
Application of the Naive Bayes Method for Determining the Quality of Crude Palm Oil (CPO)

Dimas Raka Prananda^{1)*}, Mhd Furqan²⁾

^{1*,2)} Departement of Computer Science, Faculty Science and Technology, Universitas Islam Negeri Sumatera Utara, Indonesia.

^{1*}dimasrkprananda@gmail.com, ²⁾mfurqan@uinsu.ac.id

ABSTRACT

The palm oil industry is a vital pillar of Indonesia's economy, with Crude Palm Oil (CPO) as one of its leading commodities. The quality of CPO significantly impacts its competitiveness and market price internationally. As a prominent CPO processing company, faces challenges in consistently maintaining product quality. Key factors affecting CPO quality include moisture content, free fatty acids, and impurity levels, which are difficult to manage manually. To address these challenges, this study applies the Naive Bayes method as an efficient and fast classification tool for determining CPO quality. Naive Bayes was chosen for its simplicity in probability calculations and its ability to handle data classification with reasonable accuracy. The data used in this study include moisture content, free fatty acids, and impurity levels measured between February and June 2024. The data was split into training data (80%) and testing data (20%) and analyzed using RapidMiner software. The results show that the Naive Bayes method achieved an accuracy rate of 66.6%, with precision and recall values of 50% each. Although the accuracy could be improved, the application of this method has significantly enhanced the efficiency of determining CPO quality. Thus, the implementation of the Naive Bayes method in determining CPO is an effective step towards improving operational efficiency, classification accuracy, and decision-making quality related to product standards, ultimately supporting the company's competitiveness in the global market.

Keywords: Crude Palm Oil (CPO), Naive Bayes, Quality Classification, Data Mining, RapidMiner

INTRODUCTION

The palm oil industry has become one of the most important economic sectors in Indonesia is the largest producer of palm oil in the world, making it one of the leading producers of vegetable oil. One of the main commodities in this industry is crude palm oil (CPO). Essentially, the quality of CPO greatly affects its market value in the global market. Therefore, determining the quality of CPO is a crucial aspect of palm oil plantation operations (Efendi et al., 2023).

PT Perkebunan Nusantara (PTPN) 2 is one of the leading palm oil plantation companies in Indonesia. PTPN 2 owns several palm oil plantations located across different regions, one of which is the Sawit Seberang plantation. At this plantation, the CPO processing is carried out using various methods and technologies to ensure that the quality of the CPO produced meets the standards set for both domestic and international markets (Nurhasanah et al., 2023).

However, in practice, determining CPO quality manually is not always easy, as it involves many factors and variations. Some factors that affect CPO quality include moisture content, free fatty acids, and impurities. To address these challenges, a systematic and efficient approach is needed. One commonly used approach in classification and quality determination is the Naive Bayes classification method.

The Naive Bayes method is a classification technique based on Bayes' theorem (Suryani et al., 2021). This method is frequently used in various applications such as text classification, pattern recognition, and image classification. The main advantages of the Naive Bayes method are its simplicity, speed, and efficiency in use. Although it makes the very simple assumption of independence between features, in many cases, it produces satisfactory results.

In the context of determining CPO quality, applying the Naive Bayes method can be an effective solution to aid decision-making. By utilizing relevant CPO quality features, the Naive Bayes model can classify CPO quality into various categories, such as export quality, industrial quality, or other standards.

In line with this, this research aims to apply the Naive Bayes method to determine CPO quality at PTPN 2 Sawit Seberang. By implementing this method, it is expected that a classification model can be developed to assist in the quick and efficient determination of CPO quality. Moreover, the model is anticipated to improve accuracy and consistency in quality determination, supporting better decision-making at the operational level of the company.

Overall, the application of the Naive Bayes method in determining CPO quality at PTPN 2 Sawit Seberang

* Corresponding author



represents an important and strategic step toward increasing efficiency and accuracy in the company's operational processes. By utilizing advanced technology and data analysis methods, a more effective and efficient system for managing product quality can be established, contributing to the growth and development of the palm oil plantation industry in Indonesia.

LITERATURE REVIEW

Recent studies have shown that the Naive Bayes method is a popular choice for predicting and classifying the quality of crude palm oil (CPO). One study found that the Naive Bayes algorithm is used to calculate probabilities based on various criteria related to CPO, which are then optimized to predict product quality through these probability calculations (Sidik, 2024).

