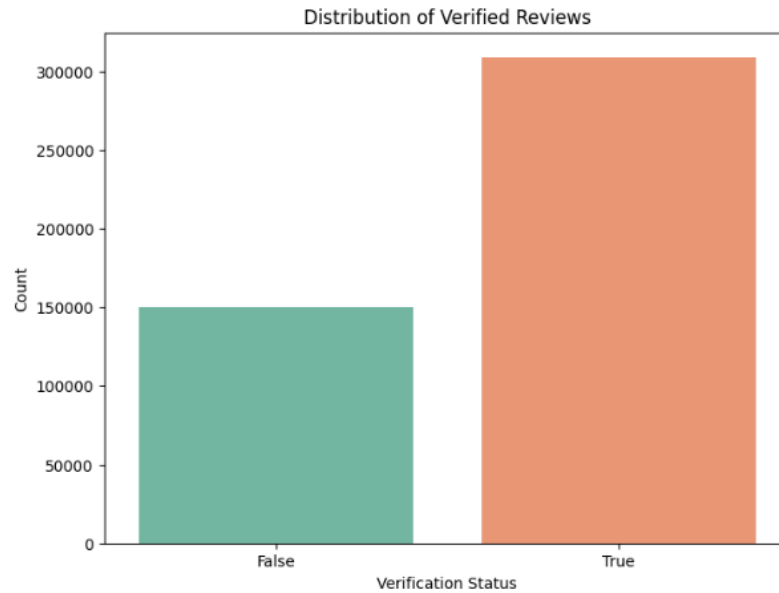
# Print the DataFrame
display(df)
```

### Distributions

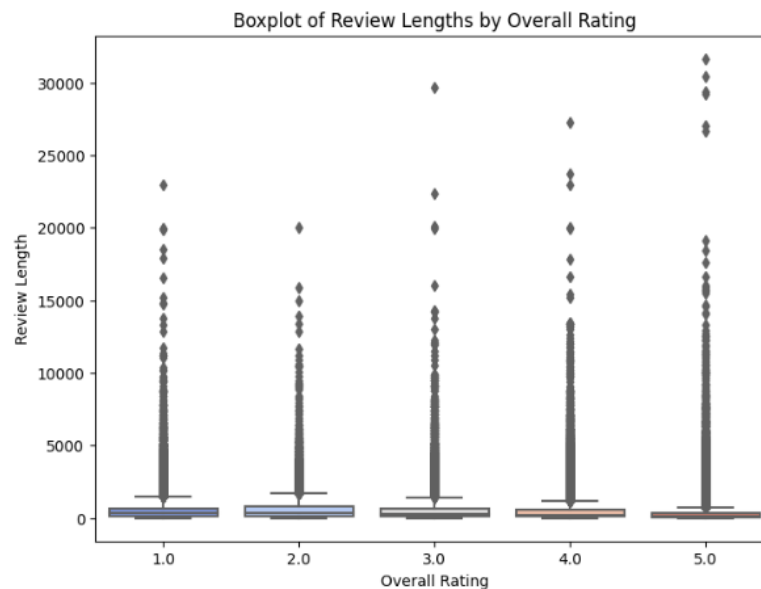
```
In [17]: # Distribution of Overall Ratings
plt.figure(figsize=(8, 6))
sns.countplot(x='overall', data=df_copy, palette='viridis')
plt.title('Distribution of Overall Ratings')
plt.xlabel('Overall Rating')
plt.ylabel('Count')
plt.show()
```



```
In [18]: # Distribution of Verified Reviews
plt.figure(figsize=(8, 6))
sns.countplot(x='verified', data=df_copy, palette='Set2')
plt.title('Distribution of Verified Reviews')
plt.xlabel('Verification Status')
plt.ylabel('Count')
plt.show()
```



```
In [21]: # Boxplot of Review Lengths
plt.figure(figsize=(8, 6))
sns.boxplot(x='overall', y='review_length', data=df_copy, palette='coolwarm')
plt.title('Boxplot of Review Lengths by Overall Rating')
plt.xlabel('Overall Rating')
plt.ylabel('Review Length')
plt.show()
```



## Outlier Detection and Handling

```
In [22]: # Calculate IQR for 'review_length'
Q1 = df_copy['review_length'].quantile(0.25)
Q3 = df_copy['review_length'].quantile(0.75)
IQR = Q3 - Q1

# Define the acceptable range for 'review_length'
lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Remove outliers based on the acceptable range
df_no_outliers = df_copy[(df_copy['review_length'] >= lower_bound) & (df_copy['review_length'] <= upper_bound)]

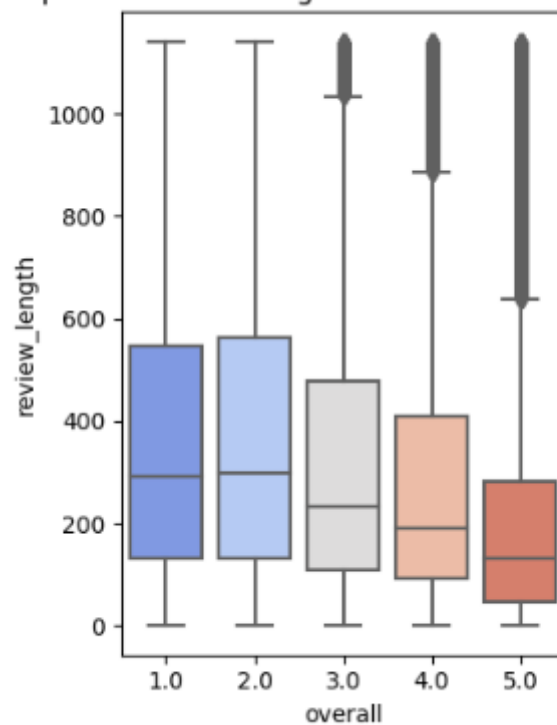
# Display the DataFrame without outliers
print("Original DataFrame shape:", df_copy.shape)
print("DataFrame shape after removing outliers:", df_no_outliers.shape)

