# Sentiment analysis on social media against public policy using multinomial naive bayes

ZULFIKAR, Wildan Budiawan, ATMADJA, Aldy Rialdy and PRATAMA, Satrya Fajri

## Published version

## Copyright and re-use policy

# Sentiment Analysis on Social Media Against Public Policy Using Multinomial Naive Bayes

**Wildan Budiawan Zulfikar[1*], Aldy Rialdy Atmadja[2], Satrya Fajri Pratama[3]**

[1, 2]Department of Informatics, Faculty of Science and Technology,UIN Sunan Gunung Djati Bandung, Indonesia

[3]Department of Computing, College of Business, Technology and Engineering, Sheffield Hallam University, United Kingdom

**Abstract.**

**Purpose:** The purpose of this study is to analyze text documents from Twitter about public policies in handling COVID-19 that are currently or have been determined. The text documents are classified into positive and negative sentiments by using Multinomial Naive Bayes.

**Methods:** In this research, CRISP-DM is used as a method for conducting sentiment analysis, starting from the business understanding process, data understanding, data preparation, modelling, and evaluation. Multinomial Naive Bayes has been applied in building classification based on text documents. The results of this study made a model that can be used in classifying texts with maximum accuracy.

**Result:** The results of this research are focused on the model or pattern generated by the Multinomial Naive Bayes Algorithm. The classification results of social media users' tweets against the new normal policy obtained good results with an accuracy value of 90.25%. After classifying the tweets of social media users regarding the new normal policy, the results show that more than 70% agreed and supported the new normal policy.

**Novelty:** This study resulted in how classification can be done with Multinomial Naive Bayes and this algorithm can work well in recognizing text sentiments that generate positive or negative opinions regarding public policies handling COVID-19. So, the research provided conclusions about the views of people around the world on new normal public policy.

**Keywords**: Classification, COVID-19, Multinomial Naïve Bayes, Sentiment, Social Media, Public Policy

**Received** November 2022 / **Revised** January 2023 / **Accepted** January 2023

## INTRODUCTION

The new normal is a program policy taken by various countries in the face of the COVID-19 pandemic. This program is implemented as a new order to adapt to COVID-19 or can also be said as a new way of life during the coronavirus pandemic whose recovery rate is increasing. One of the considerations for implementing the new normal is related to the impact of the pandemic on the economy which is considered difficult to worry about. All countries in the world are affected by this pandemic. In living in the new normal era, we learn to carry out activities according to the COVID-19 Health protocol. This situation gave birth to new policies in this era, including the policy of working from home [1]–[5].

Currently, most people in the world are busy talking about things related to this global pandemic on social media. The new normal was one of the things that came up in conversation. People express their thoughts in writing; both positive and negative opinions predominate [6]. It is undeniable that many people are still living a new normal life because the number of positive cases of COVID-19 is still increasing. Based on observations on social media, "new normal", "work from home" and "remote work" are some of the topics that often become trending topics in social media [7]–[9].

A large amount of data is very supportive to search for information through one study, namely sentiment analysis. Sentiment analysis is a computational study that is taken from people's opinions, sentiments and emotions through the entities and attributes they have which are represented in the form of text [10]–[17]. Sentiment analysis will classify the polarity of the text in a sentence or document to find out whether the

opinions expressed in the sentence or document are positive or negative. Therefore, the urgency of this research is how to classify public opinion around the world against the new normal policy into positive and negative opinions. The results of this study can instantly provide conclusions about the views of people around the world on public policy [8].

Multinomial Naive Bayes which is a development of the Naive Bayes algorithm used in this study. Several related studies have used this method on classification models with text documents. As was done in previous research, which is about Emotion Detection among Facebook Diabetes Community with 75.7% for average F-score from 800 documents [18]. In other works, Multinomial Naïve Bayes was implemented to classify text document bases such as email filtering, actor specification of software development process with dial or multi label classification [19]–[21]. These results indicate that the multinomial naive bayes algorithm can work well in classifying data in the form of text [22]–[26].

**METHODS**
**Business Understanding**
In this study, business understanding refers to public opinion regarding the new normal policy. At this stage, it is necessary to understand the background and objectives of business processes related to public policy regarding the new normal. The business purpose of this sentiment analysis is to determine whether the new normal policy should be implemented based on public opinion on Twitter.

This research is related to the atmosphere or trending topic on twitter. This program is a program created to classify documents regarding the new normal. In this program, the processed data is related to someone's twitter data or tweets containing opinions that refer to the new normal policy. The purpose of this study is to provide conclusions about people's responses or reactions to public policies based on social media. This is done by grouping documents based on data interest in sample data or training data on public opinion regarding the new normal policy by crawling using the hashtags "newnormal", the hashtags "remotework", and "workfromhome".

