



Sentiment analysis on Bengali comments to predict emotions

Papia Akter¹, Rabeya Sultana²

¹ Lecturer, Department of Computer Science and Engineering, Prime University, Bangladesh

² Senior Lecturer, Department of Computer Science and Engineering, Prime University, Bangladesh

Abstract

Sentiment Analysis, as the term implies, is the process of determining the point of view or sentiment underlying a scenario. It simply involves analyzing and determining the emotion or intent behind a piece of writing, speech, or any other form of communication. Sentiment analysis is one of the fastest-developing areas of computer science research nowadays. Using sentiment analysis, a text may be separated into several different emotions. Because there have been few studies on emotion analysis in the Bengali language, it is recognized as an essential study subject in Bengali language analysis. This study addresses three types of emotions: negative, positive, and neutral. In this paper, we trained the dataset using various machine learning techniques and obtained accuracy of 70.87% and 73.95% for Linear and RBF Support Vector Machine (SVM), 74.60% for Logistic Regression (LR), 62.46% for Decision Tree (DT), 69.90% for Random Forest (RF), 75.40% for Naive Bayes (MNB), and 67.48% for K-Nearest Neighbors (KNN).

Keywords: Sentiment analysis, machine learning, TF-IDF, MNB

Introduction

Sentiment analysis, often known as opinion mining, is the study of analyzing human sentiments through systematic detection, extraction, measurement, and interpretation of the text using Natural Language Processing (NLP). It is, in turn, a technique for evaluating the emotional nature of a set of words, which is then utilized to provide a more accurate understanding of the attitudes, ideas, and sentiments stated in an online mention. Detecting underlying emotions has become essential in today's fast-paced environment. We may gain significant benefits in a variety of industries by recognizing the reason behind a human choice, which can range from purchasing a certain product to criminal motivation, employee-employer interactions, and so on. Sentiment analysis is a highly automated approach for understanding a point of view from spoken language, papers, or even video streams. Human emotions are biochemical states caused by neurophysiological changes in the nerve systems [1, 3] that are connected to behavioral reactions, sensations, and thoughts. Psychologists have sought to classify the numerous types of emotions that humans experience. A few distinct notions have emerged to define and characterize people's feelings. Six primary emotions have been recognized by psychologists: grief, happiness, anxiety, disgust, fury, and surprise. Over the last two decades, research on human reactions has blossomed, with contributions from a wide range of disciplines, including psychology, psychiatry, affective biology, sociology of feelings, computer science, medicine, history, and medicine. However, when academics tried to recognize these fundamental human emotions from human-written text, voice, and facial expressions, the phrase sentiment analysis was developed in computer science. The first academic studies of public opinion were conducted during and after WWII, with a primarily political motivation. However, modern sentiment research did not become popular until the mid-2000s and is mostly focused on online

product ratings [4]. Following the epidemic, sentiment analysis has been used in a variety of other fields, including stock market forecasts and terrorist attack reactions. Furthermore, sentiment analysis is beneficial for social media monitoring since it allows us to observe how the general population feels about a certain topic. It has several uses and is quite efficient. Organizations all across the world are seizing the chance to get insight from social data. However, after reviewing several research works on human emotion analysis, we discovered that, due to the language's complexity, there have been comparatively few studies on emotion analysis in the Bangla language, although multiple works have been done for English and other languages [5], [6]. With around 250 million native speakers globally, Bangla is the world's fourth-most widespread language. We prioritized utilizing Bangla to extract emotion because it is one of the most popular languages and about 250 million people use it to communicate their feelings. Detecting slang or harsh terms in Bangla, on the other hand, is relatively more difficult due to a lack of resources than in English [7, 8]. Furthermore, religious differences exist among Bengali speakers. People are using social media to malign many religious beliefs, and religious violence is on the rise. People generally share their thoughts and ideas through social media platforms such as Facebook, Twitter, and others, and a substantial quantity of data on numerous subjects is available [9]. In this research, we offer a technique for predicting sentiment from Bengali text using the SVM algorithm, inspired by the use of sentiment analysis. The remainder of our study will be organized as follows: Section 2 evaluated the majority of the relevant literature reviews for our investigation. Section 3 then provided a brief overview of the suggested methodology as well as an explanation of subject modeling techniques. Section 4 then shows the results of the investigations conducted for this study. Finally, Section 5 brings the piece to a close by bringing it to a close.

Related Works

Among the many study fields in the field of computer science, the analysis of sentiment is one of the most rapidly developing. According to a study, the beginnings of sentiment analysis may be traced to surveys on public opinion conducted at the turn of the 20th century and text analysis undertaken by a computational linguistics group in the 1990s^[4]. 118-124,119 M. S. S. Khan *et al.*/JEA Vol. 02 (03) 2021 With the availability of messages on the internet, computer-based sentiment analysis came closer to the epidemic. This pandemic had a significant influence on a variety of professions by determining the underlying mood in any text or voice transmission. Many articles have addressed the subject of sentiment analysis in various ways. Mahmudun M. *et al.*^[9] employed the concepts of TF (term frequency) and IDF (inverse document frequency) values to achieve a better answer, and they extracted the various features of negative, positive, or neutral terms from Bangla text to reach a more accurate result. Das. A.^[10] *et al.* proposed a hybrid approach for extracting views from text (Bangla and English) by combining rule-based and automated systems. This system was built using Natural Language Processing (NLP) and SVM. Das D.^[11] suggested an emotion monitoring system for a specific subject or event that utilized SentiWordNet for both Bangla and English text and sense-based affect scoring algorithms. To create their approaches for identifying sentiment from Bangla text, Hasan K.A. *et al.*^[12] employed SentiWordNet and WordNet to evaluate the polarity and meaning of terms in the text. Chowdhury S. *et al.*^[13] employed the SVM and MaxEnt (Maximum Entropy) algorithms for automatically extracting emotions from Bangla Microblog (Twitter) messages, regardless of text polarity. Go A. *et al.*^[14] employed three machine learning algorithms to identify the moods of Twitter tweets using emoticons: SVM, MaxEnt, and Naive Bayes. Pandey P. *et al.*^[15] suggested a technique for