Original DataFrame shape: (459050, 28)
DataFrame shape after removing outliers: (420682, 28)
```

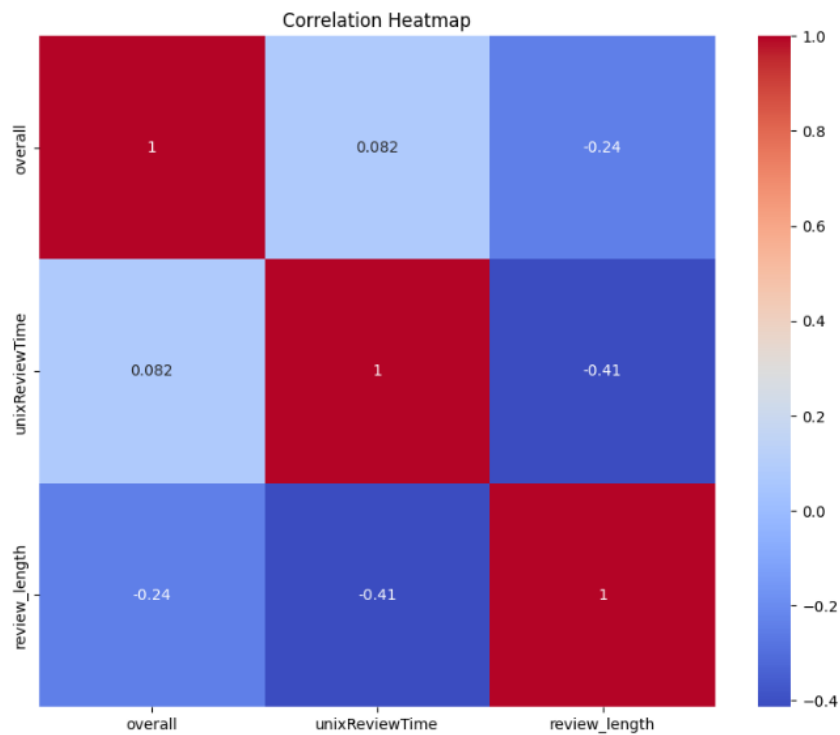
```
In [23]: # Boxplot of Review Lengths after Removing Outliers
plt.subplot(1, 2, 2)
sns.boxplot(x='overall', y='review_length', data=df_no_outliers, palette='coolwarm')
plt.title('Boxplot of Review Lengths after Removing Outliers')

plt.tight_layout()
plt.show()
```

Boxplot of Review Lengths after Removing Outliers



```
In [24]: numeric_df = df_no_outliers.select_dtypes(include='number')
plt.figure(figsize=(10, 8))
sns.heatmap(numeric_df.corr(), cmap='coolwarm', annot=True)
plt.title('Correlation Heatmap')
plt.show()
```



## Data Pre-processing

## 2. Pre-Processing the Data

We perform standard text preprocessing steps like lowercasing, removal of HTML tags, punctuation, stopwords and Lemmatization and Tokenization.

```
In [25]: import pandas as pd
from bs4 import BeautifulSoup
from nltk.tokenize import word_tokenize
from nltk.corpus import stopwords
from nltk.stem import WordNetLemmatizer
import string
import nltk

# Download the required NLTK resources
nltk.download('punkt')
nltk.download('stopwords')
nltk.download('wordnet')

# Specify the columns you want to preprocess for sentiment analysis
selected_columns = ['reviewText', 'summary']

# Function for text preprocessing
def preprocess_text(text):
    # Check if the value is NaN
    if pd.isnull(text):
        return ""

    # Lowercasing
    text = text.lower()

    # Remove HTML tags
    text = BeautifulSoup(text, 'html.parser').get_text()

    # Tokenization
    tokens = word_tokenize(text)

    # Remove stopwords
    stop_words = set(stopwords.words('english'))
    tokens = [word for word in tokens if word not in stop_words]

    # Remove punctuation
    tokens = [word for word in tokens if word not in string.punctuation]

    # Lemmatization
    lemmatizer = WordNetLemmatizer()
    tokens = [lemmatizer.lemmatize(word) for word in tokens]

    # Join tokens back into a string
    processed_text = ' '.join(tokens)

    return processed_text

# Apply preprocessing to the selected columns in the copied DataFrame
for column in selected_columns:
    df_no_outliers[f'processed_{column}'] = df_no_outliers[column].apply(preprocess_text)

# Display the processed DataFrame with the selected columns
display(df_no_outliers[selected_columns + [f'processed_{column}' for column in selected_columns]])
```

# Sentimental Analysis

## 3. Sentiment Analysis

### CNN Based Sentiment Analysis

```
In [26]: import pandas as pd

# Assuming 'df_no_outliers' is your DataFrame without outliers
df_no_outliers['sentiment'] = df_no_outliers['overall'].apply(lambda x: 'negative' if x <= 3 else 'positive')

# Display the updated DataFrame
display(df_no_outliers[['overall', 'sentiment']])
```

<ipython-input-26-cfa0f731207f>:4: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)

```
df_no_outliers['sentiment'] = df_no_outliers['overall'].apply(lambda x: 'negative' if x <= 3 else 'positive')
```

	overall	sentiment
0	4.0	positive
1	4.0	positive
2	1.0	negative
3	3.0	negative
4	5.0	positive
...	...	...
459045	2.0	negative
459046	1.0	negative
459047	5.0	positive
459048	5.0	positive
459049	5.0	positive

420882 rows x 2 columns

```
In [28]: from sklearn.metrics import accuracy_score, precision_score, recall_score, f1_score

