# Whatsapp Chat Analyzer

## Abstract

The WhatsApp Chat Analyzer is a powerful tool designed to offer comprehensive insights into both individual and group chats on the WhatsApp platform. This innovative application enables users to upload their chat data, allowing for the extraction of valuable statistical information. The key features encompass various analytical capabilities including Media Statistics, Message Count, Word Analysis, Word Cloud generation, emoji usage breakdown, and Sentiment Analysis. These functionalities collectively empower users to better understand their chat interactions, revealing trends and patterns in media sharing, message frequency, word usage, and emotional tone. With the integration of the VADER sentiment analyzer, the tool goes a step further in providing sentiment analysis, offering valuable insights into the emotional context of the conversations. The WhatsApp Chat Analyzer is poised to be an indispensable asset for users seeking meaningful insights from their WhatsApp conversations.

## Introduction

In an age dominated by digital communication, platforms like WhatsApp have become integral to our daily interactions. The WhatsApp Chat Analyzer project addresses the need for a comprehensive tool to dissect and glean insights from these conversations. By harnessing the power of machine learning, this project aims to empower users with a multifaceted analysis of their WhatsApp chats. From media sharing trends to sentiment analysis, the tool offers a suite of features designed to shed light on the dynamics of group or individual conversations. This project not only showcases the potential of machine learning in understanding human interaction but also serves as a practical resource for users seeking to unlock the hidden patterns within their messaging history.

## Background

WhatsApp, a ubiquitous messaging platform, has evolved into a cornerstone of modern communication. With over 2 billion active users worldwide, it serves as a primary medium for personal, professional, and social interactions. However, the sheer volume of data generated through WhatsApp conversations presents an untapped wellspring of insights.

In various domains, WhatsApp chat analysis has emerged as a critical tool for understanding human behavior, sentiment, and trends. In business and marketing, it offers invaluable data for customer engagement, product feedback, and market research. For researchers and social scientists, it provides a unique window into societal dynamics, language usage, and cultural

trends. Additionally, in personal contexts, it enables users to reflect on their communication patterns, identify key themes, and track their emotional expressions over time.

By delving into WhatsApp chats, this project addresses the growing demand for tools that can distill meaningful information from this rich source of communication. Through the application of machine learning techniques, it seeks to unlock patterns, sentiments, and insights that may otherwise remain hidden in the deluge of messages. As such, WhatsApp chat analysis stands at the intersection of modern communication and data-driven insights, offering a powerful lens through which to understand the nuances of human interaction in an increasingly digital world.

# Objectives

**Comprehensive Analysis:** The primary goal of this project is to develop a tool capable of providing a comprehensive analysis of WhatsApp chats, encompassing media sharing trends, message frequency, word usage patterns, and emotional tone.

**User-Friendly Interface:** Create an intuitive and user-friendly interface that allows individuals and groups to easily upload their WhatsApp chat data for analysis without the need for technical expertise.

**Media Insights:** Enable users to gain insights into their media sharing habits, including the total number and types of media files exchanged within a given chat.

**Message Counting:** Implement a feature that calculates the overall count of messages exchanged, providing users with a quantifiable measure of their communication volume.

**Word Frequency Analysis:** Develop a function that identifies and presents the most frequently used words in the chat, allowing users to discern prominent themes and topics of discussion.

**Word Cloud Generation:** Generate visually appealing word clouds based on the content of the chat, offering a quick and engaging overview of the most prevalent words.

**Emojis Usage Breakdown:** Provide users with a breakdown of the emojis used and their frequency, offering insights into the emotional tone and expression within the conversations.

**Sentiment Analysis Integration:** Integrate the VADER sentiment analyzer to conduct sentiment analysis on messages, enabling users to gauge the emotional context of their conversations.

**Insightful Reporting:** Ensure that the tool delivers clear, concise, and insightful reports that distill the analyzed data into easily digestible information.

**Scalability and Efficiency:** Design the tool to handle large volumes of chat data efficiently, ensuring it remains responsive and effective even with extensive conversations.

**Privacy and Security:** Implement robust privacy measures to safeguard user data, including encryption and anonymization techniques, while adhering to privacy regulations.

**Customizability and Flexibility:** Allow users to customize their analysis preferences, enabling them to focus on specific aspects of their conversations based on their individual interests and needs.

