

## Research Article

# Nature-Inspired-Based Approach for Automated Cyberbullying Classification on Multimedia Social Networking

N. Yuvaraj,<sup>1</sup> K. Srihari,<sup>2</sup> Gaurav Dhiman ,<sup>3</sup> K. Somasundaram,<sup>4</sup> Ashutosh Sharma ,<sup>5</sup> S. Rajeskannan,<sup>4</sup> Mukesh Soni,<sup>6</sup> Gurjot Singh Gaba ,<sup>7</sup> Mohammed A. AlZain ,<sup>8</sup> and Mehedi Masud <sup>9</sup>

<sup>1</sup>Training and Research, ICT Academy, Chennai, India

<sup>2</sup>Department of Computer Science and Engineering, SNS College of Technology, Coimbatore, India

<sup>3</sup>Department of Computer Science, Government Bikram College of Commerce, Patiala-147001, Punjab, India

<sup>4</sup>Dept of Computer Science and Engineering, Chennai Institute of Technology, Chennai, India

<sup>5</sup>Southern Federal University, Rostov-on-Don, Russia

<sup>6</sup>Dept of Computer Science and Engineering, Jagran Lakecity University, Bhopal, India

<sup>7</sup>School of Electronics and Electrical Engineering, Lovely Professional University, Phagwara 144411, India

<sup>8</sup>Department of Information Technology, College of Computers and Information Technology, Taif University, P. O. Box 11099, Taif 21944, Saudi Arabia

<sup>9</sup>Department of Computer Science, College of Computers and Information Technology, Taif University, P. O. Box 11099, Taif 21944, Saudi Arabia

Correspondence should be addressed to Mehedi Masud; mmasud@tu.edu.sa

Received 22 December 2020; Revised 23 January 2021; Accepted 1 February 2021; Published 23 February 2021

Academic Editor: Erik Cuevas

Copyright © 2021 N. Yuvaraj et al. This is an open access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

In the modern era, the cyberbullying (CB) is an intentional and aggressive action of an individual or a group against a victim via electronic media. The consequence of CB is increasing alarmingly, affecting the victim either physically or psychologically. This allows the use of automated detection tools, but research on such automated tools is limited due to poor datasets or elimination of wide features during the CB detection. In this paper, an integrated model is proposed that combines both the feature extraction engine and classification engine from the input raw text datasets from a social media engine. The feature extraction engine extracts the psychological features, user comments, and the context into consideration for CB detection. The classification engine using artificial neural network (ANN) classifies the results, and it is provided with an evaluation system that either rewards or penalizes the classified output. The evaluation is carried out using Deep Reinforcement Learning (DRL) that improves the performance of classification. The simulation is carried out to validate the efficacy of the ANN-DRL model against various metrics that include accuracy, precision, recall, and f-measure. The results of the simulation show that the ANN-DRL has higher classification results than conventional machine learning classifiers.

## 1. Introduction

Cyberbullying (CB) is considered as a new or electronic form of traditional bullying [1]. CB is defined as a repetitive, intentional, and aggressive reaction committed by a group or an individual against another group or an individual, which is made by the utilization of Information Communication Technology (ICT) tools such as social media, Internet, and

mobile phones [2]. The entire CB incidents are carried out virtually in Internet media rather than in physical form. The CB consists of hatred messages transmitted via social networking, e-mails, etc. through personal or public computers or through personal mobile phones. This has aroused as a serious threat among nations [1]. Various privacy-preserving tools are adopted in the Internet arena to protect the data; however, most mechanisms are challenged by the process of

traffic classification [3], which is a vital workhorse for network management, where it becomes a key factor in assigning the privacy level to classify malign and benign standpoints [4]. This is true in case of testing the methods with a selected dataset on the Dark Web Forum Portal [5]. The CB consists of hatred messages transmitted via social networking, e-mails, etc. through personal or public computers or through personal mobile phones. This has aroused as a serious threat among nations [1].

The research on previous studies considers CB as a distinct variant from the traditional bullying [2, 6]. The suggested variances between the CB and traditional bullying reveal the inadequacy of CB findings from conventional bullying [7]. Evidences found in [8] reveal that there exist several features of CB that vary between its prevalence rates, protective and risk factors, risk outcomes, and strategies adopted for its prevention. The CB features are partially related and partially distinct with conventional bullying [2, 9]. CB, on other hand, has impacted the victims psychologically and physically with its increasing prevalence, where most vulnerability is reported among youths [2].

Hence, it is vital to detect the CB context and its applications to reduce the vulnerability. However, from the view of the cyber world, the application involving CB involves difficulties associated with ignorance of aggressors and their identity, lack of direct communication, and relating consequences over others [10–15].

The failure to direct communication causes partial interpretation of the significance or the nature of the message, and it leads to confusion over the individual's intentionality with exchange or interaction messages. In spite of the problems while identifying the behavioral intent of an individual, the major factor that creates transition from aggression to CB is the intention of harming oneself [16].

In the current scenario, an automated behavior of social network platforms alerts the moderators to review the reported CB contents. However, most of the frameworks lacks an automated intelligent system that alerts the moderators and detects the contents in an automated way faster than the traditional reporting system. This enables the moderator to respond on the alert and take required action on reporting the user or removing the content [17].

The major constraint existing on existing detection systems with CB research is the lack of input data. The existing research is carried out conventionally on available datasets or the surveyed data, where the perpetrators or the victims are allowed to report the impressions [18]. The other issue associated with automated CB detection is the proper operationalization on CB contents that considers only the available literatures in the CB detection field for achieving the aim of automated detection to accurately identify the events of CB. The other issue associated with automated CB detection is the proper operationalization on CB contents that considers only the available literature in the CB detection field for achieving the aim of automated detection to accurately identify the events of CB. This creates complexity in identifying the events, and hence, well-developed tools are essential in integrating the features with an automated decision model [19].

Various research studies on automated cyberbullying detection with intelligent systems are reported in [8, 18, 20–28]. These studies utilized machine learning algorithms for automated detection of CB contents utilizing several common and psychological features. These intelligent systems on CB detection are reported to be low, and it is principally limited with the comment of an individual leaving the context. An existing study has reported utilization of the user context in action that involves the characteristics of users and history of user comments to improve the performance of CB detection/classification [17].

In this paper, we utilize an integrated feature model that collects and trains the system with taking psychological features, user comments, and the context into consideration for CB detection. A classification engine using an artificial neural network (ANN) as impacted from [22] enables CB classification, and the operation on each classification is monitored by the reward-penalty model of a Deep Reinforcement Learning (DRL) engine.

The study contributes to the following in the field of CB detection:

- (a) The authors develop a series of frameworks that extracts the CB contexts from raw input messages. The study considers utilizing wide varied features to train the feature extraction module, and this involves the psychological traits, user comments, and context.
- (b) The authors develop an integrated classification engine that combines an ANN with DRL to classify the CB contents and improve the results after each iteration based on the feedback obtained from the DRL mechanism. Here, the entire classification is carried out by the ANN algorithm, and the DRL provides state-action-reward for each classified results.

