

Price Forecasting of Shallots Using the Machine Learning Approach of Random Forest Regression Supporting Price Stabilization

Muhammad Naufal Rauf Ibrahim^{1*}

¹Department of Agricultural and Biosystem Engineering, Faculty of Agroindustrial Technology, Padjadjaran University, Jl. Ir. Soekarno KM.21, Sumedang, West Java 54363 Indonesia.

*Corresponding author, email: rauf.ibrahim@unpad.ac.id

Article Info	Abstract
<p><i>Submitted: 17 July 2025</i> <i>Revised: 17 September 2025</i> <i>Accepted: 22 September 2025</i> <i>Available online: 9 October 2025</i> <i>Published: September 2025</i></p> <p>Keywords: Machine Learning; Price Forecasting; Random forest regression; Shallot.</p> <p>How to cite: Ibrahim, M. N. R. (2025). Price Forecasting of Shallots Using the Machine Learning Approach of Random Forest Regression Supporting Price Stabilization. <i>Jurnal Keteknikan Pertanian</i>, 13(3): 449-461. https://doi.org/10.19028/jtep.013.3.449-461.</p>	<p><i>Shallots (<i>Allium cepa</i> L.) are a major horticultural commodity in Indonesia, with a production of 1.98 million tons in 2022, representing 13.59% of the total national vegetable production. Accurate forecasting of agricultural commodity prices is fundamental to sustainable development in the agricultural sector and contributes to broader economic stability. This study uses the random forest regression algorithm, a supervised machine learning technique that utilizes ensemble learning to combine multiple decision trees. This approach offers advantages in modeling non-linear relationships for agricultural price prediction while also reducing the risk of overfitting, resulting in more accurate and stable forecasts compared to individual decision trees. The purpose of this research is to develop and optimize a shallot price forecasting model using random forest regression. The optimized model, using 50 decision tree estimators, successfully predicted up to 15 months ahead of monthly prices and achieved an RMSE of 2363.15 and a MAPE of 8.71% in validation, then a MAPE of 10.31% in test evaluation.</i></p>

Doi: <https://doi.org/10.19028/jtep.013.3.449-461>

1. Introduction

Shallots (*Allium cepa* L) are one of the largest horticultural commodities in Indonesia, with production reaching 1.98 million tons in 2022, accounting for 13.59% of total vegetable production (BPS, 2024). National shallot consumption in 2022 reached 890 thousand tons and continues to increase by an average of 3.88% per year (Pusdatin Kementan, 2023). The existence of shallots plays a vital role in influencing the economy and creating job opportunities in Indonesia. The survey show that each hectare of shallot farmland can create approximately 290 workdays.(Wandschneider et al., 2013).

The forecasting of agricultural commodity prices serves as an important tool for sustainable development in the agricultural economy and broader economic stability. The ability to predict price movements enables farmers to make informed decisions about when to plant and sell their crops, potentially allowing them to switch between commodities and alternative markets to ensure favorable prices and maximize income. The price of shallots at the market fluctuates every month, and farmers

often suffer losses due to falling selling prices during the harvest season. On the other hand, consumers also feel disadvantaged when the price of shallots soars during periods of low availability. The fluctuating price have a significant impact, especially on people with low incomes (Matondang et al., 2024). Price forecasting is not only important for farmers and consumers but also provides future information on agricultural commodity prices that can help the government in formulating policies to maintain the stability of shallot prices.

The purpose of this research is to develop a forecasting model for shallot prices using the random forest method. To optimize the accuracy of price forecasts and computational data requirements, tests will be conducted with various numbers of estimators to determine the optimal value. This model can later be integrated into information systems or applications, allowing farmers to access it to predict shallot prices.

2. Material and Methods

2.1 Data Acquisition and Preprocessing

Monthly data on shallots prices is obtained from the official government website PIHPS (Pusat Informasi Harga Bahan Pangan Strategis Nasional, akses: <https://www.bi.go.id/hargapangan>) from the Bank Indonesia database as of July 30, 2025 for the Bandung area. To maximize the training data, the monthly data is obtained based on the oldest data available in the system, from October 2018 to May 2025, containing 81 data points. This dataset consists of two attributes, namely date and shallot price.

