

Language Matters: A Systemic Functional Linguistics-Enhanced Machine Learning Framework for Cyberbullying Detection

Raghad Altowairgi[†], Ala Eshamwit^{††}, and Lobna Hsairi^{†††}

whyragad@gmail.com aaeshmawi@uj.edu.sa lalhabib@uj.edu.sa

[†]College of Computer Science and Engineering, University of Jeddah, Jeddah, Saudi Arabia

^{††}Faculty College of Computer Science and Engineering, University of Jeddah, Jeddah, Saudi Arabia

^{†††}Faculty College of Computer Science and Engineering, University of Jeddah, Jeddah, Saudi Arabia

Summary

Cyberbullying is a growing problem among adolescents and can have serious psychological and emotional consequences for the victims. In recent years, machine learning techniques have emerged as promising approach for detecting instances of cyberbullying in online communication. This research paper focuses on developing a machine learning models that are able to detect cyberbullying including support vector machines, naïve bayes, and random forests. The study uses a dataset of real-world examples of cyberbullying collected from Twitter and extracts features that represents the ideational metafunction, then evaluates the performance of each algorithm before and after considering the theory of systemic functional linguistics in terms of precision, recall, and F1-score. The result indicates that all three algorithms are effective at detecting cyberbullying with 92% for naïve bayes and an accuracy of 93% for both SVM and random forests. However, the study also highlights the challenges of accurately detecting cyberbullying, particularly given the nuanced and context-dependent nature of online communication. This paper concludes by discussing the implications of these findings for future research and the development of practical tool for cyberbullying prevention and intervention.

Keywords:

Cyberbullying; Machine learning; Systemic functional linguistics; Latent Dirichlet allocation

1. Introduction

Cyberbullying or cyber harassment is a form of bullying that occurs through electronic communication channels. It can happen to anyone, but it is particularly prevalent among children and adolescents who spend a lot of time online. Cyberbullying can take many forms, including sending hateful or threatening messages, posting embarrassing photos or videos, spreading rumors or lies, and creating fake profiles to impersonate or harass someone. It can happen on social media platforms, through text messages, in online chat rooms, or through email. The effects of cyberbullying can be devastating for the victim. It can lead to anxiety, depression, and other mental health issues, as well as lower self-esteem and academic performance. Victims may also experience social isolation and may be reluctant to

participate in activities they once enjoyed. In extreme cases, cyberbullying can even lead to physical harm. Preventing cyberbullying requires a multi-faceted approach that involves parents, educators, and technology companies. Parents can help by monitoring their children's online behavior and teaching them how to be responsible digital citizens. Educators can raise awareness about the dangers of cyberbullying and provide resources and support for victims. Technology companies can develop tools and policies to prevent cyberbullying on their platforms and respond quickly to reports of harassment. Together, these efforts can help to create a safer and more respectful online environment for everyone. Detecting cyberbullying using machine learning involves using computational techniques to automatically identify instances of cyberbullying in online text. Cyberbullying is a serious issue that can have negative impacts on mental health, social relationships, and academic performance, among other things [1].

Machine learning algorithms can be trained to recognize patterns in language use that are indicative of cyberbullying, allowing for more efficient and accurate detection. There are several approaches to detecting cyberbullying using machine learning. One common approach is to use supervised learning, where a machine learning model is trained on a labeled dataset of cyberbullying instances and non-cyberbullying instances. The model learns to recognize patterns in the language use of cyberbullies, such as the use of derogatory language, threats, or insults [2]. Once the model is trained, it can be used to automatically identify cyberbullying instances in new text data. Another approach is to use unsupervised learning, where the machine learning model is trained on an unlabeled dataset of text data. The model learns to identify patterns in the language use of cyberbullies without prior knowledge of what constitutes cyberbullying. This approach is useful when there is no labeled dataset available or when the nature of cyberbullying is constantly changing. Machine learning models can also be used in combination with natural language processing (NLP) techniques to improve the accuracy of cyberbullying detection. NLP techniques can be used to pre-process the text data, such as

removing stop words, stemming, and lemmatizing. This can help to reduce the noise in the text data and improve the performance of the machine learning model [3]. Overall, using machine learning to detect cyberbullying is a promising approach that has the potential to improve the efficiency and accuracy of cyberbullying detection. However, it is important to note that machine learning models are not infallible and may produce false positives or false negatives. Therefore, it is important to carefully evaluate the performance of these models and to incorporate human judgment and oversight in the cyberbullying detection process.

