



## SentimentalAnalysisOfCovid-19VaccinationTweets

3274

Dr.SarikaZaware<sup>1</sup>,SanjanaBhosale<sup>2</sup>,RajeshreeKalburgi<sup>3</sup>,ShrutiModale<sup>4</sup>,andPratikKadam<sup>5</sup>

<sup>1</sup>HOD-AISSMSInstituteofInformationTechnology  
<sup>2</sup>Student-AISSMSInstituteofInformationTechnology  
<sup>3</sup>Student-AISSMSInstituteofInformationTechnology  
<sup>4</sup>Student-AISSMSInstituteofInformationTechnology  
<sup>5</sup>Student-AISSMSInstituteofInformationTechnology

### Abstract.

In March 2020, Coronavirus disease was officially announced as a pandemic all over the world by the World Health Organization (WHO). Since then, the whole pharmaceutical world is in a state of war with COVID-19 and has a responsibility to provide its vaccine for the entire world as soon as possible. The coronavirus outbreak has brought unprecedented measures, which forced the authorities to make decisions related to the installation of lockdown in the areas most hit by the pandemic. Social media has been an important support for people while passing through this difficult period. Tweets collected, analyzed, and included in the media reports. Based on the analysis, it can be seen that most tweets are neutral, while the number of compatible tweets exceeds the number of tweets against tweets. [5, 15, 17] In terms of news, it is considered that the occurrence of tweets follows the practice of events. Moreover, the proposed method can be used in a long-term monitoring campaign that can help governments to establish appropriate communication systems and to evaluate them in order to provide clear and adequate information to the general public, which can increase public confidence in the vaccine campaign. The dataset is trained on a machine learning model to classify the opinion of people on the vaccination process. The algorithms used are BERT, SVM and Naive Bayes. [1, 7]

**Keywords:** BERT, Multinomial Naive Bayes, SVM, Covid-19, Vaccination, Sentimental analysis

DOI Number: 10.14704/nq.2022.20.11.NQ66336

NeuroQuantology 2022; 20(11): 3274-3258

## 1 INTRODUCTION

In this proposed system, we are implementing a sentiment analysis model on twitter dataset using Machine Learning algorithms like BERT, SVM and Naive Bayes which will classify the tweets into positive, negative and neutral classes and compare them. It was found that BERT gives the highest accuracy [3, 4]. The proposed approach can be used for a longer monitoring campaign that can help the government to create appropriate means of communication and to evaluate them in order to provide clear and adequate information to the general public,





which could increase the public trust in a vaccination campaign.[8] Performing sentimental analysis on vaccination tweets dataset.This study emphasizes on examining and evaluating the key topics and themes corresponding to COVID-19 vaccinerelated tweets, which were posted by many individual and community profiles, and exploring the opinions laid by the monTwitter.[11,14]

## 2 LITERATURE SURVEY

Opinions Dynamics From Tweets in the Month Following the First Vaccine An-nouncement Opinion mining is a growing area of the Natural Language Pro-cessing field commonly used to determine viewpoints towards targets of interest using computational methods.[1, 6] We look at one such popular microblog called Twitter and build models for classifying “tweets” into positive, negative and neutral sentiment. We build models for two classification tasks: a binary task of classifying sentiment into positive and negative classes and a 3-way task of classifying sentiment into positive, negative and neutral classes.[2, 5, 9] A Survey of Techniques Sentiment Analysis is a term that includes many tasks such as sentiment extraction, sentiment classification, subjectivity classification, summarization of opinions or opinion spam detection, among others. [4, 11] Vast usage of social media by people for expressing their opinions in all aspects of life produces a lot of information on the Web. Analyzers and analysts are constantly dealing with this, how they can transform this massive information available on the web to useful information. [6, 9] The COVID-19 outbreak has brought significant attention to the healthcare sector in recent times, and it has changed the concept of safety in every aspect of our lives.[6] Sentiments Analysis is a task which mainly focuses on textual data and we expect there to be a huge amount of text data.[5]

- Data Gathering
- Pre-Processing
- Feature Extracting
- Feature Selection
- Classifying Methods

Humans are defined as social beings that make us relate to each other by emotion and thoughts. Indeed, the development of the child’s knowledge is through social interaction.[3, 13] Social interaction increased in the 20th century to be a technological platform to interact with humans around the world to get the acceptance of society. Different approaches are also considered while performing and classifying the sentiments such as Deep Learning, Neural Network and also hybrid approach.[19,20]



