# A Crop Recommendation System Using Ensemble Learning Approach for Sustainable Farming

Pragyat Jyoti Baruah[1], Brajnandan Prasad[2], Arun Jyoti Nath[3] , Arnab Paul[4*], Soham Paik[5], Anand Raj[6], Vishwakarma Shempyu Shridayal[7]

[1,2,4,5,6,7] Department of Computer Science and Engineering, Assam University Silchar, India

[3] Department of Ecology and Environmental Sciences, Assam University Silchar, India

pragyat.jyoti.baruah@aus.ac.in[1], brajnandanprasad21@gmail.com[2], arun.jyoti.nath@aus.ac.in[3], arnab.paul@aus.ac.in[4], soham.paik@aus.ac.in[5], anand201712@gmail.com[6], shempyu@gmail.com[7]

**Abstract:**

Sustainable farming uses a precise crop recommendation model to optimize yield and resource use. Integrating ensembled machine learning with soil and environmental data, this study applies Principal Component Analysis to identify influential factors among 22 crops grown across the world. An ensembled model combining Random Forest, Naive Bayes and Decision Tree produces 99.24% accuracy, outperforming individual models. A website and a smartphone-based application is built on the proposed model for real-world use by farmers. This data-driven approach supports robust recommendations and enhances agricultural decision-making. This study demonstrates the role of artificial intelligence in advancing resilient and resource-efficient farming, thus providing a practical tool for sustainable crop selection and management.

**Keywords**: *Precise Agriculture, Feature selection, Classification, Soil data, Environmental data, Smartphone application.*

## 1. Introduction

Agriculture is an ancient industry crucial for sustaining the world's population. With the help of technological integration and modernization, it has grown to its maximum efficiency, engaging more participants and increasing overall quality standards[1]. In recent years, the application of machinery and technological development has substituted most of the labouring work in agricultural practices to a great extent, resulting in the overall increase in quality and efficiency[2]. India is among the developing countries with agriculture as its backbone source of revenue[3]. Agriculture produces about 14% of the GDP but significantly influences the Indian economy[4]. The United Nations' Food and Agriculture Organization (FAO) says that almost 33% of all food harvested for human use is lost annually due to numerous factors [3]. One of such factors is the absence of a fully dependable crop-recommendation model.

India is a country with six-seasons[5], which allows it to grow various types of crops throughout the year. Its primary crops include rice, wheat, and potato. The staple crops of India, like rice can be grown during three seasons[6]. In this study, we forecast the yield of 22 crop varieties (rice, maize, potato, wheat, etc.) using soil and environmental information (i.e. rainfall, humidity, pH, temperature, N, P, and K) from an open sourced dataset[7]. Climate change is one of the key drivers that may influence crop yields. The fluctuating environmental conditions, especially global warming and climate variability, have a detrimental effect on the

future of agriculture[8]. Due to such climate changes, developing an automated crop prediction mechanism using the environmental parameters is essential.
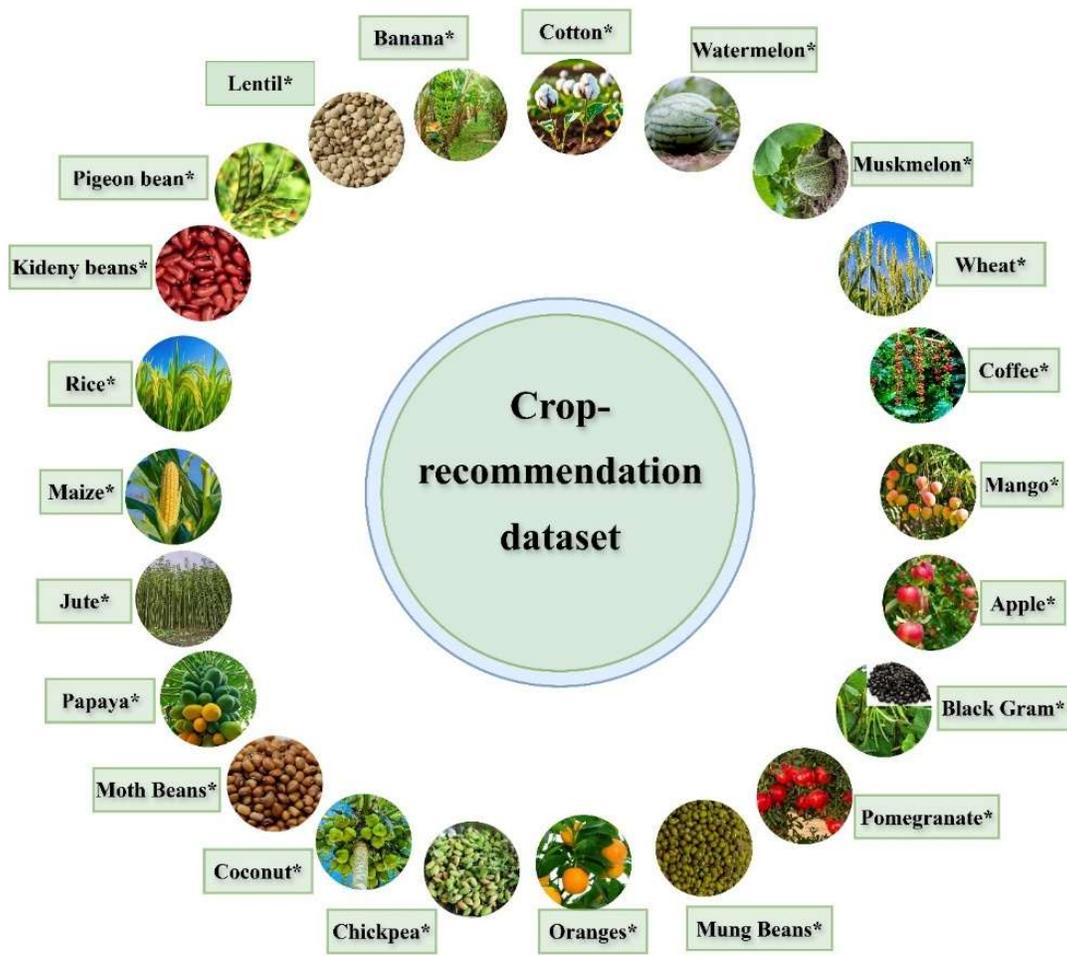
In the past, several studies attempted to employ a crop prediction model. Hasan *et al.*[8] performed prediction using various Machine Learning (ML) methods and compared the performance of an ensembled-learning method with other ML models like Support Vector Regression (SVR), Naive Bayes (NB), Catboost. They used $R^2$, MSE (Mean Square Error), root MSE, mean absolute error for result evaluation. The model in the study[8] could predict 5 crops (Aus rice, Aman rice, Boro rice, potato, and Wheat). Garanayak *et al.*[9] used RF, Logistic Regression (LR), SVR, Decision Tree regression, and polynomial regression, for crop anticipation and achieved 3.6% improvement over its previous works. The model proposed by Garanayak *et al.*[9] could recommend 5 crops (rice, ragi, gram, potato, and onion). Later, Thilakarathne *et al.*[10], employed five ML algorithm viz. Extreme Gradient Boosting, K-Nearest Neighbours (KNN), Support Vector Machine (SVM), Random Forest (RF), Decision Tree (DT) for finding the best model for crop prediction using a cloud-based service. Madhuri *et al.*[11] proposed an Improved Deep Belief Network (IDBN) model infused with Ranger Optimizer and Gaussian Restricted Boltzmann Machines for crop recommendation. Their analysis suggested that their model performs better than the traditional Deep Belief Networks for the crop prediction task across 4 crops (finger millet, sugarcane, maize, and rice). Mokarrama *et al.*[12] developed a real-time crop recommendation system named RSF (Recommendation System for Farmers) for recommending 4 crops. Their model operates by acquiring a farm's current coordinates and then recommending crops using climate and ecological data prevailing in real-time.

