



Research Article

# Feature Selection and Performance Evaluation of Buzzer Classification Model

Dian Isnaeni Nurul Afra <sup>a</sup>, Radhiyatul Fajri <sup>a,\*</sup>, Harnum Annisa Prafitia <sup>a</sup>, Ikhwan Arief<sup>b</sup>, Aprinaldi Jasa Mantau <sup>c</sup><sup>a</sup> Research Center for Data and Information Sciences, National Research and Innovation Agency, Jakarta, Indonesia<sup>b</sup> Department of Industrial Engineering, Universitas Andalas, Padang, Indonesia<sup>c</sup> Department of Computer Science and System Engineering (CSSE), Kyushu Institute of Technology, Fukuoka, Japan\* Corresponding Author: [radh001@brin.go.id](mailto:radh001@brin.go.id)

© 2024 Authors

DOI: [10.25077/josi.v23.n1.p1-14.2024](https://doi.org/10.25077/josi.v23.n1.p1-14.2024)

Submitted : November 16, 2023; Accepted : June 8, 2024; Published : July 10, 2024

## ABSTRACT

In the rapidly evolving digital age, social media platforms have transformed into battleground for shaping public opinion. Among these platforms, X has been particularly susceptible to the phenomenon of 'buzzers', paid or coordinated actors who manipulate online discussions and influence public sentiment. This manipulation poses significant challenges for users, researchers, and policymakers alike, necessitating robust detection measures and strategic feature selection for accurate classification models. This research explores the utilization of various feature selection techniques to identify the most influential features among the 24 features employed in the classification modeling using Support Vector Machine. This study found that selecting 11 key features yields a remarkably effective classification model, achieving an impressive F1-score of 87.54 in distinguishing between buzzer and non-buzzer accounts. These results suggest that focusing on the relevant features can improve the accuracy and efficiency of buzzer detection models. By providing a more robust and adaptable solution to buzzer detection, our research has the potential to advance social media research and policy. This enabling researchers and policymakers to devise strategies aimed at mitigating misinformation dissemination and cultivating an environment of trust and integrity within social media platforms, thus fostering healthier online interactions and discourse.

**Keywords:** buzzer classification, feature selection, spearman correlation, pearson correlation, chi-square test, social media

## INTRODUCTION

The rapid growth of social media platforms has enabled users to share information, news, and opinions on various topics. Social media platforms have emerged as powerful conduits for disseminating information, news, and updates rapidly. Among these platforms, X (formerly Twitter) stands out as a particularly influential medium, providing a vast repository of real-time, unfiltered, and diverse user-generated content [1], [2]. It enables people to connect, communicate, and express their thoughts to the world. Social media offers a wide range of benefits, impacting various aspects of personal, professional, and societal life, including an unprecedented opportunity to gauge public opinions and sentiments on a myriad of topics [3], [4]. In today's world, it has become commonplace for most government bodies and offices to utilize social media for regular communication or dialogue with citizens. In fact, these platforms have even become the standard for such interactions. Citizen engagement in social media can offer valuable insights for the government to monitor and evaluate public feedback. Recent research indicates that the two-way interaction facilitated by social media enables the government to disseminate messages and receive public feedback. Consequently, the utilization of social media in government administration can enhance public participation in public policy formulation [5], [6]. Moreover, researchers and policymakers employ sentiment analysis to gauge public sentiment towards critical societal issues, enabling evidence-based decision-making and policy formulation [7], [8]. However, currently, social media encounters a distinctive obstacle and creates opportunities for malicious actors to manipulate public opinion and influence online discussions known as the buzzer phenomenon.

The buzzer phenomenon in social media refers to the practice of individuals or groups who manipulate information on social media to influence public opinion, deliberately generating artificial buzz or hype around specific topics, products, or issues with the intention of manipulating public opinion and influencing online discussions. Buzzers are individuals or collectives intentionally creating artificial excitement or attention around particular subjects, products, or concerns with the intention of serving their own or their sponsors' interests [9]–[11]. They can use various strategies, such as raising issues, supporting emerging issues, creating disinformation or hoaxes, attacking a person or group, changing issues that are currently viral, and building the image of a figure [6]. They can work individually or as a team, using human resources and bot accounts [1]. Buzzers are often paid or coordinated actors who work to amplify certain messages, promote particular ideologies, or spread misinformation with the aim of promoting their own interests or those of their sponsors [9]–[11].

The buzzer phenomenon poses significant challenges for social media platforms, researchers, and users alike. The use of buzzers can lead to a post-truth situation where there is no absolute truth, but the truth is realized from the victory of the propaganda action of certain groups [6]. It can lead to the spread of misinformation, undermine trust in online information sources, and disrupt the healthy exchange of ideas. This can lead to harmful impacts if not managed properly, such as causing conflicts in society due to propaganda wars on social media.

Given the potential negative impacts, it is crucial to detect and regulate the activities of buzzers. Previous study suggests that there need to be legal rules that apply specifically to actors who play in the digital world of social media [6]. These regulations are not meant to limit creativity or democratic values, but to minimize the emergence of social conflicts that can threaten the stability and security of the country. The strategic role of buzzers in shaping public sentiment and influencing online discussions has been recognized, leading to the necessity of buzzer detection. Detecting and mitigating the impact of buzzers requires the implementation of robust content moderation systems, user education about recognizing and reporting misleading content, and the development of advanced algorithms to identify and address deceptive activities on social media [12], [13]. Finally, developing a buzzer classification model is necessary to automatically identify and classify buzzers based on the selected features. This model can help in proactively monitoring and moderating content on social media platforms, thereby fostering a more authentic and constructive environment for discussions and information sharing.

This paper concentrates on the phenomenon of buzzers on platform X. This is due to the platform's provision of real-time content, ease of capturing trending posts, and its popularity as a means for community discussions on various topics [1], [2]. The data accessible from X comprises a wide array of attributes that can serve as features when constructing models to differentiate between buzzer and non-buzzer accounts. Nevertheless, not all of these attributes truly contribute to bolstering the accuracy of the buzzer classification model. Therefore, it is necessary to do feature selection to select the most relevant and discriminative attributes from the metadata, which can significantly impact the performance of the detection algorithm. By choosing the right features, the detection model becomes more efficient, reducing computational overhead and improving the speed of analysis [14].

Feature selection is an essential step in this process. It involves identifying the most relevant features that contribute to a user being a buzzer. This step is crucial to ensure the efficiency and accuracy of the detection model. The primary objective of this study is to explore various feature selection techniques and assess their impact on the performance of the buzzer classification model through rigorous experimentation with classification algorithm. This experimentation endeavors to discern the salient features essential for identifying buzzer accounts, thereby optimizing the buzzer classification modeling process. Simultaneously, less significant features can be excluded from the modeling to enhance overall accuracy and efficiency.

