

SENTIMENT ANALYSIS AND TOPIC MODELING ON SOCIAL MEDIA AND NEWS PLATFORMS USING VADER AND RoBERTa

PROJECT REPORT

Submitted by

HARIHARAN SRIRAM

310623205044

Easwari Engineering College

Ramapuram, Chennai – 600089

Submitted as a part of:

Internship at

SIFY TECHNOLOGIES LIMITED

Tidel Park, Tharamani, Chennai - 600113

July 2025



ACKNOWLEDGEMENT

I sincerely express my gratitude to my mentors and colleagues at Sify Technologies and my academic institution who provided me continuous support and guidance during the course of this project.

I am deeply grateful for the accessible resources from Google Colab, HuggingFace Transformers, NLTK, YouTube's API and Reddit's API, which have been instrumental in the successful completion of this project.

Finally, I would also like to thank my peers and reviewers for their valuable feedback during this internship.

ABSTRACT

This project explores the implementation of **Sentiment Analysis** and **Topic Modeling** on data gathered from various digital platforms including **YouTube, Reddit, Medium, and Bing News**. The goal is to offer a consolidated tool capable of collecting, analyzing, and visualizing public sentiment and topic trends across these platforms.

For sentiment analysis, two well-established models were used:

- **VADER** (Valence Aware Dictionary and Sentiment Reasoner)
- **RoBERTa** (Robustly Optimized BERT Approach)

The project additionally integrates **Topic Modeling using NMF (Non-Negative Matrix Factorization)** to extract meaningful insights from the gathered data.

Outcomes include:

- Clean, deduplicated datasets.
- Sentiment analysis, classification and scores.
- Topic clusters with key discussion themes.
- Sentiment versus Topic intensity split-up
- Informative visualizations.

TABLE OF CONTENTS

S.NO	TOPIC	PAGE NUMBER
1.	Introduction	5
2.	Problem Statement	6
3.	Objective	7
4.	Technologies Used	8
5.	Methodology	9
6.	Data Collection	10
7.	Sentiment Analysis	12
8.	Topic Modeling	15
9.	Visualization	18
10.	Results & Discussion	24
12.	Conclusion	25
13.	Future Scope	25
14.	References	26

INTRODUCTION

In today's digital landscape, **public sentiment analysis** is critical for businesses, researchers, and governments to understand public opinion trends across social media and news platforms. This project aims to solve the gap of collecting data across multiple sources, analyzing it with robust AIML models, and presenting it in a user-friendly manner.

Platforms Scrapped:

- **YouTube** (comments)
- **Reddit** (posts)
- **Medium** (articles)
- **Bing News** (articles)

Models Used:

- **VADER** (Valence Aware Dictionary for Sentiment Reasoner)
- **RoBERTa** (Robustly Optimized BERT Approach)

This project combines **scraping**, **sentiment classification**, **topic extraction**, and **visualization** into a unified pipeline for efficient analysis.

PROBLEM STATEMENT

In the current digital era, a massive volume of user-generated content is produced daily across various online platforms such as **YouTube, Reddit, Medium, and news portals**. This content reflects valuable public opinions, sentiments, and trends which, if analyzed effectively, can provide deep insights for businesses, researchers, and individuals. However, **manually analyzing such vast amounts of textual data is not only impractical but also highly inefficient and time-consuming**.

The challenge lies in gathering reliable, deduplicated, and valid content from varied platforms and applying sophisticated models to extract actionable insights.

Traditional tools:

- Often limited to a single platform.
- Lack flexibility for deep sentiment or topic analysis.

This project bridges these gaps through:

- **Cross-platform data extraction.**
- **Advanced NLP models.**
- **Readable and insightful visualizations.**

Furthermore, the current solutions rarely provide a **seamless end-to-end pipeline that integrates scraping, deduplication, cleaning, and comparative analysis of sentiment results from multiple models**.

Additionally, platforms like Medium, Reddit, and YouTube each have their unique formats and restrictions for accessing data, making unified sentiment analysis across them challenging. Users seeking meaningful insights from these platforms often lack a unified, flexible, and easy-to-use tool that can automate this process.

OBJECTIVE

This project aims to **develop a comprehensive sentiment analysis framework** that:

- Scraps and collects data from platforms such as **YouTube, Reddit, Medium, and Bing News**.
- Cleans and deduplicates the collected data to ensure only **valid and unique articles/posts/comments are considered**.
- Applies both **VADER (rule-based)** and **RoBERTa (transformer-based)** models to perform sentiment analysis, offering users a comparative understanding of results from different approaches.
- Presents insights through **visualizations like count plots and heatmaps** to enhance interpretability.
- Provides the user with **an option to select the sentiment model** based on their specific needs, through an intuitive interface.

TECHNOLOGIES USED

➤ Languages & Tools:

- Python
- Google Colab

➤ Libraries:

NLTK, BeautifulSoup, requests, pandas, matplotlib, seaborn, transformers, tqdm, praw, sklearn, newspaper3k, torch

