



Comparison of Naive Bayes Multinomial Algorithm and Decision Tree on Tweet Emotion Classification

Arif Bijaksana Putra Negara^{1,*}, Rina Septiriana¹, Windari Oktapia Simanjuntak¹

¹ Universitas Tanjungpura, Pontianak, Indonesia

ARTICLE INFO

Article history:

Received 15 October 2025

Received in revised form 23 November 2025

Accepted 30 November 2025

Available online 9 December 2025

Keywords:

Multinomial Naive Bayes; decision tree;
bag of word; bigram; SMOTE

ABSTRACT

Social media has become a place where everyone can express their feelings and thoughts without limitations. One of the most widely used social media platforms is Twitter. Twitter has 238 million active users and provides access to search for information through specific tweets. Therefore, Twitter can be used as a source of information to analyze a person's emotions based on their writings/tweets. In analyzing the feelings of a tweet, a method is needed to classify tweets into appropriate emotion classes. The classification of tweet emotions aims to group tweets into predetermined emotion classes such as anger, joy, fear, love, and sadness. The algorithms used to build machine learning models for emotion classification are Multinomial Naive Bayes and Decision Trees. This research aims to determine which classification algorithm, Multinomial Naive Bayes or Decision Tree, is better by comparing the accuracy values of these classification algorithms. This study applies feature extraction using Bag of Words, Bigram and the SMOTE method. The research results show that the Multinomial Naive Bayes classification model, which involves feature extraction using Bag of Words and the SMOTE method, has the highest accuracy value of 67.15%.

1. Introduction

Emotion is one of the essential aspects of human life. Emotions occur due to reactions to external and internal stimuli, involving physiological changes and various thoughts [1]. Often, humans find it difficult to express their emotions verbally but tend to tell whatever they desire through social media [2]. One of the popular social media platforms for obtaining trending information and news is Twitter [3].

Twitter is a social media platform with a broad user base that allows users to post messages called Tweets, which are limited to 280 characters and can include images, videos, or blog links [4]. Twitter provides an Application Programming Interface (API) to enable users to access and search for specific Tweet information, such as specific topics, locations, current trends, and more. Thus, Twitter can be

* Corresponding author.

E-mail address: arifbpn@untan.ac.id

a data source for analysing writings/tweets. Studying emotions in tweets requires a classification method to categorise tweets into the correct emotion classes.

The classification of emotions in tweets aims to group tweets into predetermined emotion classes, such as anger, joy, fear, love, and sadness [5]. Several algorithms commonly used to build a classification model include K-Nearest Neighbors, Logistic Regression, Naive Bayes Classifier, Decision Tree, and Support Vector Machine. Each algorithm has its way of building models and its complexity [6]. In one study, the Multinomial Naive Bayes algorithm achieved the highest accuracy value of 84.60% [7]. However, comparing the Decision Tree classification algorithm obtained the highest accuracy of 99.95% compared to Naive Bayes, SVC, K-NN, and Random Forest algorithms [8].

Previous research has compared the performance of several feature extractions to improve the performance of the naive Bayes algorithm by applying Unigram, Bigram, and Trigram feature extractions. The results showed that the Naive Bayes algorithm with Bigram feature extraction achieved higher results [9]. On the other hand, another study comparing the performance of Ensemble Features and Bag of Words showed that Bag of Words had the highest accuracy [10].

One of the challenges in building a machine learning model is using an imbalanced dataset with varying data quantities among classes. SMOTE (Synthetic Minority Over-Sampling Technique) is used to address the issue of imbalanced data. This method creates replicas of minority data by finding k-nearest neighbors (nearest neighboring data points in the minority data) [11]. The Naive Bayes algorithm was optimized using the SMOTE approach in a study by Saputra *et al.*, [12], providing evidence supporting this assertion. The research obtained an accuracy value of 86.33%, better than the Naive Bayes algorithm without using SMOTE, which had an accuracy of 78.88%.

Based on the research description, this study evaluates emotion classification algorithms using Multinomial Naïve Bayes and Decision Tree algorithms. Both algorithms are compared to determine which one performs better. The parameter used for evaluation is accuracy. The study compares classification algorithms and feature extraction methods such as Bag of Words, Bigram, and SMOTE.

