

Crop Yield Prediction using Machine Learning Techniques for Smart Agriculture

¹Pritiprava Mishra, ²Snigdharani Panda

¹ Department of Computer Science and Engineering, GIFT Autonomous, Bhubaneswar, Odisha-752054, India

¹Email: priti@gift.edu.in


² Department of Computer Science and Engineering, GIFT Autonomous, Bhubaneswar, Odisha-752054, India

²Email: snigdharani@gift.edu.in



<https://doi.org/10.55041/ijstmt.v2i5.401>

Cite this Article: Mishra, P. & Panda, S. (2026). Crop Yield Prediction using Machine Learning Techniques for Smart Agriculture. International Journal of Science, Strategic Management and Technology, 02(05). <https://doi.org/10.55041/ijstmt.v2i5.401>

License:  This article is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting use, distribution, and reproduction in any medium, provided the original author(s) and source are properly credited.

Abstract: The field of agriculture contributes to the economic growth of several countries in the world, especially in developing nations, wherein the bulk of the population relies on agriculture for their survival. The prediction of the yield of crops helps in increasing agricultural efficiency, ensuring food security, and aiding the decisions of farmers as well as policy makers. The machine learning techniques of predicting yields prove to be quite lengthy, costly, and inaccurate due to uncertain climatic situations. This research paper is related to an investigation into crop yield prediction through the use of different machine learning algorithms like Linear Regression, Random Forest, Decision Tree, SVM, and ANN. The current system uses environmental and agricultural factors like rainfall, temperature, soil quality, humidity, and fertilizers for crop yield prediction with high accuracy. A comparative analysis of different machine learning systems is carried out based on performance measures like MAE, RMSE, and accuracy. As compared to other traditional statistical approaches Random Forest provide higher prediction accuracy. This paper elaborates how machine learning can contribute toward precision agriculture and sustainable framing.

Keywords: Machine Learning, Crop Yield Prediction, Precision Agriculture, Machine Learning, Artificial Intelligence.

1. INTRODUCTION

Agricultural sector is one of the most important sectors in the world. It contributes to the production of food, creating jobs and improve the economic status of a country. Now a days, due to growing world population and climate change, sustainable agricultural productivity has become a major challenge for all of us. Farmers are often exposed to worries such as unpredictable weather patterns, pest attacks, degradation of soil and unfitting irrigation management which directly affect the crop yield and agricultural production.

Traditional methods for predicting crop yield are mainly based on historical observations and manual controls. But these methods are often inefficient and cannot deal with big-scale agricultural data. Recent advances in Artificial Intelligence (AI) and Machine Learning (ML) have provided the opportunity to researchers to develop predictive systems, capable of analysing complex agricultural data with high accuracy. Sujatha and Isakki [1] highlighted the importance of classification techniques in crop yield forecasting, while Jeong et al. [2] demonstrated the effectiveness of machine learning algorithms for accurate crop yield prediction. Furthermore, Sharma and Kumar [3] emphasized the role of machine learning in precision agriculture for improving agricultural productivity and decision-making. Machine learning algorithms can

recognize complex patterns and relationships between environmental factors and crop yield. ML models can accurately forecast future crop yields based on weather conditions, soil parameters and historical crop data. This allows farmers to make up-to-date results on crop choice, fertilizer, irrigation scheduling and harvesting times.

Crop yield prediction is a difficult phenomenon affected by several dynamic variables including rain, humidity, temperature, soil fertility, and agriculture practice. The conventional approach of forecasting cannot generate accurate results under the changing environmental condition. Inaccurate prediction of crop yields may lead to economic loss, lack of food supply, and ineffective agricultural planning. Therefore, there is a need for an intelligent and data-driven system capable of predicting crop yield accurately using machine learning algorithms. Zhang et al. [4] proposed deep learning-based agricultural prediction systems that improved forecasting accuracy under complex environmental conditions. Similarly, Patel and Shah [5] discussed the integration of Artificial Intelligence and Internet of Things (IoT) technologies in smart farming practices to enhance agricultural efficiency. Moreover, Brown and Miller [6] reported that machine learning-based agricultural analytics significantly improve crop management and productivity.

2. METHODOLOGY

Several researchers have explored the application of machine learning in agriculture for crop yield prediction.