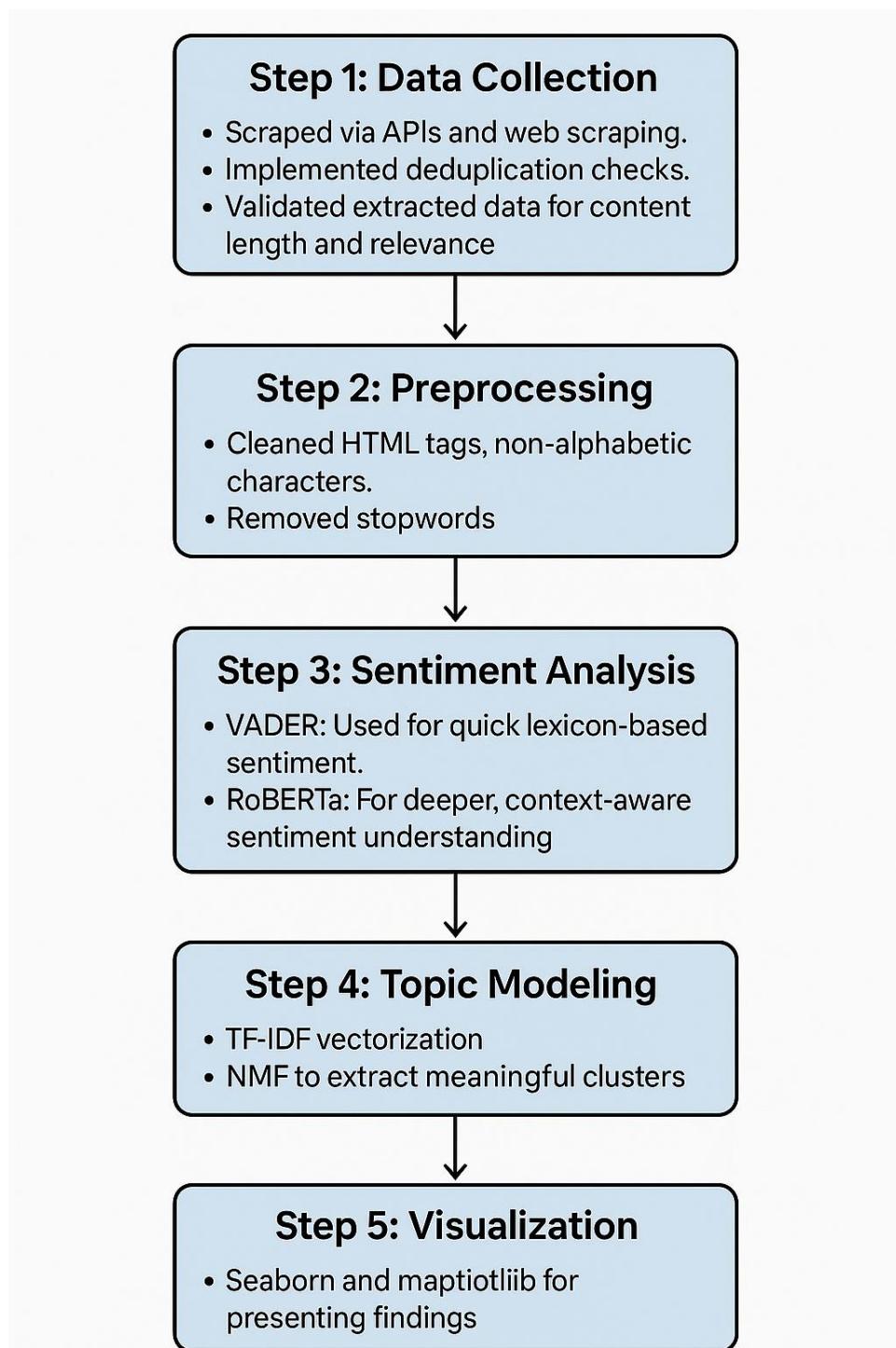
➤ Models:

- **VADER** (Lexicon-based)
- **RoBERTa** (Transformer-based via HuggingFace)

➤ APIs:

- YouTube Data API
- Reddit PRAW

METHODOLOGY



DATA COLLECTION

Platform	Type of Data	Reason for Selection
Reddit	Posts & Comments	Community-driven discussions covering diverse topics.
YouTube	Video Comments	Rich in informal user opinions, sentiment-heavy content.
Bing News	News Articles	Professional articles, global perspectives, informative tone.
Medium	Articles	Opinionated, long-form, structured articles.

Collection Methods

Reddit Data Collection

- Employed **PRAW (Python Reddit API Wrapper)** to gather posts and comments from relevant subreddits.
- Suggested subreddits based on keywords.
- Collected only **meaningful, non-duplicate content** with sufficient length.

YouTube Data Collection

- Used **YouTube Data API v3** to fetch video IDs and extract comments.
- Searched videos with **user-defined keywords** and ensured a **minimum number of comments** for quality.
- Handled API rate limits and removed duplicates.

Bing News Data Collection

- Scraped articles via **Bing Search** and **BeautifulSoup**.
- Focused on **titles, URLs, and article text**.
- Deduplicated to avoid repetitions and ensured content validity.

Medium Data Collection

- Collected articles using **Google Custom Search Engine (CSE)** and **BeautifulSoup**.
- Extracted **titles, URLs, and content bodies**.
- Ensured articles were unique and of sufficient length.

Scraping and preprocessing working (Reddit)

```
⌚ Scraping data...
⌚ Parsing 3/3 |██████████| 100%
✓ Parsing complete!
Data scrapping successfully completed

Total unique posts found: 149
✓ Reddit data saved successfully in 'reddit_full.csv'

✍ Preprocessing data...
🌐 Checking users... (this may take a few minutes)
✓ Bot filtering complete. Saved to cleaned_reddit_full.csv
```

Columns displaying the raw data (Text) and the pre-processed data (Final Text).

SENTIMENT ANALYSIS

Sentiment analysis is a core component of this project, helping to transform collected textual data into actionable insights by determining the polarity (positive, negative, or neutral) of opinions across different platforms. This project employs **two distinct approaches for sentiment analysis:**

1. **Rule-Based Analysis (VADER)**
2. **Transformer-Based Deep Learning Model (RoBERTa)**

By combining both approaches, the project ensures comprehensive coverage of sentiment detection, balancing between speed and accuracy.

VADER Sentiment Analysis

Overview:

VADER (Valence Aware Dictionary and Sentiment Reasoner) is a **lexicon and rule-based sentiment analysis tool** specifically designed for analyzing sentiments expressed in social media content. VADER is known for its simplicity and efficiency in handling informal language, slangs, and emojis common in platforms like Reddit and YouTube.

Why VADER?

- ✓ Lightweight and fast.
- ✓ Well-suited for short texts (comments, posts).
- ✓ No GPU dependency; ideal for large datasets needing quick analysis.

Implementation Details:

- **Tool Used:** vaderSentiment library in Python.
- **Process:**
 - Cleaned and pre-processed text from all platforms.
 - Used SentimentIntensityAnalyzer to compute sentiment scores.
 - Classified texts as **Positive, Negative, or Neutral** based on the compound score thresholds.
- **Output:** Each comment/post/article was tagged with a sentiment label and score.

RoBERTa Sentiment Analysis

Overview:

RoBERTa (A Robustly Optimized BERT Pretraining Approach) is a **state-of-the-art transformer-based model** from HuggingFace designed for natural language understanding tasks, including sentiment analysis. Specifically, this project used cardiffnlp/twitter-roberta-base-sentiment, a model fine-tuned for social media sentiment tasks.

Why RoBERTa?

- ✓ More accurate and nuanced understanding of complex text.
- ✓ Better at detecting sarcasm, context, and varied sentence structures.
- ✓ Suitable for both short and long texts.

Implementation Details:

- **Tools Used:** transformers library, AutoTokenizer, and AutoModelForSequenceClassification.
- **Hardware:** Enabled GPU acceleration (cuda) to ensure faster inference in Google Colab.
- **Process:**
 - Tokenized texts with a maximum length of 512 tokens.
 - Inference performed using pre-trained RoBERTa sentiment model.
 - **Softmax function** applied to logits to obtain probability distributions.
 - Sentiment determined by the class with the highest probability: **Positive, Neutral, or Negative.**
- **Output:** Sentiment labels and confidence scores stored for each text sample.

