

# International Journal of

Information Technology & Computer Engineering



Email: ijitce.editor@gmail.com or editor@ijitce.com

# Detection of Cyber bullying on Social Media Using Machine Learning

<sup>1</sup>K SUPARNA, <sup>2</sup>S. FARHEEN

<sup>1</sup>(Assistant Professor), MCA, DANTULURI NARAYANA RAJU COLLEGE(A) PG COURSES,

BHIMAVARAM ANDHRA PRADESH

<sup>2</sup>MCA, scholar, DANTULURI NARAYANA RAJU COLLEGE(A) PG COURSES, BHIMAVARAM

ANDHRA PRADESH

# **ABSTRACT**

Teens and adults alike are not immune to the pervasive online scourge of cyberbullying.Problems like sadness and suicide have resulted from it. The need for content regulation on social media platforms is on the rise. This research takes a look at two types of cyberbullying—hate speech from Twittter and comments based on personal attacks from Wikipedia model forums—to create a for detecting cyberbullying in text data using NLP and ML. In order to determine the optimal strategy, three feature extraction techniques and four classifiers were investigated. model's accuracy for Tweet data is above 90% and for Wikipedia data it's

80%. When it comes over to social interpersonal connection, networking sites are second to none. Despite the fact that social media use has skyrocketed over the years, many still engage in harmful, immoral behavior platforms. on these Sometimes this occurs between young adults and other times it occurs between teenagers. One of the worst they do is things engage cyberbullying. We have no way of knowing whether someone is speaking online for pleasure or if there is some ulterior motive behind their words..

# 1.INTRODUCTION

More than ever before, technology is intrinsic to how we live our lives. Thanks to the development of the



internet. Trending topics in social media right now. The same is true with misusers; they may appear late or early, but they will always be there. Today, cyberbullying is a prevalent problem. Social networking sites are great for people to communicate with one other. Although more and more individuals are using social media, many still post nasty, immoral, and unethical content. When this occurs, it is usually between teenagers or young adults. They engage in harmful activities such as cyberbullying. Whether someone is saying something for fun or has ulterior motives, it's hard to tell in an online setting. Saying something like, "or don't take it so seriously," usually gets people to laugh it off. A person is being bullied or targeted when they utilize technology to threaten, shame, or harass another individual. Threats to people's physical safety are a common outcome of these online disputes. A number of individuals have attempted suicide. At the outset, it is essential to halt such actions. If someone's tweet or post is deemed objectionable, for instance, it may be possible to delete or suspend their account for a certain amount of time in

#### Volume 12, Issue 3, 2024

order to prevent this. What exactly is cyberbullying, then?

Whether done in jest or with malice, cyberbullying may take many forms, including but not limited to threats, embarrassment, and harassment. Literature Review on Cyberbullying In 2018, a poll conducted by the Indian non-governmental organization Child Right and You found that 11.4% of the 720 youths questioned in the NCT DELHI had experienced cyberbullying. Almost half of those affected did not even inform their instructors, parents, or guardians about the incident. Internet users between the ages of 13 and 18 who spent three hours per day online were almost twice as likely to be victims of cyberbullying as those who spent more than four hours per day online. One must consider

Children and young adults (those between the ages of 13 and 20) have significant challenges related to their physical and mental well-being as well their capacity to make sound decisions in all areas of life, according several publications. **Scientists** should believe that all nations



investigate this issue thoroughly. Lots of kids in Russia and other countries committed suicide in 2016 after an event known as Blue Whale Challenge. A connection between a game's administrator and a player blossomed as the game expanded across several social networks. Over the course of fifty days, participants are assigned specific tasks to complete. Like getting up at 4:30 in the morning or seeing a scary movie, they're simple at first. They started off little, but eventually progressed to selfharm, which ultimately led to suicides. It was only afterwards that it was discovered that the administrators were really 12–14 year olds.

# 2.LITERATURE SURVEY

The topic of cyberbullying and its detection on social media platforms has been the subject of much study. By combining keyword matching, opinion mining, and social network analysis, Ting,IFig Hsien[1] was able to get a recall of 0.71 and a precision of 0.79 utilizing datasets from four different websites. According to the theory put forward by Patxi Gal'an-Garc'ia et al. [2], cyberbullies who use social media platforms often maintain a genuine