The outline of the study is given as follows: Section 2 discusses the related works. Section 3 provides the proposed classification engine. Section 4 evaluates the entire work. Section 5 concludes the work with possible directions of future scope;

CB	Cyberbullying
ANN	Artificial Neural Network
DRL	Deep Reinforcement Learning
SVM	Support Vector Machine
NB	Naïve Bayes
KNN	k-nearest neighbor
RF	Random Forest
LR	Logistic Regression

## 2. Related Works

Nandhini and Sheeba [20] presented a detection technique to combat CB on social media. The study extracts features such as the noun, pronoun, and adjective obtained from the text and frequency of words occurrences. These features are used to classify various activities such as Harassment, Flaming, Terrorism, and Racism using a Fuzzy logic-based

genetic algorithm. The relevant data are retrieved using the Fuzzy rule for classification, and the genetic algorithm increases the accuracy of classification by parametric optimization.

Potha et al. [21] employed a Support Vector Machine (SVM) classifier to classify the CB based on various features such as local, sentimental, contextual, and gender-specific language features. The SVM classifier combined with a tf-idf measure and linear kernel identifies the online harassment.

Kumar and Sachdeva [28] reviewed various studies and found both direct and indirect CB features have higher impacts on machine learning classification. The results of classification show that the SVM classifier has higher classification rate than other supervised/unsupervised learning methods.

Al-garadi et al. [8] used the SVM [21, 28], naïve Bayes (NB) [25, 28], k-nearest neighbor (KNN), and random forest (RF) [25] classifier with various features extracted from the Twitter data that include network, activity and user information, and tweet content. The features are selected using the information gain, c2 test, and Pearson correlation. Furthermore, the classified results are optimized using a synthetic minority oversampling approach, and classes are balanced with weight adjustment in the dataset. The result shows that the RF has higher classification accuracy.

Balakrishnan et al. [25] developed an automated detection model with Big Five and Dark Triad models for user personality determination. The classification is carried out with various machine learning classifiers, NB, RF, and J48, to detect bully, spammer, aggressor, and normal. The psychological features are selected from the twitter data for better tweet classification. The study confirmed that the user personalities on classification have higher impacts on detection than other traits.

Murnion et al. [18] developed an Artificial Intelligence-based CB detection from an automated data collection system from the chat data of online multiplayer games. The sentiment text analytics system is supported with a scoring scheme for optimal classification. The study is assigned with eight descriptive attributes including IsAbusive, IsPositive, IsNegative, HasBadLanguage, IsRacist, NoobRelated, SpecificTarget, and FilteredText for potential identification of CB. The estimation of the CB score found that the both Twinword- and Microsoft-aware sentimental analysis were poor with less classification score.

Ho et al. [27] used 90 features categorized into 10 classes and utilized it for classification using a logistic regression model. The detection is improved by training the model with 14 abusive words for reducing the false classification rate.

Balakrishnan et al. [24] used an RF classifier with multiple decision trees, where classification is finally determined based on majority of votes. The study selects 15 twitter features [23] using Big Five and Dark Triad models to find the user personalities.

Sánchez-Medina et al. [26] used ensemble classification trees with Dark Triad for identifying the personality trait. The study used psychopathy, narcissism, and abusive words and then n-grams, blacklists, and edit-distance metrics for the detection of obfuscated words. A three-layered neural

network model is used finally for classification, which acts as an unsupervised learning model. The misclassification is reduced by employing a 1.5 million nonabusive words dataset which improves the classification using neural network.

The abovementioned research used minimal features to classify the datasets, and furthermore, the CB word is treated as the seed word for DB detection. However, the CB word is a distinctive vocabulary that fails to cover all cases.

Machiavellianism for potentially detecting the CB sexual assaults in social media: Lee et al. [22] used an embedded vector representation such as skip-gram word2vec that represents the words as vectors. The cosine similarity detects the new one.

Balakrishnan et al. [24] used an RF classifier with multiple decision trees, where classification is finally determined based on majority of votes. The study selects 15 twitter features [23] using Big Five and Dark Triad models to find the user personalities.

Sánchez-Medina, et al. [26] used ensemble classification trees with Dark Triad for identifying the personality trait. The study used psychopathy, narcissism, and machiavellianism for potentially detecting the CB sexual assaults in social media.

Lee et al. [22] used an embedded vector representation such as skip-gram word2vec that represents the words as vectors. The cosine similarity detects the new abusive words and then n-grams, blacklists, and edit-distance metrics for the detection of obfuscated words. A three-layered neural network model is used finally for classification, which acts as an unsupervised learning model. The misclassification is reduced by employing a 1.5 million nonabusive words dataset which improves the classification using neural network.

The abovementioned research used minimal features to classify the datasets, and furthermore, the CB word is treated as the seed word for DB detection. However, the CB word is a distinctive vocabulary that fails to cover all cases.

### 3. Proposed Method

In the present research, the entire focus is not on a specific CB word, but the vulgarity is determined based on weight score calculation and harmfulness index estimation for the entire word sequence (optimal words chosen by the feature selection method) of the collected tweets. This reduces well the cost of training data construction and further with the dependency between the phrases. The architecture of the proposed classification model is given in Figure 1.

We consider an annotated dataset  $D = \{(x_i, \sim c_i)\}$ , where  $x_i$  are the twitter CB datasets and without label  $\sim c_i$ . The datasets are divided into smaller subset  $L \subset D$ . The aim is to detect the CB instances from the twitter data that may vary from long to short paragraphs.

**3.1. Preprocessing.** The preprocessing method uses a lexical normalization method [29] that uses various components to clean the input tweet data. It further converts the numerical variables into an equivalent text data. The spell corrector component helps to reduce the outbound vocabulary terms,

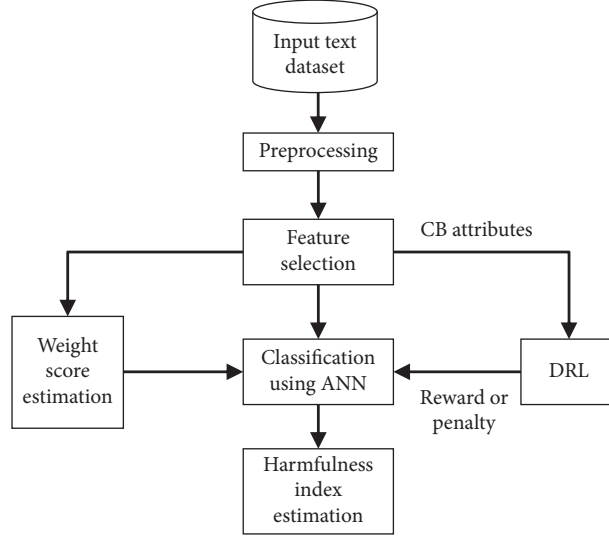


FIGURE 1: Overview of the proposed system.