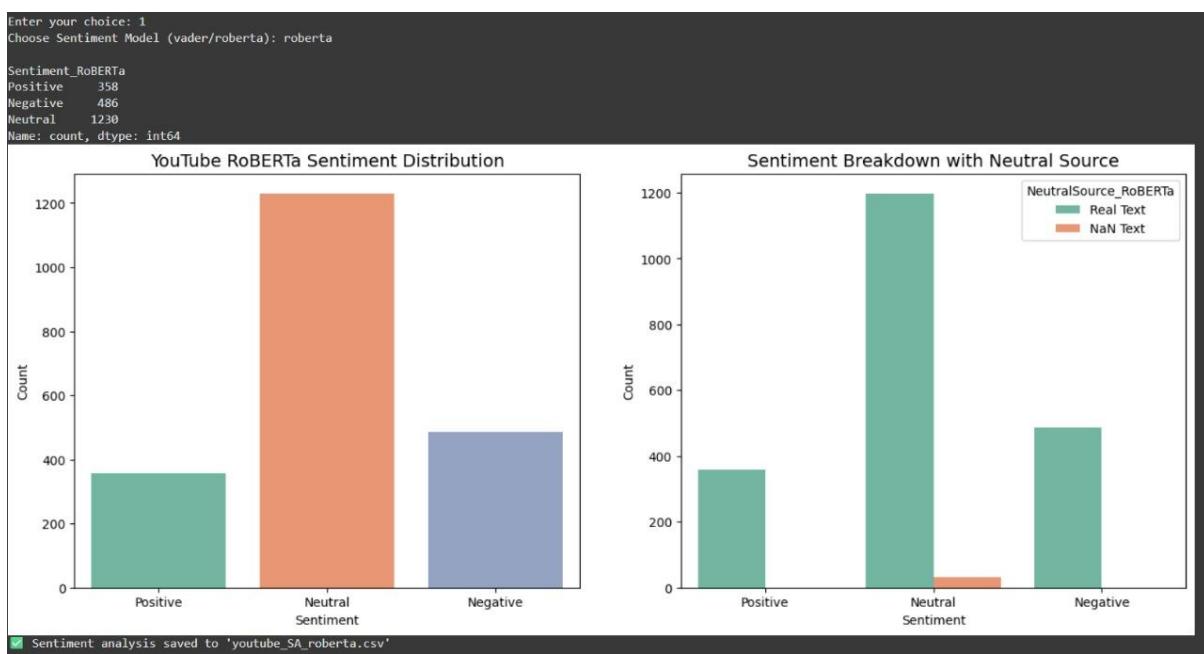
Aspect	VADER	RoBERTa
Type	Rule-Based (Lexicon)	Deep Learning (Transformer)
Ideal For	Short social media texts	Context-rich texts of varying lengths
Speed	Very Fast	Slower (requires GPU for best performance)
Accuracy	Good for simple language	Superior for complex & nuanced language
Output	Compound score & label	Class label & probability scores

Sentiment Analysis performed on the same dataset scrapped

VADER Sentiment Analysis



RoBERTa Sentiment Analysis



TOPIC MODELING

Overview:

Topic modeling was employed in this project to extract latent themes and topics from the large corpus of collected textual data. This unsupervised machine learning technique automatically identifies hidden semantic structures within the data and groups related words under coherent topics.

The goal was to understand **what people are commonly discussing** across platforms (YouTube, Reddit, Medium, Bing News) beyond just sentiment.

The project used **TF-IDF Vectorization** combined with **Non-Negative Matrix Factorization (NMF)** for topic modeling.

TF-IDF (Term Frequency-Inverse Document Frequency):

- Converts textual data into a numerical matrix, highlighting words that are **important in a specific document but less frequent globally**.
- Helps focus on meaningful and discriminative words for topic modeling.

NMF (Non-Negative Matrix Factorization):

- Decomposes the TF-IDF matrix into a predefined number of topics and associated keywords.
- Suitable for identifying **clear, non-overlapping topics** from sparse datasets like short comments and posts.

Implementation Steps:

1. Preprocessing:

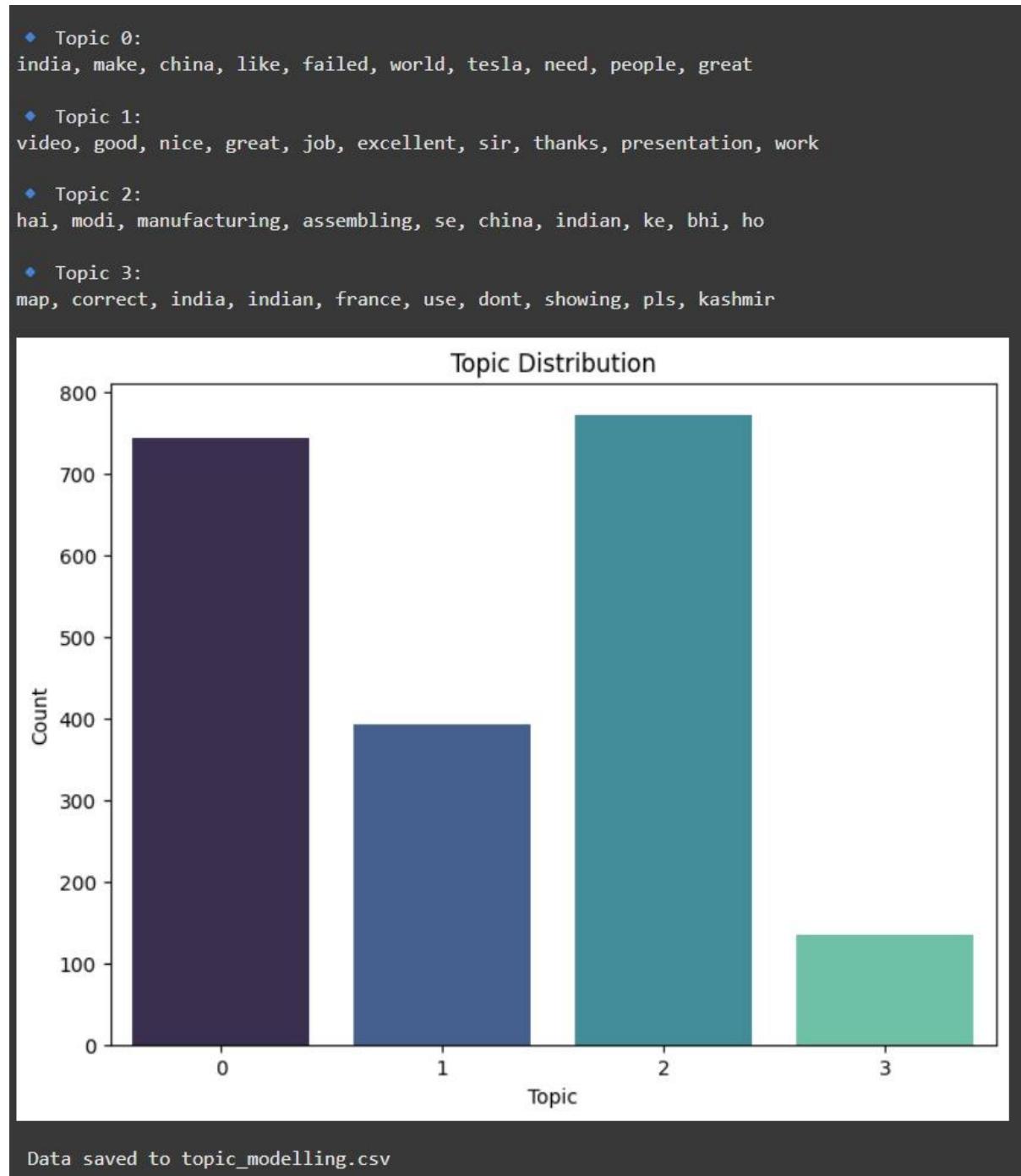
- Text cleaning: Removed punctuations, special characters, URLs, and stopwords.
- Tokenized words and performed lemmatization.

