# Text Classification for Offensive Language Detection on Social Media

**[1]Mohd Abdul Aleem,[2]Mohammed Muneef,[3]Mohammed Wajihullah**

[1]Assistant Professor, Dept of CSE-AI&ML, Lords Institute of Engineering and Technology, Hyd.

[2,3]B.E Student, Dept of CSE-AI&ML, Lords Institute of Engineering and Technology, Hyd.

aleem1234@gmail.com, muneef.mohd19@gmail.com,mdwajihofficial@gmail.com

**Abstract**: Poor verbal exchange has affected social media content. One of the nice solutions to this trouble is to apply a computational approach to isolate the fantastic content material. Also, social media users belong to one of kind groups. In this examine, we present a text category model that includes a house responsibilities model and a tokenize, 3 aggregate strategies and eight class techniques. Our check shows that it is useful for detecting suspicious messages on our data obtained from Twitter. Considering the hyper parameter optimization, the three methods AdaBoost, SVM and MLP have the very best common F1 rating of the popular joint approach of TF-IDF.

**Keywords**: cyber bullying; adolescent safety; offensive languages; social media.

## I.    INTRODUCTION

Text class is the procedure of dividing facts into predefined instructions primarily based on their content. Text type is the objective function of herbal text for the primary path. Class is the primary requirement of text retrieval systems, which retrieve the text in response to the person's question, to extract a few files and records content material knowledge, which modifications the textual content in several methods, such as developing content material, answering questions. . , pick or delete files. Paper mining has end up one of the most famous regions of generation that has covered a lot of studies, mainly in pc technology, information retrieval (IR), and statistics mining. Natural language processing (NLP) strategies are used to extract understanding from actual data gathered by means of human customers. Text mining reads unstructured records to provide accurate content models inside the shortest possible time [1]. Nowadays, social networking websites are one of the most critical information era carriers due to the fact most of the people around the sector use these websites day by day to protect the whole thing else. Social networking web sites are growing new strategies to engage people in a big community [2]. Chat lets in customers to communicate with folks that proportion

morals and values. Websites are a useful way of discussion between individuals which ends up in mutual acquaintance and sharing of understanding. On social media, it's miles becoming common to not write a sentence with correct grammar and spelling. This exercise also can lead to uncertain records together with lexical, syntactic and semantic and due to this form of statistical uncertainty it is very tough to understand the exact reality. Therefore, extracting the hypothesis with the best statistics from the needless statistical records is crucial for occasion assessment [3]. Community assessment programs had been a hit in recent years due to the close to-user boom impact associated with all different aspects of the network. The dialogue is designed as a image and the fact of the relationship is used inside the form of a large stream that can be leveraged for diverse functions. Analysis of social media texts from guides overlaying the essential length of social community analysis. This processing is quantitative, assisted by way of vital researchers on this context which permits a big choice of facts in social community mining which includes social community strategies, algorithms for seek models and content material analysis in social networks [4]. With the upward thrust of social media, humans accumulate chances

and specific patterns of information are up to date 24/7. Social networks include forums and blogs where humans can effortlessly connect them to each other. Social media is particularly described as "a miles cheaper and extra virtual tool that takes care of all bodily aspects and while having access to data, participating and sharing collectively, or growing a courting." Many researches inside the area have tried to better understand the big volumes of unfastened, customer-generated content. Research areas for e-commerce, clever transportation, smart town, cybercrime and many extra. There are no exceptions. However, extracting beneficial, actionable insights from patron-generated content is tough. As every social media provider has its personal requirements and bounds in relation to facts series. Most of the evolved reviews are used for automated processing and exploration [5]. As a end result, messages on social media are frequently quick, informal, with masses of abbreviations, jargon and slang finishing in unstructured words. In line with the above regarding the usefulness of social media structures, social media has come to be the main a part of our adventure today.

