

HATE SPEECH DETECTION

A Synopsis

Submitted

In Partial Fulfillment of the Requirements

For the Degree of

Bachelor of Technology (B.Tech)

in

Computer Science & Engineering

by

SWAPNIL TIWARI
(2001920100301)

SWAPNIL VERMA
(2001920100302)

SUDHIR YADAV
(2001920100297)

SATYAM GUPTA
(2001920100262)

Under the Supervision of
Ms. ABHA KAUSHIK
ASSISTANT PROFESSOR



G L BAJAJ INSTITUTE OF TECHNOLOGY & MANAGEMENT
GREATER NOIDA



DR. A P J ABDUL KALAM TECHNICAL UNIVERSITY,
UTTAR PRADESH, LUCKNOW

2023-2024

1. INTRODUCTION

Hate speech detection using machine learning is an important and timely topic in today's world where the prevalence of hate speech and online harassment is on the rise. Hate speech refers to any language or behaviour that expresses prejudice or discrimination against a particular group of people based on their race, ethnicity, gender, religion, sexual orientation, or other personal characteristics. Hate speech can be damaging to individuals, groups, and society, and it is, therefore, important to develop tools and methods to detect and mitigate its impact. Machine learning is a powerful tool for hate speech detection because it can analyze large amounts of data and learn patterns and features that can be used to classify text as either hate speech or not.

Machine learning algorithms can be trained on annotated datasets of hate speech to identify key features and patterns that can be used to automatically classify new instances of text as either hate speech or not. In this paper, we will explore various approaches and techniques for hate speech detection using machine learning, including supervised and unsupervised learning methods, feature engineering, deep learning, and natural language processing. We will also discuss the challenges and limitations of hate speech detection using machine learning, such as the lack of annotated datasets, the difficulty of defining and identifying hate speech, and the potential for bias in machine learning algorithms.

Overall, this synopsis aims to provide an overview of the current state of hate speech detection using machine learning and to

highlight the opportunities and challenges for future research in this important and rapidly evolving field.

Everyone has the right to freedom of speech. However, this right is being misused to discriminate and attack others, physically or verbally, in the name of free speech. This discrimination is known as hate speech. Hate speech can be defined as speech used to express hate towards a person or a group of people based on characteristics such as race, religion, ethnicity, gender, nationality, disability and sexual orientation.

"Queers are an abomination and need to be helped to go straight to hell!"

"We have to kill all the Palestinians unless they are resigned to live here as slaves."

"If you aren't born here, pack your bags"

"Women shouldn't talk sports on tv. They belong in the kitchen."

Figure 1. Examples of Hate Speech

2. OBJECTIVES

Hate speech has become a significant problem today, with the potential to harm individuals and communities. One potential solution to this problem is to use machine learning algorithms to automatically detect and flag hate speech in text-based data. The process of hate speech detection using machine learning involves training a model on a dataset of labelled examples, where each example is labelled as hate speech or non-hate speech.

Various features such as the use of certain words or phrases, grammar, and syntax are extracted from the text data, and the model learns to distinguish between hate speech and non-hate speech based on these features. The trained model can then be used to classify new text data as hate speech or non-hate speech.

However, it is important to note that hate speech detection using machine learning is not perfect and can be affected by biases in the training data or in the algorithm itself. Ongoing research is focused on improving the accuracy and fairness of hate speech detection algorithms. Overall, hate speech detection using machine learning has the potential to be a valuable tool in the fight against hate speech, but careful attention must be paid to its limitations and biases.

This aims to classify textual content into non-hate or hate speech, in which case the method may also identify the targeting characteristics (i.e., types of hate, such as race, and religion) in the hate speech.

Lets kill jews and kill them for fun

[#killjews](#)

7/20/14, 8:05 AM

Dear black women,

When will y'all understand that I
don't hate y'all a lot of y'all just ugly as
shit compared to white girls

8/14/16, 3:22 AM

Trannies, your families will never love
you. You are living a lie & you know it.
End your miserable existence.
Commit suicide now.