This paper proposes a study to promote sustainable agriculture and enhance food security amidst growing global challenges by developing an automated crop recommendation system. The approach integrates multiple environmental and soil parameters with machine learning models to predict suitable crops using a smartphone app and a website.

## 2. Dataset Description

**Table 1:** Representative samples of the dataset for the study

| N (mg/kg) | Rainfall (mm) | P (mg/kg) | Humidity (%) | K (mg/kg) | Temperature (°C) | pH | Crop |
|---|---|---|---|---|---|---|---|
| 19 | 54.73 | 55 | 87.81 | 20 | 27.43 | 7.19 | Mung Beans |
| 199 | 53.66 | 25 | 80.92 | 51 | 26.47 | 6.28 | Watermelon |
| 40 | 88.55 | 72 | 16.99 | 77 | 17.02 | 7.49 | Chickpea |
| 2 | 90.10 | 40 | 47.55 | 27 | 29.74 | 5.95 | Mango |
| 2 | 11.97 | 24 | 91.64 | 38 | 24.54 | 5.92 | Pomegranate |
| 3 | 32.68 | 49 | 64.71 | 18 | 27.91 | 3.69 | Moth Beans |
| 90 | 202.94 | 42 | 82.00 | 43 | 20.88 | 6.50 | Rice |
| 115 | 28.08 | 17 | 94.12 | 55 | 27.58 | 6.78 | Muskmelon |
| 56 | 71.89 | 79 | 63.20 | 15 | 29.48 | 7.45 | Black Gram |

* represents 100 data samples

**Figure 1:** Distribution of crop instances across the dataset

The dataset for this study is collected from Kaggle[7] data repository. Table 1 illustrates the sample data of our dataset. Figure 1 shows the data dispersion across crop species present in our dataset. It also shows that the dataset used for this study is perfectly balanced, thus preventing bias that might occur with un-balanced data while pursuing multi-class prediction using ML algorithms.
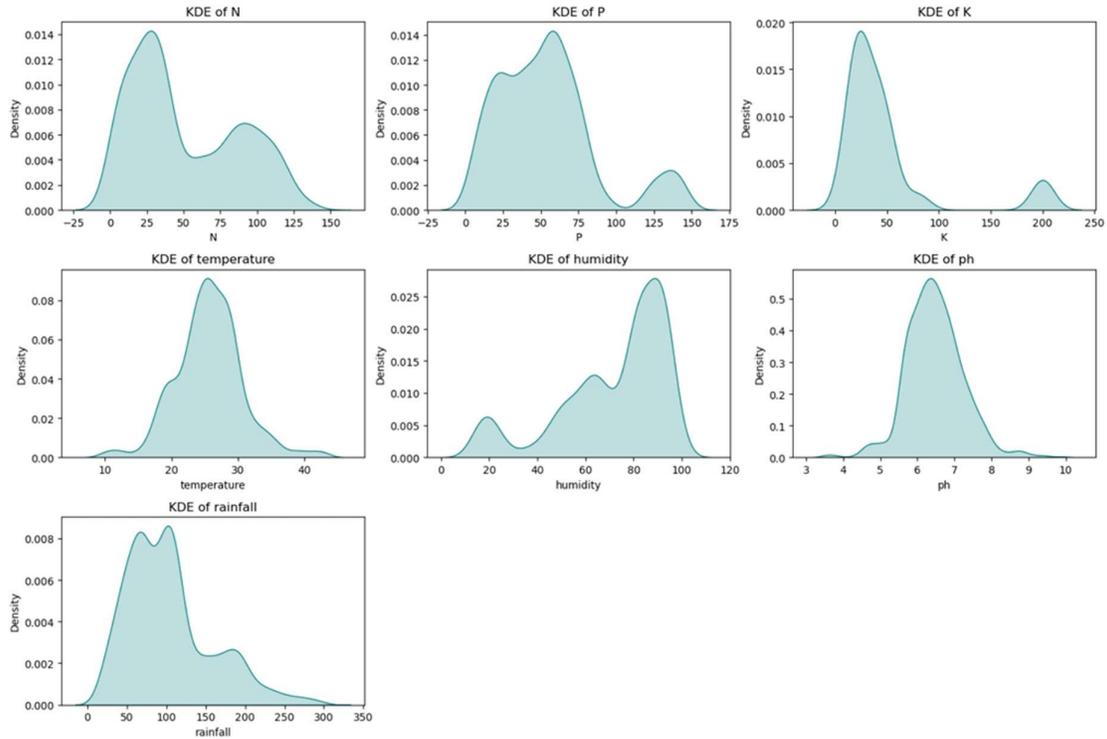
**Figure 2:** Kernel Density Estimation (KDE) of our dataset

