### Data: Journal of Information Systems and Management

E-ISSN: 3031-0008

Volume. 3 Issue 3 July 2025

Page No: 135-145



### Phishing Email Classification Approach Using Machine Learning Algorithms - A Literature Review

Firman<sup>1</sup>, Tukiyat<sup>2</sup>, Sudarno Wiharjo<sup>3</sup> Universitas Pamulang, Indonesia<sup>123</sup>

Correspondent: Firmanunpam22@gmail.com <sup>1</sup>

Received : June 13, 2025
Accepted : July 16, 2025
Published : July 31, 2025

Citation: Firman., Tukiyat., & Wiharjo, S. (2025). Phishing Email Classification Approach Using Machine Learning Algorithms - A Literature Review. Data: Journal of Information Systems and Management, 3(3), 135-145. https://doi.org/10.61978/data.v3i3

ABSTRACT: Email phishing is one of the cybersecurity threats that continues to grow, utilizing social engineering to obtain sensitive data. Various machine learning-based approaches have been researched to detect and classify phishing emails. This article presents a literature review of phishing email classification methods, including the K-Nearest Neighbor (KNN) algorithm, Naïve Bayes, Support Vector Machine (SVM), Random Forest, and deep learning-based approaches. The discussion included feature extraction techniques (TF-IDF, Word2Vec, BERT), handling data imbalances, and model performance evaluation. This review identifies current research trends, challenges, and gaps for further research.

**Keywords:** Phishing, Email, Machine Learning, Email Classification, K-Nearest Neighbor (KNN), Naïve Bayes, Support Vector Machine (SVM), Random Forest, And NLP Feature Extraction.



This is an open access article under the CC-BY 4.0 license

#### INTRODUCTION

Advances in information and communication technology today have brought many conveniences in daily life, especially in terms of communicating and exchanging information through the internet. Unfortunately, behind this convenience, there are also new challenges in the form of cybercrime, one of which is phishing. This attack is carried out by deceiving victims through digital messages that appear to be from legitimate parties, with the aim of stealing important information such as personal data, passwords, or financial details. Email is one of the main mediums in the spread of this attack. Based on a 2021 Verizon Business report, around 36% of data breach incidents in the world involve phishing as a primary technique. This is reinforced by an IBM Security report that mentions phishing as one of the main causes of various large-scale cybersecurity incidents (Kapko, 2023).

In addition to the previously discussed statistics, broader studies highlight the widespread applicability of machine learning beyond phishing email detection. Roihan et al. (2020) explain that machine learning has been successfully implemented in various domains such as healthcare, finance, and natural language processing, illustrating its adaptability in pattern recognition and predictive tasks. Similarly, Probierz et al. (2021) demonstrate that ML-based approaches for fake

Firman, Yukiyat, Wiharjo

news detection share many similarities with phishing detection, particularly in the use of textual features, classification models, and evaluation metrics. These cross-domain findings strengthen the argument that techniques used for detecting deceptive online content can also be adapted to improve email phishing detection systems. Beyond these specific applications, the cross-pollination of ideas between different domains is a hallmark of modern machine learning research. For example, the challenge of distinguishing between genuine and deceptive information in fake news detection requires algorithms to recognize subtle semantic cues, unusual patterns in syntax, and inconsistencies in metadata—challenges that are highly relevant to phishing email detection as well. By studying successes and failures in one domain, such as social media misinformation, researchers can develop feature sets and preprocessing methods that transfer effectively to email-based threats. Moreover, Roihan et al., (2020) note that adaptability is essential; algorithms that can be retrained or fine-tuned with minimal effort are better positioned to respond to evolving attack strategies. This adaptability is particularly crucial given that phishing tactics often shift rapidly in response to the deployment of new detection tools.