In practical applications, the Naive Bayes method has been applied to process data on palm oil quality in Labuhanbatu Selatan using tools like Jupyter Notebook and Python. Testing with different training-to-testing data ratios revealed that the best accuracy, 36%, was achieved with a 70% training and 30% testing ratio (Nurhasanah et al., 2023).

Another study utilizing a web-based system for CPO quality classification reported an accuracy of 82.05%, demonstrating the effectiveness of the Naive Bayes method in handling large datasets. The study also suggested that increasing the amount of test data could further improve classification performance (Suryani et al., 2021).

More broadly, the Naive Bayes algorithm has been applied in sentiment analysis on social media, alongside other algorithms like SVM and ensemble methods (Pradipta & Jayadi, 2022). Results showed that Naive Bayes performed competitively, though performance varied depending on the type of data used (Furqan & Fakhri, 2024).

Bayesian techniques have also been used for predicting CPO prices. This method compares the predictions with actual market prices, and while the focus is on price trends, it was suggested that external factors like CPO quality could be incorporated into future models to further enhance predictions (Hussin et al., 2023).

The Naive Bayes method has also proven flexible in other fields, such as diagnosing diseases in Nile tilapia, where it achieved an accuracy of 80% (Pulungan et al., 2024). Additionally, the method has been used to classify palm oil varieties, with a reported accuracy of 64.25%, supporting its effectiveness in quality classification tasks (Puspitasari et al., 2022). Outside of agriculture, the Naive Bayes algorithm has been applied to classify calorie levels in McDonald's menu items, showcasing its versatility in different industries (Manurung et al., 2022).

It has also been used to classify blogger data, achieving an accuracy of 76.27% (Widiastuti et al., 2023). In conclusion, the Naive Bayes method has demonstrated effectiveness across various fields, including CPO quality classification. While accuracy levels can vary depending on the dataset and specific application, the method remains a valuable tool for classification tasks in multiple industries.

METHOD

The research method is the fundamental structure or plan used to guide the research from start to finish. It encompasses various key elements that assist researchers in systematically organizing and conducting their study. This is crucial because the subsequent stages of the research must be presented in a coherent and logical sequence, making the creation of a research method essential. Therefore, the research method must be established beforehand, prior to initiating the research phase.



Fig. 1 Research Flow Diagram

Literature review, also referred to as a literature study, is a crucial part of the research process that involves identifying, reviewing, and analyzing relevant studies on the chosen topic or research problem. The main goal of the literature review is to understand what is already known about the subject, pinpoint gaps in

* Corresponding author

knowledge, and determine how the proposed research will contribute to the field. At this stage, researchers utilize sources such as books, journals, and other scholarly works.

1. Data collection is a systematic process of gathering relevant and accurate information or data to answer research questions, test hypotheses, or evaluate results. The methods used for data collection may vary depending on the type of research, the objectives, and the data sources.
2. Data analysis in this research is essential to ensure the study proceeds effectively. This stage involves identifying the necessary components for the research, particularly the dataset on CPO quality.
3. Data preprocessing, after collecting and analyzing the data, is the next step. It aims to remove inconsistent data to ensure clean data is available for clustering using the Naïve Bayes method. This step also includes selecting and transforming the data into numerical formats. Bayes' theorem is applied using a specific formula

$$P(H|X) = \frac{P(X|H) P(H)}{P(X)} \dots\dots\dots (1)$$

Information :

- X = Data with unknown classes
- H = The X data hypothesis is a specific class
- P(H|X) = Probability of hypothesis H based on condition x (posteriori prob.)
- P(H) = Probability of hypothesis H (prior prob.)
- P(X|H) = Probability X based on that condition
- P(X) = Probability of X

4. Design is the process of planning or organizing structures, patterns, or plans necessary to achieve a specific goal. In this research, the design process involves creating a flowchart that outlines the steps. The following is the Naïve Bayes flowchart, aimed at solving issues related to beauty products, as depicted in the following figure:

