



***VIEWER'S REVIEWS
SENTIMENT ANALYSIS
ON IMDB'S MADAM
WEB MOVIE***

Group members:

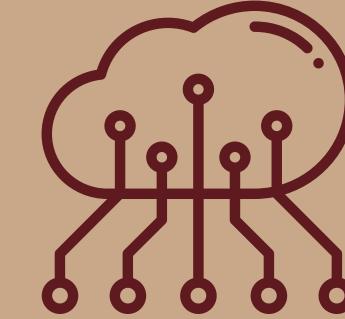
1. Muhammad Khairin Asnawi bin Rosli (ISO1082068)
2. Mohd Taufiq Eizaz bin M Jamaludin (ISO1080910)
3. Haliki Bachar Djimet (ISO1081569)
4. Ayman Fikry bin Asmajuda (ISO1081779)

Group Project (Part 2)

PREPARED FOR: TS. NUR LAILA BTE AB GHANI

[Link to Video Presentation](#)

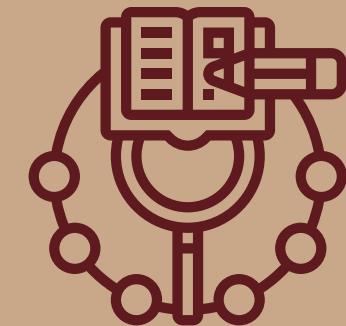
TOPIC OUTLINE:



Data Gathering



Text Preprocessing



Modelling and Evaluation.

