

DETECTING HATE SPEECH ON SOCIAL MEDIA

A. Faritha Banu^{*1}, *M. Priyadharshini*²

ABSTRACT

The Information exchange and the increasing use of social media have brought great benefits to humanity. However, this also causes some issues, such as the dissemination and exchange of hate speech messages. Therefore, to solve the issue, subsequent studies have employed a variety of feature engineering techniques and machine learning algorithms to automatically identify hate speech posts across a variety of datasets that is growing on social media sites. To our knowledge, no research has been done to contrast various feature engineering methodologies and machine learning algorithms in order to ascertain which feature engineering strategy and algorithm performs best on a widespread dataset that is publicly accessible. On a publicly available dataset with three different classes, this study examines the performance of three feature engineering strategies and eight machine learning algorithms. The experimental results show that the support vector machine approach, with an overall accuracy rate of 79%, works best when paired with bigram features. Our research has real-world applications and can serve as a benchmark in the field of automatically identifying hate speech texts. Additionally, the results of various comparisons will be used to compare upcoming studies for current automated text classification approaches using state-of-the-art methods.

Keywords: social media, text classification, NLP

I. INTRODUCTION

Hate speech is a type of communication that encourages prejudice, animosity, or violence towards certain people or groups in light of their racial, ethnic, religious, sexual, gender, or other identities or characteristics. It entails the use of offensive words, phrases, slurs, or pejorative language

with the goal to denigrate, dehumanize, or stir up hatred towards certain people or groups.

Hate speech has the ability to exacerbate social differences, reinforce stereotypes, and undercut the values of inclusivity, equality, and respect. By encouraging a climate of fear, exclusion, and marginalization, it can seriously hurt society and psychologically. Beyond the immediate targets, it has an impact on broader social beliefs and may have real-world repercussions like discrimination, harassment, or even violent attacks.

For the sake of upholding the fundamentals of human rights, encouraging social cohesion, and advancing a more inclusive and tolerant society, hate speech must be identified and dealt with. Legal frameworks, education, public awareness campaigns, and the creation of tools and technologies for detection and mitigation are just a few of the measures used in the fight against hate speech.

By comprehending the causes and effects of hate speech, society can fight to create a more secure and respected atmosphere both online and offline, where people can express themselves without fear of damage or prejudice.

II. HATE SPEECH

Hate speech is a type of communication, it may be either verbal, written or symbolic, which insults, threatens the individuals or groups based on the characteristics like race, ethnicity, religion, sexual orientation, gender identity, disability and other required characteristics. This usually involves discriminatory or offensive language or expressions needed to demean, dehumanize the hatred, violence, that is a discrimination against individuals or batch target.

Hate speech often promotes prejudice, stereotypes, and hostility towards individuals or communities, exacerbating

¹Department of Computer Technology, Karpagam Academy of Higher Education, Coimbatore.

²Department of CSE, Nalla Malla Reddy Engineering College, Hyderabad.

* Corresponding Author

social divisions and creating an environment of fear and animosity. It can occur in various forms, such as online messages, speeches, public demonstrations, or even in private conversations.

Hate speech weakens equality and leads to prejudice towards particular groups of people. Women and immigrants are frequently the major targets. Due to the refugee crisis and political changes in recent decades, there has been a sharp rise in anti-immigrant sentiment [1]. Current initiatives by some governments and authorities to address this issue include identifying and monitoring hate speech directed towards immigrants. Abuse, slander, and discrimination against women in social and professional contexts are frequent manifestations of hatred of the female gender, a prevalent and long-standing kind of discrimination. family and affiliation.

2.1. TYPES OF HATE SPEECH

Hate speech can manifest in various forms, targeting different aspects of an individual's identity [2]. Here are some common types of hate speech:

Racial Hate Speech: This involves derogatory comments, slurs, or insults based on a person's race or ethnicity. It promotes discrimination and fosters racial animosity.

Religious Hate Speech: It refers to offensive or demeaning remarks targeting individuals or groups based on their religious beliefs. This form of hate speech can fuel religious intolerance and perpetuate stereotypes.

