

## ORIGINAL RESEARCH ARTICLE

# Experiences of sexual minorities on social media: A study of sentiment analysis and machine learning approaches

Peter Appiahene<sup>1</sup>, Vijayakumar Varadarajan<sup>2,3,4,\*</sup>, Tao Zhang<sup>5</sup>, Stephen Afrifa<sup>1,5,\*</sup>

<sup>1</sup> Department of Computer Science and Informatics, University of Energy and Natural Resources, Sunyani 00233, Ghana

<sup>2</sup> School of Computer Science and Engineering, University of New South Wales, Sydney, NSW 2052, Australia

<sup>3</sup> International Divisions, Ajeenkya D. Y. Patil University, Pune 412105, India

<sup>4</sup> School of Information Technology, Swiss School of Business Management, Geneva 1213, Switzerland

<sup>5</sup> Department of Information and Communication Engineering, Tianjin University, Tianjin 300072, China

\* **Corresponding author:** Vijayakumar Varadarajan, v.varadarajan@unsw.edu.au; Stephen Afrifa, afrifastephen@tju.edu.cn

## ABSTRACT

Nowadays, social media has become a forum for people to express their views on issues such as sexual orientation, legislation, and taxes. Sexual orientation refers to individuals with whom you are attracted and wish to be engaged. In the world, many people are regarded as having different sexual orientations. People categorized as lesbian, gay, bisexual, transgender, queer, and many more (LGBTQ+) have many sexual orientations. Because of the public stigmatization of LGBTQ+ persons, many turn to social media to express themselves, sometimes anonymously. The present study aims to use natural language processing (NLP) and machine learning (ML) approaches to assess the experiences of LGBTQ+ persons. To train the data, the study used lexicon-based sentiment analysis (SA) and six distinct machine classifiers, including logistic regression (LR), support vector machine (SVM), naïve bayes (NB), decision tree (DT), random forest (RF), and gradient boosting (GB). Individuals are positive about LGBTQ concerns, according to the SA results; yet, prejudice and harsh statements against the LGBTQ people persist in many regions where they live, according to the negative sentiment ratings. Furthermore, using LR, SVM, NB, DT, RF, and GB, the ML classifiers attained considerable accuracy values of 97%, 96%, 88%, 100%, 92%, and 91%, respectively. The performance assessment metrics used obtained significant recall and precision values. This study will assist the government, non-governmental organizations, and rights advocacy groups make educated decisions about LGBTQ+ concerns in order to ensure a sustainable future and peaceful coexistence.

**Keywords:** machine learning; sentiment analysis; LGBTQ; rights; artificial intelligence; natural language processing

## ARTICLE INFO

Received: 7 May 2023

Accepted: 1 June 2023

Available online: 4 August 2023

## COPYRIGHT

Copyright © 2023 by author(s).

Journal of Autonomous Intelligence is published by Frontier Scientific Publishing.

This work is licensed under the Creative Commons Attribution-NonCommercial 4.0 International License (CC BY-NC 4.0).

<https://creativecommons.org/licenses/by-nc/4.0/>

## 1. Introduction

The globe has evolved into a global ecosystem in which information exchange has grown massive. Information is essential for communicating meaning to a recipient<sup>[1]</sup>. Information is now exchanged in cyberspace. The term “cyberspace” refers to the worldwide environment that allows electronic resources from all over the world to be shared<sup>[2]</sup>. Smartphones are used to transfer information across, for example, social media websites such as Facebook, Twitter, TikTok, and many more. Smartphone-enabling technologies aid in connecting to various platforms, allowing it to be a vital component of the Internet of Things (IoT). With the growth of technology, many individuals, regardless of race or sexual orientation, turn to social media

platforms to express their views on social, economic, and political concerns<sup>[3,4]</sup>. Sexual orientation refers to those you are attracted to and desire to be involved with<sup>[5]</sup>. Several people nowadays are classified as having different sexual orientations. The terms lesbian, gay, bisexual, transgender, and queer (LGBTQ) refer to a person's sexual orientation or gender identity. People are being attacked all around the world for who they love, how they dress, and, ultimately, for who they are<sup>[6]</sup>. Being lesbian, gay, bisexual, transgender, or queer means facing daily persecution in far too many nations. According to a research by Amnesty International<sup>[7]</sup>, between October 2017 and September 2018, at least 369 people were killed in a wave of violence against trans persons.

Many intersex persons worldwide are compelled to endure risky, intrusive, and utterly unneeded procedures that can have long-term medical and psychological consequences. Many LGBTQ persons face major gender harassment on social media platforms, which is considered cyberbullying, and this causes depression because of their sexual orientation<sup>[8]</sup>. The majority of these abuses are the result of prejudice from persons of various sexual orientations other than the LGBTQ group. The consequences of this harassment on LGBTQ people throughout the world, particularly on social media, are immense, and can lead to mental health problems. As a result, it is critical to comprehend the sentiments of the LGBTQ community and provide appropriate assistance or advice. LGBTQ persons, on the other hand, are a sexual minority who are rarely represented or spoken out in public. Individuals are more willing to express themselves on social media when they are anonymous, yet they are met with unpleasant remarks or replies. To this purpose, social media data is a suitable data resource to the study of sentiment analysis (SA) and machine learning (ML) of the LGBTQ community. This study offers a fresh perspective to the topic of SA and ML in the literature, which will aid in minimizing, avoiding, and comprehending the LGBTQ community's attitudes for informed decision making. The following are the key contributions of this study, as adapted from the study of Afrifa et al.<sup>[11]</sup> and Adu et al.<sup>[9]</sup>:

- a) To undertake SA on LGBTQ social media data using a natural language processing (NLP) technique.
- b) Use machine learning classifiers to examine the experiences of LGBTQ individuals using social media data.
- c) The research employs both qualitative and quantitative methods to present a variety of viewpoints on the topics of minimizing, avoiding, and comprehending the LGBTQ people.
- d) This study proposes a data-driven approach for policymakers to use when making choices on LGBTQ individuals.
- e) It is a fresh contribution to the literature that proposes new models of SA and ML approaches.

To be clear, the study's contribution is to use automated sentiment analysis and machine learning classifiers to analyze LGBTQ data on social media. The remainder of the research is as follows: a survey of comparable works labeled "Literature Review" in section 2. Section 3 describes the materials and procedures used in this study. Moreover, the study's findings are detailed in section 4 titled "Results and Discussion". Section 5 concludes with the study's conclusions and recommendations.

