**RESEARCH ARTICLE**

# Enhancing Agricultural Productivity: A Machine Learning Approach to Crop Recommendations

Farida Siddiqi Prity[1,2] · MD. Mehadi Hasan[1] · Shakhawat Hossain Saif[1] · Md. Maruf Hossain[1,3] · Sazzad Hossain Bhuiyan[1] · Md. Ariful Islam[1,4] · Md Tousif Hasan Lavlu[1]

**Abstract**

Agriculture constitutes the foundational pillar of the global economy, engaging a substantial segment of the workforce and making a considerable contribution to the Gross Domestic Product (GDP). However, agricultural productivity faces numerous challenges, including varying climatic conditions, soil types, and limited access to modern farming practices. Developing intelligent agricultural systems becomes imperative to address these challenges and enhance agricultural productivity. Therefore, this paper aims to present a Machine Learning (ML) based crop recommendation system tailored for the farming landscape. The proposed system utilizes historical data on climatic conditions, soil properties, crop yields, and farmer preferences to provide personalized crop recommendations. The goal of this study is to appraise the efficacy of nine distinct ML models—Logistic Regression (LR), Support Vector Machine (SVM), K-Nearest Neighbors (KNN), Decision Tree (DT), Random Forest (RF), Bagging (BG), AdaBoost (AB), Gradient Boosting (GB), and Extra Trees (ET) to generate practical recommendations for crop selection. Numerous preprocessing methods are employed to cleanse and normalize the data, thereby ensuring its appropriateness for model training. The ML models are trained using historical data sets, including temperature, rainfall, humidity, soil pH, and nutrient levels, where crop yields are correlated with environmental and agronomic factors. The models undergo fine-tuning through methods such as cross-validation to enhance their performance and ensure robustness. Among those models, Radom Forest has achieved the highest accuracy (99.31%). The proposed Machine Learning-based crop recommendation system offers a promising approach to addressing the challenges faced by the farmers. By leveraging advanced data analytics and artificial intelligence techniques, the system empowers farmers with timely and personalized recommendations, ultimately leading to improved agricultural productivity, food security, and economic prosperity.

**Keywords** Crop recommendation · Soil · Machine learning · Climate · Decision tree

## 1 Introduction

Agriculture is the oldest industry, essential for nourishing the global population. With modernization and technological integration, it has evolved to optimize efficiency, attract more participants, and elevate overall quality standards [1]. However, the looming threat of diminishing arable land due to urbanization poses a daunting challenge in agriculture. Coupled with the imperative to boost food production by over 70% by 2050 to accommodate population growth, the need for innovative solutions becomes paramount [2]. So, there is an urgent need for an automated system that can maximize harvest yield while minimizing resource inputs, ensuring agricultural sustainability in the face of evolving global demands [3].

✉ Farida Siddiqi Prity
faridasiddiqiprity@gmail.com

1   Department of Computer Science and Engineering, Shanto-Mariam University of Creative Technology, Uttara 17, Dhaka 1230, Bangladesh

2   Department of Information and Communication Engineering, Noakhali Science and Technology University, Noakhali 3814, Bangladesh

3   Department of Information and Communication Engineering, Pabna University of Science and Technology, Pabna 6600, Bangladesh

4   Institute of Information and Communication Technology, Bangladesh University of Engineering and Technology, Dhaka 1000, Bangladesh

Machine Learning has appeared as a formidable instrument in transforming the agricultural sector, proffering innovative solutions to enduring challenges such as crop yield estimation [4–6], plant species identification [7], and disease detection [8–11], thus heralding a new era of precision agriculture. By harnessing extensive datasets, encompassing variables from soil composition and meteorological patterns to crop health and yield projections, ML algorithms furnish critical insights to farmers, agronomists, and researchers alike. In recent years, a multitude of scholars have implemented ML techniques to develop crop recommendation systems grounded in agricultural parameters. In paper [12] advanced an intelligent system named Agro-Consultant, designed to aid farmers in decision-making based on soil type, area, pH, and precipitation. Chougule [13] employed ML algorithms such as Random Forest to forecast harvest yield predicated on soil categories, conditions, and rainfall. Kulkarni et al. [14] proposed a crop recommendation framework leveraging ensemble machine learning methodologies, integrating predictions from a variety of techniques including RF, Naive Bayes (NB), and Lagrangian SVM, to recommend suitable crops based on mud characteristics and parameters with heightened efficacy. These ML approaches were utilized to delineate relationships between soil nutrient levels, identifying a substantial influence of nitrogen (N) on phosphorus (P) and its subsequent effect on potassium (K), while nitrogen demonstrated a weaker correlation with potassium. Additionally, a researcher [15] utilized soil sensors interfaced with an Arduino board and ESP 8266 WiFi module to amass soil data for crop recommendation purposes.

