

# Applying AI And Semantic Analysis to Identify and Address Harmful Online Comments

Agalya M<sup>1</sup>, Ayisha Parveen A<sup>2</sup>, Kaviya V<sup>3</sup>, Keerthana J<sup>4</sup>, And Niroshini R<sup>5</sup>

<sup>1</sup>Asst. Prof. Computer Science and Engineering, Vivekanandha College of Technology for Women, Namakkal, Tamil Nadu 613 006, India

<sup>2,3,4,5</sup> Computer Science and Engineering, Vivekanandha College of Technology for Women, Namakkal, Tamil Nadu 613 006, India

**Abstract**—This project focuses on developing a robust machine learning model for toxic comment classification, aiming to enhance user experience and safety in online environments. The primary objective is to create a model that accurately classifies comments as toxic or non-toxic based on their content, thereby facilitating effective moderation and reducing the spread of harmful language. The preprocessing stage includes text normalization, tokenization, and the removal of stopwords and irrelevant characters to prepare the data for model training. Additionally, we implement techniques to address class imbalance, ensuring that the model is not biased toward the majority class. For model development, we explore a range of machine learning algorithms, including logistic regression, support vector machines (SVM), and deep learning models such as convolutional neural networks (CNN) and recurrent neural networks (RNN).

**Keywords**— Toxic comment classification, Machine learning model, Text normalization, Tokenization, Stopwords removal, Class imbalance, Logistic regression, Support vector machines (SVM), Deep learning models, Convolutional neural networks (CNN), Recurrent neural networks (RNN), Performance metrics.

## I. INTRODUCTION

Machine learning (ML) has become the cornerstone of toxic comment classification, offering an automated solution that can efficiently and accurately identify toxic language in large datasets. Traditional moderation methods, such as keyword filtering, were insufficient in capturing the complexities of toxic speech, which often involves slang, context, and cultural nuances. ML models, on the other hand, have the capacity to learn patterns from data and adapt to different forms of toxicity. Through supervised learning, algorithms are trained on labeled datasets, where comments are categorized as toxic or non-toxic. These models, after being

trained, are capable of predicting the toxicity level of new, unseen comments.

Various algorithms, including logistic regression, support vector machines (SVM), and deep learning techniques such as recurrent neural networks (RNN) and transformers like BERT, have proven effective in this task. Toxic comments—defined as those that contain offensive, inflammatory, or harmful language—are a pervasive issue on the internet today. This toxicity manifests in various forms, such as cyberbullying, racism, sexism, hate speech, personal attacks, and trolling.

Toxic comments, if left unchecked, can drive users away from platforms and contribute to a hostile environment that stifles meaningful conversation. Research has shown that exposure to toxic language can result in psychological harm, particularly for younger audiences and marginalized groups.

Machine learning enables computers to learn patterns from data, allowing them to make predictions or classifications on new, unseen data. In the context of toxic comment classification, machine learning models are trained on large datasets of user comments labeled as either toxic or non-toxic.

Machine learning models can struggle to grasp these nuances, leading to false positives (non-toxic comments flagged as toxic) or false negatives (toxic comments left unflagged). Another challenge is the inherent imbalance in toxic comment datasets. In most online communities, the majority of comments are non-toxic, while only a small fraction contain harmful language.

Explainable AI (XAI) techniques are designed to provide insights into the model's decision-making process. For example, SHAP (SHapley Additive exPlanations) values and LIME (Local Interpretable Model-agnostic Explanations) are commonly used methods to determine which words or features contributed most to a toxic classification. Toxic comment classification is a critical project for

improving the safety and inclusivity of online platforms. With ongoing advancements in machine learning and NLP, we are making significant progress toward creating automated systems that can effectively moderate harmful content.

