

Automated Detection of Cyber Bullying

Nitya Shree R, Divyashree S, Neha G, Pooja Kulkarni, Poornima K

Department of Information Science and Engineering RNS Institute of Technology, Bengaluru

Abstract:- It is the most popular channels for communication in social media. But few people use these platforms for evil intent, and "cyberbullying" is a particularly common occurrence. Cyberbullying is particularly common among young people and entails using technological methods to harass or injure others. Therefore, the aim of this study is to suggest a deeplearning algorithm-based model for identifying cyberbullying. The Long Short-Term Memory (LSTM) approach was used to forecast bullying incidents using three datasets from Facebook, Instagram, and Twitter. The outcomes showed that an efficient model for identifying cyberbullying has been developed, resolving issues with earlier methods of cyberbullying detection. For the Twitter, Instagram, and Facebook datasets, the model's accuracies were roughly 96.64%, 94.49% and 91.26%, respectively.

I. INTRODUCTION

Social media is so widely used, people are taking advantage of the anonymity and freedom of online communication, which has led to an increase in cyberbullying. With the aim to detect cyberbullying on the three main platforms of Facebook, Instagram, and Twitter, this study presents a reliable LSTM-based model. Leveraging diverse, real-time datasets, the model employs advanced feature extraction strategy to improve the perfection and adaptability of detection across multiple forms of abuse, including sexism, racism, and toxic language. It addresses limitations found in previous models, such as restricted multi-class categorization and insufficient dataset scope, by utilizing a broader range of information and adaptable time-step variations. The results indicate significant accuracy gains, highlighting this approach's ability to successfully identify and prevent detrimental online interactions, ultimately generating a safer, more positive environment for social media users.

II. LITERATURE SURVEY

It examines many analyses and tests that have been oversee in the field of interest, a literature review is crucial. It examines the findings that have previously been made public while accounting for the scope of the project and other project features. The primary objective of a literature review is to conduct a thorough analysis of the project's past, identifying shortcomings in the current configuration and emphasising issues that still require attention. In addition to illuminating the project's past, the topics covered also draw attention to the problems and weaknesses that drove the project's inception and suggested solutions.

2023

Some kids actively engage in cyberbullying, a form of online harassment of others. Many youths are unaware of the risks associated with cyberbullying, which can lead to despair, self-harm, and suicide. Cyberbullying is a major problem that must be addressed because of the significant harm it may cause to a person's mental health. This study sought to provide a technique for determining the intensity of bullying using a deep learning algorithm and fuzzy logic. In this job, Twitter data (47,733 comments) from Kaggle were processed and analysed to identify cyberbullying comments. Keras embedded remarks were fed into a long short-term memory network with four layers for categorization. The seriousness of the situation was then determined using fuzzy logic [2].

2023

Social media platforms have seen a rise in users in recent years, and although they offer several advantages, their drawbacks have also developed in tandem with their user base. Cyberbullying harm has been increased on social networking platforms in our day to day. Finding ways to identify and hold bullies responsible has become important to lessening the prevalence of cyberbullying, which has serious negative impacts on victims' mental and physical health in society. Many attempts have been made to use machine learning and deep learning algorithms with various sets of data collected from social networking sites like Facebook, Instagram, YouTube, Twitter, and others to develop models that can identify and categorize cyberbullying [9].

May.2022

One of the latest social media ills is cyberbullying. The freedom of expression is being abused as social media usage soars. According to statistics, 36.5% of people considered that they have observed cyberbullying at some point in their lives. These figures indicate that we are moving in the wrong path therefore they are more than twice as high as they were in 2007 and have increased from 2018 to 2019. The market has already seen the implementation of solutions to mitigate this problem to some degree. Nevertheless, they have usage restrictions or just don't employ effective algorithms. This project's primary objective is to look into radically novel methods for comprehending and automatically identifying instances of cyberbullying via tweets, comments, and messages on different social media platforms [4].

Jul.2021

As more people utilize the internet and social media, cyberbullying becomes more common. Cyberbullying refers to aggressive and deliberate activity by a group or individual. Negative and destructive stuff is being shared online. It causes psychiatric and emotional issues for people affected.

Developing automated ways for detecting and preventing cyberbullying is crucial. Previously, cyberbullying detection has primarily relied on text analysis. Cyberbullying incidents primarily include the use of text and images. This research introduces a deep neural model that detects cyberbullying using both text and visual data. Deep learning is now considered cutting-edge technology [7].

Sep 2021

In today's business world, Twitter is seen as an important platform for sharing ideas and information. Because of the quantity of the dataset, analysis of this data is crucial and intricate. To comprehend and examine the sentiment of such data, sentiment analysis is used. A machine learning technique is used in this study to analyse the data and determine whether the sentiment (opinion) is favourable or negative. The tweets are cleaned using several combinations of preprocessing approaches, and the feature vector of the tweets is extracted and its dimension reduced using different feature extraction techniques. Using the Sentiment140 dataset, which comprises sentiment labels and tweets, supervised machine learning models—specifically, Support Vector Machine, Naive Bayes, and Logistic Regression—are applied. Based on the outcomes of the experiment, Logistic Regression [8].