Systemic Functional Linguistics (SFL) is a linguistic theory that views language as a social semiotic system, in which meaning is created through the interaction of language with social and cultural context. It was developed by Michael Halliday in the 1960s and has since been expanded and refined by various linguists around the world. According to SFL, language has three main functions: ideational, interpersonal, and textual [4]. The ideational function refers to the representation of experience and knowledge, the interpersonal function refers to the expression of social relationships and identities, and the textual function refers to the organization of language into cohesive and coherent discourse. SFL also emphasizes the importance of context in shaping the meaning of language. Context includes not only the immediate situation in which language is used, but also the broader social and cultural context in which communication takes place. SFL analyzes language use at different levels of context, including the field (the subject matter), the tenor (the participants and their relationships), and the mode (the channel of communication). One of the key contributions of SFL is its approach to grammar, which views grammar as a resource for meaning-making. SFL analyzes grammar in terms of systems of choices that speakers make to create meaning in different contexts. These choices include things like tense, mood, voice, and modality, which are used to express different kinds of meanings and functions. SFL has been applied to a wide range of linguistic and educational contexts, including discourse analysis, language teaching, and language policy. [5] It has also been used to analyze a variety of languages, including English, Chinese, Spanish, and Arabic. In general, SFL provides a powerful framework for understanding the complex interplay between language, context, and meaning. Its emphasis on the social and cultural dimensions of language use makes it a valuable tool for analyzing and understanding communication in a variety of settings.

The Main contribution of this paper is to propose an SFL-incorporated supervised machine learning approaches to detect and limit the act of cyberbullying. This research aims to extract features that represent the ideational metafunction, since the latter is concerned about the topic of the given text. With the use of the extracted features, an

integration of Bag of Words (BoW) and ideational metafunction is made and inserted into ML models. Finally, the research compares the performance of three common classifiers in NLP after considering the SFL theory. our study contributes to the growing body of research on cyberbullying and provides a novel approach to combat cyberbullying racism using machine learning and systemic functional linguistics theory.

The remaining sections of the paper are outlined as follows: Section 2 presents the related work. Section 3 explains the theory and implementations of the proposed model. Section 4, interprets the results and discussion. Section 5 presents the conclusion and some suggestions for future work.

2. Related Work

There are several approaches that proposed systems which are automatically able to detect cyberbullying with adequate accuracy. No empirical study exists for the integration of machine learning and the theory of systemic functional linguistics within the context of cyberbullying; however, a limited number of studies have proposed the combination of machine learning algorithms and systemic functional linguistic theory in alternative contexts. In his research article, Wei Dong [6] explores the application of text mining techniques, specifically within the context of SFL theory, to detect financial statement fraud. The primary focus is on analyzing textual information in financial reports and statements to identify potential fraudulent activities or misrepresentations This research aims to harness the potential of the vast amount of often overlooked textual content within financial statements. Guided by the theoretical principles of SFL, a systematic framework is developed for analyzing textual data to detect financial statement fraud. Seven categories of information, encompassing topics, sentiments, emotions, modality, personal pronouns, writing style, and genres, are identified based on the ideational, interpersonal, and textual functions outlined in SFL. Within this analytical framework, we have integrated various methods, including the Latent Dirichlet Allocation (LDA) algorithm, computational linguistics, and the term frequency-inverse document frequency method, to synergistically extract both word-level and document-level features. These extracted features are subsequently used as input for the Liblinear Support Vector Machine classifier.

Kappagoda [7] conducted a study where they integrated SFL and text mining for military applications. In their research, they demonstrated the augmentation of words with Word-Function in the Group (WFG) tags using Conditional Random Fields (CRF) and the addition of Part-Of-Speech (POS) tags to words. These word-function tags provide insights into the systemic functional roles that words fulfil within a text. It's worth noting that this approach lacks a hierarchical structure. Nonetheless, their

innovative work successfully showcased that the labelling process can be partially automated and that the resultant tags contribute to enhanced document comprehension. The study focused on analyzing the Enron dataset, comprising over 60,000 emails, with the objective of classifying them into "official" and "non-official" categories.

Zhang [8], proposed a method that automatically detects cyberbullying text from Twitter. Using multiple machine learning classifiers on Japanese dataset containing around 2M tweets. Multiple textual features and multiple machine learning algorithms were used to construct classification models. With the experiments based on the collected tweets, the quality of automatic cyberbullying detection is evaluated and the best model performs over 90% for the four criteria: accuracy, precision, recall, and F-measure. (Logistic regression or Gradient boosting regression tree) the classification quality is unclear if a cyberbullying text only contains bullying words which do not appear in the collected learning data or even has no bullying words such as sarcastic expressions.

