

Evaluating Online Sexism Detection: A Comparative Study of Machine Learning Models using the EDOS Dataset

Abstract—Online sexism refers to gender-based discrimination and harassment that occurs in online spaces, such as social media platforms, online communities, and forums. Machine learning models can recognize and lessen online sexism by automatically detecting sexist content in social media posts. In this experimental analysis, we have evaluated the performance of five different machine learning models, including Logistic Regression, Gaussian Naive Bayes, Decision Tree Classifier, Support Vector Machine (SVM), and KNeighbors Classifier. Our objective was to detect online sexism using the Explainable Detection of Online Sexism (EDOS) data set. We preprocess the data set by cleansing the text data with a regular expression, removing null values, removing redundant columns, and vectorizing it with TfidfVectorizer. The results of our study indicate that the Logistic Regression, Gaussian Naive Bayes, Decision Tree Classifier, and Support Vector Machine (SVM) models are efficacious in detecting occurrences of online sexism. Nonetheless, the KNeighbors Classifier algorithm shows comparatively lower accuracy in this aspect. The present analysis highlights the capacity of machine learning models to identify instances of online sexism. The highest accuracy score we obtained for the Support Vector Machine (SVM) model is 94.64%.

Index Terms—online sexism, machine learning, text preprocessing, vectorization, social media, online safety

I. INTRODUCTION

Sexism is a type of gender-based abuse in which people are treated badly because of their gender. This kind of discrimination can come in many forms, like abuse, making someone feel like an object, and violence based on their gender. Persistent societal challenges include gender discrimination and sexism, exerting enduring adverse effects on individuals and society. Despite advancements toward gender equality, manifestations persist, notably in online platforms where technology is employed for discriminatory and harassing behaviors, underscoring the need for comprehensive interventions to address these pressing issues. Online sexism comes in many forms, such as trolling, and making women feel like objects. It can hurt people's mental health, self-esteem, and sense of safety in very bad ways. It can also lead to a society of sexism and inequality between men and women, where women are seen as less important than men. This study employs natural language processing to illuminate instances of sexist language on prevalent social media platforms. Methodologically, different machine learning models are applied to discern and analyze discernible patterns within the linguistic landscape,

contributing to the broader discourse on digital communication and gender dynamics. This can be very useful for social media managers and law enforcement authorities to help identify the patterns of abusive messages in their social media presence. We have used a publicly available Kaggle dataset [10] to detect such patterns.

II. RELATED WORK

Several works were done in this field. The paper [1] proposes a way of detecting email-based cyberstalking on textual data. They have followed the EBCD multi-model soft voting technique of machine learning. Here, three types of datasets are D1, D2, and D3. D1 contains the spam email subject, D2 contains the spam email body and lastly, D3 carries the documents related to cyber harassment. Initially, they combined all the classifiers of three datasets to make the detection smooth. Then they applied the EBCD model which classified normal emails and suspicious emails. In every trial, several machine learning models like random forest, naive bias, and logistic regression were used for the test. The final decision was taken by the soft-voting technique for each model. Each of the models has individual accuracy. The average accuracy was 96.3% whereas the precision, recall, and f1 score were 98.1%, 94%, and 95.9% respectively.

A similar type of paper was also written by Sule Kaya and Bilal Alatas in the ProQuest web journal. This paper introduces a hybrid model which can predict racist and sexist comments in social media. In this study, the LSTM Neural Network and Recurrent neural network were combinedly and developed. This model divides the texts into certain categories. Afterward, this hybrid model was compared with various classifier models like random forest, logistic regression, naive bias, etc. After doing all comparisons, it was found that LSTM Neural Network gave the best prediction result. The accuracy was 95.25% but the F1 score was 51.32% due to lack of dataset [2].

John Hani along with five other mates published a paper about cyberbullying and online harassment detection. The following paper uses two machine learning approaches to detect social media bullying which are Neural Network and Support Vector Machine model. The whole process was divided into three steps. They are pre-processing feature extraction, and lastly classification. In data pre-processing word tokenization,

lowerization was performed to make the data clean and noise-free. In the second part, textual data were tuned into a particular format so that it could go with the model with the help of the TF IDF model. With this model, word weights were calculated based on the sentence. Afterward, word polarity was calculated with the help of sentiment analysis. After a successful feature extraction, those features were sent to SVM and Neural Network classifiers. The accuracy of SVM and Neural Network were 90.3 and 92.8 respectively [3].