Numerous prior studies have focused on the subject of recognizing buzzers or computer bots, driven by the goal of mitigating the adverse consequences arising from the dissemination of bot or buzzer accounts across social media platforms. Kantepe and Ganiz [13] emphasized the significant issue of social media platforms, particularly X, being plagued by a considerable number of social bot accounts controlled by automated agents to engage in various malicious activities, such as spamming, spreading misinformation, recruiting individuals for illegal organizations, and blackmailing to disseminate private information. To address this concern, the authors present a bot detection method utilizing a machine learning algorithm.

The study involves the extraction of 62 features from each collected account. Subsequently, feature selection is performed using three distinct techniques: Information Gain, Mutual Information, and Chi-Square, resulting in 11 final features for bot classification modeling. The study used four machine learning algorithms—Logistic Regression, Multinomial Naive Bayes, Support Vector Machine, and Gradient Boosted Trees—to develop a bot classification model [13]. In contrast to the previous research conducted by Kantepe and Ganiz, which did not thoroughly investigate all features and their impact on the classification model's performance, this study aims to explore these aspects more comprehensively. Furthermore, in that study, the results of feature selection using Information Gain and Mutual Information showed less favorable outcomes, with the highest values being 0.38 using Information Gain and 0.29 using Mutual Information, out of a maximum value of 1. Therefore, in this research, an attempt was made to experiment with other feature selection techniques, namely Spearman and Pearson Correlation, which demonstrated favorable results in studies related to feature selection techniques in classification modeling [15], [16].

Another study [11] developed an automatic buzzer detection using machine learning algorithms before conducting political sentiment analysis. From the beginning, this research has predetermined seven features to be used in the model's development, hence feature selection methods were not applied in this study. Therefore, it is necessary to conduct further experiments and in-depth analysis to determine whether the features used are indeed correlated with the characteristics of buzzer accounts using feature selection techniques.

The research conducted by Panatra et al [17] provides an overview of the social media landscape and presents a method in the form of a process diagram for performing buzzer detection, specifically for Instagram data. The features are categorized into four sections: post times, images, hashtags, commonly followed accounts, and frequency of posts. Subsequently, these features are classified using the Naive Bayes, Support Vector Machine, and Random Forest methods. Likewise, within a study conducted by Suciati et al [18], a model for identifying buzzers is formulated utilizing a dataset derived from Twitter, specifically focusing on the context of the Indonesian presidential election. The investigation entails the application of a singular feature selection approach, denoted as Mutual Information. This strategy culminates in achieving a preeminent accuracy metric of 62.3%, a result attained through the utilization of a subset of 25 features in conjunction with the AdaBoost model. In the present inquiry, we intend to delve into the examination of alternative feature selection methodologies.

Although previous research has made significant contributions to understanding and addressing the buzzer phenomenon, there still exists a critical research gap. Existing research on buzzer detection often involves the use of too many features or employs a limited set of features without thoroughly investigating their effectiveness in capturing the distinctive behaviors of buzzer accounts. The use of an excessive number of possibly ineffective features can lead to the development of an inefficient buzzer classification model.

This study aims to fill this research gap by proposing a novel approach to buzzer detection. Our approach is characterized by its comprehensive use of various feature selection techniques. The objectives of this study are:

- Explore different techniques for feature selection: Explore different ways to select relevant features for a buzzer classification model.
- Evaluate the impact on model performance: Thoroughly experiment with the classification algorithms to evaluate the impact of these feature selection techniques on the performance of the buzzer classification model.
- Identify salient features: Identify key features essential to accurately identify buzzer accounts.
- Optimize the modeling process: Improve the accuracy and efficiency of the buzzer classification model by excluding less important features.

The novelty of our research lies in its methodological rigor and its focus on buzzer behavior in Indonesia through experiments and analysis of profile metadata and user posts on social media. This is a relatively new and specific topic that has not been extensively studied before. By providing a more robust and adaptable solution to buzzer detection, our research has the potential to significantly advance the field of social media research and policy. Platform administrators can use our findings to develop more effective algorithms for identifying and mitigating the impact of inauthentic accounts, thereby improving the integrity of their platforms. Additionally, society at large

stands to gain from enhanced transparency and trust in social media interactions, reducing the spread of misinformation and fostering a more informed public discourse. Effective buzzer detection can greatly benefit policymakers by enabling more informed decision-making regarding the regulation of social media activity. Through these contributions, our research underscores its practical implications and relevance to key stakeholders.

## METHODS

Figure 1 shows the steps of the methodology used in this research. A detailed explanation of each stage is provided below.

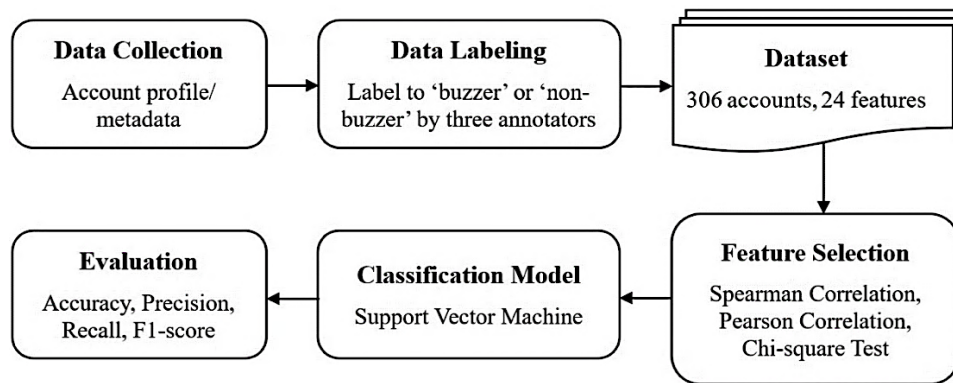


Figure 1. Workflow Diagram of the Proposed Methodology

### Data Collection and Labelling

This research utilizes a dataset about the New Indonesia's Capital City (Ibu Kota Negara Indonesia) obtained from the study by Pebiana et al [19]. This dataset contains 16,300 tweets in the Indonesian language related to the New Indonesia's Capital City. From this dataset, a list of usernames is extracted, which is subsequently used for crawling data using X API to gather profile information from these usernames. The profile data in question consists of metadata associated with each X account, such as the username's ID, tweet count, location, account creation date, URL, profile picture, and more.

Upon completing the crawling process and obtaining the account profile data, the next step involves labeling by three annotators to categorize each account into two classes: "buzzer" and "non-buzzer". Each annotator is asked to label each account based on the profile data and tweet data obtained from the previous data collection stage. The final labels used are determined through majority voting among the three assigned annotators. Based on the labeling results, a dataset is compiled, consisting of a total of 306 accounts, including 130 buzzer accounts and 176 non-buzzer accounts. Next, in the subsequent stage, which is the classification model creation, the labels for buzzers will be converted into numeric values, namely 1, while non-buzzers will be labeled as 0.