## 2. Literature review

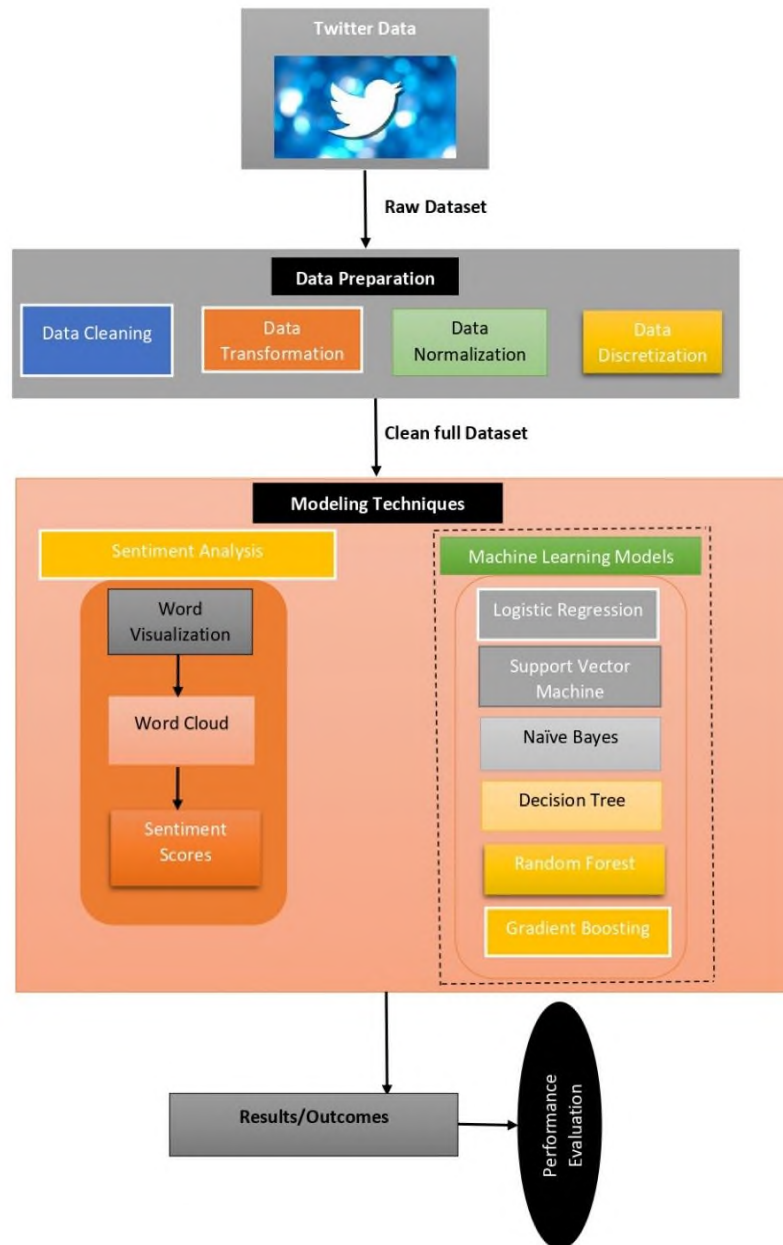
The artificial intelligence (AI) techniques have been widely applied in many different fields as AI has advanced over the years. AI has been applied to increase disease prediction accuracy<sup>[10,11]</sup>, particularly in healthcare. AI has also been used in secure networks to identify intrusions<sup>[12,13]</sup>. The natural language processing and machine learning have been widely employed in sentiment analysis of social media data. Afrifa and Varadarajan<sup>[2]</sup>, for example, exploited social media data to detect cyberbullying using NLP and ML approaches. Băroiu and Trăuşan-Matu<sup>[14]</sup> developed a method for automatically detecting sarcastic context using social media data, with the goal of developing an attention-based long short-term memory (LSTM)

architecture. Their study solely employed the deep learning method LSTM, with no sentiment scores or analysis derived from the data.

Additionally, Ainapure et al.<sup>[15]</sup> used Twitter tweets to examine the sentiments of Indian citizens towards the coronavirus disease 2019 (COVID-19) pandemic and vaccine effort. Deep learning and lexicon-based algorithms were used to classify the sentiments. The researchers concluded that the generated models can help healthcare personnel and governments make the best decisions during pandemic outbreaks in the future. Turner and Hammersjö<sup>[8]</sup>, conducted in-depth interviews with LGBTQ survivors of intimate partner violence (IPV) in Sweden to investigate the lived experiences of support-seeking. Their research produced an in-depth, phenomenological explanation of the support-seeking process, including the obstacles to, as well as the individual and societal facilitators of, seeking help. Although their study attempted to address LGBTQ support systems, it did not use SA and ML techniques to analyze the difficulties and experiences of LGBTQ persons. IFoodCloud, created by Zhang et al.<sup>[16]</sup>, automatically gathers data from more than 3100 public sources to comprehensively collect and evaluate public opinion on food safety in Greater China. Their study used numerous lexicon-based and machine learning-based algorithms coupled with IFoodCloud to create sentiment categorization models that give an unparalleled quick way of analyzing public opinion toward specific food safety issues. Their study highlighted the value of big data and machine learning in risk communication and decision-making. Furthermore, Çilgin et al.<sup>[17]</sup> conducted sentiment analysis on tweets shared by numerous people, organizations, and government agencies via Twitter during the worldwide COVID-19 epidemic using the VADER Sentiment Analysis approach. A total of 60,243,040 tweets were collected from Twitter. Their study's findings evaluated how tweets regarding COVID-19 shared at different times of the release mirrored distinct emotive scenarios. In a related work, Arcila-Calderón et al.<sup>[5]</sup> created and tested an automated detector of hate speech driven by gender and sexual orientation. The emphasis was on Twitter posts in Spanish. Their study failed to make use of sentiment analysis of Twitter data. Çilgin et al.<sup>[18]</sup> used machine learning to assess public sentiment of vaccine-related tweets gathered on Twitter in order to better understand social media users' opinions and concerns, particularly about COVID-19 vaccinations in Turkey. For the sentiment analysis studies, they used the majority voting approach in machine learning. Their study's findings demonstrated that the proposed approach is a viable and simple tool for monitoring the sensitivity of COVID-19 vaccinations using a sentiment analysis methodology via social media. Last but not least, Huynh Thai et al.<sup>[19]</sup> used YouTube user comments to conduct sentiment analysis and machine learning techniques on public attitudes on virtual tourism in the context of COVID-19. They concluded that their research met the necessity to analyze text analyses and natural language processing models using various sentiment analyses in order to attain ideal performance matrices.