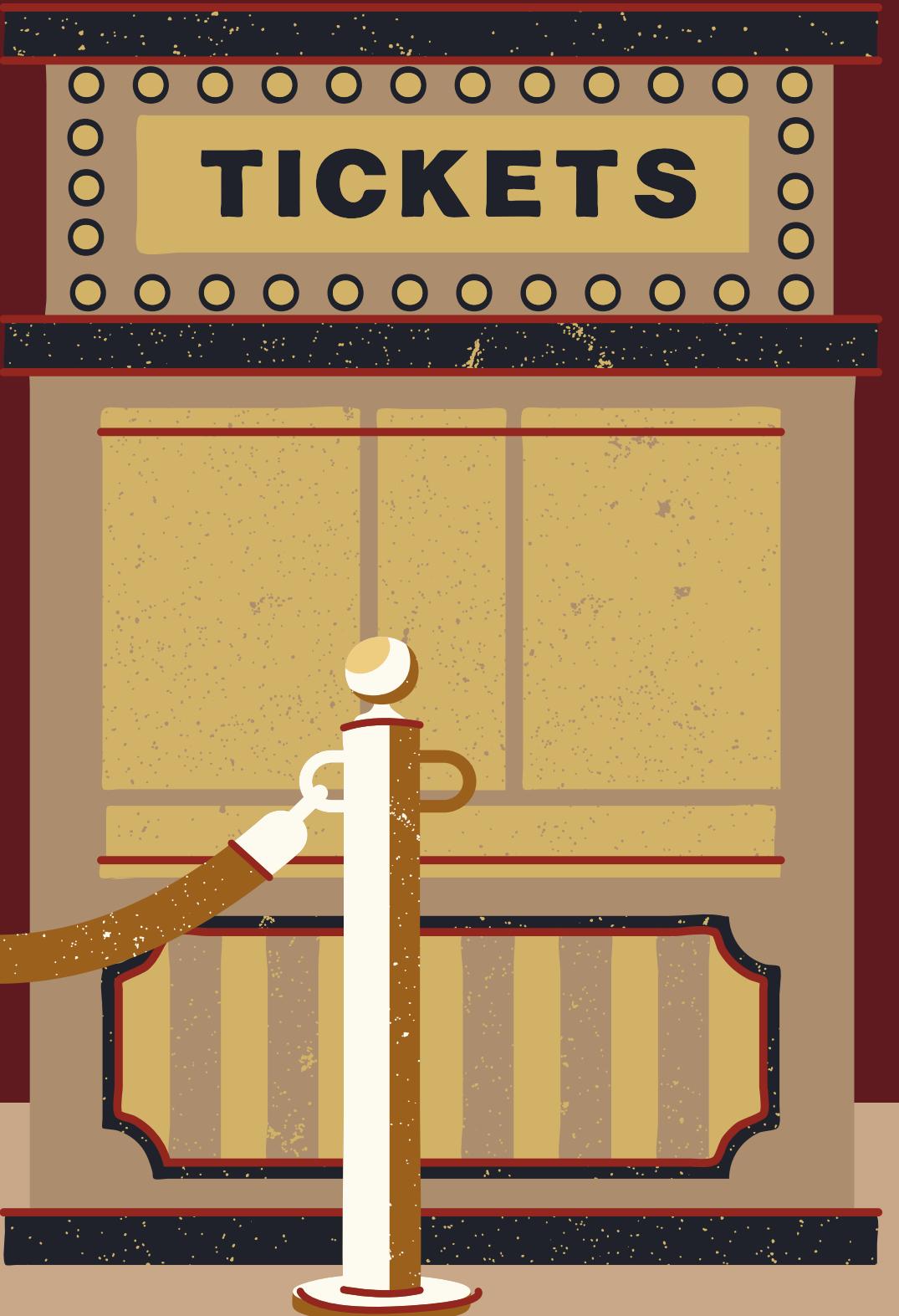


Results Visualization



Key Findings

DATA GATHERING



DATA GATHERING

1. Method used for Data Gathering:

- The method used is web scraping for data extraction

2. Target Data Source:

- Viewer reviews for the movie "Madam Web" on IMDb

3. Tool Utilized:

- Python package BeautifulSoup

4. Process:

- Parsed and collected textual data from IMDb user reviews section

5. Outcome:

- Collected a thousand plus reviews on IMDB
- These were separated by name, score, title of their review and comment underneath



TEXT PREPROCESSING

Data preprocessing steps:-

- Data Selection:
 - Selected relevant columns: Score and Text
- Data Cleaning:
 - Removed duplicate rows
 - Cleaned text by removing HTML tags, URLs, special characters, digits, and extra whitespaces
 - Converted text to lowercase
- Handling Missing Values:
 - Checked for missing values
 - Removed rows with missing 'Score' values
 - Tokenization & Stopword Removal:
 - Tokenized text into individual words using NLTK
 - Removed stop words
- Lemmatization:
 - Applied lemmatization to tokens using WordNet Lemmatizer
- Reconstruction:
 - Joined lemmatized tokens back into sentences



MODELLING AND EVALUATION



MODELLING : TEXTBLOB

```
In [26]: from textblob import TextBlob  
blob = TextBlob(mystring)  
blob.sentiment.polarity
```

Out[26]: 0.025581905683177322

```
In [27]: blob.sentiment.subjectivity
```

Out[27]: 0.5279996622967936

In [28]: blob.words

MODELLING : LEXICON-BASED APPROACH

```
In [29]: # Download the VADER Lexicon
nltk.download('vader_lexicon')

# Define a function to assign sentiment labels based on the 'Score' column
def assign_sentiment(score):

    if score >= 7:
        return 'Positive'
    elif score <= 3:
        return 'Negative'
    else:
        return 'Neutral'

# Assign sentiment labels based on the 'Score' column
data['Sentiment'] = data['Score'].apply(assign_sentiment)

# Initialize the sentiment analyzer
sid = SentimentIntensityAnalyzer()

# Calculate sentiment scores for each review
data['Lexicon_Sentiment'] = data['Preprocessed_Text'].apply(lambda x: sid.polarity_scores(x)['compound'])

# Map sentiment scores to labels
data['Lexicon_Sentiment_Label'] = data['Lexicon_Sentiment'].apply(lambda x: 'Positive' if x > 0 else ('Negative' if x < 0 else 'Neutral'))

# Evaluate the lexicon-based approach
lexicon_accuracy = accuracy_score(data['Sentiment'], data['Lexicon_Sentiment_Label'])
print("Accuracy of the Lexicon-based Approach:", lexicon_accuracy)
```

[nltk_data] Downloading package vader_lexicon to
[nltk_data] C:\Users\Administrator\AppData\Roaming\nltk_data...
[nltk_data] Package vader_lexicon is already up-to-date!