Jul.2020

Cyberbullying on social media platforms consists of harassment via insulting words or sensitive content, which is difficult to identify because of the vast volume of figures and linguistic complexity. Traditional detection methods frequently fall short, producing false positives. Recently, deep neural networks, particularly BERT (Bidirectional Encoder Representations from Transformers), have showed great promise in NLP applications due to their capacity to construct contextual embeddings. This paper provides a novel technique to cyberbullying detection based on the pre-trained BERT model and a simple linear neural network layer as a

classifier. The model is trained and assessed on two datasets—one small and one large—and outperforms existing methods in terms of detection accuracy and scalability [3].

Jul.2019

From the perspective of computer science, the scientific investigation of hate speech is relatively new. The present status of the area is arranged and described in this survey, which also offers a structured summary of earlier methodologies, including key characteristics, algorithms, and techniques. In addition to discussing the complexities of the definition of hate speech across many platforms and settings, this work offers a unified definition. Undoubtedly, this field has the potential to have an impact on society, especially in online forums and digital media platforms. A key step in improving automatic hate speech identification is the creation and organization of common resources, including algorithms, annotated datasets in many languages, and guidelines [1].

III. PROPOSED SYSTEM

The suggested automated cyberbullying detection system uses machine learning (ML) as well as natural language processing (NLP) methods to find offensive material in online conversations. To categorise information as bullying or non-bullying, it first gathers and preprocesses data from public databases, messaging apps, and social media. In the process of identifying sarcasm, irony, and emotional nuance—all these are frequently employed in cyberbullying—the algorithm also integrates contextual awareness. When a message is detected as potentially dangerous, it is blocked and automatically responds with reports or warnings. This method guarantees real-time, scalable detection on multiple web platforms. Natural language processing (NLP) and automated cyberbullying detection are combined in the suggested method.

