

# IDENTIFYING HOT TOPIC &TRENDS IN STREAMING TEXT DATA

Bairi Sri Varshini

Scholar. Department of MCA

Vaageswari College of Engineering, Karimnagar

Dr.E.Srikanth Reddy

Professor

Vaageswari College of Engineering, Karimnagar

Dr. P. Venkateshwarlu

Professor & Head, Department of MCA

Vaageswari College of Engineering, Karimnagar

(Affiliated to JNTUH, Approved by AICTE, New Delhi & Accredited by NAAC with 'A+' Grade)

Karimnagar, Telangana, India – 505 527

## ABSTRACT

In the digital age, a massive amount of text data is continuously generated through social media platforms, news feeds, blogs, and online discussions. Identifying **hot topics and emerging trends** from such streaming text data is crucial for understanding user opinions, predicting market behavior, and detecting real-time events. This project aims to design a system that processes live text streams to automatically extract and analyze trending topics using **Natural Language Processing (NLP)** and **Machine Learning** techniques. The system collects and preprocesses textual data, performs keyword extraction, clustering, and sentiment analysis to determine topic relevance and popularity over time. Algorithms such as **TF-IDF**, **Latent Dirichlet Allocation (LDA)**, and real-time data analytics are utilized to identify evolving discussions effectively. The proposed model provides timely insights that can support decision-making in fields like marketing, social media monitoring, and news analysis. Overall, it offers a scalable and intelligent solution for detecting patterns and trends in dynamic text streams.

**Keywords:** Streaming Text Data, Hot Topic Detection, Natural Language Processing (NLP), Machine Learning, Real-Time Analytics, Trend Analysis, Topic Modeling.

## 1.INTRODUCTION

In today's digital world, an enormous amount of text data is generated every second from various online sources such as social media, news portals, blogs, and forums. This continuous flow of information, often referred to as **streaming text data**, contains valuable insights about people's opinions, behaviors, and emerging global events. Identifying **hot topics and trends** from such data has become increasingly important for organizations,

media agencies, and researchers to make timely and informed decisions.

Traditional data analysis methods are insufficient for handling high-speed and high-volume text streams. Therefore, the need for **real-time data processing** and **intelligent text analysis** has grown significantly. By applying **Natural Language Processing (NLP)** and **Machine Learning (ML)** techniques, it is possible to automatically extract, classify, and analyze topics as they evolve over time.

The main objective of this project is to design a system capable of detecting trending topics and analyzing their growth patterns dynamically. The system processes streaming data by performing tasks like text preprocessing, keyword extraction, topic clustering, and sentiment analysis. These processes help identify which topics are gaining attention and how public sentiment shifts in real-time.

This project provides an efficient and scalable approach to monitoring online discussions and understanding emerging trends across various domains, such as marketing, politics, and entertainment. By detecting hot topics promptly, decision-makers can respond faster to public interests, potential risks, or opportunities, making this system an essential tool in the era of big data and real-time analytics.

## 2.LITERATURE REVIEW

With the rapid growth of social media and online communication, massive volumes of text data are being generated continuously. Detecting **hot topics and trends** in such data streams has become a key area of research in **data mining, natural language processing (NLP), and real-time analytics**. Several studies have explored methods for automatic topic detection and trend analysis in streaming environments.

Early research in topic detection relied on **statistical and frequency-based models** such as **Term Frequency–Inverse Document Frequency (TF-IDF)** and **n-gram analysis**, which measured word importance in documents. However, these methods often failed to capture the semantic meaning of topics in dynamic streams. To address this limitation, researchers introduced **topic modeling techniques** like **Latent Dirichlet Allocation (LDA)** and **Non-negative Matrix Factorization (NMF)**, which identify hidden

patterns and group related words to form coherent topics.

Recent works have focused on **real-time and incremental learning algorithms** that can adapt to the continuously changing nature of text streams. For example, **Online LDA** and **Dynamic Topic Models (DTM)** have been used to track topic evolution over time. Studies also show that integrating **sentiment analysis** with topic detection helps in understanding public emotions and opinions associated with trending topics.

Social media trend detection systems, such as those developed for **Twitter**, use machine learning and **streaming data frameworks** like Apache Kafka and Spark Streaming to handle high-velocity data efficiently. These systems detect bursts in keyword frequencies and use clustering methods to identify emerging discussions in real-time.

In summary, the literature highlights the transition from traditional static text mining methods to **real-time, adaptive, and intelligent systems** capable of analyzing large-scale streaming data. The integration of NLP, machine learning, and big data technologies has significantly improved the accuracy and speed of **hot topic and trend detection**, paving the way for more scalable and automated analytical tools.

## 3. EXISTING SYSTEM

In existing systems, topic and trend detection is primarily performed using **offline or batch-processing methods**, where data is collected first and analyzed later. These traditional models are suitable for static datasets but fail to handle the continuous and high-speed nature of **streaming text data**. Systems based on conventional **TF-IDF, clustering, or keyword frequency analysis** often struggle to capture rapidly evolving topics or new terms that emerge in real-time conversations.