Accuracy of the Lexicon-based Approach: 0.4727822580645161

MODELLING: MACHINE-LEARNING BASED APPROACH

x_train_tfidf shape: (793, 10753)

x_test_tfidf shape: (199, 10753)

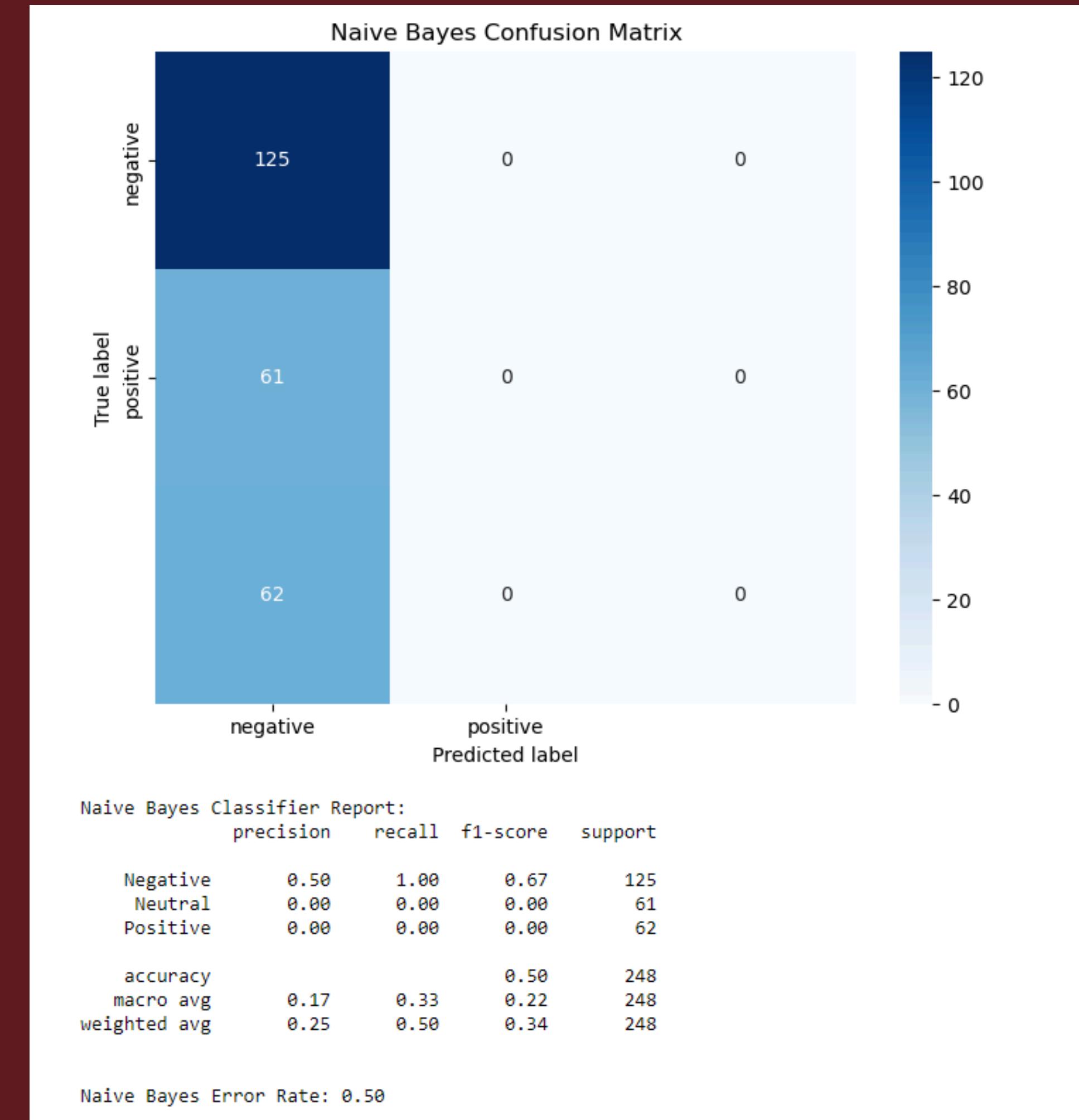
Naive Bayes Classifier Accuracy: 0.49748743718592964

	precision	recall	f1-score	support
Negative	0.50	1.00	0.66	99
Neutral	0.00	0.00	0.00	50
Positive	0.00	0.00	0.00	50
accuracy			0.50	199
macro avg	0.17	0.33	0.22	199
weighted avg	0.25	0.50	0.33	199