2. Vectorization:

- Applied TfIdfVectorizer from scikit-learn to convert preprocessed text into a weighted numerical representation.
- Set parameters such as max_features and ngram_range to optimize performance.

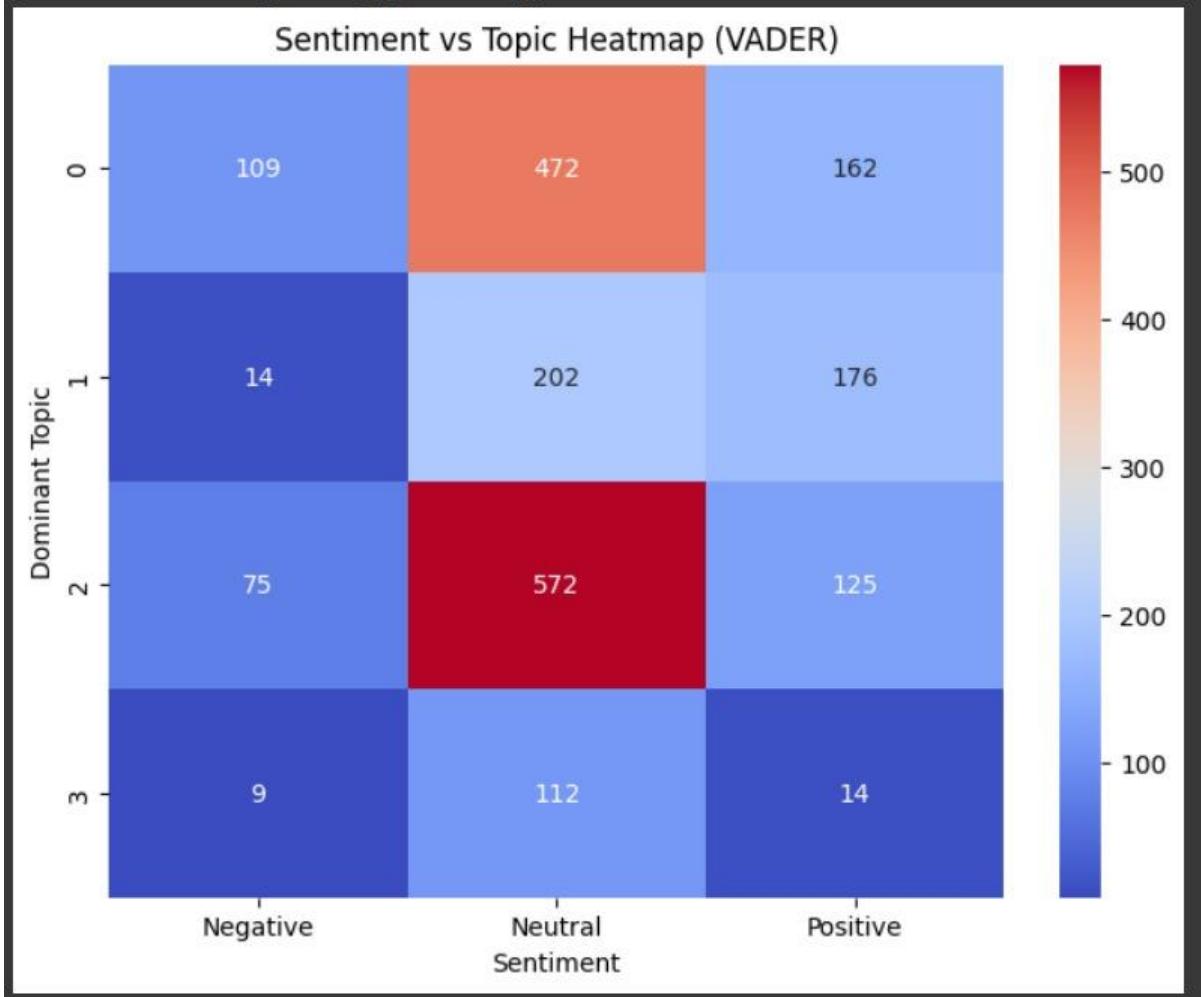
3. Modeling:

- Used NMF to extract topics and keywords.
- The number of topics was set dynamically based on user input (commonly 5-10 topics).
- Extracted **top 10 keywords per topic** to label and understand each theme.



```
Sentiment
Neutral    1358
Positive   477
Negative   207
Name: count, dtype: int64
```

```
Sentiment  Negative  Neutral  Positive
Topic
0          109      472      162
1          14       202      176
2          75       572      125
3          9        112      14
```



VISUALIZATION

❖ Sentiment Analysis:

- **Bar Graphs:** Displayed sentiment distribution (Positive, Neutral, Negative) across platforms.
- **Count Plots:** Compared sentiment categories by source (YouTube, Reddit, Medium, Bing News).

❖ Topic Modeling:

- **Bar Graphs:** Showed the frequency of each discovered topic.
- **Heatmaps:** Mapped the relationship between topics and sentiment categories, providing insights into how certain topics align with user sentiments.

