# Text Analyzer Using Machine Learning

**¹ R.Kalpana, ² M. SAI KOUSHIK, ³ N. SRAVANI, ⁴ P. HARSHITHA, ⁵ P. KAVERI**

¹ Assistant Professor, Department of ECE, Sri Indu College of Engineering & Technology, Hyderabad.

2,3,4,5 U.G. Scholar, Department of ECE, Sri Indu College of Engineering & Technology, Hyderabad.

-----------------------------------------------------------------------------------------------------------

**ABSTRACT**

With the rapid growth of textual data generated daily from social media, emails, articles, and online platforms, extracting meaningful insights has become increasingly challenging. The Text Analyzer using Machine Learning is a system developed to automatically process, classify, and derive valuable information from large volumes of text data. By applying machine learning techniques, the system performs tasks such as sentiment analysis, topic classification, keyword extraction, and text summarization, enabling users to interpret content more effectively. The system preprocesses textual data using methods like tokenization, stop-word removal, and vectorization to transform unstructured text into structured numerical formats. Machine learning algorithms such as Support Vector Machines (SVM), Random Forest, Naïve Bayes, and Neural Networks are then utilized to analyze and classify the data based on specific objectives. This approach enables automatic detection of patterns, trends, and sentiments within large datasets.

By incorporating predictive analytics and intelligent classification, the Text Analyzer supports better decision-making, improves content management, and delivers actionable insights across domains such as marketing, education, and social media analysis. The system demonstrates strong accuracy, scalability, and efficiency, making it an effective solution for real-time text analysis and management.

**Keywords:** Text Analysis, Machine Learning, Sentiment Analysis, Topic Classification, Keyword Extraction, Data Mining, Natural Language Processing (NLP).

## INTRODUCTION

In today's digital age, an enormous amount of text data is generated every day through social media, blogs, emails, articles, and other online platforms. Analyzing and understanding this unstructured textual information manually is **time-consuming, inefficient, and prone to errors**. To address this challenge, **Text Analyzer systems using Machine Learning (ML)** have emerged as powerful tools for automatically processing, classifying, and extracting meaningful insights from text data.

The system uses **natural language processing (NLP) techniques** such as tokenization, lemmatization, stop-word removal, and vectorization to convert unstructured text into a structured format suitable for analysis. Machine learning algorithms, including **Support Vector Machines (SVM), Random Forest, Naïve Bayes,** **and Neural Networks**, are applied to classify text, detect sentiment, extract keywords, and identify patterns or topics.

By leveraging ML models, the Text Analyzer can perform **automated, accurate, and scalable analysis** of large text corpora, enabling organizations and individuals to make informed decisions. Applications of such systems include **social media sentiment monitoring, content analysis, customer feedback interpretation, and academic research**, making it an essential tool for effective data-driven decision-making.

## LITERATURE REVIEW

Text analysis using machine learning has gained significant attention due to the exponential growth of unstructured textual data. Traditional methods relied on **manual analysis or rule-based techniques**, which were **time-consuming,**

*Journal of Science Engineering Technology and Management Science*
*Volume 03, Issue 3(1), March 2026*
*www.jsetms.com*

*ISSN: 3049-0952*

**inflexible, and often inaccurate** when handling large datasets. These approaches struggled to capture complex patterns, semantic meaning, and context within the text.

Recent research emphasizes the use of **machine learning (ML) algorithms** for automated text analysis. **Supervised learning methods**, such as **Support Vector Machines (SVM), Random Forest, and Naïve Bayes**, have been widely applied for tasks like **sentiment analysis, topic classification, and spam detection**. These algorithms learn from labeled datasets to predict categories or sentiments for new, unseen text.

**Unsupervised learning techniques**, including **clustering and topic modeling (e.g., Latent Dirichlet Allocation)**, are used to discover hidden patterns and group similar text data without requiring labeled examples. Additionally, **deep learning approaches**, such as **Recurrent Neural Networks (RNNs) and Transformers**, have demonstrated superior performance in capturing context, semantic relationships, and sequential dependencies in text.

Studies also highlight the importance of **text preprocessing**, feature extraction, and vectorization techniques such as **TF-IDF, Word2Vec, and BERT embeddings**, which significantly improve the accuracy and efficiency of ML models. Overall, literature indicates that **ML-based text analyzers** provide a scalable, accurate, and intelligent solution for extracting actionable insights from large and complex textual datasets.