### 3. Materials and methods

The **Figure 1** depicts the use of a proposed conceptual framework in this study. Data collection, data preparation, modeling, results/outcomes, and performance evaluation are all part of the methodology. The data used in this study was obtained from the microblogging website Twitter. Moreover, data preparation procedures such as data cleansing, data transformation, data normalization, and data discretization were used. The data preparation strategies aid in the creation of clean data and the extraction of significant features for training models<sup>[1]</sup>. To train the model, NLP and various ML classifiers were used. The outcomes or results of the various models are predicted. Lastly, several performance evaluations were used to assess the effectiveness of the various techniques employed. For data analysis, pre-processing, and modeling, the R-studio program was used.



**Figure 1.** The proposed conceptual framework.

### 3.1. The dataset

Many people currently communicate quickly owing to contemporary gadgets like smartphones and social networking websites like Twitter, Facebook, TikTok, and many more. Smartphones are important in the Internet of Things (IoT) because they can operate various IoT devices via an app on a smartphone. The dataset used in this study was collected from the microblogging website Twitter. The publicly available dataset from the FigShare data repository via the link (<https://doi.org/10.6084/m9.figshare.19787617.v1>) was accessed on 10 February 2023. The dataset includes 726,998 observations regarding lesbian, gay, bisexual, transgender, queer, and other (LGBTQ+) people. To train the proposed conceptual framework, the first 1000 Twitter comments out of 726,998 were used in this study. The entire dataset spans 8 months, from 1 October 2021 to 15 May 2022. Between 1 October 2021 and 31 December 2021, the first 1000 were collected. This section of data is used in this study to achieve a compromise between computational efficiency, model performance, and effective assessment depending on computing resources available.

## 3.2. Data preparation

The data preparation stage includes the procedures involved in obtaining clean data from raw data. It should be emphasized that data preparation is not done in any particular order, but rather repeatedly<sup>[9,20]</sup>. Data cleansing, transformation, annotation, normalization, and discretization are among the stages involved. Data preparation assists in the development of a good model, which may aid in the achievement of successful outcomes<sup>[21]</sup>.

### 3.2.1. Data cleaning

Data cleaning data is a vital step in every machine learning effort. Data cleaning detects and removes redundant instances from the data. Outlier detection and missing value imputation are two components of data cleansing<sup>[22]</sup>. Some data observations contain whitespaces (blank spaces), symbols, and uniform resource locators (URLs), making training inefficient<sup>[23]</sup>. Utilizing unclean data might also be expensive. Because of greater quality data in the decision-making process, clean data requires less time and effort to build models.

### 3.2.2. Data transformation

Data transformation has a significant impact on data mining since it helps to fill in missing values in data and brings information to the surface by developing new features to reflect trends and other ratios<sup>[24]</sup>. To label the data, data annotation was conducted on the dataset. Data annotation is the process of identifying specific bits of training data (whether text, photos, audio, or video) to assist machines in understanding what is in it and what is significant<sup>[25]</sup>. This annotated data is then utilized to train models. It is important to note that the observations (text comments) are likewise pre-processed. By converting all letters to lowercase, deleting punctuation marks, and removing stop words and errors, text pre-processing helps to remove unwanted bits of the data, or noise. To choose the desired characteristics, we use the Term Frequency-Inverse Document Frequency approach. The feature selection is depicted in Equation (1) below;

$$TF - IDF = FF * \log(N|DF) \quad (1)$$

where  $N$  represents the number of documents and  $DF$  is the number of documents that possess the feature.  $FF$  assigns a value of 0 or 1 based on the absence or existence of a feature in the document.

### 3.2.3. Data normalization

Results may be influenced when many qualities have different scales. Normalization equalizes all qualities on a single scale<sup>[26]</sup>. All attributes were scaled into smaller ranges ranging from 0 to 1. All text attributes were scaled from 0 to 1. Normalization is essential in text analysis and sentiment analysis study for cleaning and normalizing text data, establishing consistent representations, lowering vocabulary size, dealing with noisy material, and enabling accurate and relevant analysis and comparisons<sup>[27,28]</sup>. The Min-Max approach is a commonly utilized normalizing method in this study. The normalization procedure used for this study is depicted in Equation (2). Additionally, the Min-Max approach is efficient since results may be improved when there are outliers or missing values in the data.

$$X_{norm} = \frac{X - X_{min}}{X_{max} - X_{min}} \quad (2)$$

where  $X$  is the text value,  $X_{norm}$  is the normalized value,  $X_{min}$  is the minimum value, and  $X_{max}$  is the maximum value.

### 3.2.4. Data discretization

Numeric data are discretized by assigning values to interval or concept labels. Several approaches such as binning, correlation, clustering, and decision tree analysis might be used to do this. In this study, the binning approach was used to discretize data. Moreover, the equal-frequency interval-based discretization approach was used. The lowest and maximum values of all discretized characteristics are computed using the equal-

frequency interval-based procedure<sup>[29]</sup>. The values are then sorted in ascending order. The domain's technique based on the same distribution of data points overcomes the problem of equal width interval discretization. This approach also attempts to address the drawback of equal-width interval discretization. This approach was used to discretize all of the dataset attributes in this study.

### 3.3. Modeling techniques

This section describes the models that were used to train the dataset in this study. The NLP and ML approaches were utilized to train the dataset models. The subsections that follow give a full description of the various models.

#### 3.3.1. Natural language processing

The NLP is a component of AI. A computer program's ability to comprehend natural language, or human language as it is spoken and written, is known as NLP<sup>[30–32]</sup>. Sentiment analysis, often known as opinion mining, is a NLP method for identifying the positivity, negativity, or neutrality of data<sup>[33,34]</sup>. Businesses frequently do sentiment analysis on textual data to track the perception of their brands and products in customer reviews and to better understand their target market. The lexicon-based sentiment analysis in the NLP technique was employed for the study. Words in texts are categorized as positive or negative (and occasionally neutral) in lexicon-based sentiment analysis using a valence dictionary<sup>[35,36]</sup>. Once each word in the text has been classified, we can get an overall sentiment score by counting the amount of positive and negative words and computationally combining these values. A popular formula to calculate sentiment score (*SenSc*) is summarized in Equation (3).

$$SenSc = \frac{\text{number of positive words} - \text{number of negative words}}{\text{total number of words}} \quad (3)$$

In the lexicon-based sentiment analysis, the overall sentiment of the text is calculated on the fly, based solely on the dictionary employed for identifying word valence. Moreover, the word cloud and sentiment scores are predicted using the lexicon-based approach in the study.