## II. LITERATURE SURVEY

In[1] Jacob Devlin et al., (ICIEA, 2018). "Toxic Comment Classification With Bidirectional Encoder Representations From Transformers (Bert) " Toxic comment classification is a critical task in the domain of natural language processing (NLP), especially as online platforms grapple with moderating abusive or harmful content. Various machine learning and deep learning algorithms have been applied to this problem, with each demonstrating differing degrees of effectiveness when applied to the Jigsaw Toxic Comment Classification Challenge dataset. Among the most prominent and successful approaches is the use of Bidirectional Encoder Representations from Transformers (BERT), a deep learning model introduced by Jacob Devlin et al.

In[2] Yoon Kim "A Convolutional Neural Network For Toxic Comment Classification" (ICCE, 2022). In contrast to BERT, earlier deep learning architectures like Convolutional Neural Networks (CNNs) have also been applied to the task, although with slightly lower effectiveness. Yoon Kim's implementation of a CNN for toxic comment classification yielded an accuracy of 94.0%. While CNNs are typically associated with image processing, they have proven effective in text classification by capturing local features and patterns such as n-grams, which are helpful for detecting specific toxic phrases or combinations of words. The CNN model's performance, though commendable, reflects its limitations in modeling long-range dependencies and complex contextual nuances, areas where transformer-based models like BERT excel. Still, CNNs maintain advantages in terms of computational efficiency and are often easier to train compared to more complex models.

In[3] A. P. O'Neill et al. (IJERT, 2019). " Deep Learning For Toxic Comment Classification Using

Long Short-Term Memory Networks ". Another significant approach is the use of Long Short-Term Memory (LSTM) networks, a type of recurrent neural network capable of capturing temporal dependencies in sequential data. A. P. O'Neill and collaborators applied an LSTM-based model to the toxic comment classification task, resulting in an accuracy of 92.4%. LSTMs are adept at maintaining long-term dependencies in input sequences, making them suitable for processing entire comment histories and identifying toxicity trends over time.

## III. EXISTING SYSTEM

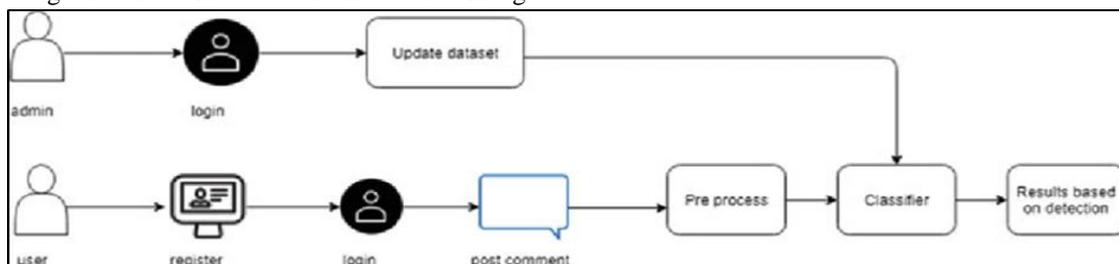
The existing system for identifying and addressing harmful online comments using AI tools and semantic analysis faces challenges in contextual understanding, detecting sarcasm and irony, and keeping pace with evolving online language. Existing systems may also perpetuate biases, struggle with scalability, and rely too heavily on keyword detection. Furthermore, they may lack nuanced understanding, adequate training data, and sufficient human oversight. To improve, it's essential to address these limitations and develop more sophisticated AI models that can effectively identify and mitigate online harm.

## IV. PROPOSED METHODOLOGY

### Architecture

An architecture diagram is a visual representation of the structure and design of a system or process. It provides a high-level view of the components, their relationships, and interactions within the system.

In various fields such as software engineering, network design, and enterprise architecture, architecture diagrams serve as essential tools for communicating complex concepts and ensuring that all stakeholders have a shared understanding of the system's structure. These diagrams typically include elements such as components, modules, interfaces, and data flows, and are crucial for both planning and implementation phases of a project.



