

Sentiment Analysis of X Platform on Viral 'Fufufafa' Account Issue in Indonesia Using SVM

Suryanto¹, Widyastuti Andriyani*²

^{1,2}Magister of Information Technology Universitas Teknologi Digital Indonesia Yogyakarta, Indonesia

e-mail: ¹suryantokantor01@gmail.com, ²widya@utdi.ac.id, rikie@utdi.ac.id

Abstract

In this study, we conducted a comprehensive sentiment analysis of users on the social media platform X concerning the viral controversy surrounding the KasKus account known as "Fufufafa." This issue attracted widespread attention and sparked varied reactions within the online community. To gain insights into public opinion on the topic, we utilized the Support Vector Machine (SVM) method, a widely recognized machine learning algorithm for classification tasks. The data for this research was gathered from various posts, comments, and public discussions on platform X, which were pre-processed to filter out irrelevant information, such as spam, unrelated topics, and non-informative content. After cleaning the data, user sentiments were categorized into three primary classes: positive, negative, and neutral. The SVM model was then trained and tested using a labeled dataset to accurately predict user sentiments based on the textual content of their interactions. Through this approach, we aimed to capture the overall mood and attitudes of the online community towards the "Fufufafa" issue. The findings of the study reveal that the majority of users on platform X expressed negative sentiments about the viral controversy, suggesting dissatisfaction or disapproval. Meanwhile, a smaller portion of the users remained neutral, while an even smaller segment displayed positive sentiment. This disparity highlights the polarized reactions within the online discourse. Our study also demonstrates the efficacy of the SVM method in analyzing large-scale social media data to understand public sentiment on viral issues. Ultimately, this research offers a valuable contribution to understanding how users on social media respond to controversial topics and trending events.

Keywords— Sentiment Analysis, X, SVM, Fufufafa, Indonesia

1. INTRODUCTION

The rise of social media has dramatically transformed the way individuals communicate, share, and disseminate information on a global scale. Platforms like Facebook, Twitter, and Instagram, and in this case, platform X, have enabled real-time interactions, allowing information to spread rapidly and influence public discourse almost instantaneously. Viral phenomena, particularly those that arise from contentious or controversial issues, can quickly spark widespread discussions and play a crucial role in shaping public opinion. One such viral issue in Indonesia is the controversy surrounding the KasKus account "Fufufafa," which drew significant attention and elicited a variety of responses from the online community.

Understanding how users respond to viral topics on social media has become increasingly important, especially for businesses, policymakers, media organizations, and researchers. Social media platforms have evolved into powerful tools for gauging public sentiment, and as viral issues

emerge, it is vital to monitor the reactions of users in order to make informed decisions. Sentiment analysis, a method that leverages natural language processing and machine learning techniques, is one of the most effective approaches to analyzing public opinion in the digital realm. By employing sentiment analysis, it is possible to determine whether users' views on a specific issue are predominantly positive, negative, or neutral. This information can offer insights into the general mood of the population concerning the topic at hand and assist in predicting future trends or reactions.

In this study, we focus on the viral controversy surrounding the KasKus account “Fufufafa,” which has become a significant topic of discussion on the social media platform X. To analyze user reactions, we utilize the Support Vector Machine (SVM) method for sentiment analysis. SVM is a supervised machine learning algorithm that has proven highly effective for text classification tasks, particularly in sentiment analysis. It operates by finding the hyperplane that best separates different classes of data, in this case, the classes representing positive, negative, and neutral sentiments. SVM’s ability to handle large, complex datasets makes it an ideal choice for this research.

Sentiment analysis, often referred to as opinion mining, is the computational investigation of individuals' opinions, feelings, or attitudes toward a particular entity. This entity could be a person, an event, or a broader topic, and sentiment analysis helps to automatically extract and process the emotional information contained in user-generated text. In the context of social media, sentiment analysis has become a critical tool for understanding public reactions to trending issues, as it allows researchers to sift through vast amounts of data to identify patterns of sentiment and public opinion.

In our research, we collected data from posts, comments, and public discussions on platform X related to the “Fufufafa” issue. This raw data was pre-processed to remove irrelevant information such as spam, non-informative content, and off-topic discussions, ensuring that only meaningful data was analyzed. After pre-processing, the dataset was labeled and categorized into three sentiment classes: positive, negative, and neutral. The labeled data was then used to train and test the SVM model, allowing it to predict sentiment with a high degree of accuracy. The aim was to identify overarching sentiment patterns within the dataset, providing insights into how social media users felt about the issue.

