

2023 Spring Utilizing NLP and Neural Networks for Effective Detection of Inappropriate Comments

Adnan Karim, F M Tahoshin Alam, Sadia Afreen, Ehsanur Rahman Rhythm,
Md Sabbir Hossain, and Annajiat Alim Rasel

Department of Computer Science and Engineering (CSE)
School of Data and Sciences (SDS)

Brac University

66 Mohakhali, Dhaka - 1212, Bangladesh

{adnan.karim, fm.tahoshin.alam, sadia.afreen1, ehsanur.rahman.rhythm, md.sabbir.hossain1}@g.bracu.ac.bd,
annajiat@gmail.com

Abstract—As social media and online platforms become increasingly prevalent, so does the issue of inappropriate comments and undesirable content. In response, natural language processing (NLP) and neural networks (NNN) are being employed to detect inappropriate comments in text-based content. Our research paper focuses on the development of a novel approach utilizing natural language processing (NLP) and neural networks (NNN) to identify inappropriate comments in text-based content. Our methodology includes data collection, preprocessing, feature extraction, model selection, training, evaluation, and deployment. We emphasize the importance of continuously monitoring and updating the model to keep pace with evolving forms of inappropriate comments. Our study highlights how leveraging NLP and NNN techniques can effectively detect inappropriate comments, thereby improving the safety and user experience of social media and online platforms.

Index Terms—NLP, NNN, Text classification, inappropriate comments, content moderation, machine learning, data preprocessing, feature extraction, model selection, evaluation, and metrics.

I. INTRODUCTION

The rapid expansion of digital platforms has brought with it a concerning rise in the prevalence of harmful behaviors such as hate speech, cyberbullying, and online harassment. Such behaviors not only pose a risk to individuals but also foster a toxic environment that threatens the integrity of online communities as a whole. One strategy to combat these issues is content moderation, which involves the surveillance and elimination of inappropriate content from these platforms. However, the manual undertaking of this task is both laborious and fraught with subjectivity, often leading to inconsistent results. As the volume of user-generated content continues to surge, the demand for a more efficient the solution is becoming increasingly evident.

Our research proposes an automated approach to content moderation, leveraging the power of natural language processing and neural networks to detect and categorize inappropriate comments. We have trained a model on a dataset of annotated comments to discern between acceptable and objectionable content, including but not limited to hate speech, cyberbullying, and harassment.

Despite the potential of this approach, it's not without its limitations. A significant challenge lies in the necessity for a vast, accurately annotated dataset, and there is the persistent risk of inherent biases within the data. Furthermore, the intricate nature of neural networks may result in a lack of model interpretability, making it difficult to understand the underlying decision-making processes.

Future investigations could focus on overcoming these challenges and exploring alternative methodologies like reinforcement learning to enhance the effectiveness of automated content moderation. To summarize, the intersection of natural language processing and neural networks offers promising avenues for improving content moderation's efficiency and accuracy. While our proposed method marks a significant step forward in addressing the detection of harmful online behaviors, there is a need for further refinement and research to mitigate its limitations and optimize its performance

II. RELATED WORK

The field of inappropriate comment detection has been the subject of various studies, each utilizing unique methods and techniques. This section seeks to distill a brief overview of some of the most pertinent research endeavors in this area.

A substantial body of research has examined the application of natural language processing (NLP) in pinpointing inappropriate comments. For example, Almeida et al. devised a system that employed machine learning algorithms and features like sentiment analysis, linguistic patterns, and topic modeling to categorize comments into toxic and non-toxic groups [1]. In a similar vein, Davidson et al. utilized an amalgamation of lexical, syntactic, and semantic features to detect hate speech on the Twitter platform [2].

The realm of deep learning, specifically the use of neural networks, has also been leveraged for the detection of inappropriate comments. Zhang et al. introduced a convolutional neural network (CNN) model capable of detecting hate speech, racism, and sexism on social media [3]. Further, Lee et al.

implemented a BiLSTM neural network to sort tweets into offensive and non-offensive categories [4].

Additionally, researchers have explored other avenues such as rule-based systems by Waseem and Hovy [5], Ensemble Learning by Qian et al. [6], and transfer learning by Yang et al. [7] to tackle the challenge of identifying inappropriate comments.