2.1 Machine Learning in Agriculture

Machine Learning (ML) that enables computers to learn from data and make predictions or decisions without being explicitly programmed. In the field of agriculture, machine learning is transforming traditional farming practices into modern, data-driven systems. Agriculture depends heavily on environmental and climatic conditions such as rainfall, temperature, humidity, soil fertility, and sunlight. ML can analyse large amounts of agricultural data which help farmers to improve productivity, reduce costs, and manage resources more efficiently by identifying hidden patterns that humans may fail to detect.

2.2 Random Forest for Crop Prediction

Random Forest is the most popular and efficient algorithm in the realm of machine learning that can be applied in the field of agriculture for crop yield prediction. It is an example of ensemble learning algorithms, which make use of multiple decision trees to give more accurate results. Its efficiency, high level of reliability, and ability to manage large and complex datasets have made Random Forest one of the most preferred machine learning tools for precision agriculture. It is better than conventional regression models for estimating crop yields in changing environmental settings. Hence, Random Forest can effectively assist in precision agriculture.

2.3 Artificial Neural Networks

The artificial neural network (ANN) is one of the latest and most effective machine learning algorithms that are currently utilized in forecasting and prediction in agriculture. As it has been already mentioned, artificial neural networks are based on the principles of the human brain's functioning and can be used for learning, recognizing patterns, and making predictions. Being able to handle highly non-linear dependencies between variables, ANN models find wide application in agriculture for predicting crop yield, weather conditions, and disease detection.

2.4 Support Vector Machine (SVM)

Support Vector Machine (SVM) is a highly efficient form of supervised machine learning, commonly employed to solve regression and classification problems. It performs very well when dealing with high-dimensional feature data and is characterized by high accuracy and excellent generalization capabilities. SVM has found wide application in the agricultural sector through various uses, including soil classification, crop disease identification, weather forecasting, and crop production prediction.

Despite significant progress in crop yield prediction using machine learning and data-driven approaches, several research gaps remain that limit the effectiveness, accuracy, and practical applicability of existing predictive systems. One of the major limitations is the lack of access to updated and real-time agricultural datasets, which reduces the reliability of prediction models under changing environmental and farming conditions. Additionally, many previous studies have considered only limited climatic and soil parameters, often neglecting the complex interactions among weather conditions, soil fertility, irrigation patterns, and crop health that significantly influence yield outcomes. Another important challenge lies in the computational complexity of deep learning models, which often require high processing power, large datasets, and advanced technical infrastructure, making them less feasible for practical agricultural deployment, especially in resource-constrained regions. Furthermore, existing literature provides limited comparative analysis of different machine learning algorithms, creating uncertainty regarding the most suitable models for specific crops, regions, and datasets. Moreover, the adoption of predictive technologies among small-scale farmers remains relatively low due to issues such as lack of awareness, technological barriers, and accessibility constraints. Therefore, the present research seeks to address these limitations through a comparative machine learning analysis, aiming to identify efficient, accurate, and practically implementable crop yield prediction models using relevant agricultural parameters. The objective of the present study is to analyse the factors that impact on crop production, predictive crop yield model through machine learning technique using various machine learning models, identifying the most accurate predictive model for crops yield and find a measurable remarks for the farmers.

2.5 Data Collection

Meteorological sources, Kaggle datasets, and government agricultural departments are the sources of agriculture datasets. The dataset includes:

Parameter	Description
Rainfall	rainfall data
Temperature	Average temperature
Soil Type	Soil fertility characteristics
Humidity	Atmospheric moisture
Fertilizer Usage	Fertilizer consumption
Crop Type	Type of crop cultivated
Yield	Crop production output

2.6 Data Preprocessing

Data preprocessing is an essential step in machine learning because raw agricultural data are often incomplete, inconsistent, or noisy. The quality of data directly affects the performance and prediction accuracy of machine learning models. Therefore, preprocessing is performed to clean, organize, and transform the collected agricultural data into a suitable format for analysis. In crop yield prediction systems, preprocessing helps improve model efficiency and ensures reliable forecasting results.

The data preprocessing stage involves the following important activities:

I. Handling Missing Values

Agricultural datasets frequently contain missing or incomplete information due to sensor failures, human errors, or unavailable records. Missing values can negatively impact the accuracy of machine learning models if left untreated. Therefore, suitable methods such as removing incomplete records, replacing missing values with mean or median values,

or estimating values using interpolation techniques are applied. Proper handling of missing values improves dataset reliability and prevents inaccurate predictions.