### Feature Selection

The process of feature selection holds paramount importance within the machine learning pipeline as it plays a crucial role in enhancing model performance. This step entails pinpointing the dataset's most informative features that substantially contribute to the model's predictive capabilities. In the context of our research, we utilize diverse feature selection techniques to discern the most pertinent and differentiating attributes from the metadata associated with X accounts. This comprehensive approach enables us to refine our understanding of the key factors influencing the model's predictive accuracy and ensures that the selected features are both relevant and impactful in capturing the nuances of the dataset.

The feature selection process in our study begins with the extraction of 24 distinct features from the account metadata. These features will be processed in the subsequent phase of our study. Two of these features are particularly

Table 1. Feature Description

No.	Feature	Description	Data Type	No.	Feature	Description	Data Type
1.	all_tweets	Total number of tweets	Numerical	13.	count_replies	Number of replies received	Numerical
2.	followers	Total number of followers	Numerical	14.	count_likes	Number of likes received	Numerical
3.	following	Total number of followings	Numerical	15.	count_retweets	Number of retweets received	Numerical
4.	bio	Indicates bio on profile	Nominal	16.	count_reply_to	Number of tweets that are replies to other tweets	Numerical
5.	bg_image	Uses background image on profile	Nominal	17.	count_mentions	Number of mentions in created tweets	Numerical
6.	profile_image	Uses profile picture on profile	Nominal	18.	count_photos	Number of photos in created tweets	Numerical
7.	location	Indicates location on profile (0: no, 1: yes)	Nominal	19.	count_hashtags	Number of hashtags in created tweets	Numerical
8.	url	Includes URL on profile (0: no, 1: yes)	Nominal	20.	count_urls	Number of URL links in created tweets	Numerical
9.	days_created	Account age since creation (days)	Numerical	21.	url_to_web	Number of URL links referring to web pages in created tweets	Numerical
10.	all_likes	Number of liked tweets	Numerical	22.	url_to_tweet	Number of URL links referring to other tweets in created tweets	Numerical
11.	all_media	Number of tweets containing media (photos, videos, etc.)	Numerical	23.	TIE_diff_interval	Time Interval Entropy	Numerical
12.	tweet	Number of tweets within the crawling time range	Numerical	24.	avg_similarity	Similarity between tweets	Numerical

noteworthy. The first, 'TIE\_diff\_interval', represents the time interval entropy. This measures the regularity of tweet timings to identify patterns that may indicate automated posting. The second feature, 'avg\_similarity', calculates the average similarity between one tweet and all other tweets from the same account. This helps us understand the level of repetition or variation in the content posted by each account. The specifics of each of these 24 features, along with their descriptions, are presented in Table 1. The data type of each feature is also explained in the table because the data type determines the feature selection technique that will be used in the next stage. Numerical data represents values that are measurable and can be expressed as numbers, while nominal data, also known as categorical data, represents values that label or categorize without implying any quantitative relationship between the categories. This type of data is used to name, label, or categorize attributes or items.

Among the abundance of features at our disposal, the process of feature selection is imperative to sift through and identify the significant and influential attributes while eliminating those deemed irrelevant. The primary goal of feature selection is to improve the overall performance of the classification model and concurrently reduce computation time. This meticulous selection process ensures that the chosen features not only contribute substantially to the model's predictive accuracy but also streamline computational efficiency, fostering a more effective and resource-efficient classification system. [14]. Several feature selection techniques are employed in this study, encompassing the following methods.

- Spearman Correlation [20], [21]. Spearman's correlation coefficient, known as the Spearman's rank correlation coefficient, is a statistical tool employed to evaluate the intensity and direction of a monotonic connection

between a pair of variables. Spearman Correlation, a non-parametric measure of statistical dependence between two variables, offers several advantages in various analytical scenarios. One notable advantage is its robustness to outliers. This property makes Spearman Correlation a suitable choice when dealing with datasets that may contain irregularities or extreme values. Additionally, Spearman Correlation is particularly useful when the relationship between variables is monotonic but not necessarily linear. It doesn't assume a linear association between variables, making it applicable to a broader range of data distributions. This flexibility is beneficial when working with real-world datasets where the underlying patterns may not strictly adhere to linear relationships. Another advantage lies in its applicability to ordinal data. Spearman Correlation is well-suited for assessing the strength and direction of monotonic relationships in variables measured on ordinal scales. This characteristic extends its utility beyond scenarios where only interval or ratio data is available, making it versatile for various types of statistical analyses. Furthermore, Spearman Correlation doesn't require the assumption of normality, making it a valuable tool in cases where the data distribution deviates from a normal distribution. This non-parametric nature allows for a more robust analysis in situations where parametric methods may not be suitable or reliable.

In summary, Spearman Correlation offers advantages such as robustness to outliers, suitability for monotonic relationships, applicability to ordinal data, and independence from the assumption of normality. These characteristics make it a valuable statistical tool in diverse analytical contexts. This technique is non-parametric in nature, implying that it avoids presumptions about data distribution and uses the rank order of the data. It is suitable for both ordinal and continuous data and doesn't assume any specific distribution of data. Spearman's rank correlation coefficient ( $\rho$ ) is calculated using the following Equation (1).

$$\rho = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2-1)} \quad (1)$$

where  $d_i$  is the difference between the ranks of corresponding data points and  $n$  is the number of data points.

The Spearman correlation coefficient spans from -1 to 1, with 1 signifying an impeccable positive monotonic relationship, wherein augmented values of one variable align with increased values of the other variable. Conversely, -1 represents an impeccable negative monotonic relationship, indicating that higher values of one variable coincide with diminished values of the other. A coefficient of 0 denotes an absence of a monotonic relationship, implying a lack of correlation between the variables.

- Pearson Correlation [22], [23]. Pearson's correlation coefficient serves as a statistical tool to numerically assess the intensity and orientation of a linear connection between two continuous variables. It ranks among the frequently employed correlation coefficients and is an integral component of parametric statistical methods. Pearson's correlation proves valuable in situations where a linear relationship between two variables is anticipated, especially when the data adheres to a fairly normal distribution pattern. One of the key advantages of Pearson Correlation is its simplicity and ease of interpretation. The coefficient provides a clear indication of the strength and direction of the linear relationship between variables. A positive correlation suggests a direct relationship, while a negative correlation implies an inverse relationship. Additionally, Pearson Correlation is sensitive to linear relationships, making it suitable for detecting and quantifying the degree of linear association between variables. This sensitivity allows researchers to identify whether changes in one variable are associated with systematic changes in another, providing valuable insights into patterns and trends. Moreover, Pearson Correlation is widely used in various fields such as economics, psychology, biology, and social sciences. Its widespread application contributes to the convenience of comparing relationships across diverse domains, facilitating interdisciplinary research and analysis. When linearity is present, Pearson Correlation is a powerful and straightforward tool for statistical analysis. Equation (2) is the formula to calculate Pearson's correlation coefficient.