**Data Preparation**
Initial data collection was carried out by crawling public opinion on social media related to research. These documents are used to obtain initial data that supports research related to the new normal policy. Crawling is done by using several related keywords, namely the new normal hashtag and other hashtags related to the policy, namely the hashtags "work from home" and "remote work".

Table 1. Sample of dataset

| No | Document |
|----|----------|
| 1 | 1 = welcome, ay 2020-2021! 📅 |
|   | with its commitment to provide quality education in any ways possible, lpu-batangas college of nursing aims to innovate learning while staying at home as we face the #newnormal. 📲 💻 happy first day of school! 😊 |
|   | #stayconnected https://t.co/ifrjwhxctw |
| 2 | my uncle passed away (non-covid). but so many people are visiting my aunt (in oklahoma), i fear someone will inflict covid on her #newnormal |
|   | straight outta balbriggan |
|   | #newnormal |
|   | #balbriggan https://t.co/pagg4sfaqb |
| 3 | welcome, ay 2020-2021! 📅 |
|   | with its commitment to provide quality education in any ways possible, lpu-batangas college of nursing aims to innovate learning while staying at home as we face the #newnormal. 📲 💻 happy first day of school! 😊 |
|   | #stayconnected https://t.co/ifrjwhxctw |
| 4 | full face mack anti-droplets anti-fog dust-proof face shield protective cover transparent face eyes protector safety accessories https://t.co/1aatwvupos |
|   | #faceshield  #newnormal  #newnormal2020  #faceshield  #facemask #facemasks #rlynedition |
| 5 | #timtalks - how to create inspired, high performance employees with jonas hansen https://t.co/yrogvwxoj3 via @dlaignite @bladt77 #socialselling #digitalselling #employeeexperience #internalcomms #leadership #newnormal https://t.co/xzrpzfrfc6 |
| 6 | "we should not go back to normal, |
|   | because normal was the problem." |
|   | #covid19 #newnormal #pandemicoutbreak PHPHPH |
| 7 | school shopping 2020 style! #newnormal #covidsucks #letsgotoschool https://t.co/89jpxlcmtk |

| No | Document |
|----|----------|
| 8 | wear your mask and wear it properly. read this. #covid19 #newnormal #who #wearamask #coronavirus source: world health organization https://t.co/st4evttwf2 |
| 9 | 😍👒 @krystledsouza instagram story 😁 @smritiiraniofficial .#krystledsouza #coronatime #newnormal #covid19 #love #life https://t.co/ufm0r3ccdz |
| 10 | dear health workers ✅ remember!! even when you're not treating patients, the risk for #covid19 continues. |

This stage includes all activities that build the final dataset (data to be included in the modelling) from the initial raw data. Data preparation includes all activities to build data sets that will be entered into modelling tools from initial raw data or create new data for data mining setups. Data preparation includes all activities to build data sets that will be processed in the modelling process using the Multinomial Naïve Bayes algorithm to classify opinions.

The documents explored to determine the relation between its meaning and the purpose of this research. After that, before entering the classification method, the document preprocessing stage must first be carried out. The purpose of this stage is to eliminate noise, uniform word form and reduce vocabulary volume. As in Figure 1, this preprocessing consists of various stages, namely case folding, tokenizing, stopword and stemming.
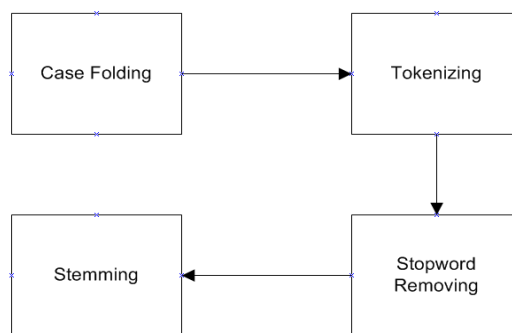


Figure 1. Text preprocessing

The initial stage is to clean the document from numbers, usernames, websites, and hashtags. Case folding is the stage of changing capital letters to lowercase letters. The letters in question are from a-z. Characters other than letters are cleared and considered as delimiters. In the tokenizing process, all words in each document are separated and punctuation is removed, and if there are symbols or anything that is not letters. Tokenizing is the process of trimming the input string based on each word that composes it and distinguishing special characters that are considered as word separators or not. Furthermore, the description of the training data is separated into chapters, paragraphs, sentences, and into words by trimming the strings of the constituent words. There are several process rules so that the results are as desired. On stopword removing, irrelevant words will be removed. In addition, words that do not have their own meaning when separated from other words and are not related to adjectives related to sentiment. Words that are included in the special stopwords list of words contained in the list of English stopwords that often appear but have no meaning are removed at this stage. Stemming is the process of searching for basic words by removing affixes [27]–[30]. In this process the words will be grouped into several groups that have the same root word, such as inaugurate, inaugurate, and inaugurate where the root words of all of them are lantik words[27], [31]–[33]. Stemming is the stage of removing word suffixes [34].