Homophobic and Transphobic Hate Speech: Hate speech targets people based on their gender identity or sexual orientation. This comprises phrases that are disparaging, offensive, or dangerous against LGBT+ (lesbian, gay, bisexual, transgender, and queer) individuals.

Sexist or Misogynistic Hate Speech: This type of hate speech degrades and belittles individuals based on their gender, particularly targeting women. It may involve

objectification, sexist slurs, or the promotion of gender-based violence.

Ableist Hate Speech: Hate speech targeting individuals with disabilities or impairments [3]. It involves derogatory language, mockery, or the marginalization of people with disabilities.

Xenophobic Hate Speech: Hate speech that exhibits prejudice, discrimination, or hostility towards individuals based on their nationality, immigration status, or foreign background. It can contribute to social divisions and promote exclusion.

Online Harassment and Cyberbullying: Hate speech in the form of targeted harassment, threats or insults directed at specific individuals or groups. This can happen on social media platforms, forums or other online spaces.

It's crucial to recognize that hate speech can overlap across these categories and take on multiple forms simultaneously. Additionally, hate speech can evolve and adapt over time, adopting new language or tactics to spread discriminatory ideologies.

III. HATE SPEECH DETECTION

Deep learning is widely used in hate speech detection because of its main feature, to effectively learn patterns and representations from large amounts of text data. Some of the ways how deep learning is commonly used in detecting hate speech.

Word Embeddings: Word embeddings, such Word2Vec or GloVe, are used in deep learning models for hate speech detection to display the words in a continuous vector space. These embeddings identify the semantic connections between words and aid the model's comprehension of contextual data.

Recurrent Neural Networks (RNNs): RNNs are frequently employed to identify hate speech and include

variations such Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU). RNNs process sequential information in text by considering the order of words and capturing dependencies over time.

Convolutional Neural Networks (CNNs): CNNs are excellent for extracting well-known patterns and textual properties. CNNs aid in the extraction of significant traits from sentences and phrases that are then utilized for categorization in the identification of hate speech.

Transformer-based Models: Transformer models with improved processing performance include BERT (Bidirectional Encoder Representation from Transformer) and GPT (Generative Pre-Trained Transformer). that includes the identification of hate speech in natural language. These models have state-of-the-art outcomes across a range of applications and contain attention methods to pick up contextual information.

Transfer Learning: Hate speech identification can be honed using some Deep Learning samples that have been pre-trained on complex language tasks like sentiment analysis or language modeling. Transfer learning enables models to leverage pre-existing knowledge, which can improve performance even with limited labeled hate speech data.

Model Ensemble: The effectiveness and reliability of hate speech detection systems will rise with the use of ensemble approaches, such as the fusion of several deep learning models and the integration of other machine learning algorithms. By aggregating predictions from multiple models, the overall performance can be improved.

The success of deep learning models is important for hate speech detection and it relies on high-quality, diverse, and well-labeled training data. Ongoing model evaluation, fine-tuning, and adapting to emerging language patterns are also crucial to address the emerging nature of speech on hate. The definition and frameworks surrounding speech on hate may vary between countries and jurisdictions. The goal

of laws and regulations pertaining to hate speech is to achieve a compromise between safeguarding the right to free speech and avoiding the negative effects of promoting prejudice, hatred, and violence.

IV. HATE SPEECH ON SOCIAL MEDIA

Speech on hate in social media platforms can take various forms and may include:

Racist Comments: Slurs are insults directed against specific people or groups based on their race or ethnicity.

Homophobic or Transphobic Remarks: Insults, harassment, or threats directed towards individuals based on their sexual orientation or gender identity.

Religious Discrimination: Offensive statements or hate speech aimed at a particular religion or religious group.

Misogyny and Sexism: Sexually explicit or demeaning language towards women, promoting stereotypes, or advocating violence against them.

Ableism: Discrimination or derogatory remarks targeting individuals with disabilities or impairments.

Xenophobic Comments: Prejudice, intolerance, or hostility towards individuals based on their nationality or immigration status.

Cyberbullying: Targeted harassment, insults, or threats towards specific individuals, often based on their personal characteristics or perceived vulnerabilities.