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (2)$$

where  $x_i$  and  $y_i$  represent individual data points for the respective variables,  $\bar{x}$  and  $\bar{y}$  denote the means (averages) of the  $x$  and  $y$  data points, and  $n$  is the number of data points.

The Pearson correlation coefficient ranges from -1 to 1. A value of  $r = 1$  signifies a flawless positive linear correlation, indicating that as one variable rises, the other simultaneously increases in proportion. On the contrary, a value of  $r = -1$  denotes a complete negative linear correlation, signifying that as one variable advances, the other correspondingly declines. A value of  $r = 0$  indicates no linear correlation between the variables; they are not related in a linear manner.

- Chi-Square Test [24], [25]. The Chi-Square test stands as a statistical technique employed to ascertain the presence of a notable link between two categorical variables. This method finds frequent application in the scrutiny of data structured within contingency tables, enabling the exploration of connections between two categorical variables. One of its primary advantages lies in its versatility and applicability to different types of data. This test does not make assumptions about the distribution of the data, making it robust in various situations. It is particularly useful when dealing with nominal data, where variables are divided into distinct categories. Furthermore, the Chi-Square Test is valuable for analyzing large datasets and complex relationships between variables. It allows researchers to explore patterns and dependencies among categorical variables within contingency tables. The test can be applied to uncover relationships in fields such as social sciences, medicine, and market research, providing insights into whether observed differences are statistically significant. Additionally, the Chi-Square Test is relatively easy to understand and implement. Its simplicity makes it accessible to researchers with varying levels of statistical expertise, and it can be used to test hypotheses in situations where other statistical methods may not be as suitable. This ease of use contributes to its widespread application in different research domains. Moreover, the Chi-Square Test is non-parametric, meaning it doesn't require assumptions about the underlying distribution of the data. This non-parametric nature makes it robust when dealing with data that may not adhere to normal distribution assumptions, providing a valuable tool in situations where parametric tests might be less appropriate.

In summary, the Chi-Square Test offers versatility, applicability to various types of data, ease of use, and the ability to uncover significant associations in categorical variables without stringent distribution assumptions, making it a valuable statistical tool in diverse research contexts. The formula to calculate the Chi-Square statistic ( $\chi^2$ ) is defined in Equation (3).

$$\chi^2 = \sum_{i=1}^k \frac{(O_i - E_i)^2}{E_i} \quad (3)$$

where  $O$  represents the observed frequency in a specific cell of the contingency table,  $E$  represents the expected frequency in the same cell under the assumption of independence, and  $k$  is the number of categories.

In Chi-Square calculations, two hypothesis formulas are present: the Null Hypothesis ( $H_0$ ), which states that there is no significant association between the variables, and the Alternative Hypothesis ( $H_a$ ), which posits a significant association between the variables. The results of the Chi-Square calculation are compared against a specific significance level of 0.05, where if the calculated Chi-Square statistic surpasses the critical value, the null hypothesis is rejected, leading to the conclusion that a significant association exists between the variables.

## Classification Model

In this study, one of the powerful classification algorithms, namely Support Vector Machines (SVM), is employed. The primary benefit of SVM is their ability to handle complex and high-dimensional data. SVM is particularly useful when dealing with datasets that have many features or dimensions. It can efficiently handle situations where the number of features is greater than the number of data points. Additionally, SVMs excel in situations where the data is not linearly separable. Through the use of kernel functions, SVMs can implicitly map data into higher-dimensional spaces, allowing for the separation of complex patterns that might not be discernible in the original feature space. This ability to handle non-linear relationships in data contributes to the flexibility and robustness of SVMs.

In conclusion, SVM offers advantages such as effectiveness in high-dimensional spaces, ability to handle non-linear data, and good generalization performance. These advantages have been substantiated by research papers that discuss the theoretical underpinnings and practical applications of SVMs in various domains. SVM has been widely employed in various research studies and has consistently exhibited excellent performance [19], [26]. Furthermore, based on several comparative reference studies, these research works utilized SVM in their buzzer classification modeling experiments and yielded satisfactory results, averaging above 80% in accuracy [27]–[29]. Therefore, for our experiment, which aims to explore features that possess distinct characteristics and strong correlation with buzzer accounts, the same classification model is employed.

## Evaluation

In evaluating the effectiveness of the buzzer classification model, standard metrics like accuracy, precision, recall, and F1-score are commonly employed [30]. Accuracy offers a general assessment of the model's performance, while precision, recall, and the F1-score provide more detailed insights into the model's efficacy concerning different classes and types of errors [6]. The adoption of a variety of evaluation metrics is intentional, aiming to prioritize specific metrics over others and consequently providing a more thorough and nuanced analysis of the model's performance. This diverse set of metrics contributes to a comprehensive understanding of the classification model's strengths and areas for improvement.

## RESULT AND DISCUSSION

### Feature Selection Result

The feature selection method utilized is tailored to the data type of each feature. For features with numerical data, Spearman and Pearson correlation feature selection techniques can be employed. On the other hand, for features with nominal data type, the Chi-Square test is utilized. The outcomes of feature selection employing Spearman and Pearson correlation are displayed in Table 2, where features enclosed in parentheses denote absolute values of negative correlations for ranking purposes. The outcomes of Spearman correlation computations on features

Table 2. The Result of Spearman and Pearson Correlation on Numerical Features

No.	Spearman		Pearson	
	Feature	Score	Feature	Score
1.	avg_similarity	0.6408	avg_similarity	0.6174
2.	count_hashtags	0.4519	count_hashtags	0.3136
3.	tweet	0.3641	count_photos	0.2854
4.	count_photos	0.3279	tweet	0.2367
5.	count_urls	0.2670	count_mentions	0.2272
6.	count_mentions	0.2566	count_reply_to	0.2061
7.	count_reply_to	0.2428	TIE_diff_interval	0.2057
8.	all_media	0.2357	(days_created)	(0.2023)
9.	count_retweets	0.2316	all_media	0.1917
10.	(days_created)	(0.2188)	count_urls	0.1801
11.	count_replies	0.1644	all_tweets	0.1378
12.	TIE_diff_interval	0.1571	following	0.1038
13.	count_likes	0.1350	all_likes	0.0794
14.	all_tweets	0.6408	(url_to_tweet)	(0.0468)
15.	(url_to_web)	0.4519	url_to_web	0.0375
16.	followers	0.3641	count_retweets	0.0276
17.	following	0.3279	(count_likes)	(0.023)
18.	all_likes	0.2670	count_replies	0.0135
19.	(url_to_tweet)	0.2566	followers	0.0102