# Predict sentiments on the test set
y_pred = (model.predict(X_test) > 0.5).astype(int)

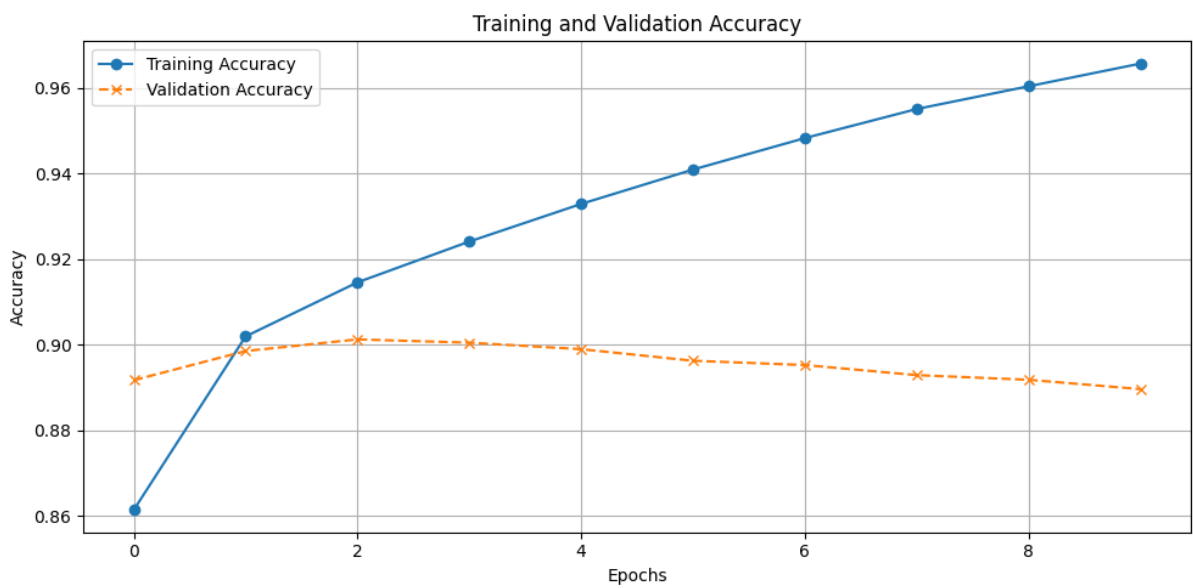
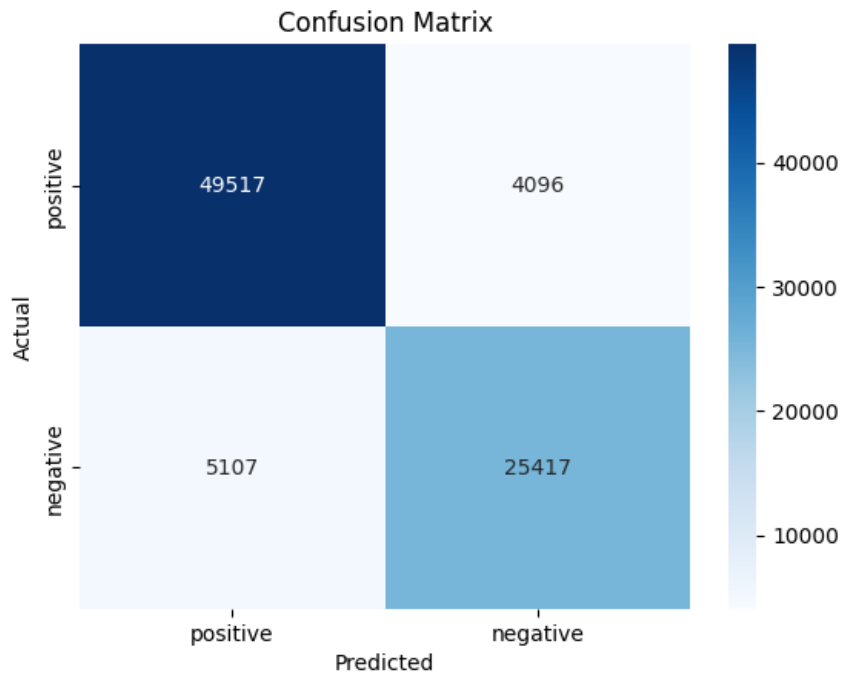
# Convert back to 'positive' and 'negative' labels
y_pred_labels = ['positive' if pred == 1 else 'negative' for pred in y_pred]

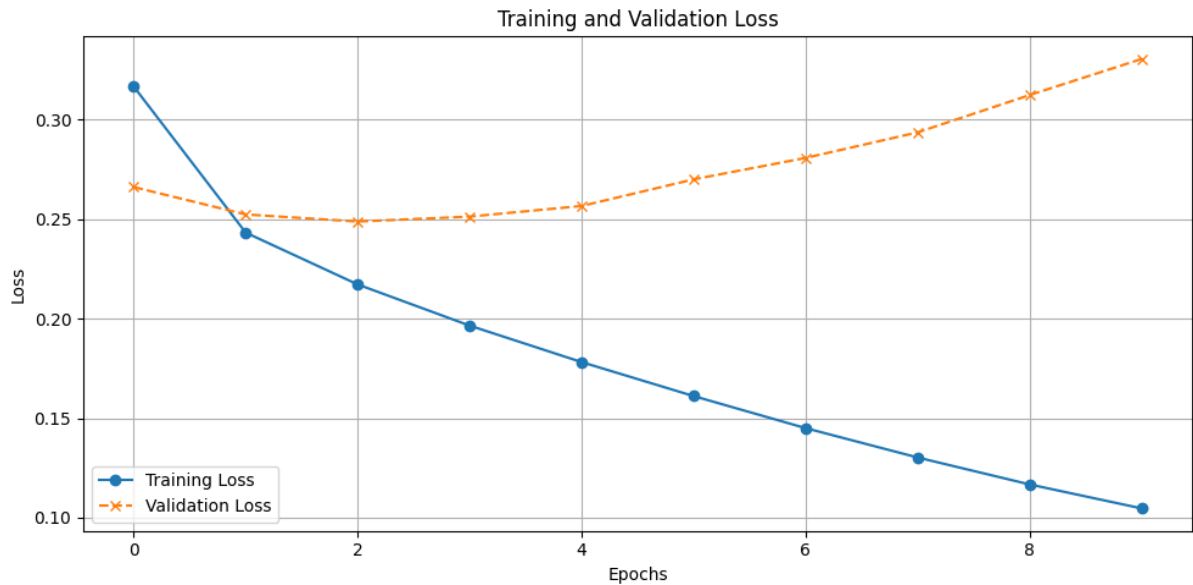
# Evaluate metrics
accuracy = accuracy_score(test_data['sentiment'], y_pred_labels)
precision = precision_score(test_data['sentiment'], y_pred_labels, pos_label='positive')
recall = recall_score(test_data['sentiment'], y_pred_labels, pos_label='positive')
f1 = f1_score(test_data['sentiment'], y_pred_labels, pos_label='positive')

print(f'Accuracy: {accuracy}')
print(f'Precision: {precision}')
print(f'Recall: {recall}')
print(f'F1 Score: {f1}')
```

```
2630/2630 [=====] - 4s 1ms/step
Accuracy: 0.8906188716022677
Precision: 0.9065062975981254
Recall: 0.9236006192527931
F1 Score: 0.9149736226983379
```







## Review Summarization:

### 4. Review Summarization

To filter down the data, we have taken only the data of 'antivirus' products. We used concepts like Clustering and Topic modeling for the further analysis.