**Modelling**
At the modelling stage, there are several things to do, among others, choosing a modelling technique, building a model, and assessing the model. Modelling is a phase that directly involves Data Mining techniques, namely by selecting Data Mining techniques and determining the algorithms that will be used according to running business processes. In this study using the Multinomial Naïve Bayes algorithm classification method. The process in implementing this system will use the term frequency to find the

pattern of the algorithm by performing the Multinomial Naïve Bayes algorithm calculation flow according to Figure 2 so that the expected results to compare document predictions are more accurate by using a dataset that has been prepared in the form of a classified document.
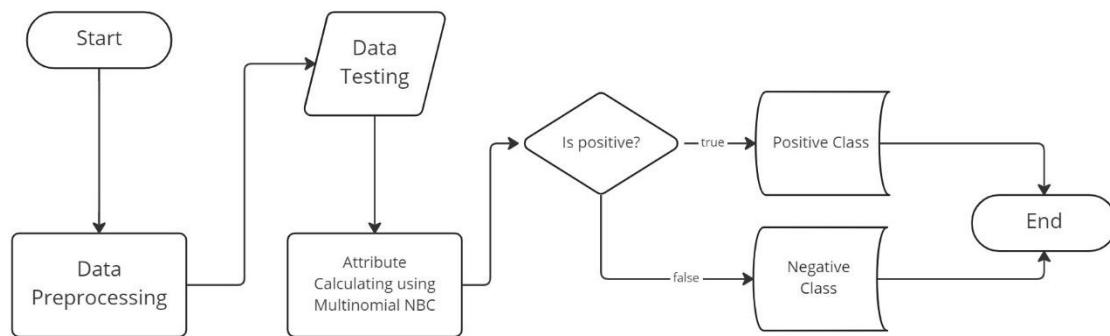


Figure 2. Multinomial naïve bayes

The data mining technique chosen is classification using the Multinomial Naïve Bayes algorithm. The Multinomial Naive Bayes algorithm is very appropriate to be used to achieve the initial goal of this study, namely to classify data based on data interest in sample data or training data on public opinion regarding the new normal policy. Data mining modelling begins with document labelling to be used as data with certain classes.

The Naïve Bayes Classifier algorithm is an algorithm used to find the highest probability value and then classify the test data in the most appropriate category [35]–[37]. In this study, the test data are documents that come from social media. There are two stages in classifying documents. The initial stage is training on documents that have been labelled class or category. While the next stage is the process of classifying documents whose class is not yet known. This system accepts input in the form of raw document data which is still a lot of noise and then preprocessing will be carried out first before working on the next stage.

In Figure 2, the raw document is a document that has not been preprocessed which consists of parsing, case folding, deletion of stopwords and stemming. Example: "What hard just keeping your fucking mask, the new normal get the F**K over—". Furthermore, a collection of documents is processed before being analyzed. This process is called preprocessing text which aims to reduce vocabulary volume, equalize words and eliminate noise. A clean document is a collection of documents that have been cleaned of noise, non-uniform words, for example: "scared the new normal living".

Multinomial Naive Bayes classification is an algorithm that classifies whether the document that is the test data includes a positive or negative sentiment document. The total of each positive and negative document is the final result of the classification. The Multinomial Naive Bayes algorithm is a classification method with probability, which predicts future opportunities based on previous experience so it is known as Bayes theorem.

Data that has gone through the text processing stage can then proceed to the next stage, namely classification with the Multinomial Naïve Bayes Classifier algorithm. Data in the form of text will appear in two text classification results containing positive and negative. The following is the calculation of the Naïve Bayes Classifier algorithm:

1. The initial stage in the Multinomial Naïve Bayes Classifier process is to calculate the probability of each class from the entire training data.

2. Testing process. This process is to determine the accuracy of the model built in the training process, generally using data called a test set to predict labels. The Multinomial Naïve Bayes Classifier method consists of two stages in the text classification process, the training stage and the classification stage. At the training stage, an analysis process is carried out on sample documents in the form of vocabulary selection, namely words that are possible to appear in a collection of sample documents that become

document representations. The next step is to determine the probabilities for each category based on a sample of documents.

The model built in this research has two aspects, namely document as *x* and label or class as *y*. These two aspects are called attributes. Attributes have values called instances. Then to facilitate the research, the classification model that we have trained is stored in memory serialized into a file in the form of a pickle. This pickle is made using Python's cPickle module, so if we want to use a model that has been trained with the same training data as before.