Table 3. The Results of Chi-Square Test on Nominal Features

No.	Feature	Score	Conclusion
1.	bio	0.04225	Dependent
2.	bg_image	0.05152	Not dependent
3.	url	0.36943	Not dependent
4.	profile_image	0.42067	Not dependent
5.	location	0.53755	Not dependent

categorized as numerical data display a similarity with the rankings of the top four features obtained through Pearson correlation computations. These features are “avg\_similarity”, “count\_hashtags”, “tweet”, and “count\_photos”, although there exists a minor disparity in the arrangement of the last two features. As for the subsequent rankings up to the last position, both correlation techniques used exhibit notably different outcomes. In the Spearman correlation results, three features have negative values, representing a negative monotonic relationship. These features are “days\_created”, “url\_to\_web”, and “url\_to\_tweet”. This implies that smaller values of these features correspond to greater correlation with accounts classified as buzzers. Conversely, in the Pearson correlation results, three features exhibit a negative linear correlation with the buzzer label: “days\_created”, “url\_to\_tweet”, and “count\_likes”. The negative relationship refers to the pattern observed between two variables when, as one variable increases, the other variable consistently decreases. It means that for the feature “days\_created”, for example, the higher the value of “days\_created” or the longer an account has been created, the smaller the likelihood that the account will be detected as a buzzer account. Furthermore, these results also indicate that buzzer accounts tend to exhibit the characteristic of having more tweets that either refer to other websites or reference other tweets.

Observing the values obtained from the aforementioned correlation techniques, the highest value is achieved by “avg\_similarity”, with a value exceeding 0.6, which can be considered strongly correlated. Despite “count\_hashtags” being the second-ranking feature in both correlation techniques, its value falls within the category of moderate correlation. Overall, other features generally exhibit weak correlations due to values below 0.3. The high correlation of “avg\_similarity” suggests that buzzer accounts tend to post tweets that are very similar to each other, possibly indicating a lack of originality or diversity in their content. This is similar to what was found in previous research, which also discovered that the similarity among tweets from buzzer accounts resulted in high correlation values [13]. Moreover, these findings align with Panatra's research [17], which defines buzzers as having the characteristic of posting repeatedly more than three times, even if the tweets are posted on different accounts. The moderate correlation of “count\_hashtags” implies that buzzer accounts use more hashtags than non-buzzer accounts, possibly to increase their visibility or relevance to certain topics. The moderate correlation of “tweet” implies that buzzer accounts tend to post more tweets than non-buzzer accounts. This aligns with previous research [13], [17] that defines one characteristic of buzzer accounts as making numerous posts in a day, with a majority of their content being reposts from other accounts. The negative correlation of “days\_created” indicates that buzzer accounts are relatively newer than non-buzzer accounts, possibly reflecting their short-term or disposable nature [17]. The low or negligible correlation of other features suggests that they are not very useful or discriminative for distinguishing between buzzer and non-buzzer accounts.

Table 3 presents the outcomes of feature selection using the Chi-Square Test on features with nominal data type. The results indicate that only one feature exceeds the significance threshold of 0.05, which is the “bio” feature. This implies that for the “bio” feature, the null hypothesis is rejected, leading to the conclusion that the “bio” feature is dependent and has a significant association with the “buzzer” variable. Meanwhile, the remaining four features are considered to lack a substantial association. The dependence of “bio” suggests that buzzer accounts have a different pattern of using bio on their profiles than non-buzzer accounts. This could be related to the purpose or intention of the buzzer accounts, such as promoting a certain product, service, or ideology. The independence of other features suggests that buzzer and non-buzzer accounts have similar tendencies of using background image, url, profile image, and location on their profiles, and these features do not help to differentiate them.

Table 4. Performance of SVM Model

Feature	Spearman				Pearson				Feature	Spearman				Pearson			
	A	P	R	F	A	P	R	F		A	P	R	F	A	P	R	F
19	82.26	73.91	77.27	80.82	82.26	73.91	77.27	80.82	<b>19 + bio</b>	<b>87.10</b>	<b>79.17</b>	<b>86.36</b>	<b>86.18</b>	87.10	79.17	86.36	86.18
18	80.65	70.83	77.27	79.26	82.26	72.00	81.82	81.16	18 + bio	82.26	72.00	81.82	81.16	82.26	73.91	77.27	80.82
17	75.81	62.96	77.27	74.69	87.10	85.00	77.27	85.60	<b>17 + bio</b>	77.42	66.67	72.73	75.81	<b>88.71</b>	<b>89.47</b>	<b>77.27</b>	<b>87.25</b>
16	77.42	64.29	81.82	76.54	80.65	72.73	72.73	78.86	16 + bio	79.03	65.52	86.36	78.35	80.65	72.73	72.73	78.86
15	72.58	59.26	72.73	71.32	77.42	70.00	63.64	74.80	15 + bio	75.81	64.00	72.73	74.30	77.42	70.00	63.64	74.80
14	82.26	70.37	86.36	81.44	83.87	87.50	63.64	81.03	14 + bio	80.65	72.73	72.73	78.86	83.87	87.50	63.64	81.03
<b>13</b>	<b>87.10</b>	<b>79.17</b>	<b>86.36</b>	<b>86.18</b>	83.87	87.50	63.64	81.03	13 + bio	85.48	78.26	81.82	84.30	87.10	88.89	72.73	85.24
12	83.87	73.08	86.36	83.00	74.19	61.54	72.73	72.81	12 + bio	83.87	73.08	86.36	83.00	79.03	68.00	77.27	77.73
<b>11</b>	82.26	72.00	81.82	81.16	<b>88.71</b>	<b>85.71</b>	<b>81.82</b>	<b>87.54</b>	11 + bio	85.48	76.00	86.36	84.58	83.87	75.00	81.82	82.72
10	80.65	69.23	81.82	79.61	79.03	66.67	81.82	78.07	10 + bio	79.03	66.67	81.82	78.07	80.65	72.73	72.73	78.86
9	74.19	63.64	63.64	71.82	79.03	66.67	81.82	78.07	9 + bio	80.65	70.83	77.27	79.26	77.42	64.29	81.82	76.54
8	75.81	62.07	81.82	75.02	79.03	71.43	68.18	76.86	8 + bio	75.81	62.96	77.27	74.69	75.81	66.67	63.64	73.30
7	75.81	66.67	63.64	73.30	77.42	70.00	63.64	74.80	7 + bio	77.42	68.18	68.18	75.34	79.03	71.43	68.18	76.86
6	80.65	72.73	72.73	78.86	77.42	70.00	63.64	74.80	6 + bio	80.65	72.73	72.73	78.86	79.03	71.43	68.18	76.86
5	80.65	72.73	72.73	78.86	80.65	72.73	72.73	78.86	5 + bio	80.65	72.73	72.73	78.86	80.65	72.73	72.73	78.86
4	80.65	72.73	72.73	78.86	80.65	72.73	72.73	78.86	4 + bio	80.65	72.73	72.73	78.86	80.65	72.73	72.73	78.86
3	75.81	66.67	63.64	73.30	77.42	68.18	68.18	75.34	3 + bio	75.81	66.67	63.64	73.30	77.42	68.18	68.18	75.34
2	77.42	70.00	63.64	74.80	77.42	70.00	63.64	74.80	2 + bio	77.42	70.00	63.64	74.80	77.42	70.00	63.64	74.80
1	75.81	66.67	63.64	73.30	75.81	66.67	63.64	73.30	1 + bio	77.42	68.18	68.18	75.34	77.42	68.18	68.18	75.34