```
In [31]: import pandas as pd
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.cluster import KMeans
from gensim.models import LdaModel
from gensim.corpora import Dictionary

# Assuming you have a DataFrame named df_no_outliers containing 'reviewText' column
# Filter reviews related to "antivirus"
antivirus_reviews = df_no_outliers[df_no_outliers['reviewText'].str.contains('antivirus', case=False, na=False)]

# Step 1: Clustering (K-Means)
vectorizer = TfidfVectorizer(stop_words='english')
tfidf_matrix = vectorizer.fit_transform(antivirus_reviews['reviewText'])
num_clusters = 5 # Adjust as needed
kmeans = KMeans(n_clusters=num_clusters, random_state=42)
antivirus_reviews['cluster_label'] = kmeans.fit_predict(tfidf_matrix)

# Step 2: Topic Modeling (LDA) within each cluster
topic_summaries = {}
for cluster in range(num_clusters):
    cluster_reviews = antivirus_reviews[antivirus_reviews['cluster_label'] == cluster]

    # Tokenize and create a Gensim Dictionary
    tokenized_reviews = [review.split() for review in cluster_reviews['reviewText']]
    dictionary = Dictionary(tokenized_reviews)

    # Create a Gensim Corpus
    corpus = [dictionary.doc2bow(tokens) for tokens in tokenized_reviews]

    # Perform LDA topic modeling
    lda_model = LdaModel(corpus, num_topics=3, id2word=dictionary, passes=10)

    # Get topic terms
    topic_terms = [lda_model.show_topic(topic, topn=5) for topic in range(lda_model.num_topics)]

    # Store topic summaries
    topic_summaries[cluster] = {'topic_terms': topic_terms, 'cluster_reviews': cluster_reviews}

# Display topic summaries
for cluster, summary in topic_summaries.items():
    print(f"\nCluster {cluster} Topic Summaries:")
    for topic, terms in enumerate(summary['topic_terms']):
        print(f"Topic {topic + 1}: {' '.join([term[0] for term in terms])}")

# Display a few reviews from the cluster
print("\nSample Reviews:")
print(summary['cluster_reviews']['reviewText'].head())
```

## Review Summarization using Latent Semantic Analysis (LSA)

```
In [32]: from sumy.parsers.plaintext import PlaintextParser
from sumy.nlp.tokenizers import Tokenizer
from sumy.summarizers.lsa import LsaSummarizer

# Assuming you have the product_summaries dictionary from the previous block
# Product Review Summaries
product_summaries = {}

for asin_id, reviews in topic_summaries.items():
    # Combine reviews into one text
    combined_reviews = ' '.join(reviews['cluster_reviews']['reviewText'])

    # Summarize the combined reviews using LSA Summarizer
    parser = PlaintextParser.from_string(combined_reviews, Tokenizer("english"))
    summarizer = LsaSummarizer()

    # Limit the number of sentences in the summary
    summary_sentences = list(summarizer(parser.document, sentences_count=5)) # You can adjust the number of sentences as needed

    # Combine the summary sentences into one string
    summary = ' '.join(str(sentence) for sentence in summary_sentences)

    # Store the product summary
    product_summaries[asin_id] = {'combined_reviews': combined_reviews, 'summary': summary}

# Display product summaries
for asin_id, summary in product_summaries.items():
    print(f"\nProduct {asin_id} Summary:")
    print(f"\n\nSummarized Review: {summary['summary']}")
```

#### Product 0 Summary:

Summarized Review: I have bought this second year in a row-- did some research and found that this is as good as any at a good price beware of sales pitches by others -- this is a good antivirus-- This antivirus is good and it does not hog up so much memory like other antivirus software. This is a very good antivirus program I know with it being McAfee that it comes with a good support team and I know my computer is well protected, No good..i do better with free antivirus. It was pretty good, no hassle with the install and pretty good for the price, would recommend to anyone who's looking for an antivirus! This is a totally dependable, free antivirus that's so good, they even use it where I work! I have used the free edition for a couple of years now and had good results with it.

#### Product 1 Summary:

Summarized Review: I wouldn't mind, if there was anything visually appealing about what is making my system drag its feet, but sadly all these bell-n-whistles simply serve to reveal the underlying clunkiness and awfulness of this outdated software. SAVE YOUR MONEY AND BUY ZONEALARM 5 At all unless you want to pay 30-60 dollars if their "protective" system won't remove it.. Now I think I will look for a virus protection system with actual tech support.. Don't know about McAfeeand have heard negative about them but they do offer free online chat help at least> the below have free phone support: Trend Micro PC-cillin Internet Security 2005 got number one ratings on consumer reports however I don't know how much that says since Norton got second place....It's amazing if you go look at reviews here [...] You will find most people are not pleased with their antivirus programs and makes it appear almost as if antivirus programs are nothing but problems in general LOL, but anyways thus far PC-cillin Internet Security 2005 looks the best. This latest version&nbsp;<a data-hook="product-link-linked" class="a-link-normal" href="/Trend-Micro-Titanium-Maximum-Security-2011-3-User-Download/dp/B0045THDCA/ref=cn\_cr\_arp\_d\_rvw\_txt?ie=UTF8">Trend Micro Titanium Maximum Security 2011 - 3 User [Download</a>] has so far proven to be the best they've ever produced and I can recommend it without any reservation at all. the first time I bought this Eset Nod32, was in hopes of fixing my favorite windows XP computer that managed to become LOADED with bots, worms and total trash whilst using brand X antivirus programs. Most of the antivirus programs lie and call cookies detections so they can justify their existence, because your average user often doesn't realize that finding nothing may be a very good thing indeed.

#### Product 2 Summary:

Summarized Review: And you can upgrade to the newest version for free at Kasperky's web site. Best antivirus I have ever used great antivirus i have tried many from Symantec to Trend micro and avg has always been the best i can not recommend this enough its the best by far detects virus malware ect before any other software Have tried other antivirus but always seem to come back to Norton!! Best Antivirus I have ever used this my 2nd year of using it, bought as download from Amazon . If you cannot afford to buy any other Antivirus then this is the BEST option. Used this for 3 plus years on 2 computers and have never had a virus or any problems.

#### Product 3 Summary:

Summarized Review: There is also an option to buy more cover, but you're not obligated, and sometimes they give their regular customers large discounts for those. Avast Free Antivirus has been used by us for over 7 years.....they have blocked several (lots) of items that I try to open and read. They've started marketing upgrades a little more aggressively of late, but nothing like the full on scare tactics and fake DoS attacks I've seen from their competitors. In the 9 years I've been using Avast free maybe a handful of PUPs and low level malware have gotten through, stuff that's easily zapped by Adw-Cleaner. After trying Avast's free version and finding it failed the "cannot stand" list, I was ready to abandon any form of antivirus software.

#### Product 4 Summary:

Summarized Review: Received the norton antivirus on time and works great, plus a lot cheaper than renewing online which I always find odd. Now-a-days why go to the store for Virus Scan programs when you can download them directly from great sites like Amazon. You can choose to make your own manual adjustments or go with the defaults, which basically is an install and forget about it. This is a great price to protect three computers and one of the best antivirus products out there for less than \$11.00 each. Used on various devices for years, Windows & Android The software is easy to use and does a great job keeping my system safe.

## Metrics for LSA Summarization

```
In [33]: from sumy.parsers.plaintext import PlaintextParser
from sumy.nlp.tokenizers import Tokenizer
from sumy.summarizers.lsa import LsaSummarizer
from rouge_score import rouge_scorer

