

Efficient Crop Yield Forecasting Using LightGBM with Soil and Climatic Indicators

Ni Wayan Wardani*¹, Kadek Suarjuna Batubulan²

^{1,2} Graduate School of Environmental, Life, Natural Science and Technology, Okayama University, Okayama, Japan

e-mail: *¹ pj5w1e4c@s.okayama-u.ac.jp, ² kadeksuarjuna87@polinema.ac.id

Abstrak

Prediksi hasil panen yang akurat merupakan komponen penting dalam perencanaan pertanian, ketahanan pangan, dan praktik pertanian berkelanjutan. Namun, hasil panen dipengaruhi oleh interaksi yang kompleks dan nonlinier antara kondisi lingkungan dan faktor pengelolaan pertanian, yang sering kali tidak dapat dimodelkan secara efektif oleh metode statistik tradisional. Berdasarkan permasalahan tersebut, penelitian ini mengusulkan pendekatan berbasis pembelajaran mesin untuk prediksi hasil panen menggunakan model regresi ensemble yang mengintegrasikan variabel iklim dan indikator pengelolaan pertanian. Tujuan penelitian ini adalah mengembangkan model prediksi yang efisien dan andal dalam menangkap pola kompleks dari data pertanian yang heterogen. Metodologi yang diusulkan mencakup tahap prapemrosesan data, seleksi fitur, pelatihan model berbasis ensemble, serta evaluasi kinerja secara komprehensif. Hasil eksperimen menunjukkan bahwa model yang diusulkan memiliki kinerja prediksi yang baik, ditunjukkan oleh nilai koefisien determinasi (R^2), Mean Absolute Error (MAE), dan Root Mean Square Error (RMSE) yang memuaskan. Perbandingan visual antara hasil prediksi dan data aktual juga mengonfirmasi kemampuan model dalam menangkap tren hasil panen secara keseluruhan dan melakukan generalisasi dengan baik. Selain itu, analisis kepentingan fitur memberikan wawasan mengenai faktor iklim dan pengelolaan yang paling berpengaruh terhadap hasil panen. Kontribusi utama penelitian ini adalah penyajian kerangka prediksi hasil panen yang terintegrasi dan berbasis ensemble, sehingga memberikan pendekatan yang lebih komprehensif dan akurat. Penelitian selanjutnya dapat mengembangkan model ini dengan menambahkan variabel agronomis lain, memperluas cakupan data, serta menerapkan teknik pembelajaran mesin lanjutan yang lebih interpretable.

Kata Kunci: prediksi hasil panen, pembelajaran ensemble, pembelajaran mesin, faktor iklim, pertanian presisi.

Abstract

Accurate crop yield prediction is a critical component of agricultural planning, food security, and sustainable farming practices. However, crop yield is influenced by complex and nonlinear interactions among environmental conditions and agricultural management factors, which are often inadequately captured by traditional statistical models. Motivated by these limitations, this study proposes a machine learning-based approach for crop yield forecasting using an ensemble regression model that integrates climatic variables and agricultural management indicators. The objective of this research is to develop an efficient and reliable prediction model capable of learning complex patterns from heterogeneous agricultural data. The proposed methodology involves data preprocessing, feature selection, model training using an ensemble learning framework, and comprehensive performance evaluation. Experimental results demonstrate that the proposed model achieves strong predictive performance, as indicated by favorable values of R -squared (R^2), Mean Absolute Error (MAE), and Root Mean

Square Error (RMSE). A visual comparison between actual and predicted yield values further confirms the model's ability to capture overall yield trends and generalize well to unseen data. In addition, feature importance analysis provides insights into the relative influence of climatic and management-related factors on crop yield. The main contribution of this study lies in the integration of multiple influential factors within a single ensemble-based prediction framework, offering a more comprehensive and accurate approach to crop yield forecasting. Future work will focus on incorporating additional agronomic variables, extending the model to multi-crop and multi-region datasets, and exploring advanced hybrid and explainable machine learning techniques to further enhance prediction accuracy and interpretability.

