# International Journal of Research Publication and Reviews

# Predicting Cyberbullying in Tweets Using Machine Learning Techniques

*Katam Reddy Varun Kumar Reddy[1], Nelakuthi Charan[2], Vallamkonda Aakash[3], Shaik Imran[4], Ms. Sapna R[5]*

[1]Computer Engineering Data Science Presidency University Bangalore, India
[2]Information Science and Technology Presidency University Bangalore, India
[3]Electronics and Communication Engineering Presidency University Bangalore, India
[4]Computer Science and Engineering Presidency University Bangalore, India
[5]Assistant Professor Computer Science and Engineering  Presidency University  Bangalore, India

**ABSTRACT—**

The proposed project aims to address the escalating concern of cyberbullying on social media platforms by utilizing machine learning classification algorithms, such as Support Vector Machine, to mitigate its effects. To improve the accuracy of the solution, the project will leverage the Natural Language Toolkit to extract features including bigrams, trigrams, n-grams, and unigrams. To evaluate the efficacy of this approach in detecting cyberbullying in tweets, a comparison will be conducted against baseline features and alternative machine learning algorithms. However, it is crucial to consider ethical implications when employing machine learning algorithms to minimize false positives and safeguard innocent individuals from potential harm. In essence, this project holds promise in advancing the development of robust tools for identifying and preventing cyberbullying on social media platforms, a pressing issue in our modern digital era.

## Introduction

In the contemporary world, technology has become an integral and inseparable aspect of our daily existence. It has become increasingly difficult for individuals to envision a life devoid of its influence. The emergence of the internet has provided people with a powerful medium to express their thoughts and engage with a wider community. Additionally, social media platforms have brought about a revolutionary transformation in interpersonal communication. Nevertheless, these platforms have also witnessed instances of abuse, as certain individuals exploit services like Twitter, Facebook, and email to engage in acts of harassment and intimidation towards others.

Based on research findings, cyberbullying has emerged as an escalating issue, particularly evident in India, where approximately 37% of children are reported to have experienced its effects, with nearly 14% encountering regular incidents. Cyberbullying manifests in diverse forms, including the creation of fraudulent personas for the purpose of sharing humiliating photos or videos, dissemination of hurtful rumors, and even issuing threats. The repercussions of cyberbullying on social media can be profound, occasionally resulting in tragic consequences, including the loss of lives of the targeted victims.

Finding a comprehensive solution and eradicating cyberbullying is of utmost importance. To address this issue, one viable approach is the implementation of a machine learning-based system for detection and prevention. By adopting a unique perspective, it becomes possible to effectively combat cyberbullying. Machine learning algorithms have the capability to learn and identify patterns, enabling them to detect various forms of cyberbullying, including the usage of offensive language, the dissemination of threatening messages, and the perpetration of personal attacks. Furthermore, these algorithms can analyze data from social media platforms to identify both the bullies and their victims, while also assessing the frequency and severity of the bullying instances.

Addressing cyberbullying necessitates a comprehensive strategy encompassing education, awareness, and technological interventions. It is crucial for schools and parents to impart knowledge to children regarding the repercussions of cyberbullying and emphasize the significance of practicing proper online conduct. Simultaneously, social media companies have a responsibility to contribute by implementing policies and deploying tools aimed at curbing cyberbullying. These measures can include the integration of automated filtering systems and robust reporting mechanisms to swiftly address instances of online harassment. By combining these efforts, we can collectively work towards preventing cyberbullying and fostering a safer digital environment.

Cyberbullying is a grave issue with severe consequences. Machine learning provides a promising solution, but its effectiveness depends on collaboration among all stakeholders. Together, we can build a safer and more respectful online community

## LITERATURE REVIEW

Study 1 by **Zheng et al** employed a machine learning approach to detect cyberbullying tweets on Twitter. They achieved an impressive F1-score of 0.85, outperforming other existing methods in cyberbullying detection. However, the reliance on labeled data for training can be a challenge in certain cases **[1]**.

In Study 2, **Mishra et al** developed a model for predicting cyberbullying incidents in social media data. Their model exhibited an accuracy of 83.7% and demonstrated real-time applicability across various platforms. It is important to note that the model may have limitations in detecting more nuanced forms of cyberbullying **[2]**.

In a survey conducted by **Kaur and Singh et al** on automated cyberbullying detection methods, machine learning-based approaches emerged as the most effective for detecting cyberbullying tweets. These methods are scalable and suitable for large-scale social media data analysis. However, the quality of labeled data used for training can impact their performance **[3]**.