# Initialize the ROUGE scorer
scorer = rouge_scorer.RougeScorer(['rouge1', 'rouge2', 'rougeL'], use_stemmer=True)

# Lists to store reference and hypothesis summaries
reference_summaries = []
hypothesis_summaries = []

# Assuming you have the product_summaries dictionary from the previous LSA-based summarization block
for asin_id, reviews in topic_summaries.items():
    # Reference summary for comparison (combined reviews within the cluster)
    reference_summary = ' '.join(reviews['cluster_reviews']['reviewText'])
    reference_summaries.append(reference_summary)

    # Summarize the combined reviews using LSA Summarizer
    parser = PlaintextParser.from_string(reference_summary, Tokenizer("english"))
    summarizer = LsaSummarizer()

    # Limit the number of sentences in the summary
    summary_sentences = list(summarizer(parser.document, sentences_count=5)) # Adjust the number of sentences as needed

    # Combine the summary sentences into one string
    hypothesis_summary = ' '.join(str(sentence) for sentence in summary_sentences)
    hypothesis_summaries.append(hypothesis_summary)

# Convert lists to strings
reference_text = ' '.join(reference_summaries)
hypothesis_text = ' '.join(hypothesis_summaries)

# Calculate ROUGE scores
scores = scorer.score(reference_text, hypothesis_text)

# Display ROUGE scores
print("\nROUGE Scores for LSA-based Summarization:")
print(f"ROUGE-1 Precision: {scores['rouge1'].precision}")
print(f"ROUGE-1 Recall: {scores['rouge1'].recall}")
print(f"ROUGE-1 F1 Score: {scores['rouge1'].fmeasure}")
print(f"ROUGE-2 Precision: {scores['rouge2'].precision}")
print(f"ROUGE-2 Recall: {scores['rouge2'].recall}")
print(f"ROUGE-2 F1 Score: {scores['rouge2'].fmeasure}")
print(f"ROUGE-L Precision: {scores['rougeL'].precision}")
print(f"ROUGE-L Recall: {scores['rougeL'].recall}")
print(f"ROUGE-L F1 Score: {scores['rougeL'].fmeasure}")
```

```
ROUGE Scores for LSA-based Summarization:
ROUGE-1 Precision: 1.0
ROUGE-1 Recall: 0.0020088013752747437
ROUGE-1 F1 Score: 0.004009548364281089
ROUGE-2 Precision: 0.9904534606205251
ROUGE-2 Recall: 0.0019872576084432113
ROUGE-2 F1 Score: 0.003966556670386309
ROUGE-L Precision: 1.0
ROUGE-L Recall: 0.0020088013752747437
ROUGE-L F1 Score: 0.004009548364281089
```

## Review Summarization using BERT

```
In [34]: from transformers import pipeline

# Load the summarization pipeline
summarizer = pipeline("summarization")

# Assuming you have the product_summaries dictionary from the previous block
# Product Review Summaries
product_summaries_bert = {}

# Display available product indices
product_indices = list(topic_summaries.keys())
print("Available Product Indices:", product_indices)

# Get user input for the desired product index
selected_product_index = int(input("Enter the product index you want to view (e.g., 0, 1, 2, ...): "))

# Check if the selected index is valid
if selected_product_index in product_indices:
    # Retrieve reviews and summary for the selected product index
    reviews = topic_summaries[selected_product_index]['cluster_reviews']['reviewText']
    combined_reviews = ' '.join(reviews)

    # Split the combined reviews into chunks of suitable length
    max_chunk_length = 512 # Adjust as needed based on the model's maximum sequence length
    chunks = [combined_reviews[i:i + max_chunk_length] for i in range(0, len(combined_reviews), max_chunk_length)]

    # Summarize each chunk separately
    summary_chunks = []
    for chunk in chunks:
        summary = summarizer(chunk, max_length=100, min_length=50, length_penalty=2.0, num_beams=4, early_stopping=True)
        summary_text = summary[0]['summary_text'] if 'summary_text' in summary[0] else summary[0]['summary']
        summary_chunks.append(summary_text)

    # Combine the summaries of chunks into one summary
    summary = ' '.join(summary_chunks)

    # Store the product summary
    product_summaries_bert[selected_product_index] = {'combined_reviews': combined_reviews, 'summary': summary}

    # Display product summary
    print(f"\nProduct {selected_product_index} Summary:")
    print(f"\n\nCombined Reviews: {product_summaries_bert[selected_product_index]['combined_reviews']}")
    print(f"\n\nSummarized Review: {product_summaries_bert[selected_product_index]['summary']}")
else:
    print("Invalid product index. Please select a valid index.")
```

## Metrics for Summarization using BERT

```
In [35]: from nltk.translate.bleu_score import sentence_bleu
from nltk.translate.bleu_score import SmoothingFunction
import nltk

nltk.download('punkt')

def calculate_bleu_score(reference, hypothesis):
    reference = [reference.split()]
    hypothesis = hypothesis.split()
    smoothing_function = SmoothingFunction().method4
    score = sentence_bleu(reference, hypothesis, smoothing_function=smoothing_function)
    return score