### Data Collection

In the realm of toxic comment classification, gathering relevant and diverse data is critical for building a reliable machine learning model. The data primarily comes from online platforms where user-generated content is prevalent, such as social media sites, forums, comment sections on news websites, and public discussion platforms.

Platforms like Twitter, Reddit, and YouTube are common sources of toxic content as they host large communities where users engage in conversations, debates, and discussions. Some platforms have shared publicly available datasets for academic and research purposes. For instance, Kaggle has hosted a well-known Jigsaw Toxic Comment Classification Challenge, providing a dataset with labeled comments that are categorized into various toxicity classes.

Another valuable source includes Reddit, where moderators often label and remove toxic comments, making it easier to access labeled data by scraping toxic comments. The first step in preprocessing is text cleaning, which involves removing irrelevant elements from the comments, such as URLs, email addresses, HTML tags, and special characters. These elements do not contribute to the meaning of the text and may introduce noise into the model if left untreated. Another important step is case normalization, where all text is converted to lowercase to ensure uniformity, as "Toxic" and "toxic" would be treated as the same word by the model.

### Data Preprocessing

Data preprocessing also involves reducing noise, removing duplicates, and performing feature engineering to create new, more informative variables. Overall, the goal of data preprocessing is to ensure that the data is clean, consistent, and in a format that enhances the performance and accuracy of the machine learning model. Data preprocessing is a pivotal step in toxic comment classification, involving various techniques to transform raw text data into a format suitable for analysis and model training.

The first crucial phase is text cleaning and normalization, which ensures that the dataset is devoid of unnecessary noise and inconsistencies that could interfere with the model's performance. Text

cleaning involves several operations, beginning with the removal of irrelevant elements that do not contribute to the comment's meaning.

### EDA Techniques

#### Descriptive Analysis

Descriptive analysis is a critical first step in understanding the characteristics of the dataset used in toxic comments classification. It involves summarizing and interpreting the features of the dataset to gain insights into its structure, distribution, and patterns. The context of toxic comments classification, descriptive analysis helps in identifying key aspects of the data, such as the frequency and distribution of toxic comments, the characteristics of the comments, and potential imbalances in the dataset.

#### Null Value Handling

In toxic comments classification, null values (or missing values) can present significant challenges, impacting the quality and accuracy of machine learning models. Null values in text data typically refer to instances where information is absent or incomplete, which can occur in various forms: missing text fields, empty comments, or placeholder values that indicate lack of data.

The first step in handling null values is to identify and quantify their occurrence in the dataset. This involves examining the dataset to detect missing or incomplete values and understanding their distribution across different features and records. Identification of null values can be performed using exploratory data analysis (EDA) techniques. For text data, null values might be represented as empty strings, NaN (Not a Number), or NULL entries in the dataset.

#### TF-IDF Vectorization

TF-IDF (Term Frequency-Inverse Document Frequency) is a fundamental technique in text mining and natural language processing (NLP) used for vectorizing text data. It transforms text into a numerical representation that can be used as input for machine learning algorithms.

The goal of TF-IDF is to convert raw text into a structured format that captures the importance of terms within documents and across a corpus. This technique is widely employed in various

applications, including document classification, information retrieval, and text clustering. TF-IDF operates on the principle that certain words are more informative than others. Term Frequency (TF) measures how often a term appears in a document, reflecting its significance within that particular document

## V. MODEL BUILDING

Model building in machine learning is a multifaceted process that involves several key stages, each critical to the development of an effective model. The first step in model building is to clearly define the problem you are trying to solve. This involves understanding the specific goals of the machine learning task, such as classification, regression, clustering, or reinforcement learning. Techniques for feature selection include statistical tests, recursive feature elimination, and using algorithms such as LASSO (Least Absolute Shrinkage and Selection Operator) that perform feature selection as part of the modeling process.

### Train Test Split

In a toxic comments classification project, the train-test split is a crucial step in developing a machine learning model that accurately classifies comments as either toxic or non-toxic. The primary purpose of this split is to ensure that the model is evaluated on data that it has not seen during the training phase, which helps assess its generalization capability.