TABLE 1: Selected attributes to classify the Tweets.

Attributes	Class	Format
Noun	CB/non-CB	Text
Pronoun	CB/non-CB	Text
Adjective	CB/non-CB	Text
Local features	The basic features extracted from a tweet	Text
Contextual features	Professional, religious, family, legal, and financial factors specific to CB	Text
Sentiment features	Positive or negative (foul words specific to CB) or direct or indirect CB	Text
Emotion features	Polite words, modal words, unknown words, number of insults and hateful blacklisted words, harming with detailed description, power differential, any form of aggression, targeting a person, targeting two or more persons, intent, repetition, one-time CB, harm, perception, reasonable person/witness, and racist sentiments	Text
Gender-specific language	Male/female	Text
User feature	Network information, user information, his/her activity information, tweet content, account creation time, and verified account time	Text/numeric
Twitter basic features	Number of followers, number of mentions, and number of following, favorite count, popularity, number of hash tags, and status count	Numeric
Linguistic features	Other languages words, punctuation marks, and abbreviated words rather than abusive sentence judgments	Text

and in prior, the entire redundant or missing variables are cleaned that involve spelling errors, wrong punctuations, etc.

**3.2. Feature Selection.** The selection of features (given in Table 1) from the input twitter datasets involves three different methods including Information Gain [30], chi-square  $\chi^2$  [31], and Pearson correlation [32]. These methods are employed to select the features from the preprocessed datasets.

**3.2.1. Information Gain.** Decision tree algorithm is utilized to implement the feature extraction using information gain. The information gain is defined as the measure of entropy that is used widely in the machine learning domain. It acts as a statistical method that assigns the weights of features based on the correlation between the categories and the features.

We consider a dataset  $S(s_1, s_2, \dots, s_n)$ , which is regarded as the collection of varying instances, say  $n$  s. t.  $A(A_1, A_2, \dots, A_p)$  is the attributes set for  $p$ , where  $C(c_1, c_2, \dots, c_m)$  is regarded as the collection of different label categories  $m$ .  $p(c_i)$  represents the  $i^{\text{th}}$ -class label proportion with  $c_i$  ( $i = 1, 2, \dots, m$ ) in  $S$ . The dataset entropy is, thus, represented as

$$H(C) = - \sum_{i=1}^m p(c_i) \log_2(p(c_i)). \quad (1)$$

The information gain on each feature is defined used for classification of input data, where  $A_q(a_{q1}, a_{q2}, \dots, a_{qk})$  represents the  $q^{\text{th}}$  attribute ( $q = 1, 2, \dots, p$ ). The conditional entropy for an attribute  $A_q(a_{q1}, a_{q2}, \dots, a_{qk})$  is, thus, represented as

$$H(C|A_q) = - \sum_{j=1}^k p(a_{qj}) \sum_{i=1}^m p(c_i|a_{qj}) \log_2(p(c_i|a_{qj})), \quad (2)$$

where  $a_{qj}$ - $A_q$  is the attribute value with a  $k$  value,  $p(a_{qj})$  is the probability of categorical variable  $C$ , and  $p(c_i|a_{qj})$  is the conditional probability of  $C$  after the value of  $A_q$  is fixed.

Then, information gain is estimated as the difference between the value  $H(C)$  and  $H(C|A_q)$ , and this offers the attribute value  $A_q$  as stated below:

$$IG(A_q) = H(C) - H(C|A_q). \quad (3)$$

Usually, the higher the information gain is, the more vital the feature is then considered for classification.

If the value of information gain is high, the feature is considered to be vital for the purpose of classification.

**3.2.2. Chi-Square  $\chi^2$ .** The chi-square statistics is used in feature extraction as an information theory function that helps in extraction of elements, say  $t_k$  over a class  $c_i$ . These elements are considered to be distributed widely and differently in sets of negative and positive examples of  $c_i$ .

$$\chi^2(t(k), c(i)) = \frac{N(A D - CB)^2}{(A + C)(B + D)(A + B)(C + D)}, \quad (4)$$

where  $N$ - total documents;  $A$ - total documents in  $c_i$  containing  $t_k$ ;  $B$ - total documents containing  $t_k$  other than  $c_i$ ;  $C$ - total documents in  $c_i$  without  $t_k$ ; and  $D$ - total documents without  $t_k$  other than  $c_i$ .

The next step is the assignment of scores for each  $c_i$  as discussed in the abovementioned equation, and the collective scores are summed into a single final score. The final score helps in classification of attributes, and the top score is selected.

**3.2.3. Pearson Correlation.** The Pearson correlation coefficient in the present study is used for the estimation of optimal features by calculating the degree of linear correlation between the extracted class and original class.

$$\text{sim}_i = \frac{\sum_{j=1}^N (X_j - \bar{X})(Y_{ij} - \bar{Y})}{\left[ \sum_{j=1}^N (X_j - \bar{X})^2 \right] \left[ \sum_{j=1}^N (Y_{ij} - \bar{Y})^2 \right]}, \quad (5)$$

where  $\text{sim}_i$ - similarity between the  $i^{\text{th}}$  class and original class of a dataset;  $X_j$  and  $Y_{ij}$ - selected attribute data to be tested on the  $i^{\text{th}}$  class,  $\bar{X}$ - and  $\bar{Y}$ -average value of selected attribute data, and with the original class of a dataset, and finally, the entire attribute data are normalized.

**3.3. ANN.** Artificial neural networks [33] are trained with weights of input features as in Figure 2(a), and furthermore, it is trained by proper reduction of an error function. The selection of a reduced error function helps in classification in terms of reduced cross-entropy error as follows:

$$E = \sum_{i=1}^n y \log o_N + (1 - y) \log (1 - o_N). \quad (6)$$

The size of the input twitter dataset  $D$ , for an ANN classification model  $P(y|x)$  is influenced by the selection of

CB from  $D$ . The challenge of model building is to summarize the underlying distribution from the specific instance  $D$  of the samples. The problem with the memory of the dataset is known as overfitting rather than identifying the dataset distribution.

An activation feature is considered as a real function that determines the value of the neuron returned. The present study uses inverse trigonometric functions as the activation function.

Multilayer perceptron is the most frequent architecture of a feedforward neural network. The input layer, output layer, and hidden layer consist of at least three layers (Figure 2(b)). The deep neural network (DNN) is a multi-layered MLP. More precisely using fewer neurons, additional layers and, therefore, connections enable the modelling of rare dependencies in the training data [4]. Nevertheless, the DNN learning process can result in overfitting and declining performance [5].

In the theory of ANN, the universal approximation theorem says that a single hidden layer of MLP is enough to estimate, with a certain accuracy, all compactly supported continuous real functions. In many cases, however, DNN predictions are more exact, as research shows [3], compared to those obtained by ANN networks.

ANN changes weights depending on the degree of an error function during the training process to minimize the error. There are several different algorithms for training purposes. Depending on a particular problem, the algorithms may vary in performance [34].