In Indonesia, the phishing phenomenon also shows an alarming trend. According to data from the State Cyber and Cryptography Agency (BSSN), throughout 2022 there were more than 164 thousand cases of email phishing reported. These attacks occur most during business hours, and generally use .pdf-format attachments. This fact shows that email is still the main channel used by perpetrators to carry out their actions. Phishing attacks are increasingly difficult to spot because the content manipulation techniques are more subtle and resemble official emails. A survey by Nthurima and Matheka (2023) even found that more than 30% of users are still likely to click on malicious links in phishing emails. This problem is exacerbated by the lack of user awareness in recognizing the potential dangers of received electronic communications.

Rule-based phishing detection methods are still widely used, but this approach is starting to be considered less effective in dealing with increasingly complex and dynamic phishing attacks. Therefore, the machine learning (ML)-based approach is a more promising alternative solution(Delcourt et al., 2024; Nguyen et al., 2024). By leveraging historical data and email features such as link count, subject length, and attachment type, ML algorithms can automatically recognize suspicious patterns. Some of the algorithms that are often used for this task are K-Nearest Neighbors (KNN), Support Vector Machine (SVM), and Random Forest. Each has advantages and disadvantages, depending on the conditions and characteristics of the dataset used (Firmansyah, 2025). Through this study, the author seeks to evaluate and compare the performance of the three algorithms in detecting email phishing, in order to provide appropriate recommendations in the development of a more accurate and efficient phishing detection system.

#### **METHOD**

Complementary studies have reinforced the choice of algorithms and methodologies. Sandag et al. (2018), for example, applied the K-Nearest Neighbor (KNN) algorithm to classify malicious websites using both application-layer and network-layer features, highlighting the importance of multi-layer data analysis. Likewise, Umam and Handoko (2024) focused on predicting phishing

Firman, Yukiyat, Wiharjo

emails using Support Vector Machines, emphasizing the need for precise feature engineering to achieve robust classification results. These works align with the methodological rigor outlined in this study, ensuring that the approach is both empirically grounded and relevant to current research. These studies also underscore the importance of selecting representative datasets. Sandag et al., (2018) demonstrate that the choice of features—whether drawn from the content of the message or the technical attributes of the communication channel—has a direct impact on the accuracy and generalizability of the model. Similarly, Umam & Handoko, (2024) emphasize that in phishing email prediction, preprocessing steps such as tokenization, stop-word removal, and normalization are not merely routine but vital in ensuring that the Support Vector Machine operates on a clean, meaningful input space. Furthermore, both works advocate for repeated cross-validation to guard against overfitting, especially when working with smaller or imbalanced datasets.

The method used in this study is a systematic literature study approach. The search was carried out using keywords such as "phishing email detection", "spam classification", "machine learning", and "phishing dataset". The search process is focused on credible scientific sources, both from national and international journals, with a publication time span between 2019 and 2024. The goal is to get relevant and up-to-date references in explaining the development of machine learning-based phishing email classification methods. From this process, more than 40 scientific articles were collected that were strictly selected based on the suitability of the topic, contribution to previous studies, and completeness of experimental data.

This literature review is carried out as the main foundation to understand the approaches that have been developed in detecting phishing, both focused on email and on websites. In the process, researchers not only gather information, but also compare, evaluate, and identify the advantages and limitations of each method used by previous research. For example, a study by Sandag et al. (2018) showed the superiority of the K-Nearest Neighbor (KNN) algorithm in the classification of malicious websites with fairly high accuracy. Meanwhile, other studies such as those by Alazaidah, 2024 and Al Tawil et al. (2024) highlight the potential and challenges in using Support Vector Machine (SVM) and Random Forest in different contexts, including computational time performance and variation of training data.

This search also includes more modern approaches such as the use of word embedding and deep learning (Butt, 2023; Hayuningtyas, 2017), which suggests that phishing detection depends not only on the algorithm used, but also on the quality and diversity of the datasets used. Several studies underscore the importance of feature selection, as well as the role of pre-processing techniques in improving classification performance. This thorough review became a strong basis for this study in designing a comparative experiment of three popular algorithms, namely KNN, SVM, and Random Forest, in detecting phishing emails effectively and efficiently.