2. Research Method

The research method used consists of datasets, text preprocessing, feature extraction, SMOTE, data splitting, modeling, and evaluation. Figure 1 illustrates the flow of the research methodology.

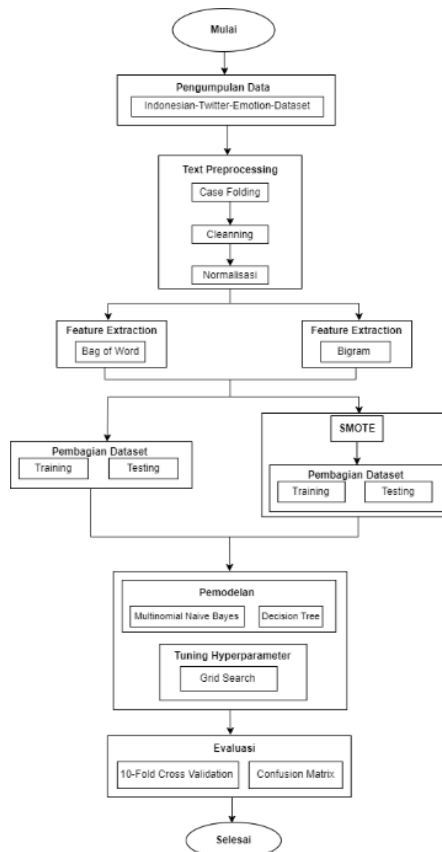


Fig. 1. Research methodology flow

2.1 Dataset

The dataset in the study is from previous research [13]. The entire dataset consists of 4,401 Tweets that have been categorized into five emotion classes: anger, happiness, sadness, love, and fear. Table 1 presents a sample of Tweet data from the five emotion classes.

Table 1

Data sample tweet

Class	Tweet
Anger	Tipikal cowo menye menye bgt ni. Paling ngga bisa diajak argumen. Menyalahkan orang dan mengeluh terus ni manusia maha sempurna. Pen ngatain goblok tapi puasa.
Fear	Jinhyuk menunduk, menatap sepatunya. Merasa takut untuk membalas permintaan sang puan. "M-maaf, saya tidak membawa alat lukis saat ini..." [URL]
Happy	[USERNAME] Hahaha samaa, aku bgitu ngliat oppa2 kore lgsg pgen bljr korea, nntn drama barat lgsg pgen bsa lncr english,ngliat bny ama mewart lgsg pgen bljr bhsa thai wkwkwlw
Love	ENGGA. KARNA SAYANG ITU TIDA MEMANDANG UMUR. KARNA AKU BAKALAN BUCININ PASANGAN AKU, JADI MAU TUA ATAU MUDA NYA DIA BAKALAN AKU IKUTIN TERUSSD
Sadness	Wktu masih jaman sekolah kalau ada laki2 yg suka boro2 mau nanggepin ..yg ada malu nya minta ampun smpe tuh laki malah d musuhin!! Tp bocah jman skrng jomblo bentar galau nya minta ampun ..gmn klau mereka kaya gua yg jomblo 7 tahun #OldMoneyGakNgerasain

2.2 Text Processing

Text preprocessing is a stage to clean the text obtained from social media, which contains unstructured sentences, non-standard words, and noise [14]. There are three stages of text preprocessing carried out in this research, which are case folding, cleaning, and normalisation [15]:

- Case Folding: converting all letter characters in the text to lowercase.
- Cleaning: removing symbols, hashtags, mentions, and unnecessary numbers from the text.
- Normalisation: identifying and replacing non-standard words with words that conform to the standard Indonesian dictionary (KBBI).