❖ Tools Used:

- **Matplotlib**
- **Seaborn**
- **TQDM** (Progress Visualization)

❖ Purpose of Visualization:

- Simplified complex analytical outcomes.
- Provided actionable insights through clear, engaging visuals.
- Enhanced interpretability for end users and stakeholders.

Sample Output – Full Pipeline run (YouTube)

```
Welcome to Policy Insights Chatbot!

● Know About Sentiment Analysis Models!
🔗 Click here to learn about VADER and RoBERTa Sentiment Models

Available Platforms for Analysis :
1. Reddit
2. YouTube
3. Bing Articles
4. Medium
5. Sentiment Distribution across platforms
6. Exit

Enter your choice : 2

Enter keywords (comma separated): make in india

You have opted for YOUTUBE Data Analysis
Enter maximum count of videos : 10
Enter maximum comments per video : 1000

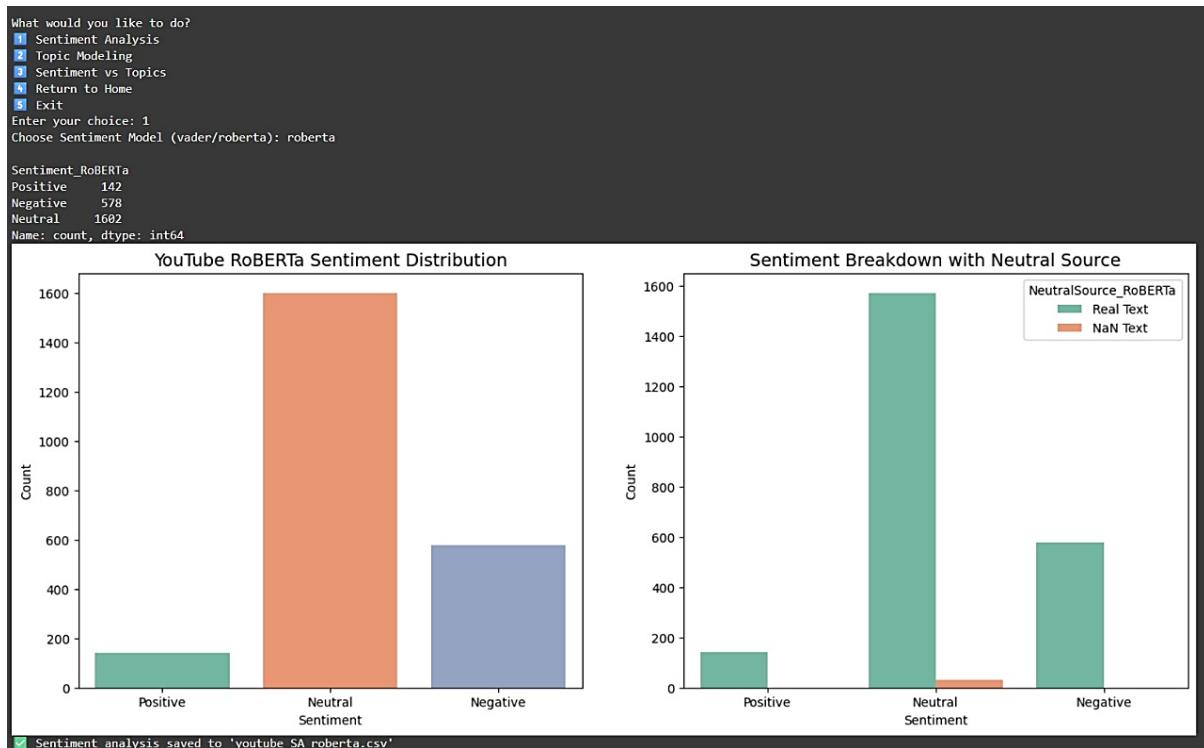
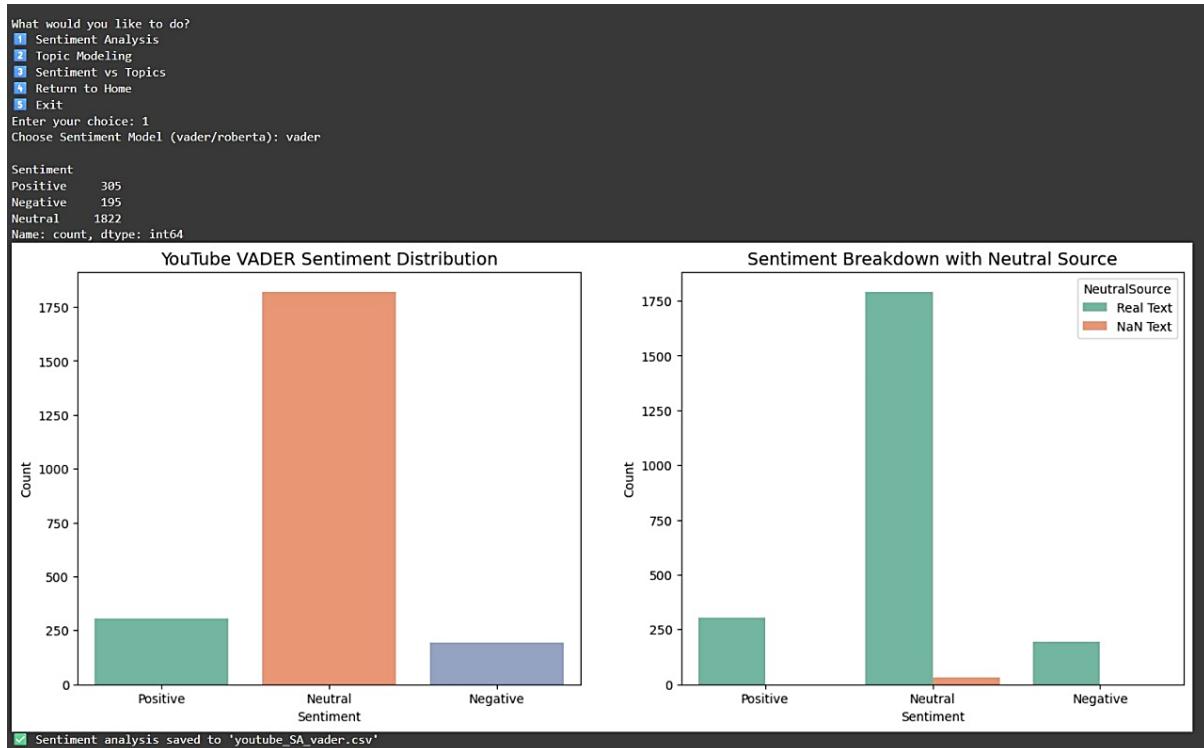
Hold on! We're collecting the content and tidying it up for you. This might take a little while--thank you for your patience!
```

```
✖ Scraping YouTube data...
✖ Starting search with keywords: ['make in india'] and max_results: 60
👉 Searching for keyword: make in india

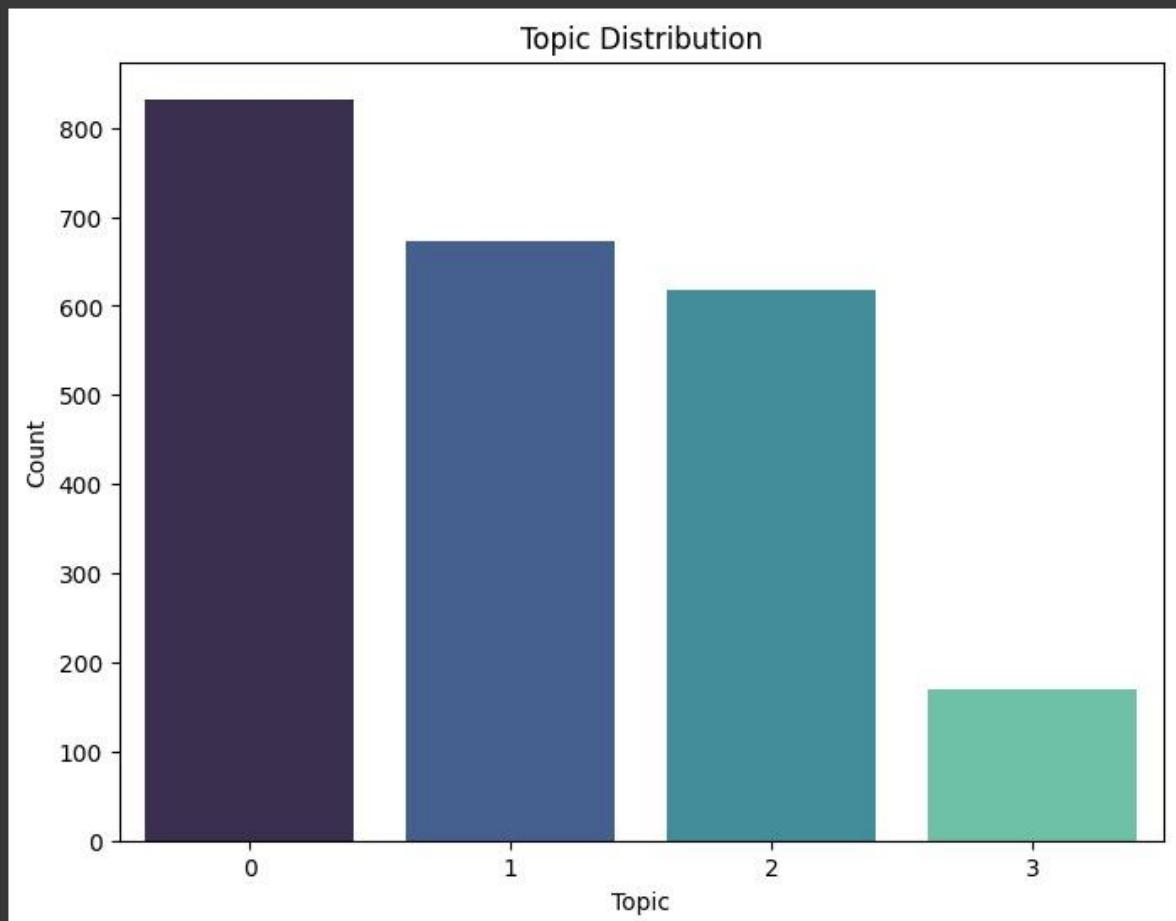
● Found 50 potential videos for keywords: ['make in india']
✖ Fetching comments for video ID: osy7Ds03J2Q
✓ Retrieved 59 comments
✓ Video ID osy7Ds03J2Q accepted (>= 10 comments)
✖ Fetching comments for video ID: TMW-iWeq08k
✓ Retrieved 0 comments
⚠ Video ID TMW-iWeq08k skipped (less than 10 comments)
✖ Fetching comments for video ID: lhjEsDRaz60
✓ Retrieved 1000 comments
✓ Video ID lhjEsDRaz60 accepted (>= 10 comments)
✖ Fetching comments for video ID: SaJ0SvoN961
✓ Retrieved 2 comments
⚠ Video ID SaJ0SvoN961 skipped (less than 10 comments)
✖ Fetching comments for video ID: TdJ4cbV6T3w
✓ Retrieved 149 comments