Modelling is done by testing the test data with training data that has been previously processed using the multinomial Naive Bayes algorithm based on predetermined rules. Labelling is done automatically by the program based on the calculation of the multinomial Naive Bayes algorithm with reference to the model that has been created. Table 2 shows an example of manual calculation of algorithms and classification.

Table 2. Sample of data training

| Set | Document ID | Document | Class |
|---|---|---|---|
| *Training* set | D1 | *scared new normal living* | Negative |
| | D2 | *already new normal peace* | Positive |
| | D3 | *new normal wild shit* | Negative |
| | D4 | *hate new normal living* | Negative |
| *Test* set | D5 | *hate new normal bull shit* | ? |

Based on the training data, the important appearance features include scared, peace, shit, and hate. After obtaining the required features, the classification process is carried out using the Naive Bayes algorithm. Below are the calculation results for the test data on Equation (1):

$$\gamma(\beta) = \frac{\gamma(\alpha) * \gamma(\beta|\alpha)}{\gamma(\beta)} \tag{1}$$

The process of calculating prior probabilities of positive and negative classes is as follows:

y(pos) = ¼
y(neg) = ¾

Then calculate the maximum likelihood value using Equation (2):

$$\gamma(\alpha) = \frac{count(\omega,\alpha)+1}{count(\alpha)+|V|} \tag{2}$$

Table 3. Multinomial naïve bayes calculation (2)

| Category | Document | | | |
|---|---|---|---|---|
| | Scared | Shit | Hate | Peace |
| Positive | $\frac{(0+1)}{(4+4)}$ | $\frac{(0+1)}{(4+4)}$ | $\frac{(0+1)}{(4+4)}$ | $\frac{(1+1)}{(4+4)}$ |
| $\gamma(\alpha)$ | $\frac{1}{8}$ | $\frac{1}{8}$ | $\frac{1}{8}$ | $\frac{1}{4}$ |
| Negative | $\frac{(1+1)}{(12+4)}$ | $\frac{(1+1)}{(12+4)}$ | $\frac{(0+1)}{(4+4)}$ | $\frac{(0+1)}{(12+4)}$ |
| $\gamma(\alpha)$ | 1/8 | 1/8 | 1/8 | 1/16 |

After getting the value of the training data, the occurrence of words in the test data will be seen in the probability model to look for possibilities in each category as in Table 4. The results of the calculation show that the negative class in D5 has the highest value, so the D5 class has a negative class.

Table 4. Classification Result

| Document | Class | |
|---|---|---|
| | Positive | Negative |
| D5 | $\frac{3}{4} \times \frac{1}{8} \times \frac{1}{8}$ | $\frac{1}{4} \times \frac{1}{8} \times \frac{1}{8}$ |
| Result | 0.00390625 | 0.01171857 |

**RESULT AND DISCUSSION**

The evaluation in this study is more focused on the model or pattern generated by the Multinomial Naive Bayes algorithm. The resulting model is analyzed to determine whether the resulting pattern is in accordance with the classification standards on the training data. If the resulting pattern is not appropriate, then further analysis of the resulting pattern needs to be carried out so that it can produce recommendations for implementing public policies, which are expected to help the government's success in implementing further policies. This study conducted 5 test scenarios based on the percentage of training data and testing data. In addition, a confusion matrix which uses 4 indicators as a reference namely precision, recall, F1-score, and accuracy is applied to each test scenario. The results of the testing process are shown in detail in Table 5.

Table 5. Evaluation results

| Training Data | | Testing Data | | Precision (%) | | Recall (%) | | F1-Score (%) | | Accuracy |
|---|---|---|---|---|---|---|---|---|---|---|
| % | Data | % | Data | 0 | 1 | 0 | 1 | 0 | 1 | % |
| 90 | 910 | 10 | 102 | 71 | 82 | 73 | 80 | 72 | 81 | 77 |
| 80 | 809 | 20 | 203 | 77 | 86 | 84 | 79 | 80 | 82 | 81 |
| 70 | 708 | 30 | 304 | 77 | 70 | 66 | 79 | 71 | 74 | 73 |
| 60 | 607 | 40 | 405 | 73 | 73 | 67 | 78 | 70 | 75 | 73 |
| 50 | 506 | 50 | 506 | 67 | 75 | 73 | 70 | 70 | 73 | 71 |
| Mean of Accuracy | | | | 73 | 77,2 | 89,9 | 77,2 | 72,6 | 77 | 75 |