Sarcasm and concealed meanings within text frequently pose challenges in accurate text classification. Vosoughi's [9] research introduced a pivotal development—a sentiment classifier incorporating a selection of unigram and bigram features. Their study presented a supervised Twitter speech act classifier, which introduced a taxonomy consisting of six speech act categories: Recommendation, Assertion, Expression, Request, Question, and Miscellaneous. To enhance classification accuracy, the researchers harnessed various linguistic features, including syntactic elements such as punctuation usage, Twitter-specific characters, abbreviations, dependency subtrees, part-of-speech (POS) patterns, and semantic attributes like opinionated terms, explicit language, emoticons, speech-act verbs, and n-grams. The dataset employed in this investigation was sourced from the Twitter public API, comprising several thousand tweets, each annotated with one of the aforementioned speech act categories. To train and evaluate the classifier's performance, five different classification algorithms were employed: Naïve Bayes, Logistic Regression, Support Vector Machine, Decision Tree, and the Baseline max classifier. Notably, Logistic Regression emerged as the most effective classifier, achieving a weighted average F1 score of 70%.

Dalvi and colleagues [10] introduced an innovative methodology designed for the identification and prevention of Twitter cyberbullying through the implementation of supervised binary classification machine learning algorithms. The assessment of their model encompassed an evaluation utilizing Support Vector Machine and Naïve Bayes classifiers. Additionally, for the extraction of relevant features, they harnessed the TF-IDF vectorizer. The outcomes of their study demonstrate the model's proficiency in detecting cyberbullying content, particularly notable for the Support Vector Machine, where an

impressive accuracy rate of approximately 71.25% was achieved, surpassing that of the Naïve Bayes approach. In a related study, Dinakar and co-authors [11] aimed to detect explicit instances of cyberbullying language associated with topics including sexuality, race, culture, and mental capacities. Their dataset was meticulously collected from the comment section of YouTube. Employing Support Vector Machine and Naïve Bayes classifiers, the Support Vector Machine classifier exhibited a noteworthy accuracy rate of 66%.

3. Methodology

This section serves as the architectural framework of our research, outlining the systematic approach employed to investigate and address the core objectives of this study. In this section, we elucidate the rigorous methods and procedures undertaken to ensure the reliability, validity, and comprehensiveness of our findings.

In the phase of data collection, the proposed approach is evaluated on a cyberbullying dataset from Mendely website which was collected and labeled by the author Elsafouri [12], it is publicly available from Twitter social platform, and is used for the purpose of automatic detection of cyberbullying. The data contain different types of cyberbullying such as hate speech, aggression, racism, insults and toxicity. NLP aims to interpret human language with the help of computers. To achieve this, the first step is to convert textual data into a numerical representation suited for machine learning models, after being cleaned from any unwanted characters, such as punctuation marks, extra white spaces, empty texts (null values) ... etc.

The most widely used numerical representation for textual data are: BoW, Term frequency inverse document frequency (TF-IDF). As in this proposed work, LDA is used as the main feature to represent the ideational metafunction. LDA is a statistical model that builds a topic per document model and words per topic model, modeled as Dirichlet distributions. Later in the model building phase Naïve Bayes (NB), Support Vector Machine (SVM), and Random Forest (RT) are used in this research. The evaluation is done using the machine learning performance metrics precision, recall, and F-1 score. Figure 1 below presents the main phases pipeline of cyberbullying classification for Twitter texts (tweets), which contains five main phases.

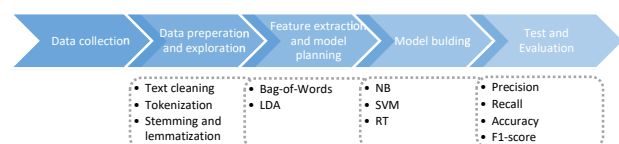


Fig. 1 Proposed Approach.

During the data preprocessing phase, the noise and unnecessary text are eliminated, streamlining the dataset.