✓ Video ID TdJ4cbV6T3w accepted (>= 10 comments)
✖ Fetching comments for video ID: lXa8P0oggyc
✓ Retrieved 0 comments
⚠ Video ID lXa8P0oggyc skipped (less than 10 comments)
✖ Fetching comments for video ID: -aHePySasFY
✓ Retrieved 45 comments
✓ Video ID -aHePySasFY accepted (>= 10 comments)
✖ Fetching comments for video ID: q0ZAr_Bvul0
✓ Retrieved 500 comments
✓ Video ID q0ZAr_Bvul0 accepted (>= 10 comments)
✖ Fetching comments for video ID: lphdtOC0h3s
✓ Retrieved 73 comments
✓ Video ID lphdtOC0h3s accepted (>= 10 comments)
✖ Fetching comments for video ID: pdjWJMfy0QA
✓ Retrieved 156 comments
✓ Video ID pdjWJMfy0QA accepted (>= 10 comments)
✖ Fetching comments for video ID: pMjD5H8fI0
✓ Retrieved 182 comments
✓ Video ID pMjD5H8fI0 accepted (>= 10 comments)
✖ Fetching comments for video ID: ycGBihoELOo
✓ Retrieved 24 comments
✓ Video ID ycGBihoELOo accepted (>= 10 comments)
✖ Fetching comments for video ID: 5kqlu8IK7iY
✓ Retrieved 134 comments
✓ Video ID 5kqlu8IK7iY accepted (>= 10 comments)

✓ All valid comments saved to youtube_full.csv
✓ Total valid videos collected: 10 / 10
```

```
👉 Preprocessing YouTube data...
🕒 Loading data for preprocessing...
✓ Total Comments: 2322
🔴 Cleaning comments...
✓ Preprocessing completed. Cleaned data saved to youtube_full_cleaned.csv
```



```
What would you like to do?  
1 Sentiment Analysis  
2 Topic Modeling  
3 Sentiment vs Topics  
4 Return to Home  
5 Exit  
Enter your choice: 2  
Enter number of topics : 4  
  
• Topic 0:  
india, make, china, dont, hygiene, tesla, world, manufacturing, cheap, people  
  
• Topic 1:  
ice, dirty, cheap, water, floor, drink, dont, food, like, use  
  
• Topic 2:  
hai, china, ki, ko, se, aur, ka, modi, bhi, ke  
  
• Topic 3:  
rusty, steel, molds, copper, lol, thats, moulds, rust, arent, like
```