There are several things that can be described from Table 5 including: there is an interesting fact that when viewed from row 2 to row 5 it appears that the more training data, the higher the accuracy. The higher the percentage of training data, the higher it is, but in line 1 where the highest amount or percentage of training data is 90%, it actually produces a lower accuracy value than line 2 which produces an accuracy value of 81%. Based on Table 5, row 2 with the composition of training data and testing data of 80% and 20% is used to conduct sentiment analysis on 964 shrunken documents on social media related to one of the policies, namely the new normal. The results of the classification process are represented in the form of a confusion matrix as described in Figure 3.
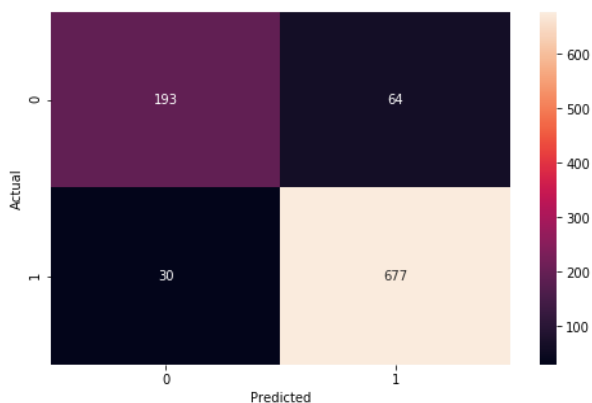


Figure 3. Visualization of confusion matrix result

Figure 3 explains that there are 4 conditions regarding the classification results which are detailed as follows:

1. True Positive, namely tweets from social media users with positive sentiments and were declared positive by the classification results of 677 documents.
2. True Negative, namely tweets from social media users with negative sentiments and declared negative by the classification results of 193 documents.
3. False positive, namely tweets from social media users with negative sentiments but were declared positive by the classification results of 64 documents.
4. False negative, namely tweets from social media users who have positive sentiments but are declared negative in the classification results of 30 documents.

Based on the true positive, true negative, false positive, and false negative values, it can be calculated to calculate the performance aspects of the classification model, namely accuracy, precision, and recall as follows:

1. *Accuracy* is the ratio of correct predictions (positive and negative) to the overall data of 90.25% as detailed in the formula below*:*

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

$$= \frac{(677+193)}{(677+64+30+193)}$$

$$= 0.9025$$

2. *Precision* is the ability of the classification model to identify only the relevant data points. Precision is the ratio of correct positive predictions to the overall positive predicted outcome. The precision obtained is 91.36% with the following calculation.

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

$$= \frac{677}{(677+64)}$$

$$= 0.9136$$

3. *Recall* is the ability of the classification model to identify all relevant data points in the dataset. Recall or sensitivity is the ratio of true positive predictions compared to the overall true positive data. The recall obtained is 95.76% with the following calculation.

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

$$= \frac{677}{(677+30)}$$

$$= 0.9676$$

Recall is the ability of the classification model to identify all relevant data points in the dataset. Recall or sensitivity is the ratio of true positive predictions compared to the overall true positive data. The recall obtained is 95.76% with the following calculation.

According to data collected from social media, the data collection includes 964 records. Figure 4 shows that 73,3% of respondents concur that the government's actions to address COVID-19 have a good impact on public perception. During this time, 26.7% disapproved of or had a poor impression of the measures adopted in response to the COVID-19 epidemic.
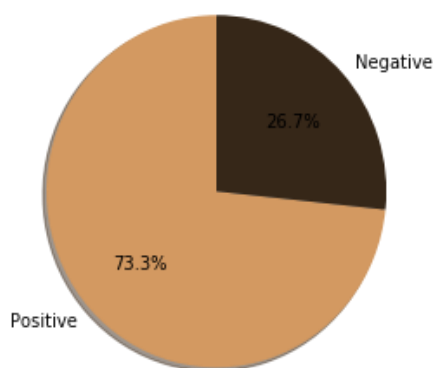


Figure 4. Visualization of the sentiment of social media users to the new normal policy

Based on the experimental results, the Multinomial Naïve Bayes algorithm can achieve a good result for the classification of social media users' tweets against the new normal policy. This classification started with four text preprocessing stages such as case folding, tokenizing, stop word removing and stemming [34]. The results showed a 90.25% accuracy, a 91.36% precision, and a 96.76% recall value. When compared to similar studies, the findings of this investigation provided quite positive outcomes [22]. After classifying the tweets of social media users regarding the new normal policy, the results show that 73.3% of 964 people tend to agree and support the new normal policy. However, there are an obstacles where recognizing sentences to be processed in the sentiment classification process are when they encounter the initial sentence which gives the perception that the sentence is a positive sentence then at the end of the sentence gives the perception that the sentence is a negative sentence. So, the high probability of term that exist in the training data has an impact on the classification outcomes. This can be overcome with finding semantic relationship between the words [38].