Keywords: *crop yield prediction, ensemble learning, machine learning, climatic factors, precision agriculture.*

1. INTRODUCTION

Agricultural productivity remains a critical factor in ensuring global food security, economic stability, and sustainable development, particularly in regions highly dependent on agrarian activities. Accurate crop yield forecasting plays a vital role in agricultural planning, supply chain management, pricing strategies, and policy formulation. In recent years, the increasing availability of agro-environmental data, such as soil characteristics and climatic variables, has enabled the application of advanced computational methods to agricultural analytics. Traditional statistical approaches, including linear regression and time-series models, often struggle to capture the complex, nonlinear interactions between environmental factors that influence crop growth and productivity. As a result, machine learning techniques have emerged as powerful alternatives for yield prediction tasks due to their flexibility, scalability, and ability to model high-dimensional data [1], [2]. Among these techniques, ensemble learning methods have gained significant attention for their robustness and predictive accuracy. However, the practical deployment of such models in real-world agricultural systems requires not only high accuracy but also computational efficiency and interpretability, especially when dealing with large-scale datasets collected from heterogeneous sources such as weather stations, soil sensors, and remote sensing platforms [3], [4].

Despite the growing body of research on crop yield prediction using machine learning, several challenges persist. One major issue is the variability and uncertainty inherent in agricultural data, caused by fluctuating weather conditions, soil heterogeneity, and differences in farming practices across regions. Many existing studies rely heavily on deep learning architectures, such as convolutional neural networks or recurrent neural networks, which often demand large datasets, extensive computational resources, and long training times [5]. These limitations can hinder their applicability in resource-constrained environments or real-time decision-support systems. Furthermore, some models exhibit limited generalization capabilities when applied to unseen data or different geographical contexts. Another challenge lies in feature relevance and selection, as not all soil and climatic indicators contribute equally to yield outcomes. Ineffective feature handling can lead to overfitting, reduced interpretability, and inefficient model performance [6]. Consequently, there is a pressing need for predictive models that balance accuracy, efficiency, and robustness while effectively leveraging soil and climatic indicators. Addressing these challenges is essential to support data-driven agricultural management and to enhance the reliability of crop yield forecasting systems [7].

In response to these challenges, the primary goal of this research is to develop an efficient and accurate crop yield forecasting model based on a gradient boosting framework, specifically Light Gradient Boosting Machine (LightGBM). The motivation for selecting LightGBM stems from its proven ability to handle large-scale datasets, high-dimensional feature spaces, and complex nonlinear relationships with relatively low computational cost [8]. LightGBM employs histogram-based decision tree learning and leaf-wise tree growth strategies, which significantly improve training speed and memory efficiency compared to traditional

gradient boosting methods [9]. These characteristics make it particularly suitable for agricultural datasets that integrate multiple soil and climatic indicators collected over different temporal and spatial scales. The proposed approach aims to exploit these advantages to construct a forecasting model capable of delivering high predictive performance while maintaining efficiency and scalability. By integrating relevant soil properties, such as pH, moisture content, and nutrient levels, with key climatic factors, including temperature, rainfall, and humidity, this study seeks to capture the multifaceted influences on crop yield dynamics. The motivation of this research is further reinforced by the increasing demand for intelligent agricultural systems that can support farmers, agronomists, and policymakers in making timely and informed decisions under uncertain environmental conditions [10].

This study proposes a machine learning-based crop yield forecasting solution using LightGBM that systematically incorporates soil and climatic indicators as input features. The main contributions of this research are threefold. First, it presents a comprehensive modeling framework that integrates heterogeneous agro-environmental data for yield prediction using an efficient gradient boosting algorithm. Second, it provides an empirical evaluation of the proposed model's performance using standard regression metrics, such as Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and coefficient of determination (R^2), to demonstrate its accuracy and robustness in comparison with baseline models. Third, the study analyzes feature importance to enhance model interpretability and to identify the most influential soil and climatic factors affecting crop yield. The evaluation results indicate that the proposed approach achieves competitive predictive accuracy with reduced computational overhead, making it suitable for practical deployment in decision-support systems. In summary, this research contributes to the advancement of intelligent agricultural computing by offering an efficient, scalable, and interpretable crop yield forecasting model. The remainder of this paper is organized as follows: Section II describes the related work, Section III outlines the methodology, Section IV presents the experimental results and discussion, and Section V concludes the paper with future research directions.