#### Volume 12, Issue 3, 2024

profile alongside their phony one, in order to gauge how others perceive it. To find these profiles, they suggested using a machine learning method. As part of the identifying procedure, we looked at profiles that are similar to them. The procedure included picking profiles to analyze, collecting data from tweets, choosing attributes to utilize from profiles, and finally, using ML to identify the tweet's source. A total of 1900 tweets from 19 distinct accounts were analyzed. When it came to author identification, it was 68% accurate. Afterwards, it was used in a Case Study at a Spanish school, where the system was successful in identifying the true owner of a profile among many pupils who were suspected of cyberbullying. A few problems remain with the following approach. Such programs or experts may modify writing styles and behaviors such that no patterns are identified, for instance, in a scenario when the trolling account does not have an actual account. More effective algorithms are required for modifying writing styles. In their collaborative detection technique, Mangaonkar et al. [3] suggested using a network of interconnected detection



nodes, each of which might utilize a unique or shared algorithm to generate findings. The authors are P. Zhou and colleagues. A 93% success rate was achieved by Banerjee et al.[4] by using KNN with updated embeddings.A dataset called Formpring, proposed by Kelly Reynolds, April Kontostathis, and Lynne Edwards [6], uses machine learning algorithms and oversampling to achieve a recall of 78.5%. This is achieved despite the fact that there is an imbalance in the dataset when it comes to cyberbullying postings. The most recent Google language model, BERST, which produces contextual embeddings for categorization, was used by Jaideep Yadav, Kumar, and Chauhan. With data from form spring, the model produced an F1 score of 0.94; with data from Wikipedia, the score was 0.81. In their study, Sweta Agrawal and Amit Awekar [5] used the same datasets to train Deep

However, their main emphasis was on using curse words as features for the job. In doing so, they identified platform-specific differences in the lexicon for such models. Building on previous work by Yasin N. Silva, Christopher Rich,

Neural Networks.

#### Volume 12, Issue 3, 2024

and Deborah Hall[6], BullyBlocker is a smartphone app that alerts parents when their kid is the target of cyberbullying on Facebook. The program takes into account warning indicators and susceptibility characteristics to determine the likelihood of cyberbullying occurring.

# 3. EXISTING SYSTEM

Hsieh [1] obtained a recall of 0.71 and a precision of 0.79 using datasets obtained from four websites by utilizing a technique that included opinion mining, social network analysis, and keyword matching. According to the theory put forward by Patxi Gal'an-Garc'ıa et al. [2], cyberbullies who use social media platforms often maintain a genuine profile alongside their phony one, in order to gauge how others perceive it. To find these profiles, they suggested using a machine learning method. As part of the identifying procedure, we looked at profiles that are similar to them.

The procedure included picking profiles to analyze, collecting data from tweets, choosing attributes to utilize from profiles, and finally, using ML to identify the tweet's source. A total of 1900 tweets from 19



distinct accounts were analyzed. When it came to author identification, it was 68% accurate. Afterwards, it was used in a Case Study at a Spanish school, where the system was successful in identifying the true owner of a profile among many pupils who were suspected of cyberbullying. A few problems remain with the following approach

As an example, consider a scenario where a trolling account lacks a legitimate account, allowing it to deceive systems or experts that may alter writing styles and behaviors to avoid detection of patterns. More effective algorithms are required for modifying writing styles. In their collaborative detection technique, Mangaonkar et al. [3] suggested using data and outcomes from many interconnected detection nodes, each of which might employ a different or same algorithm. A B-LSTM was proposed by P. Zhou et al. [4]. method that relies on mental focus. According to Banerjee et al. [5]. used KNN together with updated embeddings to achieve a 93% level of accuracy. The downsides It is not possible to construct a vocabulary from every text. Either every word (token) in every document or only the most frequently used ones could make up the vocabulary. Even though it gets its

### Volume 12, Issue 3, 2024

vocabulary characteristics in the same manner as the bag of words model, the Tf-Idf technique is distinct from it.

New Approach In this research, we tackle the subject of cyberbullying detection as a binary classification problem. Specifically, we are looking for two main types of cyberbullying: hate speech on Twitter and personal assaults on Wikipedia. We are trying to determine whether these content types include cyberbullying or not.

In tokenization, we break down raw text into smaller, more manageable pieces called tokens. Tokenizing the sentence "we will do it" into its component parts ('we,' "will," "do," and "it") is one example. Word tokenization and phrase tokenization are two different approaches to tokenization. Regex Tokenizer is one of several versions of tokenization that we utilize in this project. The decision-making process for tokens in a regex tokenizer is governed by rules, specifically regular expressions. We choose tokens that match the given regular expression, for example, The regular expression '\w+' is used to extract all the alphanumeric tokens. The term "stemming" refers to the method of breaking a word



down into its component parts. For instance, the stem of the triliteral words "eating," "eats," and "eaten" is "eat." It is reasonable to assume that the three terms that stem from the root "eat" mean the same meaning. Porter, Lancaster, Snowball, and Regexp stemmers are the four varieties available from NLTK. One project that makes use of PorterStemmer is this one.

