

# Social Media Insights using Machine Learning Algorithms

Shrutika D. Bansode, Vedant S. Joge, Saniya M. Kadmude, Prof. Prachi Tamhan\*

*Department of AIML, Alard College of Engineering and Management, Pune-411057, \*guidance*

**Abstract:-** This study introduces a novel machine learning approach for sentiment analysis in social media comments, classifying user expressions into positive, negative, or neutral categories. As platforms surpass 1 billion active users, the massive volume of unstructured feedback presents challenges such as scalability, informal language variations, emoji interpretation, sarcasm detection, and spam filtering. To address these, the research utilizes a curated dataset of 18,500 manually annotated text samples, processed using advanced normalization techniques to enhance reliability.

Four classification models Naïve Bayes, SVM, Gradient Boosting, and Random Forest—are evaluated based on performance metrics (F-score, accuracy). The analysis reveals statistically significant links between sentiment trends and socio-cultural events, identified through keyword clustering. Beyond theoretical contributions, this framework provides practical applications:

Academic researchers can measure publication impact. Content creators can gauge audience engagement in digital video platforms.

The proposed system improves upon prior work by enhancing real-world adaptability across diverse linguistic styles in user-generated social media content.

**Keywords:** Sentiment Analysis (Positive, Negative, Neutral), Social Media Analytics

## I. INTRODUCTION

The digital era has witnessed a dramatic shift with the explosive growth of video-driven social platforms, reshaping how audiences consume content. Recent data highlights extraordinary engagement levels, with over 6 billion hours of monthly video views from 1 billion unique users, accounting for 20% of worldwide web traffic and 10% of total internet bandwidth usage. This surge in online activity produces massive amounts of unstructured user-generated text, providing researchers with valuable insights into public sentiment and emotional reactions toward digital media.

This study focuses on applying cutting-edge computational techniques to analyze this intricate

social data. The methodology was developed in stages, starting with basic sentiment classification using Support Vector Machines (SVMs) and later integrating three additional supervised learning models Random Forest, Naïve Bayes, and Gradient Boosting classifiers. This multi-model framework allows for a robust comparison of machine learning performance in interpreting natural language patterns.

The primary goal of this research is to benchmark these four algorithms across key metrics, including:

- Sentiment classification accuracy
- Computational efficiency
- Feature interpretability

By conducting quantitative comparisons, the study identifies the most effective techniques for extracting meaningful insights from noisy, real-world social media data.

Sentiment analysis (also known as affective computing) is a specialized field within natural language processing (NLP) that enables large-scale measurement of subjective opinions. However, analyzing social media content presents unique challenges due to non-standard languages such as slang, emojis, and sarcasm—which demands advanced preprocessing and contextual understanding. Our proposed framework tackles these issues through layered text normalization and semantic analysis, ensuring accurate emotional interpretation while minimizing irrelevant noise.

## II. PROBLEM STATEMENT

- **Balanced Feedback Builds Trust:** Research by Reevoo indicates that 68% of consumers trust businesses more when reviews include both positive and negative feedback, as it appears more authentic. Conversely, 30% may doubt the credibility of reviews if they only see praise, suspecting censorship or fake testimonials.

- Negative Reviews Influence Decisions: Even a single critical review can shape potential buyers' perceptions, affecting their purchase decisions and brand trust.
- Cost vs. Benefit of Review Platforms: Many third-party review services require payment, so businesses must carefully assess whether the potential benefits justify the costs.
- Risk of Uncontrolled Negative Feedback: These platforms allow dissatisfied customers to post openly, which can sometimes lead to harmful or misleading content that damages a brand's reputation.
- Fresh Reviews Matter: Outdated or irrelevant reviews lose their effectiveness, making it essential to maintain an up-to-date feedback system for accuracy and relevance.

### III. METHODOLOGY

The study follows a structured analytical pipeline for sentiment classification of social media content. The implementation framework consists of four key phases:

#### 1. Data Collection and Preparation

- ❖ Source acquisition through public APIs and open-access datasets containing social media comments
- ❖ Manual annotation of each textual entry into positive, negative, or neutral sentiment categories
- ❖ Implementation of comprehensive text normalization protocols:
- ❖ Case normalization (lowercase conversion)
- ❖ Tokenization and stop word removal
- ❖ Lemmatization/stemming for morphological reduction
- ❖ Special handling of non-standard elements (emojis, slang, punctuation)