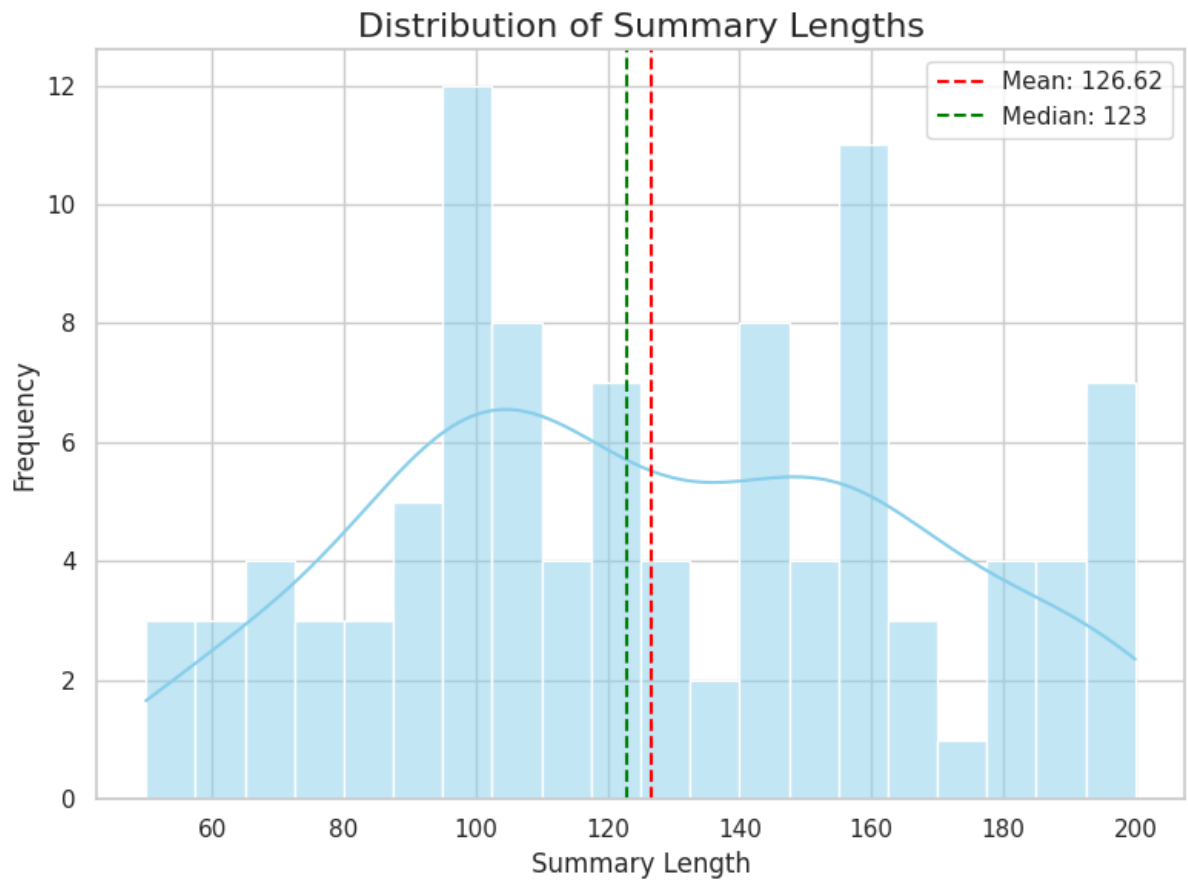
# Reference summary
reference_summary = combined_reviews

# Summarized review
summarized_review = summary

# Calculate BLEU score
score = calculate_bleu_score(reference_summary, summarized_review)

# Print BLEU score
print("BLEU Score:", score)
```

BLEU Score: 0.23845552768565806



Word Cloud:



## Product Comparison:

### 5. Product Comparison

```
In [38]: from difflib import SequenceMatcher
from textblob import TextBlob # Make sure to install the library: pip install textblob

def get_sentiment(text):
    analysis = TextBlob(text)
    return analysis.sentiment.polarity

def compare_products(product_summaries, index1, index2):
    # Get the summaries for the chosen products
    summary1 = product_summaries.get(index1, {}).get('summary', '')
    summary2 = product_summaries.get(index2, {}).get('summary', '')

    # Calculate sentiment scores
    sentiment_score1 = get_sentiment(summary1)
    sentiment_score2 = get_sentiment(summary2)

    # Use SequenceMatcher to find similarities
    matcher = SequenceMatcher(None, summary1, summary2)
    match_ratio = matcher.ratio()

    # Print the summaries
    print(f"\nProduct {index1} Summary:")
    print(f"\n{summary1}\n")
    print(f"Sentiment Score: {sentiment_score1:.2f}")

    print(f"\nProduct {index2} Summary:")
    print(f"\n{summary2}\n")
    print(f"Sentiment Score: {sentiment_score2:.2f}")

    # Print the comparison results
    print("\nComparison Results:")
    print(f"Similarities: {match_ratio:.2%}")
    print("Differences:")
    for op, i1, i2, j1, j2 in matcher.get_opcodes():
        if op == 'equal':
            continue
        elif op == 'insert':
            print(f" - Added: {summary2[j1:j2]}")
        elif op == 'delete':
            print(f" - Removed: {summary1[i1:i2]}")
        elif op == 'replace':
            print(f" - Replaced: {summary1[i1:i2]} with {summary2[j1:j2]}")

    # Compare sentiment scores
    if sentiment_score1 > sentiment_score2:
        print(f"\nProduct {index1} has a higher positive sentiment.")
    elif sentiment_score1 < sentiment_score2:
        print(f"\nProduct {index2} has a higher positive sentiment.")
    else:
        print("\nBoth products have similar sentiment scores.")

# Example usage:
# Choose two product indices
index_product1 = int(input("Enter the index of the first product: "))
index_product2 = int(input("Enter the index of the second product: "))

