

# Social Media based Hate Speech Detection using Machine Learning

Dr. Nisha Auti<sup>1</sup>, Shreeraj Ghadge<sup>2</sup>, Rajdatta Jadhav<sup>3</sup>, Prajwal Jagtap<sup>4</sup>, Sumit Ranaware<sup>5</sup>

<sup>\*2,3,4,5</sup>Student, Department of Computer Engineering, JSPM Narhe Technical Campus, Pune, Maharashtra, India

<sup>1</sup>HOD, Associate Professor, Department of Computer Engineering, JSPM Narhe Technical Campus, Pune, Maharashtra, India

## ABSTRACT

Hate speech is a crime that has been increasing in recent years, not only in person but also online. There are several causes for this. There is tremendous growth in social media that promotes full freedom of expression through anonymity features. Freedom of expression is a human right, but hate speech directed at individuals or groups on the basis of race, caste, religion, ethnicity or nationality, gender, disability, gender identity, etc. is a violation of that sovereignty. Freedom of expression is a human right, but hate speech directed at individuals or groups on the basis of race, caste, religion, ethnicity or nationality, gender, disability, gender identity, etc. is a violation of that sovereignty. It promotes violence and hate crimes, creates social imbalances, and undermines peace, trust and human rights. Revealing hate speech in social media discourse is a very important but complex task. On the one hand, the anonymity provided by the Internet, especially social networks, makes people more likely to engage in hostile behavior. On the other hand, the desire to express one's thoughts on the Internet has increased, leading to the spread of hate speech. Governments and social media platforms can benefit from detection and prevention technologies, as this kind of bigoted language can wreak havoc on society. We help resolve this dilemma by providing a systematic overview of research on this topic in this survey. This project aims to accurately predict various forms by addressing different categories of hate individually and examining a set of text mining functions. Hate speech detection.

**Keywords:** Hate Speech, Machine Learning, SocialMedia, Social Network, Multi-Class Hate Speech, Natural Language Processing, Hate Speech Classification, Social Media Microblogs, Multi-Class Hate Speech Dataset, Twitter Hate Speech, Text Mining, Features Exploration

## Article Info

### Publication Issue :

Volume 8, Issue 6

November-December-2022

**Page Number :** 443-450

## Article History

Accepted: 20 Nov 2022

Published: 05 Dec 2022

## I. INTRODUCTION

Hate speech is a crime that has been on the rise in recent years, not just in face-to-face contacts but also online.

Social media is exploding in popularity, and its anonymity aspect fully fosters freedom of expression. Hate speech directed at an individual or group based on race, caste, religion, ethnic or national origin, sex, handicap, gender identity, or other factors is an abuse of this sovereignty.

It actively promotes violence or hate crimes and disrupts society by jeopardizing peace, credibility, and human rights, among other things.

Detecting hate speech in social media discourse is crucial, but it's a difficult undertaking. This study aims to address the quality of datasets, which is a major concern raised by many of the problems that have been brought to light. This paper also addresses the second issue, which is that the best characteristics for hate speech identification must be investigated and determined before developing a suitable classifier. For this reason, datasets tend to fall into one of these categories. The work is divided into two parts:

Hate speech tweets are categorized into 8 types.

- Abusive
- Comparison
- Passing judgment
- Religious
- Sarcasm
- Vulgar
- Spam
- Non spam

Tweet is classified into one of these types or as normal tweets.

## II. LITRATURE SURVEY

1. Muhammad sabih Et.al. presented "Un-Compromised Credibility: social media based MultiClass Hate Speech Classification for Text" In

this paper their work was that to identify the problem which is had speech towards a person or a group because of that it promotes violence or hate crimes and create and imbalance in society. Addressing different categories of head separately this paper aims to correctly predict their forms.