**3.4. DRL Algorithm for Reward-Penalty Decision.** DRL [35, 36] consists of agents that access its actions and observations at a time to either reward or penalize the actions, i.e., the classification. The detailed steps are given in Algorithm 1, where DRL compares the classified results of the ANN with features extracted in the repository. If the observed and the original class are the same, then the classifier is rewarded, and vice versa.

The executions of Algorithm 1 are sent to the ANN that determines whether the unsupervised learning at each iteration is of a reward or a penalty one. This ensures that the classification of ANN-DRL is accurate and precise. Finally, the estimation of the harmfulness index [37] helps in the estimation of the CB detection as accurate or not.

## 4. Results

In this section, we present the details of the experiments using the collected datasets and the performance metrics. The study has selected 30,384 tweets collected from the twitter datasets [4]. The tweets contain both CB and non-CB tweets, where automated labelling or tagging is carried out using feature selection methods. The tagging of CB and non-CB is made based on various attributes as mentioned in Table 1, which is a common trait used in online communication over social networks. The input tweet data are, hence, classified as CB and non-CB, where the former indicates the vulnerable behavior and the latter indicates genuine behavior. Out of

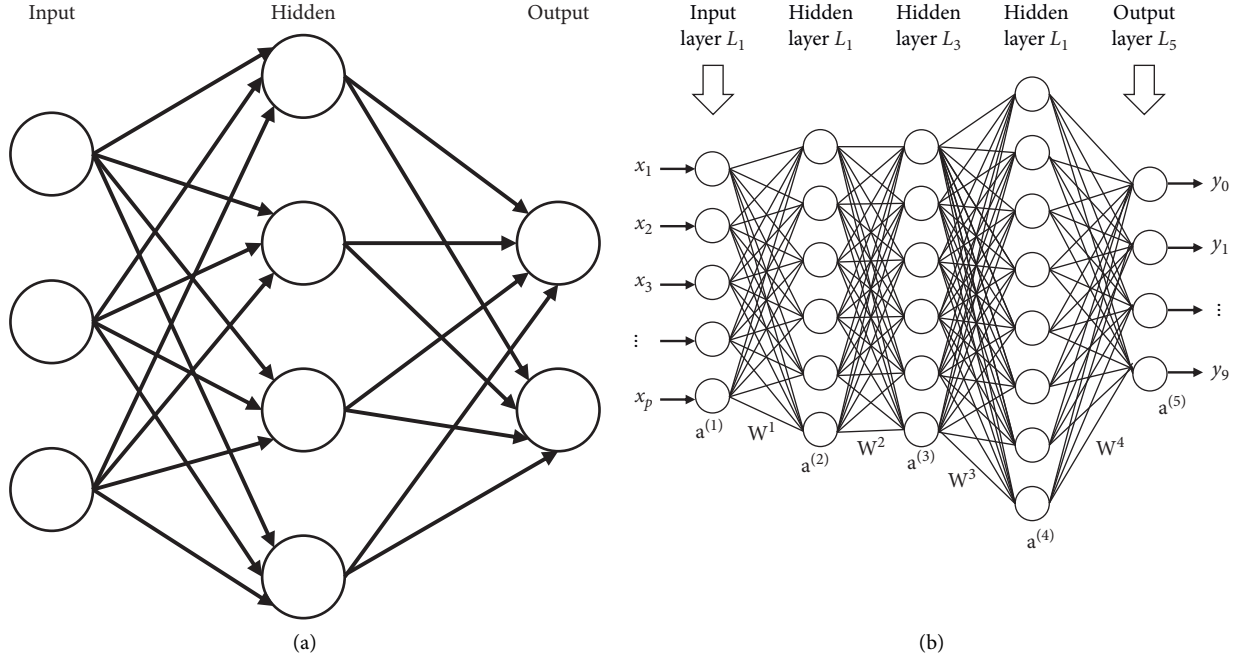


FIGURE 2: (a) ANN architecture. (b) 3-layered ANN architecture

30,384, more than 1252 tweets are classified as CB datasets; however, the labelled data are not used to train the classifier. These labelled data act as an input for the DRL method, which rewards or penalizes the ANN mechanism. The entire datasets have more imbalanced classes that penalize the unsupervised ANN with inaccurate results in identifying the relevant instances. The ANN, on other hand, with imbalanced classes, ignores minor classes, and it performs well with major classes.

The weight adjustment approach helps to avoid over-sampling of the minority class, i.e., abnormal class and undersampling the majority class, i.e., the normal class. The entire set of experiments is conducted with the topmost algorithms performed well in existing methods that include the ANN, SVM, RF, and LR. These existing methods are compared with ANN-DRL to find the classification accuracy. As in [8], the present study utilized three feature selection methods, namely, information gain,  $\chi^2$ , and Pearson correlation techniques. A 10-fold cross validation is conducted, and the proposed classifier is tested individually with all three feature selection methods.

The performance is estimated against various metrics that include accuracy, F-measure, geometric mean (G-mean), percentage error, precision, sensitivity, and specificity. The details of the metrics are given below.

Accuracy is defined as the total number of predictions required to ensure that the system works correctly. It is estimated as the ratio of the total number of correct predictions and the total predictions; .

$$\text{Accuracy} = \frac{\text{TP} + \text{TN}}{\text{TP} + \text{TN} + \text{FP} + \text{FN}} \quad (7)$$

Here, TP is the true positive cases, where the model classifies the CB classes correctly. TN is the true negative cases, where

the model classifies the non-CB classes correctly. FP is the false positive cases, where the model wrongly classifies the CB classes correctly. FN is the false negative cases, where the model wrongly classifies the non-CB classes correctly.

F-measure is the weighted harmonic mean of the recall and precision values, which ranges between zero and one. Higher value of F-measure refers to higher classification performance.

$$F - \text{measure} = \frac{2\text{TP}}{2\text{TP} + \text{FP} + \text{FN}} \quad (8)$$

G-mean is defined as the aggregation of sensitivity and specificity measure, which intends to maintain the trade-off between them, especially when the dataset is found to be imbalanced. This is measured as follows:

$$G - \text{mean} = \sqrt{\frac{\text{TP}}{\text{TP} + \text{FN}} \times \frac{\text{TN}}{\text{TN} + \text{FP}}} \quad (9)$$

Mean Absolute Percentage error (MAPE) is defined as the measure of prediction accuracy that measures the total loss while predicting the actual classes. It is measured as the ratio of the difference between the actual ( $A_t$ ) and predicted class ( $F_t$ ), and the actual class. The entire value is multiplied by 100% and divided by the fitted points ( $n$ ). The formula for the percentage error is defined as follows:

$$\text{MAPE} = \frac{100}{n} \sum_{t=1}^n \left| \frac{A_t - F_t}{A_t} \right| \quad (10)$$

Sensitivity is defined as the ability of the deep learning model to identify correctly the true positive rate.

$$\text{Sensitivity} = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad (11)$$