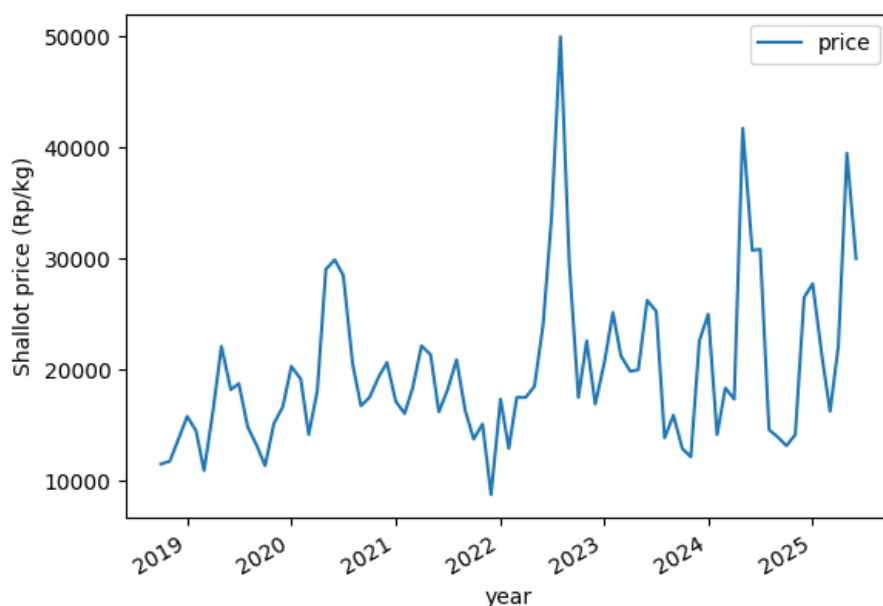


Figure 1. Shallot monthly prices from October 2018 to May 2025.

The system's default data format is in .xls and the price data is still in text format. Once the data is obtained, data pre-processing is necessary to ensure smooth processing. The data pre-processing steps include 1) checking for missing data, 2) standardizing the numbering format (for example, converting text to numbers) and simplifying the names of row indices, 3) converting to .csv, and 4) converting dates to indices. Pre-processing is carried out using Excel and Python. The ready-to-use monthly shallot price data is shown in Figure 1 above.

2.2 Training, Validation, and Testing of the Random Forest Regression Model

Random forest regression is a powerful ensemble learning technique that combines the predictions of multiple decision trees to improve the accuracy and stability of estimates (Breiman, 2001). This method is very useful for modeling complex non-linear relationships in data, which are common in agricultural price forecasting. Random Forest Regression is a part of machine learning algorithms, which is a branch of artificial intelligence. This algorithm is a supervised learning algorithm that combines several decision trees to produce predictions that are more accurate and stable compared to those achieved by individual trees through ensemble learning (Fitri, 2023). Ensemble learning randomly samples training data using bootstrap aggregation and combines multiple models, which are decision trees. The predictions from all these models are then averaged to minimize errors caused by outliers and overfitting. In addition to its ability to handle outliers and reduce the likelihood of overfitting, random forest regression can capture both linear and non-linear relationships in agricultural commodity price data, where weather, climate change, and market demand influence price fluctuations. (Zhang et al., 2020). Figure 2 illustrates how the Random Forest Regression method resembles a forum of experts, each analyzing and predicting using different random data subsets, with the final decision made based on a vote or the average conclusion of the experts.

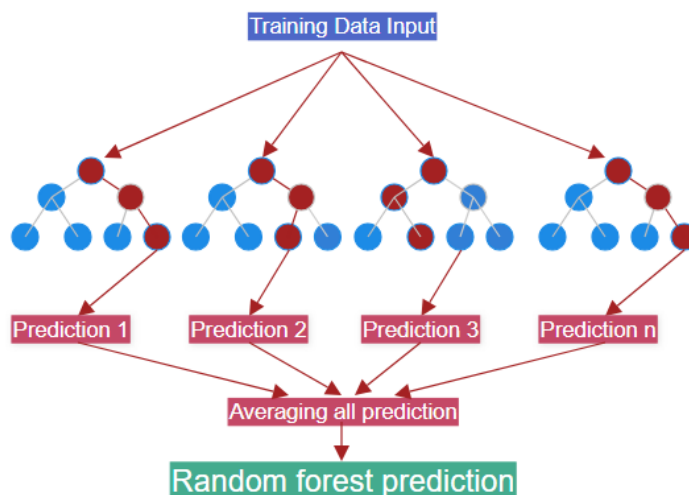


Figure 2. The random forest algorithm process for making predictions.