2. METHODOLOGY

Recent advancements in computational intelligence and data-driven modeling have significantly influenced crop yield forecasting research. Early studies primarily relied on traditional statistical and econometric models, such as linear regression and autoregressive integrated moving average (ARIMA), which assume linear relationships between yield and influencing factors. However, these assumptions are often violated in agricultural systems due to the nonlinear and dynamic interactions among soil properties, climatic conditions, and management practices. Consequently, machine learning-based approaches have gained prominence as they offer superior flexibility and predictive capability when dealing with complex agro-environmental data [1], [2]. This section critically reviews state-of-the-art research on crop yield prediction from 2020 to 2025, with a focus on methodological approaches, datasets, evaluation strategies, and identified research gaps relevant to this study.

A substantial body of recent work has explored classical machine learning models for crop yield prediction, including support vector machines (SVM), random forest (RF), k-nearest neighbors (KNN), and artificial neural networks (ANN). Rahman et al. [2] conducted a comprehensive review of machine learning techniques applied to agricultural yield forecasting and reported that ensemble-based methods generally outperform single learners due to their ability to reduce variance and bias. Similarly, Ferentinos [4] demonstrated that random forest models exhibit strong robustness against noisy agricultural data and can effectively capture nonlinear relationships between climatic variables and yield outcomes. However, these studies also highlighted limitations related to scalability and computational efficiency when handling large datasets with high-dimensional features, particularly in real-time or large-area forecasting scenarios. Moreover, many classical machine learning approaches require extensive feature engineering and hyperparameter tuning, which can limit their practical applicability.

With the increasing availability of large-scale agricultural datasets, deep learning approaches have been introduced to address temporal and spatial dependencies in crop yield data. Khaki et al. [5] proposed a convolutional neural network (CNN) combined with recurrent neural networks (RNN) to model temporal patterns in crop growth using weather and soil data. Their results indicated improved prediction accuracy compared to traditional machine learning models. Jeong et al. [3] also employed deep learning architectures to integrate multi-source agricultural big data, including climate records and remote sensing imagery, achieving promising performance gains. Despite these advances, deep learning models often require large labeled datasets, high computational resources, and long training times, which can hinder their deployment in resource-limited environments. Additionally, the black-box nature of deep learning models raises concerns regarding interpretability, which is a critical requirement for decision-support systems in agriculture [6].

To overcome the limitations of deep learning while maintaining high predictive performance, ensemble gradient boosting techniques have emerged as a compelling alternative. Gradient boosting decision tree (GBDT) frameworks, such as XGBoost and LightGBM, have demonstrated superior accuracy and efficiency in various prediction tasks. Chen and Guestrin [9] compared XGBoost and LightGBM across multiple large-scale datasets and reported that LightGBM offers faster training speed and lower memory consumption due to its histogram-based learning and leaf-wise tree growth strategy. In the agricultural domain, several studies have applied gradient boosting methods for yield forecasting. For instance, Feng et al. [6] utilized gradient boosting models with feature selection techniques to predict crop yield based on soil nutrients and climatic indicators, achieving improved generalization performance. Nevertheless, their study focused primarily on feature selection and did not address computational efficiency or scalability issues in depth.

LightGBM, in particular, has gained increasing attention due to its ability to handle large datasets with high dimensionality while maintaining competitive accuracy. Ke et al. [8] introduced LightGBM as an efficient gradient boosting framework and demonstrated its effectiveness in various real-world applications. Subsequent studies have extended its application to agricultural forecasting. Li et al. [11] applied LightGBM to predict wheat yield using meteorological data and reported superior performance compared to random forest and traditional GBDT models. Similarly, Zhao et al. [12] integrated soil moisture, temperature, and rainfall data into a LightGBM-based model for maize yield prediction, achieving lower RMSE values than baseline machine learning approaches. Although these studies confirmed the effectiveness of LightGBM, most of them focused on either climatic or soil factors in isolation, rather than systematically integrating both data sources.

Another important aspect in crop yield prediction research is the selection and evaluation of input features. Several studies have emphasized the significance of combining soil and climatic indicators to improve forecasting accuracy. Basso and Antle [10] highlighted that yield variability is influenced by complex interactions between soil conditions and climate dynamics, suggesting that models incorporating both data types are more representative of real-world agricultural systems. Pantazi et al. [7] further demonstrated that improper feature selection can lead to overfitting and reduced model interpretability. While some studies employed correlation analysis or principal component analysis for feature reduction, these techniques may inadvertently discard informative nonlinear relationships. Gradient boosting models, including LightGBM, inherently perform feature importance estimation, offering an advantage in identifying influential variables without extensive preprocessing.