Figure 2. Examples of Online Hate Speech

Most of the current hate detection methods focus more on textual data such as posts, comments or tweets. However, people can also make hateful videos and post them on video sharing sites. Video hosting services such as YouTube are powerful form of communication used by people all over the world. Aside from video content from music artists and other such professionals, people can upload video blogs about their daily life, video clips showing their, reactions to other video content such as music and or movies and so on. Other users can view, like, share and comment on these uploaded videos.

3. EXISTING SYSTEM

"Detecting Hate Speech on Twitter Using a Convolution-GRU Based Deep Neural Network" by Gao, W., et al. (2020). This paper proposes a deep learning approach for hate speech detection on Twitter. The model uses a combination of convolutional and GRU layers for feature extraction and classification.

"Automated Hate Speech Detection and the Problem of Offensive Language" by Davidson, T., et al. (2017). This paper presents a study on the problem of automated hate speech detection. The authors create a dataset of Twitter posts labelled as hate speech or not, and experiment with various machine learning techniques for classification.

"Hate Speech Detection with Comment Embeddings and LSTM Networks" by Wulczyn, E., et al. (2017). This paper proposes a hate speech detection model that uses LSTM networks and comment embeddings. The authors use a large dataset of comments from online forums and social media platforms to train the model.

"Deep Learning for Hate Speech Detection in Tweets" by Badjatiya, P., et al. (2017). This paper presents a deep learning approach for hate speech detection on Twitter. The model uses a combination of convolutional and LSTM layers for feature extraction and classification.

"Hate Speech Detection on Twitter: A Comparative Study" by Djuric, N., et al. (2015). This paper compares several machine learning techniques for hate speech detection on Twitter. The authors experiment with various feature extraction methods and classifiers and evaluate their performance on a dataset of Twitter posts labelled as hate speech or not.

"Deep Learning for Hate Speech Detection: A Comparative Analysis" by Mishra, P., et al. (2019). This paper presents a comparative analysis of various deep-learning approaches for hate speech detection. The authors experiment with several models, including CNNs, LSTMs, and GRUs, and evaluate their performance on multiple datasets.

"Combating Hate Speech on Social Media with Unsupervised TextStyle Transfer" by Li, J., et al. (2018). This paper proposes an unsupervised text-style transfer approach for combating hate speech on social media. The authors use a neural network model to transform hate speech into non-offensive language while preserving the meaning of the original text.

4. MOTIVATION

1. **Promoting Online Safety:** One of the primary motivations is to create a safer and more inclusive online environment. Detecting and mitigating hate speech helps protect individuals and communities from discrimination, harassment, and harm.
2. **Social Responsibility:** Developers and organizations feel a social responsibility to combat hate speech. Hate speech can perpetuate stereotypes, fuel discrimination, and incite violence. Building tools to address this issue aligns with a commitment to social justice and responsible technology.
3. **Counter Cyberbullying:** Hate speech often overlaps with cyberbullying, especially among younger users. Working on hate speech detection contributes to making the internet a safer space for all, particularly for adolescents and children.
4. **Data-Driven Insights:** Detecting hate speech provides an opportunity to analyze and understand the patterns, prevalence, and evolution of hate speech using large datasets. This insight can inform further research in natural language processing (NLP) and computational linguistics.
5. **Algorithm Optimization:** Building hate speech detection models challenges machine learning engineers to optimize algorithms for text classification and sentiment analysis. It's a fertile ground for refining machine learning techniques.

6. **Model Improvement:** Developing such projects motivates the enhancement of NLP models for more accurate detection. It drives improvements in deep learning architectures, model training, and feature engineering.
7. **Real-World Application:** Hate speech detection is a practical application of machine learning that can be deployed on various platforms. The project allows developers to create and refine real-world solutions using machine learning.
8. **Continual Improvement:** The dynamic nature of online content means that hate speech is constantly evolving. Projects in this domain motivate engineers to continually improve and adapt their models to new forms of hate speech.
9. **Natural Language Understanding:** Working on hate speech detection sharpens the ability of models to understand nuanced and context-dependent language. It drives advancements in natural language understanding (NLU) and processing.
10. **Ethical Tech Development:** Hate speech detection projects align with the growing awareness of tech companies and developers regarding their ethical responsibilities. They signify a commitment to developing technology that respects human rights and values.