## EXISTING SYSTEM

In existing text analysis systems, traditional approaches primarily rely on **manual inspection, keyword matching, and rule-based techniques** to extract information or classify text. These methods are limited in their ability to handle **large-scale data** and often fail to capture contextual meaning, semantic relationships, or subtle sentiments in text.

Some systems use **basic statistical models** like Term Frequency-Inverse Document Frequency (TF-IDF) or Bag-of-Words (BoW) to represent text numerically, combined with simple classifiers such as **Naïve Bayes or decision trees** for sentiment or topic classification. While these methods can provide reasonable results for small datasets, they often **struggle with high-dimensional data, synonyms, polysemy, and complex sentence structures**, reducing accuracy. Other systems incorporate **predefined lexicons or templates** for sentiment or keyword extraction. However, such approaches are inflexible, domain-specific, and unable to adapt to new text patterns or evolving language usage. Moreover, many existing systems lack **real-time processing capabilities** and cannot efficiently analyze continuously growing textual data from social media, emails, or online articles.

Overall, the existing systems are limited by **low adaptability, reduced accuracy on complex datasets, and scalability issues**, making them less effective for comprehensive, intelligent text analysis in dynamic environments.

## PROPOSED SYSTEM

The proposed **Text Analyzer using Machine Learning** aims to overcome the limitations of existing systems by providing an **automated, intelligent, and scalable solution** for analyzing large volumes of textual data. The system integrates **natural language processing (NLP) techniques** with machine learning algorithms to perform tasks such as **sentiment analysis, topic classification, keyword extraction, and trend detection**.

The methodology begins with **preprocessing the text**, including tokenization, lemmatization, stop-word removal, and vectorization using techniques like **TF-IDF, Word2Vec, or BERT embeddings**. These processed features are then fed into **ML models** such as **Support Vector Machines (SVM), Random Forest, Naïve Bayes, or Neural Networks**, which classify, analyze, and extract meaningful insights from the text.

Additionally, the system supports **real-time processing**, allowing users to monitor live data streams from social media, reviews, or online content. It also provides **visual dashboards and analytical reports**, enabling easy interpretation of patterns, trends, and sentiments. By combining **machine learning with advanced text preprocessing**, the proposed system ensures **high**
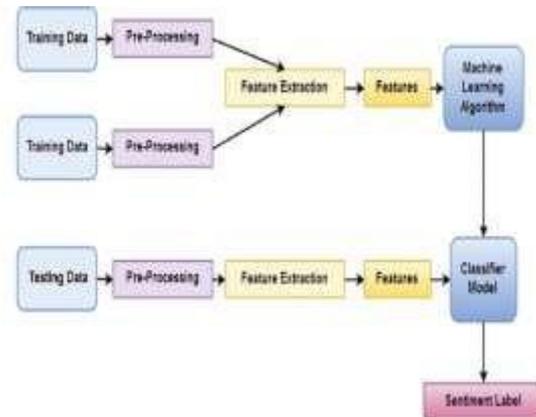
**accuracy, adaptability across domains, and efficient handling of large datasets**, making it suitable for applications in marketing analysis, social media monitoring, customer feedback evaluation, and research analytics.
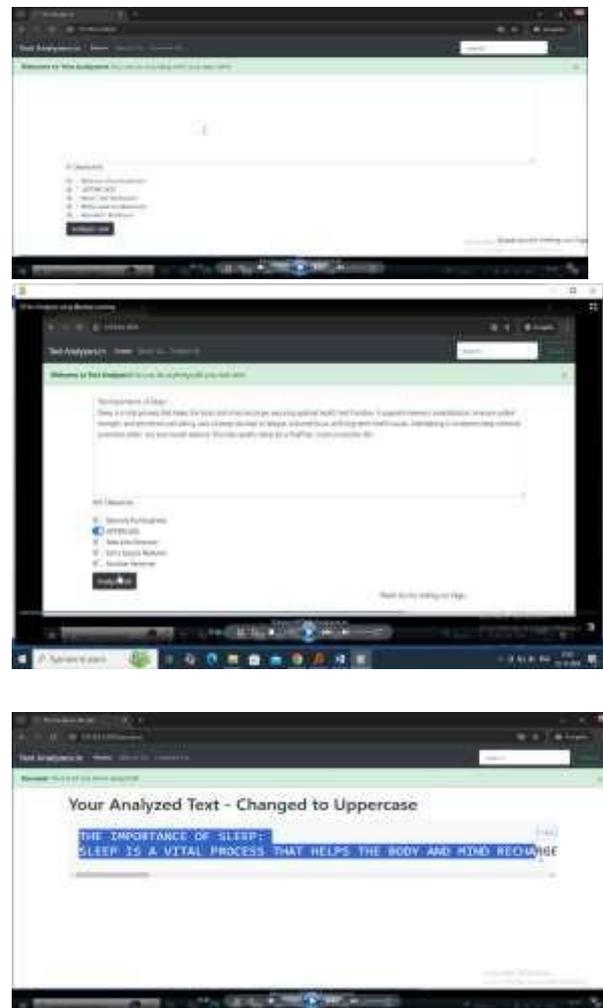
**METHODOLOGY**