sentiment analysis of Hindi movie reviews using HindiSentiWordNet (HSWN) and the Synset replacement algorithm. Tuhin R. A. *et al.*^[16] suggested two machine learning algorithms for extracting emotion from any Bangla text: the Naive Bayes Classification Algorithm and the Topical Method. The strategies proposed have been utilized at the article and sentence levels. Tembhumkar S. D. *et al.*^[17] created a sentiment analysis model based on LDA, which they utilized to score tweets in terms of popularity. Umamaheswari K. *et al.*^[18] identified the perspective from the IMDB movie analysis dataset using SVM and LDA.

Methodology

In this work, supervised learning—one of the three major categories of machine learning—is utilized. The best strategy is to use the supervised learning method for text classification and multiclass-labeled datasets. Due to the fact that supervised text classification requires a pre-labeled dataset with accurate values, free text, image, or video data are again incomprehensible to machines. It can only read 1s and 0s. The dataset must be changed or encoded to make it machine-readable. In this instance, the Term Frequency-Inverse Document Frequency (TF-IDF) function from the Python library "Scikit-learn" is utilized. Therefore, this study's approach calls for cleaning the dataset and vectorizing it. The Support Vector Machine (SVM) classifier is used to teach the computer a dataset. SVM produces findings that are far superior to those of other classifiers utilized in this study, including K-Nearest Neighbors, Naive Bayes, Random Forest, and Logistic Regression. Two measures are typically used in supervised learning-based sentiment analysis to classify documents. It learns a model using the training data, as depicted in Fig. 1. By predicting the target class of ambiguous test data, the trained classifier assesses the model's accuracy during testing.

