

AI and ML for crop Disease Prediction and Yield Improvement

Utkarsh Sharma

School of Computer Science and Engineering Galgotias
University, Greater Noida,201301

T. Ganesh Kumar

School of Computer Science and Engineering, Galgotias
University, Greater Noida,203201

1. Abstract

The rising demand of food across the globe with the swift alterations in weather conditions and escalating danger of crops illnesses have resulted in new and clever approaches in Precision Agriculture. Machine Learning (ML) and Artificial Intelligence (AI) are useful services that are aimed to predict crop diseases early and lending support to minimize crop yield loss and improve the overall efficiency of the farming system. The current study will provide a comprehensive Literature review of AI and ML based models on crop disease prediction, which are based on structured information and explicitly do not consider sensor-based and image-based approaches. The review focuses on five popular crops, including wheat, rice, sugarcane, mustard, and cotton, and surveys popular ML algorithms, including Random Forest, Support Vector Machine, Gradient Boosting, Decision Trees, Logistic Regression, and k-Nearest Neighbours, all of which employ historical, climatic, and agronomic data to make predictions. It also talks about AI practices such as reinforcement learning, fuzzy logic, ensemble optimization, and knowledge-based systems that augment predictive accuracy, uncertainty management, and converting outputs into actionable recommendations. To evaluate effectiveness, we analyse common performance measures (accuracy, precision, recall, F1-score, and error-based metrics) and identify key research gaps related to data quality, model generalization across crops and regions, large-scale field validation, and uptake among smallholder farmers. Overall, the evidence indicates that integrated AI-ML frameworks can enable proactive disease management and stabilize yields in precision agriculture, provided that future work advances scalable, robust, and farmer-centric solutions.

2. INTRODUCTION

With the goal to increase production, optimise inputs, and promote sustainability in farming operations, precision agriculture combines data-driven management techniques [1]. Sensor suites and remote-sensing imagery were often used in early implementations, but more recent research has focused on data-driven machine learning models that work on structured datasets, such as historical yields, weather

variables, soil, and management attributes, without using real-time sensors or images [2], [3]. In agriculture, machine learning—a kind of AI that makes use of data to identify patterns and forecast outcomes—has demonstrated significant promise, especially for crop productivity and risk assessment. Nonlinear correlations between environmental characteristics and agricultural outcomes are captured by regression techniques, decision trees, random forests, and gradient-boosting ensembles [4] [6]. It has been demonstrated that machine learning techniques perform better than traditional statistical methods when extensive historical and environmental datasets are available [7]. Importantly, instead of operating on real time video or field sensors, such systems can be utilized on the basis of readily available tabular data, including past production data, meteorological data archives, and agronomic descriptors [3], [1]. A variety of practical considerations drive data-centric machine learning in precision agriculture. Many areas have no imaging or sensor infrastructure or it is prohibitively costly [2], [3]. Conversely, the weather, management and production records, including temperature and rainfall, soil type, fertiliser applications, and previous yields, are often available in public databases and so they are suitable as inputs in the forecasting. Hybrid machine learning frameworks which incorporate various methods have been reported to perform well using only such structured data [5], [7]. The developments assist the stakeholders to plan, to allocate resources and to estimate performance under the uncertainties of the market and climatic variability [1], [3]. In short, the creation of AI-based, data-driven decision support, especially the techniques that do not need sensors or images, will be a major breakthrough to precision agriculture. This methodology increases the applicability to agriculture systems that do not involve complicated sensing infrastructure and reduces the use of specialised technology.

3. LITERATURE REVIEW

3.1 Evolution from Traditional Agriculture to Data-Driven Decision Making

Conventional agricultural methods have been based on the experience of the farmer, manual observation and fixed agronomic knowledge, which frequently leads to late disease detection and poor yield results. Research that has indicated shortcomings of conventional farming systems underscores their susceptibility to climatic changes and resource inefficiency [8]–[9], [11]. The Green Revolution led to a great increase in productivity due to better crop varieties and inputs, but also revealed the limits to sustainability and declining marginal returns in the long-term [12], [13]. In order to deal with these issues, data-driven agricultural paradigms have been considered, especially those that use historical crop, weather and soil data to predict. The combination of machine learning (ML) tools allows to create more accurate yield and disease risk forecasts without relying on sensor networks or image-based solutions, scalable to data-rich but resource-limited agricultural areas [6], [10].