Zahraa Jihad Berjawi from the American University Of Beirut published a paper that describes the way of detecting sexism in text via NLP and the Machine Learning Model. Multiple models were used in this study. For example, the linear-based model (logistic regression), distance-based model (K-neighbors and SVM), probabilistic model (Naive Bias classifiers), and lastly tree-based model (Decision tree and random forest). Among all of these, Random forest showed the best performance in the case of cross-validation. After doing all of the testing we got the cross-validation accuracy and feature extracting accuracy 0.720943 and 0.726125 respectively [4].

Francisco Rodríguez-Sánchez published a paper in IEEE where several analyses were done to determine how sexism is found in our day-to-day conversation, particularly on Twitter. Here a fully automatic system was presented based on the Machine Learning model. It shows that several models were compared based on the tf-idf feature but BERT shows the best performance among all of those with an accuracy of 74% [6].

Djuric, N. et al. assessed their method using an extensive dataset of user feedback gathered from the Yahoo Finance website. This dataset consists of 56,280 comments that include hate speech, as well as 895,456 clean comments. These comments were produced by 209,776 anonymous users and were collected and labeled by editors during a span of six months [8].

Xu and Zhu introduced a paper that focuses on offensive language in online community text messages and introduces a novel automatic sentence-level filtering method. This approach leverages grammatical relations between words to semantically eliminate offensive content. Compared to existing methods, the proposed approach yields filtering results that closely resemble manual filtering. The researchers curated a dataset through meticulous screening of 11,000 YouTube text comments. Experimental outcomes demonstrate a filtering concurrence exceeding 90%, affirming the efficacy of the proposed methodology in comparison to manual curation. [9].

Zeeraq Waseem and Dirk Hovy introduced a dataset containing tweets labeled for instances of racism and sexism. Employing lexical modeling and bootstrap techniques, they constructed a classifier capable of identifying tweets exhibiting sexist and racist content, achieving an F1 score of 0.69. The study explores different feature combinations for a model's training. They use character n-grams and gender information, finding a slight score increase. Another set involves tweet description lengths and n-gram lengths. Gender + location combines location and n-grams. This slightly worsens classifier performance. Gender + location + length includes gender,

location, and length features. The F1 scores for these combinations are 73.89 (gender), 73.66 (length), 73.62 (gender + location), and 73.47 (gender + location + length) [11].

An assessment conducted by S Zimmerman, U Kruschwitz, and C Fox involved an ensemble of deep learning methods applied to the Twitter dataset. Their results surpassed the performance of individual state-of-the-art deep learning classifiers trained for the identical purpose, achieving an impressive accuracy level of 94%. This outcome underscores the potency of deep learning in elevating the accuracy standards of hate speech detection systems [12].

Gitari et al. analyzed sentences from notable U.S. "hate sites," classifying them as strongly hateful (SH), weakly hateful (WH), or non-hateful (NH). Utilizing semantic attributes and grammatical patterns, they applied the classification to a separate test dataset, yielding a 65.12% F1-score [13].

Kwok and Wang focused on finding hurtful tweets directed at black people. They looked at individual words in the tweets and were able to accurately tell whether the tweets were mean or not about 76% of the time. This shows their method worked well to identify offensive content related to the black community on social media [14].

K. S. Alam et al. conducted research using different machine-learning methods to identify cyberbullying in tweets. They tried various techniques to extract features from the data and used different combinations of n-gram analyses. In their study, Logistic Regression and Bagging ensemble models stood out as the top performers in spotting cyberbullying. The proposed SLE and DLE models demonstrated superior performance with a remarkable 96% accuracy, leveraging TF-IDF (Unigram) feature extraction in conjunction with K-Fold cross-validation. This outcome underscores the efficacy of the devised methodology in enhancing accuracy in the studied context. [15].

III. MODEL APPROACHES

A. Logistic Regression

Logistic Regression, a robust statistical method, delineates the relationship between a dependent binary variable and an independent non-binary variable. It serves as a vital tool for researchers, providing systematic insights into complex datasets and contributing to a nuanced understanding of variable interactions.