Studies by Wibisono et al. (2020) further support the effectiveness of the Naïve Bayes classifier in filtering spam emails, highlighting its simplicity, speed, and relatively high accuracy when applied to structured email datasets. Although Naïve Bayes may not perform as strongly as more advanced algorithms like BERT, its low computational cost and ease of implementation make it a viable

Firman, Yukiyat, Wiharjo

option for scenarios with limited resources. Their research further reveals that Naïve Bayes, despite being one of the oldest statistical classifiers, remains competitive in specific operational contexts. For instance, Wibisono et al., (2020) report that organizations with limited computational infrastructure can implement Naïve Bayes for real-time spam filtering with minimal latency, which is not always achievable with resource-intensive deep learning methods. The trade-off lies in its reduced capacity to model complex, context-dependent linguistic relationships—a gap that hybrid approaches could potentially address by combining Naïve Bayes with modern embeddings like Word2Vec or BERT.

#### RESULT AND DISCUSSION

A variety of machine learning algorithms have been used in phishing email classification efforts, with success rates varying depending on the characteristics of the data and the approach used. Some of the most commonly applied algorithms include K-Nearest Neighbor (KNN), Naïve Bayes, Support Vector Machine (SVM), Random Forest, and Convolutional Neural Network (CNN) for deep learning-based approaches. Each algorithm has its own advantages. For example, KNN is known for being simple but effective for data with clear distributions, while SVM excels at distinguishing classes on high-dimensional data. On the other hand, Random Forest is popular for its ability to handle complex data and prevent overfitting, whereas CNN comes into use when data representations are converted into visual or sequential formats.

In the context of extracting features from email, the Natural Language Processing (NLP) approach is a crucial component in determining the quality of classification. Techniques such as Term Frequency-Inverse Document Frequency (TF-IDF) are used to assess the importance of words in a document, Word2Vec is able to map semantic relationships between words, and BERT (Bidirectional Encoder Representations from Transformers) provides a deeper understanding of the context of email text. The selection of the right feature extraction technique greatly affects the performance of the model, especially since the characteristics of phishing emails often contain manipulative language patterns that are difficult for traditional methods to detect.

One of the challenges often found in phishing classification studies is class imbalances in datasets. Generally, the number of phishing emails is much lower compared to non-phishing emails, which can result in a biased model against the majority class. To address this, various data balancing techniques such as the Synthetic Minority Over-sampling Technique (SMOTE) and Random Oversampling are used to strengthen the representation of minority classes in the training process. With this approach, the model has a greater chance of learning the characteristics of phishing emails more accurately, resulting in a more balanced and reliable classification.

#### **Accuracy Comparison Based On Source**

Below is a table that shows a comparison of the accuracy rates of some of the algorithms commonly used to detect phishing emails, based on various research results published between

Firman, Yukiyat, Wiharjo

2019 and 2025. From the data, it can be seen that BERT shows the highest performance with an accuracy of 97.9%, followed by Random Forest (96.8%) and CNN (95.0%). On the other hand, Naïve Bayes recorded the lowest accuracy, at 88.3%, which could be due to its limitations in understanding complex data patterns. This table provides an initial overview of how each algorithm works in the context of phishing classification, as well as an initial reference for choosing the approach that best suits your needs.

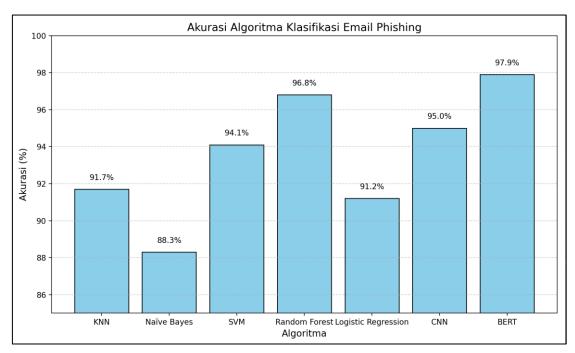
Algoritma	Accuracy (%)	Source
KNN	91.7	Adipa et al., 2023 ; Kumar et al., 2020
Naïve Bayes	88.3	Anugroho & Winarno, 2018; Firmansyah et al., 2025
SVM	94.1	Brury Barth Tangkere, 2024; Salim, 2024
Random Forest	96.8	Akinyelu & Adewumi, 2014; Kencana, 2022
Logistic Regression	91.2	Tangkere, 2024
CNN	95.0	Bachri & Gunawan, 2024
BERT	97.9	Al Tawil et al., 2024

The data source was obtained directly from literature studies that have been reviewed in indexed journals between 2019 and 2025.