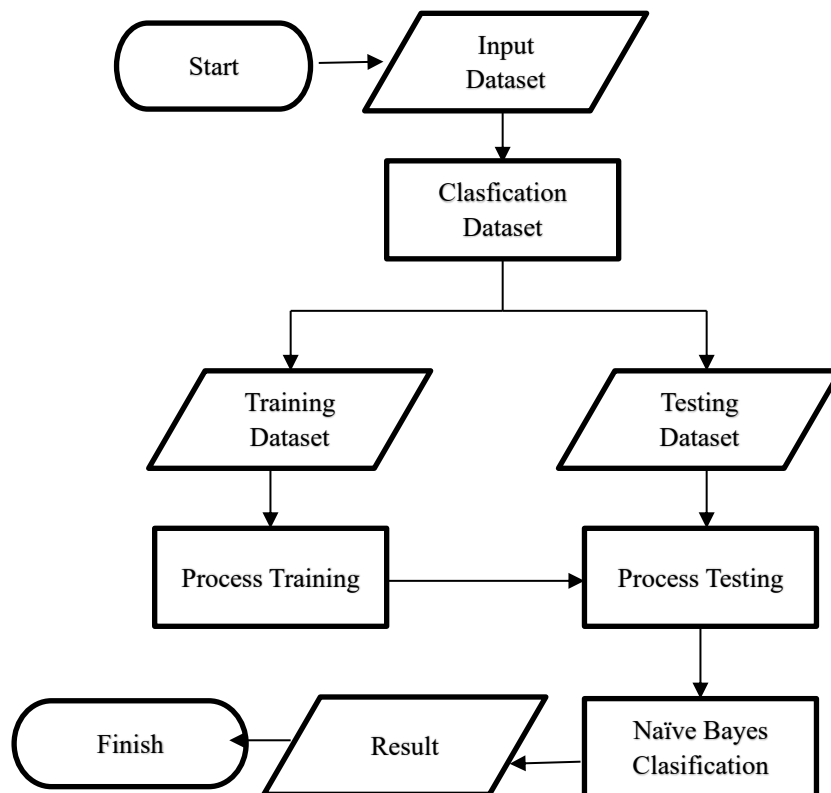


Fig. 2 Naive Bayes Classification Process

* Corresponding author



[Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.](https://creativecommons.org/licenses/by-nc-sa/4.0/)

From the design of the system process above, it has been aligned with the stages of implementing the Naïve Bayes algorithm using RapidMiner. The first step taken is to input the dataset obtained from the cooperative, after which the dataset will be classified into training data and testing data. Once the training data has been processed, the Naïve Bayes function will be called in RapidMiner. This is intended to calculate the probabilities of the testing data, resulting in predictions that will serve as a reference for determining the feasibility of palm oil.

- The testing stage is a stage to find out whether the *rapidminer application* that has been used is in accordance with its functions and outputs.

RESULT

In this study, data testing was conducted using the Naive Bayes algorithm to determine the quality of Crude Palm Oil (CPO) based on factors such as ALB levels, moisture content, and impurity levels. The test results for ALB levels from February to June are presented in Table 1.

Tabel 1 Data Testing Crude Palm Oil

No	ALB Levels					Water levels					Dirt levels				
	Feb	Mar	Apr	May	Jun	Feb	Mar	Apr	May	Jun	Feb	Mar	Apr	May	Jun
1	-	-	-	-	-	-	-	-	-	-	0,000	0,000	0,000	0,000	0,000
2	-	3,45	3,45	-	3,43	-	0,193	0,193	-	0,194	0,000	0,018	0,018	0,000	0,018
3	-	-	3,50	3,54	3,31	-	-	0,193	0,194	0,192	0,000	0,000	0,018	0,018	0,019
4	3,50	-	-	3,35	3,33	0,193	-	-	0,193	0,193	0,018	0,000	0,000	0,017	0,018
5	-	-	-	-	3,32	-	-	-	-	0,193	0,000	0,000	0,000	0,000	0,018
-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-
26	-	-	-	3,42	3,31	-	-	-	0,193	0,194	0,000	0,000	0,000	0,018	0,020
27	3,31	3,51	3,51	3,41	3,36	0,193	0,193	0,193	0,193	0,197	0,018	0,018	0,018	0,018	0,022
28	-	3,52	3,50	3,43	3,52	-	0,193	0,193	0,195	0,198	0,000	0,018	0,018	0,019	0,024
29	-	-	-	3,42	3,53	-	-	-	0,195	0,198	0,000	-	0,000	0,019	0,024
30	3,40	-	3,50	3,41	3,51	0,193	-	0,193	0,193	0,198	0,018	-	0,018	0,019	0,021
31	3,40	-	3,50	-	3,51	0,193	-	0,193	-	0,198	0,018	-	0,018	0,000	0,021

After data is determined, the next step is to calculate the number of quality classifications of palm oil products based on the palm oil product quality data used. The calculation of prior probability is essential in determining the quality of palm oil. In determining each probability for each criterion, the quality of each part of the criteria used is calculated. Thus, the calculation of the probability of each criterion. Each criterion in numerical form will facilitate the calculation in producing the calculation of the probability of each criterion.