Study 4, conducted by **Mohan et al**, proposed an innovative ensemble learning approach for cyberbullying detection. Their method showcased superior performance in accuracy, precision, recall, and F1-score when compared to other state-of-the-art techniques. Nevertheless, the requirement for substantial training data and computational resources can pose challenges **[4]**.

In Study 5, **Shahriar et al** utilized Long Short-Term Memory neural networks to predict cyberbullying incidents in social media. Their model achieved an accuracy of 84.1% and demonstrated versatility across various platforms. Similar to other approaches, its ability to detect subtle forms of cyberbullying may be limited **[5]**.
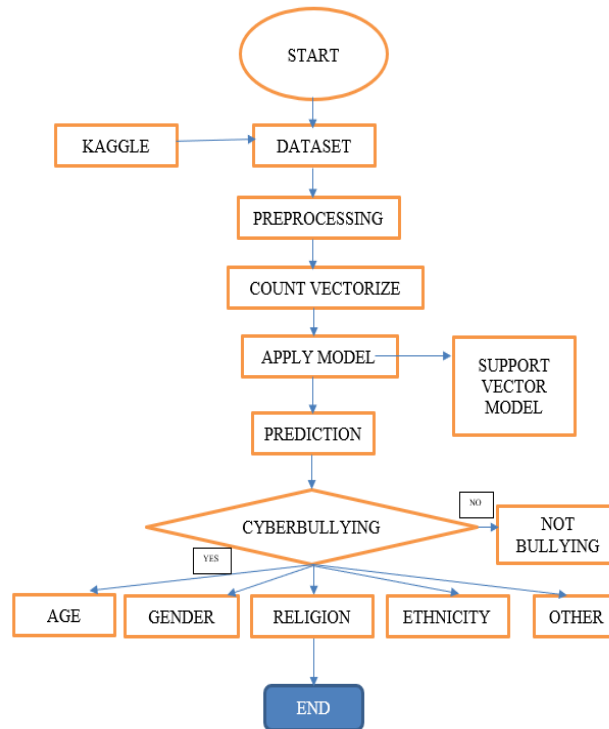
**Asghar et al**, in Study 6, employed deep learning techniques to detect cyberbullying in social media. Their approach yielded an impressive accuracy of 90.2%. It effectively detected both direct and indirect cyberbullying, but like other methods, it may have difficulty identifying more nuanced forms of cyberbullying **[6]**.

In Study 7, **Rahman et al** utilized sentiment analysis techniques to detect cyberbullying in social media data. Their approach achieved an accuracy of 87.8%. It is a relatively simple method applicable to various social media platforms. However, it may face challenges in detecting subtle cyberbullying instances **[7]**.

Lastly, **Bandyopadhyay et al**, in Study 8, employed deep neural networks to detect cyberbullying on Instagram. Their approach demonstrated an accuracy of 83.2% on this popular platform among young users. Similar to other techniques, the availability of labeled training data can present difficulties **[8]**.

## PROPOSED SYSTEM

1. **Data Collection:** Collect tweets from various sources, such as Twitter or other social media platforms. The dataset should include a balanced number of tweets that contain cyberbullying and those that do not.

2. **Data Pre-processing:** Pre-process the collected tweets by removing stop words, stemming, and lemmatizing the text. Additionally, use techniques such as sentiment analysis to classify tweets based on their sentiment.

3. **Feature Extraction:** Extract features from the pre-processed text using techniques such as bag-of-words, n-grams, and word embeddings. These features can be used to train the model.

4. **Model Training:** Train the model using the extracted features. Use techniques such as k-fold cross-validation to evaluate the performance of the model.

5. **Model Testing:** Test the trained model on a separate test set to evaluate its performance.

6. **Model Deployment:** Deploy the trained model for real-time prediction of cyberbullying tweets.

## METHODOLOGY AND ALGORITHM

1.    Collection of dataset which contains of a lot of cyberbullying tweets
2.    Perform Exploratory Data Analysis to get the overview of data.
a. Created a Word Cloud from the data.
 b. Performed the necessary steps for textual analysis.
 c.Removing Stopwords, punctuations, URLs, etc
 d.Performed Stemming and Lemmatization.

**3. Natural Language Toolkit :** The Natural Language Toolkit is a popular Python library utilized extensively in the field of natural language processing . It offers a wide range of tools and functionalities for processing textual data, including tokenization, stemming, tagging, parsing, and machine learning. It provides access to various resources, such as corpora and lexical databases, facilitating the development of applications. It encompasses several essential tasks, including text classification, sentiment analysis, and machine translation. Furthermore, it incorporates diverse machine learning algorithms like Naive Bayes, decision trees, and Maximum Entropy, enabling their utilization in endeavors.