Several previous studies have examined the forecasting of agricultural commodity prices using methods such as LR (Linear Regression) (Ge & Wu, 2020), ARIMA (Weng et al., 2019), SARIMA (Chen et al., 2025), and machine learning techniques such as LSTM (Long and Short Term Memory) (Choudhary et al., 2025; Murugesan et al., 2022), ANN (Artificial Neural Network) (Anggraeni et al., 2018; Mahto et al., 2021), and SVR (Support Vector Regression) (Liu et al., 2019; Sun et al., 2023). The regression method seeks linear relationships within the data to describe the movement of agricultural commodity prices; however, accuracy decreases if the data has non-linear relationships and a significant number of outliers. A sudden price increases due to seasonal events, such as Ramadan, greatly affects the accuracy of regression methods. Machine learning methods have the advantage of capturing non-linear relationships and seasonal events to forecast agricultural commodity prices (Paul et al., 2022). In addition, machine learning can be combined with other non-linear parameters such as rainfall data, national commodity production, and temperature as supporting data to improve price forecasting accuracy. Despite these advantages, machine learning has several drawbacks, such as requiring large datasets, which in turn demands greater computational power. For forecasting, machine learning divides the data into two parts: training data and test data. Machine learning is at risk of overfitting if the dataset is not properly processed and tuned beforehand. Random forest regression is a part of machine learning that has a lower risk of overfitting compared to other machine learning methods, while also being able to effectively capture non-linear relationships that are usually not achievable with standard regression methods (Wang et al., 2016).

Although it possesses significant advantages, research on the use of random forest algorithms in the field of agricultural commodity price forecasting, especially shallots, has not been widely explored. In fact, random forest has advantages over other algorithms for short- and medium-term price forecasting. Short-term forecasting is crucial for farmers as it serves as the basis for determining when to sell, while medium-term forecasting can help farmers plan shallot planting schedules for the upcoming season or throughout the year.

The number of decision trees or estimators used in decision making affects the accuracy of price forecasting. The more decision trees there are, the higher the accuracy will be, but data processing requirements will also increase, which can reduce computational speed (Salman et al., 2024). Optimization of the number of decision trees (estimators) is necessary to achieve maximum forecast accuracy with minimal computational processing. In this study, sk-learn from the Python library was used to train and validate the random forest model.

The random forest model is initialized by specifying the number of decision trees (n estimators) in the ensemble. Each tree is grown using a random subset of the training data and a random subset of features, and then each estimator predicts the outcome based on the training data. Finally, the final prediction is determined by averaging the results from all the estimators. For random forest regression with M trees, the final prediction is calculated as:

$$\hat{y} = \frac{1}{M} \sum_{m=1}^M h_m(x) \quad (1)$$

Where \hat{y} is the final prediction value and $h_m(x)$ is prediction of m-th estimators.

In this study, 25, 50, 75, 100, and 125 estimator hyperparameters were selected based on a balance between computational efficiency and model robustness, as referenced in a previous study by Kaewchada et al. (2023). The hyperparameter tuning method applied in this study involves utilizing a simple search technique to identify the optimal number of estimators from 25–125. This method systematically evaluates hyperparameter combinations to determine the best-performing model configuration. Estimators in the forest are built by selecting random samples from the dataset with replacement, a method known as bootstrapping. This implementation has the advantage of generating diverse estimators that can capture various aspects of the data. Ultimately, the performance of each configuration based on the number of estimators will be evaluated.

The historical value of the target variable or predictor, namely the lag feature, is implemented on the training data. In random forest regression, lag features are primarily used in time series analysis to help the model capture patterns and dependencies over time as features. In this study, a 4-month lag was created for the training data, consisting of 4 data sets that describe prices from 1 to 4 months prior, because the shallot planting cycle is 90-120 days. In this case, the model will predict the price based on the prices from 1 to 4 months prior. The price data is transformed into each lag feature data using the sliding window method. For example, on the date December 2023, Lag 1 data is the price in November 2023, then Lag 2 is represented by the price in October 2023, and so on. The table below is an example of the method applied using lag features in this research.

Table 1. Training data from lag feature transformation using a sliding window. The arrow indicates the shift movement of the lag data (tn).

Date	Price	Lag (monthly)			
		Lag-1 (value t-1)	Lag-2 (value t-2)	Lag-3 (value t-3)	Lag-4 (value t-4)
Oct-18	11500				
Nov-18	11750	11500			
Dec-18	13850	11750	11500		
.	.				
.	.				
.	.				
Jan-24	22600	12150	12900	15900	13850
Feb-24	17340	22600	12150	12900	15900

Training data and test data are separated using an 80/20 ratio. The training data consists of data from October 2018 to February 2024, while the test data comes from March 2024 to May 2025. The training dataset is used to train the random forest model. The test data is then used after the model is trained to assess the model's performance.