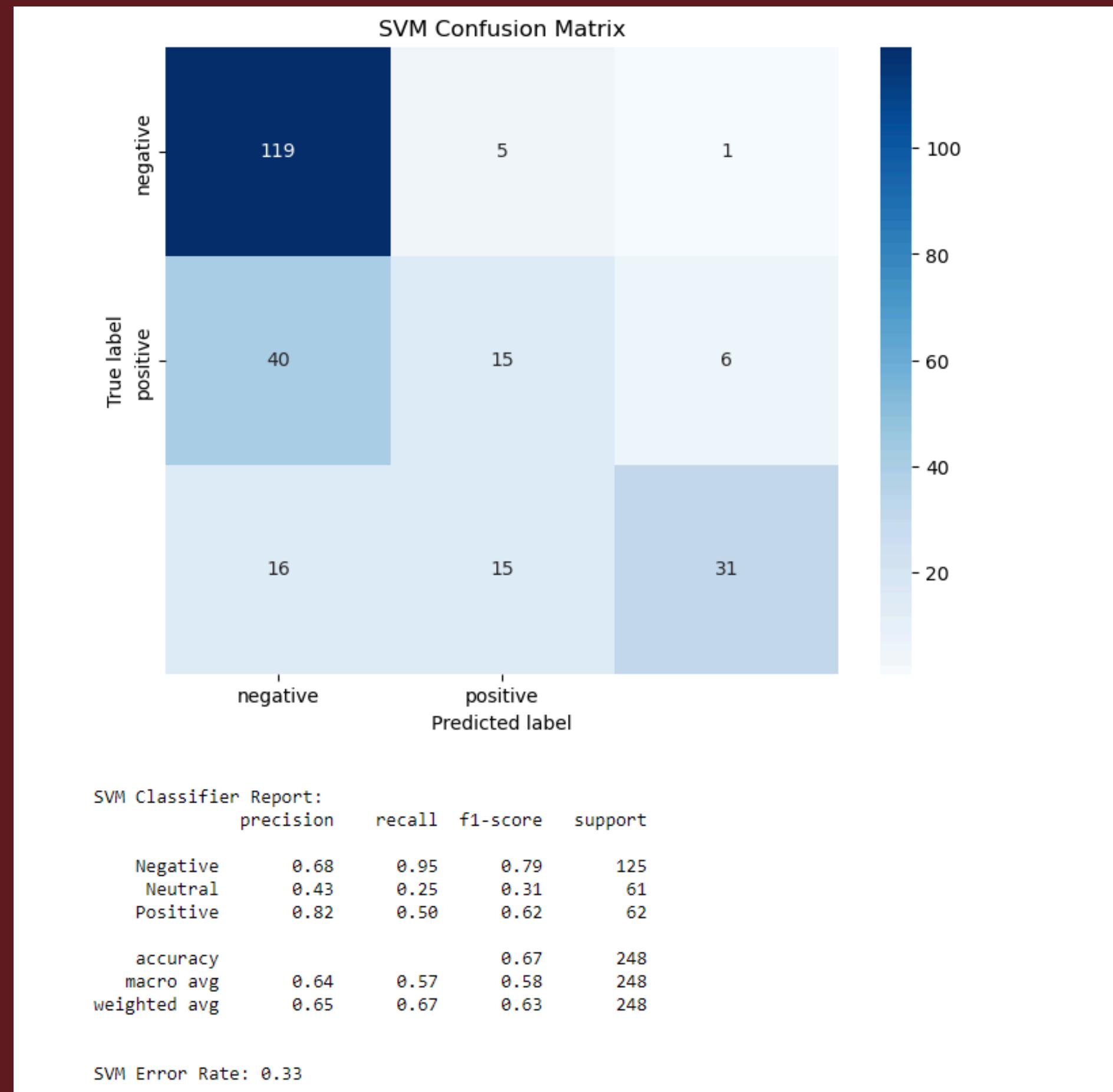
Support Vector Machine (SVM) Classifier Accuracy: 0.7185929648241206

	precision	recall	f1-score	support
Negative	0.73	0.97	0.83	99
Neutral	0.55	0.32	0.41	50
Positive	0.79	0.62	0.70	50
accuracy			0.72	199
macro avg	0.69	0.64	0.65	199
weighted avg	0.70	0.72	0.69	199