The preprocessing procedure encompasses the following steps:

- (i) Text cleaning; this step involves removing unwanted elements from the raw text such as HTML tags, punctuation, numbers, and stop words. Stop words are the commonly used words like "the", "a", and "an" which do not carry significant meaning.
- (ii) Tokenization; this step involves breaking the cleaned text into individual words or phrases, called tokens, that can be analyzed further. This can be done using various techniques whitespace-based tokenization or regular expression-based tokenization.
- (iii) Stemming and lemmatization; are two common techniques used in NLP to reduce words to their base or root form. Stemming is the process of reducing words to their stem or root form by removing suffixes. For example, the stem of the words "jumping", "jumps", and "jumped" is "jump". While lemmatization is the process of reducing words to their base or dictionary form, called a lemma, by taking into account the word's context and part of speech. For example, the lemma of the words "jumping", "jumps", and "jumped" is "jump", and the lemma of the word "was" is "be". Lemmatization can be more accurate than stemming since it considers the context and part of speech of the word. However, it can also be more computationally expensive and complex than stemming, since it requires access to a dictionary or knowledge base of the language. In general, stemming is faster and simpler, while lemmatization is more accurate but more complex.

The next phase is features extraction, In the feature extraction phase, we employed LDA to transform the raw text data into a structured and meaningful representation. LDA is a powerful probabilistic model that allows us to identify underlying topics within a corpus of text. First, the text data underwent pre-processing, including tasks such as tokenization, stop-word removal, and stemming to enhance the quality of the input. Subsequently, the LDA algorithm was applied to generate a Bag of Words representation, wherein each document was represented as a vector of topic probabilities. These topic probabilities encapsulate the distribution of latent topics within the text, enabling a compact and informative representation of the original content. The utilization of LDA not only facilitated dimensionality reduction but also provided insights into the underlying thematic structure of the text data. This feature extraction step served as a pivotal foundation for subsequent analysis, enabling us to uncover hidden patterns and relationships within the dataset."

After that comes the next phase which is the model building. After the successful extraction of BoW and LDA features from the textual data, our next pivotal step involved the construction of machine learning models to tackle the classification task. Python libraries, including scikit-learn and NLTK, were harnessed to facilitate this process. We employed a diverse ensemble of classifiers, each renowned for its distinct strengths and characteristics. The SVM classifier, implemented using scikit-learn, exhibited its prowess in handling complex data structures and non-linear relationships. NB, another reliable scikit-learn classifier, was selected for its simplicity and efficiency in text classification tasks. Furthermore, the Random Forest classifier, a robust ensemble model, was leveraged to harness the power of decision trees for accurate predictions. These classifiers were meticulously trained on the feature-rich BoW and LDA representations of our textual data. The Python libraries offered a seamless and efficient environment for model development, parameter tuning, and cross-validation, ensuring that our machine learning models were primed for the subsequent evaluation and validation phases of this study.

In the final phase of our research, we meticulously evaluated the performance of our constructed machine learning models to ascertain their efficacy in the context of cyberbullying detection. To gauge the models' effectiveness, we relied on a comprehensive set of performance metrics. These metrics included accuracy, precision, recall, and the F1-score. Accuracy provided an overall measure of the model's correctness in classifying cyberbullying instances, while precision quantified the proportion of true positives among all positive predictions, illuminating the model's ability to make accurate positive identifications. Recall, on the other hand, measured the proportion of true positives detected among all actual positive instances, highlighting the model's ability to capture genuine cyberbullying content. Lastly, the F1-score, which strikes a balance between precision and recall, served as a consolidated measure of the model's overall effectiveness. In the typical assessment of classifiers, multiple evaluation metrics are employed, contingent upon the information provided by the confusion matrix. Within this array of criteria, one finds metrics such as accuracy, precision, recall, and the F1-score. These metrics are computed through the following equations:

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad (1)$$

$$Precision = \frac{TP}{TP+FP} \quad (2)$$

$$Recall = \frac{TP}{TP+FN} \quad (3)$$

$$F1 - score = \frac{2 \times (precision \times recall)}{precision + recall} \quad (4)$$

Where 'TP' denotes the count of true positives, 'TN' signifies the count of true negatives, 'FP' represents the count of false positives, and 'FN' indicates the count of false negatives.

4. Experimental Results and Analysis

In this study, we investigated the effectiveness of LDA features for cyberbullying detection. We evaluated three different classifiers: Naïve Bayes, SVM, and Random Forest. The performance of each classifier was measured in terms of train and test accuracy, precision, recall, and f1-score. We used a publicly available dataset of social media posts that were labeled as either racism or non-racism. The dataset was preprocessed by removing stop words, stemming the remaining words, and converting the text into LDA features. Prior to the extraction of LDA features as a representation of the SFL Ideational metafunction, an initial phase involved the application of Python-based machine learning classifiers in conjunction with BoW techniques to categorize the data into two distinct classes: racism and non-racism. This preliminary step was undertaken to facilitate a comparative analysis between the outcomes achieved when employing SFL-based features and those attained in the absence of such linguistic features. We randomly split the dataset into a training set (80%) and a test set (20%). We trained each classifier on the training set and evaluated its performance on the test set. We used the scikit-learn library to implement the classifiers.