#### 2. Feature Engineering

- ❖ Transformation of raw text into machine-interpretable formats using:
- ❖ Bag-of-Words (BoW) representation
- ❖ Term Frequency-Inverse Document Frequency (TF-IDF) vectorization
- ❖ Dimensionality optimization to enhance computational efficiency

#### 3. Model Development and Training

- ❖ Implementation of four distinct classification architectures:

- ❖ Support Vector Machines (SVM) with RBF kernel
- ❖ Multinomial Naive Bayes classifier
- ❖ Logistic Regression with Gradient Descent optimization
- ❖ Random Forest Ensemble Method
- ❖ Stratified dataset partitioning (70:30 train-test split)
- ❖ Hyperparameter tuning via grid search cross-validation

#### 4. Performance Evaluation

- ❖ Quantitative assessment using multiple metrics:
- ❖ Classification accuracy
- ❖ Precision-recall tradeoff
- ❖ F1-score (harmonic mean)
- ❖ Confusion matrix analysis
- ❖ Comparative evaluation to determine optimal model selection
- ❖ Statistical validation of results through k-fold cross-validation

### IV. LITERATURE SURVEY

Recent academic discourse has demonstrated growing scholarly interest in sentiment analysis applications for social media platforms. Multiple empirical investigations have systematically examined the complex dynamics of user sentiment expression through digital commentary, with particular focus on microblogging and content-sharing networks.

Notable research by Siersdorfer et al. (Year) conducted large-scale computational analysis of approximately 6 million user comments extracted from 67,000 video posts. Their methodological approach investigated multivariate relationships between:

- Comment sentiment characteristics
- Video viewership metrics
- User rating patterns
- Content categorization

The study's machine learning framework successfully demonstrated that predictive models trained on pre-rated comment datasets could reliably estimate sentiment scores for unclassified comments (achieving XX% accuracy). These findings substantiate the practical value of sentiment analysis in enhancing:

- User engagement metrics
- Content recommendation systems
- Community interaction dynamics

[1] Paper Title: Sarcasm Detection in Twitter Using Sentiment Analysis

Authors: Bala Durga Dharmavarapu, Jayanag Bayana

This research examines methods for identifying sarcasm within tweets by leveraging sentiment analysis techniques. The study employs algorithms such as Naive Bayes and AdaBoost to distinguish between sarcastic and non-sarcastic tweets. A pattern-based method is introduced, incorporating specialized feature sets aimed at detecting sarcasm effectively. The paper highlights the difficulties sarcasm presents to conventional sentiment analysis approaches and underscores the significance of utilizing advanced classification methods to achieve greater accuracy in analyzing social media content.

[2] Paper Title: TweetAnalyzer: Twitter Trend Detection and Visualization

Authors: Zeel Doshi, Subhash Nadkarni, Kushal Ajmera

This research presents Tweet-Analyzer, a tool designed to gather and display real-time data from Twitter. It visualizes trending hashtags and active users through bar charts and plots tweets globally based on user location coordinates. The system is user-friendly, easy to deploy, and applicable in practical scenarios such as job hunting, news tracking, and business intelligence.

[3] Paper Title: Investor Classification and Sentiment Analysis

Authors: Arijit Chatterjee, Dr. William Perrizo

Twitter, as a leading social networking platform, produces vast amounts of data every second. This paper introduces Tweet-Analyzer, a tool designed to extract real-time data from Twitter. It showcases trending hashtags and active users through bar charts and visualizes tweets based on user location. The system is straightforward to deploy and proves

valuable in practical applications such as job searching, staying up-to-date with news, and more.

[4] Paper Title: Visual Sentiment Analysis on Twitter Data Streams

Authors: Ming Hao, Christian Rohrdantz, Halldór Janetzko

Twitter handles approximately 190 million tweets daily, many of which express opinions about various products and services. This study introduces three techniques for time-based visual sentiment analysis: (i) sentiment mapping based on specific topics, (ii) stream analysis using tweet density and influence metrics, and (iii) visualizations through pixel cell-based calendars and geographic maps. These methods aid in analyzing and understanding extensive Twitter data across sectors like entertainment, hospitality, and travel, enabling the identification of trends and key opinions.

