

Sentiment Analysis of Nepali COVID19 Tweets Using NB, SVM AND LSTM

Milan Tripathi

Student, Department of Computer Engineering, Advanced College of Engineering and Management, Institute of Engineering, Tribhuvan University, Nepal

E-mail: milan.tripathi@acem.edu.np

Abstract

The government's months-long total lockdown in response to the COVID19 outbreak has resulted in a lack of physical connection with others. This resulted in a massive increase in social media communication. Twitter has become one of the most popular places for people to communicate their thoughts and opinions. As a result, massive amounts of data are created every day. These data can assist businesses in making better judgments. In the case of Nepal, there has been relatively little investigation into the text's analysis. Because few researchers are working in the field, development is slow. In this study, Four language-based models for sentiment analysis of Nepali covid19 tweets are designed and evaluated. Because the number of individuals using social media is expected to skyrocket in the next few days, companies will benefit from an AI-based sentiment analysis system. It will greatly assist firms in adapting to the changing climate.

Keywords: COVID19, social media, Twitter, sentiment analysis, artificial intelligence

1. Introduction

With the surge of COVID19 cases and the government's months-long complete lockdown. Citizens from all around the world were severely impacted. To keep the country's economy afloat, educational institutions, companies, and government offices must rely on modern technologies to carry out daily duties. Technology has been utilized to its utmost

potential. In the past, the use of social media to communicate was on the rise, and this shutdown provided it a big push. Twitter is one of the fastest-growing social media sites where individuals share their opinions on a variety of topics. As a result, a massive quantity of data has been created on this platform. As a result, comprehending this data will assist companies in making better judgments. Any stakeholder should be able to deduce the mood and significance of the fine-grained user evaluations.

COVID-19 attitudes changed quickly in response to the World Health Organization's (WHO) pandemic announcement and government-led measures to combat the illness. During the COVID-19 pandemic, this research [1] looked at global patterns in four emotions (fear, anger, sorrow, and joy) as well as the narratives that underpin those feelings. Using the search terms "Wuhan", "corona", "nCov", and "covid" researchers gathered almost 20 million social media tweets generated during the early stages of the COVID-19 epidemic from 28 January to 9 April 2020. Throughout the epidemic, public emotions moved dramatically from fear to rage, with grief and joy also surfacing. Fears about COVID-19 test shortages and medical supply constraints were more common conversation themes, according to word clouds. Anger moved from xenophobia at the start of the epidemic to a discussion over stay-at-home warnings as the virus progressed. Issues linked to sadness included the loss of friends and family members, while topics relating to pleasure included thankfulness and excellent health. Over a few weeks, worldwide COVID-19 opinions have changed dramatically. The findings point to the emergence of emotion-driven communal concerns centered on shared public distress experiences of the COVID-19 pandemic, which include large-scale social isolation and human life loss. The constant growth in social worries reflected in negative emotions must be monitored and regulated by combining frequent crisis communication with strategic public health communication aimed at balancing public psychological well-being. In the wake of the Corona outbreak, Twitter has become one of the most popular social networking sites. We conducted sentiment analysis on tweets using the BERT model, based on people's similar enthusiasm and opinions to better comprehend their mental state.

The researchers in this study [2] do sentiment analysis on two data sets: one data set comprises tweets from individuals all over the globe, while the other data set contains tweets from people in India. The accuracy of the emotion categorization from the GitHub source has been verified by researchers. Researchers used a supervised machine learning technique to analyze Covid-19 tweet sentiment in this study [3]. The ability to identify Covid-19 emotions from tweets would enable more educated judgments to be made about how to best handle the present pandemic scenario. The utilized dataset was retrieved from Twitter using the IEEE data port's IDs. Tweets are retrieved using the Tweepy library by an in-house crawler. Preprocessing methods are used to sanitize the data, and the TextBlob library is used to extract sentiments. The performance evaluation of several machine learning classifiers utilizing their suggested feature set is the contribution of this study. Concatenating the bag-of-words and the term frequency-inverse document frequency yields this collection. Positive, neutral, and negative tweets are categorized. The accuracy, precision, recall, and F-1 score of classifiers are used to evaluate their performance.