# Call the function to compare the chosen products
compare_products(product_summaries, index_product1, index_product2)
```



Enter the index of the first product: 1  
Enter the index of the second product: 2

#### Product 1 Summary:

I wouldn't mind, if there was \_anything\_ visually appealing about what is making my system drag its feet, but sadly all these bell-n-whistles simply serve to reveal the underlying clunkiness and awfulness of this outdated software. SAVE YOUR MONEY AND BUY ZONEALARM 5 At all unless you want to pay 30-60 dollars if their "protective" system won't remove it.. Now I think I will look for a virus protection system with actual tech support.. Don't know about McAfee and have heard negative about them but they do offer free online chat help at least> the below have free phone support: Trend Micro PC-cillin Internet Security 2005 got number one ratings on consumer reports however I don't know how much that says since Norton got second place.....It's amazing if you go look at reviews here [...] You will find most people are not pleased with their antivirus programs and makes it appear almost as if antivirus programs are nothing but problems in general LOL, but anyways thus far PC-cillin Internet Security 2005 looks the best. This latest version  
<a data-hook="product-link-linked" class="a-link-normal" href="/Trend-Micro-Titanium-Maximum-Security-2011-3-User-Download/dp/B0045THDCA/ref=cm\_cr\_arp\_d\_rvw\_txt?ie=UTF8">Trend Micro Titanium Maximum Security 2011 - 3 User [Download</a>] has so far proven to be the best they've ever produced and I can recommend it without any reservation at all. the first time I bought this Eset Nod32, was in hopes of fixing my favorite windows XP computer that managed to become LOADED with bots, worms and total trash whilst using brand X antivirus programs. Most of the antivirus programs lie and call cookies detections so they can justify their existence, because your average user often doesn't realize that finding nothing may be a very good thing indeed.

Sentiment Score: 0.24

#### Product 2 Summary:

And you can upgrade to the newest version for free at Kasperky's web site. Best antivirus I have ever used great antivirus I have tried many from Symantec to Trend micro and avg has always been the best I can not recommend this enough its the best by far detects virus malware ect before any other software Have tried other antivirus but always seem to come back to Norton!! Best Antivirus I have ever used this my 2nd year of using it, bought as download from Amazon . If you cannot afford to buy any other Antivirus then this is the BEST option. Used this for 3 plus years on 2 computers and have never had a virus or any problems.

Sentiment Score: 0.43

#### Comparison Results:

Similarities: 1.07%

Differences:

- Added: And you can upgrade to the newest version for free at Kasperky's web site. Best antivirus
- Added: have ever used great antivirus I have tried many from Symantec to Trend micro and avg has al
- Replaced: ou with ays been the best I can not recommend this enough its the best by far detects virus ma
- Replaced: dn't mind with ware ect before any other software Have tried other antivirus but always seem to come back to Norton!! Best Antivirus I have ever used this my 2nd year of using it
- Replaced: if there with bought as do
- Replaced: as \_anything\_ visua with n
- Replaced: ly a with oad from Amazon . If you cannot afford to buy any other Antivirus then this is the BEST o
- Replaced: peeling about what is making my system drag its feet, but sadly all these bell-n-whistles sim with tion. Used th is for 3
- Replaced: y serve to revea with us years on 2 computers and have never had a virus or any prob
- Replaced: the underlying clunkiness and awfulness of this outdated software with ems
- Removed: SAVE YOUR MONEY AND BUY ZONEALARM 5 At all unless you want to pay 30-60 dollars if their "protective" system won't remove it.. Now I think I will look for a virus protection system with actual tech support.. Don't know about McAfee and have heard negative about them but they do offer free online chat help at least> the below have free phone support: Trend Micro PC-cillin Internet Security 2005 got number one ratings on consumer reports however I don't know how much that says since Norton got second place.....It's amazing if you go look at reviews here [...] You will find most people are not pleased with their antivirus programs and makes it appear almost as if antivirus programs are nothing but problems in general LOL, but anyways thus far PC-cillin Internet Security 2005 looks the best. This latest version  
<a data-hook="product-link-linked" class="a-link-normal" href="/Trend-Micro-Titanium-Maximum-Security-2011-3-User-Download/dp/B0045THDCA/ref=cm\_cr\_arp\_d\_rvw\_txt?ie=UTF8">Trend Micro Titanium Maximum Security 2011 - 3 User [Download</a>] has so far proven to be the best they've ever produced and I can recommend it without any reservation at all. the first time I bought this Eset Nod32, was in hopes of fixing my favorite windows XP computer that managed to become LOADED with bots, worms and total trash whilst using brand X antivirus programs. Most of the antivirus programs lie and call cookies detections so they can justify their existence, because your average user often doesn't realize that finding nothing may be a very good thing indeed.

Product 2 has a higher positive sentiment.

```
In [39]: from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity

# Assuming 'df_no_outliers' is the DataFrame containing processed reviewText and product information

# Selecting relevant columns for product comparison
selected_columns = ['asin', 'processed_reviewText']

# Extracting product IDs and their processed review texts
product_reviews = df_no_outliers[selected_columns]

# Grouping reviews by product (ASIN)
grouped_reviews = product_reviews.groupby('asin')['processed_reviewText'].apply(' '.join).reset_index()

# Applying TF-IDF Vectorization to processed review texts
vectorizer = TfidfVectorizer(stop_words='english')
tfidf_matrix = vectorizer.fit_transform(grouped_reviews['processed_reviewText'])

# Calculate cosine similarity between products based on their reviews
cosine_sim = cosine_similarity(tfidf_matrix, tfidf_matrix)

# Calculate metrics for product comparison
# Assuming 'asin' column contains product IDs
product_indices = grouped_reviews['asin'].tolist()

# Choose a specific product index for comparison (e.g., 0 for the first product)
selected_product_index = 0

# Retrieve similarity scores for the selected product compared to others
selected_product_similarities = cosine_sim[selected_product_index]

# Sort similarity scores in descending order and exclude the similarity to the same product
similar_products_indices = sorted(range(len(selected_product_similarities)), key=lambda i: selected_product_similarities[i], reverse=True)

# Select top similar products and their similarity scores
num_top_similar_products = 5 # Adjust as needed
top_similar_products = [(product_indices[i], selected_product_similarities[i]) for i in similar_products_indices[:num_top_similar_products]]

# Display top similar products and their similarity scores
print(f"Top {num_top_similar_products} Products Similar to Product {product_indices[selected_product_index]}:")
for product, similarity_score in top_similar_products:
    print(f"Product: {product} | Similarity Score: {similarity_score}")
```

Top 5 Products Similar to Product 0030672120:

Product: 0007N6W22	Similarity Score: 1.0000000000000002
Product: 1597461296	Similarity Score: 0.8077273088781065
Product: 0012P4Z00G	Similarity Score: 0.7301176106226555
Product: 000XGG7RC0	Similarity Score: 0.6948635066087674
Product: 1602552304	Similarity Score: 0.4438649935038219

### Mean Average Precision

```
In [40]: def calculate_average_precision(similarity_scores):
# Assuming you have ground truth relevance scores (e.g., manually labeled relevance)
# Here, set all relevance scores to 1 for demonstration purposes
relevance_scores = [1] * len(similarity_scores) # Binary relevance scores (1 for relevant, 0 for irrelevant)

# Sort similarity scores and corresponding relevance scores in descending order
sorted_scores_and_labels = sorted(zip(similarity_scores, relevance_scores), key=lambda x: x[0], reverse=True)

# Calculate average precision
num_relevant_items = sum(relevance_scores)
average_precision = 0.0
num_retrieved_relevant_items = 0

for i, (score, label) in enumerate(sorted_scores_and_labels):
    if label == 1:
        num_retrieved_relevant_items += 1
        precision_at_k = num_retrieved_relevant_items / (i + 1) # Precision at the current position
        average_precision += precision_at_k

if num_relevant_items == 0:
    return 0.0 # Handle the case when there are no relevant items

average_precision /= num_relevant_items
return average_precision

# Calculate Average Precision for the selected product's similarity scores
average_precision_score = calculate_average_precision(selected_product_similarities)
print(f"Mean Average Precision (MAP) for the selected product: {average_precision_score}")
```

Mean Average Precision (MAP) for the selected product: 1.0

#### Normalized Discounted Cumulative Gain (NDCG):

```
In [42]: from sklearn.metrics import ndcg_score

def calculate_ndcg(relevance_scores, similarity_scores):
    # Check if the range of relevance scores is zero
    if max(relevance_scores) - min(relevance_scores) == 0:
        normalized_relevance = [1.0] * len(relevance_scores) # Assigning 1.0 if the range is zero
    else:
        # Normalize scores to range [0, 1]
        normalized_relevance = [(score - min(relevance_scores)) / (max(relevance_scores) - min(relevance_scores)) for score in relevance_scores]