Table 2

Text processing

Text	Text Preprocessing
Kepingin gudeg mbarek Bu hj. Amad Foto dari google, sengaja, biar teman-teman jg membayangkannya. Berbagi itu indah.	Text
kepingin gudeg mbarek bu hj. amad foto dari google, sengaja, biar teman-teman jg membayangkannya. berbagi itu indah.	Case Folding
kepingin gudeg mbarek bu hj amad foto dari google sengaja biar teman teman jg membayangkannya berbagi itu indah	Cleaning
ingin gudeg mbarek bu hj amad foto dari google sengaja biar teman teman juga membayangkannya berbagi itu indah	Normalization
ingin gudeg mbarek bu hj amad foto dari google sengaja biar teman teman juga membayangkannya berbagi itu indah	Preprocessing Result

2.3 Frequency Distribution

Frequency Distribution is a graphical or tabular representation that shows the number of occurrences of the most frequently used words in a document [16]. In frequency distribution, the data is broken down into individual words and then grouped based on their frequency of occurrence. Table 3 presents the frequency distribution results with a total vocabulary of 17.083 in the document and 125.408 tokens. Figure 2 shows a graph of the frequency of occurrence of the top 40 most frequently appearing word sequences found in the dataset.

Table 3

Frequency distribution

No	Vocabulary	Token
1	yang	3359
2	aku	2599
3	tidak	2252
4	dan	1895
5	di	1757
6	sudah	1168

No	Vocabulary	Token
...
17083	sembuhkanlah	1
Vocabulary dataset		17083
Token		125408

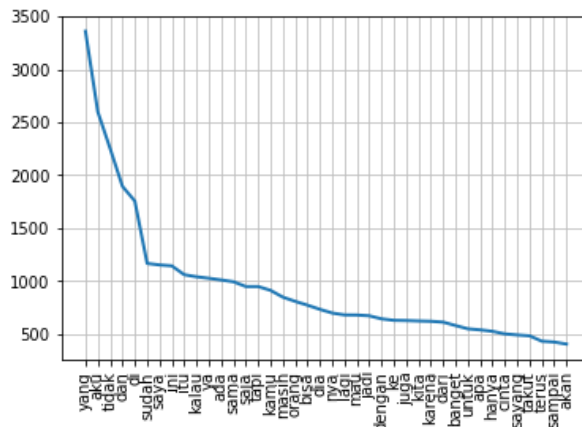


Fig. 2. Graph frequency distribution

2.4 Feature Extraction

Feature extraction aims to convert text into vector representations [17]. In this stage, the Tweet data will be tokenised into individual words, then transformed into numeric features, and the frequency of each word will be counted [18]. This research employs two feature extraction methods, Bag of Words and Bigram, which are implemented using the Sklearn CountVectorizer library.

a. Bag of Word

Bag of Word represents text as a vector of the number of words frequently appearing in documents [19]. This method works by cutting documents into individual words that are collected as a vocabulary and then counting the number of occurrences of each term [20]. The vocabulary from the Bag of Word results is presented in Figure 3, and the frequency results calculation is shown in Figure 4.

Vocabulary content: {'soal': 14721, 'jalan': 6185, 'jatibaru': 6261, 'polisi': 12652, 'tidak': 15931, 'bisa': 1935, 'gertak': 4987, 'gubernur': 5188, 'memang': 9055, 'nya': 11153, 'ikut': 5828, 'pembahasan': 11872, 'jangan': 6218, 'berpoliti': 1681, 'pengaturan': 12043, 'wilayah': 16812,

(0, 14721) 2
(0, 6185) 1
(0, 6261) 1
(0, 12652) 2
(0, 15931) 2
(0, 1935) 1
(0, 4987) 1
(0, 5188) 2
(0, 9055) 1
(0, 11153) 1
(0, 5828) 1
(0, 11872) 1
(0, 6218) 1
(0, 1681) 1
(0, 12043) 1
(0, 16812) 1

Fig. 3. Vocabulary bag of word

Fig. 4. Bag of word

b. Bigram

A bigram is part of the N-gram method, where this process cuts multiple words or two words collected in the form of a vocabulary and counts the number of occurrences [21]. The vocabulary from the Bigram results is presented in Figure 5, and the frequency results calculation is shown in Figure 6.