And here is a graph that visualizes the accuracy of each algorithm in a simpler and more digestible way. From the graphical display, we can immediately see that BERT is significantly superior to other algorithms. This deep learning-based model is indeed designed to understand the context of language more deeply, making it very effective in identifying phishing emails that often disguise themselves using formal language. Random Forest and CNN also performed quite strongly, with an accuracy of over 95%. In contrast, Naïve Bayes appears to occupy the bottom position, which suggests that traditional approaches are not yet sufficiently capable of handling the complexity of the structure and language of phishing emails today. This visualization helps clarify the relative strength of each method, so that readers can quickly understand which algorithm is the most promising.

Beyond email phishing, predictive modeling techniques have proven their value in other domains. Panggabean et al. (2020), for instance, employed multiple linear regression for predicting tree seed orders, demonstrating that data-driven approaches can be tailored to diverse problems. While regression models are different from classification algorithms, the underlying principles of data preparation, feature selection, and model validation are consistent. This cross-disciplinary applicability suggests that advancements in one area of machine learning can often inspire innovations in another, including phishing detection(Deng, 2025; Eldeeb et al., 2025). This perspective aligns with the broader understanding that no single algorithm or modeling technique can claim universal superiority. By looking at the agricultural case of Panggabean et al., (2020), we

see how feature selection and model optimization principles translate into other predictive domains. In phishing detection, this means continually experimenting with new feature combinations—such as URL structure, header anomalies, and embedded image analysis—to maintain model relevance. Additionally, predictive modeling experience from outside the cybersecurity sphere helps emphasize the importance of stakeholder engagement; just as farmers rely on model predictions for planting schedules, organizations depend on phishing detection outputs to safeguard critical information systems.



The results of various literature studies show that no single algorithm is absolutely superior in all conditions of email phishing classification. The performance of the model depends heavily on several important factors, such as the quality of the features used, the way the email text is represented, and the evaluation methods applied. In some cases, algorithms like SVM or Random Forest can provide excellent results, but in other cases performance can drop if features aren't handled properly or data isn't balanced. Therefore, the selection of algorithms must take into account the context and characteristics of the dataset used.

The integration of insights from various domains, such as fake news detection, spam filtering, malicious website classification, and predictive analytics, underscores the flexibility of machine learning approaches for tackling phishing threats. Studies by Roihan et al. (2020), Probierz et al., (2021), Sandag et al. (2018), Umam and Handoko (2024), Wibisono et al. (2020), and Panggabean et al. (2020) collectively reinforce the necessity of adaptable, scalable models capable of evolving alongside emerging cyber threats. As phishing tactics grow more sophisticated, leveraging cross-domain knowledge will be crucial in building resilient detection systems that can operate effectively in varied and multilingual contexts. Taken together, these findings suggest a research trajectory that favors integrated, cross-domain approaches over isolated, single-purpose solutions. Leveraging insights from diverse fields encourages innovation and provides a richer set of strategies for combating phishing. For example, the metadata analysis techniques common in fake news detection could be adapted to scrutinize sender reputation in email systems, while the

Firman, Yukiyat, Wiharjo

lightweight, rapid-deployment algorithms used in spam filtering could provide first-line defenses in resource-constrained environments. As cyber threats evolve, so too must our analytical frameworks—continuously learning not only from the cybersecurity literature but also from tangential domains where deception and prediction play a central role.