Eliminating superfluous words: In English, superfluous words like "what," "is," "at," and "a" are examples of stop words. You may eliminate these words since they are unnecessary. You may filter out all the tweets using NLTK's collection of English stop words. When training deep learning and machine learning models, it is common practice to delete stop words from text input. This is done because the information that stop words give is useless to the model and may actually improve its performance. The Benefits The Common Bag of Words model takes in a set of input words and makes a word prediction depending on the surrounding text. You may enter a single word or a string of words. The CBOW model averages the input words' contexts, but each word might have two different interpretations selected. in other words, we

#### Volume 12, Issue 3, 2024

may anticipate two Apple vectors. The first one is for the fruit, while the second one is for the firm.

# 4. OUTPUT SCREENS

#### User:

# **Userlogin:**



# **Registration:**



# **PredictCyberbullying:**



# Viewyourprofile:



**Viewallremoteusers:** 





ServiceProvider:

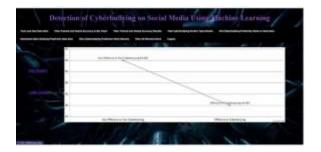
**Adminlogin:** 



**Trainandtestdatasets:** 



Viewcyberbullyingpredictionratioresults(l inechart):





5. CONCLUSION

As a result of the serious consequences it causes (e.g., suicide, depression, etc.), it is imperative that cyberbullying be curbed. As result, identifying cyberbullying on social media is crucial. improved More data and user information classification means more opportunities for cybercriminals launch a wide range of assaults. Social media platforms may use cyber bullying detection systems to block users that engage in cyberbullying. In this study, we provide a solution to the problem by outlining an architecture that identify cyberbullying. Two data sets were covered in our discussion: hate speech data from Twitter and personal assaults data from Wikipedia. Tweets with hate speech were readily identifiable due to the prevalence of profanity, therefore Natural Language



Processing approaches using simple Machine Learning algorithms were successful with accuracies of over 90%. For this reason, BOW and TF-IDF models provide superior outcomes compared to Word2Vec models. While all three feature selection approaches worked similarly, it was particularly challenging to utilize the same model to identify personal assaults in comments that lacked a common emotion that could be learnt.Word2Vec models that take feature context into account performed well in both datasets, producing comparable results with less features when paired with Multi Layered Perceptrons.

# **6.REFERENCES**

- [1] M. Di Capua, E. Di Nardo and A. Petrosino, Unsupervised cyberbullying detection in social networks, ICPR, pp. 432-437, doi: 10.1109/ICPR.2016.7899672. (2016)
- [2] D. Poeter. (2011) Study: A Quarter of Parents Say Their Child Involved in Cyberbullying. pcmag.com. [Online]. Available: <a href="http://www.pcmag.com/article2/0,2817,2388540">http://www.pcmag.com/article2/0,2817,2388540</a>, 00.asp

#### Volume 12, Issue 3, 2024

- [3] J. W. Patchin and S. Hinduja, "Bullies move Beyond the Schoolyard; a Preliminary Look at Cyberbullying," Youth Violence and Juvenile Justice, vol. 4, no. 2, pp. 148–169,2006
- [4] Anti Defamation League. (2011) Glossary of Cyberbullying

Terms.adl.org.[Online].Available:
http://www.adl.org/educati on/curriculum
connections/cyberbullying/glossary.pdf

- [5] N. E. Willard, Cyberbullying and Cyberthreats: Responding to the Challenge of Online Social Aggression, Threats, and Distress. Research Press, 2007.
- [6] D. Maher, "Cyberbullying: an Ethnographic Case Study of one Australian Upper Primary School Class," Youth Studies Australia, vol. 27, no. 4, pp. 50–57, 2008.
- [7] D. Yin, Z. Xue, L. Hong, B. D. Davison, A. Kontostathis, and L. Edwards, "Detection of Harassment on Web 2.0," in Proc. Content Analysis of Web 2.0 Workshop (CAW 2.0), Madrid, Spain, 2009.
- [8] K. Dinakar, R. Reichart, and H. Lieberman, "Modeling the Detection of Textual Cyberbullying," in Proc. IEEE International Fifth International AAAI Conference on Weblogs and Social Media (SWM'11), Barcelona, Spain, 2011.
- [9] I. H. Witten and E. Frank, Data Mining: Practical Machine Learning Tools and



Techniques, Second Edition. San Francisco, CA: Morgan Kauffman, 2005.