### 3 SYSTEM ARCHITECTURE

The first step is to collect a COVID-19 vaccination dataset.[1] A monthly combined dataset has been manually annotated as POSITIVE, NEUTRAL or NEG-ATIVE vaccination, in order to be used in the training phase of the classification algorithms.[5]

3251

COVID-19 Vaccine Dataset: In this stage, A machine learning approach is chosen to detect the stance of tweets in the dataset which is obtained by labeling the dataset to train the models.[3]

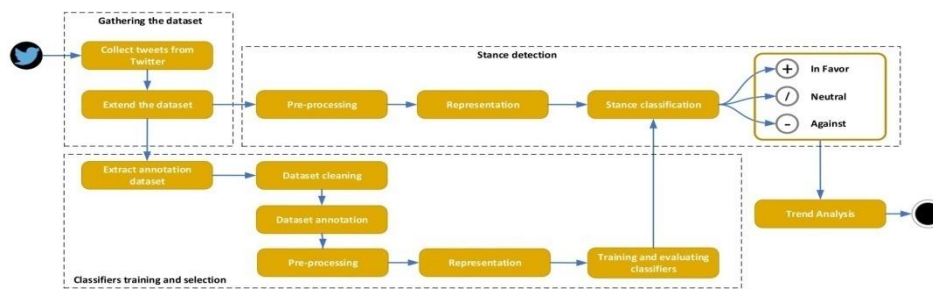


Fig.1. System architecture

#### 3.1 Data Collection

- Several large datasets including collections of tweets related to the COVID has been taken for the analysis.
- In order to collect a dataset of COVID-19 vaccination, a large dataset from kaggle was taken.

#### 3.2 Data Annotation

To maintain the quality of the annotated dataset, which will be used for training the machine learning algorithms, duplicated tweets are removed and the retweets have been easily identified due to the presence of the “RT” symbol and deleted.[2,4,5]

#### 3.3 COVID-19 Vaccine Sentiment Detection

The main components of the stance detection process are the pre-processing, the feature extraction and the machine learning classification.[3]

- Pre-processing: As seen the social media messages are many a times written using casual language, a preprocessing phase is used to filter out the unwanted tweets from the dataset containing non-english language.



- Features: In order to use machine learning algorithms for text classification, the text content has to be first converted into numerical feature vectors.
- Learning Algorithms: Machine learning algorithms are used to predict the sentiments from the tweets. Some popular machine learning algorithms are used like: Multinomial Naive Bayes (MNB), Support Vector Machine (SVM), Bidirectional Encoder Representations from Transformers (BERT).

## 4 LEARNING ALGORITHMS

### 4.1 MNB (Multinomial Naive Bayes)

Multinomial Naive Bayes is most popular supervised learning classification that is used for the analysis of categorical text data. We have selected the first algorithm as MNB because our model is predicting more than one class and Multinomial Naive Bayes is used to solve Multi-Class classification problems. [6] It is a learning method that is mostly used in Natural Language Processing (NLP). The algorithm is based on the Bayes theorem. It calculates the probability of each tag for a given sample and then gives the tag with the highest probability as output. Mathematical Expression Bayes theorem, formulated by Thomas Bayes, calculates the probability of an event occurring based on the prior knowledge of conditions related to an event. [4] It is based on the following formula:

$$P(A|B) = P(A) * P(B|A) / P(B) \quad (1)$$

Where we are calculating the probability of class A when predictor B is already provided.

$P(B)$  = prior probability of B  
 $P(A)$  = prior probability of class A

$P(B|A)$  = Occurrence of predictor B given class A probability.

This formula helps in calculating the probability of the tags in the text. In MNB the dataset is divided into 80% for training purposes and 20% for testing. Multinomial Naive Bayes gives an accuracy of 61.62% which is very less so it is not suitable for regression. Naive Bayes algorithm is only suitable for textual data classification and cannot be used to predict numeric values. As this algorithm has many disadvantages, we shifted toward the SVM algorithm.

### 4.2 SVM (Support Vector Machine)

SVM is a supervised machine learning algorithm used for both classification and regression challenges. SVM performs classification by finding the hyperplane that differentiates the classes we plotted in n-dimensional space. [10] It draws the hyperplane with the help of kernels. They are less prone to overfitting and can perform multiclass classification. Mathematical Explanation The aim of the SVM is to find the hyperplane which maximizes the margin. Here the question



arises how to find the optimal hyperplane.