#### 3.3.2. Machine learning classifiers

In this study, various ML techniques are utilized to train the data. ML models have been used in the domain of IoT<sup>[1,37]</sup>, health sector<sup>[10]</sup>, and many other sectors of the economy. The dataset has two sections: training (80%) and testing (20%). The dataset serves as the basis for training the various machine learning classifiers. The models' detection is validated using 10-fold cross validation. Cross-validation is a testing approach that includes training several ML models on portions of the available input data and then assessing them on the complementary subset<sup>[10]</sup>. The following are the various ML classifiers utilized in this study:

##### 1) Logistic regression:

Supervised learning is demonstrated using logistic regression (LR). Logistic regression is used to compute or forecast the likelihood of a binary (yes/no) event occurring<sup>[38]</sup>. The most typical use of logistic regression is to address classification issues. This study employs the classification of the dataset to train the models. In logistic regression, the logistic function or sigmoid function is employed to compute probability. The logistic function is a straightforward S-shaped curve that converts input into a value between 0 and 1. The Equation (4) below summarizes the probability computation.

$$h\theta(X) = \frac{1}{1 + e - (\beta_0 + \beta_1 X)} \quad (4)$$

where  $h\theta(X)$  is output of logistic function  $0 \leq h\theta(X) \leq 1$ ,  $\beta_1$  is the slope,  $\beta_0$  is the y-intercept, and  $X$  is the independent variable. It must be emphasized that  $(\beta_0 + \beta_1 * X)$  is derived from equation of a line  $Y (\text{predicted}) = (\beta_0 + \beta_1 * X) + \text{Error value}$ .

## 2) Support vector machine:

Support vector machines (SVMs) may do binary separation but are typically built for multiclass classification. The support vector machine algorithm seeks a hyperplane in an  $N$ -dimensional space ( $N$ —the number of features) that distinguishes between data points. Hyperplanes are decision boundaries that contribute in the classification of data items<sup>[39]</sup>. Data points on either side of the hyperplane can belong to distinct classes. Support vectors are data points that are closer to the hyperplane and have an effect on its location and orientation<sup>[40]</sup>.

## 3) Naïve bayes:

The naïve bayes (NB) algorithm is a supervised learning method for classification issues that is based on the Bayes theorem<sup>[41]</sup>. The Bayes theorem, commonly referred to as Bayes' law or Bayes' rule, is used to calculate the likelihood of a hypothesis given certain previous information. The Equation (5) below shows the formula for the Bayes theorem. The NB is mostly employed in text classification tasks involving large training datasets<sup>[42]</sup>. Moreover, NB classifier is one of the most straightforward and efficient classification algorithms, aiding in the development of rapid machine learning models capable of making accurate predictions.

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (5)$$

where  $P(A|B)$  is posterior probability,  $P(B|A)$ ,  $P(B|A)$  is likelihood probability,  $P(A)$  is prior probability, and  $P(B)$  is marginal probability.

As a probabilistic classifier, the NB makes predictions based on the likelihood that an object will occur.

## 4) Decision tree:

The supervised learning algorithm family includes the decision tree (DT) algorithm. The decision tree approach may be used to resolve regression and classification issues, unlike other supervised learning techniques<sup>[43]</sup>. With a decision tree, the objective is to develop a training model that can be used to forecast the class or value of the target variable by learning straightforward decision rules inferred from prior data (training data)<sup>[44]</sup>. Decision trees are often designed to resemble how people think while making decisions, making them simple to comprehend. Since it displays a structure like a tree, the decision tree's reasoning is simple to comprehend.

## 5) Random forest:

The simple, adaptable, and quick random forest (RF) technique is based on the decision tree concept. Both classification and regression issues in machine learning may be addressed with the random forest. It is built on the idea of ensemble learning, which is a method of integrating many classifiers to address difficult issues and enhance model performance<sup>[2]</sup>. Some decision trees may predict the proper output, while others may not, since the random forest combines numerous trees to forecast the class of the dataset<sup>[21]</sup>. Yet, when all the trees are combined, they forecast the right result. When compared to other methods, the random forest requires less time.

## 6) Gradient boosting:

Gradient boosting is an effective approach for developing predictive models. The concept of boosting arose from the question of whether a weak learner may be transformed to become a better learner<sup>[45]</sup>. Gradient boosting involves three elements:

- A loss function to be optimized.
- A weak learner to make predictions.
- An additive model to add weak learners to minimize the loss function.

The type of loss function utilized is determined on the issue being addressed. The loss function must be differentiable, however there are numerous common loss functions available as well as the ability to define your own. In gradient boosting, decision trees serve as the weak learner. Regression trees are utilized specifically because they provide actual values for splits and can be combined together, allowing future model outputs to be added and “correct” the residuals in the predictions. Existing trees in the model are not modified, and new trees are inserted one at a time. While adding trees, a gradient descent approach is utilized to reduce loss.

### 3.4. Performance evaluation

Several performance assessment criteria are used to evaluate the machine learning classifiers used in this study. This aids in determining how well your machine learning model performs on a dataset it has never seen before<sup>[46,47]</sup>. The accuracy, recall, and precision performance measures were utilized in this study to evaluate the various machine learning classifiers. Accuracy is defined as the ratio of correct sample predictions to total number of predictions. The recall of the model assesses its ability to recognize positive samples. The recall is also known as true positive rate (*TPR*). The more positive samples identified, the larger the recall. Additionally, precision is used to calculate the number of accurately categorized predicted positive cases by the algorithm. All assessment measure is based on one of four classifications: true positives (*TP*), true negatives (*TN*), false positives (*FP*), and false negatives (*FN*). The accuracy, recall, and precision are presented in Equations (6)–(8), respectively.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (6)$$

$$Recall (TPR) = \frac{TP}{TP + FN} \quad (7)$$

$$Precision = \frac{TP}{TP + FP} \quad (8)$$

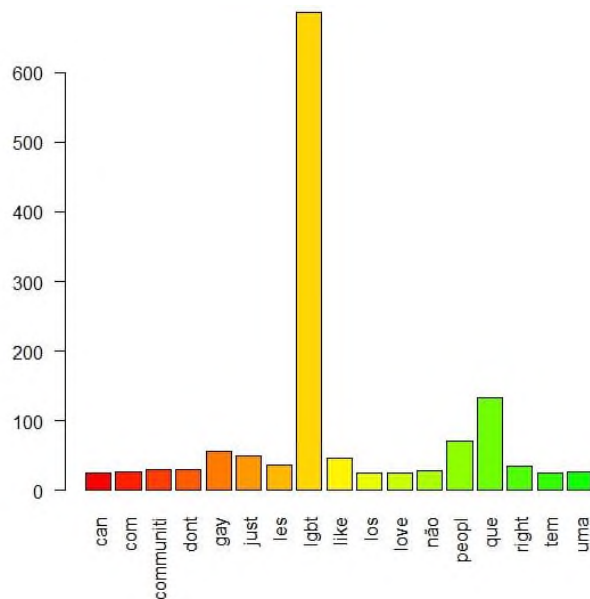
## 4. Results and discussion

The study’s findings are presented in this section. It should be noted that the study trained the data using SA and ML classifiers.