## CONCLUSION

The classification results of social media users' tweets against the new normal policy have reached an accuracy value of 90.25%. Those outcomes were tested using 964 documents. Additionally, 73.3% of people support the government's recommended policies, demonstrating how effective the Multinomial Naïve Bayes algorithm to recognize sentiment classification towards the Covid-19 epidemic. There are challenges where both positive and negative terms are present in one sentence as in the following document the document "This policy is good to implement but there are still many concerns because the number of patients with COVID-19 is still increasing". In manual data classification, the sentence is included in the category of negative sentiment, but when the classification uses the system, the sentence is included in the category of positive sentiment because at the beginning of the sentence there are words that contain positive sentiment. In this case, the results of the classification are impacted by the high likelihood of terms that are available on social media text. To solve this problem, the research can be improved by finding semantic relationship between the words. The findings of this study can also serve as one of the arguments for why the new normal policies that various nations will be implementing or already have implemented are successful. For further research, it is recommended to add classes such as very negative, negative, neutral, positive and very positive so that more categories can be used as assessment parameters to achieve the goal of a better business understanding.

## REFERENCES

[1]     J. A. Pacheco, "The 'New Normal' in Education," *Prospect. 2020 511*, vol. 51, no. 1, pp. 3–14, 2020, doi: 10.1007/S11125-020-09521-X.

[2]     S. Zhang, M. J. Ventura, and H. Yang, "Network Modeling and Analysis of COVID-19 Testing Strategies," in *Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc., EMBS*, 2021, pp. 2003–2006, doi: 10.1109/EMBC46164.2021.9629754.

[3]     I. Aygun, B. Kaya, and M. Kaya, "Aspect Based Twitter Sentiment Analysis on Vaccination and Vaccine Types in COVID-19 Pandemic With Deep Learning," *IEEE J. Biomed. Heal. Informatics*, vol. 26, no. 5, pp. 2360–2369, 2022, doi: 10.1109/JBHI.2021.3133103.

[4]     Y. Pathak, P. K. Shukla, and K. V. Arya, "Deep Bidirectional Classification Model for COVID-19 Disease Infected Patients," *IEEE/ACM Trans. Comput. Biol. Bioinforma.*, vol. 18, no. 4, pp. 1234–1241, 2021, doi: 10.1109/TCBB.2020.3009859.

[5]     S. Zhang, S. Yang, and H. Yang, "Statistical Analysis of Spatial Network Characteristics in Relation to COVID-19 Transmission Risks in US Counties," in *Proc. Annu. Int. Conf. IEEE Eng. Med. Biol. Soc., EMBS*, 2021, pp. 2278–2281, doi: 10.1109/EMBC46164.2021.9629892.

[6]     R. Jayapermana, A. Aradea, and N. I. Kurniati, "Implementation of Stacking Ensemble Classifier for Multi-class Classification of COVID-19 Vaccines Topics on Twitter," *Sci. J. Inform.*, vol. 9, no. 1, pp. 8–15, 2022, doi: 10.15294/sji.v9i1.31648.

[7]     D. J. Pradipta, I. N. P. Pradnyana, and T. Raharjo, "The New Normal Strategy for Project Management in Directorate General of Customs and Excise," in *2020 3rd Int. Conf. Comput. Inform. Eng. IC2IE 2020*, 2020, pp. 249–254, doi: 10.1109/IC2IE50715.2020.9274568.

[8]     R. P. Aluna, I. N. Yulita, and R. Sudrajat, "Electronic News Sentiment Analysis Application to New Normal Policy during the Covid-19 Pandemic Using Fasttext and Machine Learning," in *2021 Int. Conf. Artif. Intell. Big Data Anal. ICAIBDA 2021*, 2021, pp. 236–241, doi: 10.1109/ICAIBDA53487.2021.9689756.

[9]     D. A. Azarov, D. M. Nazarov, and Y. P. Silin, "Fuzzy Assessment of the Russian Federation

Military-Industrial Complex Economic Influence in the Context of a 'New Normal,'" in *Proc. 2017 20th IEEE Int. Conf. Soft Comput. Meas. SCM 2017*, 2017, pp. 856–858, doi: 10.1109/SCM.2017.7970745.

[10] S. Yuliyanti, T. Djatna, and H. Sukoco, "Sentiment Mining of Community Development Program Evaluation Based on Social Media," *TELKOMNIKA (Telecommun. Comput. Electron. Control.*, vol. 15, no. 4, pp. 1858–1864, 2017, doi: 10.12928/TELKOMNIKA.V15I4.4633.

[11] T. Shaik, X. Tao, C. Dann, H. Xie, Y. Li, and L. Galligan, "Sentiment Analysis and Opinion Mining on Educational Data: A Survey," *Nat. Lang. Process. J.*, vol. 2, p. 100003, 2023, doi: 10.1016/J.NLP.2022.100003.