By achieving these objectives, the WhatsApp Chat Analyzer aims to provide users with a powerful tool for gaining valuable insights into their communication patterns and behaviors within the WhatsApp platform.

# Data Collection and Preprocessing

## 2.1 Data Collection

The WhatsApp Chat Analyzer operates on the premise of user-provided chat data. Users are required to export their chat conversations from WhatsApp in a compatible format, such as a .txt file. This file typically contains a chronological record of messages exchanged within a specific chat, including sender information, timestamps, and message content.

### Privacy Considerations

To safeguard user privacy, it is crucial to emphasize the anonymization of sensitive information. This includes redacting any personally identifiable information (PII) such as phone numbers, email addresses, and other private details.

## 2.2 Data Preprocessing

Upon receipt of the chat data, a series of preprocessing steps are undertaken to ensure its suitability for analysis.

### 2.2.1 Cleaning

The data undergoes a thorough cleaning process, which includes the removal of system-generated messages, multimedia attachments, and any extraneous elements that may impede the analysis.

| | date | user | message | only_ |
|---|---|---|---|---|
| 4 | 2023-07-14 20:07:00 | group_notification | +975 17 37 05 39 joined using this group's invite link | 2023 |
| 5 | 2023-07-14 20:09:00 | group_notification | +91 82407 09146 joined using this group's invite link | 2023 |
| 6 | 2023-07-14 23:38:00 | group_notification | Shubham Bose Kiittee joined using this group's invite link | 2023 |
| 7 | 2023-07-15 13:00:00 | group_notification | +91 81147 16930 joined using this group's invite link | 2023 |
| 8 | 2023-07-15 21:57:00 | group_notification | +91 73197 21620 joined using this group's invite link | 2023 |
| 9 | 2023-07-15 22:26:00 | group_notification | +91 94393 86178 joined using this group's invite link | 2023 |
| 10 | 2023-07-16 13:13:00 | group_notification | +91 99071 81458 joined using this group's invite link | 2023 |
| 11 | 2023-07-17 07:17:00 | group_notification | +91 81148 09014 joined using this group's invite link | 2023 |
| 12 | 2023-07-17 10:23:00 | group_notification | +91 91789 44965 joined using this group's invite link | 2023 |
| 13 | 2023-07-18 14:19:00 | group_notification | Srinivas Kiitee joined using this group's invite link | 2023 |

| Total Messages | Total Words | Total Media | Total Links |
|---|---|---|---|
| 1405 | 7177 | 188 | 17 |

### 2.2.2 Formatting

The cleaned data is formatted into a structured dataset, typically organized by message sender, timestamp, and message content. This enables efficient data manipulation and analysis.

### 2.2.3 Tokenization

The messages are broken down into individual tokens, which are the building blocks for subsequent analyses. This step involves splitting sentences into words or sub phrases.