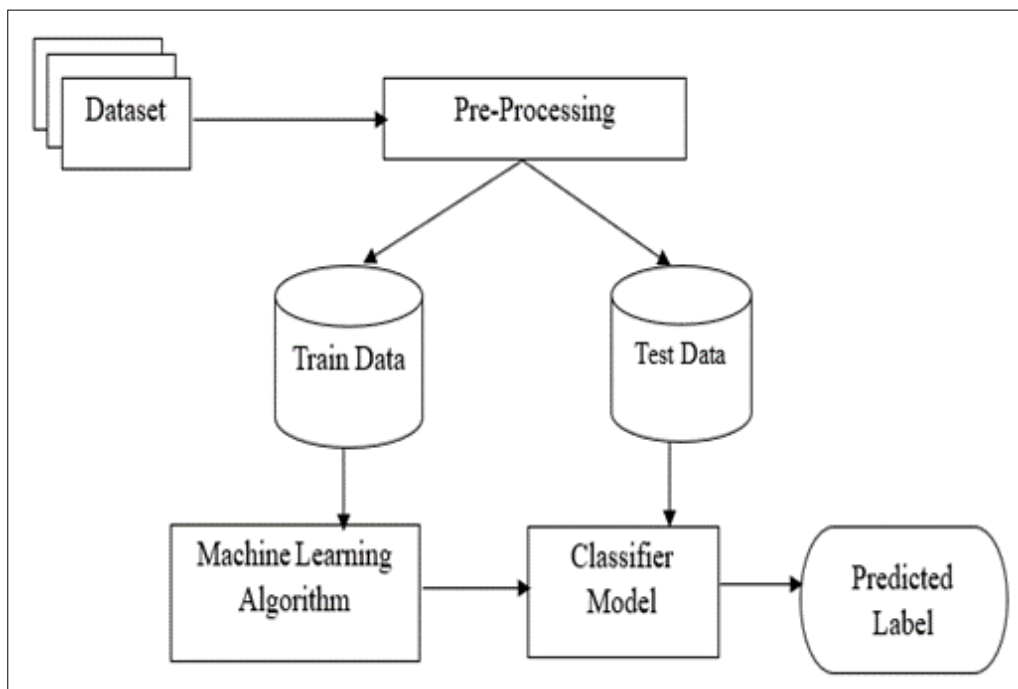


Fig 1: Working Mechanism of proposed method

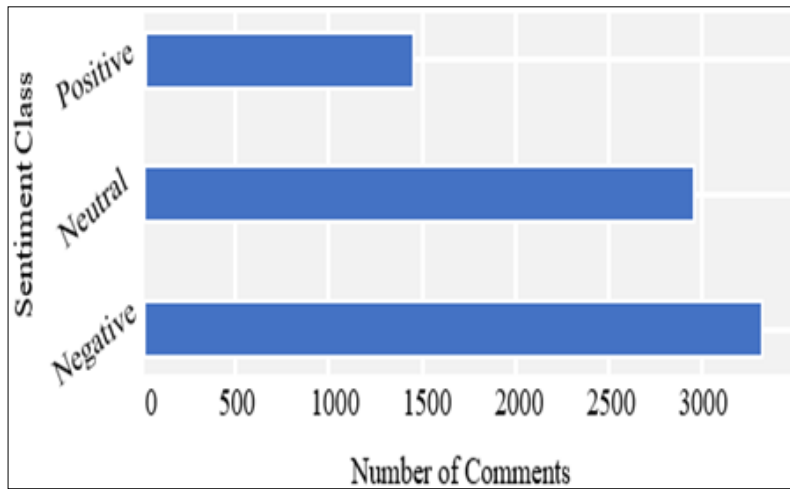


Fig 2: Dataset Distribution

Table 1: Dataset Sample

	Name	Comment Link	Tag	Comment
0	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/pos...	Negative	আমার টাকা চোর, আমার টাকা ফেরত দেন,, ইভ্যালিতে...
1	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/pos...	Positive	তুমি অসম্ভব সুন্দর, এটা বলতে কোন দিধা নেই.....
2	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/pos...	Neutral	মনের মধ্যে কারো প্রতি ক্ষোভ রেখে নিজেকে কষ্ট দ...
3	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/pos...	Negative	চুল ছোট রাখাকে তোমরা ফ্যাশন বানাই ফেলছ আর এই চুল...
4	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/pos...	Neutral	যার কথা ভাসে, মেঘলা বাতাসে তবু সে দূরে তা মানি.....
5	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/pos...	Neutral	খারাপ মন্তব্যে রিপ্লাই না দেয়াই উত্তম তাতে খার....
6	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/vid...	Neutral	আপু, চেহারায় একটু মলিনতা দেখা যাচ্ছে। জন্মদিনে...
7	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/vid...	Positive	শুভ জন্মদিন জন্মদিনের শুভেচ্ছা ও ভালোবাসা অভির.....
8	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/vid...	Neutral	ভালো মানুষের পেছনে অবশ্যই খারাপ মানুষ থাকবেই.....
9	Sabnam Faria	https://www.facebook.com/FariaSabnamTripti/vid...	Positive	অনেক সুন্দর লাগছে। আমি সাউথ আফ্রিকা থেকে।