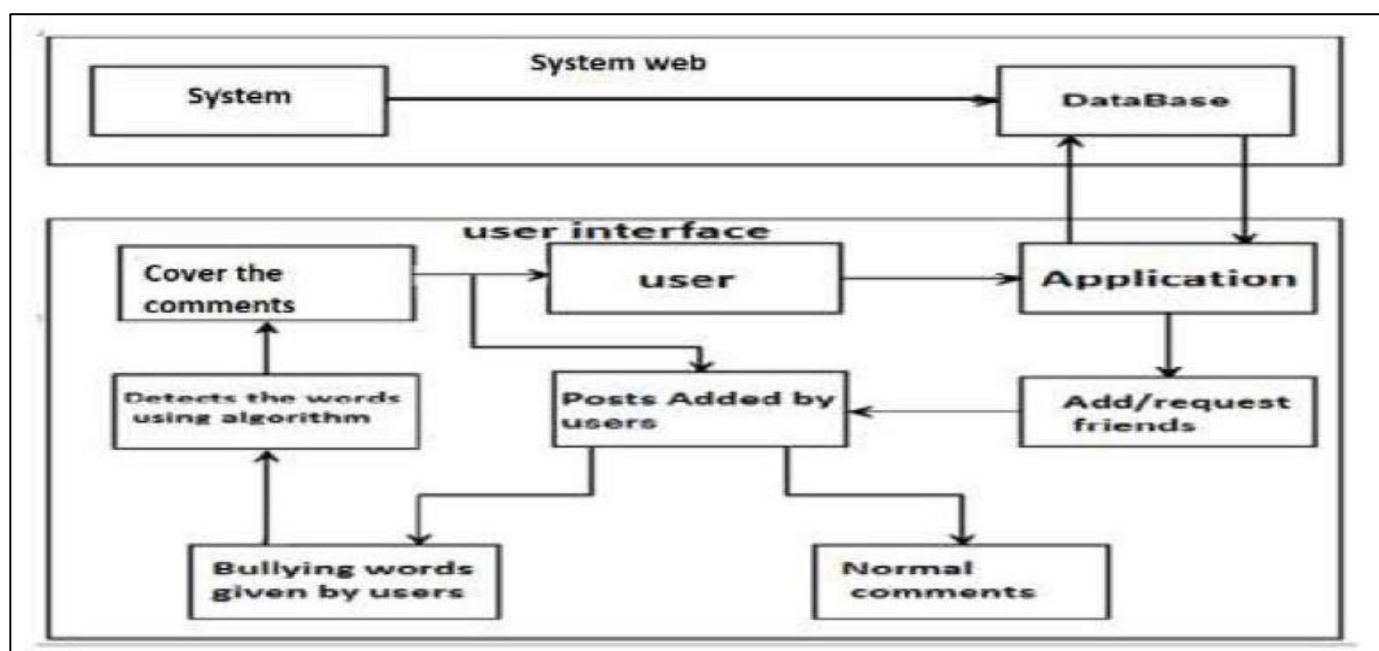


Fig 1: Proposed System

➤ *Advantages of Proposed System*

- **Efficiency and Scalability:** The system's ability to process massive amounts of data in real-time allows for quick detection and intervention across the multiplicity of online platforms, that makes them economical and scalable.
- **Contextual and Nuanced Recognition:** The system can precisely detect subtle types of cyberbullying, such as sarcasm and irony, by utilising sophisticated natural language processing (NLP) and machine learning techniques (e.g., sentiment analysis, sarcasm recognition, BERT, and LSTM).

Regular feedback and retraining improve the system's accuracy over time, allowing it to adjust to new language trends and bullying strategies. Automated detection guarantees the consistent and objective identification of harmful information.

IV. METHODOLOGY

- **Model Setup:** Using Keras, a sophisticated neural network API that is connected with TensorFlow, a sequential model was implemented, enabling efficient deep learning capabilities specifically designed for the detection of cyberbullying.
- **NLP Library:** NLTK, a well-known Python NLP package, was used. It offers crucial text processing capabilities for examining, tokenising, and comprehending linguistic patterns in online conversation.
- **Development Environment:** Code management, version control, and the integration of crucial machine learning libraries were made easier with the usage of PyCharm throughout the development and debugging process.
- **Algorithm Stages:** A four-step algorithm for detecting cyberbullying was created and improved via a great deal of testing to attain the best accuracy and performance in a range of scenarios.

V. CONCLUSION

Through the utilization of numerous datasets and deep learning for feature extraction, the study created an effective cyberbullying detection model that beat earlier methods, increasing detection flexibility. LSTM performance was improved over sigmoid by ReLU activation. To further enhance detection capabilities, future research proposes to include picture, video, and multilingual datasets.

REFERENCES

- [1]. P. Fortuna and S. Nunes, "A survey on automatic detection of hate speech in text," *ACM Comput. Surv.*, vol. 51, no. 4, pp. 1–30, Jul. 2019, doi: 10.1145/3232676.
- [2]. M. H. Obaid, S. K. Guirguis, and S. M. Elkaffas, "Cyberbullying detection and severity determination model," *IEEE Access*, vol. 11, pp. 97391–97399, 2023, doi: 10.1109/ACCESS.2023.3313113.
- [3]. J. Yadav, D. Kumar, and D. Chauhan, "Cyber bullying detection using pre trained BERT model," in *Proc. Int. Conf. Electron. Sustain. Commun. Syst. (ICESC)*, Jul. 2020, pp. 1096–1100, doi:10.1109/ICESC48915.2020.9155700.
- [4]. K. Shah, C. Phadtare, and K. Rajpara, "Cyberbullying detection in hinglish languages using machine learning," *Int. J. Eng. Res. Technol.*, vol. 11, no. 5, pp. 439–447, May 2022.
- [5]. C. Iwendi, G. Srivastava, S. Khan, and P. K. R. Maddikunta, "Cyberbullying detection solutions based on deep learning architectures," *Multimedia Syst.*, vol. 29, no. 3, pp. 1839–1852, Jun. 2023, doi: 10.1007/s00530-020-00701-5.
- [6]. A. Dewani, M. A. Memon, and S. Bhatti, "Cyberbullying detection: Advanced preprocessing techniques & deep learning architecture for Roman Urdu data," *J. Big Data*, vol. 8, no. 1, pp. 1–20, Dec. 2021, doi: 10.1186/s40537-021-00550-7.
- [7]. V. V. and H. P. D. Adolf, "Multi modal cyber bullying detection using hybrid deep learning algorithms," *Int. J. Appl. Eng. Res.*, vol. 16, no. 7, p. 568, Jul. 2021, doi: 10.37622/ijaer/16.7.2021.568-574.
- [8]. M. W. Habib and Z. N. Sultani, "Twitter sentiment analysis using different machine learning and feature extraction techniques," *Al-Nahrain J. Sci.*, vol. 24, no. 3, pp. 50–54, Sep. 2021.
- [9]. N. Haydar and B. N. Dhannoon, "A comparative study of cyberbullying detection in social media for the last five years," *Al-Nahrain J. Sci.*, vol. 26, no. 2, pp. 47–55, 2023.
- [10]. V. Balakrishnan, S. Khan, and H. R. Arabnia, "Improving cyberbullying detection using Twitter users' psychological features and machine learning," *Comput. Secur.*, vol. 90, Mar. 2020, Art.no. 101710, doi: 10.1016/j.cose.2019.101710.