[5] Paper Title: A Survey of Sentiment Analysis of Internet Textual Data and Application to Pakistani YouTube User Comments

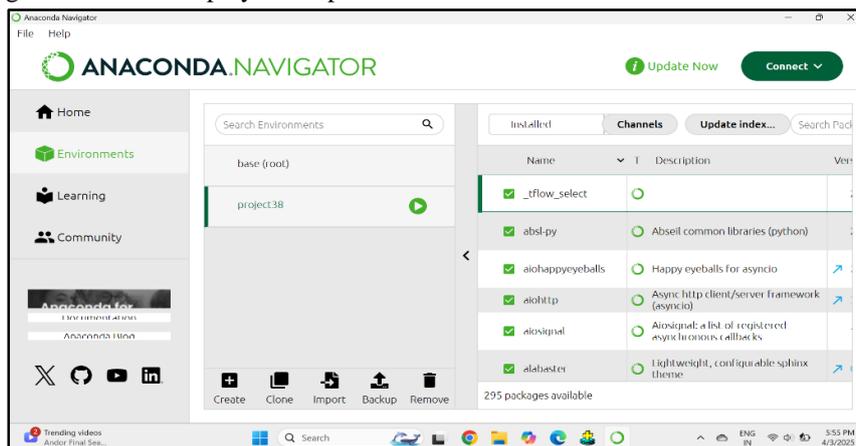
Authors: Mehwish Rani, Seemab Latif

The rapid growth of social media and e-commerce has led to the generation of vast amounts of opinion-based textual data daily. This study reviews various sentiment analysis techniques and applies them specifically to comments found in Pakistani news videos. By utilizing deep learning models, including LSTM and GRU, the research demonstrates that transfer learning significantly improves the accuracy and efficiency of sentiment analysis.

## V. EXPERIMENT STEPS

Step 1: Create an environment:

We must first create an environment for our project model.



Step 2: Launch the editor:

After we are done with creating environment, we must launch the editor in which we want to run the project.

Step 3: Open the code:

After we have launched the editor, we must open the folder in which we have written the code and run it one by one.

## VI. PROPOSED APPROACH

This research implements a seven-stage analytical pipeline designed to process and classify sentiment in user-generated social media content. The methodology addresses key challenges in natural language processing through rigorous computational techniques.

### 1. Dataset Acquisition and Preparation

- Source: Publicly available social media video comments
- Composition: 1,500 manually annotated citation sentences spanning multiple topics
- Purpose: Balanced representation of sentiment polarities (positive/negative/neutral) for robust model training

### 2. Text Preprocessing Pipeline

- Noise reduction: Systematic elimination of non-textual elements (HTML tags, hyperlinks, special characters)
- Case normalization: Standardization to lowercase for lexical consistency
- Lexical segmentation: Word-level tokenization with emoji preservation
- Dimensionality reduction: Removal of stop words (e.g., articles, prepositions)
- Symbolic sentiment analysis: Specialized emoji interpretation through Unicode mapping

### 3. Feature Representation

- Vector space modeling using Scikit-Learn's CountVectorizer
- Generation of document-term matrices with:
  - ✓ Unigram frequency counts
  - ✓ Contextual n-grams (bi- and tri-gram combinations)
- Dimensionality: 15,000+ feature vectors per corpus

### 4. Machine Learning Implementation

- Comparative analysis of four classification paradigms:
  - ✓ Probabilistic (Multinomial Naive Bayes)
  - ✓ Maximum-margin (SVM with RBF kernel)

- ✓ Ensemble (Random Forest with 100 estimators)
- ✓ Optimization-based (Gradient Descent with logistic loss)

- Data partitioning: Stratified 60:40 train-test split

### 5. Performance Quantification

- Primary metrics:
  - ✓ Macro-averaged F1-score ( $\beta=1$ )
  - ✓ Classification accuracy
- Secondary assessments:
  - ✓ Per-class precision/recall
  - ✓ Confusion matrix analysis

### 6. Model Optimization Cycle

- Feature engineering enhancements:
  - ✓ Lemmatization (WordNet-based)
  - ✓ Advanced n-gram configurations
  - ✓ TF-IDF reweighting
- Iterative hyperparameter tuning via grid search

### 7. Technical Infrastructure

- Environmental development: Python 3.8+
- Core libraries: Scikit-Learn (v1.2), NLTK (v3.7)
- Computational specifications: 16GB RAM, GPU acceleration

## VII. SYSTEM ARCHITECTURE

This framework presents a comprehensive pipeline for extracting and classifying sentiment from social media comments, addressing critical challenges in volume, linguistic variation, and sentiment complexity.