It's crucial to keep in mind that social media platforms often have policies in place to fight hate speech. These policies may include measures like removing or flagging objectionable content, suspending or deleting user accounts, or offering reporting tools for users to report instances of hate speech.

4.1. HATE SPEECH DETECTION IN DEEP LEARNING

Detection of Hate speech in deep learning involves using machine learning algorithms to automatically identify and classify text or speech which may be hate speech or non-hate speech. The general overview of the process involves:

Collection of Data: A numerous amount of dataset of labeled examples containing both speech on hate and speech on non-hate instances is gathered. The deep learning model's training data is comprised of these instances.

Preprocessing: The text data is preprocessed to remove noise, punctuation, and special characters. It may also involve tokenization, stemming, or lemmatization to standardize the text.

Word Embeddings: Word embeddings are used to represent words in a numerical form. Popular techniques include Word2Vec, GloVe, or FastText. By capturing the semantic connections between words, these embeddings help the deep learning model comprehend the context.

Model Architecture: The detection of hate speech frequently makes use of deep learning models like recurrent neural networks (RNNs), convolutional neural networks (CNNs), or transformer-based models (like BERT and GPT). To create predictions, these models analyze correlations and patterns in text data.

Training: The training and validation sets are separated from the labeled dataset. The validation set is used to track the performance of the deep learning model and avoid overfitting after it has been trained using the training set. Iterative adjustments are made to the model's parameters using methods like gradient descent and backpropagation.

Evaluation: After training, a separate test data set is used to evaluate the model's performance using measures including accuracy, precision, recall, and F1 score. This assessment aids in determining how well the model recognizes hate speech.

Deployment: The trained model can be used to examine new, unexplored text data and categorize it as hate speech or non-hate speech after demonstrating adequate performance.

It is important to note that the success of hate speech detection models heavily relies upon quality and diversity of the training data, as well as ongoing monitoring and refinement to keep up with evolving language and different forms that arise in hate speech.

4.2. NLPIN DETECTION OF HATE SPEECH

Natural Language Processing (NLP) plays a vital role in hate speech detection with the help of enabling the analysis along with understanding of text data. Here are some key components of NLP applied to hate speech detection:

Text Preprocessing: The text data is preprocessed using NLP techniques before modeling. This involves handling special characters and punctuation, deleting lowercase letters, stop words (common words like "the", "and", etc.), and tokenization (breaking text into separate words or tokens).

Feature Extraction: NLP makes it possible to extract useful features from text data. Word frequency and relevance can be represented in a text document using methods like Bag-of-Words (BoW) and Term Frequency-Inverse Document Frequency (TF-IDF). Machine learning or deep learning models use these features as input.

Language Modeling: NLP models such as word embeddings (Word2Vec, GloVe, etc.) and context embeddings (BERT, ELMO, etc.) capture the semantic meaning of words and phrases in a language. These models represent words as dense vectors, capturing relationships and contextual information, which helps in understanding the meaning and nuances of hate speech.

Sentiment Analysis: Sentiment analysis techniques are used to identify the sentiment or tone of a given text, which can be useful in detecting hate speech [4,5]. Hate speech is often characterized by negative sentiment or derogatory language.

Sentiment analysis models can help identify the underlying sentiment of a text, aiding in the identification of potentially offensive content.

Text Classification: Text is categorized into groups like hate speech and non-hate speech using NLP-based text classification models that have been trained. The detection of hate speech frequently makes use of supervised machine learning techniques like Naive Bayes, support vector machines (SVMs), or deep learning models like CNNs and LSTMs. These models learn from labeled data to identify patterns and make predictions on unseen text.

Named Entity Recognition (NER): NER is the process of identifying and classifying named entities in text, such as: Name, location, organization, or racial, ethnic, or religious group. Identifying such entities can help detect hate speech, as hate speech mainly targets specific groups or individuals based on its features.

Dependency Parsing: In order to understand the links between words, dependency parsing examines the grammatical structure of sentences. It focuses on detecting complex syntactic patterns or identifying specific linguistic constructs that may be indicative of hate speech.