As shown in Table 1, the SVM classifier achieved the highest test accuracy (92.5%) followed by Random Forest (91%) and Naïve Bayes (82.75%). All three classifiers achieved high precision, recall, and f1-score, indicating that the application of these classification models enriched with the topic-related information yielded acceptable results in detecting cyberbullying.

Table 1: Performance of the classifiers

Classifier	Performance Metrics				
	Train accuracy	Test accuracy	Precision	Recall	F1-score
NB	92.20%	82.75%	91%	83%	85%
SVM	93%	92.5%	92%	93%	92%
RF	93%	91%	93%	93%	93%

The outcome of this study, rooted in the utilization of LDA to represent the Ideational metafunction and machine learning classifiers for cyberbullying tweet classification,

yielded valuable insights into the effectiveness of this approach for topical categorization of social media content. The analysis of results encompassed several facets, including cross-validation scores, training accuracy, and test accuracy, providing a comprehensive assessment of the model's performance. During the training phase, the NB model achieved an impressive accuracy rate of 92.20%. This high training accuracy highlights the model's capacity to effectively learn and categorize tweets based on the topics derived from LDA. It also suggests that the models successfully captured the underlying patterns within the training data, indicating a robust grasp of the inherent relationships between words and topics. In evaluating the model's performance on previously unseen tweets, the test accuracy was calculated at 82.75% for NB model. This test accuracy score demonstrates the model's proficiency in classifying cyberbullying tweets into their respective topics based on the information derived from the LDA topic modeling process. While the test accuracy may exhibit a slight reduction compared to the training accuracy, it remains commendably high, affirming the model's capability to generalize effectively to new and unseen tweet data. Figure 2 shows the performance metrics when LDA isn't employed as features into the classifiers.



Fig. 2 Performance of the three classifiers without LDA features.

Figure 3 visually presents the performance metrics of NB, SVM and RF classifiers employing LDA features.

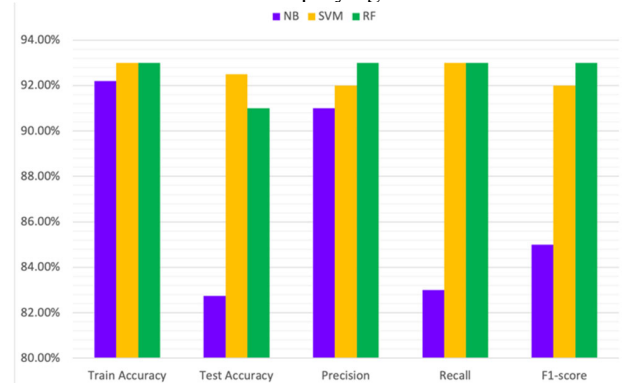


Fig. 3 Performance of the three classifiers with LDA features.

The change in performance metrics between the two sets of results, one without using LDA features and the other with the application of LDA features, likely stems from the incorporation of LDA features into the cyberbullying detection models. LDA features capture latent topics within the text data. By incorporating LDA features, we introduce additional information about the content of the text, which may enhance the model's ability to discriminate between cyberbullying and non-cyberbullying instances. LDA features can potentially highlight patterns and topic distributions that are more indicative of cyberbullying behavior. This can lead to improved precision, recall, and F1-score for both classes. The effectiveness of LDA features can also depend on the specific characteristics of your dataset. If the dataset contains textual patterns that align well with the topics extracted by LDA, you are likely to see a positive impact on performance. It's crucial to note that the extent of improvement may vary based on the quality of the LDA features, the complexity of the dataset, and other factors. The observed changes suggest that incorporating LDA features has enhanced the model's ability to classify cyberbullying and non-cyberbullying content, resulting in improved accuracy, precision, recall, and F1-score.