[12] L. Shang, H. Xi, J. Hua, H. Tang, and J. Zhou, "A Lexicon Enhanced Collaborative Network for targeted financial sentiment analysis," *Inf. Process. Manag.*, vol. 60, no. 2, p. 103187, 2023, doi: 10.1016/J.IPM.2022.103187.

[13] H. Li, B. X. B. Yu, G. Li, and H. Gao, "Restaurant Survival Prediction Using Customer-Generated Content: An Aspect-Based Sentiment Analysis of Online Reviews," *Tour. Manag.*, vol. 96, p. 104707, 2023, doi: 10.1016/J.TOURMAN.2022.104707.

[14] N. Leelawat *et al.*, "Twitter Data Sentiment Analysis of Tourism in Thailand During the COVID-19 Pandemic Using Machine Learning," *Heliyon*, vol. 8, no. 10, p. e10894, 2022, doi: 10.1016/J.HELIYON.2022.E10894.

[15] H. T. Ismet, T. Mustaqim, and D. Purwitasari, "Aspect Based Sentiment Analysis of Product Review Using Memory Network," *Sci. J. Inform.*, vol. 9, no. 1, pp. 73–83, 2022, doi: 10.15294/SJI.V9I1.34094.

[16] S. Fransiska, R. Rianto, and A. I. Gufroni, "Sentiment Analysis Provider By.U on Google Play Store Reviews with TF-IDF and Support Vector Machine (SVM) Method," *Sci. J. Inform.*, vol. 7, no. 2, pp. 203–212, Nov. 2020, doi: 10.15294/SJI.V7I2.25596.

[17] Jumanto, M. A. Muslim, Y. Dasril, and T. Mustaqim, "Accuracy of Malaysia Public Response to Economic Factors During the Covid-19 Pandemic Using Vader and Random," *J. Inf. Syst. Explor. Res.*, vol. 01, no. 01, pp. 49–70, 2023.

[18] V. Balakrishnan and W. Kaur, "String-based Multinomial Naïve Bayes for Emotion Detection among Facebook Diabetes Community," *Procedia Comput. Sci.*, vol. 159, pp. 30–37, 2019, doi: 10.1016/J.PROCS.2019.09.157.

[19] V. K. Vineetha and P. Samuel, "A Multinomial Naïve Bayes Classifier for Identifying Actors and Use Cases from Software Requirement Specification documents," in *2022 2nd Int. Conf. Intell. Technol. CONIT 2022*, 2022, pp. 1–5, doi: 10.1109/CONIT55038.2022.9848290.

[20] S. Kadam, A. Gala, P. Gehlot, A. Kurup, and K. Ghag, "Word Embedding Based Multinomial Naive Bayes Algorithm for Spam Filtering," in *Proc. - 2018 4th Int. Conf. Comput. Commun. Control Autom. ICCUBEA 2018*, 2018, pp. 1–5, doi: 10.1109/ICCUBEA.2018.8697601.

[21] R. A. Pane, M. S. Mubarok, N. S. Huda, and Adiwijaya, "A Multi-Lable Classification on Topics of Quranic Verses in English Translation Using Multinomial Naive Bayes," in *2018 6th Int. Conf. Inf. Commun. Technol. ICoICT 2018*, 2018, pp. 481–484, doi: 10.1109/ICOICT.2018.8528777.

[22] A. R. Susanti, T. Djatna, and W. A. Kusuma, "Twitter's Sentiment Analysis on GSM Services using Multinomial Naïve Bayes," *TELKOMNIKA (Telecommun. Comput. Electron. Control.*, vol. 15, no. 3, pp. 1354–1361, 2017, doi: 10.12928/TELKOMNIKA.V15I3.4284.

[23] L. Jiang, S. Wang, C. Li, and L. Zhang, "Structure Extended Multinomial Naive Bayes," *Inf. Sci. (Ny).*, vol. 329, pp. 346–356, 2016, doi: 10.1016/J.INS.2015.09.037.

[24] N. Shiri Harzevili and S. H. Alizadeh, "Mixture of Latent Multinomial Naive Bayes Classifier," *Appl. Soft Comput.*, vol. 69, pp. 516–527, 2018, doi: 10.1016/J.ASOC.2018.04.020.

[25] P. Bermejo, J. A. Gámez, and J. M. Puerta, "Improving the Performance of Naive Bayes Multinomial in E-mail Foldering by Introducing Distribution-based Balance of Datasets," *Expert Syst. Appl.*, vol. 38, no. 3, pp. 2072–2080, 2011, doi: 10.1016/J.ESWA.2010.07.146.