### 4.1. Outcomes of sentiment analysis classification

Individuals utilize smartphones and smart gadgets to communicate their views on taxation, laws, and other topics via social media. On 10 February 2023, 1000 text data (tweets) were retrieved from the entire data accessible. The data was pre-processed, and sentiment analysis was performed on it. It should be noted that this study leverages the lexicon-based approach’s Valence Aware Dictionary for sEntiment Reasoning (VADER). VADER provides more than just a vocabulary that distinguishes between positive, neutral, and negative sentiment polarity<sup>[48]</sup>. While VADER is a sentiment analysis tool that use a pre-built lexicon (dictionary) of words, it is capable of analyzing more than just three sentiment categories<sup>[2]</sup>. VADER’s vocabulary includes terms that capture sentiment intensity, sentiment modifiers, and negations in addition to positive and negative words. VADER can detect sentiment nuances and grasp the strength of sentiment expressions, making it an effective tool for studying sentiment across domains and capturing subtleties that go beyond binary positive or negative classification. To determine the frequency in text data (tweets), a term document matrix (TDM) was created. The word frequency effect describes the finding that high-frequency words are processed faster than low-frequency terms. **Figure 2** depicts the word frequency from Twitter data. The word frequency analyses the significance of words in a text or group of texts by counting the number of times specific words appear.





**Figure 2.** Word count visualization of the data.

The word frequency in the data (**Figure 2**) shows that “lgbt” had the greatest mention among the words after data pre-processing. As seen in **Figure 3** below, the word cloud illustrated how the words were distributed over the dataset. The objective is to discover patterns in the mentions of LGBTQ+ persons in the communities in which they live. Together with the search term “lgbt”, other terms were discovered, including “gay”, “people”, “right”, “communiti”, and “just”. It can be shown that people were quite worried about human rights concerns, such as the word “right”, with others even making bigoted comments to harm LGBTQ+ persons, such as, “poor”, and “necio” in Spanish which means foolish in English. The term “just” indicates that individuals were appealing for justice for the LGBTQ+ community, who felt their lives were at danger. Nonetheless, the word “love” showed in the word cloud, indicating that a portion of the population shares the LGBTQ+ community’s love emotions.

Additionally, LGBTQ+ persons feel endangered in certain places where they reside; nevertheless, in societies where laws protect them, it is a normal to live with other people without harm. **Figure 4** depicts the outcomes of LGBTQ tweets sorted into various sentimental scores or effects.

The following terms were used to determine the sentiment ratings in the text data (tweets): anger, anticipation, disgust, fear, joy, sadness, surprise, trust, negative, and positive. According to the lexicon, the basic emotions are anger, anticipation, disgust, fear, joy, sadness, surprise, and trust, with two polarities (negative and positive). **Figure 4** shows that positive sentiments (“positive”, “trust”, “joy”) outnumber negative sentiments (“negative”, “disgust”, “anger”). It can also be seen that the sentiment “fear” reached a higher value 130 as compared to the sentiment “joy”. The results show that individuals are optimistic about LGBTQ concerns; nonetheless, bigotry and harsh comments towards the LGBTQ population remain in many places where they live, according to the negative sentiment scores. The sentiment ratings for each term are summarized in **Table 1**. Positive sentiments account for 487 of the total number of tweets into various emotion categories, including 251 for positive, 90 for joy, and 146 for trust. Moreover, negative sentiments account for a total of 428 comments, with negative accounting for 227, disgust accounting for 78, and anger accounting for 118. The negative sentiments imply that many continue to hold homophobic feelings against LGBTQ+ people. Although the negative comments are similar to the positive comments, it is crucial to educate individuals about their perceptions of LGBTQ+ persons in the areas in which they live. A more secure and peaceful society promotes development toward a more sustainable future and peaceful cohabitation.

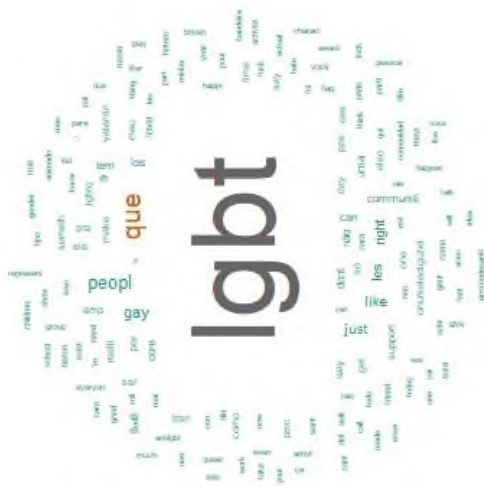


Figure 3. Word cloud analysis of the study.

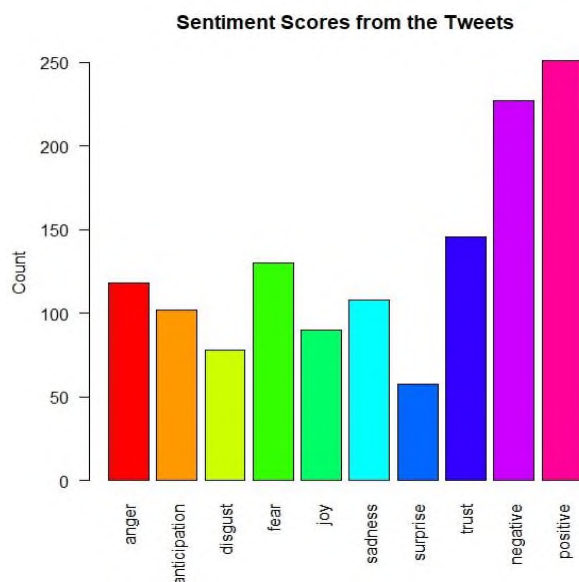


Figure 4. Sentiment scores and effects from the text data (tweets).

Table 1. Sentiment rating counts for each term.