In the Nepali digital arena, one may discover Nepali-language material, reviews, and comments. Document-level, Sentence-level, and Aspect-level sentiment analysis are the three methods for assessing sentiments from texts. For Nepali texts, research has been done on identifying emotions [4] and categorizing sentiments [5] at the document level. Sentence-level sentiments offer an emotion for each sentence, whereas Aspect-level sentiments provide sentiments for phrases or entities inside the sentence. Aspects, which are commonly found as nouns, noun groups, or proper nouns, and adjectives, are linked to feelings. Aspect-based Sentiment Analysis [6] aids in the comprehension of the opinions of the connected entities, resulting in improved service or product quality. Using Machine Learning (ML) classifier algorithms such as Support Vector Machine (SVM) and Nave Bayes, a model is constructed to identify aspect-based sentiment in Nepali text (NB). The system gathers Nepali text data and uses Part of Speech (POS) tagging to extract the necessary aspect and emotion characteristics. Each sentence is manually labeled to determine the emotion of the statement. The significance of the terms is calculated using the Term Frequency-Inverse Document Frequency (TF-IDF)

method. The feature vectors are then used to forecast and categorize the text using the Classifier algorithms. The SVM classifier has a 76.8% accuracy, whereas Bernoulli NB has a 77.5 percent accuracy.

Other researches have utilized deep learning as well. This study [14] looks at efficient energy procedures for wireless sensor networks (WSN). The numerous protocol types are provided by the newly suggested taxonomic categorization and comparison. The simulation results of several protocols on the NS-simulator demonstrate that the routing task must be based on a variety of intelligent technologies in order to increase network life and provide greater sensory area coverage. Computing technology advancements have proven to be very effective, resulting in the generation of huge amounts of data that must be evaluated. However, there is much worry about the data's privacy protection. As a result, this work [20] offers a solution to such problems by employing an effective perturbation method that makes use of large data via optimum geometric transformation. Capsule networks that use structured data perform well in visual inference domains. In this work [21], classification of hierarchical multi-label text is achieved utilizing a simple capsule. The study [15] used Twitter to gather tweets on a specific disease and treatment combination, which were then processed to extract attitudes. The tweets obtained on the illness and therapy combination were subjected to the Nave Bayes algorithm. Fake fingerprint identification has recently become a difficult problem in the cyber-crime industry in any industrialized country. The implementation and assessment of appropriate machine learning methods to identify fingerprint liveness are the subject of this paper [16]. A comparison of the Ridge-let Transform (RT) and Machine Learning (ML) approaches is also included. Emotion Analysis for movie reviews detects a commentator's overall assessment or sentiment toward a film. The sentiment analysis model is being pruned by a number of analysts. When compared to classifiers like Gini index with SVM, correlation with random forest, and information with random forest, the proposed study[17] shows that using features optimized by Gini index feature selection and then concatenating with Machine Learning classification gives better results in terms of accuracy and class details parameters. Deep learning algorithms will be used in a variety of applications to aid humans in the near future. Deep learning

algorithms have a proclivity towards allowing a machine to work on its own assumptions. This paper [18] lays out the building pieces needed to implement a deep learning-based method. In addition, the study examines the importance of the preprocessing phase in a number of deep learning-based applications. In this work [19], researchers show the results of a machine learning method utilizing R and Rapid Miner to categorize the sentiment of Twitter tweets. The tweets are retrieved and pre-processed before being categorized into neutral, negative, and positive attitudes, and ultimately the findings are summed up. The Naive Bayes algorithm was used to classify the attitudes expressed in recent tweets from various airlines.

Only machine learning algorithms were utilized in the previous paper. In one article, the entire phrase is used for the training model, but in the other, only one element is employed. The goal of this article is to solve these difficulties by training the model in both machine learning and deep learning based models, and using both aspect and full sentence training. This will allow for a more in-depth examination.