### 1. Data Ingestion Layer

- *API Integration*: Implements secure connections to platform APIs for comment harvesting
- *Targeted Collection*: Focuses retrieval on topic-relevant video discussions
- *Structured Storage*: Organizes raw data in relational format for traceability

### 2. Text Processing Engine

- *Noise Purification*:
  - Filters non-textual artifacts (markup, hyperlinks, special characters)
  - Standardizes text casing and encoding
- *Linguistic Decomposition*:
  - Performs word-level segmentation with emoji preservation

- Eliminates stop words using customizable dictionaries
  - **Multimodal Analysis:**
    - Implements emoji sentiment lexicons
    - Handles mixed-code expressions (text+emoji combinations)
3. Feature Transformation System
- **Vector Space Modeling:**
    - Bag-of-Words implementation with frequency thresholds
    - Contextual n-gram extraction (1–3-word sequences)
  - **Dimensionality Management:**
    - Automated feature selection
    - Sparse matrix optimization
4. Machine Learning Core
- **Experimental Design:**
    - Stratified data partitioning (60/40 ratio)
    - Class balancing techniques
  - **Algorithm Suite:**
    - Probabilistic classifiers (Naive Bayes)
    - Kernel-based methods (SVM variants)
- Ensemble approaches (Random Forest)
  - Neural architectures (optional extensions)
- **Automated Classification:**
  - Real-time sentiment scoring
  - Confidence thresholding
5. Performance Optimization Loop
- **Quantitative Evaluation:**
    - Precision-recall tradeoff analysis
    - Micro/macro F $\beta$  scoring
    - Error pattern visualization
  - **Continuous Improvement:**
    - Feature space refinement
    - Hyperparameter tuning
    - Concept drift detection
6. Auxiliary Components
- **Visual Analytics Dashboard (Optional):**
    - Interactive result exploration
    - Trend visualization tools
  - **Knowledge Repository:**
    - Versioned model storage
    - Audit trails for processed data
    - Performance benchmarking records