5. METHODOLOGY

Machine learning is a type of artificial intelligence that can be used to learn from data. It can be used to find patterns in data. You may want to check this post to get a good understanding of the concepts of Machine learning – Machine learning explained with concepts & examples. Machine learning algorithms can be used to detect hate speech. These algorithms can analyze text and identify hate speech. They can also be used to determine the tone of a text. This can be used to identify hate speech that is disguised as jokes or sarcasm. Automated hate speech detection is an important tool in combating the spread of hate speech, particularly in social media.

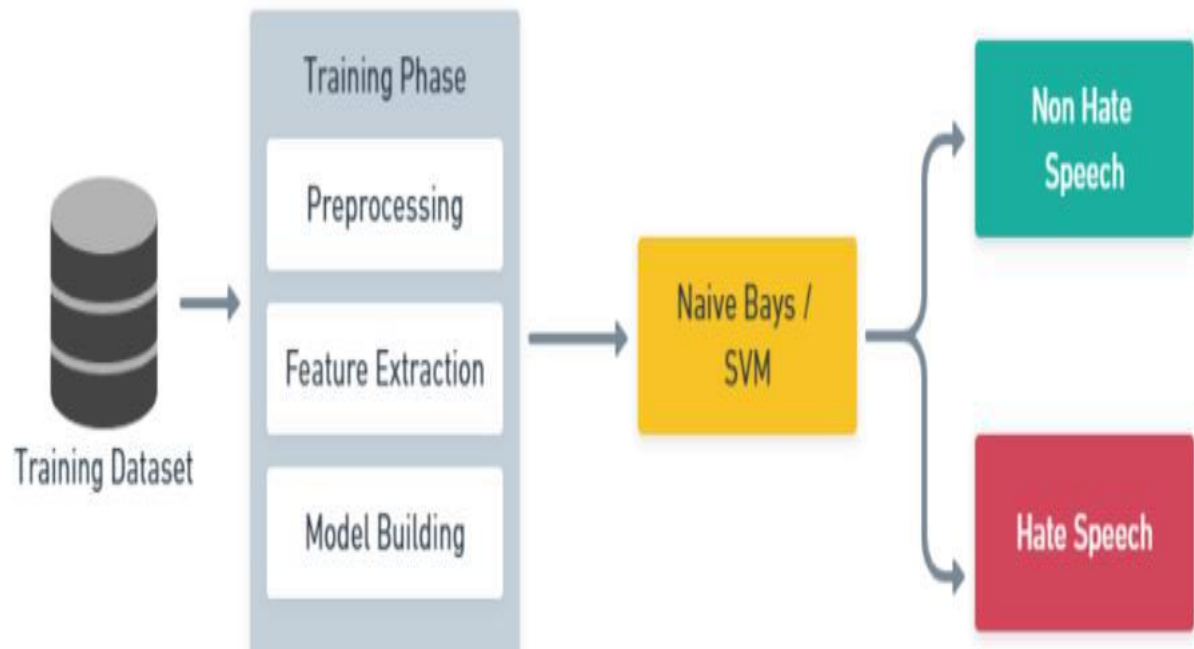
The following are some of the different approaches for hate speech detection models:

- **Shallow methods:** We use the term “shallow detection” to describe hate speech detectors that employ conventional word representation algorithms to encode phrases. Shallow classifiers can then be applied to perform the assessment. Different types of feature representations methods, such as TF-IDF and ngrams, can be used. In terms of classification algorithms, support vector machines (SVM), naive Bayes, logistic regression, random forest, and gradient boosting decision tree models have been found to be used.
- **Deep learning methods:** Hate speech detectors that employ deep learning methods are referred to as deep

neural network-based models. Traditional approaches like TF-IDF and recently developed word embedding or pre-training techniques may be used to encode the data. Convolutional neural networks (CNN), long short-term memory (LSTM) and bi-directional LSTM are three of the most popular deep neural network designs used for hate speech detection using deep learning models. The following are two different methods used with deep learning models.