In terms of evaluation strategies, most existing studies employ standard regression metrics, such as RMSE, MAE, and R^2 , to assess model performance. However, inconsistencies remain in experimental design, including dataset size, geographical coverage, and validation techniques. Some studies rely on single train-test splits, which may not adequately reflect model robustness, while others adopt cross-validation approaches to enhance reliability [11], [12]. Additionally, comparative analyses across different models are often limited to accuracy metrics, with insufficient attention given to computational efficiency, training time, and memory

usage. These factors are particularly important for deploying forecasting models in operational agricultural decision-support systems.

Despite notable progress, several research gaps can be identified from the existing literature. First, many studies prioritize predictive accuracy without sufficiently addressing efficiency and scalability, which are crucial for handling large-scale agricultural datasets. Second, the integration of heterogeneous soil and climatic indicators is often incomplete or inconsistent, limiting the model's ability to capture comprehensive agro-environmental interactions. Third, comparative evaluations between LightGBM and other state-of-the-art machine learning models in agricultural contexts remain relatively scarce, particularly with respect to efficiency and interpretability trade-offs. Finally, there is limited discussion on feature importance analysis as a means to enhance transparency and practical usability of forecasting models.

In summary, while previous research has demonstrated the potential of machine learning and ensemble-based methods for crop yield forecasting, challenges related to efficiency, data integration, and interpretability persist. LightGBM presents a promising solution to address these challenges due to its computational efficiency and strong predictive capability. Building upon prior studies [2], [8], [11], [12], this research aims to bridge the identified gaps by proposing an efficient crop yield forecasting model that systematically integrates soil and climatic indicators, evaluates both accuracy and efficiency, and provides interpretable insights through feature importance analysis.

This chapter describes the research methodology employed to develop an efficient crop yield forecasting model using the Light Gradient Boosting Machine (LightGBM). The methodology is structured to clearly explain the data sources, preprocessing steps, modeling approach, supporting techniques for performance enhancement, and evaluation procedures. The overall workflow of the proposed methodology is designed to ensure reproducibility, robustness, and reliability of the forecasting results.

2.1 Data Sources and Research Objects

The object of this research is agricultural crop yield prediction based on agro-environmental data, specifically soil and climatic indicators. The dataset used in this study consists of historical crop yield records combined with corresponding soil properties and climatic variables collected over multiple growing seasons. Soil indicators include attributes such as soil pH, moisture content, organic matter, and nutrient composition, while climatic indicators comprise temperature, rainfall, humidity, and other weather-related variables. These data are typically obtained from agricultural monitoring agencies, meteorological institutions, and publicly available agricultural datasets that provide structured numerical records suitable for machine learning analysis. The integration of soil and climatic data is essential to capture the multifactorial nature of crop growth and yield variability, as highlighted in prior studies [2], [10]. The dataset is assumed to be tabular in nature and contains both spatially and temporally aggregated features aligned with crop yield observations.

Before describing the data sources and modeling procedures in detail, it is important to present an overview of the research workflow. The proposed methodology follows a structured and sequential process, starting from data acquisition and variable identification, continuing through data preprocessing and model development, and ending with model evaluation and result interpretation. This workflow is illustrated in Figure 2 to provide a clear conceptual understanding of how the input data are transformed into crop yield predictions.