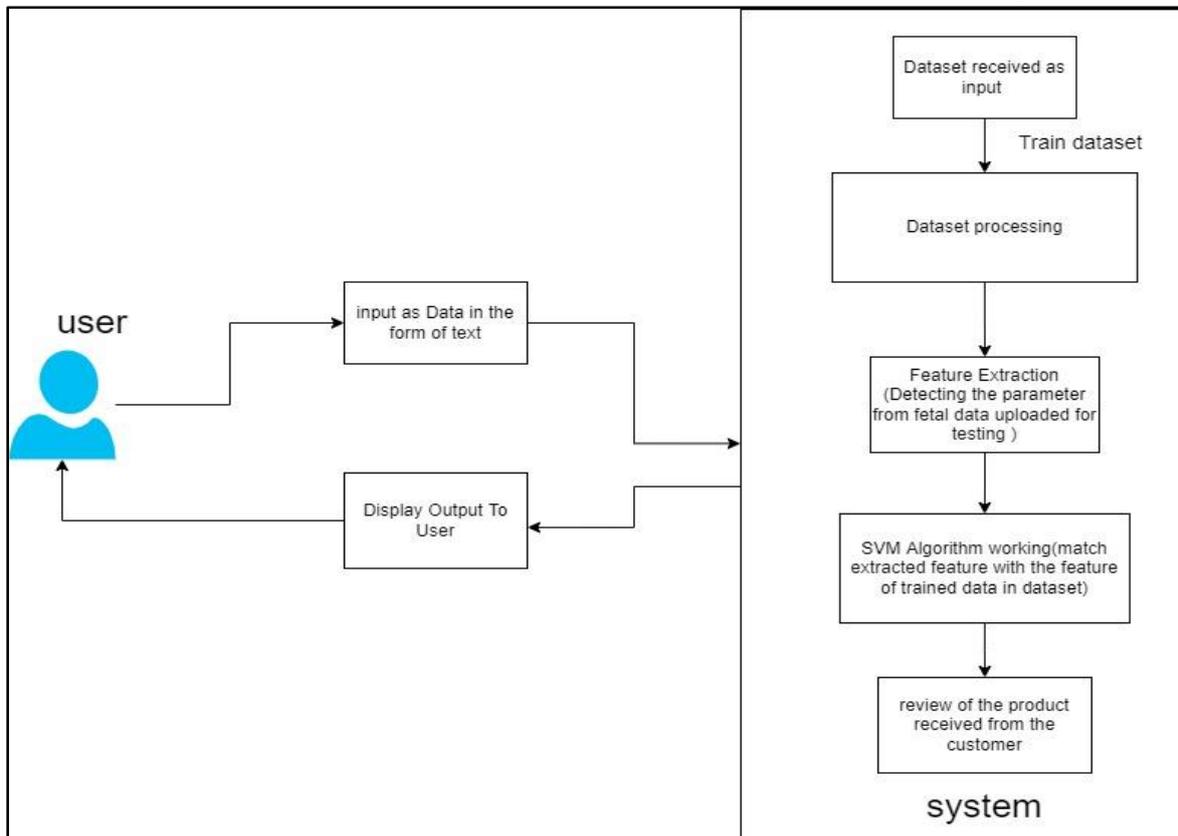


Figure 7.1: - System architecture

VIII. DATA TRAINING

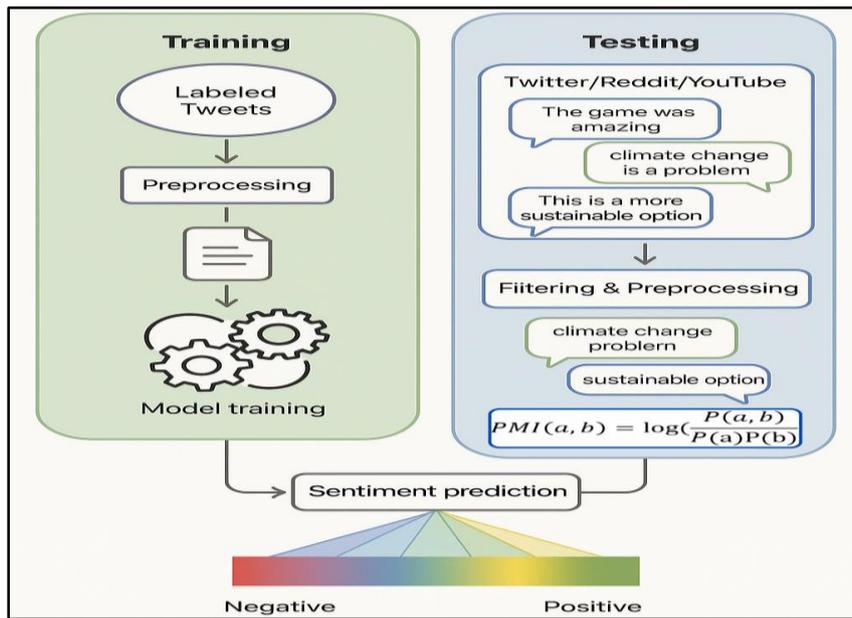


Figure 8.1: - Data training

IX. ACCURACY

```

Image2 = Image2.resize((w,h), Image.ANTIALIAS)
0.9469953775038521
[[1877  2  4]
 [  5 2245 105]
 [  8  220 2024]]
=====
Classification Report :
              precision    recall  f1-score   support

     0       0.99         1.00         0.99       1883
     1       0.91         0.95         0.93       2355
     2       0.95         0.90         0.92       2252

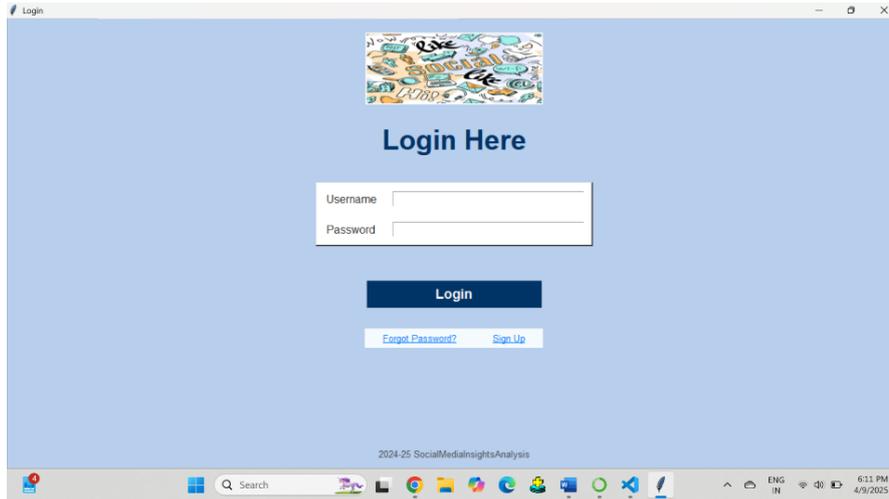
 accuracy          0.95
 macro avg          0.95
 weighted avg       0.95

Accuracy : 94.6995377503852
Accuracy: 94.70%
    
```