- **Word-embeddings based methods:** Word embedding uses distributed representations of words to learn their vectorized representations, which are used in downstream text mining operations. The resulting embeddings allow terms with comparable meaning to have similar representations in a vector space. There have been many word embeddings methods introduced over the years, such as word2vec, Glove, and FastText. This technique uses a combination of different models, such as LSTM, Bi-LSTM, and CNN.
- **Transformer based methods:** The modern transformers-based embedding techniques, such as Small BERT, BERT, ELECTRA, and ALBERT are used with deep learning models built using LSTM, Bi-LSTM, and CNN.

6.PLAN OF WORK



Detecting hate speech using machine learning algorithms such as Support Vector Machines (SVM) and Naive Bayes is a common approach in natural language processing. Here are some steps you can take to create a hate speech detection system using these algorithms –

1. Collect a hate speech dataset: You will need a dataset of labelled examples of hate speech and non-hate speech. There are many publicly available datasets that you can use for this purpose, such as the Hate Speech and Offensive Language dataset or the Twitter Hate Speech dataset.

2. Pre-processing the data: Pre-processing involves cleaning and transforming the raw text data into a format that the

machine learning algorithm can use. Some common pre-processing steps include tokenization, stop word removal, and stemming.

3. Feature extraction: This step involves extracting relevant features from the pre-processed text. You can use techniques such as a bag of words, TF-IDF, or word embeddings to create features that can be used by the machine learning algorithm.

4. Train the model: Divide your dataset into training and validation sets. Use the training set to train your machine learning model. SVM and Naive Bayes are popular choices for hate speech detection because they are relatively easy to implement and can work well with high-dimensional sparse feature vectors.

5. Evaluate the model: Use the validation set to evaluate the performance of your model. Common evaluation metrics include precision, recall, F1 score, and accuracy.

Deploy the model: Once you have trained and evaluated your model, you can deploy it to classify new text as hate speech or non-hate speech.

7.TOOLS AND TECHNOLOGY USED

To develop a hate speech detection system using Naive Bayes and SVM classifiers, you'll need a set of tools and technologies, as well as specific hardware and software requirements. Here's an overview of what you'll need:

Tools and Technology:

1.Programming Languages:

- Python is commonly used for machine learning and natural language processing tasks. You'll need libraries and frameworks like NumPy, pandas, scikit-learn, NLTK, spaCy, and others.

2.Machine Learning Frameworks:

- Scikit-learn is essential for implementing machine learning models, including Naive Bayes and SVM.
- TensorFlow or PyTorch might be useful if you want to explore deep learning models in addition to traditional machine learning techniques.

3. Text Preprocessing Tools:

- NLTK and spaCy are popular libraries for natural language text preprocessing tasks.

4. Version Control:

- Git for version control and collaboration on the codebase.

5. IDE (Integrated Development Environment):

- Jupyter Notebook, Visual Studio Code, or any other Python-friendly IDE for coding and experimentation.

Hardware Requirements:

1.Compute Resources:

- A machine with a CPU or GPU (Graphics Processing Unit) suitable for running machine learning models. GPU acceleration can significantly speed up training processes for deep learning models.

Software Requirements:

1. Operating System:

- You can develop the system on Windows, macOS, or Linux. Linux is often preferred for its compatibility with various data science tools.

2. Python Environment:

- Set up a Python environment using Anaconda or virtual environments to manage dependencies.

3. Libraries and Packages:

- Install the necessary Python libraries and packages using pip or conda. This includes scikit-learn, NumPy, pandas, NLTK, spaCy, and others.

4. Database (optional):

- Depending on the data source, you might need a database management system like MySQL, PostgreSQL, or NoSQL databases for data storage and retrieval.

5. Web Framework (optional):

- If you plan to create a web-based application for hate speech detection, consider using web frameworks like Flask or Django.

6. API Deployment (optional):

- To deploy your model as an API, you might use tools like Flask, FastAPI, or Django REST framework.