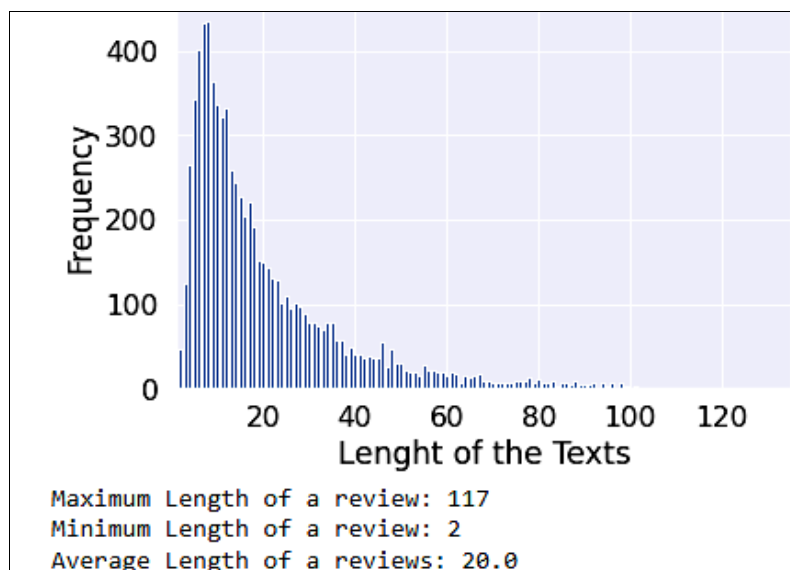


Fig 3: Length-Frequency Distribution

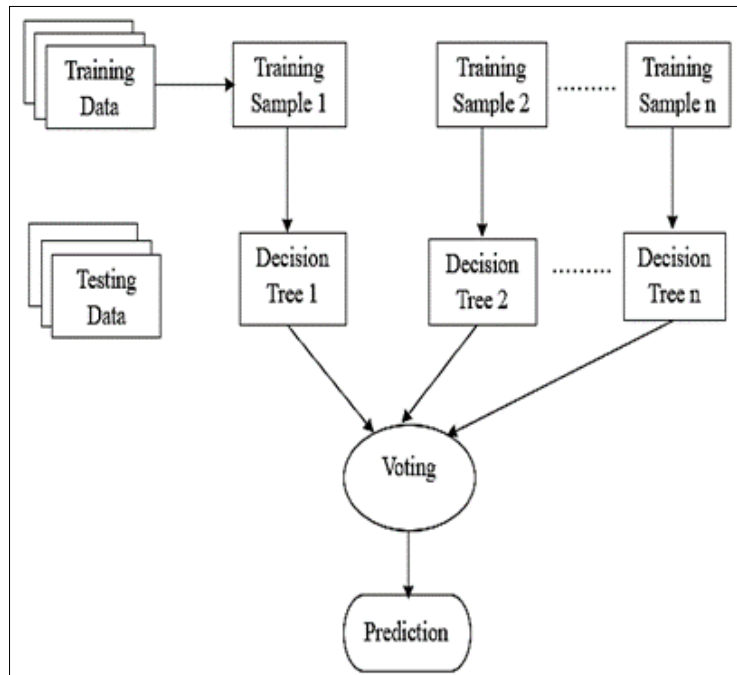


Fig 4: Working flow of Random Forest

1. Data Collection

Data was gathered from the Facebook profiles of several Bangladeshi celebrities, including a cricketer, an actor, an actress, a Youtuber, and others. The data for this study came from broad public comments on celebrity Facebook sites. Facepager was used to extract comments from some pages' postings since it can retrieve data that is publicly accessible from Facebook, YouTube, Twitter, and other websites via a JSON-based API. The access token is needed to extract data using a facepager. When a user profile is logged in using this tool, the token is produced automatically. A website named 'Graph API Explorer: Facebook Developer' is another alternative for collecting access tokens. We can use the Graph API to obtain and send data to and from the Facebook social network. Facepager extracts data by obtaining the ID of the required pages or profiles. It retains the raw data in an SQLite database after extraction, and the data may subsequently be exported as a CSV Excel document [9]. As raw data, we collected 10,000 comments. Table 1 displays several examples from my dataset, representing a brief scenario of the dataset.

2. Data Cleaning

The collected raw data contains a variety of comments that need just pure Bangla text. It signifies that a paragraph or comment should be written entirely in Bengali typefaces, with no other characters. People's remarks in various languages or a blend of different language words in a single comment are common in social media. Because just Bangla comments are required, extraneous data from the input dataset should be cleaned and removed as these data would increase noise and degrade model accuracy. The English and multilingual comments were the first to be deleted, leaving around 7,647 remarks. Following this procedure, the dataset contained either pure Bangla text or Bangla text combined with certain characters. Second, the remaining comments were cleaned using criteria such as each character's Unicode value in a comment should be between 2432 and 2559 or 32. The Unicode values 2432 to 2559 are for Bangla typefaces, and 32 is for space. In this study, the

dataset was manually tagged based on 7,647 data points. There are three types of people: negative, positive, and neutral.

3. Data Transformation

The dataset has been transformed or encoded and features have been extracted using TF-IDF. This is a typical text-to-feature vectors conversion approach that generates a numerical representation of the pattern that may be used to adapt machine learning algorithms for prediction [19]. A word's TF-IDF score is determined by multiplying two separate measures [20]. The 'Term Frequency' calculates how many times the term (wi) appears in a text using Eq. (1).

$$TF_{i,j} = \frac{\text{word } i \text{ frequency in sentence } j}{\text{total words in sentence } j} \tag{1}$$

The 'Inverse Document Frequency' indicates how often or unusual a term is across the whole dataset. It is computed using the following equation (2):

$$DF_i = \log\left(\frac{\text{total number of sentence}}{\text{number of sentences contain } i}\right) \tag{2}$$

where the logarithm of the total number of documents or sentences divided by the number of sentences containing a word (wi) is the IDF value of a word.

$$Score_{i,j} = TF_{i,j} \times IDF_i \tag{3}$$

where score i.j is a number value indicating the importance of word (wi). As a result, the significance of a keyword phrase may be evaluated by comparing its frequency in huge collections of texts. Tokens are generated from each comment or sentence in a dataset via TF-IDF, and each unique token is assigned a feature index. Finally, each phrase is converted into a vector, and the weighted integers of each vector indicate the feature score.

4. Proposed Classifiers

4.1 Support Vector Machine (SVM)

based models. This classifier divides n-dimensional space into classes by creating the decision border, allowing fresh data points to be readily classified in the future [21]. SVM is a two-class model by definition, which implies it can discriminate between two separate classes. There are several hyperplanes that may be used to split two classes.

The best hyperplane has the greatest distance between two classes. Because this study involved seven classes, the SVM used the One vs. Rest (OVR) strategy. The multiclass dataset is divided into many binary classification issues. A binary classifier is trained for each binary classification task, and predictions are produced using the most confident model [22]. My classes, for example, were labeled as Negative, Positive, and Neutral. Negative versus [Positive and Neutral] is one binary categorization. For each classifier, one class is fitted against all the other classes. It gives the multiclass problem a large computing benefit [9].

4.2 Naive Bayes (NB)

Classification with Naive Bayes is straightforward and based on event probability. Despite its simplicity, it excels at text categorization issues such as sentiment analysis [9]. It is a statistical classification approach based on the Bayes Theorem, which states:

$$P(X|Y) = P(X) / P(Y|X) P(Y) \quad (4)$$

Where

$P(X)$ =probability of X (hypothesis)

$P(Y)$ =probability of Y (data)

$P(Y|X)$ =probability of Y given X that the hypothesis X is true

$P(X|Y)$ =probability of X given Y

The Multinomial NB classifier is employed in this study since it is appropriate for this type of research. The multinomial distribution normally requires integer feature counts. In reality, fractional counts such as TF-IDF can also be used [23].

4.3 Random Forest (RF)

The Random Forest method produces decision trees based on data samples, forecasts each one, and votes on the best answer. It's an ensemble strategy that outperforms a single decision tree by averaging the results to reduce over-fitting [24]. The Random Forest classifier's workflow is depicted in Fig 4.

4.4 K-Nearest Neighbors (KNN)

The KNN method classifies data by finding the K closest matches in training data and then predicts based on the class label of the closest matches. The Euclidean distance has traditionally been employed to locate the closest match. KNN requires a numerical representation of words to compute distance and make predictions, which may be acquired via TF-IDF. The distance between all of these feature vectors created by TF-IDF in the data set will be calculated using the unlabeled text data feature vector. The K-nearest vectors will be selected from amongst them, and the class with the greatest amount of frequencies will be tagged to the unlabeled results.