## II. REVIEW OF LITERATURE

**Offensive language detection**

The analysis of cyber bullying, aggression, hate speech, toxic speech and negative feedback in social media has long received attention from the research community. There is a lot of public information available to show the classifiers for these activities. However, there is no general information or training course that can be combined to get a more powerful system. Kumar et al. (2018) provide information and results of joint work on attack detection. The data provided includes 15,000 Face book comments and comments in English and Hindi. The goal is to distinguish between our expressions: not aggressive, covertly aggressive and overly aggressive. The chemistry group's comments have been criticized on Kaggle. Several methods have been evaluated for this test on a database consisting of users with Wikipedia comments. These comments are divided into 6 types: toxic, poisonous, obscene, random, insulting, and racist. Regarding the identification of hate speech, Davidson et al. (2017) presented the latest hate speech data with more than 24,000 English-language tweets belonging to three categories: non-offensive, hate speech, and hate speech. Mandl et al. (2019) pointed out the shared responsibility regarding the violence of speech where our data is taken from Twitter and Facebook and made available

for Hindi, German and English. In addition, Zampieri et al., 2019, Zampieri et al., 2020 presented a number of results of the search for ambiguous language in different languages received from the SemEval competition group.

**Multilingual text category**

Multilingual text classification is a phenomenon in text type. However, very little work has been done in this area. At first, Lee et al. (2006) proposed a method for categorizing multilingual textual content using latent semantic indexing techniques. This method offers a multi-lingual monolingual presentation of English and Chinese datasets. In all other tables, Prajapati et al. (2009) added a method based on translating the data into familiar words, then creating categories. They include information on using WordNET to map sentences to patterns and then classify text, using the Rocchio linear classifier and probabilistic Naïve Bayes and K-Nearest Neighbour (KNN). Amini et al. (2010) studied MTC by combining semi-supervised acquisition techniques, including joint and consensus-based self-training. They trained native speakers on Reuters Corpus Volume 1 and a pair of (RCV1/RCV2) containing five languages: English, German, French, Italian and Spanish. The authors analyzed

their methods using six strategies: Boost, co-regularized boosting, boosting with self-learning, Support Vector Machine (SVM) with self-education, co-regularization + self-education, and boosting with self-complete education. Training. Bentaallah and Malki (2014) compare WordNet-based methods for categorizing multilingual text. First rely on the machine translator to enter the WordNet immediately and use a non-confrontational method to remember the meaning of most of the content to be efficient. While the second does not include translation and search WordNet related to all languages. Mittal and Dhyani (2015) discussed multi-language classroom content based on N-gram techniques. They study MTC in Spanish, Italian and English. They are done by predicting the words of the data and using Naïve Bayes in the classification section. Recently, Kapila and Satvika (2016) solved the MTC problem of Hindi and English using special tools for learning algorithms, including SVM, KNN, decision tree, self-map, and genetic algorithms. They improve the accuracy of the approach by using various selection techniques.

Recently, deep neural networks and context-aware embeddings have been proposed in the field of textual content for English (Liu and Guo, 2019 and many others) .

In short, despite the wide range of work on various language topics, MTC is almost ignored and few studies.

## III. RESEARCH METHODOLOGY

In this study, we aim for a modular data delivery pipeline including a modular maintenance level and tokenize three integration methods, and eight classifiers. The experiment performed in this study is based on Twitter and the data has been carefully edited. Although we do not guarantee that our framework will be effective on all social networks, it can provide future research to guide researchers and businesses. The broader implications of this article may relate to investigations of online crime on social media platforms. Additionally, due to the individual characteristics of social media platforms, it is impossible to generalize the model to all platforms. For example, it shows that training a classifier on Reddit is more difficult than Gab due to the average post length.

This section briefly describes the steps followed to clean and prepare the data as well as to perform the experiments. Additionally, Figure 1 shows a graphical

representation of these steps, discussed below.

## A) Data Preparation

Data Preparation is the first step for training bi nary classifiers. The strategies for data preparation, which need to be carefully conducted, are described as below:

• **Basic cleaning methods**: We need to clean the data as (i) extracting the pure text from the dataset, removing duplicates, and NaNs (ii) transforming to lowercase (iii) expanding the abbreviations.

• **Slangs:** Given the micro-blogging style of Twitter, using slangs are typical. Slangs bring difficulties to text mining approaches, especially for those emerging lately and thus do not have an updated entry in any dictionaries. So, we plan to transform the text into a canonical form using the reference dictionary1 for slangs and abbreviations.

• **Removing methods**: Using hashtags, user references, links, and emojis are typical on social media platforms. Therefore, preprocessing the data and selectively removing the typical pat terns are essential to normalize the text.

• **TF-IDF:** One way to represent words into vectors is to count the occurrence of words seen in the whole documents. One caveat of this method is the overemphasizing the frequent words in the dataset. In contrast with the word counting method, TF-IDF distributes the weight of frequent words by their relative frequency.