The equation of a hyperplane is:

$$w^T x = 0 \quad (2)$$

3253

Let's look at the steps to find the equation of the hyperplane  
Step 1: You have a dataset  $D$  and you want to classify it

$$D = \{(x_i, y_i) \mid x_i \in \mathbb{R}^p, y_i \in \{+1, -1\}\}^n \quad i=a \quad (3)$$

Step 2: You need to select two hyperplanes separating the data with no points between them

$$y_i(w \cdot x_i + b) \geq 1 \text{ for all } 1 \leq i \leq n \quad (4)$$

Step 3: Maximize the distance between the two hyperplanes

$$m = 2 \div \|w\| \quad (5)$$

Step 4: Final equation of hyperplane

$$w^T x = 0 \quad (6)$$

In SVM the dataset is divided into 80% for training purposes and 20% for testing. Support Vector Machine gives an accuracy of 83.49% which is little less. It is not suitable for large datasets. It does not execute very well when the data has more noise. Due to these disadvantages, we are shifted towards the BERT algorithm.

### 4.3 BERT

BERT (Bidirectional Encoder Representations from Transformers) is a transformer-based machine learning technique for Natural Language Processing (NLP) pre-training developed by Google. Historically, language models could only read text input sequentially (either left-to-right or right-to-left) but couldn't do both at the same time. [1, 12] BERT is different because it is designed to read in both directions at once. [7, 16] Using this bidirectional capability of BERT, it is pre-trained on two different, but related, Natural Language Processing tasks. These special tokens used by BERT for fine-tuning and specific task training. These are the following:

- CLS: The first token of every sequence. A classification token which is normally used in conjunction with a softmax layer for classification tasks. For anything else, it can be safely ignored.
- SEP: A sequence delimiter token which was used at pre-training for sequence-pair tasks (i.e. Next sentence prediction). Must be used when sequence pair tasks are required.
- MASK: Token used for masked words. Only used for pre-training.



In BERT the dataset is divided into 80% for training purposes and 20%(half for testing and other half for evaluation). After applying the BERT algorithm on the vaccination tweets, we can observe that BERT gives the highest accuracy rate among all the algorithms as it has the ability to process larger amount of text and language. It gives an accuracy of 91.24%. While calculating the accuracy of the model, accuracy for masked words is considered.

## 5 RESULTS

By applying MNB (Multinomial Naive Bayes), SVM (Support Vector Machine) and BERT (Bidirectional Encoder Representations from Transformers) on tweets regarding COVID-19 vaccination we can observe that Multinomial Naive Bayes gives 61.62%, Support Vector Machine 83.49% and BERT gives 91.24%.

The results achieved by the BERT are better than the ones obtained in the case of the other two algorithms in terms of classification report (i.e accuracy, recall, precision and F1-score.)

**Table 1. Result**

Accuracy	Precision	Recall	F1-score
61.62%			
83.49%			
91.24%			



**Fig.2. Initial Data**



```
[12] df_train.head()
```

	date	tweet	replies_count	retweets_count	likes_count	hashtags	clean_tweet
0	07-05-2021	@lm_ankit_14 @DC_Hazaribag The whole vaccinati...	0	0	2	[]	vaccination drive given 3rd Year MBBS students...
1	10-05-2021	@cmohry @mikhattar Sir Mera naam Rahul Gupta...	0	0	0	[]	Sir Mera naam Rahul Gupta hai contact 78279736...
2	07-05-2021	@NewsNationTV Good governance, follow covid gu...	0	0	0	[]	Good governance follow covid guidelines vaccin...
3	07-05-2021	@WeAreVasai @VasaiViranMcorp vaccination cent...	1	0	1	[]	vaccination center social distance super sprea...
4	07-05-2021	@htTweets @narendramodi @drharshvardhan Vaccin...	0	0	1	[]	Vaccines controlled PM fight corona shouldnt f...

Fig.3.Pre-processedData

	polarity	subjectivity	sentiment
0	0.000000	0.500000	1
1	0.000000	0.000000	0
2	0.550000	0.500000	1
3	0.088889	0.277778	1
4	0.133333	0.350000	1
...	...	...	...
6274	0.000000	0.000000	0
6275	0.000000	0.000000	0
6276	0.000000	0.000000	0
6277	0.266667	0.812963	1
6278	0.000000	0.625000	1