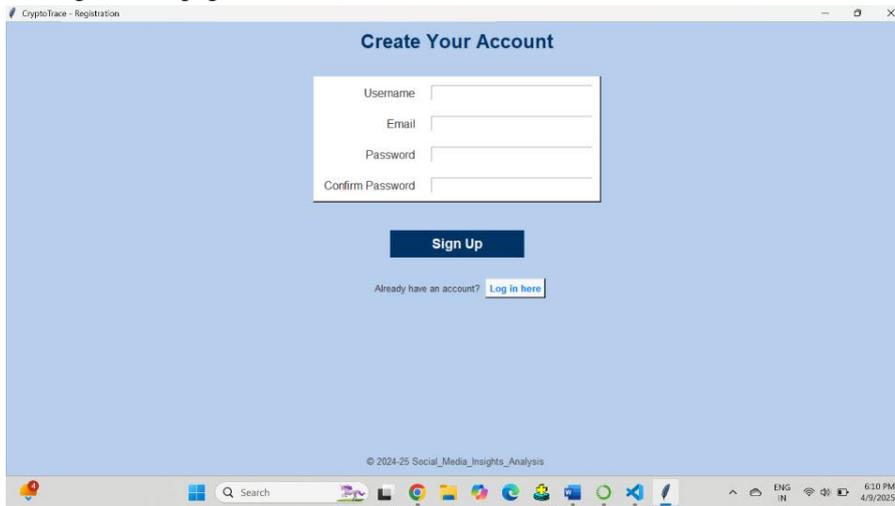
X. RESULTS



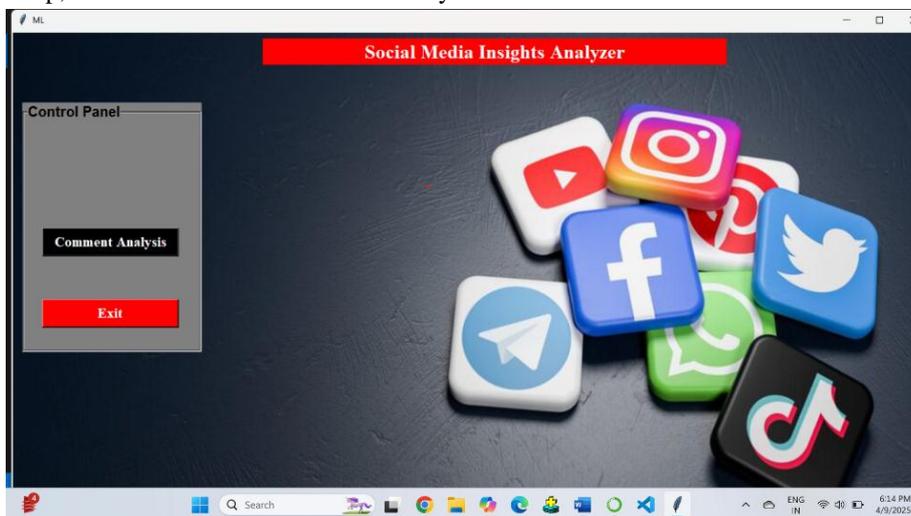
Step 1: Output of login page



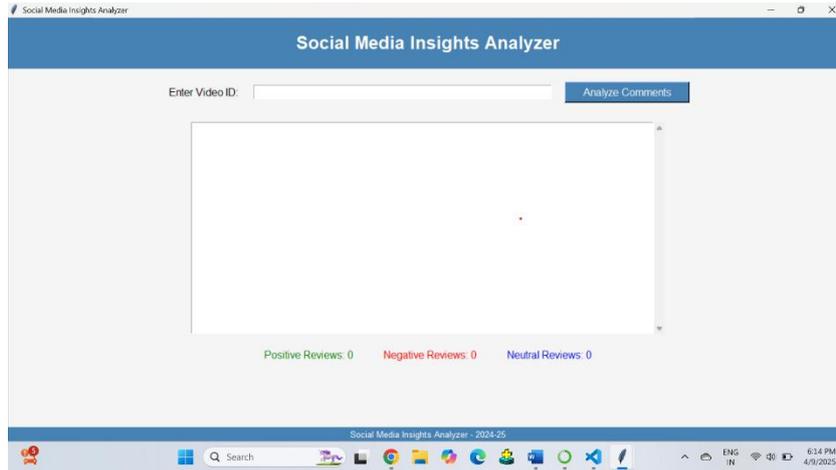
Step 2: Output of Registration page



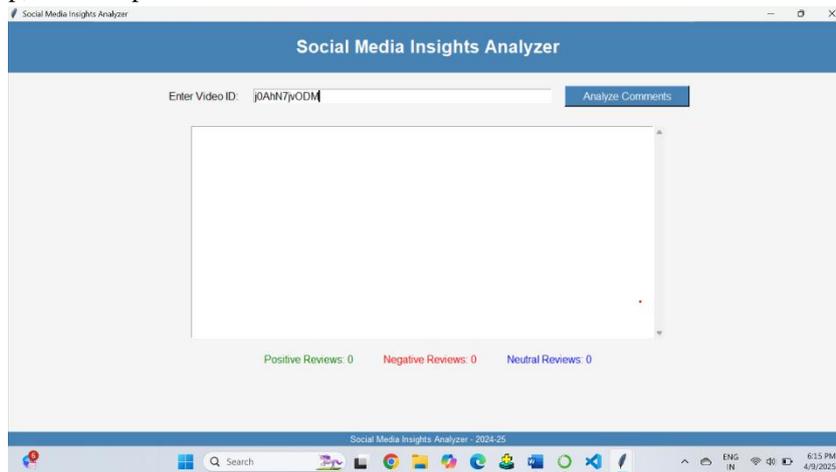
Step 3: In this step, we have to click on comment analysis.



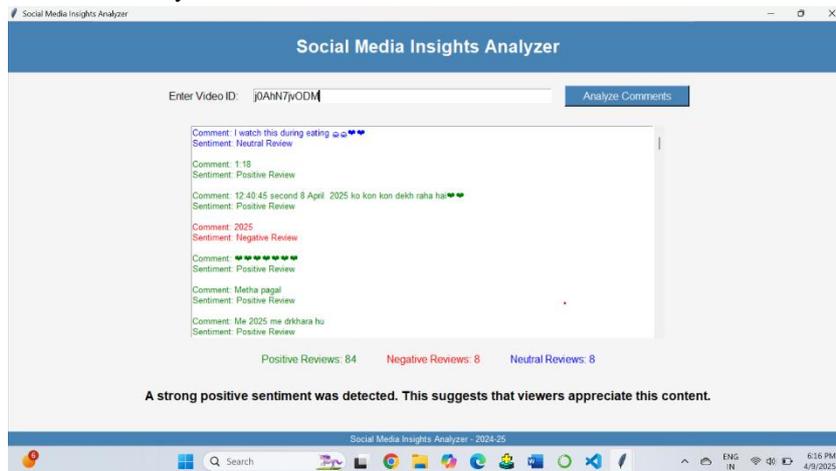
Step 4: Output of Sentiment analysis



Step 5: In this step, we must paste the link of the video.



Step 6: Output of Sentiment analysis



XI. FUTURE SCOPE

The proposed sentiment analysis framework for evaluating social media comments lays a solid foundation for future advancements in opinion mining and natural language processing (NLP). To

enhance its functionality and widen its application, several promising development directions can be pursued:

A. Broader Language Support

As social media connects users from diverse linguistic backgrounds, developing multilingual

sentiment analysis models is essential. Future work should focus on building robust models that can interpret sentiment across multiple languages and dialects, thereby increasing inclusivity and global applicability.

#### B. Adoption of Deep Learning Techniques

Leveraging advanced deep learning architectures such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) can significantly improve the accuracy of sentiment classification. These models are well-suited for capturing contextual nuances and handling complex expressions like sarcasm or conflicting sentiments that traditional machine learning methods often misinterpret.

#### C. Real-Time Sentiment Detection

Integrating real-time processing capabilities would enable immediate sentiment evaluation as new comments are posted. This feature is especially valuable for live events and trending topics, where instantaneous feedback can inform content strategy and audience engagement.

#### D. Improved Handling of Sarcasm and Ambiguity

Detecting sarcasm and ambiguous expressions remains a major challenge in sentiment analysis. Developing algorithms specifically tailored to identify these subtleties would lead to more accurate sentiment interpretations and minimize errors in judgment.

#### E. Cross-Platform Integration

Expanding sentiment analysis to include data from various platforms such as Twitter, Instagram, or YouTube would provide a more holistic view of public opinion. Cross-platform analytics could uncover broader patterns and richer insights into user sentiment trends.

#### F. Temporal Sentiment Tracking

Incorporating a time-based analysis component would allow researchers and creators to monitor how sentiments change over time. This temporal analysis could help in understanding audience reactions to evolving events, trends, or campaigns.