Sentiment scores/effects	Counts
Anger	118
Anticipation	102
Disgust	78
Fear	130
Joy	90
Sadness	108
Surprise	58
Trust	146
Negative	227
Positive	251

## 4.2. Performance of the machine learning classifiers

The accuracy, recall, and precision performance metrics are used to evaluate empirically the experiences of LGBTQ+ persons using various machine learning algorithms. The logistic regression, support vector machine, naïve bayes, decision tree, random forest, and gradient boosting models were used to classify the sentiment analysis based on the dataset. The dataset was separated into two parts: training and testing. We utilized the 20% testing sets to evaluate performance after training the ML classifiers with 80% of the data. The **Table 2** below illustrates the performance of the machine learning classifiers applied in this study under the training stage. The decision tree classifier had the greatest accuracy with 100% accuracy, recall of 1.0, and precision of 1.0, followed by logistic regression with 97% accuracy. The support vector machine classifier came in third with 96% accuracy, and the naïve bayes classifier performed the worst during the training stage. The DT's performance scores (precision and recall) were substantial. The DT is clearly the highest performing algorithm among the several ML classifiers.

**Table 2.** Performance of the ML classifiers under the training stage.

Machine learning classifier	Accuracy (%)	Recall	Precision
Logistic regression	97.0	1.00	0.97
Support vector machine	96.0	1.00	0.96
Naïve bayes	88.0	0.87	1.00
Decision tree	100	1.00	1.00
Random forest	92.0	1.00	0.92
Gradient boosting	91.0	1.00	0.91

**Table 3** summarizes the performance of the ML classifiers in comparison to the testing set. The results show that the various ML classifiers used in this study produced significant outcomes. The test results indicate how well the model was trained, which is based on the amount of data, the predicted value, and the model features.

**Table 3.** Performance of the ML classifiers under the test stage.

Machine learning classifier	Accuracy (%)	Recall	Precision
Logistic regression	87.0	1.00	0.87
Support vector machine	87.0	1.00	0.87
Naïve bayes	70.0	0.75	0.89
Decision tree	88.0	0.97	0.90
Random forest	87.0	1.00	0.87
Gradient boosting	87.0	1.00	0.87

**Tables 4** and **5** compare the proposed method, which is lexicon-based and machine learning-based sentiment analysis, to similar prior research.

The results in **Tables 4** and **5** clearly show that the suggested approach outperforms strategic approaches. The machine learning classifiers used in this work are more accurate in determining the polarity of tweets. Furthermore, the lexicon-based SA more precisely classifies the emotions in tweets to examine the experiences of LGBTQ individuals based on the data. Even when different datasets are used, comparing results from

several text mining studies may provide substantial insights and aid in determining the generalizability of findings, as demonstrated in studies by Afrifa and Varadarajan<sup>[2]</sup> and Costola et al.<sup>[23]</sup>.

**Table 4.** A comparison of the lexicon-based approach.

Research	Year	Topic classification	Approach for SA
Ainapure et al. <sup>[15]</sup>	2023	Positive, negative, neutral	VADER and NRLex
Thangavel and Lourdasamy <sup>[49]</sup>	2023	Positive, neutral, negative	VADER
Velu et al. <sup>[50]</sup>	2023	Positive, negative, neutral	VADER
Proposed method		Positive, neutral, negative	Lexicon-based VADER

**Table 5.** A comparison of the ML classifier approach.

Research	Year	Approach for SA	Result (%)
Mutinda et al. <sup>[51]</sup>	2023	CNN	88.20
Kaur and Sharma <sup>[52]</sup>	2023	Random forest	82
Paramesha et al. <sup>[53]</sup>	2023	Logistic regression	55.3
Proposed method		Logistic regression, support vector machine, naïve bayes, decision tree, random forest, gradient boosting	97, 96, 88, 100, 92, 91

## 5. Conclusion and future works

With the increased usage of social media platforms, individuals are communicating with one another more frequently to share their views on topics such as sexual orientation, taxation, and many others. This study aims to accurately apply sentiment analysis on social media data from the perspective of lesbian, gay, bisexual, transgender, queer, and other (LGBTQ+) persons. The study used natural language processing and machine learning approaches to better understand the experiences of LGBTQ+ persons. Six different machine learning classifiers and lexicon-based sentiment analysis were utilized. For the study, 1000 tweets were retrieved from a total of 726,998 publicly available datasets. The models were trained using the dataset. Individuals are positive about LGBTQ concerns, according to the sentiment analysis results; yet, intolerance and harsh statements against the LGBTQ community persist in many regions where they dwell, according to the negative sentiment ratings. Furthermore, the results of the machine learning classifiers show that the suggested method significantly enhances sentiment analysis when compared to existing research. Additionally, the performance of the proposed machine learning classifiers is compared to state-of-the-art approaches for accuracy, recall, and precision performance assessment metrics. Based on the results, the performance evaluation measures are shown to significantly improve. The study's findings can help governments, non-governmental organizations, and right advocacy groups make educated decisions on LGBTQ+ issues. The proposed model in this study aims to do the following in the future: (1) perform sentiment analysis on the entire 726,998 publicly available dataset, (2) train the entire 726,998 tweets of the publicly available dataset using the proposed machine learning classifiers, (3) investigate the proposed model's performance using diverse datasets, (4) train the entire 726,998 tweets dataset using deep learning technologies, (5) employ ensemble techniques of stacking and voting classifier methods of the proposed models to train the entire dataset, and (6) explore more techniques in the present study domain.

## Author contributions

Conceptualization, VV, TZ, and PA; methodology, SA; software, SA; validation, VV, TZ and PA; formal analysis, SA and PA; investigation, VV, TZ, and PA; resources, VV; data curation, SA; writing—original draft

preparation, SA; writing—review and editing, TZ and SA; visualization, SA; supervision, VV, TZ, and PA; project administration, VV.

## Funding

No funding was received for conducting this study.

## Acknowledgments

The authors express their profound gratitude to Adwoa Afriye and Eric Afrifa for their encouragement throughout the study.

## Conflict of interest

The authors declare no conflict of interest.

## Data availability statement

The data used for the study are available in the FigShare data repository Afrifa, Stephen (2022): LGBTQ.csv. FigShare. Dataset. <https://doi.org/10.6084/m9.figshare.19787617.v1>.

## Ethical approval

Not applicable.