Vocabulary content: {'soal jalan': 71029, 'jalan ja
tibusu': 30841, 'jatibusu polisi': 31351, 'polisi t
idak': 59798, 'tidak bisa': 77242, 'bisa gertak': 1
1592, 'gertak gubernur': 24440, 'gubernur memang':
24947, 'memang nya': 46813, 'nya polisi': 53906, 't
idak ikut': 77361, 'ikut pembahasan': 27841, 'pempa
hasan jangan': 57452, 'jangan berpolitik': 31069,

Fig. 5. Vocabulary Bigram

(0, 71029)	1
(0, 30841)	1
(0, 31351)	1
(0, 59798)	2
(0, 77242)	1
(0, 11592)	1
(0, 24440)	1
(0, 24947)	1
(0, 46813)	1
(0, 53906)	1
(0, 77361)	1
(0, 27841)	1
(0, 57452)	1
(0, 31069)	1

Fig. 6. Bigram

2.5 SMOTE

SMOTE is a technique used to address class imbalance problems in a dataset. SMOTE works by finding k-nearest neighbours in the minority class data, then creating synthetic data equal to a percentage of the minority data duplication (percentage oversampling), and the k-nearest neighbours are randomly selected [11]. One of the advantages of this method is that it does not cause any information loss since there is no reduction in data, as done in undersampling ways [22]. The total dataset before SMOTE consists of 4401 data points, and Figure 7 displays the distribution of the dataset before SMOTE.

After the SMOTE process, the total dataset increased to 5505 data points, and the distribution of the dataset can be seen in Figure 8.



Fig. 7. Data before SMOTE



Fig. 8. Data after SMOTE

2.6. Data Splitting

Data splitting is the process of dividing data into training data and testing data [23]. The training data is used to build a machine-learning model by learning the patterns in the data. The testing data is used to evaluate the machine learning model. This study divides the dataset into a ratio of 90% for training data and 10% for testing data.

2.7 Multinomial Naive Bayes

Multinomial Naive Bayes is a supervised learning algorithm for classifying or predicting class probabilities or data labels. This algorithm is often used because it is speedy in learning training and predictions [24]. The basic concept of Multinomial Naive Bayes is to study parameters by looking at each feature individually and calculating each feature's probability for each class [25]. Multinomial Naive Bayes has an alpha parameter to control the complexity and smoothing of the model. The higher the alpha value, the more smoothing is applied, and the simpler the model becomes [6].

2.8 Decision Tree

Decision Tree is a popular algorithm in classification cases because it simplifies the representation of patterns [26]. This algorithm uses the tree concept, where nodes represent each attribute, leaves represent each class, and branches hold the values of each category [27]. Decision Trees are susceptible to overfitting, so research uses parameters to control the complexity of the Decision Tree model by setting `max_depth` (the depth of the decision tree) and adjusting the value of `min_samples_leaf` (the minimum number of samples required to be at a leaf node) [28].

2.9 Scenario

There are four scenarios in this research using the Multinomial Naive Bayes and Decision Tree algorithms. The scenarios involve applying two feature extraction techniques, Bag of Words and Bigram, and the SMOTE method on the data. Table 4 displays the testing design for this study.

Table 4
The scenarios

Testing Design	<i>Bag of Word</i>	<i>Bigram</i>	SMOTE
Scenario 1	✓	✗	✗
Scenario 2	✗	✓	✗
Scenario 3	✓	✗	✓
Scenario 4	✗	✓	✓

- Scenario one, the Multinomial Naive Bayes classification algorithm and Decision Tree are trained using a dataset that has implemented the Bag of Word feature extraction technique with a total of 4401 datasets.
- Scenario two, the Multinomial Naive Bayes classification algorithm and Decision Tree are trained using datasets that have implemented the Bigram feature extraction technique with a total of 4401 datasets used.

- Scenario three, the Multinomial Naive Bayes classification algorithm and Decision Tree are trained using datasets that have been processed using the Bag of Word feature extraction technique and the SMOTE method. A total of 5505 datasets are used.

2.10 Model

Both algorithms have undergone tuning of their hyperparameters. Hyperparameter Tuning aims to find the optimal values for the hyperparameters to improve the model's performance. One commonly used method is Grid Search. Grid Search is an optimization technique that explores combinations of hyperparameter values by testing each combination of model parameters [29]. For the Multinomial Naive Bayes algorithm, the hyperparameter used is alpha. Meanwhile, the Decision Tree algorithm uses two hyperparameters, namely max_depth and min_samples_leaf. Table 5 presents the results of the Grid Search for Multinomial Naive Bayes, and Table 6 shows the Grid Search for Decision Tree results.