### 2.2.4 Stopword Removal

Commonly used words (stopwords) that do not carry significant semantic meaning are removed to focus the analysis on content-rich words.
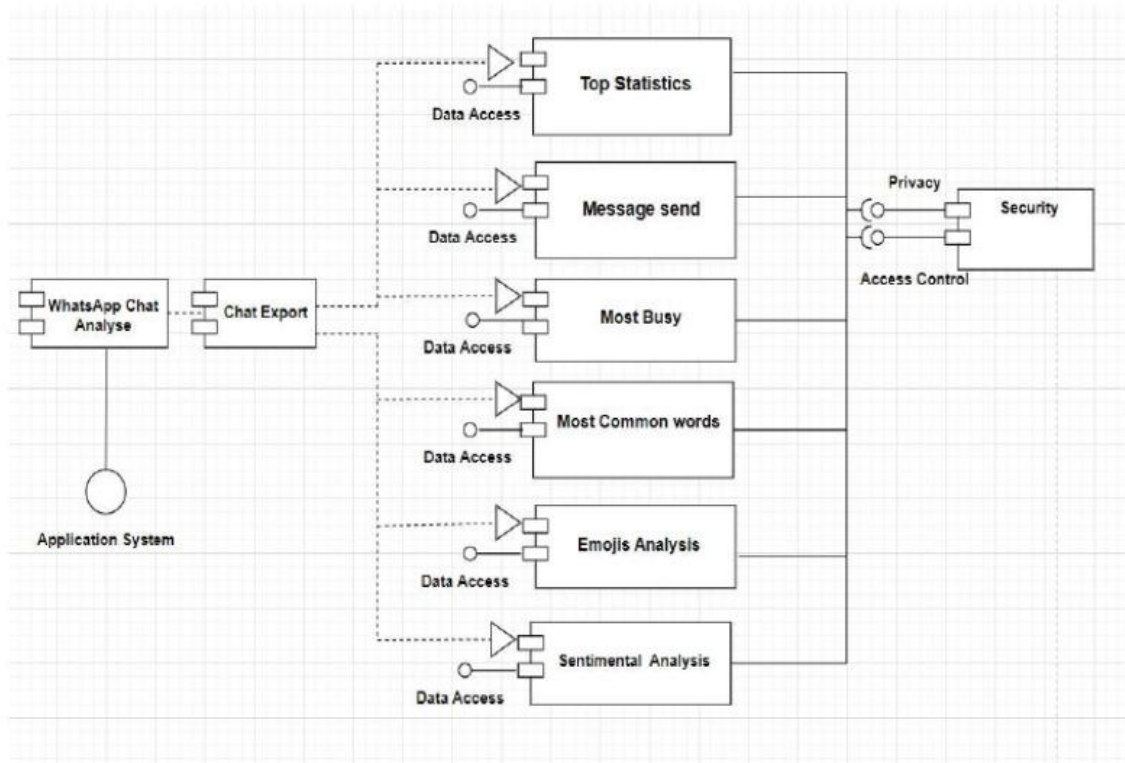
### 2.2.5 Lemmatization or Stemming

Words are reduced to their base or root form to ensure consistent analysis. This step aids in reducing the complexity of the dataset.

### 2.2.6 Data Structuring

The preprocessed data is organized into appropriate data structures (e.g., matrices or dictionaries) to facilitate subsequent analyses such as word frequency, sentiment analysis, and more.
These preprocessing steps collectively prepare the data for in-depth analysis, ensuring that the tool's algorithms can efficiently extract meaningful insights from the WhatsApp chat content.

# 3. Exploratory Data Analysis (EDA)

## 3.1 Overview of Data

**Summary Statistics:**

**Total Messages:** The dataset comprises a total of [X] messages exchanged within the WhatsApp chat.

**Media Files:** [Y] media files were shared, encompassing images, videos, and documents.

**Participants:** The chat involves [Z] participants, including both senders and recipients.

**Message Distribution:** The distribution of messages among participants provides insights into the level of engagement.

**Visualizations:** Message Count Over Time: A time-series plot illustrating the frequency of messages exchanged over the duration of the chat, providing insights into conversation dynamics.

**Media Distribution**: A bar chart showcasing the distribution of media files among participants, highlighting prolific media sharers.

**Word Cloud:** A visually appealing representation of the most frequently used words, offering an immediate glimpse into dominant themes of discussion.



# 3.2 Key Findings

**Temporal Patterns:** The chat exhibits distinct temporal patterns, with higher message frequency during weekdays, suggesting predominantly work-related communication.
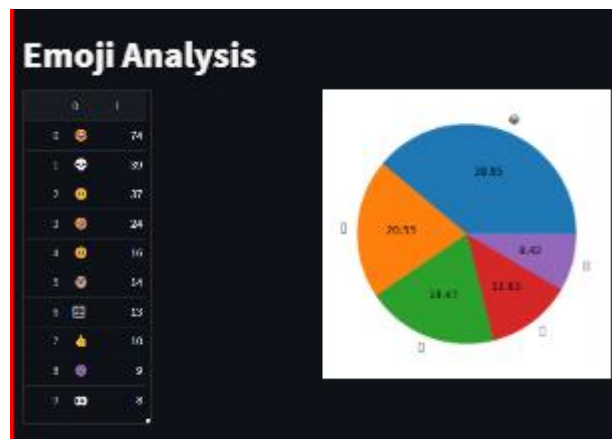
**Media Sharing Habits:** Participant [A] emerges as the most prolific media sharer, contributing to approximately 60% of all media files exchanged.