4.5 Logistic Regression (LR)

Despite its name, logistic regression is a classification model rather than a regression model. For binary and linear classification issues, logistic regression is a simpler and more efficient technique. It is a classification model that is simple to implement and delivers excellent results with linearly separable classes. It is a widely used classification method in business. Like the Adaline and perceptron, the logistic regression model is a statistical approach for binary classification that may be applied to multiclass classification. Scikit-learn includes a highly efficient logistic regression implementation that can handle multiclass classification tasks.

4.6 Decision Tree (DT)

The computational power of any decision tree classifier is based on the notion of splitting criteria. Decision trees are shown in a tree form, similar to a flow chart, with instances categorized based on their feature values. A node in a decision tree represents an instance, a branch reflects the results of the test, and the leaf node represents the class name.

1. Results

To better illustrate the suggested classifiers, four separate metrics are used to analyze and comprehend my model's performance: Accuracy, Precision, Recall, and F1-Score. In each measure, 5 classifiers and SVM were used to see which one worked best.

1.1 Accuracy

Accuracy is defined as the ratio of accurately predicted True Positives and True Negatives to the whole test dataset. Eq. (5) may be used to compute the number of correctly identified samples out of the total samples.

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (5)$$

To carry out the experiment, the dataset was divided into two parts: 90% for training and 10% for testing. The obtained accuracy of 70.87% and 73.95% for Linear and RBF Support Vector Machine (SVM), 74.60% for Logistic Regression (LR), 62.46% for Decision Tree (DT), 69.90% for Random Forest (RF), 75.40% for Nave Bayes (MNB), and 67.48% for K-Nearest Neighbors (KNN) after training the algorithms with 3 class labeled dataset. From Fig.5, MNB worked better than the rest of the algorithms.

1.2 Precision

The precision measure represents the class's accuracy. This metric determines whether or not the positive class prediction is correct [9]. The highest score is 100 if the classifier properly identifies all positive values. Eq. (6), on the other hand, may be used to compute precision.

$$Precision = \frac{TP}{TP + FP} \quad (6)$$

Which is the proportion of positively predicted observations that were successfully anticipated to the total number of positively predicted observations? As a result, a higher accuracy score indicates a lower false positive rate.

Figure 6 demonstrates that MNB has a high overall Precision for all three classes. Decision Tree classifier, on the other hand, had a lower overall accuracy value than MNB, but it could categorize all of the classes. Precision alone is ineffective owing to its disregard for the negative class. It is frequently used in conjunction with the Recall metric.

1.3 Recall

Recall is defined as the ratio of accurately predicted positive observations to all observations in the actual class [9]. In contrast to precision, recall may be determined using Eq. (7) where FN is taken into account.

$$Recall = \frac{TP}{TP + FN} \tag{7}$$

This metric measures the model's ability to distinguish positive classifications. As a result, a greater recall score implies a reduced false-positive rate. Figure 7 shows that MNB has an overall higher recall score than other classifiers.

1.4 F1 Score

The F1 Score is the weighted average of Precision and Recall. As a result, this score considers both false positives

and false negatives. Although f1 is less intuitive than accuracy, it is more beneficial when class distribution is unequal [25, 26]. When the cost of false positives and false negatives is the same, accuracy improves. If they don't, Precision and Recall can assist, which is why F1 Score was created. Unlike precision, F1 score may be determined using Eq. (8), which takes FN into account.

$$F1\ score = 2 \times \frac{(Precision \times Recall)}{Precision + Recall} \tag{8}$$

This study employed a weighted average to construct F1 Score because the dataset was uneven, and it additionally prioritized some predictions depending on their percentage. It computes the weighted mean of the measurements. Each weight reflects the number of samples in that class. Still, this study doesn't claim which classifier is superior because a greater F1-score does not always imply a better classifier. It combines Precision and Recall to assess overall model correctness. As a result, a high f1 score implies a high Precision and Recall.

As shown in Fig. 8, F1 score 72.70, 60.68, 65.93, 73.79, 66.78, 65.52, and 71.30 for LR, DT, RF, MNB, KNN, Linear SVM, and RBF SVM respectively. MNB has a higher score than others, which means it predicted 73.79% of data correctly, considering both false positive and false negative values.