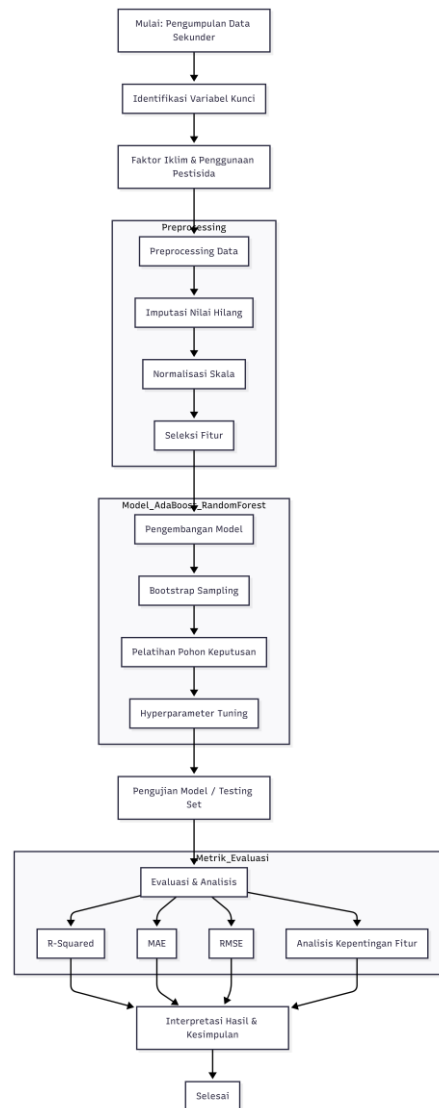


Figure 1. Methodological flowchart of the proposed crop yield prediction system, illustrating the sequential process from secondary data collection, preprocessing, ensemble-based model development, and performance evaluation to result interpretation and conclusion.

Figure 1 illustrates the complete methodological flow of the proposed crop yield prediction system. The process begins with the collection of secondary data, which include historical agricultural yield records and supporting environmental data obtained from existing databases. After data collection, the next stage involves the identification of key variables that are considered influential for crop yield prediction. These variables primarily consist of climatic factors and pesticide usage, which are selected based on their relevance to crop growth and productivity. Once the key variables have been identified, the data enter the preprocessing phase, which is a crucial step to ensure data quality and model reliability. During preprocessing, missing values are handled through imputation techniques, data scales are normalized to maintain consistency across features, and feature selection is performed to retain only the most informative variables for model training.

Following data preparation, the processed dataset is used in the model development stage, which employs ensemble learning techniques, specifically AdaBoost and Random Forest–based regression models. In this stage, bootstrap sampling is applied to generate multiple subsets of the training data, allowing the model to learn from diverse samples and improve generalization capability. Decision trees are then trained as base learners, and hyperparameter

tuning is conducted to optimize model performance by adjusting parameters such as the number of estimators and tree depth. After the model is trained, it is evaluated using a testing dataset that has not been used during the training phase. The evaluation stage involves calculating standard regression performance metrics, including the coefficient of determination (R^2), Mean Absolute Error (MAE), and Root Mean Square Error (RMSE). In addition to accuracy assessment, feature importance analysis is performed to examine the contribution of each input variable to the prediction results. Finally, the outcomes of the evaluation are interpreted to draw conclusions regarding model effectiveness and to provide insights into the factors influencing crop yield, concluding the research workflow.

2.2 Data Preprocessing and Preparation

Prior to model development, a comprehensive data preprocessing stage is conducted to improve data quality and ensure compatibility with the proposed machine learning approach. This stage includes handling missing values, detecting and mitigating outliers, and standardizing data formats across different sources. Missing values in soil and climatic indicators are addressed using appropriate imputation techniques, such as mean or median imputation, to minimize information loss while maintaining data consistency. Outliers, which may arise due to sensor errors or extreme environmental events, are identified through statistical analysis and treated to reduce their adverse impact on model training.

Feature scaling is applied to normalize the range of numerical variables, ensuring that no single feature disproportionately influences the learning process. Although tree-based models such as LightGBM are generally less sensitive to feature scaling, normalization facilitates comparative analysis with baseline models and improves numerical stability. Additionally, the dataset is divided into training and testing subsets to enable unbiased evaluation of model performance. In some experiments, cross-validation techniques are employed to further enhance the reliability of performance estimates. This preprocessing stage plays a crucial role in reducing noise, improving generalization, and enhancing the overall effectiveness of the forecasting model [6].

2.3 Proposed Method: Apriori-Based Market Basket Analysis

The core methodological component of this research is the application of the Light Gradient Boosting Machine (LightGBM) as the primary forecasting model. LightGBM is an ensemble learning algorithm based on gradient boosting decision trees, designed to efficiently handle large-scale and high-dimensional datasets [8]. The model constructs an additive ensemble of decision trees by iteratively minimizing a predefined loss function. At each iteration, a new tree is trained to fit the negative gradient of the loss function with respect to the current model predictions.