2.3 Model Evaluation

This study evaluates the training and validation models. The training model is evaluated using the coefficient of determination (R^2) to determine how well the model fits the data, and then evaluated using Root Square Mean Error (RSME) and Mean Average Percentage Error (MAPE). For validation and testing, performance will also be evaluated using RSME and MAPE. RSME and MAPE were chosen based on the usefulness of these metrics in explaining errors more easily because RSME and MAPE can intuitively convey the magnitude of the error. The following are the formulas for calculating R^2 , RSME, and MAPE:

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (2)$$

$$RSME = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3)$$

$$MAPE = \frac{100\%}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (4)$$

Where R^2 = coefficient of determination, n = total number of data points, y_i = The actual value of the target variable for data i point, \hat{y}_i = Predicted value of the target variable for the data point i , \bar{y} = The average of the actual values. The overall workflow of model development, from data preprocessing to validation, is illustrated in Figure 3.

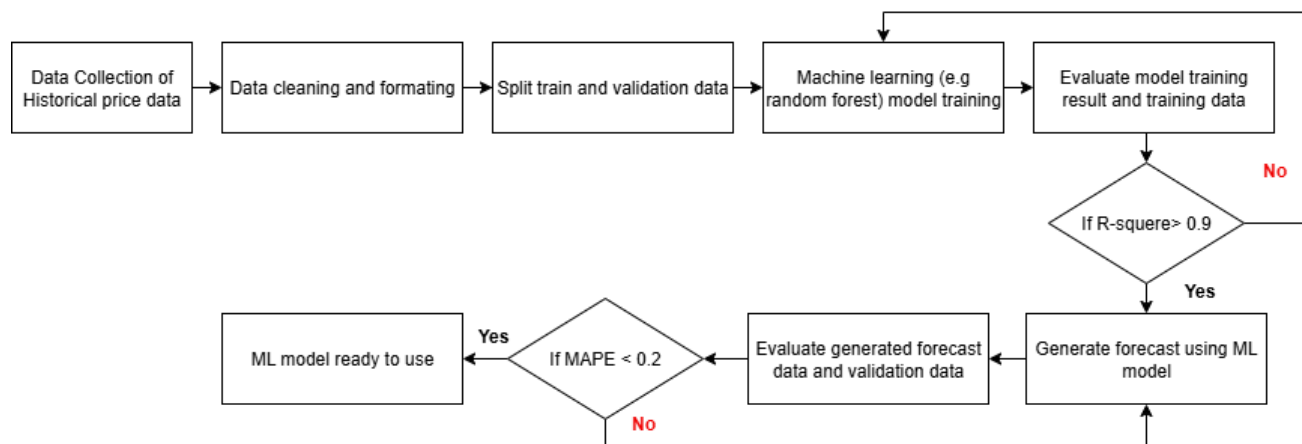


Figure 3. Workflow of developing a machine learning (ML) model for price forecasting.

3. Results and Discussion

3.1 Training and Model Deployment Performance

The random forest regression model demonstrated excellent performance during the training phase on historical shallot price data from October 2018 to February 2024. Each hyperparameter, consisting of 25 to 125 estimator models, was trained and evaluated. The model achieved an R^2 score of 0.89 to 0.9, indicating that up to 90% of the variance in shallot prices was explained by the ensemble learning approach. The training phase was also evaluated using RMSE and MAPE, with scores of 2,116.79 IDR/kg and 7.99%, respectively. This performance is consistent with previous studies on price prediction using the random forest algorithm, which have demonstrated superior predictive capabilities compared to traditional regression models.(Jui et al., 2020).

The bootstrap methodology employed in random forests has proven particularly effective for handling the inherent volatility of market prices (Kara et al., 2021). By creating multiple estimators trained on different bootstrap samples of the training data, the model successfully captured the complex nonlinear relationships between historical price patterns and market dynamics. The results of the model fitting using 100 estimators as the best result are shown in Figure 3.