- [10] R. Quinlan, C4.5: Programs for Machine Learning. San Mateo, CA: Morgan Kauffman, 1993.
- [11] W. W. Cohen, "Fast Effective Rule Induction," in Proc. Twelfth International Conference on Machine Learning (ICML'95), Tahoe City, CA, 1995, pp. 115–123.
- [12] D. W. Aha and D. Kibler, "Instance-based Learning Algorithms," Machine Learning, vol. 6, pp. 37–66, 1991.
- [13] J. C. Platt, "Fast Training of Support Vector Machines using Sequential Minimal Optimization," Advances in Kernel Methods, pp. 185–208, 1999. [Online]. Available: <a href="http://portal.acm.org/citation.cfm?id=299094.29">http://portal.acm.org/citation.cfm?id=299094.29</a>

### [14]

https://www.sciencedirect.com/topics/computer science/deep-neural-network

- [15] An Effective Approach for Cyberbullying Detection and avoidance ieee paper
- [16] Approaches to Automated Detection of Cyberbullying: A Survey ieee paper
- [17] Cyberbullying Detection System on Twitter ieee paper

#### Volume 12, Issue 3, 2024

- [18] Methods for Detection of Cyberbullying: A Survey ieee paper
- [19] Using Machine Learning to Detect Cyberbullying ieee paper
- [20] Deep Learning Algorithm for Cyberbullying Detection ieee paper
- [21] Online Social Network Bullying Detection Using Intelligence Techniques ieee paper
- [22] Livingstone S, Haddon L, Görzig A, Ólafsson K. Risks and safety on the internet: The perspective of European children. Initial Findings. London: EU Kids Online; 2010.
- [23] Tokunaga RS. Following You Home from School: A Critical Review and Synthesis of Research on Cyberbullying Victimization. Computers in Human Behavior. 2010; 26(3):277–287.
- [24] Mckenna KY, Bargh JA. Plan 9 From Cyberspace: The Implications of the Internet for Personality and Social Psychology. Personality & Social Psychology Review. 1999;4(1):57–75.
- 25] Gross EF, Juvonen J, Gable SL. Internet Use and Well-Being in Adolescence. Journal of Social Issues. 2002;58(1):75–90.
- [26] Juvonen J, Gross EF. Extending the school grounds?—Bullying experiences in cyberspace. Journal of School Health. 2008;78(9):496–505. pmid:18786042



- [27] Hinduja S, Patchin JW. Bullies Move Beyond the Schoolyard: A Preliminary Look at Cyberbullying. Youth Violence And Juvenile Justice. 2006;4(2):148–169.
- [28] Van Cleemput K, Bastiaensens S, Vandebosch H, Poels K, Deboutte G, DeSmet A, et al. Zes jaar onderzoek naar cyberpesten in Vlaanderen, België en daarbuiten: een overzicht van de bevindingen. (Six years of research on cyberbullying in Flanders, Belgium and beyond: an overview of the findings.) (White Paper). University of Antwerp & Ghent University; 2013.
- [29] Livingstone S, Haddon L, Vincent J, Giovanna M, Ólafsson K. Net Children Go Mobile: The Uk report; 2014. Available from: http://netchildrengomobile.eu/reports. [Accessed 30th March 2018].
- [30] O'Moore M, Kirkham C. Self-esteem and its relationship to bullying behaviour. Aggressive Behavior. 2001;27(4):269–283.
- [31] Fekkes M, Pijpers FIM, Fredriks AM, Vogels T, Verloove-Vanhorick SP. Do Bullied Children Get Ill, or Do Ill Children Get Bullied? A Prospective Cohort Study on the Relationship Between Bullying and HealthRelated Symptoms. Pediatrics. 2006;117(5):1568–1574. pmid:16651310

- [32] Cowie H. Cyberbullying and its impact on young people's emotional health and well-being. The Psychiatrist. 2013;37(5):167–170.
- [33] Price M, Dalgleish J. Cyberbullying: Experiences, Impacts and Coping Strategies as Described by Australian Young People. Youth Studies Australia. 2010;29(2):51–59.
- [34] Van Royen K, Poels K, Daelemans W, Vandebosch H. Automatic monitoring of cyberbullying on social networking sites: From technological feasibility to desirability. Telematics and Informatics. 2014;.
- [35] Cortes C, Vapnik V. Support-Vector Networks. Machine Learning. 1995;20(3):273–297.
- [36] Vandebosch H, Van Cleemput K. Cyberbullying among youngsters: profiles of bullies and victims. New Media & Society. 2009;11(8):1349–1371.
- [37] Vandebosch H, Van Cleemput K. Defining cyberbullying: a qualitative research into the perceptions of youngsters. Cyberpsychology and behavior: the impact of the Internet, multimedia and virtual reality on behavior and society. 2008;11(4):499–503.
- [38] Nahar V, Al-Maskari S, Li X, Pang C. Semi-supervised Learning for Cyberbullying Detection in Social Networks. In: ADC.Databases Theory and Applications.



Springer International Publishing; 2014. p. 160–171.

[39] Galán-García P, Puerta JGdl, Gómez CL, Santos I, Bringas PG. Supervised machine learning for the detection of troll profiles in twitter social network: application to a real case of cyberbullying

[40] Rehurek R, Sojka P. Software framework for topic modelling with large corpora. In: The LREC 2010 Workshop on new Challenges for NLP Frameworks. University of Malta; 2010. p. 45–50.