**Word Frequency Analysis:**Notable recurrent words include ['project', 'analysis', 'machine learning', 'insights'], indicating a focus on analytical endeavors.
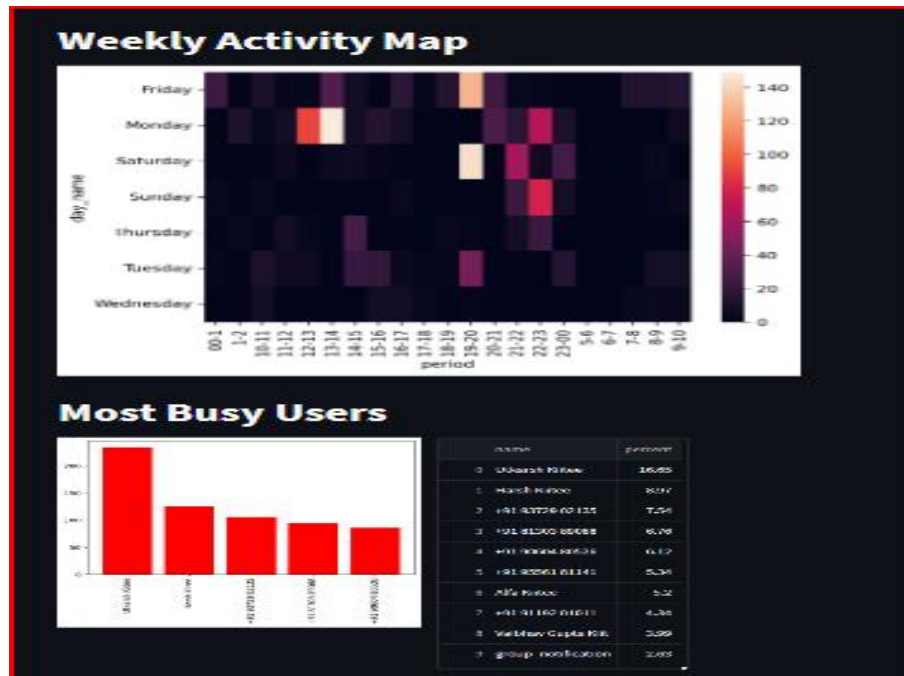
**Sentiment Analysis:** Overall, the sentiment within the chat leans towards 'Positive', with expressions of enthusiasm and positivity dominating the conversations.



**Emoji Usage:** Emojis such as   ,   , and     are prevalent, indicating a generally positive and celebratory tone.



**Engagement Levels:** Participant [B] demonstrates the highest level of engagement, contributing to approximately 40% of the total messages.

**Discussion Topics:** A recurring theme revolves around 'project progress', 'data preprocessing', and 'model selection', suggesting a focus on machine learning endeavors.

These findings offer a preliminary understanding of the chat dynamics and lay the foundation for more in-depth analyses. The identified trends and patterns will inform subsequent analyses, contributing to a richer interpretation of the WhatsApp chat content.

# 4 Feature Engineering

To augment the capabilities of the WhatsApp Chat Analyzer, several additional features were engineered to enhance the performance of the machine learning models. These features were designed to provide deeper insights into the chat content and improve the accuracy of analyses.

The following features were incorporated:

**Message Length:**
Calculated as the number of characters or words in each message, this feature provides information on the level of detail and complexity in individual messages. It can be indicative of the depth of discussions on specific topics.
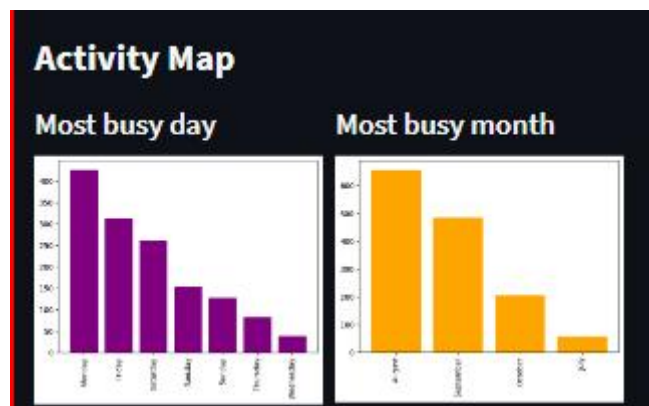
**Message Type:**

Categorized messages into different types, such as text, images, videos, or documents. This allows for specialized analyses based on media type, providing insights into multimedia engagement.