Different algorithms, including Naive Bayes, Logistic Regression, and C4.5, were applied to rainfall data, with C4.5 achieving the highest accuracy at 85.07% [16]. Gosai et al. developed a robust feature extraction framework deploying various ML algorithms such as NB, RF, SVM, DT, LR, and XGBoost. Their experimental results demonstrated that NB, RF, and XGBoost yielded the highest prediction accuracy, reaching an impressive 99% [17]. The researchers further recommended optimal crop selections based on geospatial and climatic data gathered from agricultural databases. A novel hybrid model was introduced, merging Naive Bayes and J48 with association rule mining, where the J48 algorithm achieved a commendable accuracy of 95.9% [18]. Abrougui et al. conducted a predictive analysis for organic potato yield estimation, employing soil property data. Through a comparative evaluation of Artificial Neural Networks (ANN) and Multiple Linear Regression (MLR), ANN exhibited superior performance, boasting a correlation coefficient of 0.975 [19]. Villanueva et al. carried out image acquisition of bitter melon plants from agricultural fields, subsequently analyzing their fruit-bearing potential through the Convolutional Neural Network (CNN)

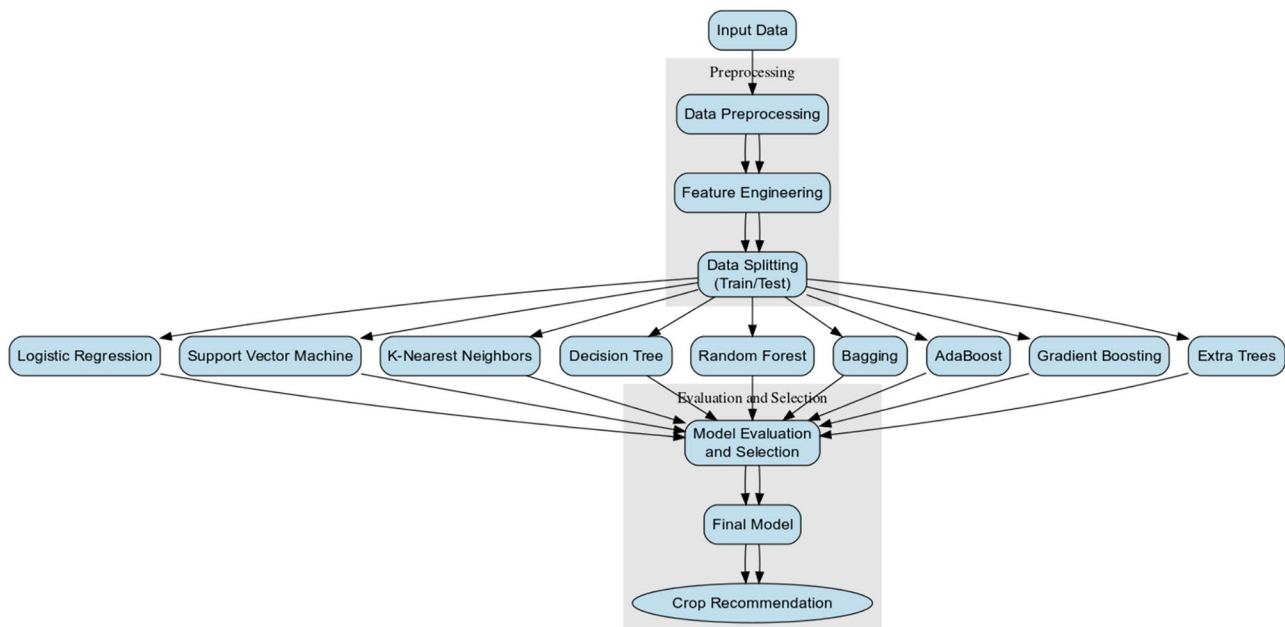methodology. Their findings affirmed that CNN reliably predicts crop yield [20].