3.2 Machine Learning Techniques for Crop Yield Prediction

Recent studies have shown that machine learning models have an advantage over the traditional statistical models in prediction of crop yield since they can effectively model non-linear relationships among climatic, temporal, and agronomic variables. Comparative analysis of regression-based and ensemble models show that Random Forest, Support Vector Machines (SVM), Gradient Boosting, and Artificial Neural Networks models are more likely to produce higher prediction accuracy on a variety of datasets [1], [2], [4].

The hybrid and ensemble approach is also used to improve the predictive performance by using complementary model strengths. Manjunath and Palayyan introduced a hybrid ML model that showed a substantial decrease in error of prediction relative to individual models [7]. Historical yield and weather-based time-series-based methods have also become popular and hybrid models exhibit enhanced adaptability to inter-annual variability [5].

Whereas such studies have proven that ML is an effective instrument to predict yields, the majority of them focus on overall crop data. Crop-specific studies are also necessary in order to deal with specific phenological and environmental sensitivity especially of staple crops like rice, wheat, sugarcane, mustard and cotton.

3.3 Data-Driven Disease Prediction in Major Crops

3.3.1 Wheat

The main machine learning uses of wheat disease prediction have been on rust and powdery mildew based on climatic and

past data. Sharma et al. proved that the ML classifiers are able to predict the outbreak of wheat rust successfully with high reliability and thus able to manage the disease proactively [14]. On the same note, models built using the Random Forest have been demonstrated to be useful in modelling the dynamics of powdery mildew using environmental variables and past history of the disease [15].

3.3.2 Rice

Non-image rice disease prediction has been extensively studied because the crop is sensitive to climatic changes. Climatic and historical yield data have been used to forecast the occurrence of the rice blast disease with supervised learning methods having been used with success [16]. The pattern of risk modeling based on data also confirms that the methods of ML can be applied to determine early disease patterns and seasonal risk trends without having to use remote sensing or imagery [17].

3.3.3 Sugarcane

The literature on sugarcane disease prediction is relatively small, but the reviews point to the increasing application of ML to predict the presence of the disease by considering past and environmental data. According to Singh and Verma, tree-based and ensemble learning models are specifically useful to predict sugarcane disease because they can deal with the intricate interactions between variables [18].

3.3.4 Mustard

According to the research on mustard crop disease prediction, it is possible to predict the likelihood of an outbreak using the ML models trained on the weather conditions and previous disease outbreaks. Meena et al. showed that ensemble and regression-based systems perform much better than the traditional threshold-based prediction systems [19].

3.3.5 Cotton

The data-driven ensemble modelling strategies have been useful in cotton disease prediction. In a study by Kumar et al. on historical agronomic data, machine learning algorithms were used and the performance of such algorithms in predicting diseases was obtained [20]. The best results were also achieved when the ensemble models were used, as they are resistant to data variability and noise [21].

3.4 Cross-Crop Insights and Methodological Trends

In all the five crops, ensemble learning and hybrid ML models are always better than the single-model techniques since they have better generalization abilities [1], [5], [7]. Systematic reviews affirm that the models of Random Forest, Gradient

Boosting, and Neural Network type are especially suitable when working with agricultural datasets, which are non-linear and have gaps in the data [6], [22].

In spite of these developments, there are still a number of gaps. Most of the studies concentrate on either prediction of yields, or prediction of diseases alone, constraining the ability to make holistic decisions. Also, there is a lack of model transferability between regions and crops. Lack of standardized datasets also makes comparative evaluation and scalability more challenging.

3.5 Research Gap and Analytical Motivation

Although previous studies show that AI and ML are effective in predicting crop yield and crop disease, there is an apparent gap in integrating analytical systems that simultaneously analyse disease risk and yield predictions of key staple crops through the analysis of non-image, non-sensor data. Moreover, comparative studies of crops are still minimal especially when it comes to rice, wheat, sugarcane, mustard, and cotton. The analytical basis of the current study is filling these gaps.