B. Gaussian Naive Bayes

Gaussian Naive Bayes addresses classification complexities by presuming normal feature distributions within classes. This method harmonizes probability and classification, shedding light on the route to informed decision-making in the realm of machine learning. The model's reliance on Gaussian distributions reflects a nuanced understanding of underlying data patterns in diverse classification scenarios.

C. KNeighborsClassifier

Utilized in classification and regression, KNN is a data analysis method assigning a point's classification or prediction

through the majority class or average of its nearest neighbors in the feature space. This approach enhances efficacy across various analytical applications, showcasing KNN's versatility and relevance in data-driven research and problem-solving scenarios.

D. DecisionTreeClassifier

The decision tree classifier partitions the input feature space iteratively, forming a tree structure. Internal nodes signify decisions based on features, guiding the classification process towards leaves with class labels. The algorithm strives to construct an optimal decision sequence for effective data point separation and classification.

E. Support Vector Machine

The SVM algorithm seeks to construct an optimal hyper-plane or decision boundary within n-dimensional space, facilitating the classification of data points into distinct categories. This delineation enhances the algorithm's ability to efficiently categorize new data points, contributing to its efficacy in machine learning applications.

IV. EXPERIMENTAL ANALYSIS

A. Dataset

The Explainable Detection of Online Sexism (EDOS) [10] dataset was used for this research, which contains 14,000 unique values. The dataset includes five columns, namely `rewire_id`, `text`, `label_sexist`, `label_category`, and `label_vector`. The `text` column contains the social media posts, the `label_sexist` column contains labels for the sexism classification, and the other columns contain additional information about the post. The main goal of this research is to develop models that can predict whether a given social media post is sexist or not. As only the sexist and not sexist classification is required, unnecessary columns such as `rewire_id`, `label_category`, and `label_vector` were dropped from the dataset.

In pursuit of the research goal, various machine learning models were trained and evaluated on the EDOS dataset, employing techniques such as data preprocessing, text embedding, and rigorous model assessment. The ultimate objective was to develop a highly accurate and interpretable model for identifying instances of online sexism in social media posts. By focusing on this critical issue, the research aimed to contribute to a safer and more inclusive online environment.

Before commencing the analytical process, an examination of the class distribution was undertaken to enhance comprehension of the dataset. Employing a pie chart, the distribution of the two classes revealed that 75.7% of posts exhibited non-sexist content, while 24.3% contained sexist elements. This observation underscores the dataset's imbalanced nature, presenting a potential impediment in developing a model capable of achieving optimal performance across both classes. Addressing this imbalance is crucial for mitigating biased model outcomes and optimizing predictive accuracy in the classification task.

Distribution of Labels of Original Dataset

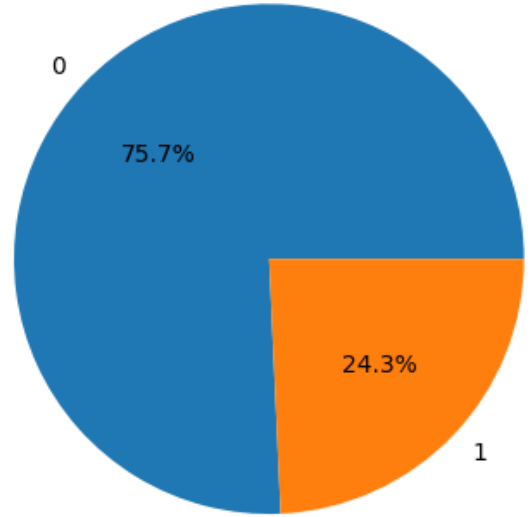


Fig. 1. Distribution of Sexist and Non-Sexist Posts in the EDOS Dataset.

Distribution of Labels of Balanced Dataset

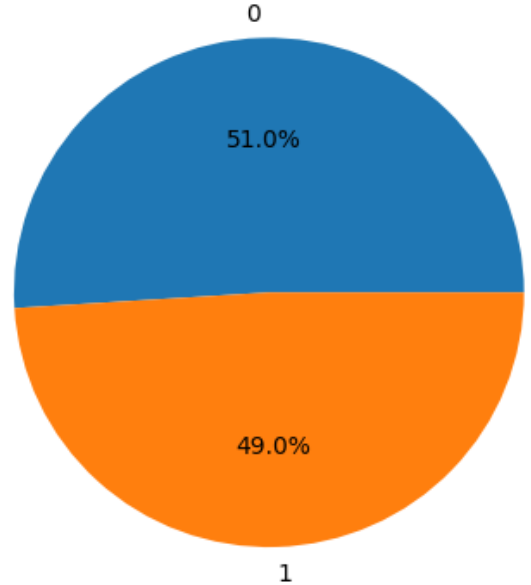


Fig. 2. Distribution of Balanced Sexist and Non-Sexist Posts in the EDOS Dataset.