Table 5
Grid search for Multinomial Naive Bayes

Model	PARAMETER	Range Value	GRID SEARCH Result
Model 1	Alpha	0.01, 0.1, 0.5, 0.9, 1.5, and 5.0	Alpha = 0.9
Model 2	Alpha	0.01, 0.1, 0.5, 0.9, 1.5, and 5.0	Alpha = 5.0
Model 3	Alpha	0.01, 0.1, 0.5, 0.9, 1.5, and 5.0	Alpha = 1.5
Model 4	Alpha	0.01, 0.1, 0.5, 0.9, 1.5, 5.0, and 5.5	Alpha = 5.0

Table 6
Grid search for decision tree

Model	PARAMETER	Range Value	GRID SEARCH Result
Model 1	Max_depth	100 - 120	Max_depth = 100
	Min_samples_leaf	11 - 20	Min_samples_leaf = 15
Model 2	Max_depth	160 - 170	Max_depth = 160
	Min_samples_leaf	11 - 20	Min_samples_leaf = 11
Model 3	Max_depth	100 - 120	Max_depth = 100
	Min_samples_leaf	2 sampai 10	Min_samples_leaf = 9
Model 4	Max_depth	160 sampai 170	Max_depth = 160
	Min_samples_leaf	11 sampai 20	Min_samples_leaf = 11

3. Results And Analysis

The evaluation is to determine how well the model performs in classifying test data and identify whether the model is experiencing overfitting (where the accuracy of the train data is better than the test data) [30].

3.1 Model Performance Multinomial Naive Bayes

The accuracy results obtained in the testing scenarios for the Multinomial Naive Bayes algorithm, as shown in Table 7, indicate that Model 1, which uses the Bag of Words feature extraction, has a 0.92% increase in accuracy on the training data and a 2.07% increase on the testing data compared to Model 3, which incorporates both the Bag of Words feature extraction and the SMOTE method. The study demonstrates that applying the SMOTE method on the Multinomial Naive Bayes algorithm using the Bag of Words feature extraction significantly improves accuracy from the testing data to the training data.

In Model 2, which utilizes the Bigram feature extraction, there is a decline of -10.09% in accuracy on the training data and -6.97% in the testing data compared to Model 4, which combines the Bigram feature extraction with the SMOTE method. Model 1 and Model 3 show that they do not suffer from overfitting (where the training data accuracy is better than the testing data)—however, Model 2 experiences overfitting with a difference of -1.10% between the training and testing. The research findings indicate that applying the Bigram feature extraction reduces the algorithm's performance.

Model 4 experienced an increase in the accuracy of the test data by 2.02% compared to the accuracy of the train data. Model 3 shows that the model performs best by obtaining test data accuracy of 67.15% and data training of 65.04%.

Table 7
Result of accuracy Multinomial Naive Bayes

MODEL	Data		Difference
	Train	Test	
1 (Bow)	64.12%	65.08%	0.96%
2 (Bigram)	52.35%	51.25%	-1.10%
3 (Bow-SMOTE)	65.04%	67.15%	2.11%
4 (Bigram-SMOTE)	42.26%	44.28%	2.02%

3.2 Model Performance Decision Tree

Table 8 shows the results of the accuracy of the Decision Tree algorithm resulting in Model 1 using the Decision Tree algorithm with the application of the Bag of Word feature extraction compared to Model 3, which is the Decision Tree algorithm with the application of the Bag of Word feature extraction and the SMOTE method experienced a decrease in accuracy of -0.59% in train data and an increase in accuracy in test data of 3.11%. Model 2 is a Decision Tree algorithm model that applies a Bigram feature extraction compared to model 4, a Decision Tree algorithm that applies a Bigram feature extraction. The SMOTE method decreases accuracy on the train data by -4.44% and -2.3% on test data. Model 1 shows that the model is not overfitting (where the accuracy of the train data is better than the test data). Model 3 shows that the model performs best by obtaining an accuracy test data result of 53.90%.