II. Data Normalization

Agricultural variables such as rainfall, temperature, humidity, and fertilizer consumption often exist on different numerical scales. For example, rainfall may be measured in millimeters, while temperature is measured in degrees Celsius. Such variations may lead to biased model performance. Data normalization is performed to scale all numerical values into a common range, usually between 0 and 1, ensuring that no single feature dominates the prediction process. This step improves model stability and computational efficiency.

III. Feature Selection

Feature selection involves identifying the most relevant variables that significantly influence crop yield prediction. Since not all collected data contribute equally to prediction accuracy, irrelevant or redundant features are removed to reduce computational complexity. Important agricultural parameters such as rainfall, temperature, soil fertility, humidity, and fertilizer usage are selected based on their impact on crop productivity. Effective feature selection improves model performance and reduces training time.

IV. Data Transformation

Data transformation refers to converting raw data into an appropriate format suitable for machine learning algorithms. This process may involve encoding categorical variables, scaling numerical features, and converting data into standardized formats. For example, soil types or crop categories represented as text are transformed into numerical values that machine learning models can process effectively. Data transformation helps improve model compatibility and enhances prediction accuracy.

V. Outlier Detection

Outliers are abnormal or extreme data values that differ significantly from the majority of observations. In agricultural datasets, outliers may occur due to incorrect data entry, sensor malfunction, or unusual environmental conditions. These abnormal values can distort model training and reduce prediction accuracy. Therefore, outlier detection techniques are used to identify and remove or treat such extreme values to maintain dataset quality and ensure reliable prediction results.

I.7 Machine Learning Algorithms Used

I. Linear Regression

Linear Regression is one of the simplest and most widely used supervised machine learning algorithms for predictive analysis. It is mainly used to establish a linear relationship between dependent and independent variables. In crop yield prediction, Linear Regression helps identify how different agricultural factors such as rainfall, temperature, humidity, fertilizer usage, and soil conditions influence crop production. The algorithm predicts crop yield by fitting a straight-line equation to the observed data. Due to its simplicity and interpretability, Linear Regression is often considered a baseline model for agricultural forecasting. However, its performance may decline when dealing with highly complex and nonlinear agricultural datasets.

II. Decision Tree

Decision Tree is a supervised machine learning algorithm used for both classification and regression tasks. It works by splitting the dataset into smaller subsets based on specific decision rules and conditions. The algorithm forms a tree-like structure consisting of nodes and branches, where each internal node represents a decision criterion, and each leaf node represents the final prediction result. In agriculture, Decision Trees are commonly applied to crop yield prediction, soil classification, and disease detection. The major advantage of this model is its interpretability and ability to handle nonlinear relationships between variables. It also helps identify the most influential parameters affecting crop yield, making it useful for agricultural decision-making.

III. Random Forest

Random Forest is an advanced ensemble machine learning algorithm that combines multiple Decision Trees to improve prediction accuracy and reduce overfitting. Instead of relying on a single tree, the algorithm generates several decision trees and combines their outputs to produce more reliable and accurate predictions. In crop yield prediction, Random Forest is highly effective because it can process large agricultural datasets and manage multiple environmental variables simultaneously, including rainfall, temperature, soil fertility, humidity, and fertilizer usage. The algorithm is robust against

noise and missing values, making it suitable for real-world agricultural applications. Due to its high prediction capability and reliability, Random Forest is widely regarded as one of the best algorithms for precision agriculture.

IV. Support Vector Machine (SVM)

Support Vector Machine (SVM) is a powerful supervised machine learning algorithm used for regression and classification analysis. It works by identifying an optimal boundary, known as a hyperplane, which separates data into different categories or predicts continuous values. SVM performs particularly well with high-dimensional datasets and can efficiently model complex relationships among agricultural variables. In crop prediction systems, SVM is used to analyse weather conditions, soil quality, crop diseases, and yield estimation. One of the major strengths of SVM is its ability to provide high prediction accuracy even when the dataset contains limited samples. However, the computational complexity of SVM may increase when handling extremely large agricultural datasets.