B. Pre-Processing

After removing the unnecessary columns, two columns were remaining, namely `text` and `label_sexist`. The data was then converted to numerical data so that it could be used

for machine learning computations. Since there are only two classes, they were mapped to 0 and 1, with the sexist class being mapped to 1 and the not_sexist class being mapped to 0.

We over-sampled the sexist data within this dataset as a measure to achieve dataset balance. Following the data over-sampling process, the dataset attains equilibrium with 49% of instances classified as 'sexist' and 51% categorized as 'not sexist' data.

Before proceeding with the analysis, it was important to check whether the dataset had any null or missing values. The dataset was checked for null values, and none were found. However, it is always recommended to perform a thorough check for null values, especially when working with a large dataset.

The text data required cleaning as there was noise in the data. Regular expressions were used to perform the following cleaning steps:

- Any text inside square brackets was removed along with the brackets from the text. Square brackets are commonly used to denote meta information or citations and are not relevant to the classification task.
- Non-word characters were replaced with a space. Non-word characters such as hashtags, at signs, or punctuation marks can be noisy and irrelevant to the classification task.
- Any URLs were removed from the text. URLs are also commonly found in social media posts but are not relevant for the classification task.
- Any HTML tags were removed from the text. HTML tags are commonly found on web pages and can be irrelevant to the classification task.
- Any punctuation marks were removed from the text. Punctuation marks such as periods, commas, or exclamation marks can be noisy and irrelevant to the classification task.
- Any newlines were removed from the text. Newlines are commonly found in text and can be irrelevant to the classification task.
- Any alphanumeric character that contained digits was removed from the text. Alphanumeric characters such as emojis or emoticons can be noisy and irrelevant to the classification task.

The dataset underwent stratified splitting into training (80%) and testing (20%) sets using the sklearn model_selection's train_test_split function. The 'text' column served as the feature variable, while the 'label_sexist' column was the target variable. A random state of 10 was applied for reproducibility, ensuring consistency in results during subsequent text data vectorization.

After cleaning the text data, it was vectorized using TfidfVectorizer from the sklearn feature extraction library. The vectorized text data was then used to train the models. Vectorization is an important step as it converts text data into numerical data, which can be used by machine learning models. The TfidfVectorizer was used as it is a commonly used

vectorizer for text data and takes into account the importance of each word in the document. Utilizing a hybrid approach, this methodology employs both term frequency (TF) within a document and inverse document frequency (IDF) across the corpus to ascertain the significance of a word. The resulting metric captures the nuanced relevance of terms, contributing to a more refined analysis of textual data in academic research.

C. Results

In this section, we present a comprehensive comparison of the performance of five different classification models on both imbalanced and balanced datasets. Table I illustrates the results obtained from the imbalanced dataset, highlighting the challenges posed by class imbalance. Subsequently, Table II showcases the performance of the same models after manual dataset balancing, revealing the impact of addressing the class imbalance issue.

Logistic Regression: In the context of imbalanced datasets, the Logistic Regression model demonstrated an initial accuracy of 81.71%. Precision and recall were observed at 82% and 65%, respectively, resulting in an F1 score of 68%. Subsequent manual rebalancing of the dataset yielded notable enhancements in model performance. The accuracy persisted at 81.78%, with precision and recall both ascending to 82%. Consequently, the F1 score exhibited a commensurate increase to 82%, underscoring the model's enhanced proficiency in discerning positive instances while upholding overall accuracy. This underscores the significance of dataset balance in optimizing the predictive capabilities of machine learning models and highlights the potential for refinement through deliberate adjustments to class distribution.

Decision Tree Classifier: For the imbalanced dataset, the Decision Tree Classifier achieved an accuracy of 75.68%. Its precision and recall were close at 67% and 68% respectively, resulting in an F1 score of 68%. Upon dataset balancing, the model's performance significantly improved, with accuracy rising to 89.04%. The precision and recall values also reached 90% and 89% respectively, leading to a balanced F1 score of 89%. This enhancement suggests that the Decision Tree Classifier was able to better generalize its learned patterns after addressing the class imbalance.