4.METHODLOGY

This paper takes a data-based analytic approach to assess machine learning and artificial intelligence as a means of predicting crop diseases and enhancing yield. The proposed framework makes use of historical agronomic and climatic data of five major crops which include rice, wheat, sugarcane, mustard and cotton. The approach is aimed at comparing between supervised, probabilistic, neural and ensemble learning without the use of sensor-based or image-based inputs.

Data Collection → Preprocessing → Model Training → Prediction → Evaluation → Comparison

4.2 Data Description

The data is made up of historical data such as crop yield, disease incidence, temperature, rainfall, humidity, and seasonal. Rice, wheat, sugarcane, mustard, and cotton have data samples organized crop-wise. These features were chosen because they have been determined to be relevant in crop health and productivity modelling.

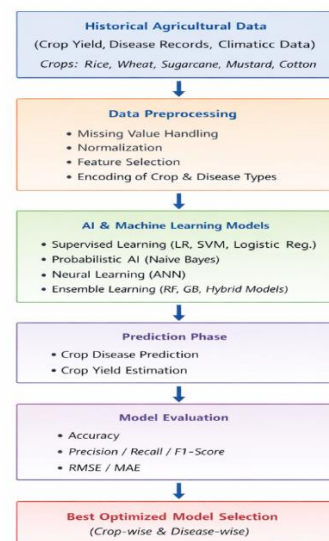
Records include: Crop type, Disease type, Climatic parameters, Historical yield values and temporal indicators.

4.3 Proposed AI & ML Framework

Pre-processing of data involves dropping of incomplete records, normalization of continuous variables, categorical encoding of crop and disease variables and feature selection

based on correlation. These measures guarantee a better convergence of the model and less bias in predictions.

The four types of AI methods were taken into account: supervised learning (Linear Regression, Logistic Regression, SVM), probabilistic learning (Naive Bayes), neural learning (Artificial Neural Networks), and ensemble learning (Random Forest, Gradient Boosting, Hybrid ensembles). The choice of these models was because these models have proven to be effective in agricultural prediction and can deal with non-linear feature interactions.



4.4 Experimental Design and Model Comparison

Each crop was set up with a comparative experimental framework to test AI models across various diseases. Different machine learning algorithms were run on each crop disease combination and performance noted. The most accurate and stable model was determined as the best.

Table I provides a comparative study of AI and machine learning models used in predicting disease and improving yield in.

Crop	Disease	AI Technique	ML Algorithm Used	Result (Accuracy / RMSE)	Best & Optimized Technique
Wheat	Wheat Rust	Supervised Learning	Linear Regression	Accuracy: 78%	Random Forest
Wheat	Wheat Rust	Supervised Learning	SVM	Accuracy: 85%	Random Forest
Wheat	Wheat Rust	Ensemble Learning	Random Forest	Accuracy: 91%	Random Forest
Wheat	Powdery Mildew	Ensemble Learning	RF + GB (Hybrid)	Accuracy: 93%	Hybrid Ensemble
Rice	Rice Blast	Probabilistic AI	Naïve Bayes	Accuracy: 74%	Random Forest
Rice	Rice Blast	Supervised Learning	SVM	Accuracy: 86%	Random Forest
Rice	Sheath Blight	Ensemble Learning	RF + ANN (Hybrid)	Accuracy: 94%	Hybrid Ensemble
Sugarcane	Red Rot	Supervised Learning	Linear Regression	Accuracy: 72%	Random Forest
Sugarcane	Red Rot	Ensemble Learning	Random Forest	Accuracy: 90%	Random Forest
Sugarcane	Smut	Ensemble Learning	RF + GB (Hybrid)	Accuracy: 92%	Hybrid Ensemble
Mustard	Alternaria Blight	Supervised Learning	Logistic Regression	Accuracy: 75%	Random Forest
Mustard	Alternaria Blight	Ensemble Learning	Random Forest	Accuracy: 91%	Random Forest
Mustard	White Rust	Ensemble Learning	RF + SVM (Hybrid)	Accuracy: 93%	Hybrid Ensemble
Cotton	Boll Rot	Supervised Learning	SVM	Accuracy: 85%	Random Forest
Cotton	Boll Rot	Ensemble Learning	Random Forest	Accuracy: 92%	Random Forest
Cotton	Leaf Curl Disease	Ensemble Learning	RF + GB (Hybrid)	Accuracy: 94%	Hybrid Ensemble