### Performance Of Buzzer Classification Model

After calculating the feature relevance level using several feature selection techniques, the next step involves conducting experiments on modeling buzzer classification using the Support Vector Machine (SVM) classification algorithm. The study conducts two sets of experiments. The first set involves modeling using 19 features ranked by Spearman and Pearson correlations. A total of 19 experiments are performed, each time eliminating the lowest-scoring feature until only one feature with the highest value remains. The optimal outcomes were achieved with 13 features based on the Spearman correlation sequence and 11 features based on the Pearson correlation sequence. The model using the 11 features as per the Pearson correlation ranking exhibited slightly better performance with an F1-score of 87.54, compared to the 13 features from the Spearman correlation ranking, which achieved an F1-score of 86.18. The second set of experiments was conducted similarly to the first, but with the addition of one feature, namely “bio”, which was deemed relevant based on the results of the Chi-Square Test. The best outcomes were achieved using the Spearman sequence, resulting in an F1-score of 86.18 (19 numerical features plus “bio” feature), and employing the Pearson sequence produced the optimal result of an F1-score of 87.25 (17 numerical

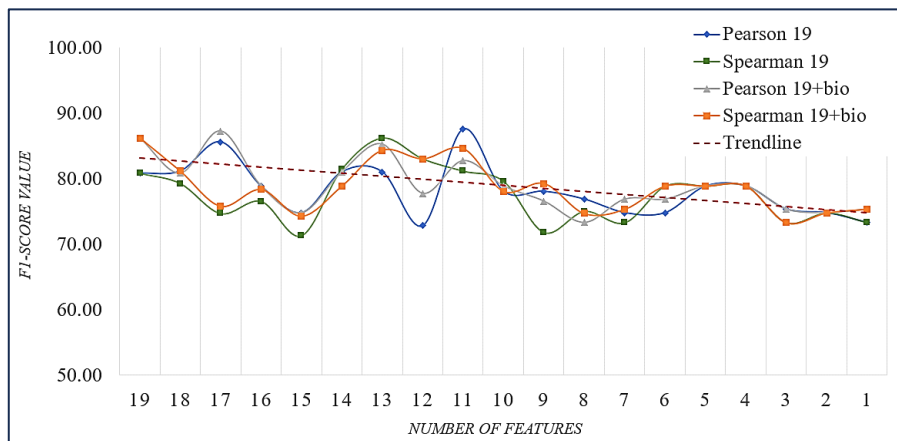


Figure 2. Trendline of SVM Model Performance (F1-Score)

features plus “bio” feature). The performance of modeling using the 19 features and 19 features + “bio” with SVM is presented in Table 4.

The study found that as the number of utilized features diminishes, the performance trend declines as shown in Figure 2. This trend is evident in the evaluation metric values for features equal to or less than 10, where no F1-score surpasses 80. Despite the reduction in the number of features, those used are deemed most relevant to the variable. However, the decreasing trend could be attributed to the fact that the obtained correlation values are not significantly strong. Moreover, in the second experiment involving the addition of one more feature, the results were not superior to the 11 features from Pearson ranking. This indicates that through diverse feature selection methods, a more efficient modeling can be achieved with a reduced number of features. Across all experiments, the best performance is demonstrated in the modeling employing 11 features in accordance with the Pearson correlation ranking, namely “avg\_similarity”, “count\_hashtags”, “count\_photos”, “tweet”, “count\_mentions”, “count\_reply\_to”, “TIE\_diff\_interval”, “days\_created”, “all\_media”, “count\_urls”, and “all\_tweets”.

## CONCLUSION

In this study, an exploration of several feature selection techniques was conducted to identify the most relevant features for enhancing the performance of the buzzer classification model in social media data. Experiments were carried out using feature selections obtained through three types of feature selection techniques: Spearman correlation, Pearson correlation, and Chi-Square test. The first experiment involved 19 numerical features, ranked based on Spearman and Pearson correlations. In the second experiment, an additional nominal feature, deemed relevant according to the Chi-Square test results, was added. The performance of the buzzer classification model using Support Vector Machine yielded the best results when employing 11 features based on Pearson correlation calculations, achieving an F1-score of 87.54. These features are “avg\_similarity”, “count\_hashtags”, “count\_photos”, “tweet”, “count\_mentions”, “count\_reply\_to”, “TIE\_diff\_interval”, “days\_created”, “all\_media”, “count\_urls”, and “all\_tweets”. These results provide insights into features that have correlations and can be utilized as distinguishing factors between buzzer and non-buzzer accounts. For the development of a buzzer classification model, the focus can be directed towards incorporating these features, aiming to achieve more accurate and efficient classification results without the need for an extensive number of features. The results of this study can be utilized by researchers and platform administrators to develop more effective algorithms for identifying and preventing buzzers on social media, thereby preventing the spread of misinformation. Policymakers can leverage these insights to develop more effective regulations and strategies to combat the spread of misinformation, ultimately fostering a healthier public discourse. By addressing these broader impacts, our research underscores its significance not only for the academic community but also for social media platform administrators, developers, policymakers, and society at large.

For future research, it is essential to explore additional feature selection methods and alternative classification algorithms. Specifically, deep learning algorithms like Long Short-Term Memory (LSTM), Convolutional Neural Networks (CNN), and Bidirectional Encoder Representations from Transformers (BERT) should be considered due to their potential to capture complex patterns in sequential and unstructured data. Further investigation into distinctive features that characterize buzzers, particularly paid ones, such as interaction behaviors and linguistic features, is also recommended. Expanding the dataset to include more buzzer and non-buzzer accounts could enhance the model's ability to capture a wider range of buzzer behaviors. By addressing these aspects, future research can build upon our findings to develop more effective and comprehensive buzzer detection systems.

## ACKNOWLEDGMENT

The authors express their sincere gratitude to all individuals and entities who contributed to the realization of this research endeavor. Special appreciation is extended to Mohammad Teduh Uliniansyah, Elvira Nurfadhilah, and Agung Santosa for their invaluable support, discussions, and feedback, which significantly enriched the quality and depth of this study. The authors also wish to acknowledge the diligent efforts of the anonymous reviewers, whose insightful feedback and constructive comments played a crucial role in enhancing the overall quality of the manuscript. Their expertise and dedication are truly appreciated.