• Data set used: -

- a. Hate Based Twitte
- b. Hat Eval
- c. Waseem A
- d. Waseem B

In this study, major challenges are identified first and the complex problem of multi-class automated hate speech classification for text is accomplished with much better results. Ten separate binary classified datasets consisting of different hate speech categories are constructed. Each dataset was annotated by experts with the strong agreement of annotators under comprehensive, clear definition and well-defined rules. Datasets were well balanced and broad. They were also supplemented with language subtleties. Compilation of such dataset was achieved as necessary requirement for filling the gap of the field. After the development of high-quality datasets, a list of effective, commonly used and recommended features extracted from related studies under the field of text mining were identified. In addition to these features our own potential features were also proposed. These features were then explored and identified with respect to their problem objective. It is found that character 2 to 4-grams, word 1 to 5-grams, dependency tuples, sentiment scores, and count of 1st, and 2nd person pronouns were very effective. There are a total of ten separate datasets compiled with binary labels each. Different features together with a different set of models are explored over each dataset. Best features are identified and ten independent models are trained. Each tweet will be passed to all ten models and therefore it may have multiple hate classes identified by each model.

2. Idris et.al published "Detecting Hate Speech on Social Media Using Deep Learning Techniques". Her work shows that hate speech is a recurring problem on social media platforms, attacking specific groups of people based on certain common characteristics. Online data is created by users so quickly that it has become a daunting task to manually moderate a user's comments, including hate speech, in order to reduce the negative impact on the platform. In our previous research, we were able to create a model that can detect hate speech with high accuracy when detecting hate speech in user comments and posts (called tweets) on her Twitter, a social media platform. rice field. two base classifiers Long-short-term memory labels were used for naive-based SVM b. classification. H. "Hate Speech" and "No Hate Speech". David Sonnet. al. has shown in his work that this has led to the "hate speech" label, which includes forms of offensive language other than hate. B. Offensive language. The ensemble model was developed using a soft-voting combination technique with his two base classifiers, NBSVM (Naive Bayes Support Vector Machine) and his LSTM (Long Short-Term Memory). We created a data representation using Facebook's Fast Text and TF-IDF (Term Frequency Inverse Document Frequency) using character n-grams known for rare detection.

3. Raquel Fernandez et al. Published "Hate Speech Corpus Research for Detecting Hate Speech and Predicting Popularity". Her research showed that, as a result, her discussion environment online can become abusive, hateful, and toxic, especially when user anonymity is added. In order to identify, study, and ultimately contain this problem, such negative environments and the language used within them are studied under the name of hate speech. In this post, the last Focus on her two points. Consider a specific hate speech corpus, the Twitter corpus collected by Waseem and Hovy (2016). This corpus is gaining momentum as a resource for training

models to detect hate speech. Manually annotated to distinguish between two types of hate speech (sexist and racist), allowing for more nuanced insight and analysis. In addition, as a Twitter corpus, we offer all sorts of analysis and research possibilities for typical characteristics of Twitter corpora, such as: B. User and Tweet metadata, user interactions, etc.

4. Eileen Kwok et al published "Localizing Hate: Detecting Tweets Against Black People". Their study found that Twitter has a sizeable Black following, but anti-Black people We've found that racist tweets are particularly damaging to the Twitter community.<sup>1</sup> In doing so, they can provide data on the sources of hate speech against black people. Nov 2012 On January 1st, a Twitter user wrote, "So my tweet caused an 11-year-old black girl to commit suicide? It has been retweeted 77 times, has 17 favourites, and the user currently has 14,959 followers. They processed this balanced training data set of 24582 tweets by removing URLs, mentions, stop words, and punctuation. lowercase; and equate alternate spellings of insults with their properly spelled equivalents. They showed that our bag-of-words model was insufficient to classify anti-Black tweets accurately. Their research is becoming increasingly relevant, such as how often tweets are woven into various conversations.

5. Mohyaddin et al. published Automatic Hate Speech Detection: A Literature Review. Their study provided a comprehensive review of the different approaches to detect hate speech on social media platforms that have been deployed in recent years, along with a brief description of their analysis. Language on Twitter. Data was collected and preprocessed using the Twitter API. We then applied a nearby classification algorithm and got an accuracy of 93%. Thus, a dangerous statement can be observed as follows: Dangerous speech is offensive speech that encourages the audience to participate in acts of violence against a particular group of people. Therefore, the most common hate speech online is

related to religion, race, sexual orientation, nationality, class, and gender. b) Hate speech may contain one of the pillars of dangerous speech. c) Dangerous speech often incites listeners to support or commit acts of violence against a particular group. The six most common calls to action in dangerous language are kill, riot, beat, loot, kick out, and discriminate. The Internet is inherently open and dynamic, but communities have their own rules that define language boundaries. there is. These boundaries vary by culture and are shaped by historical events and cultural norms.