The choice of SVM for our analysis was driven by several considerations. SVM is particularly effective in high-dimensional spaces, which is essential for text data that often features thousands of attributes. Its robustness to overfitting makes it suitable for situations where the number of features exceeds the number of observations, a common scenario in sentiment analysis. Unlike Naive Bayes, which assumes feature independence and may misrepresent correlated features, SVM maximizes the margin between classes, leading to better generalization on unseen data. Moreover, SVM offers flexibility through the use of various kernel functions, allowing it to model non-linear relationships in the data effectively. While neural networks can capture complex patterns, they typically require larger datasets and more computational resources, making SVM a more practical choice for our dataset of 320 observations. Additionally, SVM models can be more interpretable than complex neural networks, as the decision boundaries can be visualized and analyzed, aiding in understanding how sentiment is derived from specific features.

The results of our analysis indicate that the majority of users on platform X expressed negative sentiment regarding the “Fufufafa” controversy, reflecting dissatisfaction or disapproval of the situation. A smaller portion of users displayed neutral sentiments, indicating either

indifference or a balanced view of the issue. An even smaller group of users expressed positive sentiments, suggesting support or approval of the account in question. These findings underscore the polarized nature of online discussions, especially when it comes to viral topics that stir strong emotional reactions.

Through this study, we aim to contribute to the growing body of academic research focused on sentiment analysis and its applications in understanding public opinion in the digital era. By leveraging machine learning techniques like SVM, researchers and organizations can better comprehend the dynamics of public sentiment, particularly in relation to viral social media issues. This research highlights the potential of sentiment analysis to serve as a valuable tool for monitoring and analyzing public discourse, offering a deeper understanding of how viral topics influence public opinion in today's fast-paced, digitally connected world.

2. METHODS

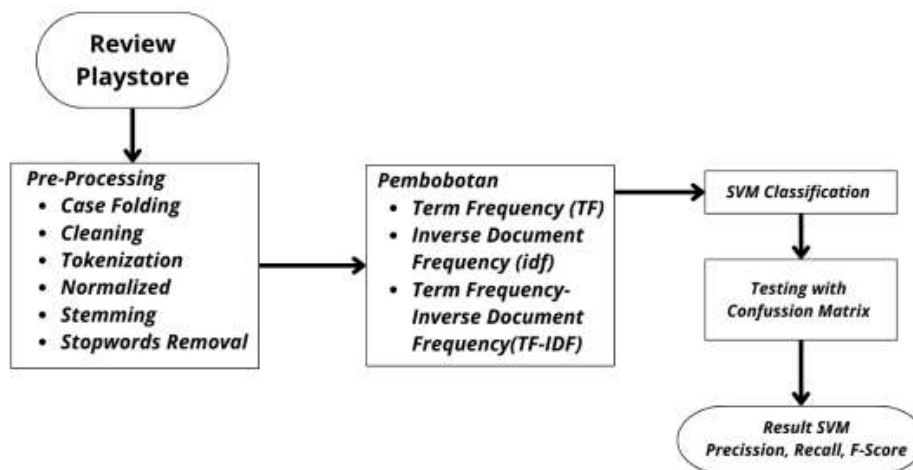


Figure 1 Research Metodology

In Figure 1, the research methodology illustrates the entire research process that has been conducted and consists of the following elements:

2.1 Data Crawling

To perform data crawling of tweets with the keyword “fufufafa,” the Python library used is Tweepy to access the Twitter API. After obtaining the API credentials, authenticate and use the Cursor method to search for tweets. The collected data is stored in a list and then converted into a Data frame for further analysis.

2.2 Preprocessing

Preprocessing is the initial processing step to select words from tweet text data. It involves choosing and removing unnecessary words to obtain a more concise representation of user emotions. To select reviews, we need to eliminate certain elements from the text. Thus, the preprocessing stages include case folding, cleansing, tokenizing, normalization, stemming, and stopword removal [1] [4]

2.3 Weighting

In this stage, each document undergoes a scoring or weighting calculation using the KBBI dictionary or a slang dictionary that contains predefined scores to match with the words present in the document. The matched words are then represented in vector form using TF-IDF, and an information table is created that includes Term Frequency (TF), Document Frequency (DF), and IDF for each term, followed by multiplying TF and IDF. The goal is to obtain labels/sentiments for each term/word in the document.