**Input:** Eligibility trace decay term  $\lambda$ , learning rate  $\alpha$ , number of objectives  $n$ , discounting term  $\gamma$ ,  $a \leftarrow$  action ( $r$ =reward or  $p$ =penalty),  $s \leftarrow$  state,  $o \leftarrow$  observer

Initialize Population

**For** all states  $s$ , actions  $a$  and objectives  $o$  do

Initialize  $Q(s, a, o)$

**Endfor**

Evaluate each member of the Population

**For** each epoch do

**For** all states  $s$  and actions  $a$  do

$e(s, a) = 0$

**Endfor**

Observe initial state  $st$

Select action  $at$  based on an exploratory policy derived from  $Q(st)$

**For** each step of the episode do

Execute action  $at$ , observe  $s'$  find the vector as reward  $r$  or penalty  $s$

Select action  $a^*$  based on a greedy policy derived from  $Q(s')$

Select action  $a'$  based on an exploratory policy derived from  $Q(s')$

**For** each objective  $o$  do

$\delta o = ro + \gamma Q(s0, a^*, o) - Q(st, at, o)$

**End for**

Set  $e(st, at) = 1$

**For** each state  $s$  and action  $a$  do

**For** each objective  $o$  do

set  $Q(s, a, o) = Q(s, a, o) + \alpha \delta o e(s, a)$

**End for**

**Endfor**

**If**  $a' = a^*$  then

set  $e(s, a) = \gamma \lambda e(s, a)$

**Else**

set  $e(s, a) = 0$

**Endif**

**Endfor**

$st = s', at = a'$

**Endfor**

ALGORITHM 1: DRL algorithm.

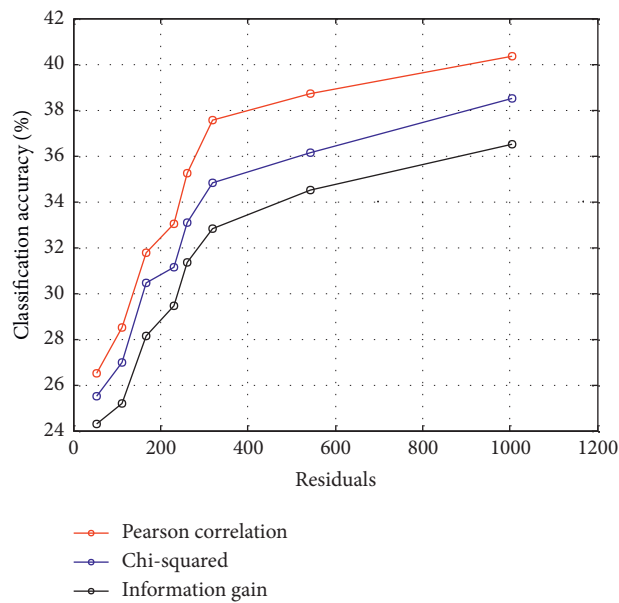


FIGURE 3: Comparison of feature selection methods with 60% training data.



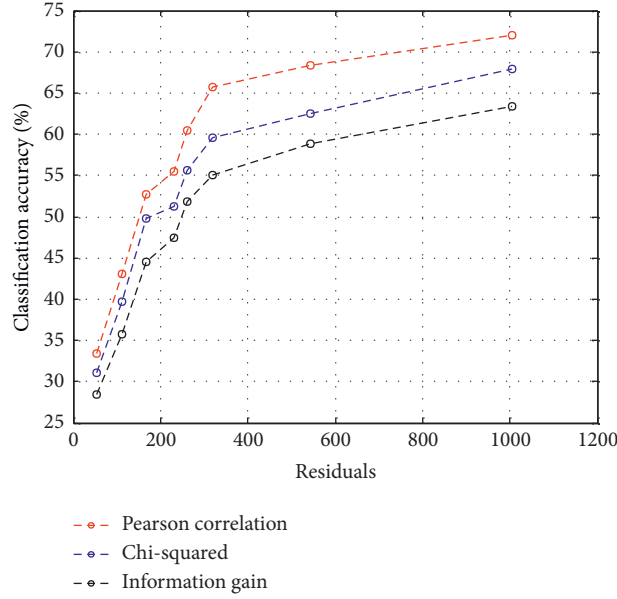


FIGURE 4: Comparison of feature selection methods with 75% training data.

Specificity is defined as the ability of the deep learning model to identify correctly the true negative rate.

$$\text{Specificity} = \frac{\text{TN}}{\text{TN} + \text{FP}} \quad (12)$$

**4.1. Analysis.** This section provides the results of classification as in following tables. The proposed ANN-DRL is validated and compared with existing methods, namely, the ANN, SVM, RF, LR, and NB. The results of predicting the CB are validated against 60%, 75%, and 90% training data with various feature extraction methods: information gain,  $\chi^2$ , and Pearson correlation techniques.

Figures 3–5 show the results of training the feature selection method with 60%, 75%, and 90% of training data and presenting the classification accuracy of the proposed classifier. The result shows that the Pearson correlation has the highest classification accuracy than information gain and  $\chi^2$ . The result further shows that, at some point, with increasing the number of residuals, the classification accuracy using information gain as a feature selection tool drops the most compared with chi-squared and Pearson correlation. Therefore, the class of CB is determined accurately with Pearson correlation and ANN-DRL as the classifier;

Tables 2–4 show the results of predicting the CB over 60%, 75%, and 90% of training data with information gain as a feature selection tool. Tables 5–7 show the results of predicting the CB over 60%, 75%, and 90% of training data with  $\chi^2$  tool. Tables 8–10 show the results of predicting the CB over 60%, 75%, and 90% of training data with Pearson correlation tool. The results of simulation show that the proposed method has higher classification accuracy than the existing classifiers. It is further inferred that the Pearson correlation has optimal selection of features that has boosted the classification accuracy with 90% training data than 75% or 60% datasets. The other metrics show optimal

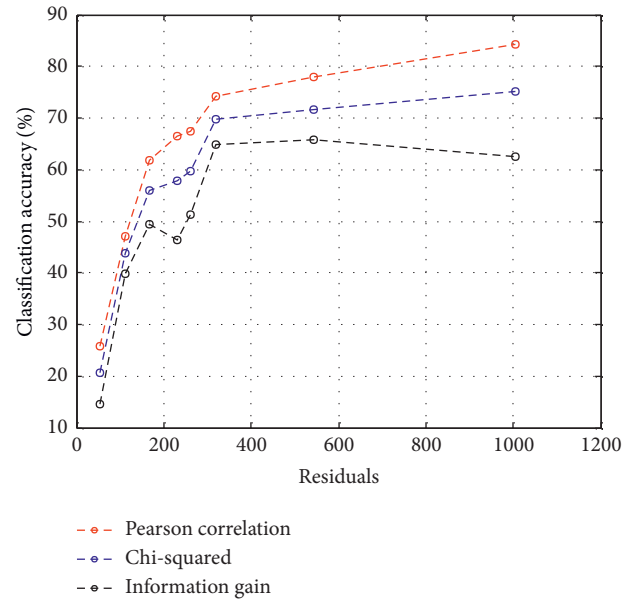


FIGURE 5: Comparison of feature selection methods with 90% training data.

performance for Pearson correlation than the other feature selection tools. Furthermore, the MAPE of the ANN-DRL is lesser than that of the other methods (Table 11).