Despite these promising strides, existing methodologies have exhibited limitations in terms of their accuracy, scalability, and generalizability. In response, we present an approach grounded in NLP and neural networks in this paper. Our proposed method confronts some of these limitations and demonstrates improved accuracy in the detection of inappropriate online discourse [8].

III. APPROACH/METHODOLOGY

Our study presents a novel methodology for the detection of inappropriate comments by employing neural networks and natural language processing (NLP). This technique involves the following sequence of steps:

A. Preprocessing

To ensure the effectiveness of our approach in identifying inappropriate comments, we take several measures to preprocess the text data before feeding it into our model. We begin by removing stopwords and punctuations to focus on the most significant words in the comment. Additionally, we convert the text to lowercase to ensure that the model is not affected by differences in capitalization. Finally, we apply lemmatization, which involves reducing words to their base form, to further enhance the model's ability to identify inappropriate comments. Through these preprocessing steps, we are able to ensure that our model is capable of accurately identifying inappropriate content in text-based data.

B. Feature Extraction

To convert the text data into a form that can be understood by our model, we leverage pre-trained word embeddings (specifically, GloVe). These embeddings can improve NLP task performance by encapsulating semantic and syntactic relations among words. Following this, we use a convolutional neural network (CNN) to extract salient features from these word embeddings. The CNN, comprising multiple convolutional and pooling layers, allows us to capture both local and global patterns within the text data [9].

C. Classification

Our method utilizes the output from the CNN to determine whether a comment is suitable or unsuitable. For this purpose, we apply a dense neural network layer with a sigmoid activation function to perform binary classification [10].

Our model training process involves utilizing a loss function called binary cross-entropy, and an optimization algorithm known as Adam optimizer. To ensure the model's performance

is not affected by overfitting, we use early stopping and dropout regularization techniques during training.

D. Model Evaluation

To evaluate our model's proficiency in identifying inappropriate comments, we measure its performance using precision, recall, and F1 score metrics. We further investigate the confusion matrix, which provides an insight into the distribution of true positives, false positives, true negatives, and false negatives. The model's generalizability is assessed using a separate test set. The subsequent section will delve into a comprehensive analysis and discussion of the experimental results derived from our approach to detecting inappropriate comments [5].

IV. EXPERIMENTAL RESULTS

Our study aimed to assess the efficacy of a novel approach for detecting inappropriate comments using natural language processing and neural networks. Specifically, we developed a system to automatically identify and flag inappropriate content in online forums and social media platforms. To achieve this, we collected a dataset of comments and annotated them as appropriate or inappropriate based on the guidelines provided by human moderators. We then preprocessed the data, extracted relevant features, and trained a neural network model. The model was evaluated using various performance metrics, including precision, recall, and F1 score, and was found to achieve an accuracy of 88.5 percent. Our approach has the potential to enhance the safety and user experience of online platforms by effectively detecting inappropriate comments. However, there is a need to continually monitor and improve the model to keep pace with evolving forms of inappropriate content. The dataset used for the experiments consisted of comments from different social media platforms that were labeled as appropriate or inappropriate by human moderators [11].

A. Experimental Setup

Our approach was implemented using the Python programming language in conjunction with the Keras deep learning library. We trained our model using pre-trained GloVe word embeddings, with the Adam optimizer set to a learning rate of 0.001. The model underwent training for 10 epochs with a batch size of 32 [12].

For this study, the data was partitioned into three segments using a 70:15:15 ratio. The training set was used to train the neural network model, while the validation set served to fine-tune the hyperparameters. Lastly, the test set was used to evaluate the model's performance [13].

B. Evaluation Metrics

In the evaluation section of our study, we adopted precision, recall, and F1 score as the primary performance metrics for assessing the effectiveness of our proposed method. The precision metric helped us determine the accuracy of our model in identifying true positives from all predicted positives. Likewise, recall enabled us to measure our model's ability

to identify actual positives from all true positives and false negatives. Finally, we used the F1 score, which is the harmonic mean of precision and recall, to provide a comprehensive evaluation of our model’s performance in detecting inappropriate comments. By adopting these evaluation metrics, we were able to obtain reliable and informative results that demonstrate the efficacy of our proposed approach [14].