2.4 SVM Classification

SVM classification is a machine learning technique used for classification and regression tasks. Its primary focus is on classification, where SVM aims to build a model that can separate two classes of data in a high-dimensional space. The objective is to find the optimal hyperplane that separates the two classes. SVM classification can be applied to various types of data, whether two-dimensional or high-dimensional. One of the advantages of SVM is its ability to naturally handle non-linearly separable data using kernel functions [5] [6]. Kernels allow data to be mapped into a higher-dimensional space, making it easier to separate using hyperplanes. As part of the sentiment analysis mentioned above, SVM can be used to classify reviews or text into positive or negative sentiment categories. By learning from labeled training data, SVM can understand patterns and build a model that can be used to predict the sentiment of new text.

2.5 Accuracy Test with Confusion Matrix

Accuracy testing with a confusion matrix is an evaluation method used in machine learning to measure how well a classification model can predict correctly. A confusion matrix provides a detailed view of the model's performance by comparing the true positive, false positive, true negative, and false negative results. By using this information, we can calculate various evaluation metrics such as accuracy, precision, recall, and F1 score, which offer deeper insights into the model's effectiveness in handling different classes or categories.

2.6 Result

The sixth section of this research discusses the results of the sentiment analysis, providing a comprehensive overview of the responses and reactions received from the collected data. This analysis illustrates the feelings and opinions that emerge from various sources, offering valuable insights into the sentiment surrounding the researched topic.

3. RESULTS AND DISCUSSION

3.1 Data Tweet "fufufafa"

The crawling process resulted in 320 tweet data points collected over the period from September 20, 2024, to September 24, 2024. An example of the documents obtained from the crawling can be seen in Table 1 and in Fig 1. Data tweet crawling.

index	id	label	tweet
199	200	0	@Cek_Kancid007 @evolerasiem @piknavigowin @SpartanTaurus @her_erni @Mamamochel @Maudy_Apan @_Ayu_ayu Woodi Fufufafa @kesferobem ki dah rindu kali miss Borbenu itu gerakannya terbang naik Jet Gutschein ign dibarkan nangis-kabong kau yaa kawatiran bro @Bera_Rc satu arah k sana tp kau hrs hrs bawa deodoran km di @malis_jenny ikut serta kuantitas ada aroma keke
200	201	0	Mih bs monevof lg stih Fufufafa ont nemanpin
201	202	0	@BosPurwa mi bukan Inggris tetapi adalah Kalo FUFUFafa dan jet pribadi itu baru banyak ingkrah apeski yg kaperlan
202	203	0	@Gustawan75_Maafkita janggg... obaknya s fufufafa mo dicongkol para obelggg... knapa ga dipatu ajaa ??
203	204	0	Dangggas sembunyikan kebetaran terkait Fufufafa Anonymous bocoran data pribadi Menkominfo Budi Ane Setiadi @budi_waras @Dobstantan @BosPurwa @denocrazymedia @dhamid_h_back @Herakobas @marina_scha @DokterTita https://t.co/Pz34Ys554

Figure 2 Data Tweet

Table 1 Example Data Tweet

No	Term
Doc 1	Kakaknya kaesang namanya gibran Apa bener yg punya fufufafa
Doc 2	Biasanya gercep pencitraan bnyk orang
Doc 3	cinta amat sama Anies nanti FUFUFABA cemburu lo

3.2 Preprocessing

1. Case Folding dan Cleansing

The processes of case folding and cleansing are critical steps in text data preprocessing. Case folding involves normalizing text by converting all letters to lowercase, while cleansing encompasses cleaning the data from unwanted characters or elements [7]. Both steps aim to ensure consistency and cleanliness of the data, preparing it for further analysis without unnecessary variations or noise. The results of case folding and cleansing can be seen in Table 2.

Table 2 Case Folding dan Cleansing

No	Term
Doc 1	kakaknya kaesang namanya gibran apa bener yang punya fufufafa
Doc 2	biasanya gerak cepat pencitraan banyak orang
Doc 3	cinta amat sama anies nanti fufufafa cemburu kamu

2. Tokenization

Tokenization is the process of breaking down text or sentences into smaller, manageable units known as tokens. These tokens can be individual words, phrases, or even characters, depending on the level of granularity required for a particular task. In the context of programming and natural language processing (NLP), tokenization serves as a foundational step in analyzing and processing textual data. By converting a continuous stream of text into distinct elements, tokenization allows for a clearer representation of the text's structure and meaning, facilitating subsequent operations such as text classification, sentiment analysis, machine translation, and more.