To test the efficacy of ANN algorithm in the proposed method, we validate the algorithm with a 3000 test dataset and present a confusion matrix. Here, the 3000 test samples are picked randomly from the overall datasets, which is not native to the trained datasets. A 10-fold cross validation is conducted to test the ANN with the DRL scheme. The result shows that the classified results have 1740 TP cases, 1030 TN cases, 160 FN, and 70 FP cases, which is evident from Table 12.



TABLE 2: Results of predicting the CB with 60% training data with information gain.

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	56.890	57.190	59.280	59.541	60.901	81.688
F-measure	39.614	41.714	52.948	53.108	55.479	84.869
G-mean	73.755	73.986	75.486	75.526	75.936	86.790
MAPE	29.550	26.609	25.209	22.628	22.048	17.346
Sensitivity	62.962	66.473	74.376	86.760	87.420	97.464
Specificity	75.396	75.586	79.097	79.117	80.488	81.328

TABLE 3: Results of predicting the CB with 75% training data with information gain.

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	97.674	98.394	98.424	98.504	98.514	98.644
F-measure	53.588	70.944	71.265	74.146	77.377	80.578
G-mean	83.099	83.949	85.570	87.130	92.172	93.692
MAPE	28.130	26.739	23.968	21.297	11.834	91.332
Sensitivity	69.914	71.305	74.076	76.756	86.210	89.811
Specificity	97.744	98.534	98.734	98.814	98.834	98.894

TABLE 4: Results of predicting the CB with 90% training data with information gain.

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	97.124	97.144	97.154	97.244	97.264	97.324
F-measure	78.597	78.727	79.247	80.328	81.008	81.298
G-mean	80.648	80.888	81.158	82.158	82.478	82.669
MAPE	32.371	32.011	31.491	29.880	29.380	29.060
Sensitivity	65.673	66.033	66.553	68.164	68.664	68.984
Specificity	95.933	95.993	96.033	97.264	97.674	98.024

TABLE 5: Results of predicting the CB with 60% training data with  $\chi^2$ .

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	57.480	60.091	62.452	63.822	67.063	85.900
F-measure	67.963	68.013	68.904	70.024	75.056	80.618
G-mean	44.725	57.650	60.681	45.926	77.307	87.200
MAPE	20.597	17.926	17.836	12.984	11.654	10.504
Sensitivity	77.447	80.128	80.218	85.059	86.400	87.550
Specificity	74.606	77.487	78.427	81.588	83.589	85.680

TABLE 6: Results of predicting the CB with 75% training data with  $\chi^2$ .

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	98.944	98.964	98.984	98.984	98.994	98.994
F-measure	90.411	91.842	91.992	92.502	92.712	93.352
G-mean	94.483	97.884	98.434	98.734	98.874	98.874
MAPE	87.860	28.240	21.477	10.504	55.849	22.228
Sensitivity	90.161	96.793	97.894	98.474	98.754	98.764
Specificity	97.974	97.994	97.994	97.994	97.994	98.634

TABLE 7: Results of predicting the CB with 90% training data with  $\chi^2$ .

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	98.584	98.584	98.664	98.664	98.684	98.734
F-measure	87.100	87.220	89.161	89.191	90.551	90.561
G-mean	95.243	95.243	95.613	95.683	96.013	96.053
MAPE	72.015	71.835	63.942	62.612	55.259	54.579
Sensitivity	91.742	91.742	92.552	92.682	93.422	93.492
Specificity	98.674	98.684	98.774	98.774	98.864	98.864

TABLE 8: Results of predicting the CB with 60% training data with Pearson correlation.

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	60.291	67.073	70.094	75.306	79.147	83.629
F-measure	70.904	71.155	71.325	71.505	76.066	81.608
G-mean	71.205	71.435	73.155	75.216	77.677	80.428
MAPE	69.334	65.693	58.970	40.854	37.993	36.142
Sensitivity	78.737	72.365	73.085	74.856	75.056	81.908
Specificity	71.615	73.495	76.566	81.798	83.109	83.519

TABLE 9: Results of predicting the CB with 75% training data with Pearson correlation.

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	95.183	95.333	95.413	95.463	95.643	95.653
F-measure	59.471	61.261	61.601	62.061	63.592	63.732
G-mean	80.358	80.868	81.398	81.628	82.739	83.199
MAPE	31.131	30.400	29.510	29.160	27.329	26.509
Sensitivity	66.913	67.643	68.534	68.894	70.714	71.545
Specificity	96.463	96.623	96.633	96.683	96.723	96.773

TABLE 10: Results of predicting the CB with 90% training data with Pearson Correlation.

Statistical parameters	NB	LR	RF	SVM	ANN	ANN-DRL
Accuracy	98.653	98.653	98.733	98.733	98.753	98.803
F-measure	87.161	87.281	89.223	89.253	90.614	90.625
G-mean	95.309	95.309	95.680	95.750	96.080	96.120
MAPE	72.065	71.885	63.987	62.656	55.298	54.617
Sensitivity	91.806	91.806	92.617	92.747	93.487	93.558
Specificity	98.743	98.753	98.843	98.843	98.933	98.933

Depending on the execution results, we found the computational complexity of the ANN-DRL is lesser than that of the existing machine learning methods on detecting the cyberbullying contents. However, the complexity increases with increased layers of the neural network and increased iterations on DRL. It is found that the ANN-DRL is  $O(nl + en + n^3 + n_{\text{layers}})$  for training and  $O(l + en + n^3 + n_{\text{layers}})$  for testing, where  $n$  is the training samples,  $l$  is the features, and  $n_{\text{layers}}$  is the total number of hidden layers with  $n$  neurons. The ANN has  $O(nl + n_{\text{layers}})$

TABLE 11: Summary of various methods on cyberbullying.

Authors	Features used	Classifier
Nandhini and Sheeba [20]	Noun, pronoun, and adjective	Fuzzy logic-based genetic algorithm
Potha et al. [21]	Local, sentimental, contextual, and gender-specific language features	SVM
Kumar and Sachdeva [28]	Direct and indirect CB features	SVM
Al-garadi et al. [8]	Network, activity and user information, and tweet content	SVM
[28]	Network, activity and user information, and tweet content	Naïve Bayes (NB)
[25]	Network, activity and user information, and tweet content	k-nearest neighbor (KNN) and random forest (RF)
Balakrishnan et al. [25]	Psychological features	NB, RF, and J48
Murnion et al. [18]	IsAbusive, IsPositive, IsNegative, HasBadLanguage, IsRacist, NoobRelated, SpecificTarget, and FilteredText	Sentiment text analytics system is supported with a scoring scheme
Ho et al. [27]	Abusive words	Logistic regression model
Balakrishnan et al. [24]	15 twitter features [23]	RF classifier
Sánchez-Medina et al. [26]	Psychopathy, narcissism, and machiavellianism	Ensemble classification trees
Lee et al. [22]	New abusive words	Three-layered neural network model

TABLE 12: Confusion matrix on a 3000 test dataset.