Figure 2 provides the Kernel Density Estimation (KDE) graphs which show the probability distributions of seven major features— humidity, N, P, temperature, K, pH humidity, and rainfall present in our dataset. The figure plots exhibit unique patterns across all the variables. N and P variables exhibit bimodal patterns peaking around 20–40 and 60–80. Whereas, K exhibits strong concentration in the 0–50 region with a secondary peak in higher values, showing high variability. Temperature shows a unimodal distribution with a peak at 25–30°C, which indicates that most values are within an optimal range for plant growth. Humidity shows a skewed multimodal distribution with a large peak at 80–90%, while pH shows a near-normal distribution with a peak at 6.5, consistent with optimal soil conditions. Rainfall also shows a right-skewed distribution with a 50–100 mm peak, reflecting normal seasonal rainfall conditions. Figure 2 shows important characteristics about the natural variation and distribution of factors determining crop recommendations in our dataset.

## 3. Research Methodology

In this study, 7 machine learning models viz. LR, SVM, KNN, DT, Naïve Bayes (NB), RF, and an ensemble learning model are used to build a crop recommendation system. These models are evaluated using accuracy, precision, recall, and F1-score metrics along with error metrices (i.e. MAE, MSE, RMSE and $R^2$ score). This study uses Principal Component Analysis (PCA) to select the most relevant features/variables impacting the prediction of the crops. This is performed to reduce the model's complexity during its prediction stage. Figure 3 gives the flowchart of the methodology followed in this study.
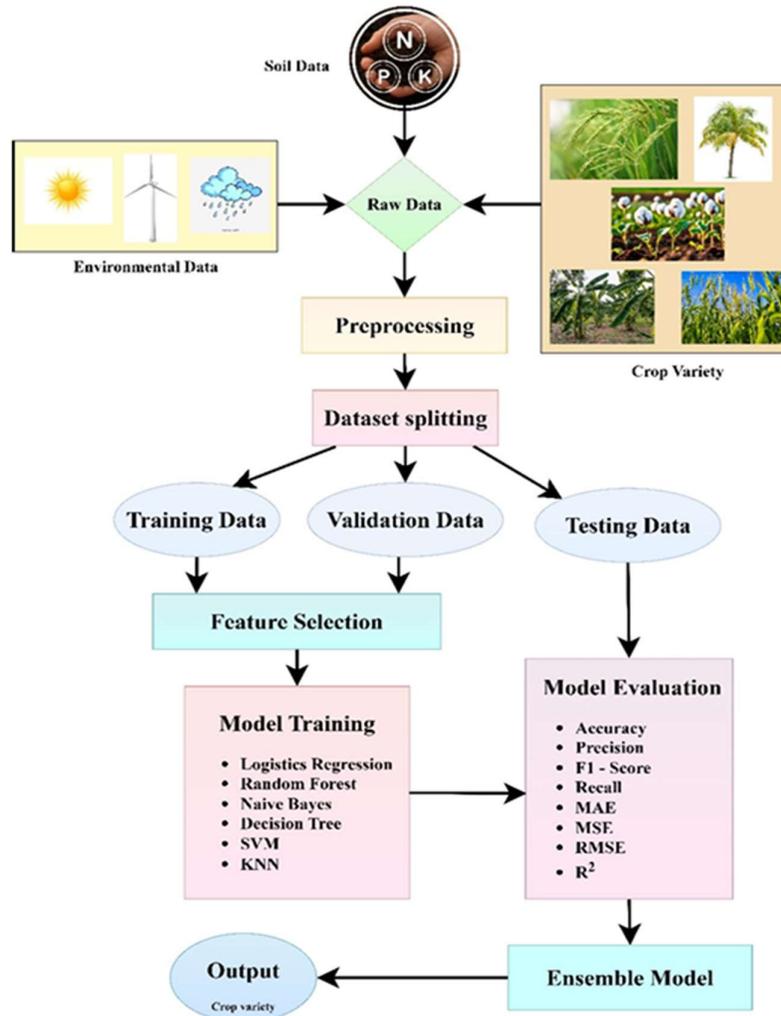
**Figure 3:** Workflow diagram of the proposed crop prediction model

### 3.1. Data Preprocessing

In preprocessing, two methods, viz. data cleaning and data transformation are performed on the raw dataset. The raw data is checked for any abnormalities like- missing values and improper formatting. Then, data cleaning is performed to handle missing-values by replacing the missing value with the mean of remaining sample data. Additionally, data-transformation is done by converting the cleaned data into suitable format.

### 3.2. Dataset splitting

The cleaned dataset is splitted into 3 sub-sets in this stage. For training, 70% data from the dataset is considered. Similarly, 15% of the dataset was considered for validation, and the remaining 15% for testing.