While these studies showcase the potential of ML in agricultural applications, a limitation observed in the literature review is the reliance on a relatively small number of ML algorithms, ranging from 2 to 6, for recommending crop strategies in intelligent farming. The narrow scope of algorithms utilized may restrict the diversity of approaches explored and the robustness of the predictive models developed. Therefore, this study has used nine ML algorithms—LR, SVM, KNN, DT, RF, BG, AB, GB, and ET for developing a crop recommendation system by collecting diverse data from the Kaggle repository on soil types, climate, historical crop yields, and farmer preferences. No study has employed all nine specified ML algorithms simultaneously on identical datasets for crop recommendation in intelligent farming applications. Numerous performance metrics, including precision, accuracy, F1-score, and recall, have been employed to gauge the effectiveness of these models. The comparison of models revealed that Random Forest exhibited the highest performance among them. The proposed system will offer a pathway towards sustainable and resilient farming practices, ensuring the long-term viability of our food systems in the face of global challenges.

## 2 Methodology

This paper endeavors to engineer a crop recommendation system by harnessing the capabilities of diverse ML algorithms. The proposed framework encompasses several critical phases: data acquisition and preprocessing, crop prediction through ML methodologies, and the subsequent stages of model training, testing, and performance evaluation. The workflow diagram of the proposed system is depicted in Fig. 1.

### 2.1 Data Collection

The dataset employed in this study was sourced from the Kaggle archives [21], accumulated by the Food and Agriculture Council of India. It comprises 2200 data points and spans 22 crops, including Rice, Maize, Jute, Cotton, Coconut, Papaya, Orange, Apple, Muskmelon, Watermelon, Grapes, Mango, Banana, Pomegranate, Lentil, Blackgram, Mungbean, Mothbeans, Pigeon peas, Kidney beans, Chickpea, and Coffee. This dataset is structured around variables related to Nitrogen, Phosphorus, Potassium, fertilizers, soil pH [22], and climate factors such as rainfall, temperature, and humidity. Initially, we import the data from datasets and check for null and duplicated values within the dataset. Subsequently, we label each crop using one-hot encoding and compile them into a dictionary.

**Fig. 1** Proposed crop recommendation system

**Table 1** Dataset description

| Attributes | Range |
|---|---|
| N | (0–139) kg/ha |
| P | (5–145) kg/ha |
| k | (5–205) kg/ha |
| Temperature | (10.78–43.36) K |
| Humidity | (14.69–98.80) F |
| ph | (3.55–7.45) |
| Rainfall | (20.21–291.29) mm |

Following this, we train the Data Distribution Testing and Scaling function using MinMaxScaler and the Data Training Model. Notably, this dataset exhibits exceptional quality by encompassing a diverse range of geographical conditions and crops, underscoring its potential utility across regions worldwide with similar environmental conditions. A detailed description of the dataset employed in this study is presented in Table 1.

The dataset encompasses information regarding various attributes pertinent to agricultural conditions. The description of each attribute is outlined below:

**N**: This attribute exhibits a range of (0–139) kg/ha, indicative of the quantity of nitrogen in the soil, measured in kilograms per hectare.

**P**: This attribute ranges from (5–145) kg/ha, representing the amount of phosphorus in the soil, measured in kilograms per hectare.

**K**: With a range of (5–205) kg/ha, this attribute denotes the quantity of potassium in the soil, measured similarly to N and P.

**Temperature**: Ranging from (10.78–43.36) K, this attribute is provided in Kelvin value, reflecting the temperature conditions.

**Humidity**: With a range of (14.69–98.80) F, this attribute can be expressed in Fahrenheit or Celsius, indicating humidity.

**$p^H$**: Spanning from (3.55–7.45), the pH attribute typically operates on a scale from 0 to 14, measuring the acidity or alkalinity of a substance, thus reflecting soil conditions.

**Rainfall**: This attribute signifies the volume of rainfall in millimetres, exhibiting a range of (20.21–291.29) mm, providing insight into precipitation levels.

## 2.2 Prediction of Crop Using ML Techniques

The study employs a range of ML algorithms to recommend crops based on various factors, including NPK fertilizer levels, soil pH, and climatic conditions. These algorithms encompass LR, SVM, KNN, DT, RF, BG, AB, GB, and ET.

### 2.2.1 Logistic Regression

Predictive LR emerges as a potent machine learning approach for crop recommendation, particularly in binary classification scenarios where the objective is to ascertain whether a specific crop should be cultivated, based on variables such as soil composition, climatic conditions, and

water resources. The methodology initiates with rigorous data preprocessing, encompassing data cleansing, imputation of missing values, and the normalization or standardization of features to maintain uniformity in scale. Throughout the training phase, Logistic Regression identifies the optimal coefficients that establish a linear decision boundary to delineate the classes. Upon completion of training, the model's performance is scrutinized using metrics. This methodical approach facilitates well-informed decisions regarding crop cultivation, utilizing Logistic Regression to enhance agricultural methodologies.