Table 8
Result of accuracy decision tree

MODEL	Data		Difference
	Train	Test	
1 (Bow)	52.73%	50.79%	-1.94%
2 (Bigram)	37.73%	38.78%	1.05%
3 (Bow-SMOTE)	52.14%	53.90%	1.76%
4 (Bigram-SMOTE)	33.29%	36.48%	3.19%

3.3 Comparison of Multinomial Naive Bayes and Decision Tree Model

The comparison of the highest accuracy results from both models is presented in Table 9. Based on the accuracy results, the model that outperforms the others is the Multinomial Naive Bayes model with the application of the Bag of Words feature extraction and the SMOTE method, achieving an accuracy of 67.15%. In contrast, the Decision Tree model applying the Bag of Words feature extraction and the SMOTE method provides the highest accuracy of 53.90%.

Table 9
Comparison of accuracy Models

Model	Accuracy
<i>Multinomial Naive Bayes (Bag of Word dan SMOTE)</i>	67.15%
<i>Decision Tree (Bag of Word dan SMOTE)</i>	53.90%

4. Conclusion

This research compares the accuracy of two classification algorithms, Multinomial Naive Bayes and Decision Tree, by applying two feature extraction techniques, Bag of Words and Bigram, and using the SMOTE method to address imbalanced class data.

The highest accuracy obtained from the Multinomial Naive Bayes model is 67.15%. This value is achieved when applying the Bag of Words feature extraction and the SMOTE method to handle imbalanced data. On the other hand, the Decision Tree model has the highest accuracy value of 53.90%.

The research findings indicate that applying the Bag of Words feature extraction yields higher accuracy than the Bigram feature extraction for both algorithms. Additionally, the SMOTE method effectively improves the performance of both algorithms when applied with the Bag of Words feature extraction, addressing the issue of imbalanced class data. However, the SMOTE method decreases the algorithms' performance when using the Bigram feature extraction.

Thus, the Multinomial Naive Bayes model with the Bag of Words feature extraction and the SMOTE method can be more effective in emotion classification.

Acknowledgment

This research was not funded by any grant.

References

- [1] Prawitasari, Johana E. "Mengenal emosi melalui komunikasi nonverbal." *Buletin Psikologi* 3, no. 1 (1995): 27-43.