Tabel 2 Criterion

No.	Quality Qualification Statement
1	Low
2	High

The following research data can be used for manual calculations of Naive Bayes. The manual calculation data is aggregated daily according to the month. Additionally, the manual calculations are performed to determine the probabilities for each criterion.

- $P(\text{Quality}) = \frac{0}{30} = 0$

- $P(\text{Quality}) = \frac{0,441}{30} = 0,015$

- $P(\text{Quality}) = \frac{0,634}{30} = 0,021$

- $P(\text{Quality}) = \frac{0,439}{30} = 0,015$

* Corresponding author



5. $P(\text{Quality}) = \frac{0,221}{30} = 0,007$
6. $P(\text{Quality}) = \frac{0,671}{30} = 0,022$
7. $P(\text{Quality}) = \frac{0,424}{30} = 0,014$
8. $P(\text{Quality}) = \frac{0,425}{30} = 0,014$
9. $P(\text{Quality}) = \frac{0,461}{30} = 0,015$
10. $P(\text{Quality}) = \frac{0,653}{30} = 0,022$

And so on until the 30th. Then, the midpoint of the prior probabilities is determined to classify the quality of CPO, where the midpoint value of the data is 0.15. It can be assumed that values <0.15 are classified as high-quality classes, while values >0.15 are classified as low-quality classes.

Tabel 3 Class Classification Assumptions

No.	Total Up	Prior Probability	Classification
1	0	0	Tall
2	0.441	0.015	Low
3	0.634	0.021	Low
4	0.439	0.015	Low
5	0.211	0.007	Tall
6	0.671	0.022	Low
7	0.424	0.014	Tall
8	0.425	0.014	Tall
9	0.461	0.015	Low
10	0.653	0.022	Low
11	0.215	0.007	Tall
12	0.443	0.015	Low
13	0.457	0.015	Low
14	0.631	0.021	Low
15	0.212	0.007	Tall
16	0.460	0.015	Low
17	0.651	0.022	Low
18	0.425	0.014	Tall
19	0.441	0.015	Low
20	0.673	0.022	Low
21	0.211	0.007	Tall
22	0.233	0.008	Tall
23	0.427	0.014	Tall
24	0.463	0.015	Low
25	0.231	0.008	Tall
26	0.425	0.014	Tall
27	0.677	0.023	Low
28	0.665	0.022	Tall
29	0.436	0.015	Low

* Corresponding author



30	0.660	0.022	Tall
31	0.448	0.015	Low

A total of 14 data points were classified as high quality, while 16 data points were classified as low quality.

$$P(\text{Quality CPO} \mid \text{Class Quality})$$

$$P(\text{Quality CPO} \mid \text{Quality High}) = \frac{14}{30} = 0,47$$

$$P(\text{Quality CPO} \mid \text{Quality Low}) = \frac{16}{30} = 0,53$$

Split Validation is a validation technique that divides data into two parts: a portion for training and a portion for testing. The prepared data for classification is split into two using random sampling techniques: 80% for training data and 20% for testing data. An example of the calculation for selecting testing data is as follows:

The total number of data points (N) for ALB content: 31

Number of testing data: 20% × 31= 6

Thus, the data used consists of the dates from the 26th to the 31st of each month from February to June. The number of training data points is 31– 6 = 2531, meaning that the data used includes the 1st to the 25th of each month as shown in the table below.

Tabel 4 Result Of The Data Split

No	ALB Levels					Water Levels					Dirt levels				
	Feb	Mar	Apr	May	Jun	Feb	Mar	Apr	May	Jun	Feb	Mar	Apr	May	Jun
1	-	-	-	3,42	3,31	-	-	-	0,193	0,194	0,000	0,000	0,000	0,018	0,020
2	3,31	3,51	3,51	3,41	3,36	0,193	0,193	0,193	0,193	0,197	0,018	0,018	0,018	0,018	0,022
3	-	3,52	3,50	3,43	3,52	-	0,193	0,193	0,195	0,198	0,000	0,018	0,018	0,019	0,024
4	-		-	3,42	3,53	-		-	0,195	0,198	0,000		0,000	0,019	0,024
5	3,40		3,50	3,41	3,51	0,193		0,193	0,193	0,198	0,018		0,018	0,019	0,021
6	3,40		3,50	-	3,51	0,193		0,193	-	0,198	0,018		0,018	0,000	0,021

The tabel 5 below shows a table containing attributes such as production date, ALB levels, impurity levels, water content, and quality. This display is shown from reading an Excel file in RapidMiner to add test data, which will later be presented on the test results page. In the image above, there are test data that have been input, with 2 high-quality entries and 4 low-quality entries.