When building a classification model, it's essential to train the model on a subset of the data (the training set) and then evaluate its performance on a separate subset (the test set). This separation helps prevent over fitting, a situation where the model performs exceptionally well on the training data but poorly on new, unseen data. Over fitting occurs because the model has learned not only the underlying patterns but also the noise and specifics of the training data, which do not generalize well to new data

## VI. MODEL TRAINING

Model training in machine learning refers to the process of teaching a machine learning model to recognize patterns, relationships, or associations within a given dataset. This is accomplished by feeding the model a set of input data and corresponding outputs, which allows the model to

adjust its internal parameters to make predictions or classifications on unseen data.

The primary goal of model training is to minimize the error or difference between the model's predictions and the actual output, optimizing it to generalize well on new, unseen data. Model training in a toxic comments classification project using machine learning is a critical and detailed process that involves several stages, each requiring careful attention to ensure the final model performs effectively in real-world applications.

## VII. CONCLUSION

Toxic comment classification using machine learning is a critical endeavor aimed at identifying and moderating harmful or offensive content across various online platforms. This task is essential for maintaining a safe and respectful online environment, where users can engage in discussions without encountering abuse, harassment, or hate speech. Machine learning models, particularly those leveraging advanced techniques such as deep learning and natural language processing (NLP), play a pivotal role in automating this process. By training models on large datasets of labeled comments, these systems can learn to differentiate between toxic and non-toxic content with increasing accuracy.

## VIII. FUTURE ENHANCEMENT

Improved Contextual Understanding: Future enhancements in toxic comment classification will benefit greatly from models with deeper contextual understanding. Current models often struggle with nuances such as sarcasm, irony, and context-dependent meanings. Advanced transformer-based models like GPT-4 and BERT have made strides in this area, but there is room for improvement.

Multilingual and Cross-Cultural Models: As the internet is a global platform, enhancing models to handle multiple languages and cultural contexts is crucial. Current systems may be effective in one language but struggle with others, especially in handling slang, idioms, or culturally specific references.

## REFERENCE

- [1] Alissa de Bruijn, Vesa Muhonen, Tommaso

- Albinonistraat, Wan Fokkink, Peter Bloem, and Business Analytics. Detecting offensive language using transfer learning. 2019.
- [2] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages 5754–5764, 2019.
- [3] Thomas A Birkland. *An introduction to the policy process: Theories, concepts, and models of public policy making*. Routledge, 2019.
- [4] Estela Saquete, David Tomas, Paloma Moreda, Patricio Martinez-Barco, and Manuel Palomar. Fighting post-truth using natural language processing: A review and open challenges. *Expert Systems with Applications*, 141:112943, 2020.
- [5] Sameer Hinduja and Justin W Patchin. Bullying, cyberbullying, and suicide. *Archives of suicide research*, 14(3):206–221, 2010.
- [6] Saleem Alhabash, Anna R. McAlister, Chen Lou, and Amy Hagerstrom. From clicks to behaviors: The mediating effect of intentions to like, share, and comment on the relationship between message evaluations and offline behavioral intentions. *Journal of Interactive Advertising*, 15(2):82–96, 2015.
- [7] Shannon D Bailey and Lina A Ricciardelli. Social comparisons, appearance related comments, contingent self-esteem and their relationships with body dissatisfaction and eating disturbance among women. *Eating behaviors*, 11(2):107–112, 2010.
- [8] Mark Hsueh, Kumar Yogeeswaran, and Sanna Malinen. Leave your comment below: Can biased online comments influence our own prejudicial attitudes and behaviors? *Human communication research*, 41(4):557–576, 2015.
- [9] Mark Hsueh, Kumar Yogeeswaran, and Sanna Malinen. Leave your comment below: Can biased online comments influence our own prejudicial attitudes and behaviors? *Human communication research*, 41(4):557–576, 2015.