MODELLING : NB CONFUSION MATRIX



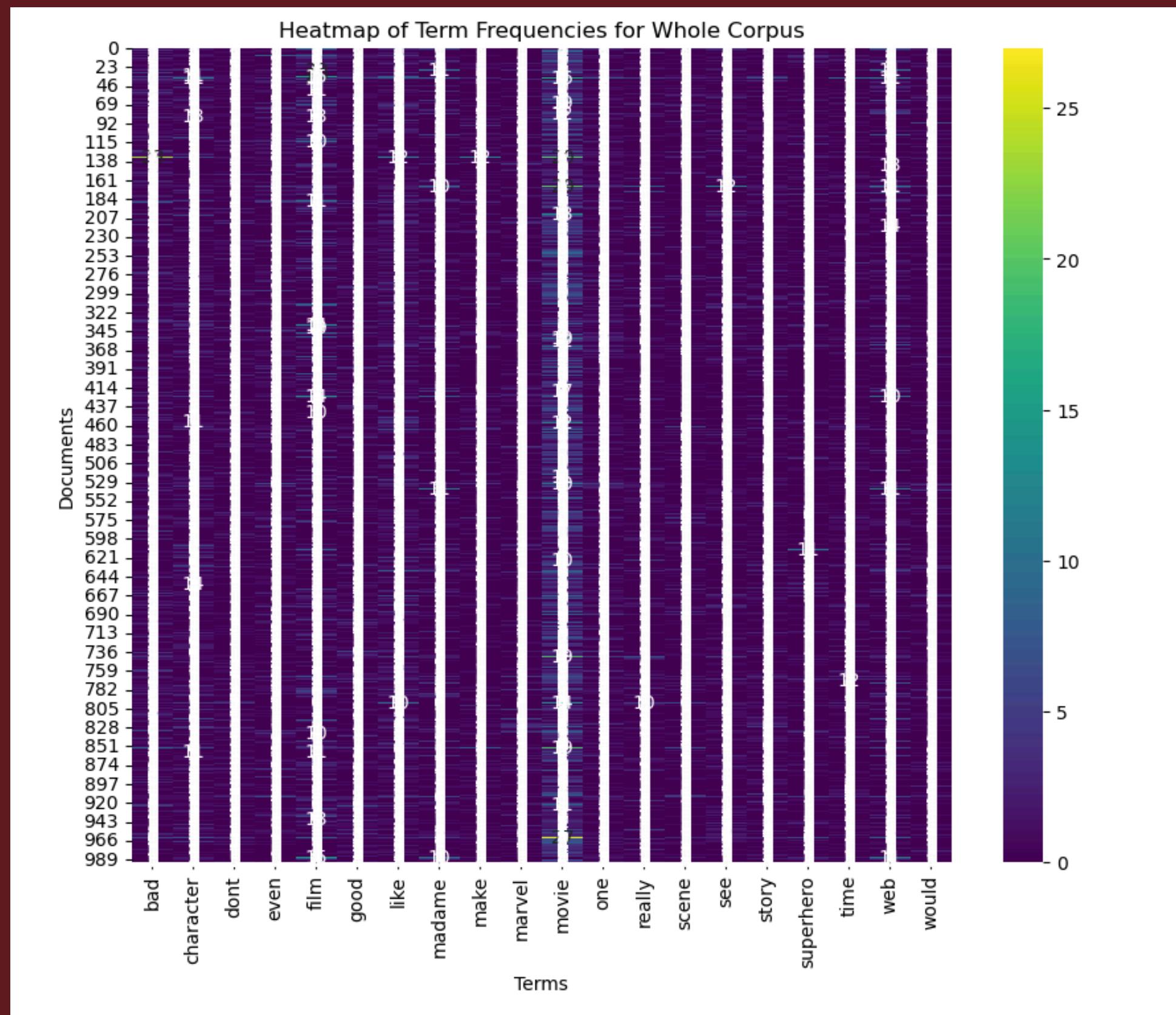
MODELLING : SVM CONFUSION MATRIX



RESULTS VISUALIZATION

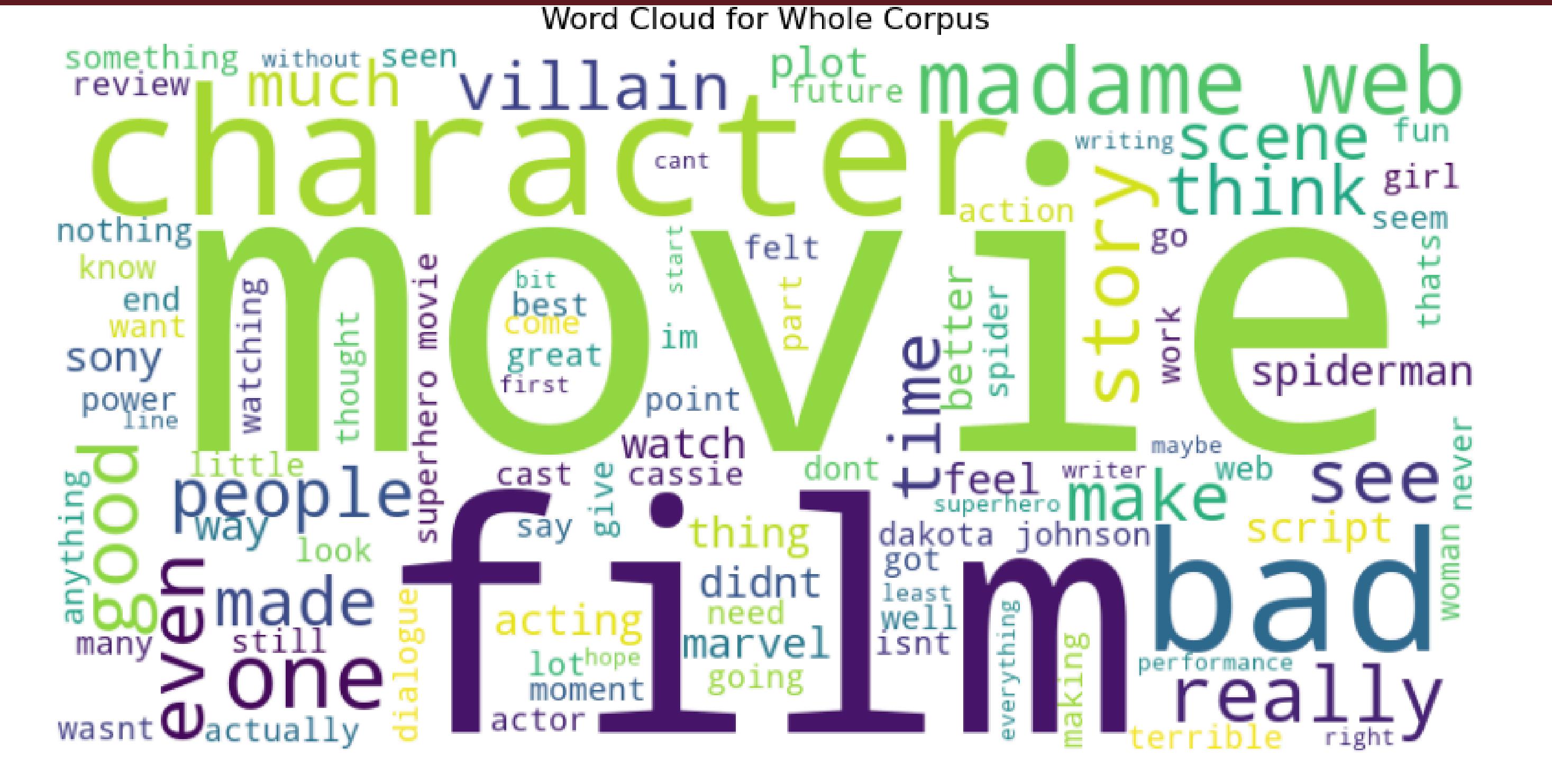


HEATMAP



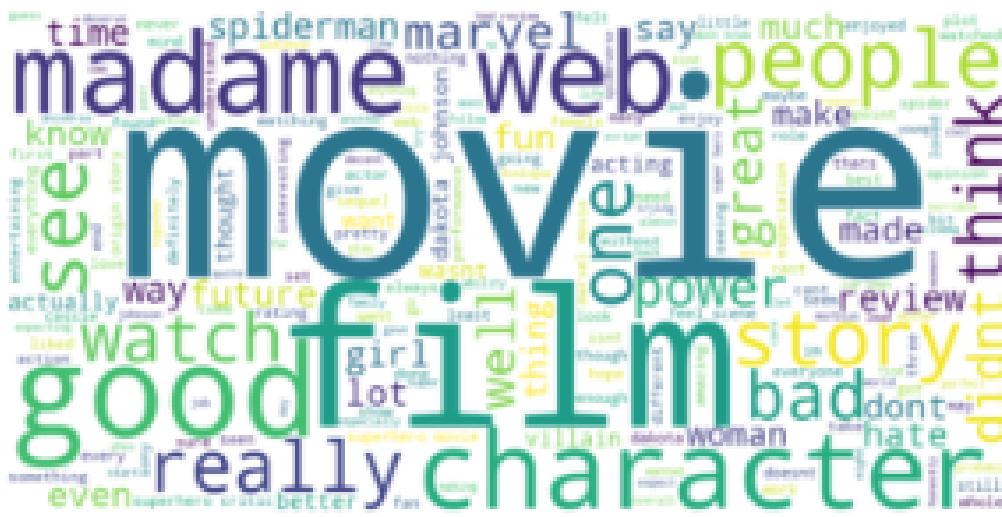
WORDCLOUD

Word Cloud for Whole Corpus



WORLD CLOUD DIVIDED BY SENTIMENTS