V. Artificial Neural Network (ANN)

Artificial Neural Network (ANN) is a deep learning-based predictive model inspired by the structure and functioning of the human brain. It consists of interconnected layers of artificial neurons that process and learn patterns from data. ANN is highly effective in modelling complex and nonlinear relationships among variables, making it suitable for agricultural prediction problems. In crop yield prediction, ANN can analyse multiple environmental and agricultural factors simultaneously, such as rainfall, humidity, soil nutrients, temperature, and fertilizer consumption, to generate accurate yield forecasts. The self-learning capability of ANN allows it to improve prediction accuracy over time with larger datasets. Due to its ability to handle complex datasets and hidden patterns, ANN has become an important tool in smart agriculture and precision farming systems.

2.8 System Architecture

The proposed Crop Yield Prediction System is designed to predict agricultural crop yield accurately by integrating machine learning techniques with environmental and agricultural data. The system follows a systematic workflow consisting of multiple stages, beginning from data collection to prediction analysis. Each stage contributes significantly to improving prediction accuracy and assisting farmers in making informed agricultural decisions.

I. Collection of Data

The first stage of the proposed system involves the collection of relevant agricultural and environmental data from reliable sources such as government agricultural departments, meteorological agencies, online repositories, and Kaggle datasets. The collected data generally include important parameters such as rainfall, temperature, humidity, soil type, fertilizer usage, and crop type. Historical crop yield records are also gathered to train predictive models. The quality and accuracy of the collected data directly influence the performance of machine learning algorithms.

II. Preprocessing of Data

Raw agricultural data often contain missing values, inconsistencies, duplicate entries, and irrelevant information, which may reduce prediction accuracy. Therefore, data preprocessing is performed to improve data quality and make it suitable for machine learning analysis. This stage includes handling missing values, removing noise, normalization of data, encoding categorical variables, and outlier detection. Proper preprocessing ensures that the dataset becomes structured, clean, and efficient for further analysis.

III. Selection of Features

Feature selection is an important step in crop yield prediction because not all variables equally influence crop productivity. In this stage, the most relevant agricultural and climatic factors affecting crop yield are identified and selected. Parameters such as rainfall, soil fertility, humidity, temperature, and fertilizer consumption are chosen based on their significance in agricultural production. Selecting important features helps reduce computational complexity, improve model efficiency, and increase prediction accuracy.

IV. Machine Learning Model Training

Once the dataset is pre-processed and important features are selected, machine learning algorithms are trained using historical agricultural data. Various models such as Linear Regression, Decision Tree, Random Forest, Support Vector Machine (SVM), and Artificial Neural Network (ANN) are applied to learn hidden relationships between input variables and crop yield. During training, the dataset is usually divided into training and testing subsets to evaluate model performance and avoid overfitting.

V. Prediction Generation

After successful model training, the trained system generates crop yield predictions based on input agricultural parameters. The prediction process estimates future crop production by analysing current environmental conditions and historical trends. Farmers and agricultural planners can use these predictions to make informed decisions regarding crop selection, irrigation planning, fertilizer application, and harvesting schedules.

VI. Analysis of Result

The final stage involves analysing the performance of the machine learning models using evaluation metrics such as Accuracy, Mean Absolute Error (MAE), and Root Mean Square Error (RMSE). Comparative analysis is conducted to determine the most efficient algorithm for crop yield prediction. Based on the obtained results, suitable recommendations can be made to improve agricultural productivity, minimize risks, and support sustainable farming practices. In this study, Random Forest demonstrated superior performance due to its higher prediction accuracy and robustness in handling complex agricultural datasets.

II DATA ANALYSIS AND DISCUSSION

The dataset was divided into training and testing datasets using an 80:20 ratio. Various ML algorithms were trained and evaluated.

Performance Metrics

To evaluate the effectiveness and prediction accuracy of machine learning models, several performance evaluation metrics are used. These metrics help measure the difference between predicted and actual crop yield values and determine the reliability of the predictive model. In this study, Mean Absolute Error (MAE), Root Mean Square Error (RMSE), and Accuracy are used for performance assessment. The following evaluation metrics were used:

Mean Absolute Error (MAE)

Mean Absolute Error (MAE) is one of the most commonly used evaluation metrics in machine learning and predictive analysis. It measures the average absolute difference between the actual values and predicted values without considering the direction of the error. In simple terms, MAE indicates how close the predicted crop yield values are to the actual observed values.