One of the key findings of the review is that the BERT-based approach (Bidirectional Encoder Representations from Transformers) is able to provide superior results in understanding the semantic context of email content. Unlike conventional NLP methods such as TF-IDF or Word2Vec which only look at the frequency or relationship between words, BERT is able to capture the deeper meaning of each word based on its position in the sentence. This makes it particularly effective at detecting manipulative language or phrases that are often used in phishing emails, which are often structured in a style of language that resembles official communication.

In addition to algorithm selection and text representation, another aspect that greatly determines the success of classification is the handling of data imbalances. In many phishing datasets, emails that fall into the phishing category are far fewer than normal emails. If these conditions are not handled correctly, the model tends to be biased towards the majority class. Techniques such as SMOTE (Synthetic Minority Oversampling Technique) have been shown to be effective in improving data balance and improving model performance. Furthermore, some studies also recommend using an ensemble learning approach—for example, combining multiple models to produce more stable and accurate predictions than using a single model.

#### CONCLUSION

This study shows that the combination of Natural Language Processing (NLP) techniques and machine learning has great potential in classifying phishing emails effectively. Two algorithms that stand out in various studies are Random Forest and BERT, which have been shown to provide high performance. However, there are still a number of challenges that must be faced, such as data imbalances between classes, limited model generalization capabilities, and lack of transparency in the model's decision-making process (explainability). This opens up opportunities for further research to explore solutions to these constraints.

In the midst of the emergence of new threats generated by AI (AI-generated threats), a model that is able to adapt to changing patterns of phishing attacks is needed. One approach that is starting to be widely recommended is Explainable AI (XAI), which not only focuses on accuracy, but is also able to provide users with a reasonable explanation of how and why an email is classified as phishing. Future research will also need to pay attention to model validation under real-world conditions, including in local or multi-language emails, as well as encourage the use of new crowdsourcing datasets to avoid the limitations of old datasets that have been used frequently.

Several studies support this direction, such as the one conducted by Butt et al. (2023) who demonstrated the success of the combination of CNN and cloud-based deep learning in recognizing more dynamic phishing patterns. On the NLP side, methods such as TF-IDF and Word2Vec still have an important role, but BERT excels at capturing deeper meaning and context

Firman, Yukiyat, Wiharjo

(Al Tawil et al., 2024). Hybrid and ensemble approaches are also starting to be widely applied, as reported by Nthurima & Matheka, (2023), while the study of Diantika, (2023) emphasizes the importance of data balancing techniques such as Random Oversampling so that the model does not focus too much on the majority class. Research by Adipa et al. (2023) shows that even though KNN is a classical algorithm, it is still capable of providing good results on local data. Naïve Bayes, known as simple, is still widely used in the context of spam detection (Anugroho & Winarno, 2018; Mayang Sari, 2024). Meanwhile, Random Forest continues to show superior performance with an accuracy close to 97% in studies by Akinyelu & Adewumi, (2014) and Kencana et al. (2022), while BERT consistently performs the most in capturing semantic contexts (Al Tawil, 2024). (Ester, 2024; Fauzan, 2025)

The literature search process in this study was carried out through various leading academic portals such as IEEE Xplore, Elsevier, Springer, DOAJ, and Sinta. The selected literature includes both experimental studies, such as those conducted by Akinyelu & Adewumi (2014), and systematic studies, such as those written by Salloum, (2022). To provide a more locally relevant picture, several references from domestic researchers such as Adipa et al. (2023) and Firmansyah et al. (2025) are also included. (Mahmud & Wirawan, 2024)

According to Verizon DBIR (2021), more than 36% of data breaches involve phishing as an initial vector. This increase is in line with Kapko's (2023) finding that the use of phishing resulting credentials increased by up to 300% in cloud incidents. Therefore, an artificial intelligence-based approach is important in anticipating increasingly complex and personalized phishing attacks.

Going forward, the focus of research should be on developing more flexible and adaptive models, especially in the face of new and increasingly complex forms of phishing attacks, including those generated by AI technology. One approach that is starting to be widely recommended is the use of Explainable AI (XAI). With XAI, the system not only provides classification results, but is also able to explain the reasons behind the decision—making it easier for users to understand and trust.