There are different approaches to tokenization. For instance, word tokenization splits a text into individual words, which is useful for tasks where word-level analysis is needed. Character tokenization, on the other hand, breaks the text into single characters, which can be useful in language models dealing with specific alphabets or characters. Phrase tokenization divides the text into multi-word units or phrases, capturing meaning beyond single words.

The tokenization process also takes into account certain factors like punctuation, special characters, and spaces to ensure that each token is meaningful and relevant for further analysis. After tokenization, the text is typically converted into a format that can be more easily processed by machine learning algorithms or other computational tools.

In this study, tokenization was applied to the dataset, and the results are presented in Table 3, which showcases how the original text was transformed into tokens for further analysis. Tokenization thus plays a critical role in preparing the data for the sentiment analysis task.

Table 3 Tokenization

No	Term
Doc 1	['kakaknya', 'kaesang', 'namanya', 'gibran', 'apa', 'bener', 'yang', 'punya', 'fufufafa']
Doc 2	['biasanya', 'gerak', 'cepat', 'pencitraan', 'banyak', 'orang']
Doc 3	['cinta', 'amat', 'sama', 'anies', 'nanti', 'fufufafa', 'cemburu', 'kamu']

3. Normalization

Normalization is the process of data processing aimed at bringing variable values into a standard range or uniform scale. The goal is to eliminate magnitude differences between variables, ensuring that each variable has a balanced impact in the analysis and supporting consistent interpretation. The results of normalization can be seen in Table 4.

Table 4 Normalization

No	Term
Doc 1	['kakak', 'kaesang', 'nama', 'gibran', 'apa', 'benar', 'yang', 'punya', 'fufufafa']
Doc 2	['biasanya', 'gerak', 'cepat', 'pencitraan', 'banyak', 'orang']
Doc 3	['cinta', 'amat', 'sama', 'anies', 'nanti', 'fufufafa', 'cemburu', 'kamu']

4. Stemming

Stemming is the process of removing prefixes or suffixes from words that include conjunctions, prepositions, and pronouns, resulting in the corresponding root words as defined by the Indonesian Dictionary (KBBI). The results of the stemming process are presented in Table 5.

Table 5 Stemming

No	Term
Doc 1	['kakak', 'kaesang', 'nama', 'gibran', 'apa', 'benar', 'yang', 'punya', 'fufufafa']
Doc 2	['biasa', 'gerak', 'cepat', 'citra', 'banyak', 'orang']
Doc 3	['cinta', 'amat', 'sama', 'anies', 'nanti', 'fufufafa', 'cemburu', 'kamu']

5. Stopwords Removal

This process involves filtering out words in the documents that are unrelated to sentiment analysis, as presented in Table 6.

Table 6 Stopwords removal

No	Term
Doc 1	['kakak', 'kaesang', 'nama', 'gibran', 'benar', 'punya', 'fufufafa']
Doc 2	['biasa', 'gerak', 'cepat', 'citra', 'banyak', 'orang']
Doc 3	['cinta', 'amat', 'anies', 'fufufafa', 'cemburu']

After completing the preprocessing process with term cleansing, the next step is to perform TF-IDF weighting, which is explained in Table 7. TF-IDF weighting.

Table 7 TF-IDF Weighting

Term	TF			DF	IDF	TF x IDF		
	D1	D2	D3			D1	D2	D3
kakak	1	0	0	1	0,47	0,47	0	0
kaesang	1	0	0	1	0,47	0,47	0	0
nama	1	0	0	1	0,47	0,47	0	0
gibran	1	0	0	1	0,47	0,47	0	0
benar	1	0	0	1	0,47	0,47	0	0
punya	1	0	0	1	0,47	0,47	0	0
fufufafa	1	0	1	2	0,18	0,18	0	0,18
biasa	0	1	0	1	0,47	0	0,47	0
gerak	0	1	0	1	0,47	0	0,47	0
cepat	0	1	0	1	0,47	0	0,47	0
citra	0	1	0	1	0,47	0	0,47	0
banyak	0	1	0	1	0,47	0	0,47	0
orang	0	1	0	1	0,47	0	0,47	0
cinta	0	0	1	1	0,47	0	0	0,47
amat	0	0	1	1	0,47	0	0	0,47
anies	0	0	1	1	0,47	0	0	0,47
cemburu	0	0	1	1	0,47	0	0	0,47

From Table 7, the TF-IDF results for the documents used as test data in this research are obtained. The TF column indicates how often a word/term appears in the test document. The IDF column represents the result of $IDF = \log \left(\frac{N}{Df} \right)$ [8], sample term is “kakak” $IDF = \log \left(\frac{3}{1} \right) = 0,477$, and the result TF-IDF with $TF \times IDF = 1 * 0,477 = 0,477$.