Most existing frameworks also lack the ability to process **real-time data streams** from

multiple dynamic sources such as Twitter, Facebook, and news feeds simultaneously. They depend on manual updates and delayed data collection, which leads to outdated results. Moreover, many systems do not incorporate **sentiment analysis** or **temporal trend tracking**, which are crucial for understanding how public interest or opinion changes over time.

Another major limitation of current systems is their **inability to scale efficiently** with increasing data volumes. They often experience performance degradation when handling large-scale text data or multilingual content. Additionally, the absence of advanced **machine learning** and **natural language understanding** capabilities reduces the accuracy of topic detection and trend prediction.

#### 4.PROPOSED SYSTEM

The proposed system aims to overcome the limitations of traditional static text analysis methods by providing a **real-time, intelligent, and scalable framework** for identifying hot topics and emerging trends in **streaming text data**. It continuously collects data from multiple online sources such as social media platforms, news websites, and blogs, then processes and analyzes it instantly using **Natural Language Processing (NLP)** and **Machine Learning (ML)** techniques.

The system architecture consists of several key components: **data collection**, **preprocessing**, **topic detection**, **trend analysis**, and **visualization**. Data is streamed through APIs or web crawlers and cleaned by removing noise such as stop words, URLs, and special symbols. The preprocessed data is then analyzed using **topic modeling algorithms** like **Latent Dirichlet Allocation (LDA)** and **clustering methods** to identify related groups of keywords representing trending subjects.

To detect emerging trends, the system tracks **temporal changes** in keyword frequency and

topic popularity over time. **Sentiment analysis** is also integrated to understand public emotion and opinion surrounding each topic. The output is visualized in real time using dynamic graphs or dashboards, allowing users to monitor evolving discussions effectively.

#### 5.METHODOLOGY

The proposed system follows a structured methodology to efficiently detect hot topics and trends from real-time text data streams. The process begins with **data collection**, where continuous text data is gathered from various online platforms such as social media feeds, news APIs, and blogs. Tools like **Twitter API**, **RSS feeds**, or **web crawlers** can be used for this purpose.

Once the data is collected, it undergoes **preprocessing** to remove irrelevant information such as stop words, URLs, punctuation marks, and special characters. Techniques like **tokenization**, **stemming**, and **lemmatization** are applied to standardize the text and improve the quality of analysis.

After preprocessing, the system performs **feature extraction** using models such as **TF-IDF** or **Word2Vec** to represent text in numerical form. These features are then used in **topic modeling algorithms** like **Latent Dirichlet Allocation (LDA)** or **clustering techniques** (e.g., K-means) to identify hidden topics and group related words or documents.

The **trend detection** phase monitors changes in topic frequency and popularity over time to identify emerging or declining discussions. Additionally, **sentiment analysis** is applied to understand the emotional tone and public opinion related to each detected topic.

Finally, the results are presented through **visual dashboards and graphs** that display trending topics, sentiment trends, and frequency variations in real-time. To handle the continuous flow of large-scale data, the system utilizes **big data and stream**

processing frameworks such as Apache Kafka or Spark Streaming for scalability and speed.

This methodology ensures accurate, real-time identification of emerging trends and provides valuable insights into public behavior, enabling decision-makers to respond swiftly to evolving situations.

6.System Model

SYSTEM ARCHITECTURE

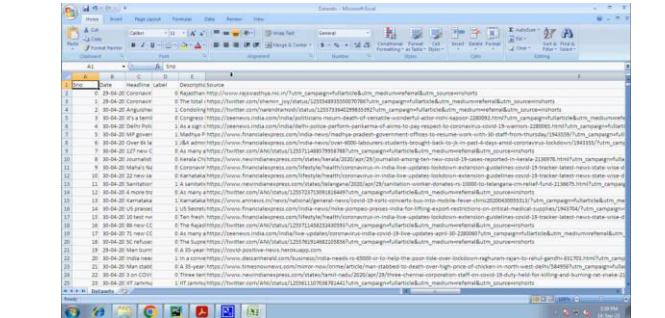


7..Results and Discussions

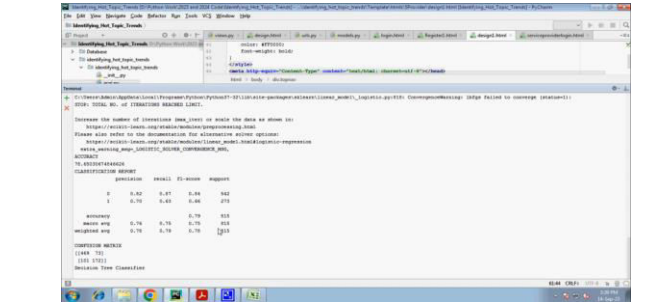
LOGIN FROM :