### 3.3. Feature Selection

This step is performed to select the best features from the set of variables present in the dataset (i.e. humidity, N, P, temperature, K, pH humidity, and rainfall). For this task,

dimensionality reduction technique named Principal Component Analysis (PCA) is used. The scatter plot in Figure 4 represents the distribution of the 22 crops varieties across the dataset based on the 7 variables, using two principal components (Principal Component 1 and 2). This plot is useful in interpreting data separability and applicability for classification task. From the image, it can be observed that the cluster of 3 crops that are produced using environmental and soil data, viz. pigeon peas, black gram and apple are significantly distinguishable from the other 19 crops varieties. However, almost all the crop varieties in the image show overlapping nature in the graph. For these characteristics, the dataset needs the use of an efficient machine learning models for the task of real-world prediction.
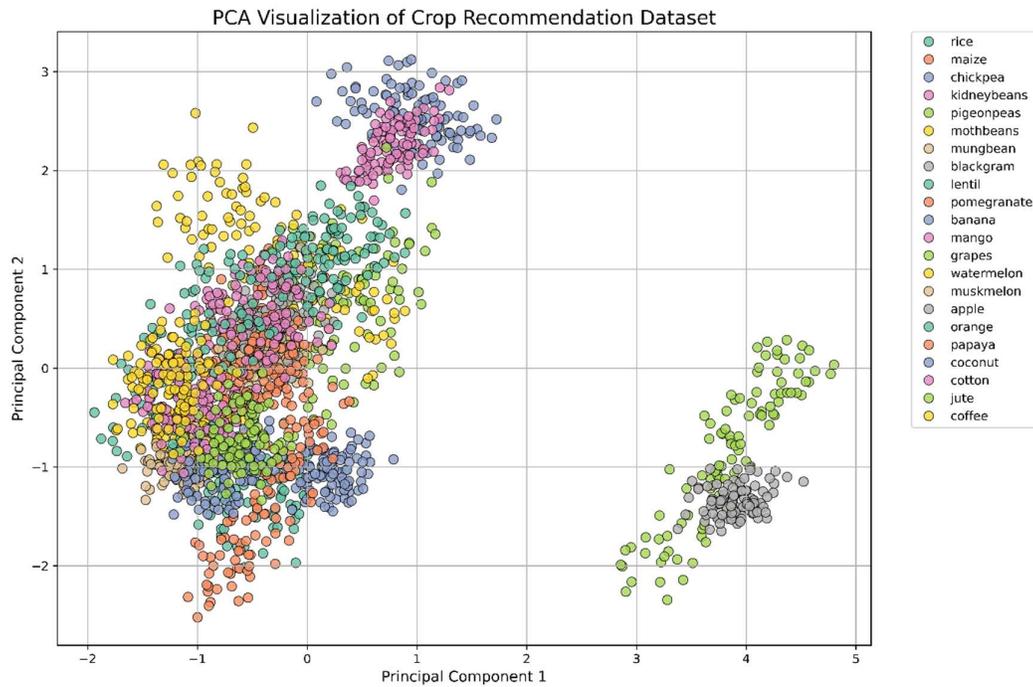


**Figure 4:** Graphical Representation of Principal Component Analysis (PCA)
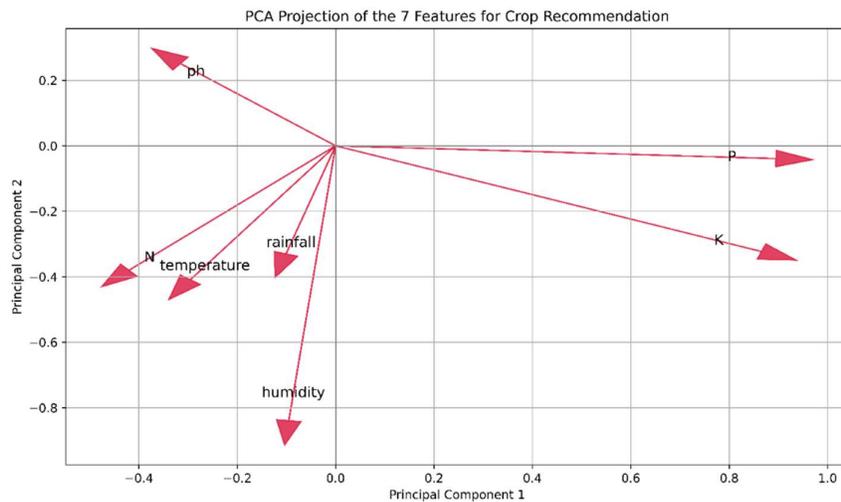


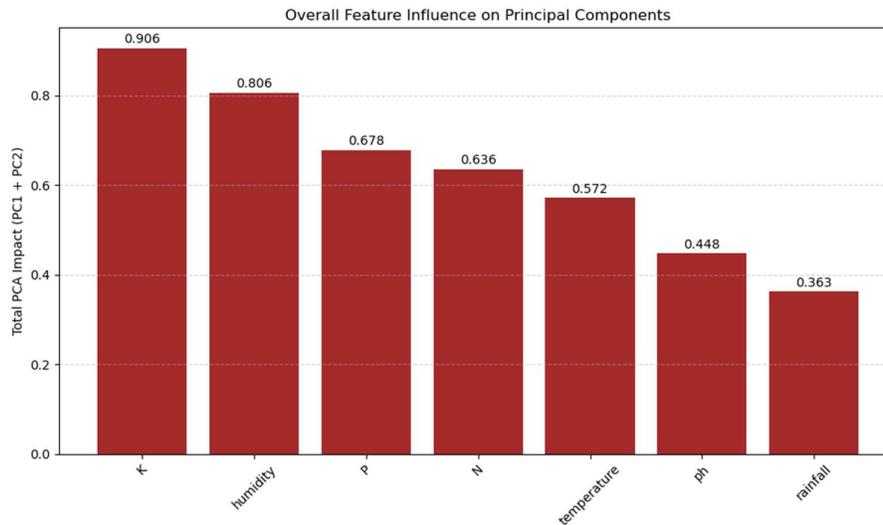**Figure 5:** PCA projections for all the features present in a dataset

**Figure 6:** This graph illustrates how attributes affect crop prediction

Figure 5 and Figure 6 indicates the impact of all the 7 features present in the dataset. This visualization identifies the most critical features to capture the variance structure in the dataset so that the feature selection and interpretation is performed for the purpose of dimensionality reduction and clustering. It can be observed that the 3 features, viz. K, humidity, and P are the most influential features with impact scores of 90.60%, 80.60 %, 67.80% respectively. For the further task of crop prediction using ML methods, the dataset variables are reduced to only K, humidity, and P.