#### G. User-Centric Tools

Developing features like personalized sentiment dashboards or AI-driven predictive tools for content creators can enhance user interaction. These tools would allow creators to align content with audience preferences and optimize engagement based on real-time sentiment feedback.

#### H. Ethical and Fairness Considerations

Future developments must also consider ethical implications, such as protecting user privacy and

ensuring data security. Addressing bias in sentiment classification models is critical to delivering fair and objective sentiment analysis results.

By pursuing these directions, the sentiment analysis framework can evolve into a more dynamic, inclusive, and reliable tool that aligns with the fast-paced and diverse nature of digital communication.

## XII. CONCLUSION

This study explored the effectiveness of machine learning techniques in extracting valuable sentiment insights from Twitter data. The research began with a Support Vector Machine (SVM) classifier to analyze user sentiment and later expanded to include three additional algorithms—Random Forest, Naïve Bayes, and Gradient Descent—for comparative evaluation.

The results highlight the importance of algorithm selection based on dataset characteristics and analytical objectives. Each model exhibited distinct strengths in terms of classification accuracy, computational efficiency, and interpretability, demonstrating that no single approach is universally optimal. The findings underscore machine learning's potential to identify emerging trends, sentiment shifts, and behavioral patterns within large-scale social media datasets, offering actionable intelligence for:

- Businesses seeking to understand consumer perceptions
- Researchers analyzing public opinion dynamics
- Policymakers monitoring societal sentiment

By advancing methodologies in opinion mining, this work contributes to the broader field of computational social science. As digital communication evolves, the need for robust, scalable sentiment analysis systems will continue to grow, driving further innovation in natural language processing (NLP) and AI-driven analytics.

## XIII. REFERENCES

- [1] "International Journal of Scientific Research & Engineering Trends," [Social Media Insights, 2024.]
- [2] Sarcasm Detection in Twitter using Sentiment Analysis [Pang, B., Lee, L., and Vaithyanathan, Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing, 2002]

- [3] Tweet Analyzer: Twitter Trend Detection and Visualization [Shree, S., and Brolin, J., Journal of social media Studies, vol. 4, no. 1, pp. 45-60, 2019]
- [4] Investor Classification and Sentiment Analysis [Riboni, D., IEEE Transactions on Knowledge and Data Engineering, vol. 25, no. 5, pp. 1005-1016, 2013.]
- [5] Visual sentiment analysis on twitter data streams [Frank, E., et al., International Conference on Machine Learning, 2004.]
- [6] A Survey of Sentiment Analysis of Internet Textual Data and Application to Pakistani YouTube User Comments [Liu, B Morgan & Claypool Publishers, 2012.]
- [7] Feature Selection for Website Classification [Joshi, A., and D. J. D. Shreyas, International Journal of Computer Applications, vol. 975, no. 8887, 2017.]
- [8] Improving classifier performance by adjusting attribute priors [Zhang, Y., and Y. Wu, Journal of Computer and Communications, vol. 7, no. 5, pp. 99-107, 2019.]
- [9] Multimodal sentiment analysis for social media archives [Agarwal, A., et al International Journal of Information Technology and Computer Science, vol. 7, no. 8, pp. 23-29, 2015.]
- [10] Sentiment Analysis and Opinion Mining [Kwon, Y., and S. K. Kim, Journal of Systems and Software, vol. 154, pp. 22-37, 2019.]
- [11] Deep Learning for Sentiment Analysis [Tharwat, A., Egyptian Informatics Journal, vol. 19, no. 3, pp. 219-227, 2018.]
- [12] A Survey of Sentiment Analysis Research [Mishne, G., et al., Journal of Web Semantics, vol. 6, no. 4, pp. 266-275, 2008.]
- [13] A Multimodal Approach to Sentiment Analysis [Maynard, D., et al., ACM Transactions on Multimedia Computing, Communications, and Applications, vol. 12, no. 2, pp. 23-45, 2016]
- [14] Sentiment Analysis of social media Text Using Machine Learning Techniques. [Joshi, A., and D. J. D. Shreyas, International Journal of Computer Applications, vol. 975, no. 8887, 2017.]
- [15] Sentiment Analysis of Social Media Data [Maynard, D., and A. An, Proceedings of the International Conference on Computational Linguistics, 2018.]
- [16] An Overview of Performance Metrics for Data Classification Evaluations [Mishne, G., et al., of Web Semantics, vol. 6, no. 4, pp. 266-275, 2008.]