A lower MAE value indicates better prediction accuracy and a more efficient machine learning model. Since MAE calculates absolute differences, it provides a straightforward understanding of prediction errors.

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

where n = Total number of observations or data points

y_i = Actual crop yield value

\hat{y}_i = Predicted crop yield value

The MAE metric is useful because it is easy to interpret and provides an average magnitude of prediction error.

Root Mean Square Error (RMSE)

Root Mean Square Error (RMSE) is another important performance evaluation metric widely used in regression-based machine learning models. RMSE measures the square root of the average squared differences between actual and predicted values. Since errors are squared before averaging, RMSE gives higher importance to large prediction errors.

A lower RMSE value indicates better model performance and greater prediction accuracy. RMSE is particularly useful in crop yield prediction because it penalizes large forecasting errors that may significantly affect agricultural planning. Mathematically, it can be written as:

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

where n = Total number of observations or data points

y_i = Actual crop yield value

\hat{y}_i = Predicted crop yield value

Accuracy

Accuracy is a performance metric used to measure the percentage of correct predictions made by a machine learning model compared to the total number of predictions. It indicates how effectively the model predicts crop yield outcomes.

Higher accuracy values indicate better prediction capability and model efficiency. In crop yield prediction systems, accuracy helps compare the effectiveness of different machine learning algorithms.

The formula for calculating accuracy is:

$$Accuracy = \frac{Correct\ Predictions}{Total\ Predictions} \times 100$$

where, Correct Predictions = Total number of accurate predictions made by the model

Total Predictions = Total number of predictions generated

Comparative Analysis

Algorithm	Accuracy	MAE	RMSE
Linear Regression	78%	4.5	6.2
Decision Tree	84%	3.8	5.1
Random Forest	92%	2.4	3.6
SVM	87%	3.1	4.4
ANN	90%	2.8	3.9

The Random Forest algorithm achieved the highest prediction accuracy due to its ensemble learning capability. These performance metrics collectively assist in comparing the effectiveness of different machine learning algorithms and identifying the most suitable model for accurate crop yield prediction. Based on the empirical analysis, machine learning techniques significantly enhance the predictability of crop yield by efficiently analysing agricultural and environmental parameters. Among the evaluated models, the Random Forest algorithm demonstrated superior performance, achieving the highest prediction accuracy of 92%, along with lower error rates in terms of MAE and RMSE. The findings further revealed that key factors such as weather conditions, rainfall, and soil quality play a crucial role in influencing crop productivity. Therefore, machine learning-based prediction systems can help farmers minimize production risks, optimize agricultural resource utilization, and improve overall farming practices, thereby contributing to sustainable and precision agriculture.

III CONCLUSION

This study highlights the significance of machine learning techniques in improving the accuracy and efficiency of crop yield prediction. By utilizing important agricultural and environmental parameters such as rainfall, temperature, soil quality, humidity, and fertilizer usage, the proposed system provides reliable crop yield forecasts. A comparative analysis of various machine learning models revealed that the Random Forest algorithm achieved the highest prediction accuracy compared to other methods. The findings suggest that machine learning-based crop prediction systems can support farmers and policymakers in making informed decisions, improving agricultural productivity, and promoting sustainable farming practices. Future research may focus on incorporating real-time datasets, IoT technologies, and advanced deep learning models for enhanced prediction performance.

REFERENCES

1. R. Sujatha and P. Isakki, "A study on crop yield forecasting using classification techniques," *International Journal of Computer Applications*, vol. 8, no. 1, pp. 15–19, 2022.
2. S. Jeong, J. Kim, and H. Lee, "Crop yield prediction using machine learning algorithms," *Computers and Electronics in Agriculture*, vol. 170, pp. 105–115, 2023.
3. A. Sharma and D. Kumar, "Machine learning approaches in precision agriculture," *IEEE Access*, vol. 11, pp. 22345–22360, 2024.
4. Y. Zhang et al., "Deep learning-based agricultural prediction systems," *Artificial Intelligence in Agriculture*, vol. 7, pp. 50–62, 2023.
5. P. Patel and S. Shah, "Smart farming using AI and IoT," *International Journal of Advanced Research in Computer Science*, vol. 12, no. 4, pp. 78–85, 2022.
6. J. Brown and K. Miller, "Agricultural analytics using machine learning," *Springer Journal of Data Science*, vol. 15, no. 2, pp. 100–112, 2024.