These NLP techniques are often combined and integrated into comprehensive hate speech detection systems, utilizing machine learning or deep learning models to achieve accurate and effective identification of offensive and harmful content.

4.3. MACHINE LEARNING IN HATE SPEECH DETECTION

Machine learning algorithms are broadly utilized in hate speech detection to routinely classify textual content as both hate speech or non-hate speech. Here are the important things concerned in the usage of machine learning for hate speech detection:

Data Collection: A diverse and representative dataset of labeled examples containing both the hate speech and non-

hate speech instances are collected. Dataset should cover various forms of hate speech, languages, and demographics.

Data Preprocessing: The text data is preprocessed to remove noise, punctuation, and special characters. You can also include tokenization, stemming, or lemmatization to standardize the text and reduce its dimensionality.

Feature Extraction: The Features are extracted from the preprocessed text to represent the input data for the machine learning model. The Bag-of-Words (BoW), Term Frequency-Inverse Document Frequency (TF-IDF), and word embeddings like Word2Vec and GloVe are popular approaches.

Model Selection: The detection of hate speech is possible using a number of machine learning methods, such as Naive Bayes, Support Vector Machines (SVM), Random Forests, and Neural Networks. The properties of the data set, the available computing resources, and the needed performance all influence the model choice.

Training and Evaluation: A training set and a test set are created from the labeled data set. Utilizing the retrieved features and their accompanying labels, a machine learning model is trained on the training set. The performance of the model is then assessed on the test set using measures like precision, recall, and F1 score.

Model Tuning and Validation: To boost performance, your model's hyperparameters must be tuned. You can use methods like grid search and cross-validation to identify the ideal hyperparameter values.

Deployment and Monitoring: Once a satisfactory model is obtained, it can be deployed to analyze new, unseen text data. Continuous monitoring is essential to assess the model's performance over time and adapt to evolving language patterns and other new forms of hate speech.

Note that the subjective character of hate speech and the continuously changing language employed present issues for hate speech identification programs. Regular updates,

feedback loops, and ongoing model improvements are crucial to keep the model effective and up-to-date with emerging forms of hate speech.

4.4. DETECTION OF HATE SPEECH IN SOCIAL MEDIA: TROUBLES

Detection of hate speech in social media poses several challenges due to the dynamic and complex nature of online communication. Some of the key challenges include:

Contextual Understanding: Hate speech can heavily rely on contextual cues, sarcasm, irony, or cultural references, making it challenging to accurately interpret the intended meaning. Algorithms may struggle to capture these nuances, leading to false positives or false negatives.

Evolving Language and New Trends: Language is constantly evolving, and hate speech adapts to new trends and expressions. Hate speech detection systems must continually update their models to understand and recognize emerging forms of hate speech and avoid becoming obsolete. **Variations in Language and Dialects:** Social media platforms have a global user base, resulting in diverse languages, dialects, and cultural contexts. Developing accurate hate speech detection models that can handle multiple languages and dialects is a significant challenge due to linguistic variations and limited labeled data for each language.

Subtle and Implicit Hate Speech: Hate speech can manifest in subtle or implicit ways, making it harder to detect. It may involve dog-whistling, coded language, or veiled references, which can be challenging for automated systems to identify accurately.

User Anonymity and User Generated Content: social media allows users to create anonymous accounts or use pseudonyms, making it difficult to track and attribute hate speech to specific individuals [5,6,7]. Additionally, the sheer volume of user-generated content makes it challenging to analyze and monitor hate speech effectively.

Legal and Cultural Differences: Different countries and regions have different legal frameworks and cultural norms regarding hate speech [8]. Developing a universal hate speech detection system that accommodates these differences while adhering to ethical and legal boundaries is a complex task.

Balancing Free Speech and Censorship: Distinguishing between hate speech and freedom of expression is a delicate balance [9,10]. Hate speech detection systems need to be careful not to overly restrict legitimate speech, emphasizing the importance of striking the right balance to avoid censorship concerns.