5. Conclusion

In this research endeavor, we embarked on a mission to address a pressing societal issue - the detection of cyberbullying, with a specific focus on identifying instances of racism in Twitter data. Our approach amalgamated the formidable capabilities of three prominent machine learning algorithms—SVM, NB, and RF—with the insightful lens of SFL, focusing on the Ideational metafunction, to comprehensively analyze and classify textual data. Our research harnessed the power of machine learning to harness patterns and nuances within Twitter data while simultaneously incorporating the linguistic insights derived from the Ideational metafunction of SFL. This amalgamation aimed to deepen our understanding of the linguistic structures within the text and subsequently employ this knowledge to differentiate between instances of cyberbullying with racial connotations and non-cyberbullying content. A pivotal step in our methodology involved the extraction of features using LDA, an innovative technique for uncovering latent topics within text. By utilizing LDA, we aimed to represent the Ideational metafunction, thereby enhancing the granularity of our analysis and facilitating the identification of pertinent themes and topics. The results of our research were highly promising. Each of the machine learning models—NB, RF, and SVM—exhibited a remarkable ability to detect instances of cyberbullying, particularly those associated with racism. Notably, Naïve Bayes achieved an impressive F1-score of 85%, while Random Forest and Support Vector Machine surpassed expectations with F1-scores of 93% and

92%, respectively. These outcomes underscore the robustness and efficacy of our approach in addressing the challenging task of cyberbullying detection within the specific context of racism. In conclusion, our research underscores the potential of integrating machine learning with linguistic theory to tackle complex societal issues. Our machine learning models, enriched by linguistic insights, demonstrated commendable performance in the detection of racism-related cyberbullying. As we move forward, the pursuit of comprehensive and ethically sound solutions to combat online harassment remains paramount, and we hope that our research serves as a catalyst for further advancements in this critical domain.

Acknowledgment

The authors would like to express their cordial thanks to Dr. Sawsan Aljahdali for her valuable advice.

References

- [1] Kowalski, R. M., Giunetti, G. W., Schroeder, A. N., & Lattanner, M. R. (2014). Bullying in the digital age: A critical review and meta-analysis of cyberbullying research among youth. *Psychological bulletin*, 140(4), 1073. Chicago
- [2] Reynolds, Kelly & Edwards, April & Edwards, Lynne. (2011). Using Machine Learning to Detect Cyberbullying. Proceedings - 10th International Conference on Machine Learning and Applications, ICMLA 2011. 2. 10.1109/ICMLA.2011.152.
- [3] Raj, M., Singh, S., Solanki, K., & Selvanambi, R. (2022). An Application to Detect Cyberbullying Using Machine Learning and Deep Learning Techniques. *SN computer science*, 3(5), 401.
- [4] Wegener, R., Cassens, J., & Butt, D. (2008). Start Making Sense. *Systemic-Functional Linguistics and Ambient Intelligence*. *Revue d'intelligence artificielle*, 22(5), 629-645.
- [5] Almurashi, W. A. (2016). An introduction to Halliday's systemic functional linguistics. *Journal for the study of English Linguistics*, 4(1), 70-80
- [6] Dong, W., Liao, S., & Liang, L. (2016). Financial statement fraud detection using text mining: A systemic functional linguistics theory perspective. In *Pacific Asia Conference on Information Systems (PACIS)*. Association For Information System.
- [7] Kappagoda, A. (2009). The use of systemic-functional linguistics in automated text mining. DEFENCE SCIENCE AND TECHNOLOGY ORGANISATION EDINBURGH (AUSTRALIA) COMMAND CONTROL COMMUNICATIONS AND INTELLIGENCE DIV.
- [8] Zhang, J., Otomo, T., Li, L., & Nakajima, S. (2019, October). Cyberbullying Detection on Twitter using Multiple Textual Features. In *2019 IEEE 10th International Conference on Awareness Science and Technology (iCAST)* (pp. 1-6). IEEE.
- [9] Vosoughi, S., & Roy, D. (2016). Tweet acts: A speech act classifier for twitter. *arXiv preprint arXiv:1605.05156*.
- [10] Dalvi, R. R., Chavan, S. B., & Halbe, A. (2020, May). Detecting A Twitter Cyberbullying Using Machine Learning. In *2020 4th International Conference on Intelligent*

- Computing and Control Systems (ICICCS) (pp. 297-301).
IEEE.
- [11] Dinakar, K., Jones, B., Havasi, C., Lieberman, H., & Picard, R. (2012). Common sense reasoning for detection, prevention, and mitigation of cyberbullying. *ACM Transactions on Interactive Intelligent Systems (TiiS)*, 2(3), 1-30.
- [12] Elsafoury, Fatma (2020), "Cyberbullying datasets", Mendeley Data, V1, doi: 10.17632/jf4pzyvnpj.1