6. Resmi-Regnathan et al. With the announcement of Hate Speech Detection in Conventional Languages on SocialMedia Using Machine Learning, the need to automate the process of classifying hate speech data arises. They also use Malayalam for hate speech. For Malayalam, we basically developed Malayalam data for the system, the system detects hate speech based on the dataset applied to the English system, and uses SVM, logistic regression, and random forest machine learning algorithms. did. This methodology describes a proposed device set up to classify speech into two specific classes, specifically "hate speech, clean speech". Suggest a perfect learning method. Specifically, as this figure shows, the research methodology consists of six major steps: dataset acquisition, pre-processing, feature extraction, model training, evaluation run, and model checking.

7. Resmi-Regnathan et al. "Detection of customary language hate speech in social media using machine learning" was published. Therefore, it becomes necessary to automate the process of classifying hate speech data. They also use Malayalam for hate speech. For Malayalam, we basically developed Malayalam data for the system, the system detects hate speech based on the dataset applied to the

English system, and uses SVM, logistic regression, and random forest machine learning algorithms. did. Methodology describes a proposed deprecated device for classifying speech into two specific classes, specifically "hate speech clean speech". Suggest a perfect learning method. Specifically, as this figure shows, the research methodology consists of six major steps: dataset acquisition, pre-processing, feature extraction, model training, evaluation run, and model checking. .

In everyday life, with the increasing use of social media, it seems that everyone thinks they can speak and write as they please. Because of this thinking, hate speech is on the rise. Hate speech can hurt individuals and communities. However, manually identifying hate speech is very difficult. Therefore, it becomes necessary to automate the process of classifying hate speech data. To simplify the process of classifying hate speech, we used a machine learning approach to detect hate speech from speech. Most machine learning algorithms require data to be formatted in a very specific way, so usually some preparation is required to get useful insights from datasets. Stemming is the process of creating morphological variants of root/base words. Remove word prefixes or suffixes. Stemming programs are commonly called stemming algorithms or stemming algorithms. The stemming algorithm reduces the words "chocolates", "chocolatey" and "choco" to the stem "chocolate" and "retrieval", "retrieved" and "retrieves" to the stem "retrieve". Stemming is an important part of pipeline processing in natural language processing. A stemmer's input is a tokenized word.

Lemmatization The process of grouping different forms of a word so that they can be analysed as a single element. It usually refers to doing things right using vocabulary and morphological analysis of words. This is usually intended to remove only inflections and return the base or lexical form of the word, known as the lemma, which is the root word

rather than the stem, which is the output of word stemming.

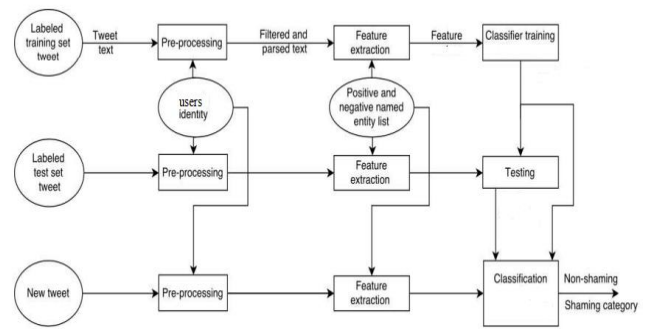
8. Aliji al-Hassan et al. "Detecting Hate Speech in Social Networks: A Multilingual Corpus Research" was published. They also distinguished various anti-social behaviors (cyberbullying, abuse, abusive language, hate speech). After Differentiation, they also published a comprehensive study on the use of text mining, examining several challenges for detecting Arabic hate speech. In this paper, they also present a table containing the different algorithms used to detect different hate categories and the accuracy of this model.

The table summarizes all the work discussed., arranged according to their chronological order. English anti-social behavior, English hate speech and finally Arabic anti-social behavior. They serve as a quick reference for the important work done in auto-discovering social media. All approaches and their respective test results are clearly listed. Consolidate all terms related to hate speech and related posts.

The Challenge of Hate Speech Detection in Arabic Hate speech detection is more than a simple keyword detection, it is a complex task with many challenges. Based on the review done in the previous section, they can identify some research challenges in automated detection of Arab hatred in social media. The first obstacle is that there is some research into hate speech detection. This can lead to high precision and high recall.