    # Calculate NDCG using scikit-learn's ndcg_score function
    ndcg = ndcg_score([normalized_relevance], [similarity_scores])
    return ndcg

# Assuming you have ground truth relevance scores and similarity scores for the selected product
# Relevance scores (manually labeled relevance) set to 1 for demonstration purposes
relevance_scores = [1] * len(selected_product_similarities)
# Use the similarity scores from the previous snippet
similarity_scores = selected_product_similarities

# Calculate NDCG for the selected product's similarity scores and relevance scores
ndcg_score = calculate_ndcg(relevance_scores, similarity_scores)
print(f"Normalized Discounted Cumulative Gain (NDCG) for the selected product: {ndcg_score}")
```

Normalized Discounted Cumulative Gain (NDCG) for the selected product: 1.0000000000000009

## Product Recommendation:

### 5. Product Recommendations

```
In [43]: from surprise import Dataset, Reader
from surprise.model_selection import train_test_split
from surprise import SVD
from surprise import accuracy

# Load the dataset
reader = Reader(rating_scale=(1, 5))
data = Dataset.load_from_df(antivirus_reviews[['reviewerID', 'asin', 'overall']], reader)

# Split the dataset into training and test sets
trainset, testset = train_test_split(data, test_size=0.2, random_state=42)

# Create the SVD (Singular Value Decomposition) model
model = SVD(n_factors=100, random_state=42)

# Train the model on the training set
model.fit(trainset)

# Make predictions for the test set
predictions = model.test(testset)

# Calculate RMSE (Root Mean Squared Error) for accuracy evaluation
rmse = accuracy.rmse(predictions)
print(f"RMSE: {rmse}")
```

RMSE: 1.4561  
RMSE: 1.4560943976070322

## Generating Recommendations for Users

```
In [44]: # Replace 'user_id' with the ID of the user for whom you want to generate recommendations
user_id = 'A2480RQ2LF9LUI'

# Get all unique product IDs (asin)
unique_products = antivirus_reviews['asin'].unique()

# Prepare a list of products not yet rated by the user
products_not_rated_by_user = []
for product_id in unique_products:
    if not trainset.knows_user(user_id) or trainset.ur[trainset.to_inner_uid(user_id)]:
        products_not_rated_by_user.append(product_id)

# Predict ratings for the products not rated by the user
predicted_ratings = [model.predict(user_id, product_id).est for product_id in products_not_rated_by_user]

# Combine the products and their predicted ratings
recommendations = list(zip(products_not_rated_by_user, predicted_ratings))

# Sort the recommendations by predicted ratings in descending order
sorted_recommendations = sorted(recommendations, key=lambda x: x[1], reverse=True)

# Display the top N recommended products
top_n = 10
top_recommendations = sorted_recommendations[:top_n]
print(f"Top {top_n} Recommendations for User {user_id}:")
for idx, (product_id, predicted_rating) in enumerate(top_recommendations, start=1):
    print(f"{idx}. Product ID: {product_id}, Predicted Rating: {predicted_rating}")

Top 10 Recommendations for User A2480RQ2LF9LUI:
1. Product ID: B00EZKNY8G, Predicted Rating: 4.58071344287636
2. Product ID: B00EZK08J0, Predicted Rating: 4.556407890948249
3. Product ID: B0056CZC2S, Predicted Rating: 4.54501131550186
4. Product ID: B00H9A6004, Predicted Rating: 4.541239333862853
5. Product ID: B005CSF154, Predicted Rating: 4.506501200888183
6. Product ID: B00EZKNYWC, Predicted Rating: 4.486671564976551
7. Product ID: B005CSF1JK, Predicted Rating: 4.4857642703819876
8. Product ID: B015IH4AZW, Predicted Rating: 4.474951855053922
9. Product ID: B003WT1KHI, Predicted Rating: 4.474219234150831
10. Product ID: B008F5THLA, Predicted Rating: 4.468026131895832
```

## Content-Based Filtering (TF-IDF and Cosine Similarity)

```
In [45]: from sklearn.feature_extraction.text import TfidfVectorizer

# Handle lists in 'description' column
def concatenate_descriptions(descriptions):
    # Ensure all values are strings and concatenate them
    cleaned_descriptions = [str(desc) for desc in descriptions if isinstance(desc, str)]
    return ' '.join(cleaned_descriptions)

# Concatenate item descriptions by 'asin'
item_descriptions = df_no_outliers.groupby('asin')['description'].apply(concatenate_descriptions).reset_index()

# Calculate cosine similarity between items based on their descriptions
item_item_similarity_tfidf = cosine_similarity(tfidf_matrix, tfidf_matrix)
```

## Combining Collaborative and Content-Based Filtering

```
In [46]: # Combine item-item similarity from both approaches (you can adjust weights based on your preference)
alpha = 0.5 # Adjust the weight for collaborative filtering
combined_similarity = alpha * item_item_similarity_tfidf + (1 - alpha) * item_item_similarity_tfidf

# Get recommendations for a specific product
selected_product_index = 0 # Index of the selected product for which you want recommendations

# Sort products by similarity scores to the selected product
similar_products_indices = combined_similarity[selected_product_index].argsort()[::-1]

# Top N similar products excluding itself (selected product)
top_n = 5
similar_products_indices = similar_products_indices[1:top_n + 1]

# Get the recommended products
recommended_products = item_descriptions.iloc[similar_products_indices]['asin'].values
print("Recommended Products:", recommended_products)

Recommended Products: ['B007N6WPZ2' '1597461296' 'B012P4ZODG' 'B00XGG7RCO' 'B0018TMV6S']
```

## Evaluation product recommendation

```
In [47]: from surprise import Dataset, Reader
from surprise.model_selection import train_test_split
from surprise import SVD
from surprise import accuracy

# Load the dataset
reader = Reader(rating_scale=(1, 5))
data = Dataset.load_from_df(antivirus_reviews[['reviewerID', 'asin', 'overall']], reader)

# Split the dataset into training and test sets
trainset, testset = train_test_split(data, test_size=0.2, random_state=42)

# Create the SVD (Singular Value Decomposition) model
model = SVD(n_factors=100, random_state=42)

# Train the model on the training set
model.fit(trainset)

# Make predictions for the test set
predictions = model.test(testset)

# Calculate RMSE (Root Mean Squared Error) for accuracy evaluation
rmse = accuracy.rmse(predictions)
print(f"RMSE: {rmse}")

# Calculate MAE (Mean Absolute Error) for accuracy evaluation
mae = accuracy.mae(predictions)
print(f"MAE: {mae}")
```

```
RMSE: 1.4561
RMSE: 1.4560943976070322
MAE: 1.1937
MAE: 1.1937164186830747
```