### 3.4. Model Training

This stage involves training 6 popular ML methods for evaluation. These models were trained using the data after the process of feature selection and reduction of the dataset. The models used are:

**Random forest**, which consists of an ensemble of decision trees. This method predicts the crop by voting the highest prediction of the crop produced by each decision tree using averaging and regression.

**Logistic Regression** is a supervised model which uses sigmoid function to predict the probable crop. The value of the sigmoid function lies between 0 and 1.

**Decision tree** is a supervised-learning algorithm which contains two major nodes: the decision node, and the leaf node. It performs probability by deciding the best split node using feature selection measures of two types of Information gain, Gini Index.

**Support Vector Machines (SVMs)** are supervised algorithms for classification and regression problems. They perform classification using hyperplanes and support vectors.

**K-Nearest Neighbours (KNN)** is a supervised-learning algorithm developed for classification tasks. It performs classification by producing clusters across the labels present in the dataset.

**Naive Bayes** is a supervised probabilistic algorithm that uses Bayes theorem to train the model for predictions.

### 3.5. Model Evaluation

This study utilizes 8 evaluation matrices for evaluation: **Precision** is the measure of exact correctness of the model. It is a useful parameter when false-positive have large impact. **Recall** measures the true positivity rates of the predicted values of a model. It is also termed as sensitivity. **F1 Score** is the measure of harmonic mean of precision and recall of a model. It combines the recall and precision values; thus, it is an important matrix when an unbalanced dataset is involved. **Accuracy** measures a model's overall correctness (i.e., the percentage of total correct prediction). The **Mean squared error (MSE)** is the difference between actual and anticipated values; a lower MSE denotes better model performance. The **Mean Absolute Error (MAE)** is the average absolute amount of the difference between the actual and anticipated values in the data sets; its value is also inversely proportion to the performance of the model. The **Root Mean Square Error (RMSE)** is the square root of the average of the squared discrepancies between the predicted labels and those actually assigned. The **Coefficient-of-Determination/$R^2$ value** measures the amount of variance in the dependent variable which is explainable by the independent variables. It is an important matrix as it indicates how well the model fits the data.

### 3.6. Ensemble Model

Ensemble Learning model combines cluster of other individual models to obtain a better prediction result than individual models. It is performed to increase the robustness and accuracy of the model. In this study, hard-voting was performed on the predictions of the top 3 individual models (based on the accuracy values) to obtain the predicted output of the ensemble model. In hard voting, the class label is predicted from the majority votes of the predicted outputs of the individual models.

## 4. Results and Discussion

**Table 2:** Evaluation Results of the individual models used for this study

| Model | Precision | Recall | MAE | F1 Score | MSE | RMSE | $R^2$ score | Accuracy (%) |
|---|---|---|---|---|---|---|---|---|
| Naïve Bayes | 0.9926 | 0.9919 | 0.0727 | 0.9929 | 0.0171 | 0.8239 | 0.9881 | 99.11 |
| Random Forest | 0.9923 | 0.9909 | 0.1091 | 0.9907 | 0.0214 | 1.0090 | 0.9774 | 99.09 |
| Decision Tree | 0.9867 | 0.9864 | 0.1485 | 0.9863 | 0.0286 | 0.9954 | 0.9715 | 98.64 |
| SVM | 0.9766 | 0.9742 | 0.2621 | 0.9738 | 0.0432 | 1.1355 | 0.9343 | 97.42 |
| KNN | 0.9651 | 0.9606 | 0.3212 | 0.9602 | 0.0547 | 0.7006 | 0.9416 | 96.06 |
| Logic Regression | 0.9604 | 0.9591 | 0.3727 | 0.9589 | 0.0568 | 0.7568 | 0.8992 | 95.91 |

Table 2 shows the testing results of the 6 individual ML algorithms used to predict crops based on the dataset. In this study, RF, NB, and decision tree performed better than the other models, with Naive Bayes achieving 99.11% accuracy (best-case). Similarly, random forest and decision tree also achieved an accuracy of 99.09% and 98.64% respectively. The F1 scores achieved by these three models also exceeded 0.98 mark, showing a better prediction rate. Naive Bayes produced the lowest MAE, MSE and RMSE values, reflecting the minimal deviation between the anticipated classes and the actual/ground-truth values. Whereas, logistic regression has the highest error rates which signifies its poor performance during fitting of the dataset. The $R^2$ score is highest of Naive Bayes, representing its superior performance on our dataset, as the model performs a good variance over the data. The results imply that Naive Bayes could be the most effective model on the given dataset for the task of crop prediction/recommendation. Whereas, random forest and decision tree also shows high accuracy and low prediction error, proving its effectiveness for use in prediction of crops.

In this study, an ensembled learning approach is introduced to enhance the performance of our model for crop recommendation. The proposed ensemble learning model (i.e. RNT) uses 3 individual models, viz. RF, NB and DT, based on their accuracy scores from Table 2.