There is growing awareness of the problem of hate spreading through social networks in the Arab region and around the world.

### III. SYSTEM ARCHITECTURE



### IV. MODULES

• **Pre-processing:**

Data pre-processing is the process of preparing raw data and making it suitable for machine learning models. This is the first critical step in creating a machine learning model. When creating machine learning projects, you don't always come across clean and formatted data. Also, it is imperative to store the data in a clean and formatted way every time you work with it. For this, we use a data pre-processing task.

Pre-processing is a machine learning term that refers to transforming raw features into data that machine learning algorithms can understand and learn from.

• **Feature Extraction:**

Feature extraction aims to reduce the number of features in a data set by generating new features from existing features and then discarding the original features. This reduced new feature set should summarize most of the information contained in the original feature set. Therefore, a condensed version of the original function can be created from the combination of the original set.

• **Classifier training:**

In data science, a classifier is a type of machine learning algorithm used to assign class labels to data inputs. An example is an image-recognition classifier to label images (e.g. "car", "truck", or "person"). The classification algorithms are trained using labelled data in the image recognition example, e.g. a classifier that receives training data for labelling images. After adequate training, the classifier can



then take the unlabelled images as input and will generate classification labels for each image. Classification algorithms use sophisticated mathematical and statistical methods to generate predictions about the likelihood that a data entry will be classified in a certain way. In the image recognition example, the classifier statistically predicts whether an image is likely to be a car, truck, or person, or some other classifier that the classifier has been trained to identify.

• **Classification:**

Classification is defined as the process of recognizing, understanding and grouping objects and ideas into predefined categories also known as "populations". By using these pre-classified training datasets, classification in machine learning programs leverages a series of algorithms to classify future datasets into corresponding categories and related.

Classification algorithms used in machine learning use input training data for the purpose of predicting the probability or probability that the following data will fall into one of the predefined categories. One of the most popular classification applications is to filter emails into "spam" or "non-spam", as used by major email service providers today.

• **Testing:**

The process of training an ML model involves feeding an ML algorithm (i.e. learning algorithm) with training data to learn. The term ML model refers to the model artifact generated by the training process. The training data must contain the correct answer, called the target or target attribute. The training algorithm finds patterns in the training data that map the attributes of the input data to the target (the response you want to predict) and creates an ML model that captures those patterns.

## V. MOTIVATION

Today, social networking sites involve billions of users around the world.

- User interactions with these social sites, like Twitter, have a huge and sometimes unwanted impact on everyday life.
- Vandals disrupt meaningful discussions in online communities by posting irrelevant comments.
- Victims receive punishment disproportionate to the extent of the crime they clearly committed

## VI. OBJECTIVE OF THE SYSTEM

- Automatically reduce and categorize Hate Speech tweets.
- Helps block haters from attacking victims on social media.
- Provides insight into embarrassing events and shameful people.
- Attempts to improve classification accuracy using machine learning and real-time Twitter data.

## VII. SYSTEM REQUIREMENT

### Software Requirement

1. Operating System - Windows 7/8/10/11
2. Application Server - Apache Tomcat 7/8/9
3. Front End - HTML, CSS  
BOOTSTRAP
4. Scripts - JavaScript.
5. Language - Python
6. Database - My SQL
7. IDE - Eclipse / Visual Studio/Anaconda

## VIII. METHODOLOGY

In the proposed systems approach, we formulate a problem classifying task to identify and mitigate the side effects of public shame on networks. Two major contributions:

- 1) Classification and automatic classification of embarrassing tweets.
- 2) Develop a web application that allows Twitter users to identify Shamers.

The goal is to automatically classify tweets into 8 categories. For each category, the labeled training and test sets undergo preprocessing and feature extraction. The training set is used to train the random forest (RM). Tweets marked as negative by all classifiers are not considered shameful.

## IX. CONCLUSION

After identifying the primary challenges, the multi-class automated hate speech categorization for text problem is solved with significantly better results. Potential solution for countering the menace of online public shaming in Twitter by categorizing shaming comments in eight types, choosing appropriate features, and designing a set of classifiers to detect it.

The propagation of hate speech on social media has been increasing significantly in recent years and it is recognized that effective counter-measures rely on automated data mining techniques. Our work made several contributions to this problem.