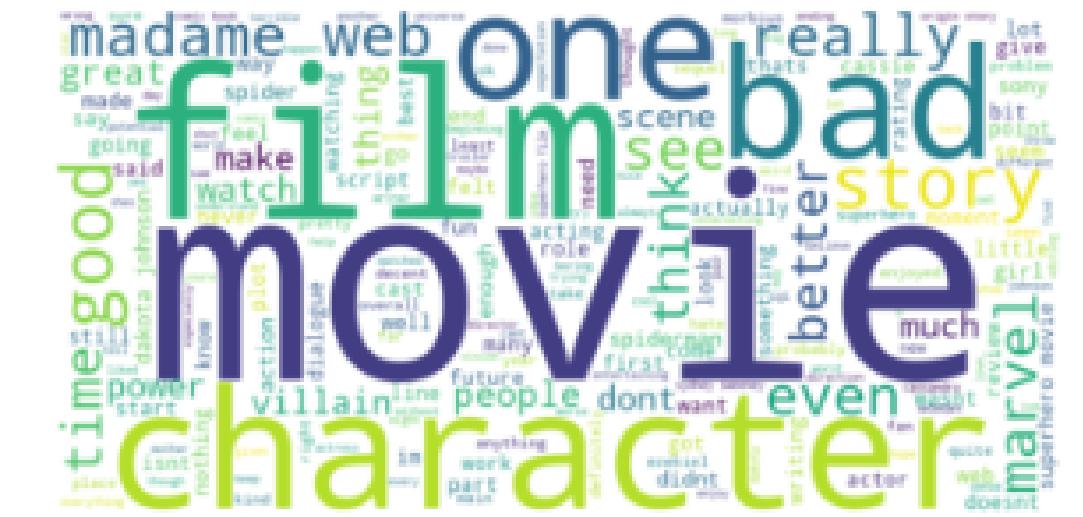
Positive Sentiment



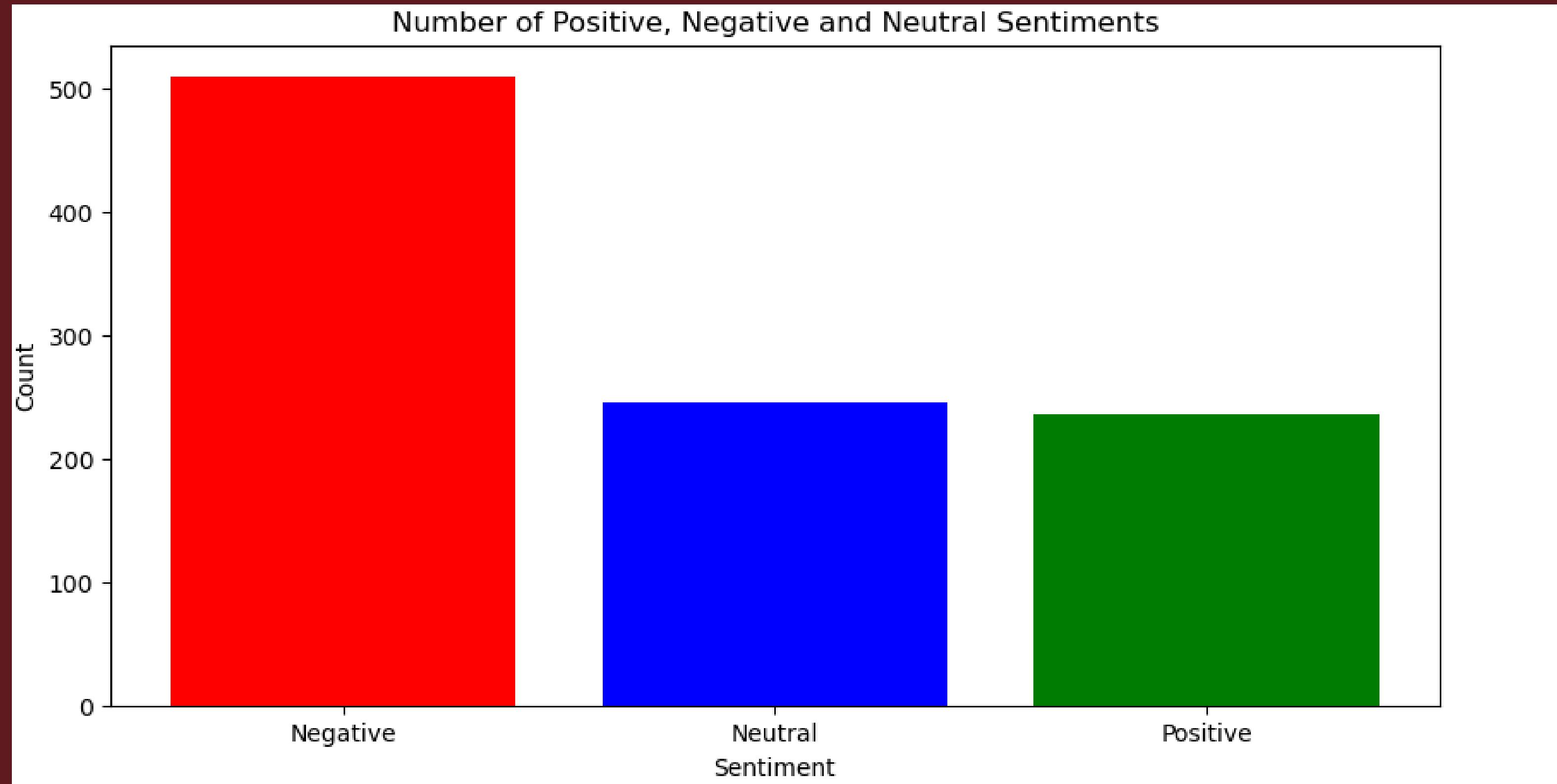
Negative Sentiment



Neutral Sentiment



BAR CHART COUNT OF DIFFERENT SENTIMENTS



KEY FINDINGS

- Data Gathering:

Method: Web scraping

Tool: BeautifulSoup

Source: IMDb user reviews for "Madam Web"

Outcome: Collected 1000+ reviews

- Text Preprocessing:

Data Selection: Selected Score and Text columns

Data Cleaning: Removed duplicates, cleaned HTML tags, URLs, special characters, digits, and converted text to lowercase

Handling Missing Values: Removed rows with missing scores

Tokenization & Stopword Removal: Tokenized text and removed stop words using NLTK

Lemmatization: Applied lemmatization and reconstructed sentences

- Modeling & Evaluation:

Models Used: Various sentiment analysis models

Evaluation: Performed cross-validation and measured accuracy

- Results:

Heat Map: Displayed sentiment distribution across reviews

THANK YOU FOR LISTENING