**Table 3:** Result comparison of our proposed RNT (RF+NB+DT) model with prior studies

| Model | Precision | RMSE | Recall | MSE | F1 Score | MAE | $R^2$ score | Accuracy (%) |
|---|---|---|---|---|---|---|---|---|
| Hasan et al.[8] | --- | 0.0130 | --- | 0.0160 | --- | 0.0230 | 0.9900 | --- |
| Garanayak et al.[9] | --- | --- | --- | --- | --- | --- | --- | 94.78 |
| Thilakarathne et al.[10] | 0.9700 | --- | 0.9700 | --- | 0.9700 | --- | --- | 97.40 |
| Madhuri et al.[11] | --- | 0.1900 | --- | 0.1000 | --- | 0.1800 | --- | 93.78 |
| Mokarrama et al.[12] | 0.8000 | --- | 0.8000 | --- | --- | --- | --- | --- |
| RNT | 0.9934 | 1.0445 | 0.9924 | 0.0909 | 0.9923 | 1.0909 | 0.9730 | 99.24 |

The comparative analysis in Table 3 shows the superior performance of the proposed RNT model with other existing works on crop recommendation. The proposed model achieved a higher accuracy of 99.24%, with leading values in precision (0.9934), recall (0.9924), and F1 score (0.9923), indicating its excellent prediction capability. In addition, it recorded low error rates (MAE = 0.0909, MSE = 1.0909, RMSE = 1.0445) and a high R² value of 0.9730 when compared with the results in Table 2, suggesting a strong correlation between predicted and actual results. The comparison of the results in Table 2 with Table 3 emphasize that blending

of more than one individual model can significantly enhance the performance, reliability and resilience of a crop recommendation system.

Hasan *et al.*[8] employed ensemble learning and reported their study using error rates, yet omitted major classification metrices like accuracy, precision, f1 and recall. They achieved modest results compared to our proposed model, but was able to predict 17 crops less than our model. Garanayak *et al.*[9] and Madhuri *et al.*[11] presented conventional ML models and an IDBN model with accuracy rates of 94.78% and 93.78% respectively. Thilakarathne *et al.*[10] reported highest result in their experiments by using RF, achieving 97.40% accuracy and 0.97 recall, precision and f1 score. The study proposed by Mokarrama *et al.*[12] recorded 0.80 precision and recall scores using their RSF model. From Table 3, it can be observed that RNT reported better results when compared to most of the models. It is also to be noted that apart from our proposed RNT and the model used by Thilakarathne *et al.*[10], all the other models could predict less than 22 crop varieties.
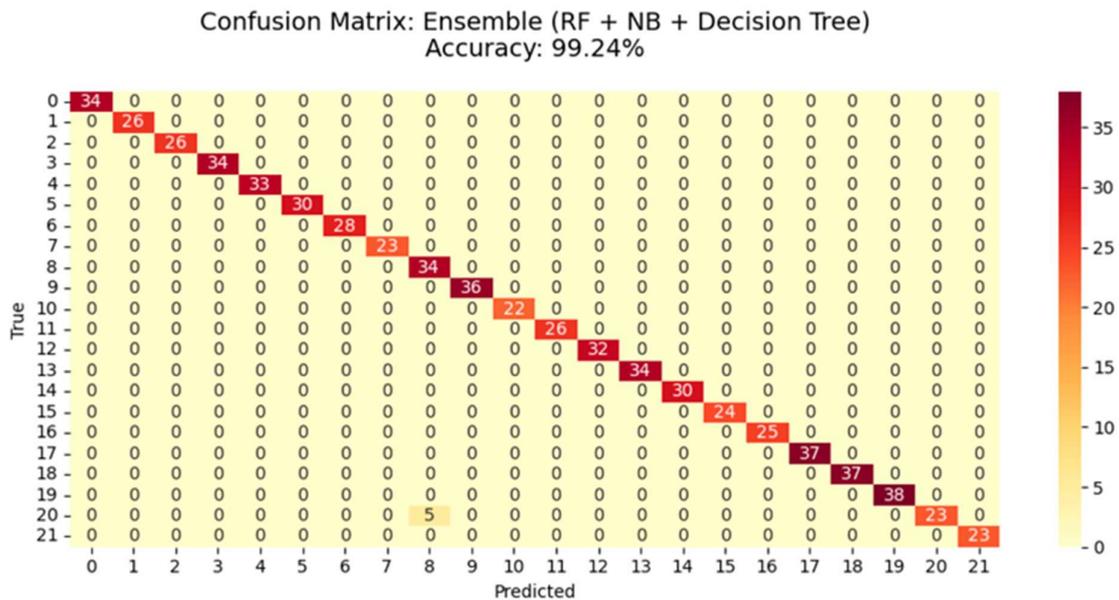


**Figure 7:** Confusion matrices for ensemble model for crop classification

Figure 7 gives heat-map of the confusion matrix produced by the ensemble model during its testing phase. The confusion matrix is represented by the actual labels (marked as **True**) and the predicted labels (marked as **Predicted**) across its axes. The 22 labels in the confusion matrix represent the 22 crop varieties in the dataset. The lables [0-21] represents [rice, maize, jute, grapes, mango, banana,…., jute, coffee]. The matrix shows a high heat signature diagonally, indicating a good prediction score.
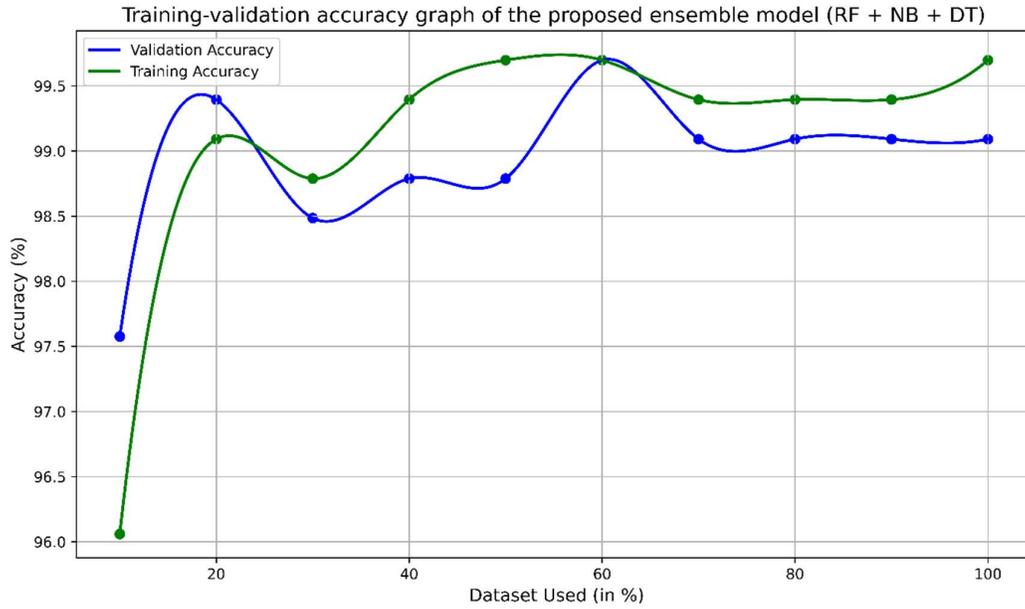
**Figure 8:** Accuracy loss graph of the proposed ensemble model