6279 rows x 3 columns

Fig.4.SentimentsAdded

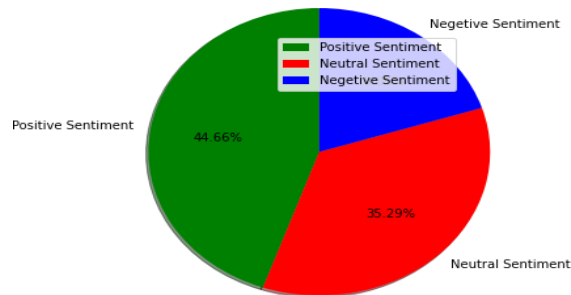


Fig.5.PieChart





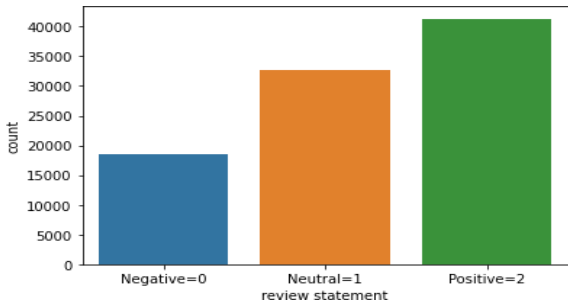


Fig.6.BarGraph

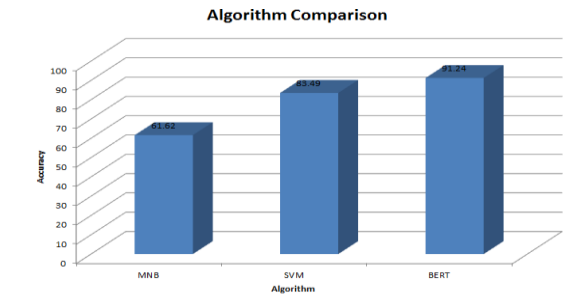


Fig.7.ComparisonGraph

## 6 CONCLUSION

Nowadays social media sites have become a strong tool for influencing people and disseminating information to the general public. However, due to the lack of contextual information in the texts, sentiment analysis for brief texts such as Twitter is particularly difficult. As a result, numerous algorithms are always being created to obtain the most accurate algorithm of all. In the current paper, a five-month period from January to May 2021 tweets has been analyzed using machine learning. The dataset used in this model is from Kaggle. In our model, we have compared three algorithms: MNB, SVM, and BERT. Based on the results, it is concluded that the majority of tweets, i.e. 44.66% were positive, 20.05% negative, and 35.29% neutral tweets. Multinomial Naive Bayes gives an accuracy of 61.62%, SVM gives 83.49%, and Bert gives 91.24%. So it concludes that Bert is the most suitable algorithm for performing sentiment analysis on a large dataset. We have classified the tweets into three main classes: positive, negative, and neutral. The count of sentiment and their percentage is calculated, and the final graph is displayed. This study can be useful for all governments to predict the sentiment of the public in order to make an effective strategy for better managing the situation. Future Twitter data can also be considered for experimentation. The impact of COVID-19 on the financial sector, employability, and personal life of individuals



may be analyzed and prediction will be performed using this machine learning techniques.

## Acknowledgment

3257

It gives us immense pleasure in presenting the project report on “SENTIMENTAL ANALYSIS OF COVID-19 VACCINATION TWEETS USING MACHINE LEARNING”. The success and the outcome of this report required a lot of guidance. We are very grateful to our guide Dr. S. N. Zaware who has provided expertise and encouragement. We thank sir who provided vision and knowledge that was very helpful throughout the research. All that we have done is only due to the great guidance. We also express our gratitude to Dr. S. N. Zaware Head of Computer Engineering Department, AISSMS's Institute of Information Technology, for the valuable support. We would also like to extend our sincere thanks to Principal Dr. P. B. Mane, for his dynamic and valuable guidance throughout the project and providing the necessary facilities that helped us to complete our dissertation work.

## References

1. LIVIU-ADRIAN COTFAS, CAMELIA DELCEA, IOAN ROXIN, CO-RINA IOANAS, DANA SIMONA GHERAI AND FEDERICO TAJARIOL “The Longest Month: Analyzing COVID-19 Vaccination Opinions Dynamics From Tweets in the Month Following the First Vaccine Announcement” February 2021 IEEE Access PP(99):1-1 doi:10.1109/ACCESS.2021.3059821
2. Kumar Rahul, Bhanu Raj Jindal, Kulvinder Singh and Priyanka Meel “Analyzing Public Sentiments Regarding COVID-19 Vaccine on Twitter” 2021 7th International Conference on Advanced Computing and Communication Systems (ICACCS), Coimbatore, India, doi:10.1109/ICACCS51430.2021.9441693
3. Zainab Tariq Soomro, Sardar Haider Waseem Ilyas and Ussama Yaqub “Sentiment, Count and Cases: Analysis of Twitter discussions during COVID-19 Pandemic”, Bournemouth, United Kingdom, doi:10.1109/BESC51023.2020.9348291
4. Maha A. Alanezi and Nabil M. Hewahi “Tweets Sentiment Analysis During COVID-19 Pandemic” 2020 International Conference on Data Analytics for Business and Industry: Way Towards a Sustainable Economy (ICDABI), Sakheer, Bahrain, doi:10.1109/ICDABI51230.2020.9325679
5. Lokesh Mandloi and Ruchi Patel “Twitter Sentiments Analysis Using Machine Learning Methods” 2020 International Conference for Emerging Technology (IN-CET) Belgaum, India, Jun 5-7, 2020 doi:10.1109/INCET49848.2020.9154183
6. Supriya Raheja and Anjani Asthana “Sentimental Analysis of Twitter Comments on Covid-19”, 2021 11th International Conference on Cloud Computing, Data Science and Engineering (Confluence 2021), Noida, India, doi:10.1109/Confluence51648.2021.9377048