## CONFLICT OF INTEREST

The author declares that there are no conflicts of interest regarding the authorship or publication of this research.

## FUNDING

The authors received no financial support for the research, authorship, and/or publication of this article.

## References

- [1] M. Arazzi, M. Ferretti, S. Nicolazzo, and A. Nocera, "The role of social media on the evolution of companies: A Twitter analysis of Streaming Service Providers," *Online Soc Netw Media*, vol. 36, Jul. 2023, doi: [10.1016/j.osnem.2023.100251](https://doi.org/10.1016/j.osnem.2023.100251).
- [2] M. Grandjean, "A social network analysis of Twitter: Mapping the digital humanities community," *Cogent Arts & Humanity*, vol. 3, no. 1, 2016, doi: [10.1080/23311983.2016.1171458](https://doi.org/10.1080/23311983.2016.1171458).
- [3] L. K. Kaye, "Exploring the 'socialness' of social media," *Computers in Human Behavior Reports*, vol. 3. Elsevier Ltd, Jan. 01, 2021. doi: [10.1016/j.chbr.2021.100083](https://doi.org/10.1016/j.chbr.2021.100083).
- [4] D. M. Boyd and N. B. Ellison, "Social network sites: Definition, history, and scholarship," *Journal of Computer-Mediated Communication*, vol. 13, no. 1, pp. 210–230, Oct. 2007, doi: [10.1111/j.1083-6101.2007.00393.x](https://doi.org/10.1111/j.1083-6101.2007.00393.x).
- [5] G. Yavetz and N. Aharony, "Social media in government offices: usage and strategies," *Aslib Journal of Information Management*, vol. 72, no. 4, pp. 445–462, Nov. 2020, doi: [10.1108/AJIM-11-2019-0313](https://doi.org/10.1108/AJIM-11-2019-0313).
- [6] I. Mergel, "Open innovation in the public sector: drivers and barriers for the adoption of Challenge.gov," *Public Management Review*, vol. 20, no. 5, pp. 726–745, 2018, doi: [10.1080/14719037.2017.1320044](https://doi.org/10.1080/14719037.2017.1320044).
- [7] E. Rosenberg et al., "Sentiment analysis on Twitter data towards climate action," *Results in Engineering*, vol. 19, Sep. 2023, doi: [10.1016/j.rineng.2023.101287](https://doi.org/10.1016/j.rineng.2023.101287).
- [8] O. Czeranowska et al., "Migrants vs. stayers in the pandemic – A sentiment analysis of Twitter content," *Telematics and Informatics Reports*, vol. 10, Jun. 2023, doi: [10.1016/j.teler.2023.100059](https://doi.org/10.1016/j.teler.2023.100059).
- [9] L. Ilias and I. Roussaki, "Detecting malicious activity in Twitter using deep learning techniques," *Appl Soft Comput*, vol. 107, Aug. 2021, doi: [10.1016/j.asoc.2021.107360](https://doi.org/10.1016/j.asoc.2021.107360).
- [10] M. T. Juzar and S. Akbar, "Buzzer Detection on Twitter Using Modified Eigenvector Centrality," in *2018 5th International Conference on Data and Software Engineering (ICoDSE)*, 2018, pp. 1–5. doi: [10.1109/ICoDSE.2018.8705788](https://doi.org/10.1109/ICoDSE.2018.8705788).
- [11] M. Ibrahim, O. Abdillah, A. F. Wicaksono, and M. Adriani, "Buzzer Detection and Sentiment Analysis for Predicting Presidential Election Results in a Twitter Nation," in *Proceedings - 15th IEEE International Conference on Data Mining Workshop, ICDMW 2015, Institute of Electrical and Electronics Engineers Inc.*, Jan. 2016, pp. 1348–1353. doi: [10.1109/ICDMW.2015.113](https://doi.org/10.1109/ICDMW.2015.113).
- [12] M. O. Ibrohim and I. Budi, "Hate speech and abusive language detection in Indonesian social media: Progress and challenges," *Heliyon*, vol. 9, no. 8. Elsevier Ltd, Aug. 01, 2023. doi: [10.1016/j.heliyon.2023.e18647](https://doi.org/10.1016/j.heliyon.2023.e18647).
- [13] M. Kantepé and M. C. Ganiz, "Preprocessing framework for Twitter bot detection," in *2017 International Conference on Computer Science and Engineering (UBMK)*, 2017, pp. 630–634. doi: [10.1109/UBMK.2017.8093483](https://doi.org/10.1109/UBMK.2017.8093483).
- [14] S. Wang, J. Tang, and H. Liu, "Feature Selection," in *Encyclopedia of Machine Learning and Data Mining*, Boston, MA: Springer US, 2016, pp. 1–9. doi: [10.1007/978-1-4899-7502-7\\_101-1](https://doi.org/10.1007/978-1-4899-7502-7_101-1).
- [15] H. I. Kuru, A. E. Cicek, and O. Tastan, "From Cell-Lines to Cancer Patients: Personalized Drug Synergy Prediction," 2023, doi: [10.1101/2023.02.13.528276](https://doi.org/10.1101/2023.02.13.528276).
- [16] R. Rodríguez-Pérez and J. Bajorath, "Feature importance correlation from machine learning indicates functional relationships between proteins and similar compound binding characteristics," *Sci Rep*, vol. 11, no. 1, p. 14245, Jul. 2021, doi: [10.1038/s41598-021-93771-y](https://doi.org/10.1038/s41598-021-93771-y).
- [17] A. J. Panatra, F. B. Chandra, W. Darmawan, H. L. H. S. Warnars, W. H. Utomo, and T. Matsuo, "Buzzer Detection to Maintain Information Neutrality in 2019 Indonesia Presidential Election," in *Proceedings - 2019*