## References

1. Afrifa S, Varadarajan V, Appiahene P, et al. Ensemble machine learning techniques for accurate and efficient detection of botnet attacks in connected computers. *Eng* 2023; 4(1): 650–664. doi: 10.3390/eng4010039
2. Afrifa S, Varadarajan V. Cyberbullying detection on twitter using natural language processing and machine learning techniques. *International Journal of Innovative Technology and Interdisciplinary Sciences* 2022; 5(4): 1069–1080. doi: 10.1515/IJITIS.2022.5.4.1069-1080
3. Ahmed C, ElKorany A, ElSayed E. Prediction of customer’s perception in social networks by integrating sentiment analysis and machine learning. *Journal of Intelligent Information Systems* 2022; 1–27. doi: 10.1007/s10844-022-00756-y
4. Choudhary D. Security challenges and countermeasures for the heterogeneity of IoT applications. *Journal of Autonomous Intelligence* 2019; 1(2): 16–22. doi: 10.32629/jai.v1i2.25
5. Arcila-Calderón C, Amores JJ, Sánchez-Holgado P, Blanco-Herrero D. Using shallow and deep learning to automatically detect hate motivated by gender and sexual orientation on Twitter in Spanish. *Multimodal Technologies and Interaction* 2021; 5(10): 63. doi: 10.3390/mti5100063
6. Westwood S. Religious-based negative attitudes towards LGBTQ people among healthcare, social care and social work students and professionals: A review of the international literature. *Health & Social Care in the Community* 2022; 30(5): e1449–e1470. doi: 10.1111/hsc.13812
7. LGBTI Rights. Available online: <https://www.amnesty.org/en/what-we-do/discrimination/lgbti-rights/> (accessed on 20 July 2023).
8. Turner R, Hammersjö A. Navigating survivorhood? Lived experiences of social support-seeking among LGBTQ survivors of intimate partner violence. *Qualitative Social Work* 2023; 0(0): 1–19. doi: 10.1177/14733250221150208
9. Adu WK, Appiahene P, Afrifa S. VAR, ARIMAX and ARIMA models for nowcasting unemployment rate in Ghana using Google trends. *Journal of Electrical Systems and Information Technology* 2023; 10: 1–16. doi: 10.1186/s43067-023-00078-1
10. Appiahene P, Asare JW, Donkoh ET, et al. Detection of iron deficiency anemia by medical images: A comparative study of machine learning algorithms. *BioData Mining* 2023; 16(1): 1–20. doi: 10.1186/s13040-023-00319-z
11. Elmira WB, Hemmak A, Senouci B. Smart platform for data blood bank management: Forecasting demand in blood supply chain using machine learning. *Information* 2023; 14(1): 31. doi: 10.3390/info14010031
12. Chaganti R, Suliman W, Ravi V, Dua A. Deep learning approach for SDN-enabled intrusion detection system in IoT networks. *Information* 2023; 14(1): 41. doi: 10.3390/info14010041
13. Giannakas F, Kouliaridis V, Kambourakis G. A closer look at machine learning effectiveness in Android malware detection. *Information* 2023; 14(1): 2–24. doi: 10.3390/info14010002

14. Băroiu AC, Trăușan-Matu S. Comparison of deep learning models for automatic detection of sarcasm context on the MUSTARD dataset. *Electronics* 2023; 12(3): 666. doi: 10.3390/electronics12030666
15. Ainapure BS, Pise RN, Reddy P, et al. Sentiment analysis of COVID-19 tweets using deep learning and lexicon-based approaches. *Sustainability* 2023; 15(3): 2573–2593. doi: 10.3390/su15032573
16. Zhang H, Zhang D, Wei Z, et al. Analysis of public opinion on food safety in Greater China with big data and machine learning. *Current Research in Food Science* 2023; 6: 100468. doi: 10.1016/j.crfs.2023.100468
17. Çilgin C, BAŞ M, Bilgehan H, Ünal C. Twitter sentiment analysis during COVID-19 outbreak with VADER. *Academic Journal of Information Technology* 2022; 13(49): 72–89. doi: 10.5824/ajite.2022.02.001.x
18. Çilgin C, Gökçen H, z Gökşen Y. Sentiment analysis of public sensitivity to COVID-19 vaccines on Twitter by majority voting classifier-based machine learning. *Journal of the Faculty of Engineering and Architecture of Gazi University* 2023; 38(2): 1093–1104. doi: 10.17341/gazimmfd.1030198
19. Thai HH, Silhavy P, Kumar Dey S, et al. Analyzing public opinions regarding virtual tourism in the context of COVID-19: Unidirectional vs. 360-degree videos. *Information* 2023; 14(1): 11. doi: 10.3390/info14010011
20. Dai A, Hu X, Nie J, Chen J. Learning from word semantics to sentence syntax by graph convolutional networks for aspect-based sentiment analysis. *International Journal of Data Science and Analytics* 2022; 14(1): 17–26. doi: 10.1007/s41060-022-00315-2
21. Appiahene P, Missah YM, Najim U. Predicting bank operational efficiency using machine learning algorithm: Comparative study of decision tree, random forest, and neural networks. *Advances in Fuzzy Systems* 2020; 2020: 8581202. doi: 10.1155/2020/8581202
22. Mohammed AFY, Sultan SM, Lee Y, Lim S. Deep-reinforcement-learning-based IoT sensor data cleaning framework for enhanced data analytics. *Sensors* 2023; 23(4): 1791. doi: <https://doi.org/10.3390/s23041791>
23. Costola M, Hinz O, Nofer M, Pelizzon L. Machine learning sentiment analysis, COVID-19 news and stock market reactions. *Research in International Business and Finance* 2023; 64: 101881. doi: 10.1016/j.ribaf.2023.101881
24. Nigam N, Yadav D. Lexicon-based approach to Sentiment Analysis of tweets using R language. In: Singh M, Gupta P, Tyagi V, et al. (editors). *Advances in Computing and Data Sciences*, Proceedings of ICACDS 2018: 2nd International Conference on Advances in Computing and Data Sciences; 20–21 April 2018; Dehradun, India. Springer; 2018. pp. 154–164.
25. Rainer J, Vicini A, Salzer L, et al. A modular and expandable ecosystem for metabolomics data annotation in R. *Metabolites* 2022; 12(2): 173. doi: 10.3390/metabo12020173
26. Zeeshan, Ali Z, Jawad, Zakira M. Research Chinese-urdu machine translation based on deep learning. *Journal of Autonomous Intelligence* 2020; 3(2): 34–44. doi: 10.32629/jai.v3i2.279
27. Sharma S, Srinivas PYKL, Balabantaray RC. Text normalization of code mix and sentiment analysis. In: Proceedings of 2015 International Conference on Advances in Computing, Communications and Informatics (ICACCI); 10–13 August 2015; Kochi, India. pp. 1468–1473. doi: 10.1109/ICACCI.2015.7275819
28. Chennafi ME, Bedlaoui H, Dahou A, Al-qaness MAA. Arabic aspect-based sentiment classification using Seq2Seq dialect normalization and transformers. *Knowledge* 2022; 2(3): 388–401. doi: 10.3390/knowledge2030022
29. Yang L, Baratchi M, van Leeuwen M. Unsupervised discretization by two-dimensional MDL-based histogram. *arXiv* 2022; arXiv:2006.01893. doi: 10.1007/s10994-022-06294-6
30. Ren D, Srivastava G. A novel natural language processing model in mobile communication networks. *Mobile Networks and Applications* 2022; 27: 2575–2584. doi: 10.1007/s11036-022-02072-9
31. Ramaswamy SL, Chinnappan J. RecogNet-LSTM+CNN: A hybrid network with attention mechanism for aspect categorization and sentiment classification. *Journal of Intelligent Information Systems* 2022; 58: 379–404. doi: 10.1007/s10844-021-00692-3
32. Khan Z, Zakira M, Slamun W, Slamun N. A study of neural machine translation from Chinese to Urdu. *Journal of Autonomous Intelligence* 2019; 2(4): 29–36. doi: 10.32629/jai.v2i4.82
33. Mat Razali NA, Malizan NA, Hasbullah NA, et al. Opinion mining for national security: Techniques, domain applications, challenges and research opportunities. *Journal of Big Data* 2021; 8: 150. doi: 10.1186/s40537-021-00536-5
34. Chakravarthi BR. Multilingual hate speech detection in English and Dravidian languages. *International Journal of Data Science and Analytics* 2022; 14: 389–406. doi: 10.1007/s41060-022-00341-0.
35. Karsi R, Zaim M, El Alami J. Assessing naive bayes and support vector machine performance in sentiment classification on a big data platform. *IAES International Journal of Artificial Intelligence (IJ-AI)* 2021; 10(4): 990–996. doi: 10.11591/ijai.v10.i4.pp990-996
36. Sivakumar M, Uyyala SR. Aspect-based sentiment analysis of mobile phone reviews using LSTM and fuzzy logic. *International Journal of Data Science and Analytics* 2021; 12: 355–367. doi: 10.1007/s41060-021-00277-x
37. Okey OD, Maidin SS, Adasme P, et al. BoostedEnML: Efficient technique for detecting cyberattacks in IoT systems using boosted ensemble machine learning. *Sensors (Basel)* 2022; 22(19): 7409. doi: 10.3390/s22197409
38. Raj C, Agarwal A, Bharathy G, et al. Cyberbullying detection: Hybrid models based on machine learning and natural language processing techniques. *Electronics* 2021; 10(22): 2810. doi: 10.3390/electronics10222810
39. Muneer A, Fati SM. A comparative analysis of machine learning techniques for cyberbullying detection on