2. Proposed Work

2.3 Architecture

Figure 1. shows the high-level architecture of the study. Initially, the data are scraped from Twitter using a crawler. The gained data are preprocessed by cleaning foreign words, symbols, and translations. Google Translator API is used to translate the sentence from Nepali to English. Afterward, the translated sentence is tokenized using NLTK library, and English stopwords are removed and iNLTK[7] library is used to tokenize Nepali sentences and Nepali stopwords are removed. NLTK[8] is mostly focused on English and other western languages while iNLTK is focused on Indic languages like Hindi, Nepali, Urdu, etc. Next, the tokens are passed through a lemmatizer where the words are converted to meaningful base form taking context into consideration. Since the Nepali pos tagger is not developed enough, the Nepali sentence is translated to English, and aspect is gained from the translated English sentence. The conversion of the sentence and extraction of the aspect is shown below.

Sentence: सम्पन्न भएकाे हाेगुठि ले गर्दा अमेरिका , बेलायत नबन्दा उपत्यका सम्पन्न भएकाे हाे ।

Sentence English: The valley was completed when the United States and the United Kingdom were not formed.

Aspect: valley state kingdom

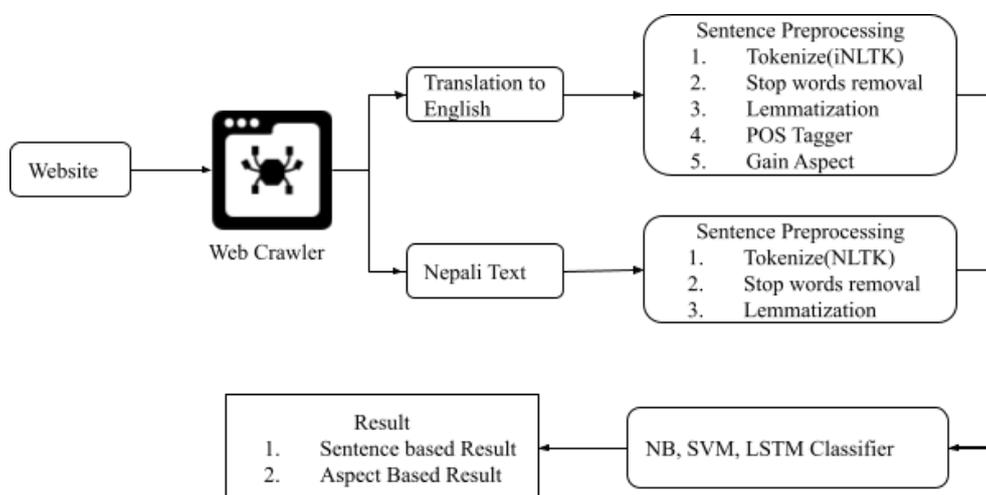


Figure 1. High-Level System Architecture

One of the most successful ways to do text-based sentiment analysis is to use aspect-based analysis. Noun words are used to indicate the aspect. The experiment is conducted using aspect and complete sentences.

CountVectorizer and TfidfTransformer were used to convert training data into vectors for classification using the ML Algorithm. The sentences are one-hot encoded and padding is applied in the case of LSTM[9]. The independent variable is the statement(sentence or aspect), whereas the dependent variable is the label, which represents the mood of the sentence (positive(0) or negative(1)). The system extracts relevant characteristics from the existing data using these vectors, allowing it to learn from them and generate predictions about future texts. Following that, distinct evaluation measure scores from the NB[10], SVM[11], and LSTM

classifiers are recorded using various hyper-parameters. The optimum model is obtained by analyzing the optimal scores from each classifier. Because the negatives outnumber the positives by a little margin, the label is stratified when the data is divided. The results gained after pos tagging a sentence are shown below. Only the aspect part is extracted from the sentence.

Sentence: The valley was completed when the United States and the United Kingdom were not formed.