Data saved to topic_modelling.csv

```

What would you like to do?
1 Sentiment Analysis
2 Topic Modeling
3 Sentiment vs Topics
4 Return to Home
5 Exit
Enter your choice: 3
Enter model type : vader
Enter number of topics : 4

• Topic 0:
india, make, china, dont, hygiene, tesla, world, manufacturing, cheap, people

• Topic 1:
ice, dirty, cheap, water, floor, drink, dont, food, like, use

• Topic 2:
hai, china, ki, ko, se, aur, ka, modi, bhi, ke

• Topic 3:
rusty, steel, molds, copper, lol, thats, moulds, rust, arent, like

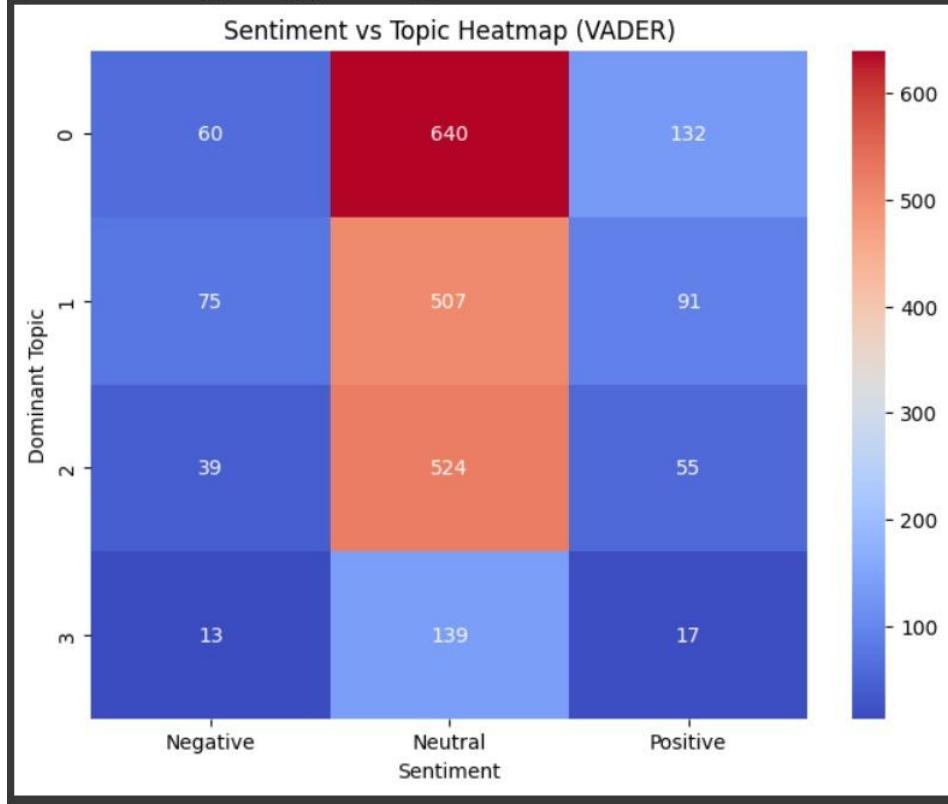
```

```

Sentiment
Neutral    1810
Positive     295
Negative     187
Name: count, dtype: int64

```

	Sentiment	Negative	Neutral	Positive
Topic				
0		60	640	132
1		75	507	91
2		39	524	55
3		13	139	17



```

What would you like to do?
1 Sentiment Analysis
2 Topic Modeling
3 Sentiment vs Topics
4 Return to Home
5 Exit
Enter your choice: 3
Enter model type : roberta
Enter number of topics : 4

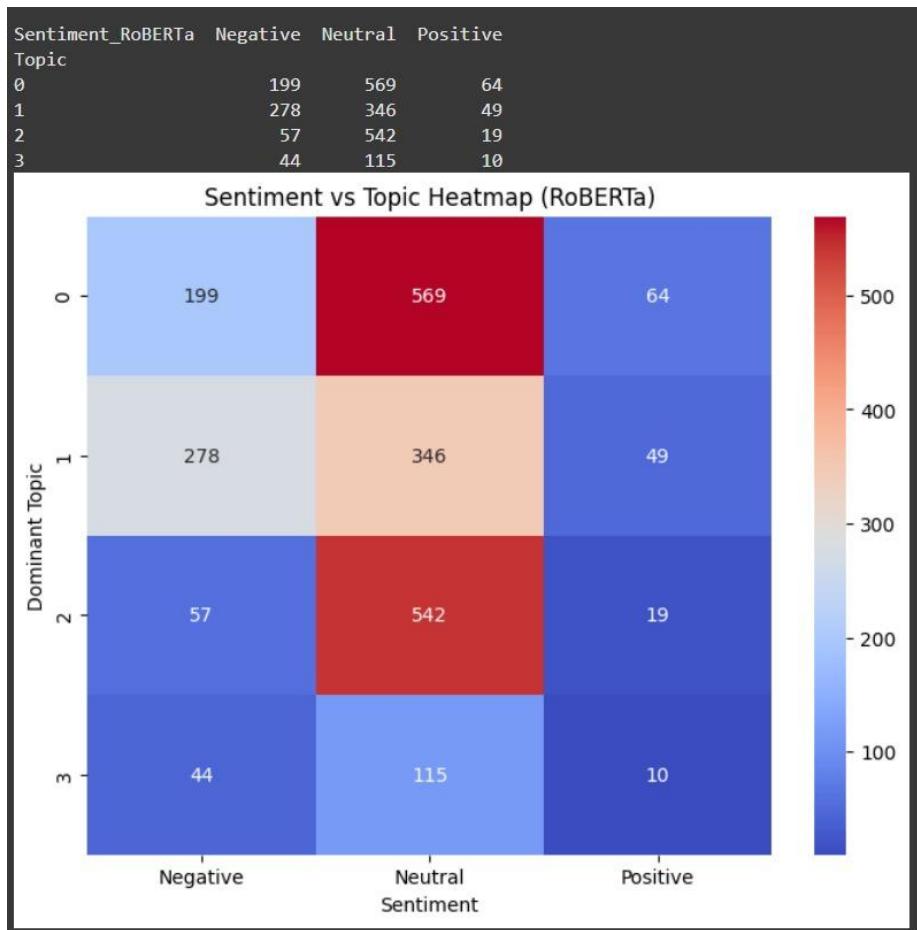
• Topic 0:
india, make, china, dont, hygiene, tesla, world, manufacturing, cheap, people

• Topic 1:
ice, dirty, cheap, water, floor, drink, dont, food, like, use

• Topic 2:
hai, china, ki, ko, se, aur, ka, modi, bhi, ke

• Topic 3:
rusty, steel, molds, copper, lol, thats, moulds, rust, arent, like