Support Vector Machine (SVM): The Support Vector Machine (SVM) demonstrated robust performance on an imbalanced dataset, achieving an accuracy of 81.86%, characterized by a noteworthy precision of 85% but a relatively lower recall of 64%, indicative of a bias toward negative class classification owing to dataset imbalance. Following dataset rebalancing, the SVM exhibited a substantial accuracy enhancement, reaching 94.64%. Notably, the model displayed outstanding precision (95%) and a commendable recall (85%), resulting in an impressive F1 score of 95%. This underscores the SVM's proficiency in discriminating between classes, particularly in mitigating the adverse impact of imbalanced data. The findings emphasize the significance of preprocessing techniques in enhancing the overall effectiveness of support vector machines in classification tasks.

Gaussian Naive Bayes: In the context of an imbalanced dataset, the Gaussian Naive Bayes model exhibited an initial accuracy of 75.50%, with commendable precision at 88% but a relatively lower recall of 51%. This imbalanced recall adversely affected the F1 score, resulting in a suboptimal value of 44%. After dataset rebalancing, the model demonstrated a modest accuracy enhancement to 80.22%. Notably, both precision and recall showed improvement, reaching 82% and 80% respectively, culminating in a more equilibrium F1 score of 80%. Despite this progress, challenges persist in accurately identifying positive instances, reflecting the enduring impact of the inherent class distribution. These findings underscore the importance of addressing imbalances to enhance model performance in real-world applications.

K-Nearest Neighbors (KNN) Classifier: The K-Nearest Neighbors (KNN) classifier demonstrated a 75.67% accuracy on an imbalanced dataset, featuring a precision of 58% and a recall of 50%, yielding an F1 score of 43%, indicative of potential enhancement opportunities. Upon dataset balancing, the accuracy diminished to 51.15%, underscoring KNN's sensitivity to alterations in class distribution. Despite a sustained precision of 73%, the recall declined to 51%, resulting in a reduced F1 score of 36%. The performance degradation post-balancing implies that KNN may not align optimally with the dataset or the specific balancing procedures applied. This underscores the need for careful consideration of classifier characteristics and dataset attributes to ascertain the most suitable model for effective classification in imbalanced scenarios.

In conclusion, our study reveals that manual balancing to address class imbalance yields significant enhancements across various models. Noteworthy improvements were observed in precision, recall, and F1 score for models such as the Support Vector Machine in Figure 3 and Decision Tree Classifier in Figure 4. These outcomes underscore the efficacy of class imbalance management. Nevertheless, caution is advised in algorithm selection for imbalanced datasets, as demonstrated by the diminished performance of the K-Nearest Neighbors model. This emphasizes the importance of thoughtful algorithmic choices in addressing class imbalance challenges.

TABLE I
ACCURACY, PRECISION, RECALL, F1 SCORE WITH ORIGINAL
IMBALANCED DATA-SET

Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)
SupportVectorMachine	81.86%	85%	64%	67%
Logistic Regression	81.71%	82%	65%	68%
DecisionTreeClassifier	75.68%	67%	68%	68%
KNeighborsClassifier	75.67%	58%	50%	43%
Gaussian Naive Bayes	75.50%	88%	51%	44%

TABLE II
ACCURACY, PRECISION, RECALL, F1 SCORE OF GIVEN VALUES WITH
MANUAL BALANCED DATA-SET

Model	Accuracy (%)	Precision (%)	Recall (%)	F1 Score (%)
SupportVectorMachine	94.64%	95%	85%	95%
DecisionTreeClassifier	89.04%	90%	89%	89%
Logistic Regression	81.78%	82%	82%	82%
Gaussian Naive Bayes	80.22%	82%	80%	80%
KNeighborsClassifier	51.15%	73%	51%	36%