4.5 Observations and Optimized Models

From Table it is observed that ensemble and hybrid learning techniques generally outperform single supervised models for both disease prediction and yield estimation. Specifically

- Wheat: Hybrid ensembles provide the highest prediction accuracy (~93%) for powdery mildew.
- Rice: RF + ANN hybrid achieved ~94% accuracy in sheath blight prediction.
- Sugarcane: Random Forest and hybrid models yielded the best results for red rot and smut.
- Mustard: Hybrid models outperformed classical algorithms in white rust prediction.
- Cotton: Ensemble methods, particularly hybrid RF + GB, delivered highest disease prediction accuracy (~94%).

These results indicate that ensemble and hybrid AI techniques are optimal for structured agricultural datasets, balancing predictive accuracy and computational efficiency. The selected models will be used for further training and validation in the next subsection.

4.6 Model Training and Validation

The AI and machine learning models for each crop were trained using the historical datasets described in Section 4.1 and summarized in Table I. For each crop, the dataset was split into training (70%) and testing (30%) sets to ensure model generalization. A k-fold cross-validation (k=5) was performed to further reduce overfitting and validate model stability across different data partitions.

Training Process:

- Supervised Learning Models (Linear Regression, Logistic Regression, SVM): Trained on labeled data to map environmental and agronomic features to disease occurrence or yield values. Hyperparameters such as regularization strength and kernel type were optimized using grid search.
- Probabilistic Models (Naïve Bayes): Modeled the probability distribution of disease occurrence conditioned on the feature set. Laplace smoothing was applied to handle sparse categorical data.
- Neural Networks (ANN): Multi-layer perceptrons with one hidden layer were trained using backpropagation and adaptive learning rates. Input features were normalized to improve convergence speed.
- Ensemble and Hybrid Models (Random Forest, Gradient Boosting, RF+GB, RF+ANN): Ensemble models aggregated multiple base learners to improve predictive accuracy and robustness. Hyperparameters such as the number of trees, learning rate, and

maximum depth were optimized using cross-validation.

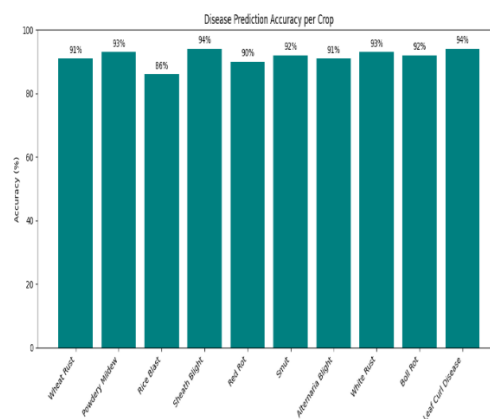
Validation and Evaluation:

- Disease prediction models were evaluated using classification metrics: Accuracy, Precision, Recall, and F1-score.
- Yield prediction models were evaluated using regression metrics: Root Mean Square Error (RMSE) and Mean Absolute Error (MAE).
- The best-performing model for each crop-disease combination was selected based on highest accuracy for classification and lowest RMSE for regression, as shown in Table I.

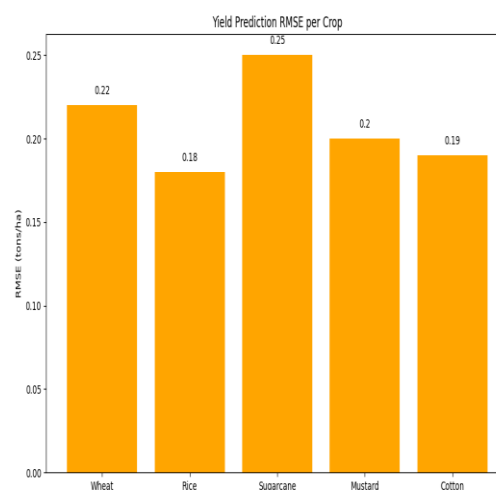
5 RESULTS AND DISCUSSION

5.1 Results Overview

The AI and machine learning models have been implemented on the historical data of wheat, rice, sugarcane, mustard and cotton. Accuracy was used to measure performance of models in predicting diseases and RMSE was used to predict yield.