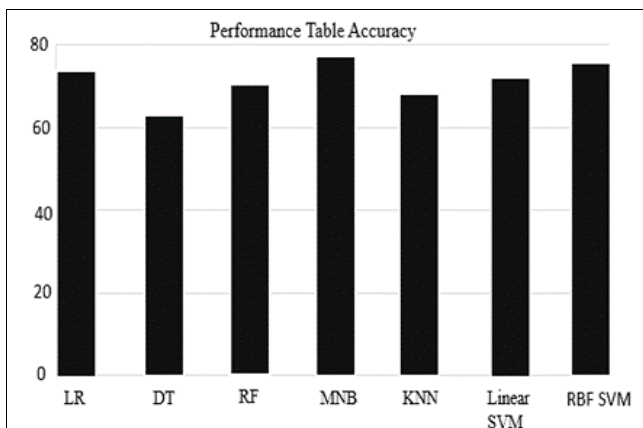


Fig 5: Accuracy comparison between classifiers

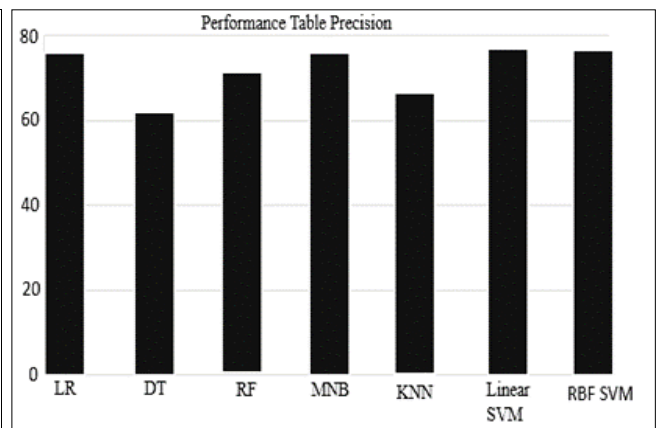


Fig 6: Precision Score for Classifiers

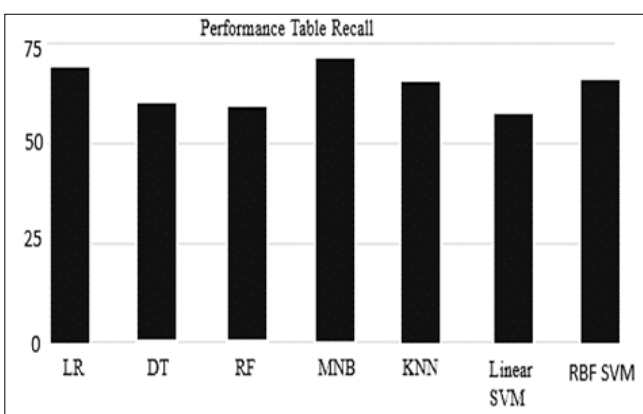


Fig 7: Recall Score for Classifiers

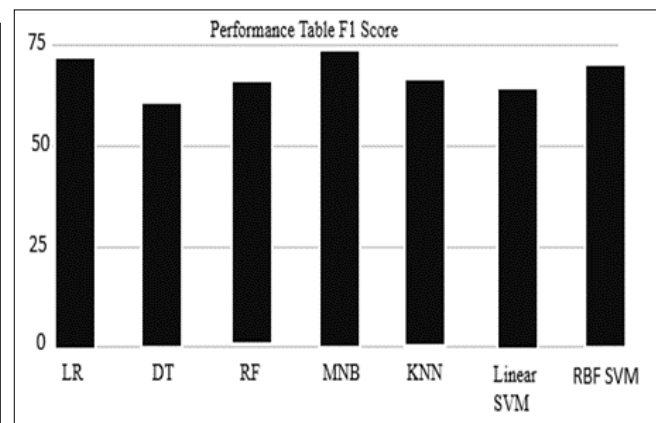


Fig 8: F1-Scores for classifiers

2. Performance Table for N-gram Feature

An n-gram model is a probabilistic language model that predicts the next item in a sequence using a (n-1) order Markov model. In probability, communication theory,

computational linguistics (for example, statistical natural language processing), computational biology (for example, biological sequence analysis), and data compression, n-gram models are currently commonly utilized. The simplicity and

scalability of n-gram models (and algorithms that employ them) are two advantages. With bigger n, a model may hold

more context with a well-understood space-time tradeoff, allowing modest experiments to scale up effectively.

Table 2: Performance table for Unigram feature

```

===== Performace Table for Unigram feature:=====
      Accuracy Precision Recall F1 Score Model Name
0      74.60    75.72  69.90   72.70      LR
1      62.46    61.51  59.87   60.68      DT
2      69.90    72.87  60.20   65.93      RF
3      75.40    76.16  71.57   73.79      MNB
4      67.48    66.01  67.56   66.78      KNN
5      70.87    76.68  57.19   65.52  Linear SVM
6      73.95    76.34  66.89   71.30      RBF SVM
    
```

Table 3: Performance table for Bigram feature

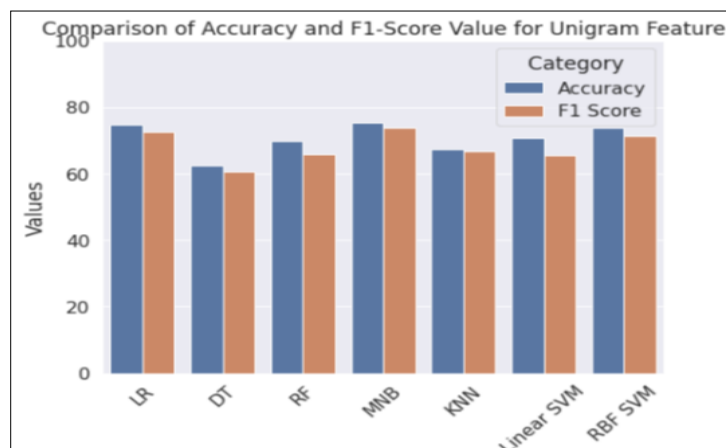
```

===== Performace Table for Bigram feature:=====
      Accuracy Precision Recall F1 Score Model Name
0      74.11    77.04  66.22   71.22      LR
1      64.89    63.31  65.22   64.25      DT
2      71.84    77.06  59.53   67.17      RF
3      74.11    74.22  71.24   72.70      MNB
4      65.70    64.36  65.22   64.78      KNN
5      60.36    96.55  18.73   31.37  Linear SVM
6      68.45    83.33  43.48   57.14      RBF SVM
    
```