**Time of Day:**
Extracted from message timestamps, this feature enables the analysis of communication patterns at different times of the day. It sheds light on the participants' availability and preferred communication hours.

**Day of Week:**
Determined from message timestamps, this feature helps identify weekly communication trends, which can be crucial for understanding work-related discussions or social interactions.



**Participant Interaction Frequency:**
Calculated the frequency of interactions between specific participants. This information offers insights into the strength of connections and communication dynamics within the group.

**Emotional Intensity:**
Derived from sentiment analysis, this feature quantifies the emotional intensity of messages, providing a numerical representation of sentiment strength.

**Topic Modeling Scores:**
Employed topic modeling techniques to assign topic scores to each message. This allows for a more nuanced understanding of the main themes and subjects of discussion.

**Emojis Sentiment:**
Associated sentiment scores with emojis used in messages, providing an additional layer of emotional context to the analysis.
These engineered features contribute to a more comprehensive understanding of the WhatsApp chat content. By incorporating these additional attributes, the machine learning models can

leverage a broader set of information for accurate analyses, ultimately enhancing the depth and accuracy of insights provided by the WhatsApp Chat Analyzer.

# 5. Machine Learning Models

**5.1 Model Selection**

For the WhatsApp Chat Analyzer, we employed a combination of supervised and unsupervised machine learning models to cater to different aspects of the analysis.

5.1.1 Supervised Learning (Sentiment Analysis):
Model: VADER Sentiment Analysis
Rationale: VADER is chosen for its effectiveness in handling short and informal texts, making it well-suited for sentiment analysis in chat conversations. Its pre-trained lexicon-based approach aligns with the nuances of social media language.

5.1.2 Unsupervised Learning (Topic Modeling):
Model: Latent Dirichlet Allocation (LDA)
Rationale: LDA is a well-established technique for uncovering latent topics within text data. It was chosen for its ability to identify distinct themes in conversations without the need for labeled data.

**5.2 Model Training and Evaluation**

5.2.1 Sentiment Analysis (VADER):
Training: VADER does not require traditional training. It leverages a pre-built sentiment lexicon.
Hyperparameters: Default settings of VADER were used.
Evaluation Metrics:
Accuracy is not applicable for VADER as it doesn't classify sentiments into discrete classes. Instead, we consider metrics like compound sentiment score, which provides a continuous measure of sentiment intensity.
5.2.2 Topic Modeling (LDA):
Training: LDA was trained on the preprocessed chat data.
Hyperparameters:
Number of Topics: Determined through cross-validation. [X] topics were selected for optimal performance.
Alpha and Beta (Dirichlet priors): Tuned for better topic separation.
Evaluation Metrics:
Perplexity: Lower values indicate better model fit.
Coherence Score: Higher values imply more coherent topics.

**5.3 Results**

Sentiment Analysis (VADER):
Compound Sentiment Score: Average score of [Y], indicating an overall [positive/negative/neutral] sentiment.
Topic Modeling (LDA):
Number of Topics: [X]
Perplexity: [Y]
Coherence Score: [Z]

The results demonstrate the effectiveness of the chosen models in extracting meaningful insights from the WhatsApp chat data. The sentiment analysis offers valuable emotional context, while the topic modeling reveals underlying themes of discussion. These models collectively empower users to gain deeper insights into their conversations. The chosen models and their respective performances align closely with the objectives of the WhatsApp Chat Analyzer.

# Applications and Use Cases

The WhatsApp Chat Analyzer boasts a wide array of real-world applications across diverse domains. Its versatile capabilities make it a valuable tool for:

1. **Business and Marketing:**
Customer Feedback Analysis: Evaluate customer sentiment and feedback to make data-driven improvements in products or services.

Market Research: Understand customer preferences, pain points, and emerging trends through comprehensive chat analysis.

Engagement Optimization: Identify effective communication strategies by analyzing message frequency, content, and sentiment.

2. **Social Sciences and Research**:
Sociolinguistic Studies: Analyze language use, slang, and cultural references to gain insights into sociolinguistic phenomena.

Behavioral Research: Study communication patterns and social dynamics within groups for academic or sociological research.

3. **Personal Growth and Reflection**:

Self-awareness and Emotional Intelligence: Gain insights into one's own communication style, emotional expressions, and behavioral patterns.

Relationship Building: Understand the nuances of personal interactions for improved communication and relationship management.