• **Word2Vec:** The word2vec method takes a corpus of text as input and returns word vectors as output. There are two model architectures to produce a distributed representation of words. The continuous bag-of-words (CBOW) architecture predicts the current word based on the context (window size), and the Skip-gram predicts surrounding words (defined window) given the current word.

• **FastText:** FastText represents a low dimensional vector text that is generated by summing vectors corresponding to the words in the text. Neural Network is being used in FastText for word embedding. FastText model is often compared to other deep learning classifiers with a higher speed and accuracy for training and evaluation.
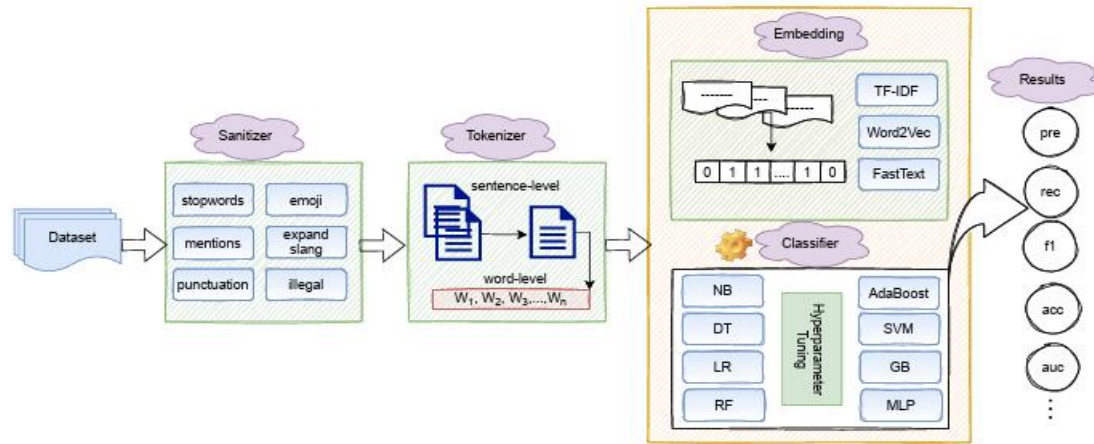
Fig.1 The modular experimental setting with the flow of data from dataset to results.

## B) Using Text Mining Techniques to Detect Online Offensive Contents

Offensive language identity in social media is a difficult project due to the fact the textual contents in such surroundings is regularly unstructured, informal, or even misspelled. While protective strategies accompanied with the useful resource of modern-day social media are not enough, researchers have studied smart methods to pick out offensive contents the use of text mining approach. Implementing text mining strategies to investigate on-line records requires the subsequent tiers: 1) facts acquisition and pre-process, 2) characteristic extraction, and three) classification. The essential worrying conditions of using text mining to detect offensive contents lie on the function selection phrase that allows you to be elaborated in the following sections.

## C) Message-degree Feature Extraction

Most offensive content material detection research extracts types of functions: lexical and syntactic capabilities. Lexical features cope with each phrase and phrase as an entity. Word patterns together with look of certain key phrases and their frequencies are frequently used to represent the language version. Early studies used Bag-of-Words (BoW) in offensiveness detection. The BoW method treats a text as an unordered series of phrases and disregards the syntactic and semantic facts. However, the use of BoW technique on my own no longer simplest yields low accuracy in diffused offensive language detection, but additionally brings in a immoderate fake fine fee particularly throughout heated arguments, shielding reactions to others' offensive posts, or even

conversations among close friends. N-gram method is considered as a complicated approach in that it brings words nearby context statistics into interest to come across offensive contents. N-grams constitute subsequences of N non-stop phrases in texts. Bi-gram and Tri-gram are the most famous N grams utilized in textual content mining. However, N-gram suffers from trouble in exploring related phrases separated by way of lengthy distances in texts. Simply developing N can alleviate the hassle but will gradual down machine processing pace and bring in extra fake positives. Syntactic features: Although lexical functions perform nicely in detecting offensive entities, without thinking about the syntactical structure of the complete sentence, they fail to distinguish sentences' offensiveness which comprise same terms but in super orders. Therefore, to remember syntactical capabilities in sentences, natural language parsers are brought to parse sentences on grammatical structures in advance than function selection. Equipping with a parser can help keep away from deciding on un-associated phrase units as skills in offensiveness detection

## D) User-level Offensiveness Detection

Most current research on detecting on-line offensive languages great attention on sentence-stage and message-stage constructs. Since no detection method is one hundred% accurate, if customers maintain connecting with the assets of offensive contents (e.g., on-line users or web sites), they're at excessive chance of constantly publicity to offensive contents. However, client-stage detection is a greater hard venture and research related to the character level of assessment is in large part lacking. There are a few confined efforts at the individual degree.