```



```

What would you like to do?
1 Sentiment Analysis
2 Topic Modeling
3 Sentiment vs Topics
4 Return to Home
5 Exit
Enter your choice: 5

👋 Exiting the chatbot. Thank you!

*****Session Terminated*****

```

RESULTS AND OBSERVATIONS

Summary of Sentiment Analysis Results:

Sentiment analysis was performed on data from **YouTube, Reddit, Medium, and Bing News** using **VADER** and **RoBERTa** models. Both provided unique insights, highlighting the difference between lexicon-based and deep learning-based approaches.

Insights from VADER Analysis:

- **VADER (Valence Aware Dictionary and sEntiment Reasoner)** is a **lexicon and rule-based model** specifically attuned to sentiments expressed in social media and informal texts.
- **Key Observations:**
 - Tends to label **neutral or mildly positive tones more favorably**.
 - Recognizes emojis, slang, and common intensifiers (e.g., "very good", "super bad").
 - Sometimes misclassifies **long, complex sentences or sarcasm** as neutral or even positive.

Insights from RoBERTa Analysis:

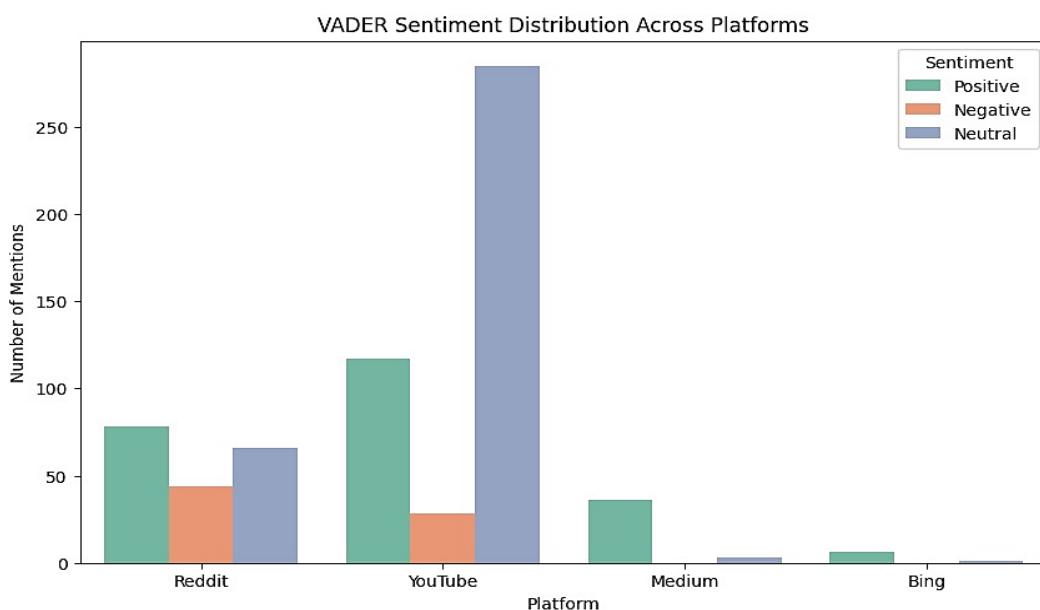
- **RoBERTa (Robustly Optimized BERT Pretraining Approach)** is a **deep learning transformer model** fine-tuned for sentiment tasks.
- **Key Observations:**
 - Evaluates sentiment based on **contextual depth**, detecting nuances and hidden tones.
 - More effective at picking up **subtle negativity or indirect sentiment expressions**.
 - Handles long texts, sarcasm, and varied sentence structures more robustly.

Aspect	VADER	RoBERTa
Strength	Fast, lightweight, social-media friendly	Deep understanding, context-aware
Weakness	Struggles with complex texts	Requires more computation, slower
Bias Tendency	Slightly optimistic	Balanced, realistic
Neutral Detection	High	Balanced

CONCLUSION

This project successfully developed a **robust and modular pipeline** for:

- ✓ Collecting data from multiple platforms (YouTube, Reddit, Medium, Bing News).
 - ✓ Efficient data cleaning and structuring for accurate sentiment analysis.
 - ✓ Performing **sentiment classification** using both **VADER** (lexicon-based) and **RoBERTa** (transformer-based) models for broader perspective.
 - ✓ Extracting meaningful **topics** through **NMF topic modeling**.
- Providing **clear, visual insights** through well-crafted graphs and breakdowns.



FUTURE SCOPE

- Integrate **real-time dashboards** for continuous monitoring of sentiments and topics.
- Deploy as a **user-friendly web application** for wider accessibility.
- Expand analysis by incorporating **emotion detection** (joy, anger, fear) for deeper understanding of public opinion.
- Incorporate trend analysis over time to observe how sentiments evolve.
- Expand data sources to include forums, blogs, and news APIs for richer insights.

REFERENCES

- **YouTube Data API v3 Documentation**
<https://developers.google.com/youtube/v3>
- **PRAW: The Python Reddit API Wrapper**
<https://praw.readthedocs.io/>
- **Google Custom Search JSON API Documentation**
<https://developers.google.com/custom-search/v1/overview>
- **Bing News Search**
<https://www.bing.com/news>
- **BeautifulSoup: Documentation**
<https://www.crummy.com/software/BeautifulSoup/bs4/doc/>
- **Transformers Library (HuggingFace)**
<https://huggingface.co/transformers/>
- **Cardiff NLP Twitter RoBERTa Base Model**
<https://huggingface.co/cardiffnlp/twitter-roberta-base-sentiment>
- **VADER Sentiment Analysis**
Hutto, C.J. & Gilbert, E.E. (2014). VADER: A Parsimonious Rule-based Model for Sentiment Analysis of Social Media Text. *Proceedings of the Eighth International Conference on Weblogs and Social Media (ICWSM-14)*.
- **Scikit-learn Documentation: NMF for Topic Modeling**
<https://scikit-learn.org/stable/modules/generated/sklearn.decomposition.NMF.html>