### 2.2.2 Support Vector Machine

SVMs can be effectively applied to crop recommendation in Machine Learning. The core concept behind using SVM for this purpose is constructing a classification model to identify the most suitable crops to grow based on various factors. The first step involves feature selection, identifying the traits that significantly influence plant growth and crop yield. After selecting the appropriate features, the SVM classifier is trained using preprocessed data [23]. This data should be cleaned, scaled, and normalized to ensure the model's robustness and to avoid overfitting. The SVM algorithm endeavors to ascertain the optimal hyperplane that demarcates the distinct classes of crops within the multidimensional feature space, thereby facilitating the precise discrimination among potential crop recommendations. Upon the successful training and rigorous evaluation of the SVM model, it becomes proficient in delivering crop recommendations by exploiting its capacity to delineate intricate decision boundaries inherent in the dataset. Users can input relevant information, including soil properties, climate conditions, and topography, and the model will predict which crops are most suitable to grow in that particular context. This approach offers a powerful tool for farmers and agricultural planners, ultimately contributing to improved farm productivity and sustainability.

### 2.2.3 K-Nearest Neighbors

K-Nearest Neighbors excels when data is non-linear or doesn't follow a specific distribution [24]. KNN classifies or predicts based on the majority class or average value of the k closest data points in the feature space. Cross-validation can optimize this, with smaller k values offering more sensitivity and larger ones providing stability. KNN is based on the principle of similarity, typically using Euclidean or Manhattan distance to find the closest neighbours. In classification, KNN assigns a class based on the majority vote among the K's Nearest Neighbours. For regression, it predicts by averaging the values of the nearest neighbours. KNN is considered a "lazy learner" because it doesn't build

a complex model during training; computation occurs during the prediction phase. This flexibility makes it a popular choice for many machine-learning tasks, especially when data structures are complex and require a simple yet effective approach. Feature scaling or normalization is recommended to ensure optimal performance with KNN.

### 2.2.4 Decision Tree

Decision Tree models decisions through a tree-like structure comprising nodes and branches, where internal nodes represent conditions or tests on specific features, and leaf nodes represent the final output or decision [25]. The algorithm constructs the tree by selecting the most significant features to split the data, often using criteria like Gini Impurity, Entropy, or Information Gain. Decision Trees are famous for their simplicity, interpretability, and ability to handle non-linear data without complex transformations. DTs are also sensitive to small data changes, leading to variations in the tree structure. Despite these limitations, they are widely used in finance, healthcare, and marketing due to their clear visualization and ease of understanding.

### 2.2.5 Random Forest

A RF is a robust supervised learning algorithm employed in both classification and regression tasks, renowned for its ensemble methodology that amalgamates multiple decision trees. Each tree within the Random Forest is generated by selecting a stochastic subset of the training dataset, alongside a randomly chosen subset of features at each decision node, fostering diversity among the trees. This aggregation process enhances model stability and generalization. Random Forests offer several advantages, including reduced overfitting, improved accuracy, and insights into feature importance, making them useful in various arenas such as healthcare, finance, and marketing. However, they are more complex than single decision trees and require more computational resources. Despite this, Random Forests are favoured for their robustness and versatility in handling large datasets and complex data structures.

### 2.2.6 Bagging

Bagging, an acronym for bootstrap aggregating, is a widely utilized ensemble learning technique in ML. This method involves training multiple models in isolation, each on a different subset of the data generated through bootstrapping. This technique is often used to reduce variance and avoid overfitting. Different crop recommendation models can be trained using different subsets of available data or different algorithms. Baggage helps create different patterns by introducing randomness into the training process. By

combining predictions from multiple models, bagging often leads to higher prediction accuracy than individual models. This can be especially beneficial in crop recommendation, where accurate predictions are crucial for maximizing yield and optimizing resource utilization. Bagging helps make the crop recommendation system more stable and robust by reducing the impact of abnormal or noisy data. Bagging can also be used for feature selection by training each model on a random subset of features. This helps identify the most relevant features for crop recommendation while reducing the risk of overfitting. Since each model in the crop ensemble is trained independently, crop algorithms can be easily parallelized, making them effective for large-scale crop recommendation systems.