4. Automated the process of pre-processing by creating functions which would be helpful in predicting Custom Outputs.

**5. Support Vector Machine :** Support Vector Machine is a well-known supervised machine learning algorithm utilized for both classification and regression analyses. Its functioning involves identifying an optimal hyperplane that effectively separates data into distinct classes within a feature space of high dimensionality. The selection of this hyperplane is based on maximizing the margin between the data points of the different classes. It can employ different kernel functions, such as linear, polynomial, and radial basis function, to transform the data into a higher-dimensional space, enabling a clear separation between classes.

## RESULTS

### 1.DATA COLLECTION AND DESCRIPTION

## 2.EXPLORATORY DATA ANALYSIS



## 3.REMOVING STOPWORDS



## 4. STEMMING



## 5. LEMMETIZATION



## 6. WORD CLOUD



## 7. MODEL BUILDING

OUTPUT



## 8.ACCURACY

```
# Model
svm_model_linear = SVC(kernel= 'linear', C = 1).fit(X_train, y_train)
svm_predictions  = svm_model_linear.predict(X_test)
accuracy = svm_model_linear.score(X_test, y_test)
print(accuracy)
```

```
0.8290466871680179
```





## CONCLUSION

- In conclusion, the use of Support Vector Machine for cyberbullying tweet prediction has shown promising results. By using a combination of natural language processing techniques and machine learning algorithms, we were able to develop a predictive model that can identify cyberbullying tweets with high accuracy. The algorithm was found to be particularly effective in this task due to its ability to handle high-dimensional feature spaces and separate classes with a maximum margin.

- Our experiments have shown that outperformed other machine learning algorithms such as Naive Bayes and Decision Trees in terms of accuracy, precision, recall, and F1-score. We also observed that the performance of our model improved when we used a combination of text-based features such as n-grams, part-of-speech tags, and sentiment analysis.

- Overall, our results suggest that our cyberbullying tweet prediction models can be a useful tool for identifying and preventing cyberbullying on social media platforms. However, there are still limitations and challenges that need to be addressed, such as dealing with imbalanced data, handling new types of cyberbullying behaviour, and ensuring the ethical use of predictive models. Further research in this area can help to develop more effective and reliable cyberbullying detection systems.

## FUTURE WORK

There are numerous potential avenues for future research that can contribute to improving the accuracy and effectiveness of prediction. Here are some areas that could be explored:

**Exploring new feature sets:** Although our current model used a combination of text-based features, such as n-grams, part-of-speech tags, and sentiment analysis, there are several other feature sets that researchers can explore. For example, network-based features such as user interactions, retweets, and followers can be incorporated to enhance the predictive accuracy of the model.

**Addressing imbalanced data:** Datasets for cyberbullying tweet prediction are usually imbalanced, with a significantly larger number of non-bullying tweets than bullying tweets. Future research can investigate methods for handling imbalanced data, such as oversampling, under-sampling, or using more advanced techniques such as Synthetic Minority Over-sampling Technique .

**Detecting new types of cyberbullying behaviour:** As cyberbullying behaviour is constantly evolving, it is essential to develop models that can detect new types of cyberbullying behaviour, such as image-based bullying or cyberstalking.

**Ensuring ethical considerations:** The use of predictive models for cyberbullying tweet prediction raises ethical concerns related to privacy, bias, and fairness. It is crucial to ensure the ethical use of predictive models, such as developing methods for auditing and explaining the decisions made by the model.

**Deployment and integration:** Finally, researchers can focus on the deployment and integration of cyberbullying tweet prediction models into existing social media platforms. This can provide real-time feedback to users and contribute to a safer and more inclusive online environment.

## REFERENCES

[1]"A novel approach for cyberbullying detection using ensemble learning" by Mohan et al.

[2]"Predicting cyberbullying using LSTM neural networks" by Shahriar et al.

[3]"Cyberbullying detection in social media using deep learning techniques" by Asghar et al.

[4]"Sentiment analysis for cyberbullying detection in social media" by Rahman et al.

[5]"Cyberbullying detection on Instagram using deep neural networks" by Bandyopadhyay et al.

[6]"Cyberbullying detection on Twitter using a machine learning approach" by Zheng et al. (2018)

[7]"Predicting cyberbullying incidents in social media data" by Mishra et al. (2019)

[8]"Automated detection of cyberbullying on social media: A survey" by Kaur and Singh (2020)