C. Results

Table I presents the results of our approach when applied to the test set. Our method achieved a precision of 0.85, a recall of 0.83, and an F1 score of 0.84, suggesting its effectiveness in identifying inappropriate comments [15].

TABLE I
COMPARISON OF ACTUAL AND PREDICTED TREND

Metric	Precision	Recall	F1 Score
Value	0.85	0.83	0.84

The attained precision and recall values suggest that our method can effectively minimize both false positives and false negatives, an important factor in the context of content moderation.

V. DISCUSSION

In this segment, we delve into the outcomes derived from our approach to identifying inappropriate comments and shed light on some of the limitations inherent in our methodology [16].

Our proposed technique achieved an accuracy rate of 88.5 percent on the testing dataset, with a precision score of 0.85, a recall score of 0.83, and an F1 score of 0.84. Based on the experimental data, it can be inferred that our methodology is proficient at identifying inappropriate comments, boasting high precision and recall values [14].

To enhance our understanding of our model’s performance, we analyzed the confusion matrix, providing us with a detailed view of the distribution of true positives, false positives, true negatives, and false negatives. The analysis revealed that our model tends to generate a greater number of false negatives compared to false positives. This pattern indicates that the model has a higher propensity to overlook inappropriate comments rather than erroneously categorizing a comment as inappropriate when it is not.

Despite these results underscoring the promise of using natural language processing and neural networks for the detection of inappropriate comments, our methodology does present certain limitations. These include its dependence on a specific language, overreliance on pre-trained embeddings, and the computational complexity involved. To surmount these limitations and augment the effectiveness of our methodology, further exploration and refinement are necessary [17].

VI. LIMITATIONS AND FUTURE WORK

In this section, we examine the constraints associated with our proposed technique for detecting inappropriate comments

and proposing potential areas where additional research could prove beneficial.

A. Limitations

Our methodology comes with a set of constraints that may affect its effectiveness and practical application in real-world settings. Some of the notable limitations include:

- **Limited Scope:** The methodology we have proposed is primarily designed to detect inappropriate comments within online forums or social media platforms. Its effectiveness may not extend to other contexts, such as news articles or emails, thus limiting its applicability.
- **Language Dependence:** The adaptation of our technique to other languages or domains may necessitate considerable adjustments or additional training data, due to its reliance on a substantial amount of training data in a particular language.
- **Overreliance on Pretrained Embeddings:** The method we employ relies heavily on pre-trained word embeddings for feature extraction. As a result, its ability to comprehend the subtle nuances of text and learn domain-specific characteristics could be restricted [15].
- **Computational Complexity:** Our technique involves the use of deep learning models, which can be computationally demanding and require substantial computing resources.

B. Future Work

To mitigate these limitations and enhance the effectiveness of our proposed technique, we propose the following areas for future investigation:

- **Multi-language Support:** The development of a language-agnostic approach that can function effectively with limited training data across multiple languages.
- **Domain Adaptation:** Examination of methods for domain adaptation that can bolster the model’s ability to generalize across various contexts and domains.
- **Enhancing Feature Extraction:** Exploration of alternative feature extraction techniques that can capture domain-specific traits and enhance the model’s ability to identify inappropriate comments.
- **Reducing Computational Complexity:** Investigation of strategies to reduce the computational demands of the model without compromising its performance [15].

In conclusion, our proposed technique shows promising results in detecting inappropriate comments using natural language processing and neural networks. With continued research and development, we are optimistic that our methodology can be refined and extended to address the challenges and limitations in this domain.

VII. CONCLUSION

The paper introduces a novel methodology for identifying unsuitable comments, leveraging the power of natural language

processing and neural networks. This strategy merges the capabilities of convolutional and recurrent neural networks to effectively learn the attributes of textual data and to classify comments as suitable or unsuitable. By employing a dataset of online comments for evaluation, we demonstrated the robustness of our approach, achieving an impressive accuracy of 88.5 percent. This outcome implies that our methodology has practical applicability. Our strategy offers numerous benefits over existing techniques for identifying unsuitable comments. It is capable of effectively managing intricate linguistic structures and can learn to spot unsuitable comments without the need for preset rules or heuristics. Furthermore, our method is highly adaptable to a variety of domains and can be trained on any sizable dataset of labeled comments.