3.3 Classification SVM

In the sentence “cinta amat sama Anies nanti FUFUFABA cemburu lo,” after undergoing the preprocessing process, it becomes ['cinta', 'amat', 'anies', 'fufufafa', 'cemburu']. In this term, it forms the matrix $K(x, x_1)$, which represents the frequency of word occurrences in the test sample with a total of 12 terms. The matrix $K(x, x_1)$ can be seen in Table 8.

Table 8 Table Matrik $K(x, x_1)$

	X1	X2
$K(x, x_1)$	0,477	0,477

From Table 8, the results to obtain the values of x and x_1 are then presented in Table 9.

Table 9 Table Matrik x,x1

0,477	0,477	0,47+0,47=0,94
1	1	

Table 9 presents the results of each term against the term weights, namely:

$$= a_1y_1k(x,x_1) + a_2y_2k(x,x_2) + b$$

$$= (1 * (1) * 0.47) + (1 * (1) * 0.47) + (0) = 0,94$$

$$\text{maka, } = f(x) \text{ sign } (0,94) = 1$$

Thus, the new sentence “cinta amat sama Anies nanti FUFUFABA cemburu lo” (1) represents a positive class.

3.4 Accuracy Test with Confusion Matrix

The testing is conducted using a confusion matrix through the SVM library, based on the previously classified training and testing models. This results in a matrix of size 2x2, representing the actual classes and predicted classes. The results from the training model using new data that has not been trained previously are presented in the confusion matrix shown in Table 10.

Table 10 Table Confusion Matrix

Data	Prediction	
	Positive	Negative
Positive	TP : 87	FP : 0
Negative	FN : 15	TN : 0

In the confusion matrix, TP, FP, TN, and FN are abbreviations for each type of prediction result in the context of classification:

1. True Positive (TP): The number of positive observations that are correctly predicted by the model. This means the model accurately identifies positive examples.
2. False Positive (FP): The number of negative observations incorrectly predicted as positive by the model. This means the model gives a positive prediction when it should have been negative.
3. True Negative (TN): The number of negative observations that are correctly predicted by the model. This means the model accurately identifies negative examples.
4. False Negative (FN): The number of positive observations incorrectly predicted as negative by the model. This means the model gives a negative prediction when it should have been positive [9] [10].

Thus, the accuracy value of the SVM method can be calculated as follows:

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN} = \frac{87}{102} = 0,85 \quad (1)$$

3.5 Result SVM

The results of the experiment with 320 data points yield values for accuracy, precision, recall, and F-1 Score, which can be seen in Figure 3.


```

/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))
/usr/local/lib/python3.10/dist-packages/sklearn/metrics/_classification.py:1531: UndefinedMetricWarning: Precision is ill-defined and being set to 0.0 in
_warn_prf(average, modifier, f"{metric.capitalize()} is", len(result))

Classification Report:
      precision    recall  f1-score   support

     0       0.85       1.00       0.92         87
     1       0.00       0.00       0.00         15

 accuracy          0.43
 macro avg          0.43
 weighted avg       0.73

Confusion Matrix:
[[87  0]
 [15  0]]

Accuracy Score: 0.8529411764705882

```

Figure 3 SVM Testing

The experiment yielded values for several key performance metrics, including accuracy, precision, recall, and F1 score, which are summarized in Figure 3. The Support Vector Machine (SVM) method, using Scikit-learn's Term Frequency-Inverse Document Frequency (TF-IDF) for feature extraction, demonstrated strong performance, achieving an accuracy score of 0.85. This high accuracy indicates that the SVM model was effective in correctly classifying the sentiment of user data in the experiment.

When compared to the alternative feature extraction method, Term-presence, the SVM model achieved notably better results. The accuracy score for SVM using Term-presence was only 0.79, demonstrating that TF-IDF provided a superior representation of the data, leading to improved model performance. This difference highlights the importance of feature selection in text classification tasks, as the choice of representation can significantly impact the effectiveness of machine learning models like SVM. Overall, the experiment underscores the strength of the SVM method combined with TF-IDF for sentiment analysis tasks..