Table 4: Performance table for Trigram feature

```

===== Performace Table for Trigram feature:=====
      Accuracy Precision Recall F1 Score Model Name
0      72.01    76.69  60.54   67.66      LR
1      62.62    62.06  58.53   60.24      DT
2      70.87    77.67  55.85   64.98      RF
3      74.11    74.39  70.90   72.60      MNB
4      65.53    64.33  64.55   64.44      KNN
5      51.78   100.00   0.33    0.67  Linear SVM
6      60.36    93.55  19.40   32.13      RBF SVM
    
```



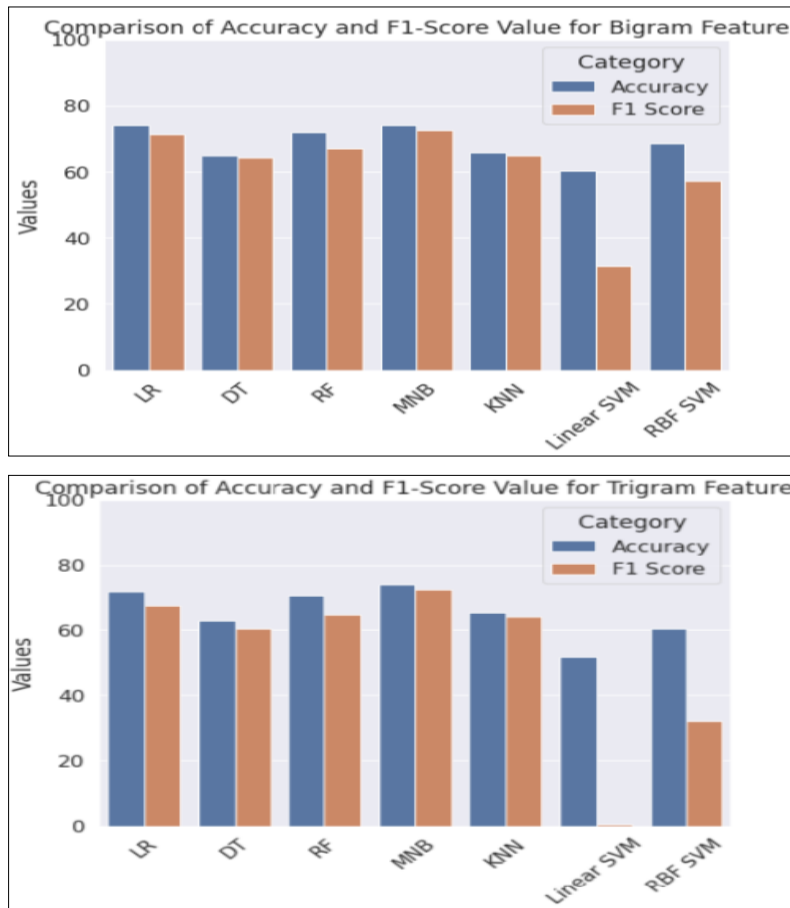


Fig 9: Comparison of Accuracy and F1-Score Value for Unigram, Bigram and Trigram Feature

From Table II, the obtained result is

Highest Accuracy achieved by MNB at = 75.4
 Highest F1-Score achieved by MNB at = 73.79
 Highest Precision Score achieved by Linear SVM at = 76.68
 Highest Recall Score achieved by MNB at = 71.57

From Table III, the obtained result is

Highest Accuracy achieved by LR at = 74.11
 Highest F1-Score achieved by MNB at = 72.7
 Highest Precision Score achieved by Linear SVM at = 96.55
 Highest Recall Score achieved by MNB at = 71.24

From table IV, the obtained result is:

Highest Accuracy achieved by MNB at = 74.11
 Highest F1-Score achieved by MNB at = 72.6
 Highest Precision Score achieved by Linear SVM at = 100.0
 Highest Recall Score achieved by MNB at = 70.89

3. Analysis

Based on the four distinct types of measurements discussed above, it is evident that MNB is the greatest fit for this research because it performed better than other classifiers. This suggests that MNB can be used in the future to anticipate human sentiment toward someone or anything. In general, the goal of this study is to extract emotion from the comments on various social media posts where individuals might remark on something. Many comments may be collected automatically from social media sites, and the suggested algorithm MNB can readily detect a human expression or a statement. Which algorithm must first be taught using sample data? As a result, in the future, detecting human sentiment may be conceivable using a fully automated approach.

Conclusion

Sentiment analysis is crucial in extracting useful information from people's thoughts. It assists us with extracting human emotions from voice, literature, and even facial expressions. As a result, extracting perspectives from text, voice, and facial expressions is a crucial task in a variety of industries such as social media monitoring, product analysis, customer reviews analysis, clinical service analysis, and so on. Inspired by these requirements, this research article focuses on extracting emotion from Bangla text. People nowadays are increasingly interested in social media, and they use it to voice their thoughts on a variety of issues.

As a result, sentiment analysis might be useful in gaining a deeper knowledge of people's perspectives in order to deliver better service. The goal of this project was to collect comments from various Facebook pages of Bangladeshi celebrities and create a model that can predict emotion based on Bengali language.