## E) Machine mastering algorithms

NaiveBayes (NB) and SVM—are used to perform the class, and 10-fold cross validation changed into carried out on this check. To fully compare the effectiveness of clients' sentence offensiveness price (LSF), style features, shape skills and content material particular functions for purchaser offensiveness estimation, we fed them sequentially into the classifiers, and get the bring about Fig.Three. The "Strong Weak" technique absolutely makes use of offensive words as the base feature to find out offensive customer. Similarly, "LSF" manner the sentence offensiveness price generated by means of LSF is used as the base characteristic.

## IV. EXPERIMENTAL RESULTS

This section describes the exclusive experiments that we done to assess the LSF for detecting offensive messages in social networks. Data Description The take a look at facts, taken from YouTube comments at the discussion board, is a spread of advertising feedback in response to the pinnacle 18 movies. Music video distribution includes thirteen classes: track, vehicles, comedy, schooling, enjoyment, films, games, fashion, information, nonprofits, animals, science and sports. Each remark segment incorporates the patron's personal records, time and content of the content. Private customers perceive the writer who posted the comment, records of time the comment became transformed right into a message and the content material statistics content consists of the person's recommendation. The database includes remarks from 2,175,474 tremendous customers.

Pre-processing Before passing the dataset to the classifier, pre-processing automatically gathers the words for each person and breaks them down into sentences. For each sentence in the sample dataset, computerized spelling and spelling correction precedes the introduction of the pattern dataset for the classifier. Using the WordNet corpus and spelling correction algorithm2, accurate spelling and grammar mistakes in incomplete sentences using

obligations which includes removing returned textual content in sentences, deleting characters unnecessary, department of lengthy words, alternative of textual content. And make modifications. Wrong letters and lacking letters inside the message. Therefore, phrases without letters, which include "spelling", are corrected to "spelling"; incorrect sentences, which include "accurate", are changed with "yes".

Ͽ� Test Locations in Sentence Crime The take a look at compares 6 sentence prediction techniques: a) Bag-of-sentences (BoW): The BoW technique ignores textual content writing patterns and word order and verifies complaints through checking whether or not or not they include all of the customer's information. Use and offensive words. This process is likewise based totally on benchmarking. B) 2-grams: The N-gram method reveals the unsatisfied sentences using the willpower of every part of n sentences in a sentence and examines whether or not or not the sentences consist of all identified and terrible sentences. In this technique, N is equal to two; it also acts as a degree. C) 3 grams: N-gram method, determines each a part of 3 words in a given sentence. It additionally follows the standard. D) five-grams: N-grams system, decide every a part of five words in a given sentence. It also follows the same old.

## V. Evaluation Metrics

In our test, the benchmarks for category in hypothesis checking out (i.e., precision, recollect, and f-rating) are used to evaluate the performance of LSF. In particular, the fact affords the share of guidelines which can be surely terrible phrases. Returns the general accuracy of the category, which represents the share of proper assaults recognized. The fake advantageous (FP) rate represents the percentage of tips that are not authentic false positives. The fake terrible (FN) fee represents the proportion of really terrible words that are not identified. The F-score represents the weighted average of the actual and the reverse, that means:

$$f-score = \frac{2(precision \times recall)}{precision + recall}$$
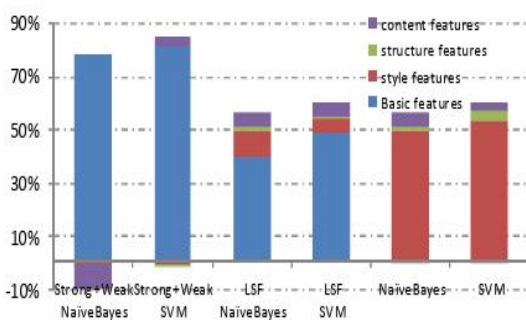
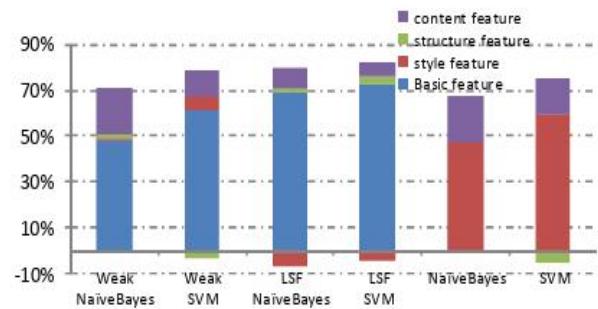