First, we introduced a method for automatically classifying hate speech on Twitter using a machine learning that empirically improve classification accuracy.

## X. REFERENCES

- [1] Muhammad Sabih “Un-Compromised Credibility: Social Media Based Multi-Class Hate Speech Classification for Text” January 2021 IEEE Access 9:109465-109477 DOI:10.1109/ACCESS.2021.3101977 License CC BY-NC-ND 4.0
- [2] Dris, David, Ogunseye, Elizabeth Oluyemisi and Akinola, Solomon Olalekan. (2020). “Detecting Hate Speech on social media Using Deep Learning Techniques”, University of Ibadan Journal of Science and Logics in ICT Research (UIJSLICTR), Vol. 5 No. 1, pp. 22 - 38. ©U IJSLICTR Vol. 5, No. 1, June 2020. [3] Filip Klubicka, Raquel Fernandez “Examining a hate speech corpus for hate speech detection and popularity prediction” arXiv:1805.04661v1 [cs.CL] 12 May 2018. [4] Irene Kwok and Yuzhou Wang “Locate the Hate: Detecting Tweets against Blacks”
- [3] Thomas Davidson, Dana Warmley, Michael Macy, Ingmar Weber “Automated Hate Speech Detection and the Problem of Offensive Language” arXiv:1805.04661v1 [cs.CL] 12 May 2017
- [4] Mohiyaddeen, Dr. Shifaula Siddiqui “Automatic Hate Speech Detection: A Literature Review” e-ISSN: 2250-0758 | p-ISSN: 2394-6962 Volume-11, Issue-2 (April 2021)
- [5] Resmi Reghunathan, Asha A S “Hate Speech Detection in Conventional Language on Social Media by using Machine Learning” International Journal of Engineering Research & Technology) <http://www.ijert.org> ISSN: 2278-0181 Vol. 11 Issue 06, June-2022
- [6] Areej Al-Hassan, Hmood Al-Dossari “Detection of hate speech in social networks: a survey on multilingual corpus” Conference Paper February 2019 DOI: 10.5121/csit.2019.90208
- [7] Sindhu Abro , Sarang Shaikh , Zafar Ali Sajid Khan , Ghulam Mujtaba “Automatic Hate Speech Detection Using Machine Learning: A Comparative Study” (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 11, No. 8, 2020.
- [8] Pete Burnap and Matthew L. Williams “Cyber Hate Speech on Twitter: An Application of Machine Classification and Statistical Modelling for Policy and Decision Making .” 1944-2866 # 2015 The Authors. Policy & Internet published by Wiley Periodicals, Inc. on behalf of Policy Studies Organization.
- [9] Sreelakshmi ka , Premjith Ba , Soman K.Pa “Detection of Hate Speech Text in Hindi English Code-mixed Data” Procedia Computer Science 171 (2020) 737–744 .

- [10] Mathew, Binny, et al. "Analyzing the hate and counter speech accounts on twitter." arXiv preprint arXiv:1812.02712 (2018).
- [11] Gaydhani, Aditya, et al. "Detecting hate speech and offensive language on twitter using machine learning: An n-gram and tfidf based approach." arXiv preprint arXiv: 1809.08651 (2018).
- [12] Watanabe, Hajime, Mondher Bouazizi, and Tomoaki Ohtsuki. "Hate speech on twitter: A pragmatic approach to collect hateful and offensive expressions and perform hate speech detection." IEEE access 6 (2018): 13825-13835.
- [13] Wich, Maximilian, Jan Bauer, and Georg Groh. "Impact of politically biased data on hate speech classification." Proceedings of the Fourth Workshop on Online Abuse and Harms. 2020.

**Cite this article as :**

Dr. Nisha Auti, Shreeraj Ghadge, Rajdatta Jadhav, Prajwal Jagtap, Sumit Ranaware, "Social Media based Hate Speech Detection using Machine Learning", International Journal of Scientific Research in Computer Science, Engineering and Information Technology (IJSRCSEIT), ISSN : 2456-3307, Volume 8 Issue 6, pp. 443-450, November-December 2022.

Available at doi :

<https://doi.org/10.32628/CSEIT228653>

Journal URL : <https://ijsrcseit.com/CSEIT228653>