PosTag: [('The', 'DT'), ('valley', 'NN'),('was', 'VBD'),('completed', 'VBN'),('when', 'WRB'),('the', 'DT'),('United', 'NNP'),('States', 'NNPS'),('and', 'CC'),('the', 'DT'),('United', 'NNP'),('Kingdom', 'NNP'), ('were', 'VBD'), ('not', 'RB'), ('formed', 'VBN')]

2.2 Data Collection

Tweets publish by Nepali users during the COVID19 period are scraped using the Twitter API. A total of 4035 sentences are scraped from Twitter in which 1899 are positive and 2136 are negative. The gained text is labeled manually based on the sentiment of the words. The positive is labeled 0 whereas the negative is labeled 1. A sample of the final preprocessed dataset is shown in Table 1. More information regarding the dataset is provided in Table 2. The complete pre-processed dataset[22] is stored in the kaggle system for further analysis.

Table 1. Sample Dataset

	sentence	sentence_english	aspect	label
31	भू माफियायो सोम भन्ने मान्छे नै भू माफिया जस्त...	The land mafia man Som felt like a land mafia.	land mafia man land mafia	1
32	सडक मा आउनुपर्छधर्म संस्कृत जोगाउन नेवार मात्र...	People from all walks of life, not just Newars...	life street religion sanskrit	0
33	देश ईसाईकरण गर्छधर्म सस्कृति सबै सकेर देश ईसाई...	The country is Christianized. The religion and...	country culture	1

Table 2. Dataset Information

Data	Label '1' (negative)	Label '0' (positive)	Total
Total	2136	1899	4035
Training Data	1691	1537	3228
Testing Data	445	362	807

2.3 Algorithm

The Naive Bayes theorem is used to create a classification method called Naive Bayes. It is widely utilized in text analysis applications such as emotion recognition, spam classification, and so forth. The naïve Bayes Theorem is represented mathematically by Equation 1 and 2.

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)} \quad \text{Eqn. 1}$$

$$P(x|c) = P(x_1|c) * P(x_2|c) * \dots * P(x_n|c) * P(c) \quad \text{Eqn. 2}$$

$P(c|x)$ represents the posterior probability, $P(x)$ represents the predictor prior probability, and $P(c)$ represents the class prior probability and $P(x|c)$ represents the likelihood in the above equation. A posterior probability is a likelihood of allocating observations to classes based on the outcomes. A prior probability is the probability of an observation belonging to a group before the data is gathered.

A Support Vector Machine (SVM) is a discriminative classifier with a separating hyperplane as its formal definition. In other words, the algorithm produces an optimum hyperplane that categorizes fresh instances given labeled training data (supervised learning). The major components of the SVM can be seen in the figure 2. The figure is gained after three-dimensional data is passed through linear kernel SVM. The equation of the supporting vectors,

separating hyperplane, is shown in the figure 2, and also the equation representing the distance between support vectors and the distance between a support plane and a support vector.

The hyperplane equation can be also written as,

$$W^T - b = 0 \quad \text{Eqn. 3}$$

where 'w' is the normal vector of the hyperplane and $\frac{b}{\|w\|}$ denotes the offset of the hyperplane from the origin along the normal vector w. $(x_1, y_1), \dots, (x_n, y_n)$ represents the training dataset of 'n' points.

In this study, ' x_n ' represents the vectorized complete sentences or aspects and ' y_n ' represents the labels that represent the sentiment of the sentences. The linear kernel is used to linearly separate the data to categorize the sentiment. The effectiveness of linear separation is attributed to the existence of increased dimensionality of instances and features.

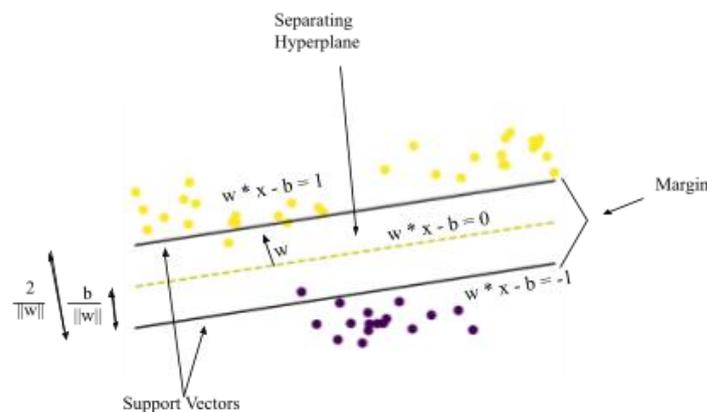


Figure 2. SVM Classifier

Long short-term memory is a type of neural network that works well with time-series data. It's primarily utilized in natural language processing for tasks like text creation and categorization. LSTM uses special units in addition to ordinary units, unlike its predecessor, RNN. An LSTM unit has a memory cell, which can retain data for a long time. Three gates