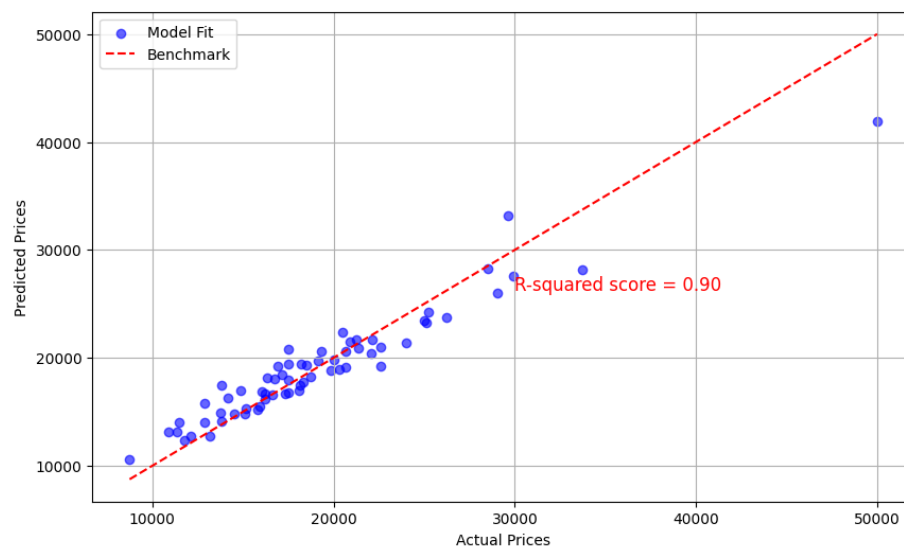


Figure 3. Model fitting in training data with n of estimators = 100.

3.2 Model Prediction

The random forest model demonstrates outstanding generalization ability when evaluated on out-of-sample test data covering the period from March 2024 to May 2025. This model achieves an excellent MAPE range of 8.71% to 9.42% on the test data, indicating consistent performance across various data periods of up to 15 months. These error metrics are comparable to other machine learning

approaches for agricultural price forecasting, where MAPE values between 10-20% are considered to have good forecasting accuracy (Kaewchada et al., 2023).

The model's ability to accurately predict extreme price fluctuations, such as the sharp increase to Rp 41,750/kg in June 2024, demonstrates the effectiveness of the random forest ensemble learning in capturing market volatility. Figure 4 shows a comparison between the test data and the model's prediction data using 50 estimators as the optimal n estimator.

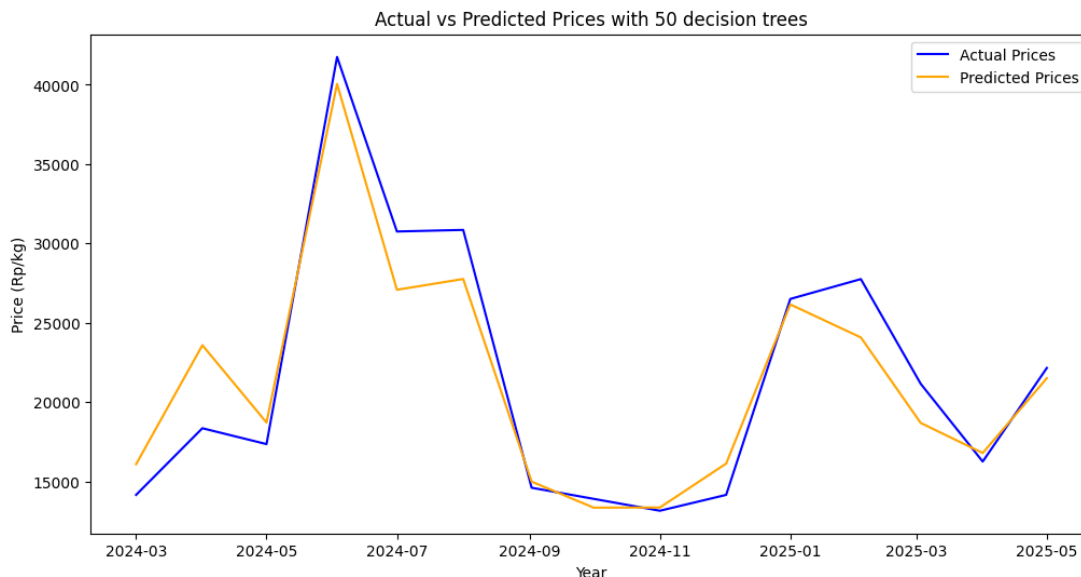


Figure 4. Test price vs predicted price with n estimators = 50.

3.3 Optimizing Number of Estimators

Evaluation of the n parameter estimator reveals important insights into the relationship between ensemble size and model performance. The analysis shows that the R^2 score increases progressively from 0.89 (25 trees) to 0.90 (125 trees) as shown in Figure 5, following the expected pattern in ensemble learning theory, where adding base learners generally improves overall performance on training data.(Mienye & Sun, 2022).

However, the RMSE performance demonstrated the most optimal results at 50 estimators (2,363), with gradual increases in the error rates as the number of trees exceeded this value. This finding suggests that while additional trees continue to improve the coefficient of determination, they may introduce slight noise or computational overhead that affects absolute error measurements and prediction errors. The MAPE similarly demonstrated that the most optimal performance was achieved with 50 estimators, resulting in an 8.71% MAPE score. This occurred despite 50 not being the optimal number of estimators in the training setting, which highlights the importance of balancing model complexity with predictive accuracy. A comparison between each estimator size is shown in Figure 6.