Actual	Predicted CB		Total
	Present	Absent	
<b>Present</b>	<b>TP</b> (1740)	<b>FN</b> (160)	1900
<b>Absent</b>	<b>FP</b> (70)	<b>TN</b> (1030)	1100
<b>Predicted CB</b>	1810	1190	3000

for training and  $O(l + n_{\text{layers}})$  for testing, and SVM has  $O(n^2l + n^3)$  for training and  $O(n_{sv}l)$  for testing, where  $n_{sv}$  is the support vectors. RF has  $O(n^2\sqrt{l}n_{\text{trees}})$  for training and  $O(ln_{\text{trees}})$  for testing, where  $n_{\text{trees}}$  is the total trees in random forest. LR has  $O(l^2n + l^3)$  for training and  $O(l)$  for testing, and NB has  $O(nl)$  for training and  $O(l)$  for testing.

However, researchers are now trying to apply their proposed methods on this problem [38–49].

## 5. Conclusions

In this paper, an integrated model using an ANN and DRL is designed for the classification of CB from raw text datasets of a social media engine. The extraction of psychological features, user comments, and the context has enabled better

classification performance, where an ANN at the initial stage performs with improved classification results. The addition of a reward-penalty system using DRL has enhanced the classification to a much greater level than the ANN model. The simulation results illustrate the improved average classification accuracy of 80.69% using ANN-DRL than existing three-layered ANN (77.40%), SVM (75.44%), RF (75.55%), LR (75.10%), and NB (75.19%). In future, the convolutional neural network can be applied on image datasets to extract the information to serve the purpose on reducing the cyberbullying. [50] [49]

## Data Availability

The data used to support the findings of this study are available from the author upon request (gdhiman0001@gmail.com).

## Conflicts of Interest

The authors declare that they have no conflicts of interest.

## Acknowledgments

The authors are thankful for the support from Taif University Researchers Supporting Project (TURSP-2020/98), Taif University, Taif, Saudi Arabia.

## References

- [1] M. Ptaszynski, F. Masui, T. Nitta et al., “Sustainable cyberbullying detection with category-maximized relevance of harmful phrases and double-filtered automatic optimization,” *International Journal of Child-Computer Interaction*, vol. 8, pp. 15–30, 2016.
- [2] N. S. Ansary, “Cyberbullying: concepts, theories, and correlates informing evidence-based best practices for prevention,” *Aggression and Violent Behavior*, vol. 50, Article ID 101343, 2020.
- [3] S. Bahassine, A. Madani, M. Al-Sarem, and M. Kissi, “Feature selection using an improved Chi-square for Arabic text classification,” *Journal of King Saud University-Computer and Information Sciences*, vol. 32, no. 2, pp. 225–231, 2020.
- [4] A. Dirksen, S. Verberne, A. Sarker, and W. Kraaij, “Data-driven lexical normalization for medical social media,” *Multimodal Technologies and Interaction*, vol. 3, no. 3, p. 60, 2019.
- [5] G. Zhang, J. Hou, J. Wang, C. Yan, and J. Luo, “Feature selection for microarray data classification using hybrid information gain and a modified binary krill herd algorithm,” *Interdisciplinary Sciences, Computational Life Sciences*, vol. 12, no. 3, pp. 288–301, 2020.
- [6] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, Seattle, WA, USA, June 2016.
- [7] C. Baral, O. Fuentes, and V. Kreinovich, “Why deep neural networks: a possible theoretical explanation,” in *Constraint Programming and Decision Making: Theory and Applications*, pp. 1–5, Springer, Berlin, Heidelberg, 2018.
- [8] G. Dhiman and V. Kumar, “Multi-objective spotted hyena optimizer: a multi-objective optimization algorithm for