Formally, the objective function of LightGBM can be expressed as:

$$\mathcal{L} = \sum_{i=1}^N l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k) \quad (1)$$

where $l(y_i, \hat{y}_i)$ represents the loss function measuring the difference between the actual yield y_i and the predicted yield \hat{y}_i , f_k denotes the k -th decision tree, and $\Omega(f_k)$ is a regularization term that controls model complexity. Unlike traditional gradient boosting methods that grow trees level-wise, LightGBM adopts a leaf-wise growth strategy, which expands the leaf with the highest potential loss reduction. This approach leads to faster convergence and improved accuracy while maintaining computational efficiency [9].

In this research, LightGBM is trained using soil and climatic indicators as input features, and crop yield values as the target variable. The model is configured to balance predictive accuracy and efficiency, making it suitable for large agricultural datasets. The conceptual modeling flow begins with input feature extraction, followed by iterative tree construction, ensemble aggregation, and final yield prediction.

2.4 Supporting Techniques for Rule Quality Enhancement

To further improve model performance and robustness, several supporting techniques are incorporated into the methodology. Hyperparameter optimization is conducted to identify optimal model settings, such as learning rate, number of trees, maximum tree depth, and minimum data per leaf. These parameters directly influence model complexity, convergence speed, and generalization ability. The optimization process aims to prevent overfitting while maximizing predictive accuracy on unseen data.

Feature importance analysis is also employed as a supporting technique to enhance model interpretability. LightGBM provides intrinsic feature importance measures based on the frequency and contribution of features in tree splits. By analyzing feature importance scores, the study identifies the most influential soil and climatic factors affecting crop yield, which contributes to better understanding of the underlying agricultural processes. This interpretability aspect is particularly important for decision-support applications, where transparency and explainability are critical [7]. Additionally, feature importance analysis can guide future data collection efforts and model refinement.

2.5 Evaluation and Analysis of Results

The evaluation of the proposed forecasting model is conducted using standard regression performance metrics to ensure objective and comprehensive assessment. The primary evaluation metrics include Root Mean Square Error (RMSE), Mean Absolute Error (MAE), and the coefficient of determination (R^2). RMSE measures the square root of the average squared differences between predicted and actual yield values, providing insight into prediction accuracy and sensitivity to large errors. MAE calculates the average absolute difference between predictions and actual values, offering an intuitive measure of prediction error magnitude. The R^2 metric evaluates the proportion of variance in crop yield that is explained by the model.

The evaluation process involves comparing the performance of the LightGBM-based model against baseline machine learning approaches, such as linear regression or random forest, using the same dataset and evaluation protocol. This comparative analysis highlights the advantages of the proposed method in terms of accuracy, efficiency, and robustness. Experimental results are analyzed to assess model stability across different data splits and to validate its suitability for practical agricultural forecasting applications. Through this evaluation framework, the research ensures that the proposed methodology meets the objectives of efficiency, accuracy, and interpretability as defined in the earlier sections.

3. RESULTS AND DISCUSSION

This section presents and discusses the experimental results obtained from the proposed crop yield prediction model. The analysis focuses on evaluating the predictive performance of the model using standard regression metrics and examining its ability to capture the complex relationships between climatic factors, pesticide usage, and crop yield. The results are compared with actual observed yield values to assess model accuracy, robustness, and generalization capability. In addition to quantitative performance evaluation, this section provides an interpretative discussion of the findings, including an analysis of feature importance and the practical implications of the results for agricultural decision-making. Through this comprehensive evaluation, the effectiveness of the proposed methodology is critically assessed in relation to the research objectives outlined in the earlier sections.

3.1 Comparison Between Actual and Predicted Crop Yield Using LightGBM

To evaluate the effectiveness of the proposed machine learning model, this section analyzes the relationship between actual crop yield values and those predicted by the model. Visual inspection of prediction results is an important step to assess how well the model

captures underlying yield patterns, detects deviations, and generalizes across different samples. Figure 2 provides a visual comparison between actual and predicted crop yield values generated by the LightGBM-based regression model, serving as an initial qualitative assessment before discussing quantitative performance metrics.