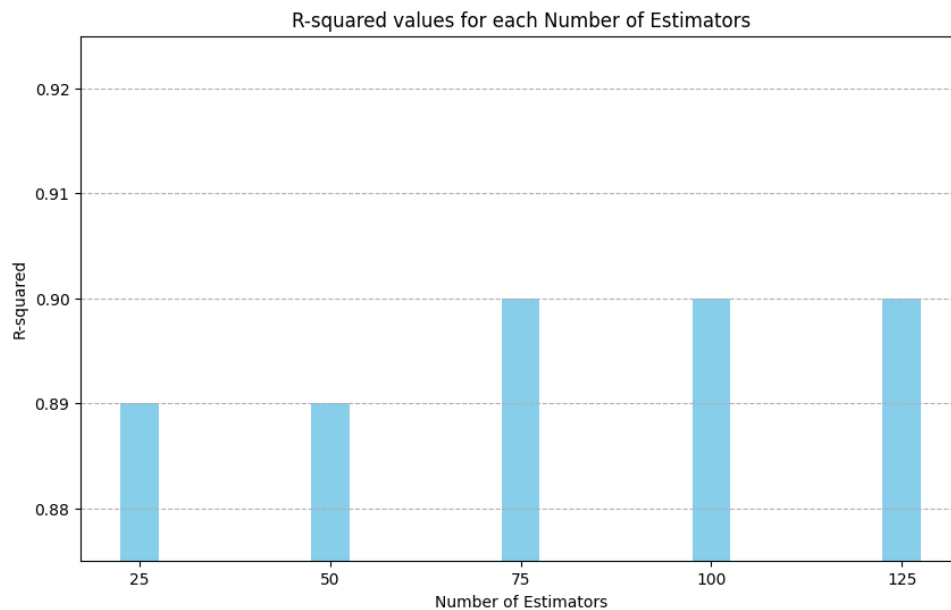


Figure 5. R^2 values for each n estimator from 25 to 125.

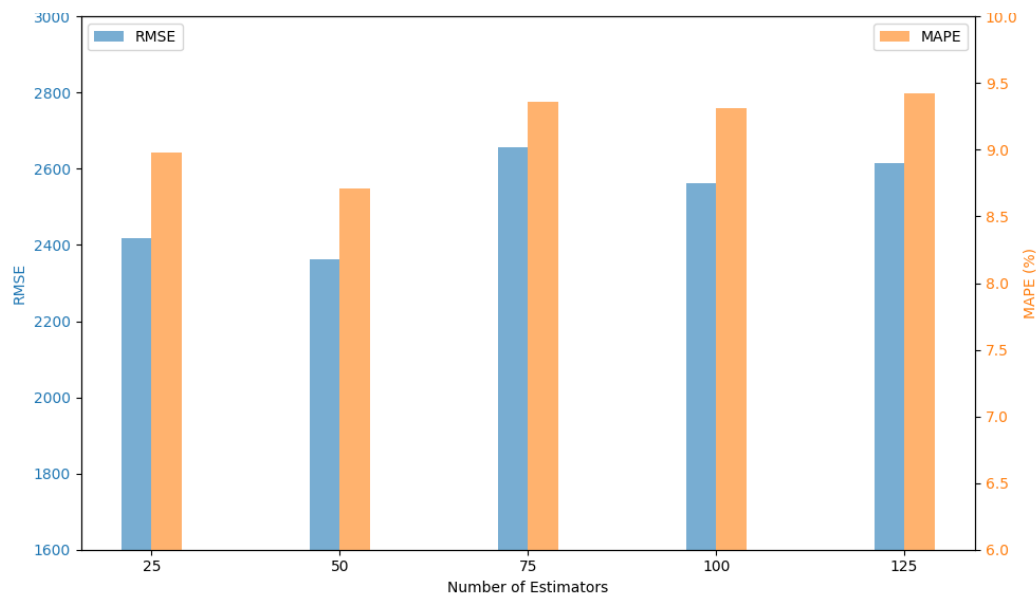


Figure 6. Error evaluation of based on each n estimator in the validation setting.

The results of hyperparameter optimization for the estimator are in line with the established guidelines for tuning random forests, where the optimal number of estimators is usually specific to each problem and depends on the complexity of the data. The effectiveness of the ensemble learning approach in capturing non-linear relationships and seasonal patterns addresses key challenges in agricultural price prediction, where traditional methods often struggle with market volatility and complex interdependencies. The bootstrap methodology provides inherent resilience against outliers

and missing data, which are common issues in agricultural datasets. This resilience is especially valuable for developing countries where data quality and availability may be limited.