Addressing these challenges requires ongoing research, collaboration with diverse communities, continuous model improvement, robust data collection, and adapting to the evolving dynamics of online discourse.

S.No	NAME	ACTIVE USERS (Approx.)
1	FACEBOOK	2.9 billion
2	LINKEDIN	0.93 billion
3	TWITTER	0.45 billion
4	INSTAGRAM	2 billion

Table 1. Active users in social media

ACTIVE USERS in Billions (Approx.)

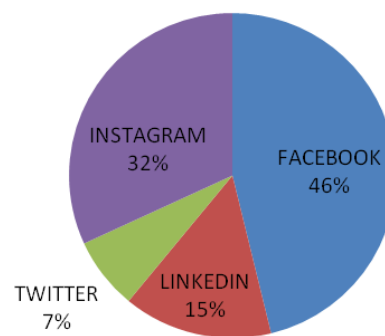


Figure 1. Number of users in social media

V. CONCLUSION

In conclusion, identifying hate speech in social media is a difficult undertaking that involves solving a number of problems. Accurately detecting hate speech requires a thorough awareness of the context, adaptation to changing linguistic trends, proficiency in a variety of languages and dialects, recognition of nuanced types of hate speech, consideration of user anonymity, and adherence to legal and cultural norms.

Machine learning and NLP techniques play a vital role in developing hate speech detection systems. Deep learning models like neural networks and transformers, along with NLP components like text preprocessing, feature extraction, and language modeling, provide the tools to analyze and classify text data effectively.

It's crucial to recognize that the discipline of detecting hate speech is active and developing. Increased accuracy, fairness, and ethical implications of hate speech detection systems depend on ongoing study, collaboration, and advances in data collecting, model building, and monitoring procedures. Finding the right balance between protecting freedom of expression and preventing the negative effects of hate speech remains a challenge, and effectively addressing these issues will require continued dialogue with diverse communities and the need for engagement is highlighted.

REFERENCES

1. C. Vania, M. Ibrahim, and M. Adriani, "Sentiment Lexicon Generation for an Under-Resourced Language," *Int. J. Comput.*, vol. 5, no. 1, pp. 59-72, 2014.
2. A. Brown, "What is hate speech? Part 1: The Myth of Hate," *Law Philos.*, vol. 36, no. 4, pp. 419-468, 2017.
3. A. Goswami and A. Kumar, "A survey of event detection techniques in online social networks," *Soc. Netw. Anal. Min.*, vol. 6, no. 1, pp. 1-25, 2016.
4. A. Assiri, A. Emam, and H. Al-Dossari, "Towards enhancement of a lexicon-based approach for Saudi dialect sentiment analysis," *J. Inf. Sci.*, vol. 44, no. 2, pp. 184-202, 2018.
5. C.-F. Tsai, "Bag-of-Words Representation in Image Annotation: A Review," *ISRN Artif. Intell.*, vol. 2012, pp. 1-19, 2012.
6. P. Burnap and M. L. Williams, "Us and them: identifying cyber hate on Twitter across multiple protected characteristics," *EPJ Data Sci.*, vol. 5, no. 1, 2016.
7. G. Xiang, B. Fan, L. Wang, J. Hong, and C. Rose, "Detecting offensive tweets via topical feature discovery over a large scale twitter corpus," *Proc. 21st ACM Int. Conf. Inf. Knowl. Manag. - CIKM'12*, p. 1980, 2012.
8. J. Lilleberg, Y. Zhu, and Y. Zhang, "Support vector machines and Word2vec for text classification with semantic features," *Proc. 2015 IEEE 14th Int. Conf. Cogn. Informatics Cogn. Comput. ICCI*CC 2015*, pp. 136-140, 2015.
9. P. Burnap and M. L. Williams, "Cyber Hate Speech on Twitter: An Application of Machine Classification and Statistical Modeling for Policy and Decision Making," *Policy & Internet*, vol. 7, no. 2, pp. 223-242, 2015.
10. T. Hua, F. Chen, L. Zhao, C.-T. Lu, and N. Ramakrishnan, "STED: semi-supervised targeted-interest event detection in twitter," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 1466-1469.