The goal of this project was to collect comments from various Facebook pages of Bangladeshi celebrities and create a model that can predict emotion based on Bengali language. Using the Multinomial Naïve Bayes(MNB) machine learning technique, we were able to identify five different emotions from Bangla text with 75.40% accuracy. Furthermore, the key rationale for selecting MNB as the principal model is that it outperforms other models discussed in the preceding section, such as Random Forest (RF), K-Nearest Neighbors (KNN), SVM and Neural Network.

References

1. Panksepp J. Affective neuroscience: The foundations of human and animal emotions. Oxford university press, 2004.
2. Damasio A R. Emotion in the perspective of an integrated nervous system. Brain research reviews,1998;26(2-3):83-86.
3. Ekman PE, Davidson RJ. The nature of emotion: Fundamental questions. Oxford University Press, 1994.
4. Viking Mäntylä M, Graziotin D, Kuutila M. The Evolution of Sentiment Analysis-A Review of Research Topics, Venues, and Top Cited Papers. arXiv e-prints, 2016, 1612.
5. Drovu MD, Chowdhury M, Uday SI, Das AK. Named entity recognition in bengali text using merged hidden markov model and rule base approach. In 2019 7th International Conference on Smart Computing & Communications (ICSCC), IEEE, 2019, 1-5.
6. Hossain MT, Hasan MW, Das AK. Bangla Handwritten Word Recognition System Using Convolutional Neural Network. In 2021 15th International Conference on Ubiquitous Information Management and Communication (IMCOM), IEEE, 2021, 1-8.
7. Mumu TF, Munnii IJ, Das A.K. Depressed people detection from bangla social media status using lstm and cnn approach. Journal of Engineering Advancements,2021;2(01):41- 47.
8. Islam J, Mubassira M, Islam MR, Das AK. A speech recognition system for Bengali language using recurrent neural network. In 2019 IEEE 4th international conference on computer and communication systems (ICCCS), IEEE, 2019, 73-76.
9. Mahmudun M, Altaf MT, Ismail S. Detecting sentiment from Bangla text using machine learning technique and feature analysis. Int. J. Comput. Appl., 975, 2016, 8887.
10. Das A, Bandyopadhyay S. Phrase-level polarity identification for Bangla. Int. J. Comput. Linguist. Appl.(IJCLA),2010;1(1-2):169-182.
11. Das D. Analysis and tracking of emotions in english and bengali texts: a computational approach. In Proceedings of the 20th international conference companion on World wide web, 2011, 343-348.
12. Hasan KA, Rahman M. Sentiment detection from bangla text using contextual valency analysis. In 2014 17th International Conference on Computer and Information Technology (ICCIT), 2014, 292-295.
13. Chowdhury S, Chowdhury W. Performing sentiment analysis in Bangla microblog posts. In 2014 International Conference on Informatics, Electronics & Vision (ICIEV), IEEE, 2014, 1-6.
14. Go A, Bhayani R, Huang L. Twitter sentiment classification using distant supervision. CS224N project report, Stanford,2009;1(12):2009.
15. Pandey P, Govilkar S. A framework for sentiment analysis in Hindi using HSWN. International Journal of Computer Applications, 2015, 119(19).
16. Tuhin RA, Paul BK, Nawrine F, Akter M, Das AK. An automated system of sentiment analysis from bangla text using supervised learning techniques. In 2019 IEEE 4th International Conference on Computer and Communication Systems (ICCCS), IEEE, 2019, 360-364.
17. Tembhurnikar SD, Patil NN. Sentiment Analysis using LDA on Product Reviews: A Survey. International Journal of Computer Applications, 975, 8887.
18. Umamaheswari K, Karthiga R. Sentiment Classification based on Latent Dirichl *et al* location. International Journal of Computer Applications, 975, 2015, 8887.
19. Rakib OF, Akter S, Khan MA, Das AK. and Habibullah, K.M. Bangla word prediction and sentence completion using GRU: an extended version of RNN on N-gram language model. In 2019 International Conference on Sustainable Technologies for Industry 4.0 (STI), IEEE, 2019, 1-6.
20. Monkeylearn.com,2021.[Online].Available:https://monkeylearn.com/blog/what-is-tf-idf/
21. www.javatpoint.com,2021.[Online].Available:https://www.javatpoint.com/machine-learning-support-vector-machine-algorithm.
22. Hossain MM, Labib MF, Rifat AS, Das AK, Mukta M. Auto-correction of english to bengali transliteration system using levenshtein distance. In 2019 7th International Conference on Smart Computing & Communications (ICSCC), IEEE, 2019, 1-5.
23. Labib MF, Rifat AS, Hossain MM, Das AK, Nawrine F. Road accident analysis and prediction of accident severity by using machine learning in Bangladesh. In 2019 7th International Conference on Smart Computing & Communications (ICSCC), IEEE, 2019, 1-5.
24. Emon EA, Rahman S, Banarjee J, Das AK, Mitra T. A deep learning approach to detect abusive bengali text. In 2019 7th International Conference on Smart Computing & Communications (ICSCC), IEEE, 2019, 1-5.
25. Ullah MR, Bhuiyan MAR, Das AK. IHEMHA: Interactive healthcare system design with emotion computing and medical history analysis. In 2017 6th International Conference on Informatics, Electronics and Vision & 2017 7th International Symposium in Computational Medical and Health Technology (ICIEV-ISCMHT), IEEE, 2017, 1-8.
26. Bhuiyan M, Rahman A, Ullah M, Das AK. iHealthcare: Predictive model analysis concerning big data applications for interactive healthcare systems. Applied Sciences,2019;9(16):3365.