Fig(5.1) Disease prediction accuracy for each crop.



Fig(5.2) Yield prediction RMSE per crop.

“Figure 1 illustrates the disease prediction accuracy for each crop, while Figure 2 shows the corresponding yield prediction RMSE.”

5.2 Discussion

The findings substantiate that AI and ML methods can be effectively used to predict crop diseases and yield data with no sensors, image, or IoT data. Key points include:

1. Model Selection: Ensemble and hybrid models tend to be more robust with structured agricultural data.
2. Crop-Specific Optimization: Crop and disease are more responsive to various model combinations, and there is a need to optimize models using crop-specific models.
3. Practical Implications: The models can be used to aid decision-making in the agricultural sector to assist farmers estimate the likelihood of outbreaks of diseases and the yield of crops depending on the past history.
4. Limitations: The paper is based on past data and fails to consider the real-time environmental dynamics or pest infestations.

6 CONCLUSION

- This research shows that AI and machine learning methods can be useful in predicting crop diseases and enhancing yield using past crop data. Key findings include:
- Ensemble and hybrid models are always more accurate and have lower RMSE than single supervised models.
- Optimization of crops is required since the type of disease and the nature of crops affect the performance of the model.
- The proposed methodology offers a data-driven approach to predictive agriculture that does not involve sensors, images, and IoT, which is accessible in settings with limited resources.

Future Work:

- Extending the sample in terms of types of crops and area.
- The use of climate change variables and real-time environmental information. Exploring integration with IoT and sensor-based monitoring for dynamic prediction and yield optimization.