- 8th International Congress on Advanced Applied Informatics, IIAI-AAI 2019, Institute of Electrical and Electronics Engineers Inc., Jul. 2019, pp. 873–876. doi: [10.1109/IIAI-AAI.2019.00177](https://doi.org/10.1109/IIAI-AAI.2019.00177).
- [18] A. Suciati, A. Wibisono, and P. Mursanto, “Twitter Buzzer Detection for Indonesian Presidential Election,” in 2019 3rd International Conference on Informatics and Computational Sciences (ICICoS), 2019, pp. 1–5. doi: [10.1109/ICICoS48119.2019.8982529](https://doi.org/10.1109/ICICoS48119.2019.8982529).
- [19] S. Pebiana et al., “Experimentation of Various Preprocessing Pipelines for Sentiment Analysis on Twitter Data about New Indonesia’s Capital City Using SVM and CNN,” in 2022 25th Conference of the Oriental COCOSDA International Committee for the Co-Ordination and Standardisation of Speech Databases and Assessment Techniques, O-COCOSDA 2022 - Proceedings, Institute of Electrical and Electronics Engineers Inc., 2022. doi: [10.1109/O-COCOSDA202257103.2022.9997982](https://doi.org/10.1109/O-COCOSDA202257103.2022.9997982).
- [20] M. Lobo and R. D. Guntur, “Spearman’s rank correlation analysis on public perception toward health partnership projects between Indonesia and Australia in East Nusa Tenggara Province,” J Phys Conf Ser, vol. 1116, p. 022020, Dec. 2018, doi: [10.1088/1742-6596/1116/2/022020](https://doi.org/10.1088/1742-6596/1116/2/022020).
- [21] P. Sedgwick, “Spearman’s rank correlation coefficient,” BMJ, p. g7327, Nov. 2014, doi: [10.1136/bmj.g7327](https://doi.org/10.1136/bmj.g7327).
- [22] F. Zinzendoff Okwonu, B. Laro Asaju, and F. Irimisose Arunaye, “Breakdown Analysis of Pearson Correlation Coefficient and Robust Correlation Methods,” IOP Conf Ser Mater Sci Eng, vol. 917, no. 1, p. 012065, Sep. 2020, doi: [10.1088/1757-899X/917/1/012065](https://doi.org/10.1088/1757-899X/917/1/012065).
- [23] P. Schober, C. Boer, and L. A. Schwarte, “Correlation Coefficients: Appropriate Use and Interpretation,” Anesth Analg, vol. 126, no. 5, pp. 1763–1768, May 2018, doi: [10.1213/ANE.0000000000002864](https://doi.org/10.1213/ANE.0000000000002864).
- [24] S. T. Nihan, “Karl Pearsons chi-square tests,” Educational Research and Reviews, vol. 15, no. 9, pp. 575–580, Sep. 2020, doi: [10.5897/ERR2019.3817](https://doi.org/10.5897/ERR2019.3817).
- [25] R. Singhal and R. Rana, “Chi-square test and its application in hypothesis testing,” Journal of the Practice of Cardiovascular Sciences, vol. 1, no. 1, p. 69, 2015, doi: [10.4103/2395-5414.157577](https://doi.org/10.4103/2395-5414.157577).
- [26] M. A. Hearst, S. T. Dumais, E. Osuna, J. Platt, and B. Scholkopf, “Support vector machines,” IEEE Intelligent Systems and their Applications, vol. 13, no. 4, pp. 18–28, 1998, doi: [10.1109/5254.708428](https://doi.org/10.1109/5254.708428).
- [27] M. Kantepe and M. C. Ganiz, “Preprocessing framework for Twitter bot detection,” in 2017 International Conference on Computer Science and Engineering (UBMK), IEEE, Oct. 2017, pp. 630–634. doi: [10.1109/UBMK.2017.8093483](https://doi.org/10.1109/UBMK.2017.8093483).
- [28] A. J. Panatra, F. B. Chandra, W. Darmawan, H. L. H. S. Warnars, W. H. Utomo, and T. Matsuo, “Buzzer Detection to Maintain Information Neutrality in 2019 Indonesia Presidential Election,” in 2019 8th International Congress on Advanced Applied Informatics (IIAI-AAI), IEEE, Jul. 2019, pp. 873–876. doi: [10.1109/IIAI-AAI.2019.00177](https://doi.org/10.1109/IIAI-AAI.2019.00177).
- [29] M. Ibrahim, O. Abdillah, A. F. Wicaksono, and M. Adriani, “Buzzer Detection and Sentiment Analysis for Predicting Presidential Election Results in a Twitter Nation,” in 2015 IEEE International Conference on Data Mining Workshop (ICDMW), IEEE, Nov. 2015, pp. 1348–1353. doi: [10.1109/ICDMW.2015.113](https://doi.org/10.1109/ICDMW.2015.113).
- [30] C. Goutte and E. Gaussier, “A Probabilistic Interpretation of Precision, Recall and F-Score, with Implication for Evaluation,” 2005, pp. 345–359. doi: [10.1007/978-3-540-31865-1\\_25](https://doi.org/10.1007/978-3-540-31865-1_25).

## AUTHORS BIOGRAPHY

**Dian Isnaeni Nurul Afra** is a researcher at the Research Center for Data and Information Sciences, National Research and Innovation Agency (BRIN), Indonesia. She obtained her Bachelor's degree in Informatics Engineering from Brawijaya University in 2015 and her Master's degree in Information Technology from the University of Indonesia in 2022. Her current research focuses extensively on natural language processing, including sentiment analysis, text processing, text classification, and machine learning. Additionally, she is actively involved in exploring various innovative applications of artificial intelligence, contributing significantly to advancements in the field.

**Radhiyatul Fajri** received her Master of Computer Science from IPB University, Indonesia, in 2022. Currently, she is an associate researcher at the Research Center for Data and Information Sciences, affiliated with the National Research and Innovation Agency (BRIN) in Indonesia. Since 2021, she has been actively engaged in research on

artificial intelligence, with a specific focus on natural language processing and computer vision. Her research includes sentiment analysis, natural language processing, and face recognition.

**Harnum Annisa Prafitia** is an engineering staff at the Center for Information and Communication Technology, Agency for the Assessment and Application of Technology from 2011 to 2021. Since 2022, she has been working at the Research Center for Data and Information Sciences, National Research and Innovation Agency (BRIN). She received her Bachelor's degree in Statistics from ITS Surabaya in 2010 and is currently pursuing a Master's degree in Information Systems Management at Bina Nusantara University. Her expertise includes Information Systems, ICT Management, and SPBE Audit.

**Ikhwan Arief**, a permanent staff member in the Department of Industrial Engineering at Andalas University, holds a Master's Degree from the University of Birmingham, England. He serves as the DOAJ Editor and Ambassador for Indonesia, is a member of the editorial team for the Scopus-indexed Journal of Jurnal Optimasi Sistem Industri, and coordinates the enhancement of scientific journals under the Institute of Industrial Engineering Higher Education Cooperation Agency (BKSTI) in Indonesia. His dedication to academic excellence is clear, particularly in his focus on Data Engineering, Business Intelligence, and Data Analysis, through which he makes significant scholarly contributions.

**Aprinaldi Jasa Mantau** received his Bachelor and Master of Computer Science degrees from the Faculty of Computer Science at the University of Indonesia in 2013 and 2014, respectively. Currently, he is a doctoral student in the Department of Computer Science and Systems Engineering at Kyushu Institute of Technology. His research interests encompass machine learning, computer vision, robotics, data mining, swarm intelligence, and swarm robotics. He is also a member of The Institute of Electrical and Electronics Engineers (IEEE).