### 2.2.7 AdaBoost

AdaBoost, short for Adaptive Boosting, is a prominent ensemble learning algorithm employed in ML for both classification and regression tasks. It functions by sequentially merging numerous weak learners—models that execute slightly better than random guessing—into a single, robust model. In each iteration, AdaBoost adjusts the weights of misclassified instances, prioritizing them in subsequent training rounds. This iterative process amplifies the contribution of weaker learners, culminating in a more accurate and resilient model.

### 2.2.8 Gradient Boosting

Gradient Boosting is a formidable ensemble learning method widely applied in ML for both classification and regression tasks. It constructs a robust predictive model through the iterative training of multiple weak learners, typically decision trees. In this process, each successive tree is trained to rectify the residual errors of its predecessor, thereby incrementally enhancing the model's overall predictive accuracy. The critical concept of Gradient Boosting is to iteratively optimize a loss function by minimizing the residual errors, allowing the model to improve with each step. The process involves adding new trees to the ensemble, each tree trained on the residuals or errors from the preceding iteration, adjusting for overfitting through regularization techniques like learning rate and tree depth control. Gradient Boosting is highly effective for achieving high accuracy, offering flexibility with various loss functions and parameters, making it adaptable to multiple tasks. While it can be computationally intensive and prone to overfitting if not adequately managed, Gradient Boosting is valued for its ability to create robust and precise predictive models, often outperforming other ensemble techniques. It's commonly used in finance, healthcare, and marketing applications, where predictive accuracy is crucial [26].

### 2.2.9 Extra Trees

Extra Trees resembles RFs but introduces greater randomness in its tree construction, leading to increased robustness and condensed risk of overfitting. Extra Trees create an ensemble of decision trees. Still, unlike Random Forests, they randomly choose both features and split thresholds rather than selecting the best split based on a specific criterion. This added randomness can accelerate training times and make the model more resilient to noisy or diverse datasets. In Extra Trees, each tree contributes to the final prediction through majority voting for classification or averaging for regression, yielding an effective ensemble model. The trade-off for this additional randomness is potentially lower accuracy in datasets with distinct patterns. Still, Extra Trees can be a highly efficient and robust method for handling large and high-dimensional datasets, making it a versatile tool in healthcare, finance, and marketing [27].

## 2.3 Training and Testing

Addressing the challenge posed by unbalanced data during the training and testing of machine learning models is crucial, as biased models often yield inaccurate predictions for minority classes. To mitigate this issue, we implemented a down-sampling strategy. Additionally, we employed early stopping techniques to prevent overfitting and enhance model generalization. Machine learning algorithms were employed to forecast the optimal cultivation strategy. The dataset was partitioned, with 80% allocated for model training and the residual 20% reserved for testing to evaluate the model's predictive efficacy.

## 2.4 Performance Metrics

This study used different performance metrics such as Confusion Matrix (CM), ROC (Receiver Operating Curve) curve, precision, recall, F1 score, and accuracy to evaluate the efficacy of the ML models.

**Confusion Matrix**: A CM is a table that visualizes a classification model's performance by comparing predicted and actual classes, highlighting True Positives (TPs), False Positives (FPs), True Negatives (TNs), and False Negatives (FNs) for each class.

**Receiver Operating Characteristic Curve**: The ROC curve illustrates the trade-off between the true positive rate (TPR) and the false positive rate (FPR) across various threshold settings.

**Precision:** Precision enumerates the accuracy of positive predictions by measuring the proportion of TPs among all predicted positives. It can be calculated using Eq. (1):

$$Precision = \frac{TP}{TP + FP} \tag{1}$$

**Recall:** Recall assesses the classifier's effectiveness by measuring the proportion of actual positives that were correctly identified. The F1 score can be calculated using the Eq. (2):

$$Recall = \frac{TP}{TP + FN} \tag{2}$$

**F1 Score:** The F1 score is the harmonic mean of precision and recall, providing a single metric that balances both aspects. It is calculated using the formula:

$$F1 - Score = \frac{2 \times Precision \times Recall}{Precision + Recall} \tag{3}$$

**Accuracy:** Accuracy measures the proportion of correctly classified instances, including both TPs and TNs, relative to the total number of instances. It can be calculated using the formula:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \tag{4}$$

## 3 Results and Discussions

This section provides an in-depth analysis of the performance metrics for all Machine Learning algorithms used in the proposed crop recommendation system. It includes detailed evaluations of ROC curve, accuracy, precision, recall, F1 score, and, confusion matrix analysis. These results enable a comprehensive assessment of each algorithm's effectiveness and appropriateness for the task.