Fig.2 F-score for different feature sets using NB and SVM



Fig.3 F-score for different feature sets using NB and SVM (without strongly offensive words)

## V. CONCLUSION

In this analysis, we examine the current content development process in investigating suspicious content for the protection of young people's online safety. In these images we show the content of the video in social media news, especially on Twitter. Our plan is to promote modular development allowing clean use of the combination of specific course contents. This document is most important if it provides new modular delivery pipeline content for easy evaluation by performing in-depth performance process analysis. Quality, performance and materials are highlighted using the new logo.

## REFERENCES

1. P. Hajibabaee, F. Pourkamali-Anaraki, and M. Hariri Ardebili, "An empirical evaluation of the t-sne algorithm for data visualization in structural

engineering," in 2021 IEEE International Conference on Machine Learning and Applications. IEEE, 2021.

2. S. Zad, M. Heidari, J. H. J. Jones, and O. Uzuner, "Emotion detection of textual data: An interdisciplinary survey," in 2021 IEEE World AI IoT Congress (AIIoT), 2021, pp. 0255–0261.

3. S. Zad, M. Heidari, J. H. Jones, and O. Uzuner, "A survey on concept-level sentiment analysis techniques of textual data," in 2021 IEEE World AI IoT Congress (AIIoT), 2021, pp. 0285–0291.

4. M. Heidari, S. Zad, B. Berlin, and S. Rafatirad, "Ontology creation model based on attention mechanism for a specific business domain," in 2021 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), 2021, pp. 1–5.

5. M. Heidari, S. Zad, and S. Rafatirad, "Ensemble of super vised and unsupervised learning models to predict a profitable business decision," in 2021 IEEE International IOT, Electronics and Mechatronics Conference (IEMTRONICS), 2021, pp.

6. A. Esmaeilzadeh, M. Heidari, R. Abdolazimi, P. Ha jibabaee, and M. Malekzadeh, "Efficient large scale nlp feature engineering with apache spark," in 2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC). IEEE, 2022.

7. R. Abdolazimi, M. Heidari, A. Esmaeilzadeh, and H. Naderi, "Mapreducepreprocess of big graphs for rapid connected components detection," in 2022 IEEE 12th Annual Computing and Communication Workshop and Conference (CCWC). IEEE, 2022.

8. M. Malekzadeh, P. Hajibabaee, M. Heidari, and B. Berlin, "Review of deep learning methods for automated sleep staging," in 2022 IEEE 12th Annual Computing and Com munication Workshop and Conference (CCWC). IEEE, 2022.

9. A. Razavi, D. Inkpen, S. Uritsky, and S. Matwin, "Offensive language detection using multi-level classification," Advances in Artificial Intelligence, vol. 6085/2010, pp. 16-27, 2010.

10. Mahmud, Ahmed, KaziZubair, and Khan, Mumit "Detecting flames and insults in text," in Proc. of 6th International Conference on Natural Language Processing (ICON' 08), 2008.

11. D. Yin, Z. Xue, L. Hong, and B. Davison, "Detection of harassment on Web 2.0," in the Content Analysis in the Web 2.0 Workshop, 2009.

12. Z. Xu and S. Zhu, "Filtering offensive language in online communities using grammatical relations," in Proceedings of The Seventh Annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference (CEAS'10), 2010

13. Prasadu Peddi and Dr. Akash Saxena (2014), "EXPLORING THE IMPACT OF DATA MINING AND MACHINE LEARNING ON STUDENT PERFORMANCE", International Journal of Emerging Technologies and Innovative Research (www.jetir.org), ISSN:2349-5162, Vol.1, Issue 6, page no.314-318, November-2014, Available: http://www.jetir.org/papers/JETIR1701B47.pdf

14. Prasadu Peddi and Dr. Akash Saxena (2015), "The Adoption of a Big Data and Extensive Multi-Labled Gradient Boosting System for Student Activity Analysis", International Journal of All Research Education and Scientific Methods (IJARESM), ISSN: 2455-6211, Volume 3, Issue 7, pp:68-73.