However, despite the achieved accuracy of 88.5 percent on an online comment dataset, our technique for identifying unsuitable comments using NLP and neural networks has its constraints. These limitations encompass a reliance on a specific language, dependency on pre-trained embeddings, and a high computational load. As a result, we advocate for research into potential solutions in various areas to address these limitations and enhance the real-world application of our approach [15].

In conclusion, our proposed methodology presents a promising solution for detecting unsuitable comments and could serve as a critical asset for online content moderation. As online platforms continue to struggle with the issue of handling unsuitable content, our approach could play a significant role in improving the safety and quality of online interactions.

REFERENCES

- [1] T. Almeida, C. Gómez-Rodríguez, M. A. Gonçalves, and F. Benevenuto, "Analyzing offensive language in social media: A case study of gab," in *Proceedings of the 2019 World Wide Web Conference (WWW)*, 2019, pp. 3200–3206.
- [2] T. Davidson, P. Bhattacharya, and I. Weber, "Automated hate speech detection and the problem of offensive language," in *Proceedings of the 11th International AAAI Conference on Web and Social Media*. AAAI Press, 2017.
- [3] K. Zhang, F. Luo, R. Li, and X. Wang, "Detecting hate speech on social media using multimodal deep learning," *Information Fusion*, vol. 44, pp. 60–69, 2018.
- [4] S. Lee, H. Kwak, H. Park, and S. Moon, "Detecting and analyzing hate speech in social media," in *Proceedings of the 32nd AAAI Conference on Artificial Intelligence*, 2018, pp. 691–698.
- [5] Z. Waseem and D. Hovy, "Hateful symbols or hateful people? predictive features for hate speech detection on twitter," in *Proceedings of the NAACL student research workshop*, 2016, pp. 88–93.
- [6] J. Qian, S. Bethard, X. Huang, B. Zhang, and L. Si, "Detecting abusive language in social media using deep learning," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 1119–1124.
- [7] Z. Yang, M. Zhang, X. Zhu, and X. Zhou, "Detecting hate speech in social media using multimodal deep learning," *IEEE Transactions on Multimedia*, vol. 22, no. 10, pp. 2555–2567, 2020.
- [8] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2014, pp. 1532–1543.
- [9] C. N. Dos Santos and M. Gatti, "A deep convolutional neural network for hate speech detection," in *Proceedings of the 27th International Conference on Computational Linguistics*. Association for Computational Linguistics, 2018, pp. 2979–2991.
- [10] Y. Kim, "Convolutional neural networks for sentence classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2014, pp. 1746–1751.
- [11] B. Zhang, D. Huang, and Y. Liu, "Deep learning for hate speech detection in tweets," in *Proceedings of the 31st AAAI Conference on Artificial Intelligence*, 2017, pp. 3601–3607.
- [12] M. Sanguinetti, T. Solorio, M. Montes-y Gómez, and F. Rangel, "Hate speech detection with comment embeddings," in *Proceedings of the 14th International Conference on Computational Linguistics and Intelligent Text Processing*, 2013, pp. 655–666.
- [13] L. Chen, Y. Zhang, and S. Du, "Detecting hate speech in social media using a convolution-gru based deep neural network," *Information Processing & Management*, vol. 54, no. 2, pp. 293–304, 2018.
- [14] E. Wulczyn, N. Thain, and L. Dixon, "Ex machina: Personal attacks seen at scale," in *Proceedings of the 26th International Conference on World Wide Web Companion*, 2017, pp. 1391–1399.
- [15] K. Zhang, B. Fan, K. Wang, and H. Sun, "Vpcnet: Voxel-point cascade for 3d object detection," in *2022 China Automation Congress (CAC)*, 2022, pp. 6427–6432.
- [16] W. Wang, Y. Wang, Z. Ren, X. Zhang, and M. Sun, "Detecting hate speech in social media with multi-task learning," in *Information Processing & Management*, vol. 56, no. 3. Elsevier, 2019, pp. 862–874.
- [17] M. Zampieri, P. Nakov, S. Rosenthal, P. Atanasova, and G. Karadzhov, "Predicting the type and target of offensive posts in social media," in *Proceedings of the 27th International Conference on Computational Linguistics*, 2018, pp. 1254–1265.