make up an LSTM. Input Gate, Forget Gate, and Output Gate are the three gates. Input, Forget, and Output Gates are mathematically represented by Equations 4, 5, and 6, respectively.

$$i_t = \sigma(W_i[h_{t-1}, x_t] + b_i) \quad \text{Eqn. 4}$$

$$f_t = \sigma(W_f[h_{t-1}, x_t] + b_f) \quad \text{Eqn. 5}$$

$$o_t = \sigma(W_o[h_{t-1}, x_t] + b_o) \quad \text{Eqn. 6}$$

where " σ " represents a sigmoid function, " w_i ," " w_f " and " w_o " represents the weight for the input, forget and output neuron gate, " w_t " is the input at timestamp 't' and " b_i ," " b_f " and " b_o " represents biases for input, forget and output gate respectively. Fig.3. represents the basic Architecture of a Block of LSTM.

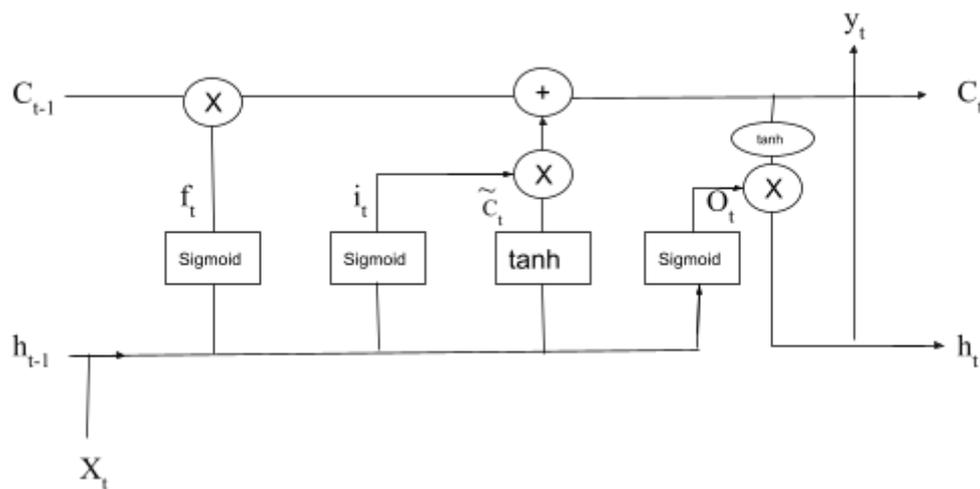


Figure 3. Basic Architecture of a Block of LSTM

In the above figure, " \bar{c}_t ", " c_t ", " h_{t-1} ", and " h_t " are the candidate cell state at timestamp (t), cell state at timestamp(t), the input of the previous hidden state, and input for the next hidden state. The mathematical equations for the candidate cell state, cell state, and the final output is shown below.

$$\bar{c}_t = \tanh(W_c[h_{t-1}, x_t] + b_c) \quad \text{Eqn. 7}$$

$$c_t = f_t * c_{t-1} + i_t * \bar{c}_t \quad \text{Eqn. 8}$$

$$h_t = o_t * \tanh(c^t) \quad \text{Eqn. 9}$$

3. Results and Discussion

The experimental system for this study is shown in Table 3.