Tabel 5 Quality Result

ALB Levels	Water Levels	Dirt Levels	Quality
0	0	0	High
3.31	0.193	0.018	Low
3.31	0.193	0.018	Low
3.4	0.193	0.018	Low
3.4	0.193	0.018	Low
3.52	0.193	0.018	High
3.52	0.193	0.018	High
0	0	0	Low
0	0	0	High

* Corresponding author



Accuracy: Accuracy = $(TP + TN) / (TP + TN + FP + FN) = (5 + 15) / (5 + 5 + 5 + 15) = 20/30 = 0.666 = 66.6\%$
Thus, the error accuracy is:

Error Accuracy = $100\% - 66.6\% = 33.3\%$

Precision: Precision = $TP / (TP + FP) = 5 / (5 + 5) = 5/10 = 0.5 = 50\%$

Recall: Recall = $TP / (TP + FN) = 5 / (5 + 15) = 5/20 = 0.25 = 25\%$

Specificity: Specificity = $TN / (TN + FP) = 15 / (15 + 5) = 15/20 = 0.75 = 75\%$

Tabel 6 Predictions Naïve Bayes

	Predictions	
	High Quality	Low Quality
High Quality	TF = 5	TN = 5
Low Quality	PP = 5	FN = 15

Thus, the system's accuracy in classifying CPO quality using the Naive Bayes method is 66.6%, with a precision of 50%, a recall of 50%, and a specificity of 75%.

DISCUSSIONS

In the discussion of this study, it is evident that the application of the Naive Bayes algorithm has provided a systematic approach for determining the quality of Crude Palm Oil (CPO) based on several key factors, including ALB levels, moisture content, and impurity levels. The data from February to June, as presented in the results, show that the Naive Bayes model offers a practical solution for classifying CPO quality efficiently. However, the accuracy of the model, which was 66.6%, indicates that there is room for improvement. This could be due to various factors, such as the limited size of the dataset used or the nature of the features selected for classification. While the model achieved a reasonable balance between true positives (high-quality CPO) and true negatives (low-quality CPO), the precision and recall values were relatively low, at 50% each. This suggests that the model may not always perform well in distinguishing between high and low-quality CPO, particularly when the features overlap or are not strongly indicative of quality.

The specificity of 75% indicates that the model performed well in identifying low-quality CPO, but the recall value reflects its challenges in identifying high-quality CPO. This disparity could point to a need for further refinement in data preprocessing or feature selection, potentially by incorporating additional relevant features or improving data cleaning processes to remove inconsistencies. The use of the Naive Bayes method, despite its simplicity and speed, is inherently limited by its assumption of feature independence. In practice, the factors affecting CPO quality, such as ALB levels, moisture content, and impurities, may not be entirely independent of each other, which could explain some of the performance limitations observed in this study.

Nonetheless, the implementation of this method has proven effective in streamlining the process of quality classification, and its integration into operational workflows can help improve decision-making in the palm oil industry. Further improvements could include exploring other machine learning algorithms or combining Naive Bayes with additional techniques to enhance overall performance and accuracy in predicting CPO quality.

CONCLUSION

In implementing the Naive Bayes method for determining CPO quality at PTPN 2 Sawit Seberang, data collection is conducted for attributes or variables related to CPO quality, such as ALB levels, water content, and impurity levels. The system testing is then carried out using RapidMiner. The determination of CPO quality for training and testing the Naive Bayes model is divided into 80% training data, consisting of 21 data points (from February to June), and 20% testing data, consisting of 6 data points (from February to June). The interpretation of applying Naive Bayes classification in the decision-making context related to CPO quality determination yielded an accuracy of 66.6%, a precision of 50%, a recall of 50%, and a specificity of 75% for the system's classification of CPO quality using the Naive Bayes method.

REFERENCES

Alita, D., Sari, I., Isnain, A. R., & Styawati, S. (2021). Penerapan Naïve Bayes Classifier Untuk Pendukung Keputusan Penerima Beasiswa. *Jurnal Data Mining Dan Sistem Informasi*, 2(1), 17–23.