[26] N. Chirawichitchai, "Sentiment Classification by a Hybrid Method of Greedy Search and Multinomial Naïve Bayes Algorithm," in *Int. Conf. ICT Knowl. Eng.*, 2013, pp. 1–4, doi: 10.1109/ICTKE.2013.6756285.

[27] A. Kraal, P. W. van den Broek, A. W. Koornneef, L. Y. Ganushchak, and N. Saab, "Differences in Text Processing by Low- and High-Comprehending Beginning Readers of Expository and Narrative Texts: Evidence from Eye Movements," *Learn. Individ. Differ.*, vol. 74, p. 101752, 2019, doi: 10.1016/J.LINDIF.2019.101752.

[28] P. K. Jayasekara and K. S. Abu, "Text Mining of Highly Cited Publications in Data Mining," in *IEEE 5th Int. Symp. Emerg. Trends Technol. Libr. Inf. Serv. ETTLIS 2018*, 2018, pp. 128–130, doi:

10.1109/ETTLIS.2018.8485261.

[29]    S. Jain, S. C. Jain, and S. Vishwakarma, "Enhanced Text Classification Methods to Improve the Performance of the Various Text Mining Processes using Rapid Miner," in *Proc. 2021 IEEE Int. Conf. Mach. Learn. Appl. Netw. Technol. ICMLANT 2021*, 2021, pp. 1–5, doi: 10.1109/ICMLANT53170.2021.9690551.

[30]    T. Matsumoto, W. Sunayama, Y. Hatanaka, and K. Ogohara, "Data Analysis Support by Combining Data Mining and Text Mining," in *Proc. - 2017 6th IIAI Int. Congr. Adv. Appl. Inform. IIAI-AAI 2017*, 2017, pp. 313–318, doi: 10.1109/IIAI-AAI.2017.165.

[31]    S. Bhattacharjee, D. Delen, M. Ghasemaghaei, A. Kumar, and E. W. T. Ngai, "Business and Government Applications of Text Mining & Natural Language Processing (NLP) for Societal Benefit: Introduction to the Special Issue on 'Text Mining & NLP,'" *Decis. Support Syst.*, vol. 162, p. 113867, 2022, doi: 10.1016/J.DSS.2022.113867.

[32]    A. Motz, E. Ranta, A. S. Calderon, Q. Adam, F. Alzhouri, and D. Ebrahimi, "Live Sentiment Analysis Using Multiple Machine Learning and Text Processing Algorithms," *Procedia Comput. Sci.*, vol. 203, pp. 165–172, 2022, doi: 10.1016/J.PROCS.2022.07.023.

[33]    D. Zhang, J. Hyönä, L. Cui, Z. Zhu, and S. Li, "Effects of Task Instructions and Topic Signaling on Text Processing Among Adult Readers with Different Reading Styles: An Eye-tracking Study," *Learn. Instr.*, vol. 64, p. 101246, 2019, doi: 10.1016/J.LEARNINSTRUC.2019.101246.

[34]    N. P. Ririanti and A. Purwinarko, "Implementation of Support Vector Machine Algorithm with Correlation-Based Feature Selection and Term Frequency Inverse Document Frequency for Sentiment Analysis Review Hotel," *Sci. J. Inform.*, vol. 8, no. 2, pp. 297–303, 2021, doi: 10.15294/sji.v8i2.29992.

[35]    N. Umar and M. A. Nur, "Application of Naïve Bayes Algorithm Variations On Indonesian General Analysis Dataset for Sentiment Analysis," *J. RESTI (Rekayasa Sist. dan Teknol. Informasi)*, vol. 6, no. 4, pp. 585–590, 2022, doi: 10.29207/RESTI.V6I4.4179.

[36]    A. Falasari and M. A. Muslim, "Optimize Naïve Bayes Classifier Using Chi Square and Term Frequency Inverse Document Frequency For Amazon Review Sentiment Analysis," *J. Soft Comput. Explor.*, vol. 3, no. 1, pp. 31–36, 2022, doi: 10.52465/joscex.v3i1.68.

[37]    T. L. Nikmah, M. Z. Ammar, Y. R. Allatif, R. M. P. Husna, P. A. Kurniasari, and A. S. Bahri, "Comparison of LSTM , SVM , and Naive Bayes for Classifying Sexual Harassment Tweets," *J. Soft Comput. Explor.*, vol. 3, no. 2, pp. 131–137, 2022, doi: https://doi.org/10.52465/joscex.v3i2.85.

[38]    P. Wang, B. Xu, J. Xu, G. Tian, C. L. Liu, and H. Hao, "Semantic Expansion Using Word Embedding Clustering and Convolutional Neural Network for Improving Short Text Classification," *Neurocomputing*, vol. 174, pp. 806–814, 2016, doi: 10.1016/J.NEUCOM.2015.09.096.