Table 3. Model HyperParametera

Experiment Parameter	Value
System	<i>Kaggle (Cloud)</i>
CPU	<i>Intel Xeon</i>
RAM	<i>25.3 GB</i>
Programming Language	<i>Python 3.5</i>
Torch	1.3.1
iNLTK library	<i>Default</i>
NLTK library	3.2.4
Pandas	1.2.2
Numpy	1.19.5
Scikit-learn	0.24.1
Keras	5.3.1
Tensorflow	2.4.1

During the experiments, ML and DL algorithms such as NB, SVM, and LSTM were used. Evaluation measures like as accuracy, precision, recall, and f1-score are generated to assess the model's performance.

$$Accuracy = \frac{TP+TN}{TP+TN+FP+FN} \quad \text{Eqn. 10}$$

$$Precision = \frac{TP}{TP+FP} \quad \text{Eqn. 11}$$

$$Recall = \frac{TP}{TP+FN} \quad \text{Eqn. 12}$$

$$F1 - Score = 2 * \frac{precision*recall}{precision+recall} \quad \text{Eqn. 13}$$

where, TP = True Positive, TN = True Negative, FP = False Positive, FN = False Negative

Table 3, 4 and 6 show the results obtained using Bernoulli Naive Bayes (BernoulliNB), SVM, and LSTM algorithm respectively.

The BernoulliNB [13] classifier is trained in different with complete sentences and aspects of the sentence in varying alpha values chosen between 0 and 1. Alpha is the additive smoothing parameter. Table 4 shows the result obtained.

Table 4. BernoulliNB Classifier

Trained On	Alpha	Accuracy	Precision	Recall	F1-Score
Complete Score	0.1	70	70	70	70
	0.5	71	72	71	71
	0.9	71	71	71	71
Aspect	0.1	76	80	76	76
	0.5	67	66	67	66
	0.9	71	73	71	68

The SVM classifier is trained in different with complete sentences and aspects of the sentence in varying alpha values chosen between 0 and 1. Alpha is the constant which

multiplies the regularization term. Regularization becomes greater when the value is increased. Table 5 shows the result obtained.

Table 5. SVM Classifier

Trained On	Alpha	Accuracy	Precision	Recall	F1-Score
Complete Score	0.1	69	71	69	67
	0.5	63	65	64	63
	0.9	55	66	55	49
Aspect	0.1	64	77	64	53
	0.5	63	67	63	52
	0.9	63	67	63	57

The SVM classifier is trained in different with complete sentences and aspects of the sentence. The model is hyper tuned and Table 6 represents the hyperparameters gained after tuning. Table 7 shows the result obtained.

Table 6. Model Hyperparameter

Hyperparameters	Value
Embedding Vector Feature	40
Epochs	10
Batch Size	32
Loss	Binary CrossEntropy
Optimizer	Adam
Metrics	Accuracy

Table 7. LSTM Classifier

Trained On	Accuracy	Precision	Recall	F1-Score
Complete Sentence	79	80	79	79
Aspect	70	70	70	70

Finally, figure .4. represents the graphical of the performance of the best model in terms of metrics for all models. It can be clearly seen that the LSTM based model trained on full test performs the best.