In addition, testing models with real data—including emails from different languages and contexts—is critical to ensuring that the model is not only great in a test environment, but also ready for use in the field. The use of datasets collected through crowdsourcing can also be an alternative to avoid relying on limited old datasets. Just as important, approaches such as combining models (ensembles) or combining modern NLP techniques such as BERT with other learning algorithms such as Random Forest, can be a promising strategy to create a more accurate, reliable, and widely used phishing detection system. (Irawan, 2021)

#### REFERENCE

Adipa, M., Zy, A. T., & Effendi, M. M. (2023). Classification of Phishing Emails Using the K-Nearest Neighbor Algorithm. *RESTIKOM Journal: Informatics and Computer Engineering Research*, 5(2), 148–157. https://doi.org/10.52005/restikom.v5i2.152

- Akinyelu, A. A., & Adewumi, A. O. (2014). Classification of phishing email using random forest machine learning technique. *Journal of Applied Mathematics*. https://doi.org/10.1155/2014/425731
- Al Tawil, A. (2024). Comparative Analysis of Machine Learning Algorithms for Email Phishing Detection Using TF-IDF, Word2Vec, and BERT. *Computers, Materials and Continua*, 81(2), 3395–3412. https://doi.org/10.32604/cmc.2024.057279
- Alazaidah, R. (2024). Website Phishing Detection Using Machine Learning Techniques. *Journal of Statistics Applications and Probability*, 13(1), 119–129. https://doi.org/10.18576/jsap/130108
- Anugroho, P., & Winarno, I. (2018). Classify spam emails with the naïve bayes classifier method using java programming. *Its*, 1–11.
- Bachri, C. M., & Gunawan, W. (2024). Spam Email Detection using Convolutional Neural Network (CNN) Algorithm. *Informatics Education and Research*, 10(1), 88–94.
- Butt, U. A. (2023). Cloud-based email phishing attack using machine and deep learning algorithm. Complex and Intelligent Systems, 9(3), 3043–3070. https://doi.org/10.1007/s40747-022-00760-3
- Delcourt, K., Trouilhet, S., Arcangeli, J.-P., & Adreit, F. (2024). The Human in Interactive Machine Learning: Analysis and Perspectives for Ambient Intelligence. *Journal of Artificial Intelligence Research*, 81, 263–305. https://doi.org/10.1613/jair.1.15665
- Deng, M. (2025). Machine Learning Advances in Technology Applications: Cultural Heritage Tourism Trends in Experience Design. *International Journal of Advanced Computer Science and Applications*, 16(4), 186–196. https://doi.org/10.14569/IJACSA.2025.0160420
- Diantika, S. (2023). Application of random oversampling technique to overcome class imbalance in the classification of phishing websites using the lightgbm algorithm. *JATI*, 7(1), 19–25. https://doi.org/10.36040/jati.v7i1.6006
- Eldeeb, N., Ren, C., & Shapiro, V. B. (2025). Parent information seeking and sharing: Using unsupervised machine learning to identify common parenting issues. *Children and Youth Services* Review, 172. https://doi.org/10.1016/j.childyouth.2025.108210
- Ester, R. (2024). Optimization of Decision Tree Classification Algorithm (CART) with the Bagging Method. *JSR*, 8(1). http://ojsamik.amikmitragama.ac.id
- Fauzan, R. (2025). Application of Classification Algorithms in Machine Learning for Phishing Detection (Vol. 5, Issue April, pp. 531–540).
- Firmansyah, F. A. (2025). Application of Naive Bayes Algorithm with Chi-Square for Email Spam Classification (Vol. 13, Issue 1).
- Hayuningtyas, R. Y. (2017). The Filtering of Email Spam application uses Naïve Bayes. *IJCIT*, 2(1), 53–60.