- [2] Hamzah, Radja Erland, and Citra Eka Putri. "Analisis Self-Disclosure Pada Fenomena Hyperhonest Di Media Sosial." *Jurnal Pustaka Komunikasi* 3, no. 2 (2020): 221-229. <https://doi.org/10.2514/6.2008-383>
- [3] Nuryawan, Azril Tazidan Octa, Mamun Hasbullah, Miftahul Rizal, Muhamad Fauzan Rajab, and Nova Agustina. "ALGORITMA DECISION TREE UNTUK ANALISIS SENTIMEN PUBLIC TERHADAP MARKETPLACE DI INDONESIA." *Naratif: Jurnal Nasional Riset, Aplikasi dan Teknik Informatika* 5, no. 1 (2023): 18-25. <https://doi.org/10.53580/naratif.v5i1.186>
- [4] Imam, Imam Santoso, and Imam Santoso. "ANALISIS SENTIMEN PADA TWITTER TERHADAP GAGALNYA PELAKSANAAN PIALA DUNIA DI INDONESIA MENGGUAKAN METODE NAÏVE BAYES." *Jurnal IKRAITH-INFORMATIKA Vol 7*, no. 2 (2023): 145.
- [5] Negara, Arif Bijaksana Putra, Hafiz Muhandi, and Fahmi Sajid. "Perbandingan algoritma klasifikasi terhadap emosi tweet berbahasa indonesia." *Jurnal Edukasi dan Penelitian Informatika* 7 (2021): 242-249. <https://doi.org/10.26418/jp.v7i2.48198>
- [6] Müller, Andreas C., and Sarah Guido. *Introduction to machine learning with Python: a guide for data scientists*. "O'Reilly Media, Inc.", 2016.
- [7] Yang, Ang, Jun Zhang, Lei Pan, and Yang Xiang. "Enhanced twitter sentiment analysis by using feature selection and combination." In *2015 International Symposium on Security and Privacy in Social Networks and Big Data (SocialSec)*, pp. 52-57. IEEE, 2015. <https://doi.org/10.1109/SocialSec2015.9>
- [8] Maulidah, Mawadatul, Windu Gata, Rizki Aulianita, and Cucu Ika Agustyaningrum. "Algoritma Klasifikasi Decision Tree Untuk Rekomendasi Buku Berdasarkan Kategori Buku." *E-Bisnis: Jurnal Ilmiah Ekonomi dan Bisnis* 13, no. 2 (2020): 89-96. <https://doi.org/10.51903/e-bisnis.v13i2.251>
- [9] Trianto, Rahmawan Bagus, Andri Triyono, and Dhika Malita Puspita Arum. "Klasifikasi Rating Otomatis pada Dokumen Teks Ulasan Produk Elektronik Menggunakan Metode N-gram dan Naïve Bayes." *Jurnal Informatika Universitas Pamulang* 5, no. 3 (2020): 295-301. <https://doi.org/10.32493/informatika.v5i3.6110>
- [10] Permatasari, Rosy Indah, Mochammad Ali Fauzi, Putra Pandu Adikara, and Eka Dewi Lukmana Sari. "Analisis Sentimen Film pada Twitter Berbahasa Indonesia Menggunakan Ensemble Features dan Naïve Bayes." *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer* 2, no. 11 (2018): 5921-5927. <http://j-ptiik.ub.ac.id>
- [11] Sofyan, Sabiq, and Achmad Prasetyo. "Penerapan Synthetic Minority Oversampling Technique (SMOTE) Terhadap Data Tidak Seimbang Pada Tingkat Pendapatan Pekerja Informal Di Provinsi DI Yogyakarta Tahun 2019." In *Seminar Nasional Official Statistics*, vol. 2021, no. 1, pp. 868-877. 2021. <https://doi.org/10.34123/semnasoffstat.v2021i1.1081>
- [12] Saputra, Dedi Dwi, Windu Gata, Nia Kusuma Wardhani, Ketut Sakho Parthama, Hendra Setiawan, Sularso Budilaksono, Dimas Yogatama et al. "Optimization Sentiments of Analysis from Tweets in myXLCare using Naïve Bayes Algorithm and Synthetic Minority over Sampling Technique Method." In *Journal of Physics: Conference Series*, vol. 1471, no. 1, p. 012014. IOP Publishing, 2020. <https://doi.org/10.1088/1742-6596/1471/1/012014>
- [13] Saputri, Mei Silviana, Rahmad Mahendra, and Mirna Adriani. "Emotion classification on indonesian twitter dataset." In *2018 International Conference on Asian Language Processing (IALP)*, pp. 90-95. IEEE, 2018. <https://doi.org/10.1109/IALP.2018.8629262>
- [14] Haddi, Emma, Xiaohui Liu, and Yong Shi. "The role of text pre-processing in sentiment analysis." *Procedia computer science* 17 (2013): 26-32. <https://doi.org/10.1016/j.procs.2013.05.005>
- [15] Hidayatullah, Ahmad Fathan, C. I. Ratnasari, and S. Wisnugroho. "The influence of stemming on Indonesian tweet sentiment analysis." In *Proceeding of International Conference on Electrical Engineering, Computer Science and Informatics (EECSI 2015)*, pp. 127-132. 2015. <https://doi.org/10.11591/eecsi.v2i1.79120>
- [16] Manikandan, S. "Frequency distribution." *Journal of pharmacology & pharmacotherapeutics* 2, no. 1 (2011): 54. <https://doi.org/10.4103/0976-500X.77120>
- [17] Salim, Emil, and Mohammad Syafrullah. "Analisis Sentimen Pada Ulasan Pelayanan Suku Dinas Kependudukan Dan Pencatatan Sipil Kota Administrasi Jakarta Barat Menggunakan Algoritme K-Nearest Neighbor." *Bit (Fakultas Teknologi Informasi Universitas Budi Luhur)* 20, no. 1 (2023): 58-65. https://kemsalim.space/ulasan_dukcapil/
- [18] Zulkifli, Zulkifli, Agung Toto Wibowo, and Gia Septiana. "Pembobotan Fitur Ekstraksi Pada Peringkasan Teks Bahasa Indonesia Menggunakan Algoritma Genetika." *eProceedings of Engineering* 2, no. 2 (2015). <https://libraryproceeding.telkomuniversity.ac.id/index.php/engineering/article/view/2782>
- [19] Permana, Adhitya Prayoga, Totok Chamidy, and Cahyo Crysdian. "Klasifikasi ulasan fasilitas publik menggunakan metode Naïve Bayes dengan seleksi fitur Chi-square." *JISKA (Jurnal Informatika Sunan Kalijaga)* 8, no. 2 (2023): 112-124. <https://doi.org/10.14421/jiska.2023.8.2.112-124>
- [20] Deepu, S., Raj Pethuru, and S. Rajaraajeswari. "A Framework for Text Analytics using the Bag of Words (BoW) Model for Prediction." *International Journal of Advanced Networking & Applications (IJANA)* 2, no. 1 (2016): 320-323. <https://archive.ics.uci.edu/ml/datasets/Bag+of+Words>