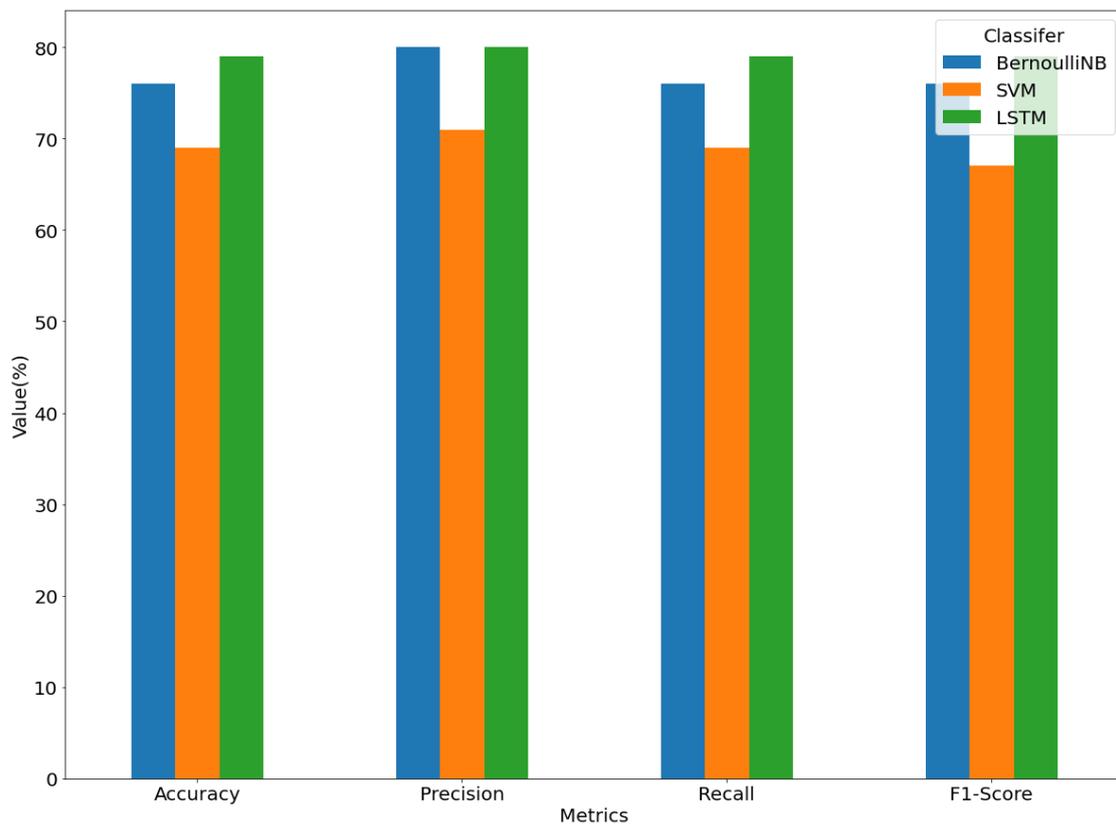


Figure 4. Models Performance

Very little research has been done in Nepali text analysis. So, few papers have been published in journals. Table 8 shows the comparison of the models with past papers used in similar tasks.

Table 8. Model Comparison

Article	Accuracy	Algorithm
Braja Gopal Patra, Dipankar Das, and Amitava Das (2018)	56.9%	SVM
Sujan Tamrakar ¹ , Bal Krishna Bal and Rajendra Bahadur Thapa (2020)	77.5%	Bernoulli Naive Bayes
Milan Tripathi (2021)	79%	LSTM

4. Conclusion

For most organizations, capturing consumer sentiment is critical. This leads to a more detailed understanding of consumer observation, which improves customer happiness. The results of algorithms such as NB, SVM, and LSTM in detecting sentiments in user opinions centered on a whole phrase and a specific entity are shown in this paper. The BernoulliNB Classifier performs best on aspect-based classification with an alpha value of 0.1, whereas the LSTM Classifier performs best on entire sentences on a restricted dataset (4035 phrases and 483 aspects). The algorithm may be evaluated against a larger dataset for future improvements, which might result in increased accuracy.

References

- [1] Lwin, M. O., Lu, J., Sheldenkar, A., Schulz, P. J., Shin, W., Gupta, R., & Yang, Y. (2020). Global sentiments surrounding the COVID-19 pandemic on Twitter: analysis of Twitter trends. *JMIR public health and surveillance*, 6(2), e19447.