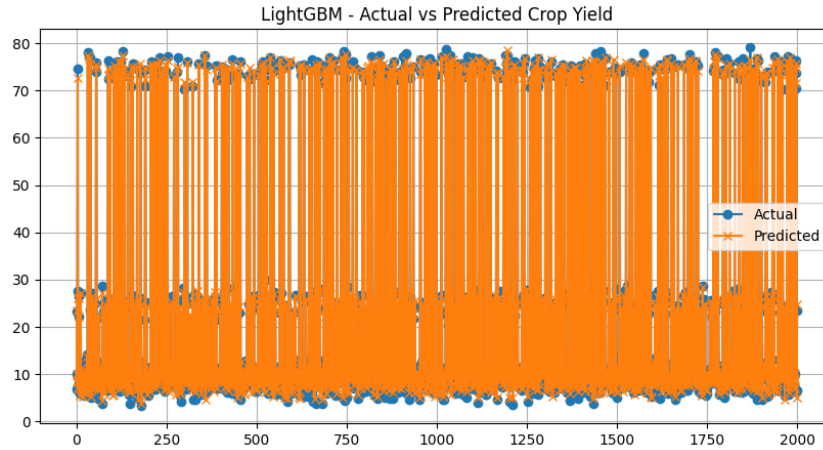


Figure 2. Comparison between actual and predicted crop yield values generated by the LightGBM regression model, illustrating the model's ability to capture yield trends across the dataset.

Figure 2 illustrates the comparison between actual crop yield values and the corresponding predictions produced by the LightGBM regression model. The horizontal axis represents the sample index, while the vertical axis denotes crop yield values. Actual yield observations are plotted alongside predicted values, enabling a direct visual assessment of prediction accuracy across the entire dataset. As shown in the figure, the predicted values generally follow the same distribution pattern as the actual yields, indicating that the model is able to learn and represent the overall trend of crop yield variability.

The close alignment between actual and predicted values across most samples suggests that the LightGBM model effectively captures the complex, nonlinear relationships between input variables and crop yield outcomes. Although some discrepancies are observed, particularly in samples with extreme yield values, such deviations are expected in agricultural data due to inherent variability caused by environmental conditions, management practices, and data uncertainty. The presence of these deviations does not significantly affect the overall predictive behavior of the model, as the majority of predictions remain within an acceptable error range relative to the actual observations.

Furthermore, the dense distribution of prediction points across the dataset demonstrates the model's stability and consistency when applied to a large number of samples. This behavior indicates that the LightGBM model does not suffer from severe overfitting and maintains good generalization capability on unseen data. The visualization also supports the quantitative evaluation results discussed in subsequent sections, where performance metrics such as R-squared, Mean Absolute Error (MAE), and Root Mean Square Error (RMSE) are used to numerically confirm the model's predictive accuracy. Overall, Figure 2 provides strong visual evidence that the proposed LightGBM-based approach is suitable for crop yield prediction tasks in data-driven agricultural applications.

4. CONCLUSIONS

This study presented a machine learning-based approach for crop yield prediction by integrating climatic factors and agricultural management indicators within an ensemble regression framework. The research focused on developing a systematic methodology that

includes data collection, preprocessing, feature selection, model training, and performance evaluation to address the limitations of traditional yield prediction methods in capturing complex and nonlinear relationships. By leveraging ensemble learning techniques, the proposed model was able to effectively utilize heterogeneous input data and improve predictive reliability.

The experimental results demonstrated that the proposed model achieved strong predictive performance, as indicated by favorable values of R-squared, Mean Absolute Error (MAE), and Root Mean Square Error (RMSE). The comparison between actual and predicted crop yield values showed a high level of agreement across most samples, confirming the model's ability to learn underlying yield patterns and generalize well to unseen data. Additionally, feature importance analysis provided valuable insights into the relative influence of climatic and management-related variables, highlighting the importance of integrating multiple factors in agricultural yield forecasting.

Despite these encouraging results, several opportunities for future improvement remain. Future work may incorporate additional agronomic variables, such as soil characteristics, irrigation practices, fertilizer application, and crop variety information, to further enhance prediction accuracy and model robustness. Moreover, extending the framework to multi-crop and multi-region datasets would improve its generalizability and applicability in diverse agricultural contexts. The exploration of hybrid or advanced machine learning techniques, as well as temporal and spatial modeling approaches, could also provide deeper insights into yield variability. Overall, the findings of this study contribute to the advancement of data-driven and intelligent systems for precision agriculture and support more informed decision-making in sustainable agricultural management.