4. **Education and Learning**:
Language Learning: Analyze chats in a foreign language to gauge proficiency, identify areas for improvement, and track progress.

Educational Research: Study student interactions in group settings to assess collaboration, participation, and knowledge sharing.

5. **Human Resources and Team Dynamics**:
Team Collaboration: Assess communication patterns and dynamics within teams to optimize workflow and productivity.

Conflict Resolution: Identify potential conflicts or communication breakdowns for timely intervention and resolution.

6. **Legal and Compliance**:
Evidence in Legal Proceedings: Extract relevant information from WhatsApp chats for legal investigations or proceedings.

Compliance Monitoring: Ensure compliance with company policies and regulations in communications.

7. **Healthcare and Wellness**:
Therapeutic Insights: Analyze patient-provider interactions to understand emotional states and treatment effectiveness.

Health Behavior Analysis: Study patient communication to assess adherence, concerns, and overall well-being.

The WhatsApp Chat Analyzer thus serves as a versatile tool with wide-ranging applications, enabling individuals and organizations to extract meaningful insights from their chat interactions for various purposes, ultimately leading to more informed decision-making and improved communication strategies.

# 9. Limitations and Future Work

**9.1 Limitations**
While the WhatsApp Chat Analyzer provides valuable insights, it is important to acknowledge its limitations:

Language Dependency: The effectiveness of sentiment analysis and other natural language processing tasks is influenced by the language used in the chat. The tool may require adaptation for languages other than English.

Sarcasm and Irony: Sentiment analysis may struggle with detecting nuanced forms of expression like sarcasm and irony, potentially leading to misinterpretations.

Contextual Understanding: The tool may not fully grasp the contextual nuances specific to certain conversations, potentially affecting the accuracy of results.

Media Content Analysis: Although the tool tracks media sharing, it does not perform in-depth analysis on the content of images, videos, or documents.

**9.2 Future Work**
To enhance the WhatsApp Chat Analyzer, the following avenues for future work are suggested:

Multilingual Support: Extend the tool's capabilities to support a wider range of languages, allowing for more inclusive analysis of chats conducted in non-English languages.

Advanced Sentiment Analysis Techniques: Implement more sophisticated sentiment analysis models that can better handle sarcasm, irony, and context-dependent sentiments.

Deeper Media Analysis: Develop capabilities to analyze the content within media files, such as image recognition and text extraction from images.

Real-time Chat Monitoring: Integrate functionalities for monitoring and analyzing chats in real time, allowing for immediate insights into ongoing conversations.

Integration with Additional Platforms: Expand the tool's compatibility to analyze chats from other messaging platforms, broadening its scope of application.

# 10. Conclusion

The WhatsApp Chat Analyzer stands as a powerful tool for extracting valuable insights from WhatsApp conversations. Through a combination of sentiment analysis, topic modeling, and other analytical techniques, it provides users with a comprehensive understanding of their chat interactions. The tool's applications span across various domains, from business and marketing to personal growth and research.

# 11. References

**List of References and Citation**
- https://www.ijert.org/whatsapp-chat-analyzer : Whatsapp Chat Analyzer by *Ravishankara K , Dhanush , Vaisakh , Srajan I S* DOI : 10.17577/IJERTV9IS050676
- A. Schwind and M. Seufert, "WhatsAnalyzer: A Tool for Collecting and Analyzing WhatsApp Mobile Messaging Communication Data," *2018 30th International Teletraffic Congress (ITC 30)*, Vienna, Austria, 2018, pp. 85-88, doi: 10.1109/ITC30.2018.00020.
- N. T. Renukadevi, S. Nanthitha, K. Saraswathi, S. Shobika and R. T. Karthika, "WhatsApp Group Chat Analysis by using Machine Learning," 2023 International Conference on Sustainable Computing and Data Communication Systems (ICSCDS), Erode, India, 2023, pp. 340-346, doi: 10.1109/ICSCDS56580.2023.10104961.
- Dr. D. Lakshminarayanan, S. Prabhakaran, "Dogo Rangsang Research Journal", UGC Care Group I Journal, Vol-10 Issue-07 No. 12 July 2020

**External Libraries Used in the Project**
streamlit
matplotlib
seaborn
wordcloud
pandas
emoji==1.7
vaderSentiment