- [2] Singh, M., Jakhar, A. K., & Pandey, S. (2021). Sentiment analysis on the impact of coronavirus in social life using the BERT model. *Social Network Analysis and Mining*, 11(1), 1-11.
- [3] Rustam, F., Khalid, M., Aslam, W., Rupapara, V., Mehmood, A., & Choi, G. S. (2021). A performance comparison of supervised machine learning models for Covid-19 tweets sentiment analysis. *Plos one*, 16(2), e0245909.
- [4] Gupta, C. P., & Bal, B. K. (2015). Detecting Sentiment in Nepali texts: A bootstrap approach for Sentiment Analysis of texts in the Nepali language. 2015 International Conference on Cognitive Computing and Information Processing (CCIP), (pp. 1-4).
- [5] Thapa, L. B., & Bal, B. K. (2016). Classifying sentiments in Nepali subjective texts. 2016 7th International conference on information, intelligence, systems & applications (IISA), (pp. 1-6).
- [6] Tamrakar, S., Bal, B. K., & Thapa, R. B. (2020). Aspect Based Sentiment Analysis of Nepali Text Using Support Vector Machine and Naive Bayes. *Technical Journal*, 2(1), 22-29.
- [7] Stigler, S. M. (1983). Who discovered Bayes's theorem?. *The American Statistician*, 37(4a), 290-296.
- [8] Hearst, M. A., Dumais, S. T., Osuna, E., Platt, J., & Scholkopf, B. (1998). Support vector machines. *IEEE Intelligent Systems and their applications*, 13(4), 18-28.
- [9] McCallum, A., & Nigam, K. (1998, July). A comparison of event models for naive bayes text classification. In *AAAI-98 workshop on learning for text categorization (Vol. 752, No. 1, pp. 41-48)*.
- [10] Patra, B. G., Das, D., & Das, A. (2018). Sentiment analysis of code-mixed indian languages: An overview of sail_code-mixed shared task@ icon-2017. *arXiv preprint arXiv:1803.06745*.
- [11] Smys, S., Abul Bashar, and Wang Haoxiang. (2021). Taxonomy Classification and Comparison of Routing Protocol Based on Energy Efficient Rate. *Journal of ISMAC* 3, no. 02 : 96-110.

- [12] Meena, R., V. Thulasi Bai, and J. Omana.(2019).Sentiment Analysis on Tweets for a Disease and Treatment Combination. In International Conference On Computational Vision and Bio Inspired Computing, pp. 1283-1293. Springer, Cham.
- [13] Adam, Edriss Eisa Babikir.(2021). Evaluation of Fingerprint Liveness Detection by Machine Learning Approach-A Systematic View. Journal of ISMAC 3, no. 01: 16-30.
- [14] Kaur, Manpreet.(2019).An Approach for Sentiment Analysis Using Gini Index with Random Forest Classification. In International Conference On Computational Vision and Bio Inspired Computing, pp. 541-554. Springer, Cham.
- [15] Ranganathan, G.(2021).A Study to Find Facts Behind Preprocessing on Deep Learning Algorithms. Journal of Innovative Image Processing (JIIP) 3, no. 01 : 66-74.
- [16] Kumar, G. Ravi, K. Venkata Sheshanna, and G. Anjan Babu.(2020).Sentiment Analysis for Airline Tweets Utilizing Machine Learning Techniques. In International Conference on Mobile Computing and Sustainable Informatics, pp. 791-799. Springer, Cham.
- [17] Haoxiang, Wang, and S. Smys.(2021).Big Data Analysis and Perturbation using Data Mining Algorithm. Journal of Soft Computing Paradigm (JSCP) 3, no. 01: 19-28.
- [18] Manoharan, J. Samuel.(2021).Capsule Network Algorithm for Performance Optimization of Text Classification. Journal of Soft Computing Paradigm (JSCP) 3, no. 01: 1-9.
- [19] Arora, G. (2020). iNLTK: Natural language toolkit for indic languages. arXiv preprint arXiv:2009.12534.
- [20] Bird, S. (2006, July). NLTK: the natural language toolkit. In Proceedings of the COLING/ACL 2006 Interactive Presentation Sessions (pp. 69-72).
- [21] Hochreiter, S., & Schmidhuber, J. (1997). Long short-term memory. Neural computation, 9(8), 1735-1780.
- [22] original | Kaggle. (n.d.). Retrieved July 21, 2021, from <https://www.kaggle.com/milan400/original> (Original | Kaggle, n.d.)

Author's Biography



Milan Tripathi received a Bachelor's degree in computer engineering from Advanced College of Engineering and Management, Tribhuvan University, Kathmandu, Nepal, in 2019. He is currently working as an AI Researcher. He is guiding bachelor's and master's students in their project and thesis papers. To date, he has guided two master's students in their thesis projects. His research interests are computer vision, deep learning, and image processing.