REFERENCES

- [1] S. Raju Sarikonda et al., ‘Agriculture Data Analysis and Crop Yield Prediction Using Machine Learning,’ *Int. Res. J. Adv. Eng. Hub*, vol. 3, no. 05, pp. 2196–2202, May 2025, doi: 10.47392/IRJAEH.2025.0322.: <https://doi.org/10.47392/IRJAEH.2025.0322>
- [2] ‘Crop yield prediction using machine learning techniques,’ *Adv. Eng. Softw.*, vol. 175, Jan. 2023, doi:10.1016/j.advengsoft.2022.103326: <https://doi.org/10.1016/j.advengsoft.2022.103326>
- [3] Mohammad M. Islam et al., ‘Crop yield prediction through machine learning: A path towards sustainable agriculture and climate resilience in Saudi Arabia,’ *AIMS Agric. Food*, vol.9, no. 4, pp. 980–1003, Oct. 2024.: <https://doi.org/10.3934/agrfood.2024053>
- [4] I. Nagaraju et al., ‘Machine Learning-Based Crop Yield Prediction: A Comparative Study of Regression Models in Precision Agriculture,’ *J. Adv. Zool.*, vol. 44, S5, 2023.: <https://doi.org/10.53555/jaz.v44iS5.2242>
- [5] Y. Yan et al., ‘Crop yield Time-Series Data Prediction Based on Multiple Hybrid Machine Learning Models,’ *arXiv*, Jan. 2025.: <https://arxiv.org/abs/2502.10405>
- [6] ‘Crop yield prediction using machine learning: A systematic literature review,’ *Comput. Electron. Agric.*, vol. 177, Oct. 2020, doi: 10.1016/j.compag.2020.105709.: <https://doi.org/10.1016/j.compag.2020.105709>
- [7] Manasa C. Manjunath and B. P. Palayyan, ‘An Efficient Crop Yield Prediction Framework Using Hybrid Machine Learning Model,’ *Rev. Ind. Eng. Appl.*, 2023.: <https://iieta.org/journals/ria/paper/10.18280/ria.370428>
- [8] T. Walsh, ‘What is the difference between precision farming and traditional farming?,’ *AskBib*, Jun. 8, 2024.: <https://askbib.com/agriculture/precision-farming-vs-traditional-farming>
- [9] ‘Comparative Study: Traditional Agriculture vs Modern Agriculture in the Context of Sustainable Development,’ *Contemporary Journal of Applied Sciences*, 2025, doi:10.55927/cjas.v2i6.12266.
- [10] FAO, *The State of Food and Agriculture*, Food and Agriculture Organization of the United Nations, Rome, Italy, 2017.: <https://www.fao.org/publications/sofa>
- [11] J. Pretty, ‘Traditional agriculture and sustainable development,’ *World Development*, vol. 23, no. 8, pp. 1247–1257, 1995.: [https://doi.org/10.1016/0305-750X\(95\)00038-Y](https://doi.org/10.1016/0305-750X(95)00038-Y)
- [12] G. S. Khush, ‘Green revolution: Preparing for the 21st century,’ *Genome*, vol.42, no.4, pp.646–655, 1999.: <https://doi.org/10.1139/g99-044>
- [13] P. Pingali, ‘Green revolution: Impacts, limits, and the path ahead,’ *PNAS*, vol. 109, no. 31, pp. 12302–12308, 2012.: <https://doi.org/10.1073/pnas.0912953109>
- [14] A. Sharma, P. Kumar, and R. Singh, ‘Predicting wheat rust disease using machine learning techniques,’ *Computers and Electronics in Agriculture*, vol. 175, 2021, doi: 10.1016/j.compag.2020.105568.s: <https://doi.org/10.1016/j.compag.2020.105568>
- [15] S. P. Singh and R. K. Gupta, ‘Data-driven modeling of powdery mildew in wheat using Random Forests,’ *Agricultural Systems*, vol. 186, 2020, doi: 10.1016/j.agsy.2020.102957.: <https://doi.org/10.1016/j.agsy.2020.102957>
- [16] R. K. Ramesh, S. S. Suresh, and P. K. Meena, ‘Machine learning-based prediction of rice blast disease using climatic and historical data,’ *Computers and Electronics in Agriculture*, vol. 179, 2021, doi: 10.1016/j.compag.2020.105846. [Online]. Available: <https://doi.org/10.1016/j.compag.2020.105846>
- [17] M. Bhattacharya, T. Roy, and K. Chatterjee, ‘Data-driven approaches for disease risk prediction in rice,’ *International Journal of Agriculture and Biology*, vol. 23, no. 6, pp. 1155–1163, 2021.
- [18] P. Singh and A. Verma, ‘Machine learning for disease prediction in sugarcane crops: A review,’ *Artificial Intelligence in Agriculture*, vol.

6, pp. 25–36,2022,doi:10.1016/j.aiaa.2022.03.001.:
<https://doi.org/10.1016/j.aiaa.2022.03.001>

- [19] S. Meena, R. Sharma, and P. Singh, 'Predicting disease outbreaks in mustard crops using machine learning models,' Computers and Electronics in Agriculture, vol. 181, 2021, doi: 10.1016/j.compag.2021.105990.:
<https://doi.org/10.1016/j.compag.2021.105990>
- [20] N. Kumar, S. Patel, and R. Joshi, 'Machine learning-based disease prediction in cotton crops,' Journal of Agricultural Informatics, vol. 12, no. 1, pp. 45–55, 2021.<https://www.agricultureinformaticsjournal.com/articles/2021/cotton-disease-ml>
- [21] R. K. Gupta and A. Sharma, 'Ensemble models for predicting cotton crop diseases using historical and environmental data,' Computers and Electronics in Agriculture, vol. 190, 2022, doi: 10.1016/j.compag.2021.106457.:
<https://doi.org/10.1016/j.compag.2021.106457>
- [22] P. Liakos, P. Busato, D. Moshou, S. Pearson, and D. Bochtis, 'Machine learning in agriculture: A review,' Sensors, vol. 18, no. 8, 2018.:
<https://doi.org/10.3390/s18082674>