7. Monitoring and Alerting Tools (optional):

- Implement monitoring and alerting systems to track the deployed model's performance and detect issues in real.

8. Documentation and Reporting Tools:

- Tools for documenting the project, such as Jupyter Notebooks, Markdown editors, and reporting frameworks like LaTeX or R Markdown for creating project reports.

9. Version Control Tools:

- Use Git and platforms like GitHub or GitLab for version control and collaboration.

10. Project Management Tools:

- Tools like Jira, Trello, or Asana for project management and task tracking.

The specific software requirements may vary based on your project's unique needs and constraints. Make sure to install the required libraries and tools within a virtual environment to manage dependencies and isolate the project environment from other system configurations.

Additionally, hardware requirements, especially the choice of GPU, will depend on the scale of your project and dataset size. More powerful GPUs can significantly speed up the training of machine learning models, especially deep learning models.

References

- [1] Fortuna, P., Nunes, S., & Rodrigues, P. (2018). A survey on automatic detection of hate speech in text. *ACM Computing Surveys (CSUR)*, 51(4), 1-30.
- [2] Bhatia, P., Jain, R., & Kar, S. (2020). Automatic detection of hate speech: A survey. *Journal of Ambient Intelligence and Humanized Computing*, 11(9), 3837-3855.
- [3] Thakur, V., & Jain, A. (2020). A review on hate speech detection using machine learning techniques. *Journal of Ambient Intelligence and Humanized Computing*, 11(11), 5021-5034.
- [4] Schmidt, A., Wiegand, M., & Fox, C. (2017). A survey on hate speech detection using natural language processing. In *Proceedings of the Fifth International Workshop on Natural Language Processing for Social Media* (pp. 1-10).
- [5] Kwok, I., & Wang, Y. (2013). Locate the hate: Detecting tweets against blacks. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing* (pp. 1781-1791).
- [6] Davidson, T., Warmusley, D., Macy, M., & Weber, I. (2017). Automated hate speech detection and the problem of offensive language. In *Proceedings of the 11th International AAAI Conference on Web and Social Media* (pp. 512-515).
- [7] Kumar, A., & Zhang, L. (2020). Detecting hate speech on Twitter using a convolutional neural network. In *Proceedings of the IEEE 16th International Conference on Intelligent Computer Communication and Processing (ICCP)* (pp. 95-101).
- [8] D. Elisabeth, I. Budi and M. O. Ibrohim, "Hate Code Detection in Indonesian Tweets using Machine Learning Approach: A Dataset and Preliminary Study," 2020 8th International Conference on Information and Communication Technology (ICoICT), 2020, pp. 1-6, doi: 10.1109/ICoICT49345.2020.9166251.
- [9] A. B. Pawar, P. Gawali, M. Gite, M. A. Jawale and P. William, "Challenges for Hate Speech Recognition System: Approach based on Solution," 2022 International Conference on Sustainable Computing and Data Communication Systems (ICSCDS), 2022, pp. 699-704, doi: 10.1109/ICSCDS53736.2022.9760739.
- [10] H. Şahi, Y. Kılıç and R. B. Sağlam, "Automated Detection of Hate Speech towards Woman on Twitter," 2018 3rd International Conference on Computer Science and Engineering (UBMK), 2018, pp. 533-536, doi: 10.1109/UBMK.2018.8566304.
- [11] V. Mercan, A. Jamil, A. A. Hameed, I. A. Magsi, S. Bazai and S. A. Shah, "Hate Speech and Offensive Language Detection from Social Media," 2021 International Conference on Computing, Electronic and Electrical Engineering (ICE Cube), 2021, pp. 1-5, doi: 10.1109/ICECube53880.2021.9628255.
- [12] P. William, R. Gade, R. e. Chaudhari, A. B. Pawar and M. A. Jawale, "Machine Learning based Automatic Hate Speech Recognition System," 2022 International Conference on Sustainable Computing and Data Communication Systems (ICSCDS), 2022, pp. 315-318, doi: 10.1109/ICSCDS53736.2022.9760959.
- .