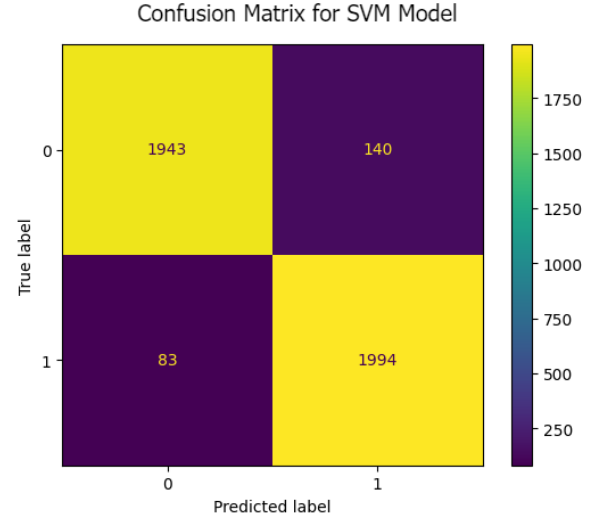


Fig. 3. Confusion matrix for SVM model in the balanced dataset.

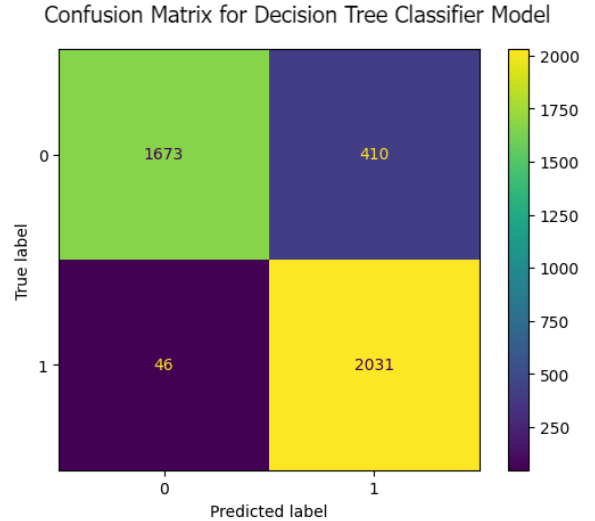


Fig. 4. Confusion matrix for Decision Tree Classifier model in the balanced dataset.

V. CONCLUSION AND FUTURE WORK

In this experimental analysis, we developed and evaluated machine learning models for detecting online sexism in social media posts.

In the imbalanced dataset scenario (Table I), all models showed reasonable accuracy but demonstrated varying precision, recall, and F1 scores. Both Logistic Regression and Support Vector Machine exhibited balanced performance. The Decision Tree and Naive Bayes achieved average precision and recall, while K-Nearest Neighbors struggled with lower recall and F1 scores. However, in the manually balanced dataset (Table II), models exhibited general improvements across the board. The Decision Tree excelled in precision and recall, the Support Vector Machine outperformed in accuracy,

and Naive Bayes achieved more balanced metrics. Despite this, K-Nearest Neighbors continued to face challenges due to its sensitivity to imbalance. The utilization of a balanced dataset substantially improved model efficacy, as evidenced by heightened recall and F1 score metrics. This underscores the critical importance of mitigating class imbalance to enhance the performance of machine learning models.

In summary, our study investigated the performance of five classification models on imbalanced and balanced datasets. The results from the imbalanced dataset (Table I) underscored the challenges of handling minority class instances, with varying impacts on the accuracy, precision, recall, and F1 scores. Upon manual dataset balancing (Table II), models like Decision Tree Classifier and Support Vector Machine exhibited substantial performance improvements, showcasing their adaptability to balanced data. Gaussian Naive Bayes achieved a more balanced performance, while K-Nearest Neighbors struggled post-balancing. These findings emphasize the importance of tailored approaches for imbalanced data and the need to consider model suitability. The study contributes insights into effectively navigating imbalanced scenarios, prompting further research for more robust machine learning solutions in real-world applications.

This study underscores the efficacy of machine learning models in identifying online sexism, fostering a safer and more inclusive digital ecosystem. Nonetheless, additional research is imperative to mitigate the study's inherent limitations, ensuring a comprehensive understanding and implementation of these models for sustained advancements in online social dynamics. For instance, the EDOS dataset is relatively small, and there is a need for more diverse datasets to improve the generalizability of the models. Additionally, the models could benefit from more advanced techniques, such as deep learning, to capture the complex and subtle nature of online sexism.