Figure 8 represents the training and validation-accuracy graph of the ensemble model. This graph gives the algorithm's trend during its training and validation phases. The figure presents a gradual increase and then displays a stable pattern during training and validation stage. A sudden increase of both lines can be observed until 20% of the data is used. Then, both the lines inhibited a stable rise until training and validation with the whole data. Both the lines followed a similar patter with close gaps between them. This indicates a close and efficient learning pattern followed by the model. The overall patterns were not over-fitting as the cent percent accuracy mark was not touched.

**Table 4:** Execution time of the models used in the study

| Model | Execution Time (sec) |
|---|---|
| RF | 0.0873 |
| NB | 0.0005 |
| DT | 0.0024 |
| KNN | 0.0102 |
| SVM | 0.0116 |
| Logistic Regression | 0.0512 |
| RNT | 0.1852 |

Table 4 presents the execution time for each model used in this study. Among individual models, Naive Bayes and Decision Tree are the fastest and best suited for real-time or low-resource environments, while SVM and KNN are relatively slower due to their complexity and computation heavy nature, making them ideal only for high-resource or offline systems. The

ensemble model produced the highest execution time. The performance of the ensemble model is better when compared to other individual models. However, due to its complexity and its heavy utilization of computational resources, the ensemble model must sacrifice its execution time.
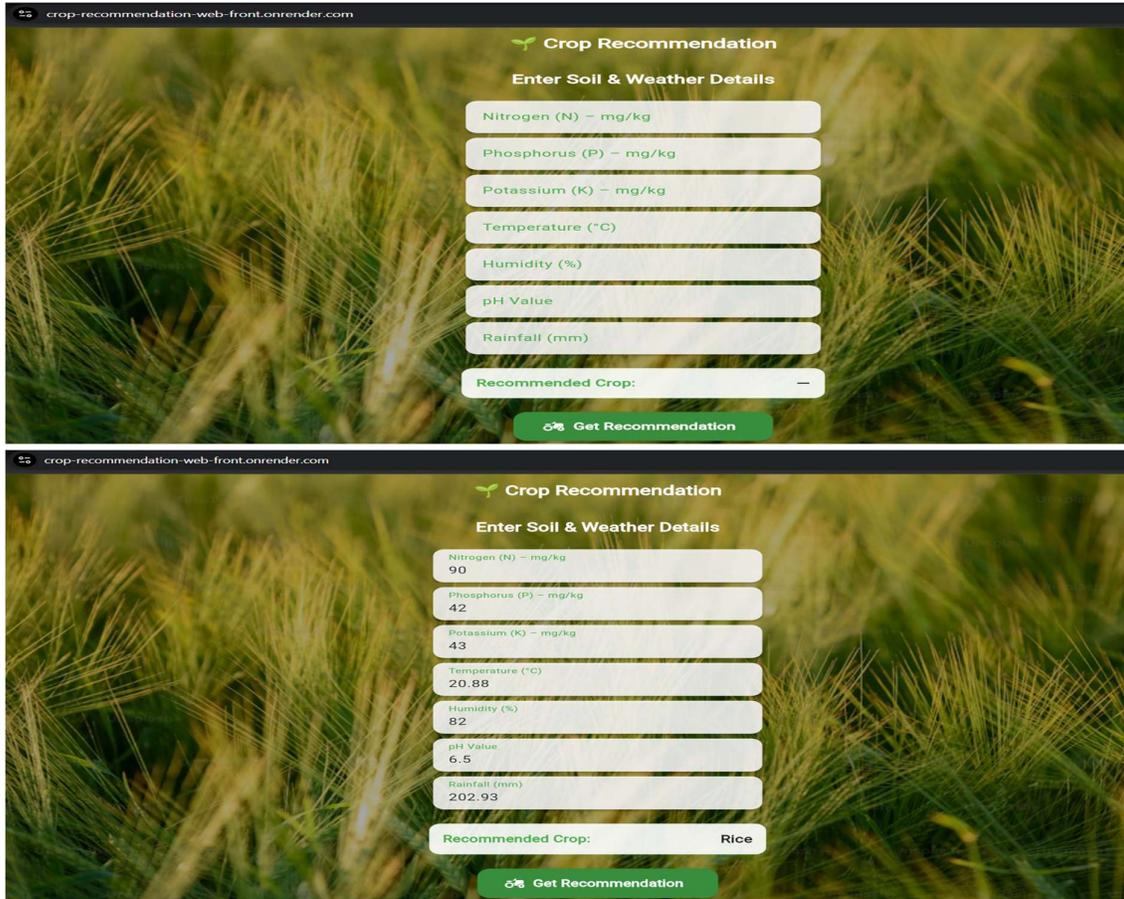


**Figure 9:** The website[1] user-interface of our proposed model

Figure 9 and Figure 10 shows the user-interface of our proposed ensemble model. The development of a website and a smartphone application is done to help farmers recommend the most suitable crop based on the 7 soil and environmental factors as shown. This will provide the farmers to get real-world access to our model.