### 3.1 Nutrient Requirements and Environmental Factors of Agricultural and Horticultural Crops

Nutrient requirements, environmental factors affecting agricultural and horticultural crops, and the correlation of crop nutrients play crucial roles in determining crop health, growth, and yield potential. Figure 2 depicts the correlation of crop nutrients, illustrating the interrelationships between various nutrient factors crucial for crop health and productivity. The plot matrixes represent NPK requirements, temperature, humidity, rainfall, and soil $p^H$ for growing various agricultural and horticultural crops. The crop temperature range is generally narrow; However, some cultures, like rice, require high rainfall and prefer warm temperatures. The crops mentioned thrive under specific environmental conditions: acidic soils (pH 5.0–6.5) are ideal for growth, while maize, jute, cotton, coconut, papaya, orange, muskmelon, watermelon, grapes, banana, mango, pomegranate,

mungbean, mothbeans, blackgram, pigeon peas, kidney beans, chickpea, and coffee flourish in warm climates with moderate to high rainfall. These crops are adaptable to a broad pH range but perform optimally in slightly acidic to neutral soils (pH 5.8–7.0). On the other hand, apple and lentil prefer cooler temperatures and moderate rainfall, thriving in well-drained soils with a slightly acidic to neutral pH (6.0–7.0).

Figure 3 displays the histogram of nitrogen for the proposed system, depicting the distribution of recommended nitrogen levels across various crops. This histogram visualizes the recommended nitrogen levels, with nitrogen levels depicted along the x-axis and the frequency or proportion of recommendations along the y-axis. Each bar within the histogram represents a specific range of recommended nitrogen levels, with the height of the bar indicating how frequently recommendations fall within that particular range. Analyzing this data offers valuable insights into the typical nitrogen requirements of different crops, thereby aiding farmers in making informed decisions regarding fertilization practices.
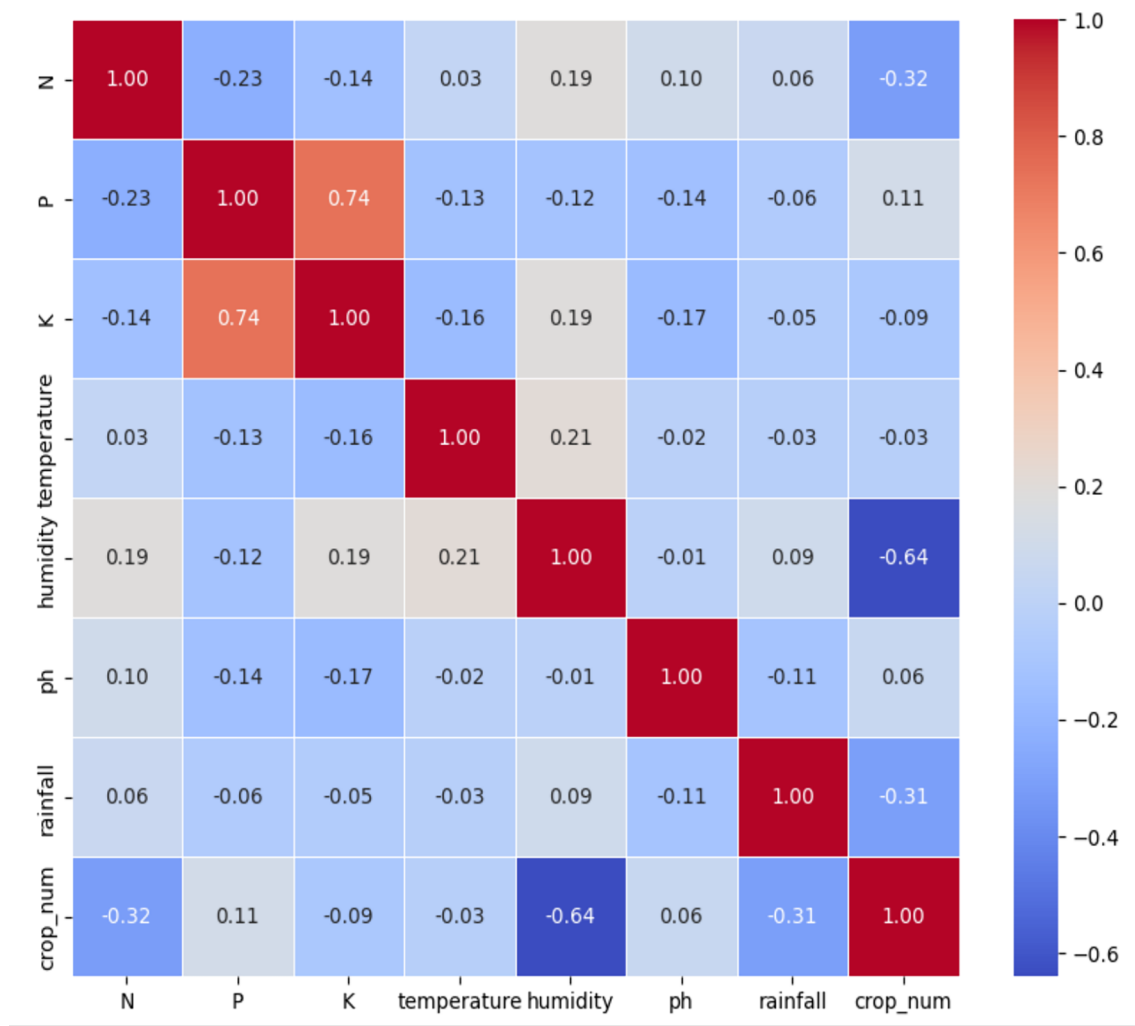
### 3.2 Visualization of Model Responses for Crop Recommendations