- [21] Fanesya, Fera, Randy Cahya Wihandika, and Indriati Indriati. "Deteksi Emosi pada Twitter Menggunakan Metode Naive Bayes dan Kombinasi Fitur." *Jurnal Pengembangan Teknologi Informasi dan Ilmu Komputer* 3, no. 7 (2019): 6678-6686.
- [22] Wijayanti, N. P. Y. T., Eka N. Kencana, and I. Wayan Sumarjaya. "SMOTE: potensi dan kekurangannya pada survei." *E-Jurnal Matematika* 10, no. 4 (2021): 235. <https://doi.org/10.24843/mtk.2021.v10.i04.p348>
- [23] Pambudi, Agung, Z. Abidin, and P. Permata. "Penerapan Crisp-Dm Menggunakan Mlr K-Fold Pada Data Saham Pt. Telkom Indonesia (Persero) Tbk (Tlkm)(Studi Kasus: Bursa Efek Indonesia Tahun 2015-2022)." *Jurnal Data Mining dan Sistem Informasi* 4, no. 1 (2023): 1. <https://doi.org/10.33365/jdmsi.v4i1.2462>
- [24] Herwanto, Herwanto, Nuke L. Chusna, and Muhammad Syamsul Arif. "Klasifikasi SMS Spam Berbahasa Indonesia Menggunakan Algoritma Multinomial Naïve Bayes." *JURNAL MEDIA INFORMATIKA BUDIDARMA* 5, no. 4 (2021): 1316. <https://doi.org/10.30865/mib.v5i4.3119>
- [25] Berliana, Garnis, Shaufiah Shaufiah, and Siti Saâ. "Klasifikasi Posting Tweet Mengenai Kebijakan Pemerintah Menggunakan Naïve Bayesian Classification." *eProceedings of Engineering* 5, no. 1 (2018).
- [26] Handayani, Putri Kurnia. "Penerapan Principal Component Analysis untuk Peningkatan Kinerja Algoritma Decision Tree pada Iris Dataset." *Indonesian Journal of Technology, Informatics and Science (IJTIS)* 1, no. 2 (2020): 55-58. <https://doi.org/10.24176/ijtis.v1i2.4939>
- [27] Robianto, Robianto, Sampe Hotlan Sitorus, and Uray Ristian. "Penerapan Metode Decision Tree Untuk Mengklasifikasikan Mutu Buah Jeruk Berdasarkan Fitur Warna Dan Ukuran." *Coding: Jurnal Komputer dan Aplikasi* 9, no. 01 (2021): 76-86. <https://dx.doi.org/10.26418/coding.v9i01.45907>.
- [28] Sembodo, J. Eka, E. Budi Setiawan, and Z. Abdurahman Baizal. "Data Crawling Otomatis pada Twitter." In *Indonesian Symposium on Computing (Indo-SC)*, pp. 11-16. 2016. <https://doi.org/10.21108/indosc.2016.111>
- [29] Rifai, Nur Azizah Komara. "Klasifikasi Penyakit Diabetes Retinopati Menggunakan Support Vector Machine dengan Algoritma Grid Search Cross-validation." *Jurnal Riset Statistika* (2023): 79-86.
- [30] Lorento, Chavin, Arif Bijaksana Putra Negara, and Rudy Dwi Nyoto. "Implementasi Sentimen Masyarakat berdasarkan Tweet terkait Kebijakan Kemendikbud di Masa Pandemi Covid-19." *JUSTIN (Jurnal Sistem dan Teknologi Informasi)* 10, no. 3: 294-302. <https://doi.org/10.26418/justin.v10i3.54243>