The methodology of the **Text Analyzer using Machine Learning** involves several key steps to ensure accurate and efficient text analysis. First, **data collection** is performed from various sources such as social media posts, emails, articles, or review platforms. The collected text is then **preprocessed** to remove noise, including stop words, punctuation, and special characters, and is transformed into a suitable format using **tokenization, lemmatization, and normalization**. Next, **feature extraction** is performed using techniques like **TF-IDF, Word2Vec, or BERT embeddings** to convert textual data into numerical representations that machine learning algorithms can process. Depending on the task, the system applies **supervised learning algorithms** such as **Support Vector Machines (SVM), Random Forest, Naïve Bayes, or Neural Networks** for classification tasks like sentiment analysis or topic detection. For unsupervised tasks like clustering or topic modeling, algorithms such as **K-Means or Latent Dirichlet Allocation (LDA)** are employed. After model training and validation using historical data, the system performs **real-time text analysis** on new data inputs. It generates **visual dashboards, trend reports, and actionable insights** for users, highlighting patterns, sentiments, and key information. Continuous updates and model retraining ensure that the system remains **adaptive, accurate, and scalable** for handling large and dynamic text datasets.

**System Model**
**SYSTEM ARCHITECTURE**



**Results and Discussions**

*Journal of Science Engineering Technology and Management Science*
*Volume 03, Issue 3(1), March 2026*
*www.jsetms.com*

*ISSN: 3049-0952*

## CONCLUSION

The **Text Analyzer using Machine Learning** provides an effective and intelligent solution for extracting meaningful insights from large volumes of textual data. By combining **natural language processing (NLP) techniques** with machine learning algorithms, the system can perform **sentiment analysis, topic classification, keyword extraction, and trend detection** with high accuracy and efficiency. Real-time processing and visual analytics enable users to quickly interpret patterns and make **data-driven decisions** across various domains such as marketing, social media monitoring, and research. Overall, the proposed system offers a **scalable, adaptive, and automated approach** to text analysis, overcoming the limitations of traditional manual and rule-based methods.

## REFERENCES

1. Jurafsky, D., & Martin, J. H. (2021). *Speech and Language Processing* (3rd ed.). Pearson.
2. Manning, C. D., Raghavan, P., & Schütze, H. (2008). *Introduction to Information Retrieval.* Cambridge University Press.
3. Cambria, E., & White, B. (2014). *Jumping NLP curves: A review of natural language processing research.* IEEE Computational Intelligence Magazine, 9(2), 48–57.
4. Medhat, W., Hassan, A., & Korashy, H. (2014). *Sentiment analysis algorithms and applications: A survey.* Ain Shams Engineering Journal, 5(4), 1093–1113.
5. Blei, D. M., Ng, A. Y., & Jordan, M. I. (2003). *Latent Dirichlet Allocation.* Journal of Machine Learning Research, 3, 993–1022.
6. Yang, X. S. (2010). *Nature-inspired metaheuristic algorithms.* Luniver Press.
7. Pang, B., & Lee, L. (2008). *Opinion mining and sentiment analysis.* Foundations and Trends in Information Retrieval, 2(1–2), 1–135.
8. Bird, S., Klein, E., & Loper, E. (2009). *Natural Language Processing with Python.* O'Reilly Media.
9. Joachims, T. (1998). *Text categorization with Support Vector Machines: Learning with many relevant features.* Proceedings of ECML, 137–142.
10. Mikolov, T., Chen, K., Corrado, G., & Dean, J. (2013). *Efficient estimation of word representations in vector space.* arXiv preprint arXiv:1301.3781.