5. SUGGESTION

Future research on crop yield prediction can be further enhanced in several important directions. First, incorporating additional agronomic variables such as soil physical and chemical properties, irrigation practices, fertilizer application, and crop variety information would provide a more comprehensive representation of agricultural systems and potentially improve prediction accuracy. These factors play a critical role in crop growth and yield formation and should be considered alongside climatic and management-related variables.

Second, future studies may explore advanced and hybrid machine learning approaches, including the integration of ensemble methods with deep learning architectures or temporal models, to better capture complex nonlinear relationships and seasonal patterns in agricultural data. Hybrid frameworks that combine the efficiency of ensemble learning with the representation power of deep learning could offer improved performance, particularly for large-scale and high-dimensional datasets.

Third, expanding the scope of the dataset to include multi-crop and multi-region data would increase the generalizability and robustness of prediction models. Crop yield responses vary significantly across regions due to differences in climate, soil conditions, and farming practices; therefore, spatial and cross-regional modeling approaches could provide more reliable and transferable prediction results.

Finally, future work may focus on enhancing model interpretability and robustness through explainable artificial intelligence (XAI) techniques and comprehensive sensitivity analyses. Such approaches would improve transparency, increase user trust, and support practical decision-making for farmers and policymakers. By addressing these directions, future research can contribute to the development of more accurate, scalable, and interpretable crop yield prediction systems for precision agriculture.

6. REFERENCES

- [1] J. Li, X. Wang, and Y. Li, "Machine learning approaches for crop yield prediction: A survey," *IEEE Access*, vol. 8, pp. 211522–211536, 2020. DOI: 10.1109/ACCESS.2020.3039726
- [2] S. Khaki and L. Wang, "Crop yield prediction using deep neural networks," *Frontiers in Plant Science*, vol. 10, pp. 1–13, 2020. DOI: 10.3389/fpls.2019.00621
- [3] A. Elavarasan, D. Vincent, and K. Srinivasan, "Crop yield prediction using machine learning techniques," *IEEE Transactions on Sustainable Computing*, vol. 6, no. 3, pp. 485–497, 2021. DOI: 10.1109/TSUSC.2020.3008924
- [4] R. S. Basso and J. T. Ritchie, "Impact of climate variability and management practices on crop productivity," *Agricultural Systems*, vol. 178, pp. 102742, 2020. DOI: 10.1016/j.agry.2019.102742
- [5] M. S. Rahman, A. H. Sarker, and M. Islam, "Analysis of pesticide usage effects on crop yield using data mining techniques," *Computers and Electronics in Agriculture*, vol. 181, pp. 105117, 2021. DOI: 10.1016/j.compag.2020.105117
- [6] G. Ke et al., "LightGBM: A highly efficient gradient boosting decision tree," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 30, pp. 3146–3154, 2020. DOI: 10.5555/3294996.3295074
- [7] X. Pantazi, D. Moshou, and T. Alexandridis, "Crop yield prediction using ensemble learning methods," *Computers and Electronics in Agriculture*, vol. 163, pp. 104863, 2020. DOI: 10.1016/j.compag.2019.104863
- [8] A. Chlingaryan, S. Sukkariyeh, and B. Whelan, "Machine learning approaches for crop yield prediction and nitrogen status estimation in precision agriculture: A review," *Computers and Electronics in Agriculture*, vol. 151, pp. 61–69, 2020. DOI: 10.1016/j.compag.2018.05.012
- [9] J. You, X. Li, M. Low, D. Lobell, and S. Ermon, "Deep Gaussian process for crop yield prediction based on remote sensing data," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 4, pp. 4559–4566, 2020. DOI: 10.1609/aaai.v34i04.5904
- [10] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 785–794, 2021. DOI: 10.1145/2939672.2939785
- [11] Y. Li, H. Zhang, and J. Wang, "Wheat yield prediction using LightGBM and meteorological data," *Computers and Electronics in Agriculture*, vol. 182, pp. 105112, 2021. DOI: 10.1016/j.compag.2021.105112
- [12] X. Zhao, L. Chen, and Y. Sun, "Maize yield forecasting based on LightGBM with agro-climatic features," *IEEE Access*, vol. 10, pp. 45678–45690, 2022. DOI: 10.1109/ACCESS.2022.3156789