The performance evaluation of the crop model is visualized through the analysis of the confusion matrix in Fig. 4, providing a comprehensive understanding of the model's classification accuracy across different classes. Additionally, the ROC curves, depicted in Fig. 5, offer insights into the model's performance by illustrating the trade-off between TP and FP rates across various classification thresholds. Together, these visualizations enable a thorough assessment of the model's predictive capabilities and effectiveness in recommending crops based on diverse agricultural and environmental factors.
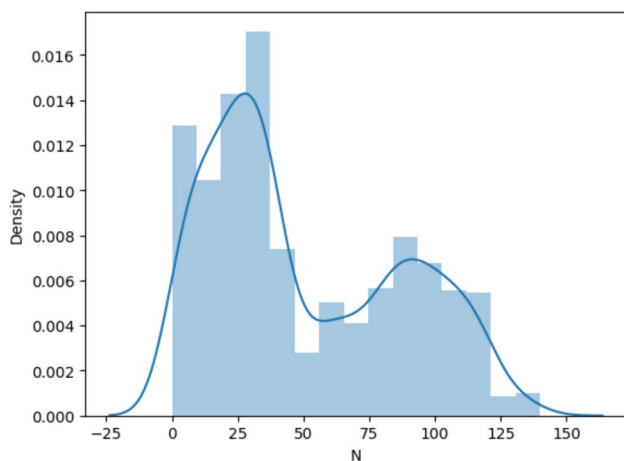
Figure 6 illustrates the correlation matrix of all the crops, serving as a tool to comprehend the interrelationships between different crops concerning various factors such as environmental conditions, soil type, climate, and agricultural practices. Each crop is systematically denoted with its name reiterated in a structured format ("Rice to rice, Maize to maize, Jute to jute…"), indicating self-correlation. The correlation matrix was established by computing Pearson correlation coefficients between crop pairs, derived from their corresponding feature sets. This matrix provides valuable insights into the degree and nature of associations between different crops, aiding in understanding their dependencies and potential interactions within agricultural ecosystems.

### 3.3 Comparative Analysis of ML Algorithms

The study aimed to recommend crops based on multiple factors, employing nine ML algorithms, including LR, SVM,
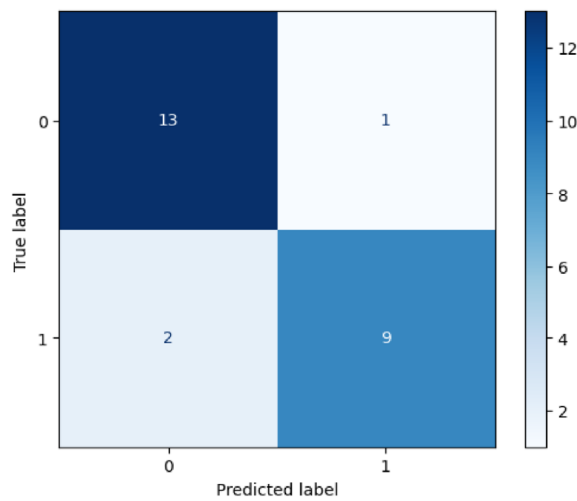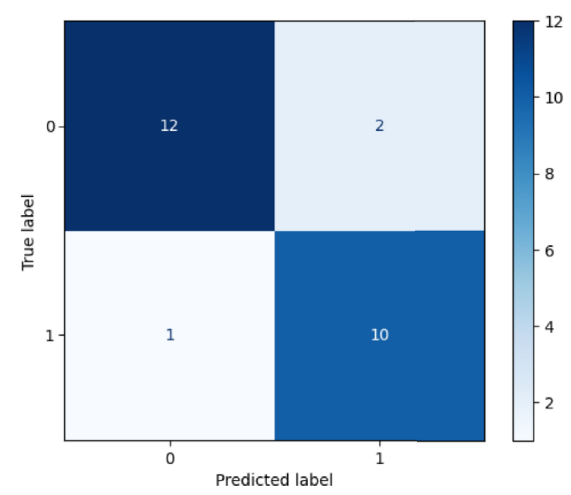
**Fig. 2** Correlation of crop nutrients