Finally, the model's performance in forecasting from June to July 2025 is evaluated using MAPE, with data obtained from PIHPS. Since the model requires 4 lag features as input, prediction data from the previous $n-1$ months is adopted as input for the lag features. Based on the table below, the model achieves a MAPE score of 10.31%, which is considered an acceptable result for prediction performance.

Table 2. Model evaluation for forecasting the testing phase of Jun-25 and Jul-25.

Date	Actual Value	Predict	Error
Jun-25	39500	34407	-5093
Jul-25	30000	27682	-2318
MAPE			10.31%

The findings regarding optimal hyperparameter settings provide practical guidance for implementing the random forest model in agricultural applications. Identifying the use of 50 estimators as the optimal range offers a balance between predictive accuracy and computational efficiency, making this approach accessible for various implementation scenarios. The proposed price forecasting model can assist farmers in deciding when to plant shallots, based on future prices. This has the potential to help stabilize farmers' incomes and, consequently, stabilize future market prices.

Future research should explore the integration of additional variables such as weather data, market indicators, and policy variables to further enhance the model's predictive capabilities and expand its application to various agricultural commodities, supporting farmers' decision-making across different crops.

4. Conclusion

This study provides an application of the random forest regression machine learning algorithm to predict shallot prices up to 15 months ahead by training on approximately 5 years of data. The training, validation, and prediction of the random forest showed promising results with MAPE values of 8.71%, 9.39%, and 10.31% respectively, which represent excellent performance for forecasting. The algorithm not only predicts based on historical price movements in the training data, but also captures the volatility of shallot prices by successfully anticipating the price surge in June 2024 during the validation setting. The finding that 50 estimators offer an optimal performance balance has important implications for computational efficiency and the implementation of the model in real-world agricultural price forecasting systems. Future research should address the integration of additional

features, other agricultural commodities, and practical utility for agricultural stakeholders (farmers, consumers, and policymakers).

5. References

- Anggraeni, W., Mahananto, F., Rofiq, M. A., Andri, K. B., Sumaryanto, Zaini, Z., & Subriadi, A. P. (2018). Agricultural Strategic Commodity Price Forecasting Using Artificial Neural Network. 2018 International Seminar on Research of Information Technology and Intelligent Systems (ISRITI), 347–352. <https://doi.org/10.1109/ISRITI.2018.8864442>
- BPS, I. (2024). Produksi Tanaman Sayuran—Tabel Statistik. <https://www.bps.go.id/id/statistics-table/2/NjEjMg==/produksi-tanaman-sayuran.html>
- Breiman, L. (2001). Random Forests. *Machine Learning*, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
- Chen, D., Cheng, B., Li, J., Shen, Q., & Liu, J. (2025). A study on the forecasting of agricultural product price trends based on SARIMA. 2025 IEEE 8th Information Technology and Mechatronics Engineering Conference (ITOEC), 8, 441–445. <https://doi.org/10.1109/ITOEC63606.2025.10968760>
- Choudhary, K., Jha, G. K., Jaiswal, R., & Kumar, R. R. (2025). A genetic algorithm optimized hybrid model for agricultural price forecasting based on VMD and LSTM network. *Scientific Reports*, 15(1), 9932. <https://doi.org/10.1038/s41598-025-94173-0>
- Fitri, E. (2023). Analisis Perbandingan Metode Regresi Linier, Random Forest Regression dan Gradient Boosted Trees Regression Method untuk Prediksi Harga Rumah. *Journal of Applied Computer Science and Technology*, 4(1), 58–64. <https://doi.org/10.52158/jacost.v4i1.491>
- Ge, Y., & Wu, H. (2020). Prediction of corn price fluctuation based on multiple linear regression analysis model under big data. *Neural Computing and Applications*, 32(22), 16843–16855. <https://doi.org/10.1007/s00521-018-03970-4>
- Jui, J. J., Imran Molla, M. M., Bari, B. S., Rashid, M., & Hasan, M. J. (2020). Flat Price Prediction Using Linear and Random Forest Regression Based on Machine Learning Techniques. In M. A. Mohd Razman, J. A. Mat Jizat, N. Mat Yahya, H. Myung, A. F. Zainal Abidin, & M. S. Abdul Karim (Eds.), *Embracing Industry 4.0* (pp. 205–217). Springer. https://doi.org/10.1007/978-981-15-6025-5_19
- Kaewchada, S., Ruang-On, S., Kuhapong, U., & Songsri-in, K. (2023). Random forest model for forecasting vegetable prices: A case study in Nakhon Si Thammarat Province, Thailand. *International Journal of Electrical and Computer Engineering (IJECE)*, 13(5), 5265. <https://doi.org/10.11591/ijece.v13i5.pp5265-5272>
- Kara, M., Atici, K. B., & Ulucan, A. (2021). Price and Volatility Forecasting in Electricity with Support Vector Regression and Random Forest. In A. B. Dorsman, K. B. Atici, A. Ulucan, & M. B. Karan (Eds.), *Applied Operations Research and Financial Modelling in Energy: Practical Applications*