4. CONCLUSIONS

This study aims to evaluate the effectiveness of the Support Vector Machine (SVM) method in conducting sentiment analysis using a dataset consisting of 320 data points. Upon examining the dataset during the data processing stage, it was observed that the data distribution was highly imbalanced. Specifically, 85.71% of the data represented negative sentiment, while only 14.29% of the data was classified as positive sentiment. This imbalance is a critical aspect of the dataset's characteristics and provides an important foundation for assessing the model's ability to handle skewed data distributions in sentiment classification tasks.

In sentiment analysis, the evaluation of the model's performance is crucial for understanding its accuracy in predicting sentiment categories. Key performance metrics used in this study include precision, recall, F1 score, and overall accuracy. These metrics offer a comprehensive view of how well the SVM model performs, not only in terms of accuracy but also in its ability to correctly identify true positives and minimize false positives and false negatives.

The results indicate that the SVM model achieved an accuracy score of 85%, which is considered a strong performance, especially given the imbalanced nature of the dataset. This high accuracy suggests that the SVM model is capable of effectively identifying sentiment, even when the majority of the data falls into one category.

Additionally, a comparative analysis was performed against the Term Presence method, which is another feature extraction technique used in text classification. The results showed that SVM outperformed Term Presence, which only achieved an accuracy of 79%. This comparison underscores the superiority of SVM, particularly when paired with robust feature extraction techniques like TF-IDF, in handling sentiment analysis tasks. The findings highlight the potential of SVM in achieving high performance, even when working with imbalanced datasets, making it a valuable tool for sentiment analysis.

REFERENCES

- [1] C. J. C. A. K. Santra, "Genetic Algorithm and Confusion Matrix for Document Clustering," *IJCSI*, vol. 9, no. Bharathiar University, Coimbatore, pp. 322-328, 2012.
- [2] M. D. C. Z. A. Salappa, "Feature selection algorithms in classification problems: an experimental evaluation," *Optimization Methods and Software*, 2005.
- [3] S. M. F. B. G. R. Y. P. K. Abdul Manan Iddrisu, "International Journal of Information Management Data Insights A sentiment analysis framework to classify instances of sarcastic sentiments within the aviation sector," *International Journal of Information Management Data Insights*, vol. 3, no. 2, 2022.
- [4] B. H. I. G. P. A. J. Z. Abd. Samad Hasan Basari, "Opinion Mining of Movie Review using Hybrid Method of Support Vector Machine and Particle Swarm Optimization," *Procedia Engineering*, vol. 53, no. Universiti Malaysia, pp. 453-462, 2013.
- [5] E. S. M. A. U. I. M. T. A. M. I. S. & A. M. H. Ahmad, "Challenges, comparative analysis and a proposed methodology to predict sentiment from movie reviews using machine learning," *ICBDAC*, vol. 10, pp. 86-91, 2017.
- [6] A. B. W. D. H. Anto Satriyo Nugroho, "Support Vector Machine Teori dan Aplikasinya," *Bioinformatika*, 2003.
- [7] S. B. Atang Saepudin, KOMPARASI ALGORITMA SUPPORT VECTOR MACHINE DAN K-NEAREST NEIGHBOR BERBASIS PARTICLE SWARM OPTIMIZATION PADA ANALISIS SENTIMEN FENOMENA TAGAR #2019GANTIPRESIDEN, JAKARTA: STMIK NUSA MANDIRI, 2018.
- [8] A. S. T. B. Azhagusundari, "Feature Selection based on Information Gain," *International Journal of Innovative Technology and Exploring Engineering*, vol. 2, pp. 18-21, 2013.
- [9] J. C. M. W. Y. & P. A. Chou, "Expert Systems with Applications Optimizing parameters of Support Vector Machine using fast messy genetic algorithm for dispute classification," *Expert Systems With Applications*, vol. 41, 2014.
- [10] Dawson, "Introduction to Research Methods: A Practical Guide for Anyone Undertaking a Research Project," *Oxford: How to Books*, 2009.
- [11] R. M. H. J. A. & S. Y. Dehkharghani, "Expert Systems with Applications Sentimental causal rule discovery from Twitter.," *EXPERT Systems With Applications*, vol. 41, 2014.