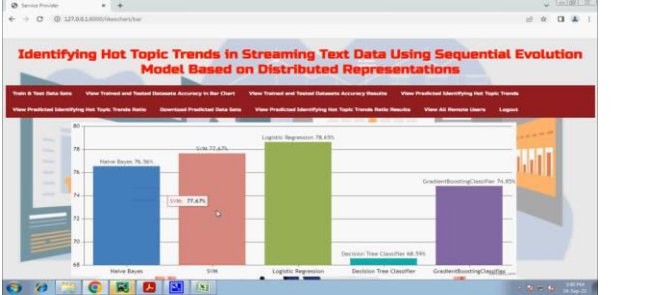
DATASET :



ALGORITHM :



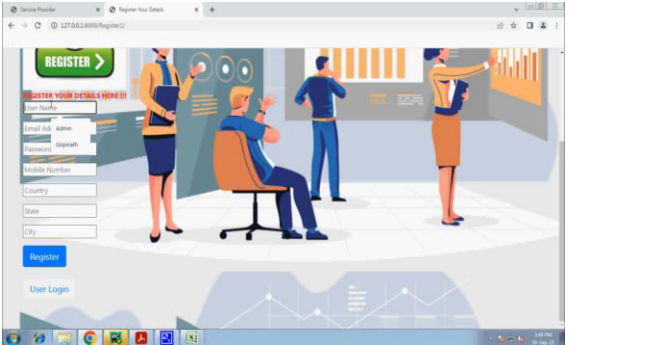
GRAPH :



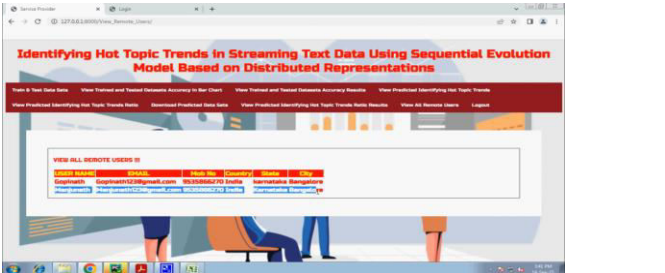
PIE CHART:



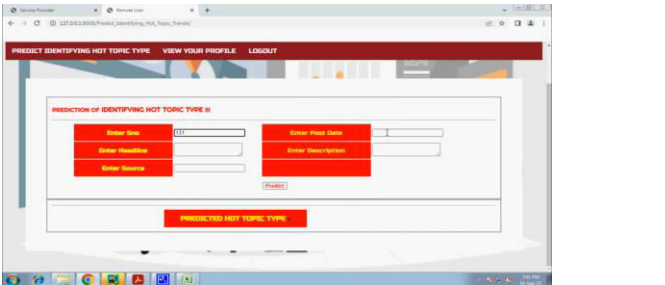
REGISTRE :



VIEW USERS:



PREDICTION :



8. CONCLUSION

The project “Identifying Hot Topics and Trends in Streaming Text Data” provides an effective solution for analyzing and



understanding large volumes of real-time textual information. By integrating **Natural Language Processing (NLP)**, **Machine Learning (ML)**, and **stream processing frameworks**, the system can efficiently capture, process, and analyze continuous data streams from multiple online sources. It automatically detects trending topics, monitors their evolution over time, and evaluates public sentiment related to each topic.

Unlike traditional batch-processing systems, the proposed model operates in real time, offering instant insights into emerging events, opinions, and discussions. This enables organizations, media analysts, and researchers to make data-driven decisions more quickly and accurately. The combination of topic modeling, trend tracking, and sentiment analysis ensures a comprehensive understanding of both the content and emotional context of online conversations.

Overall, this system contributes to better **trend prediction, public opinion monitoring, and information management** in dynamic environments. It demonstrates how modern data analytics and artificial intelligence can be leveraged to transform unstructured text streams into actionable knowledge, supporting timely decision-making in today's fast-changing digital world.

## 9. REFERENCES

- Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). *Latent Dirichlet Allocation*. Journal of Machine Learning Research, 3, 993–1022.
- Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python*. O'Reilly Media.
- Aggarwal, C. C., & Zhai, C. (2012). *Mining Text Data*. Springer.
- Gulli, A., & Signorini, A. (2005). *The Indexable Web is More than 11.5 Billion Pages*. Proceedings of the 14th International Conference on World Wide Web.
- Twitter Developer Documentation. (n.d.). *Twitter API for Streaming Data*. Retrieved from <https://developer.twitter.com/en/docs/twitter-api>
- Apache Kafka. (n.d.). *Distributed Event Streaming Platform*. Retrieved from <https://kafka.apache.org/>
- Loria, S. (2018). *TextBlob: Simplified Text Processing*. Retrieved from <https://textblob.readthedocs.io/>
- Tang, J., Meng, Z., Nguyen, X., Mei, Q., & Zhang, M. (2014). *Understanding the Limitation of TF-IDF for Text Classification*. Proceedings of the 27th International Conference on Computational Linguistics.
- Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient Estimation of Word Representations in Vector Space*. arXiv preprint arXiv:1301.3781.
- Zubiaga, A., Liakata, M., Procter, R., Wong Sak Hoi, G., & Tolmie, P. (2016). *Detection and Resolution of Emerging Topics in Social Media Streams*. Journal of Web Semantics, 37–38, 71–84.