* Corresponding author



- Boy, A. F. (2020). Implementasi Data Mining Dalam Memprediksi Harga Crude Palm Oil (CPO) Pasar Domestik Menggunakan Algoritma Regresi Linier Berganda (Studi Kasus Dinas Perkebunan Provinsi Sumatera Utara). *Journal of Science and Social Research*, 3(2), 78–85.
- Efendi, R., Faurina, R., & Hamimmah, T. S. (2023). Implementasi Metode Naïve Bayes Pada Penentuan Mutu CPO (Crude Palm Oil). *JSAI (Journal Scientific and ...)*
- Furqan, M., & Fakhri, A. ab. N. (2024). Big data approach to sentiment analysis in machine learning-based microblogs: perspectives of religious moderation public policy in indonesia. In *journal of applied engineering and technological science* (vol. 5, issue 2).
- Harahap, R. R., & Furqan, M. (2024). Sentiment Analysis towards the 2024 Vice Presidential Candidate Debate Using the Support Vector Machine Algorithm. *Sinkron: jurnal dan penelitian teknik informatika*, 8(3), 1783-1794.
- Hussin, M., Ismail, Z., & Ilias, I. S. C. (2023). Bayesian Network Design for Crude Palm Oil (CPO) Price Prediction Driven by Fluctuation Patterns and Trends. *Journal of Advanced Research in Applied Sciences and Engineering Technology*, 31(2), 117–129. <https://doi.org/10.37934/araset.31.2.117129>
- Manurung, M., Siti Dzulhijjah Nur Ammarah, N., & Nisa Sofia Amriza, R. (2022). The Application Naive Bayes Algorithm Determines Calorie Level Of The Mcdonald's Menu. In *JTSI* (Vol. 3, Issue 2).
- Nurhasanah, D., Lestari, D. A., & Simatupang, S. (2023a). Pemilihan Kualitas Produk Kelapa Sawit Menggunakan Metode Naive Bayes Di Labuhanbatu Selatan. *Jurnal Teknisi*.
- Nurhasanah, D., Lestari, D. A., & Simatupang, S. (2023b). Pemilihan Kualitas Produk Kelapa Sawit Menggunakan Metode Naive Bayes Di Labuhanbatu Selatan. *Jurnal Teknisi*, 3(1), 24. <https://doi.org/10.54314/teknisi.v3i1.1254>
- Pradipta, R., & Jayadi, R. (2022). The Sentiment Analysis Of The Indonesian Palm Oil Industry In Social Media Using A Machine Learning Model. *Journal of Theoretical and Applied Information Technology*, 30(12). www.jatit.org
- Pulungan, R. W., Sriani, S., & Armansyah, A. (2024). Implementation of Naïve Bayes Method Diagnosing Diseases Nile Tilapia. *Journal of Computer Networks, Architecture and High Performance Computing*, 6(2), 817–828. <https://doi.org/10.47709/cnahpc.v6i2.3834>
- Puspitasari, N., Rosmasari, R., Pratama, F. W., & Sulastri, H. (2022). Quality Classification of Palm Oil Varieties Using Naive Bayes Classifier. *Digital Zone: Jurnal Teknologi Informasi Dan Komunikasi*, 13(1), 11–23. <https://doi.org/10.31849/digitalzone.v13i1.9773>
- Sidik, A. (2024). Data Mining Menggunakan Metode Naive Bayes Untuk Menetapkan Standar Untuk Produk Minyak Sawit Mentah. *EJECTS: Journal Computer, Technology, and ...*
- Suryani, D., Yulianti, A., Maghfiroh, E. L., & Alber, J. (2021). *SISTEMASI: Jurnal Sistem Informasi Klasifikasi Kualitas Produk Kelapa Sawit Menggunakan Metode Naïve Bayes Quality Classification of Palm Oil Products Using Naïve Bayes Method*. <http://sistemasi.ftik.unisi.ac.id>
- Widiastuti, N., Hermawan, A., & Avianto, D. (2023). IMPLEMENTASI METODE NAÏVE BAYES UNTUK KLASIFIKASI DATA BLOGGER. *JUPI (Jurnal Ilmiah Penelitian Dan Pembelajaran Informatika)*, 8(3), 985–994. <https://doi.org/10.29100/jipi.v8i3.3713>

* Corresponding author



[Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.](https://creativecommons.org/licenses/by-nc-sa/4.0/)