- Twitter. *Future Internet* 2020; 12(11): 187. doi: 10.3390/fi12110187
40. Borg A, Boldt M. Using VADER sentiment and SVM for predicting customer response sentiment. *Expert Systems with Applications* 2020; 162: 113746. doi: 10.1016/j.eswa.2020.113746
  41. Saba T, Khan SU, Islam N, et al. Cloud-based decision support system for the detection and classification of malignant cells in breast cancer using breast cytology images. *Microscopy Research and Technique* 2019; 82(6): 775–785. doi: 10.1002/jemt.23222
  42. Azeez NA, Idiakose SO, Onyema CJ, Van Der Vyver C. Cyberbullying detection in social networks: Artificial intelligence approach. *Journal of Cyber Security and Mobility* 2021; 10(4): 745–774. doi: 10.13052/jcsm2245-1439.1046
  43. Ahmed MT, Rahman M, Nur S, et al. Natural language processing and machine learning based cyberbullying detection for Bangla and Romanized Bangla texts. *TELKOMNIKA Telecommunication Computing Electronics and Control* 2022; 20(1): 89–97. doi: 10.12928/TELKOMNIKA.v20i1.18630
  44. Sarailidis G, Wagener T, Pianosi F. Integrating scientific knowledge into machine learning using interactive decision trees. *Computers & Geosciences* 2022; 170: 105248. doi: 10.1016/j.cageo.2022.105248
  45. Murorunkwere BF, Ihrwe JF, Kayijuka I, et al. Comparison of tree-based machine learning algorithms to predict reporting behavior of electronic billing machines. *Information* 2023; 14(3): 140. doi: 10.3390/info14030140
  46. Afrifa S, Zhang T, Appiahene P, Vijayakumar V. Mathematical and machine learning models for groundwater level changes: A systematic review and bibliographic analysis. *Future Internet* 2022; 14(9): 259. doi: 10.3390/fi14090259
  47. Junior MA, Appiahene P, Appiah O. Forex market forecasting with two-layer stacked Long Short-Term Memory neural network (LSTM) and correlation analysis. *Journal of Electrical Systems and Information Technology* 2022; 9: 14. doi: 10.1186/s43067-022-00054-1
  48. Abiola O, Abayomi-Alli A, Tale OA, et al. Sentiment analysis of COVID-19 tweets from selected hashtags in Nigeria using VADER and Text Blob analyser. *Journal of Electrical Systems and Information Technology* 2023; 10: 5. doi: 10.1186/s43067-023-00070-9
  49. Thangavel P, Lourdusamy R. A lexicon-based approach for sentiment analysis of multimodal content in tweets. *Multimedia Tools and Applications* 2023; 82: 24203–24226. doi: 10.1007/s11042-023-14411-3
  50. Velu SR, Ravi V, Tabianan K. Multi-lexicon classification and valence-based sentiment analysis as features for deep neural stock price prediction. *Sci* 2023; 5(1): 8. doi: 10.3390/sci5010008
  51. Mutinda J, Mwangi W, Okeyo G. Sentiment analysis of text reviews using Lexicon-Enhanced Bert Embedding (LeBERT) model with convolutional neural network. *Applied Sciences* 2023; 13(3): 1445. doi: 10.3390/app13031445
  52. Kaur G, Sharma A. A deep learning-based model using hybrid feature extraction approach for consumer sentiment analysis. *Journal of Big Data* 2023; 10: 5. doi: 10.1186/s40537-022-00680-6
  53. Paramesha K, Gururaj HL, Nayyar A, Ravishankar KC. Sentiment analysis on cross-domain textual data using classical and deep learning approaches. *Multimedia Tools and Applications* 2023; 82: 30759–30782. doi: 10.1007/s11042-023-14427-9