7. ApoorvAgarwal, BoyiXie, Ilia Vovsha, Owen Rambowand Rebecca Pas-sonneau“Sentiment Analysis of Twitter Data” Columbia University New York,NY10027USA
8. VishalA.KhardeandS.S.Sonawane“Sentiment Analysis of Twitter Data:ASurveyofTechniques”InternationalJournalofComputerApplications(0975–8887)Volume139–No.11, April2016[9]AmanKhakharia,VruddhiShahandPragyaGupta“SentimentanalysisofCOVID-19vaccineweetsusingMachineLearning” 3258
9. NawSafrin Sattar and Shaikh Arifuzzaman“COVID-19 Vaccination Awareness andAftermath: Public Sentiment Analysis on Twitter Data and Vaccinated PopulationPredictionintheUSA”
10. Nijhum Paul and Swapna S. Gokhale“Analysis and Classification of Vaccine Di-alogueintheCoronavirusEra”2020IEEEInternationalConferenceonBigData(BigData)doi:10.1109/BigData50022.2020.9377888
11. Aman Khakharia, Vruddhi Shah and Pragya Gupta“Sentiment analysis of COVID-19vaccineweetsusingMachineLearning”
12. A. H. Alamoodi et al., “Sentiment analysis and its applications in fighting COVID-19 and infectious diseases: A systematic review,” Expert Systems with Applica-tions,p.114155,Oct.2020,doi:10.1016/j.eswa.2020.114155.
13. N.ÖztürkandS.Ayvaz,“SentimentanalysisisonTwitter:Atextminingapproach to the Syrian refugee crisis,” Telematics and Informatics, vol. 35, no. 1, pp. 136–147, Apr.2018,doi:10.1016/j.tele.2017.10.006.
14. G.A.Ruz,P.A.Henríquez,andA.Mascareño,“SentimentanalysisofTwit- ter data during critical events through Bayesian networks classifiers,” Fu-tureGenerationComputerSystems,vol.106,pp.92–104,May2020,doi:10.1016/j.future.2020.01.005.
15. P.Tiwarietal.,“SentimentAnalysisforAirlinesServicesBasedonTwitterDataset,”inSocialNetwo rkAnalytics,N.Dey,S.Borah,R.Babo, and A. S.Ashour,Eds.AcademicPress,2019,pp.149–162.
16. Himanshu Batra,Narinder Singh Punn,Sanjay Kumar Sonbhadra and Sonali Agar-wal“BERT-BasedSentimentAnalysis:ASoftwareEngineeringPerspective”arXiv:2106.02581v3,2July2021
17. Muhammad Abbas, Kamran Ali,Saleem Memon and Abdul Jamali “MultinomialNaiveBayesClassificationModelforSentimentAnalysis”IJCSNSInternationalJournalOfComputerSci enceandNetworkSecurity,VOL.19No.3, March2019
18. Kai-Xu Han , Wei Chien , Chien-Ching Chiu and Yu-Ting Cheng “Application ofSupport Vector Machine (SVM) in the Sentiment Analysis of Twitter DataSet”Appl.Sci.2020,10,1125;doi:10.3390/app10031125,7February2020
19. JatlaSrikanth,AvulaDamodaram,YuvarajaTeekaraman,RamyakuppusamyandAmruth Ramesh Thelkar “Sentiment Analysis on COVID-19 Twitter Data StreamsUsingDeepBeliefNeuralNetworks”HindawiComputationalIntelligenceandNeu-rosienceVolume2022,ArticleID8898100,6May2022
20. ChetanpalSingh ,Tasadduq Imam , Santoso Wibowo and SrimannarayanaGrandhi “A Deep Learning Approach for Sentiment Analysis of COVID-19 Re-views ” Appl. Sci. 2022, 12, 3709. <https://doi.org/10.3390/app12083709> ,7 April2022