## 5. Conclusion

An automated crop recommendation model can be a valuable tool for farmers, as it reduces the risks of crop failure. Such models can help farmers predict and recommend crop varieties most suitable for the cultivation area. The environmental and soil parameters play an important role, as the prediction model relies on their values to recommend a crop based on expert knowledge. This model can help farmers optimize fertilizers and other resources, thus promoting sustainable farming practices. It can also help in promoting better land management and balanced nutrient usage.
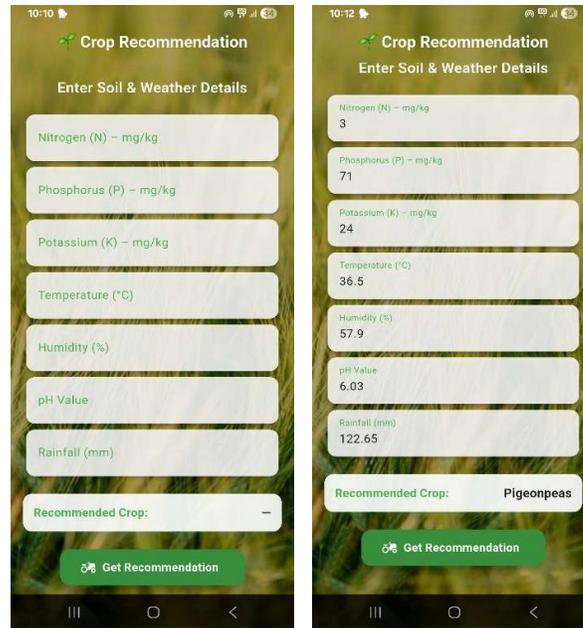
---

[1] https://crop-recommendation-web-front.onrender.com/

**Figure 10:** The smartphone application user-interface of our proposed model

## Conflict of Interest

The authors declare no conflict of interest.

## References

[1]   S. W. Wang, W.-K. Lee, and Y. Son, "An assessment of climate change impacts and adaptation in South Asian agriculture," *IJCCSM*, vol. 9, no. 4, pp. 517–534, Aug. 2017, doi: 10.1108/IJCCSM-05-2016-0069.

[2]   K. Shinde and Auricle Technologies Pvt. Ltd., "Web Based Recommendation System for Farmers," *IJRITCC*, vol. 3, no. 3, pp. 1444–1448, 2015, doi: 10.17762/ijritcc2321-8169.1503119.

[3]   D. C. Tsouros, A. Triantafyllou, S. Bibi, and P. G. Sarigannidis, "Data Acquisition and Analysis Methods in UAV- based Applications for Precision Agriculture," in *2019 15th International Conference on Distributed Computing in Sensor Systems (DCOSS)*, Santorini Island, Greece: IEEE, May 2019, pp. 377–384. doi: 10.1109/DCOSS.2019.00080.

[4]   S. Pudumalar, E. Ramanujam, R. H. Rajashree, C. Kavya, T. Kiruthika, and J. Nisha, "Crop recommendation system for precision agriculture," in *2016 Eighth International Conference on Advanced Computing (ICoAC)*, Chennai, India: IEEE, Jan. 2017, pp. 32–36. doi: 10.1109/ICoAC.2017.7951740.

[5]   Sharma, S., Jindal, S., Gosain, T., & Sharma, H. (2025, February). Deep Learning-Based Crop Recommendation System. In International Conference On Innovative Computing And Communication (pp. 489-505). Singapore: Springer Nature Singapore..

[6]  C. Campbell, S. Sands, C. Ferraro, H.-Y. (Jody) Tsao, and A. Mavrommatis, "From data to action: How marketers can leverage AI," *Business Horizons*, vol. 63, no. 2, pp. 227–243, Mar. 2020, doi: 10.1016/j.bushor.2019.12.002.

[7]  Bhola, A., Shrivastava, G., Sharma, H., & Kumar, P. (2025, February). Harnessing Digital Innovations for Sustainable Agriculture in India: Technology-Driven Smart Farming Framework. In International Conference On Innovative Computing And Communication (pp. 501-512). Singapore: Springer Nature Singapore.

[8]  M. Hasan *et al.*, "Ensemble machine learning-based recommendation system for effective prediction of suitable agricultural crop cultivation," *Front. Plant Sci.*, vol. 14, p. 1234555, Aug. 2023, doi: 10.3389/fpls.2023.1234555.

[9]  M. Garanayak, G. Sahu, S. N. Mohanty, and A. K. Jagadev, "Agricultural Recommendation System for Crops Using Different Machine Learning Regression Methods," *International Journal of Agricultural and Environmental Information Systems*, vol. 12, no. 1, pp. 1–20, Jan. 2021, doi: 10.4018/IJAEIS.20210101.oa1.

[10] Bhola, A., Sharma, H., Sagar, A. K., & Kumar, P. (2024, November). Pre-Harvest to Post-Harvest: A Review of AI and IoT Applications in Smart Agriculture and the Prospects of 6G-Enabled IoT Framework. In 2024 27th International Symposium on Wireless Personal Multimedia Communications (WPMC) (pp. 1-6). IEEE.

[11] J. Madhuri, B. Ramana Reddy, G. Kalnoor, M. Indiramma, and N. Nagarathna, "Optimizing Crop Recommendations With Improved Deep Belief Networks: A Multimodal Approach," *IEEE Access*, vol. 13, pp. 31762–31773, 2025, doi: 10.1109/ACCESS.2025.3542284.

[12] M. J. Mokarrama and M. S. Arefin, "RSF: A recommendation system for farmers," in *2017 IEEE Region 10 Humanitarian Technology Conference (R10-HTC)*, Dhaka: IEEE, Dec. 2017, pp. 843–850. doi: 10.1109/R10-HTC.2017.8289086.