- engineering problems,” *Knowledge-Based Systems*, vol. 150, pp. 175–197, 2018.
- [9] M. Cogswell, F. Ahmed, R. Girshick, L. Zitnick, and D. Batra, “Reducing overfitting in deep networks by decorrelating representations,” 2015, <https://arxiv.org/abs/1511.06068>.
  - [10] S. Dreiseitl and L. Ohno-Machado, “Logistic regression and artificial neural network classification models: a methodology review,” *Journal of Biomedical Informatics*, vol. 35, no. 5-6, pp. 352–359, 2002.
  - [11] J. Schmidhuber, “Deep learning in neural networks: an overview,” *Neural Networks*, vol. 61, pp. 85–117, 2015.
  - [12] T. T. Nguyen, N. D. Nguyen, and S. Nahavandi, “Deep reinforcement learning for multiagent systems: a review of challenges, solutions, and applications,” *Institute of Electrical and Electronics Engineers Transactions on Cybernetics*, 2020.
  - [13] B. Wang, Y. Li, W. Ming, and S. Wang, “Deep reinforcement learning method for demand response management of interruptible load,” *Institute of Electrical and Electronics Engineers Transactions on Smart Grid*, 2020.
  - [14] T. Aricak, S. Siyahhan, A. Uzunhasanoglu et al., “Cyberbullying among Turkish adolescents,” *Cyber Psychology & Behavior*, vol. 11, no. 3, pp. 253–261, 2008.
  - [15] G. Dhiman, “ESA: a hybrid bio-inspired metaheuristic optimization approach for engineering problems,” *Engineering with Computers*, vol. 37, pp. 1–31, 2019.
  - [16] D. Olweus and S. P. Limber, “Some problems with cyberbullying research,” *Current Opinion in Psychology*, vol. 19, pp. 139–143, 2018.
  - [17] T. Vaillancourt, R. Faris, and F. Mishna, “Cyberbullying in children and youth: implications for health and clinical practice,” *The Canadian Journal of Psychiatry*, vol. 62, no. 6, pp. 368–373, 2017.
  - [18] G. Dhiman and V. Kumar, “Seagull optimization algorithm: theory and its applications for large-scale industrial engineering problems,” *Knowledge-Based Systems*, vol. 165, pp. 169–196, 2019.
  - [19] M. W. Savage and R. S. Tokunaga, “Moving toward a theory: testing an integrated model of cyberbullying perpetration, aggression, social skills, and Internet self-efficacy,” *Computers in Human Behavior*, vol. 71, pp. 353–361, 2017.
  - [20] G. Dhiman and V. Kumar, “Spotted hyena optimizer: a novel bio-inspired based metaheuristic technique for engineering applications,” *Advances in Engineering Software*, vol. 114, pp. 48–70, 2017.
  - [21] G. Dhiman and V. Kumar, “Emperor penguin optimizer: a bio-inspired algorithm for engineering problems,” *Knowledge-Based Systems*, vol. 159, pp. 20–50, 2018.
  - [22] G. Dhiman and A. Kaur, “STOA: a bio-inspired based optimization algorithm for industrial engineering problems,” *Engineering Applications of Artificial Intelligence*, vol. 82, pp. 148–174, 2019.
  - [23] S. Kaur, L. K. Awasthi, A. L. Sangal, and G. Dhiman, “Tunicate Swarm Algorithm: a new bio-inspired based metaheuristic paradigm for global optimization,” *Engineering Applications of Artificial Intelligence*, vol. 90, p. 103541, 2020.
  - [24] I. Cuadrado-Gordillo and I. Fernández-Antelo, “Adolescents’ perception of the characterizing dimensions of cyberbullying: differentiation between bullies’ and victims’ perceptions,” *Computers in Human Behavior*, vol. 55, pp. 653–663, 2016.
  - [25] M. Dadvar, D. Trieschnigg, R. Ordelman, and F. de Jong, “March). Improving cyberbullying detection with user context,” in *European Conference on Information Retrieval*, pp. 693–696, Springer, Berlin, Heidelberg, 2013.
  - [26] S. Murnion, W. J. Buchanan, A. Smales, and G. Russell, “Machine learning and semantic analysis of in-game chat for cyberbullying,” *Computers & Security*, vol. 76, pp. 197–213, 2018.
  - [27] H. Rosa, N. Pereira, R. Ribeiro et al., “Automatic cyberbullying detection: a systematic review,” *Computers in Human Behavior*, vol. 93, pp. 333–345, 2019.
  - [28] B. S. Nandhini and J. I. Sheeba, “Online social network bullying detection using intelligence techniques,” *Procedia Computer Science*, vol. 45, pp. 485–492, 2015.
  - [29] N. Potha, M. Maragoudakis, and D. Lyras, “A biology-inspired, data mining framework for extracting patterns in sexual cyberbullying data,” *Knowledge-Based Systems*, vol. 96, pp. 134–155, 2016.
  - [30] M. A. Al-garadi, K. D. Varathan, and S. D. Ravana, “Cyber-crime detection in online communications: the experimental case of cyberbullying detection in the Twitter network,” *Computers in Human Behavior*, vol. 63, pp. 433–443, 2016.
  - [31] H.-S. Lee, H.-R. Lee, J.-U. Park, and Y.-S. Han, “An abusive text detection system based on enhanced abusive and non-abusive word lists,” *Decision Support Systems*, vol. 113, pp. 22–31, 2018.
  - [32] I.-K. Peter and F. Petermann, “Cyberbullying: a concept analysis of defining attributes and additional influencing factors,” *Computers in Human Behavior*, vol. 86, pp. 350–366, 2018.
  - [33] V. Balakrishnan, S. Khan, T. Fernandez, and H. R. Arabnia, “Cyberbullying detection on twitter using Big five and Dark Triad features,” *Personality and Individual Differences*, vol. 141, pp. 252–257, 2019.
  - [34] W. Feng, Q. Zhu, J. Zhuang, and S. Yu, “An expert recommendation algorithm based on Pearson correlation coefficient and FP-growth,” *Cluster Computing*, vol. 22, no. 3, pp. 7401–7412, 2019.
  - [35] V. Balakrishnan, S. Khan, and H. R. Arabnia, “Improving cyberbullying detection using twitter users’ psychological features and machine learning,” *Computers & Security*, vol. 90, Article ID 101710, 2020.
  - [36] A. J. Sánchez-Medina, I. Galván-Sánchez, and M. Fernández-Monroy, “Applying artificial intelligence to explore sexual cyberbullyingbehaviour,” *Heliyon*, vol. 6, no. 1, pp. 1–9, 2020.
  - [37] S. M. Ho, D. Kao, M.-J. Chiu-Huang, W. Li, and C.-J. Lai, “Detecting cyberbullying “hotspots” on twitter: a predictive analytics approach,” *Forensic Science International: Digital Investigation*, vol. 32, p. 300906, 2020.
  - [38] G. Dhiman and M. Garg, “MoSSE: a novel hybrid multi-objective meta-heuristic algorithm for engineering design problems,” *Soft Computing*, vol. 24, pp. 1–20, 2020.
  - [39] G. Dhiman, K. K. Singh, A. Slowik et al., “EMoSQA: a new evolutionary multi-objective seagull optimization algorithm for global optimization,” *International Journal of Machine Learning and Cybernetics*, vol. 12, pp. 1–26, 2020.
  - [40] G. Dhiman, D. Oliva, A. Kaur et al., “BEPO: a novel binary emperor penguin optimizer for automatic feature selection,” *Knowledge-Based Systems*, vol. 211, Article ID 106560, 2021.
  - [41] G. Dhiman, K. K. Singh, M. Soni et al., “MOSOQA: a new multi-objective seagull optimization algorithm,” *Expert Systems with Applications*, Article ID 114150, 2020.
  - [42] H. Kaur, A. Rai, S. S. Bhatia, and G. Dhiman, “MOEPO: a novel Multi-objective Emperor Penguin Optimizer for global optimization: special application in ranking of cloud service providers,” *Engineering Applications of Artificial Intelligence*, vol. 96, Article ID 104008, 2020.

- [43] M. Dehghani, Z. Montazeri, A. Dehghani et al., "DM: dehghani Method for modifying optimization algorithms," *Applied Sciences*, vol. 10, no. 21, Article ID 7683, 2020.
- [44] A. Sharma, P. K. Singh, A. Sharma, and R. Kumar, "An efficient architecture for the accurate detection and monitoring of an event through the sky," *Computer Communications*, vol. 148, pp. 115–128, 2019.
- [45] A. Sharma, P. K. Singh, and Y. Kumar, "An integrated fire detection system using IoT and image processing technique for smart cities," *Sustainable Cities and Society*, vol. 61, Article ID 102332, 2020.
- [46] D. Kumar, A. Sharma, R. Kumar, and N. Sharma, "Restoration of the network for next generation (5G) optical communication network," in *Proceedings of the 2019 International Conference on Signal Processing and Communication (ICSC)*, pp. 64–68, IEEE, Noida, India, March 2019.
- [47] A. Sharma, R. Sarishma, R. Tomar, N. Chilamkurti, and B.-G. Kim, "Blockchain based smart contracts for internet of medical things in e-healthcare," *Electronics*, vol. 9, no. 10, Article ID 1609, 2020.
- [48] A. Sharma and R. Kumar, "Computation of the reliable and quickest data path for healthcare services by using service-level agreements and energy constraints," *Arabian Journal for Science and Engineering*, vol. 44, no. 11, pp. 9087–9104, 2019.
- [49] <https://www.kaggle.com/data/35739Twitter> datasets.
- [50] A. Kumar and N. Sachdeva, "Cyberbullying detection on social multimedia using soft computing techniques: a meta-analysis," *Multimedia Tools and Applications*, vol. 78, no. 17, pp. 23973–24010, 2019.