In the future, researchers can expand on this study by incorporating more advanced techniques to improve the accuracy and generalizability of the machine learning models. For instance, incorporating contextual features such as the identity of the user or the specific social media platform could enhance the models' performance. Additionally, future work can focus on evaluating the models' performance on more diverse and representative datasets to address potential biases. Moreover, developing explainable models that provide insights into the features that contribute to the model's predictions can aid in understanding the complex nature of online sexism. The advancement of precise and dependable machine learning models holds promise for fostering a secure and inclusive digital ecosystem.

REFERENCES

- [1] Gautam, A.K. and Bansal, A. (2023) "Email-based cyberstalking detection on textual data using multi-model soft voting technique of machine learning approach," *Journal of Computer Information Systems*, pp. 1–20. Available at: <https://doi.org/10.1080/08874417.2022.2155267>.
- [2] Kaya, S. and Alatas, B. (2022) "A new hybrid LSTM-RNN deep learning based racism, Xenomy, and Genderism Detection Model in online social network," *International Journal of Advanced Net-*
- working and Applications*, 14(02), pp. 5318–5328. Available at: <https://doi.org/10.35444/ijana.2022.14201>.
- [3] Hani, J. et al. (2019) "Social Media Cyberbullying Detection Using Machine Learning," *International Journal of Advanced Computer Science and Applications*, 10(5). Available at: <https://doi.org/10.14569/ijacsa.2019.0100587>.
- [4] American University of Beirut Benevolent Sexism Detection in text: A ... (no date). Available at: https://scholarworks.aub.edu.lb/bitstream/handle/10938/23599/Berjawi-Zahraa_2022.pdf (Accessed: April 15, 2023).
- [5] Grosz, D., Conde-Cespedes, P. (2020). Automatic Detection of Sexist Statements. Commonly Used at the Workplace. arXiv (Cornell University). <https://doi.org/10.48550/arxiv.2007.04181>
- [6] Rodriguez-Sanchez, F., Carrillo-de-Albornoz, J. and Plaza, L. (2020) "Automatic classification of sexism in social networks: An empirical study on Twitter data," *IEEE Access*, 8, pp. 219563–219576. Available at: <https://doi.org/10.1109/access.2020.3042604>.
- [7] Katie Bouman (no date) Katie Bouman aka Katherine L. Bouman. Available at: <http://users.cms.caltech.edu/~klbouman/> (Accessed: April 15, 2023).
- [8] Djuric, N., Zhou, J., Morris, R., Grbovic, M., Radosavljevic, V., & Bhamidipati, N. (2015). Hate Speech Detection with Comment Embeddings. Proceedings of the 24th International Conference on World Wide Web - WWW '15 Companion. doi:10.1145/2740908.2742760
- [9] Xu, Z., & Zhu, S. (2010). Filtering Offensive Language in Online Communities using Grammatical Relations. ResearchGate. <https://ceas.cc/2010/papers/Paper%2010.pdf>
- [10] Asad, M.U. (2023) Explainable detection of online sexism (EDOS), Kaggle. Available at: <https://www.kaggle.com/datasets/maifeeulasad/explainable-detection-of-online-sexism-edos> (Accessed: April 16, 2023).
- [11] Waseem, Z., & Hovy, D. (2016). Hateful Symbols or Hateful People? Predictive Features for Hate Speech Detection on Twitter. Association for Computational Linguistics. <https://doi.org/10.18653/v1/n16-2013>
- [12] Zimmerman, S. M., Kruschwitz, U., & Fox, C. (2018). Improving Hate Speech Detection with Deep Learning Ensembles. Language Resources and Evaluation. <https://www.aclweb.org/anthology/L18-1404.pdf>
- [13] N. D. Gitari, Z. Zuping, H. Damien and J. Long, "A lexicon-based approach for hate speech detection", *Int. J. Multimedia Ubiquitous Eng.*, vol. 10, no. 4, pp. 215-230, Apr. 2015.
- [14] I. Kwok and Y. Wang, "Locate the hate: Detecting tweets against blacks", *Proc. AAAI*, pp. 1621-1622, Jul. 2013.
- [15] K. S. Alam, S. Bhowmik and P. R. K. Prosun, "Cyberbullying Detection: An Ensemble Based Machine Learning Approach," 2021 Third International Conference on Intelligent Communication Technologies and Virtual Mobile Networks (ICICV), Tirunelveli, India, 2021, pp. 710-715, doi: 10.1109/ICICV50876.2021.9388499.