- and Implications (pp. 101–124). Springer International Publishing. https://doi.org/10.1007/978-3-030-84981-8_6
- Liu, Y., Duan, Q., Wang, D., Zhang, Z., & Liu, C. (2019). Prediction for hog prices based on similar sub-series search and support vector regression. *Computers and Electronics in Agriculture*, 157, 581–588. <https://doi.org/10.1016/j.compag.2019.01.027>
- Mahto, A. K., Alam, M. A., Biswas, R., Ahmed, J., & Alam, S. I. (2021). Short-Term Forecasting of Agriculture Commodities in Context of Indian Market for Sustainable Agriculture by Using the Artificial Neural Network. *Journal of Food Quality*, 2021(1), 9939906. <https://doi.org/10.1155/2021/9939906>
- Matondang, M. R., Krisnamurthi, B., & Herawati, H. (2024). Price Fluctuations And Volatility Of National Strategic Food Commodities In Indonesia. *Agrisocionomics: Jurnal Sosial Ekonomi Pertanian*, 8(1), 134–146. <https://doi.org/10.14710/agrisocionomics.v8i1.17753>
- Mienye, I. D., & Sun, Y. (2022). A Survey of Ensemble Learning: Concepts, Algorithms, Applications, and Prospects. *IEEE Access*, 10, 99129–99149. <https://doi.org/10.1109/ACCESS.2022.3207287>
- Murugesan, R., Mishra, E., & Krishnan, A. H. (2022). Forecasting agricultural commodities prices using deep learning-based models: Basic LSTM, bi-LSTM, stacked LSTM, CNN LSTM, and convolutional LSTM. *International Journal of Sustainable Agricultural Management and Informatics*, 8(3), 242–277. <https://doi.org/10.1504/IJSAMI.2022.125757>
- Paul, R. K., Yeasin, M., Kumar, P., Kumar, P., Balasubramanian, M., Roy, H. S., Paul, A. K., & Gupta, A. (2022). Machine learning techniques for forecasting agricultural prices: A case of brinjal in Odisha, India. *PLOS ONE*, 17(7), e0270553. <https://doi.org/10.1371/journal.pone.0270553>
- Pusdatin Kementan. (2023). Outlook Bawang Merah 2023.
- Salman, H. A., Kalakech, A., & Steiti, A. (2024). Random Forest Algorithm Overview. *Babylonian Journal of Machine Learning*, 2024, 69–79. <https://doi.org/10.58496/BJML/2024/007>
- Sun, F., Meng, X., Zhang, Y., Wang, Y., Jiang, H., & Liu, P. (2023). Agricultural Product Price Forecasting Methods: A Review. *Agriculture*, 13(9), Article 9. <https://doi.org/10.3390/agriculture13091671>
- Wandschneider, T., Andri, K. B., Ly, K., Puspadi, K., Gniffke, P., Harper, S., & Kristedi, T. (2013). Shallot Value Chain Study Executive Summary. Australian Center for International Agricultural Research.
- Wang, L., Zhou, X., Zhu, X., Dong, Z., & Guo, W. (2016). Estimation of biomass in wheat using random forest regression algorithm and remote sensing data. *The Crop Journal*, 4(3), 212–219. <https://doi.org/10.1016/j.cj.2016.01.008>
- Weng, Y., Wang, X., Hua, J., Wang, H., Kang, M., & Wang, F.-Y. (2019). Forecasting Horticultural Products Price Using ARIMA Model and Neural Network Based on a Large-Scale Data Set

Collected by Web Crawler. IEEE Transactions on Computational Social Systems, 6(3), 547–553.
<https://doi.org/10.1109/TCSS.2019.2914499>

Zhang, D., Chen, S., Liwen, L., & Qiang, X. (2020, January). Forecasting Agricultural Commodity Prices Using Model Selection Framework With Time Series Features and Forecast Horizons – DOAJ. <https://doaj.org/article/b9093ef10e8c4218afa7ced6ac611284>