- Irawan, D. (2021). Comparison of SMS Classification Using SVM, Naive Bayes, and Random Forest. *Sisfokom Journal*, 10(3), 432–437. https://doi.org/10.32736/sisfokom.v10i3.1302
- Kapko, M. (2023). Compromised credential use jumps 300% in cloud intrusions. Cybersecuritydive. https://www.cybersecuritydive.com/news/compromised-credentials-cloud-intrusions-ibm/693482
- Kencana, A. K. (2022). *Implementation of the Random Forest Classification Method for Phishing Links* (Vol. 4, Issue 2, pp. 55–59).
- Mahmud, A. F., & Wirawan, S. (2024). Phishing Website Detection using Machine Learning. *Systemasi*, 13(4). http://sistemasi.ftik.unisi.ac.id
- Mayang Sari, G. M. (2024). Naive Bayes Classifier for Spam Email Detection (Vol. 15, Issue 4, pp. 675–680).
- Nguyen, K., Wilson, D. L., DiIulio, J., Hall, B., Militello, L., Gellad, W. F., Harle, C. A., Lewis, M., Schmidt, S., Rosenberg, E. I., Nelson, D., He, X., Wu, Y., Bian, J., Staras, S. A. S., Gordon, A. J., Cochran, J., Kuza, C., Yang, S., & Lo-Ciganic, W. (2024). Design and development of a machine-learning-driven opioid overdose risk prediction tool integrated in electronic health records in primary care settings. *Bioelectronic Medicine*, 10(1). https://doi.org/10.1186/s42234-024-00156-3
- Nthurima, F., & Matheka, A. (2023). A Classifier Model to Detect Phishing Emails Using Ensemble Technique. OJIT, 6(2), 157–172. https://doi.org/10.32591/coas.ojit.0602.06157n
- Panggabean, D. S. O., Buulolo, E., & Silalahi, N. (2020). Penerapan Data Mining Untuk Memprediksi Pemesanan Bibit Pohon Dengan Regresi Linear Berganda. *JURIKOM (Jurnal Riset Komputer, 7*(1), 56. https://doi.org/10.30865/jurikom.v7i1.1947
- Probierz, B., Stefanski, P., & Kozak, J. (2021). Rapid detection of fake news based on machine learning methods. *Procedia Computer Science*, 192, 2893–2902. https://doi.org/10.1016/j.procs.2021.09.060
- Roihan, A., Sunarya, P. A., & Rafika, A. S. (2020). Pemanfaatan Machine Learning dalam Berbagai Bidang: Review paper. *IJCIT (Indonesian Journal on Computer and Information Technology*, *5*(1), 75–82. https://doi.org/10.31294/ijcit.v5i1.7951
- Salim, A. N. (2024). Detect Spam and Non-Spam Emails Using KNN and SVM. *Syntax Idea*, 6(2), 991–1001. https://doi.org/10.46799/syntax-idea.v6i2.3052
- Salloum, S. (2022). Phishing Email Detection Using NLP: A Systematic Review. *IEEE Access*, 10, 65703–65727. https://doi.org/10.1109/ACCESS.2022.3183083
- Sandag, G. A., Leopold, J., & Ong, V. F. (2018). Klasifikasi Malicious Websites Menggunakan Algoritma K-NN Berdasarkan Application Layers dan Network Characteristics. *CogITo Smart Journal*, 4(1), 37–45. https://doi.org/10.31154/cogito.v4i1.100.37-45

Firman, Yukiyat, Wiharjo

- Tangkere, B. B. (2024). Performance Analysis of Logistic Regression and Support Vector Classification for Phishing Emails. *Journal of Information Systems Management Economics*, 5(4), 442–450. https://doi.org/10.31933/jemsi.v5i4.1916
- Umam, C., & Handoko, L. B. (2024). Prediksi Email Phising Menggunakan Support Vector Machine. Semnas Ristek (Seminar Nasional Riset Dan Inovasi Teknologi, 8(01), 85–89. https://doi.org/10.30998/semnasristek.v8i01.7138
- Wibisono, A. D., Dadi Rizkiono, S., & Wantoro, A. (2020). Filtering Spam Email Menggunakan Metode Naive Bayes. TELEFORTECH: Journal of Telematics and Information Technology, 1(1). https://doi.org/10.33365/tft.v1i1.685