**Fig. 3** Histogram of Nitrogen

KNN, DT, RF, BG, AB, GB, and ET. These models underwent training and optimization with specific parameters outlined in the methodology section. A comparative analysis of all models is presented in Table 2.
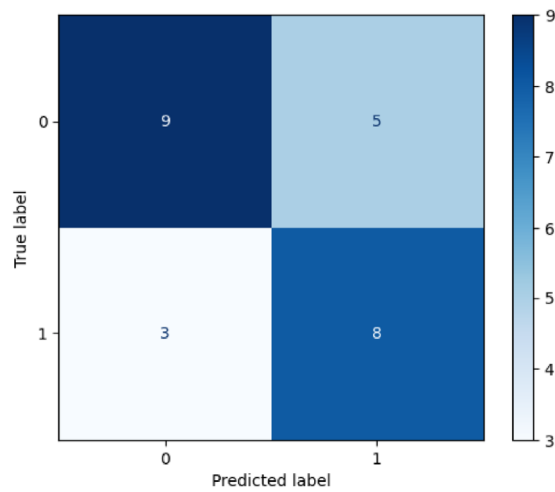
Table 2 provides a comprehensive overview of various ML classifiers and their corresponding precision, recall, and F1 scores. LR achieved a precision of 0.88, recall of 0.87, and an F1-score of 0.88. SVM attained a precision, recall, and F1-score of 0.88. KNN exhibited a precision of 0.68, a recall of 0.69, and an F1-score of 0.68. DT resulted in a precision, recall, and F1-score of 0.72. RF demonstrated a precision of 0.92, recall of 0.93, and an F1-score of 0.92. Bagging showed a precision, recall, and F1-score of 0.92. AB achieved a precision and recall of 0.78, with an F1-score of 0.76. GB yielded a precision of 0.85, a recall of 0.83, and an F1-score of 0.83. ET achieved a precision and recall of 0.96, with an F1-score of 0.96. Comparing the models, Random Forest exhibited the highest test accuracy of 99.31%, while
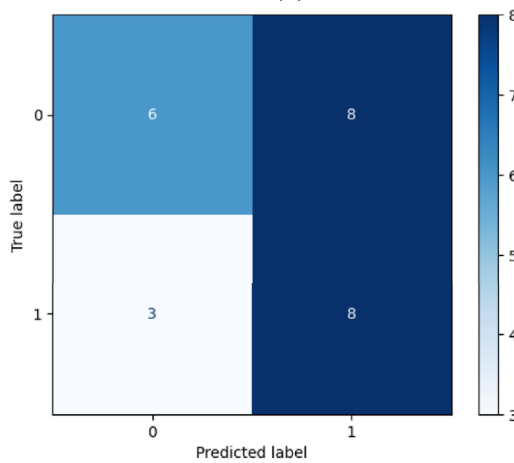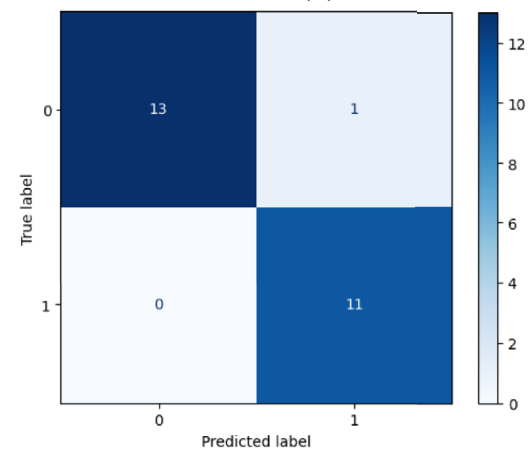
(a) LR

(b) SVM

(c) KNN

(d) DT

(e) ET

(f) RF

**Fig. 4** Confusion matrix of the ML models
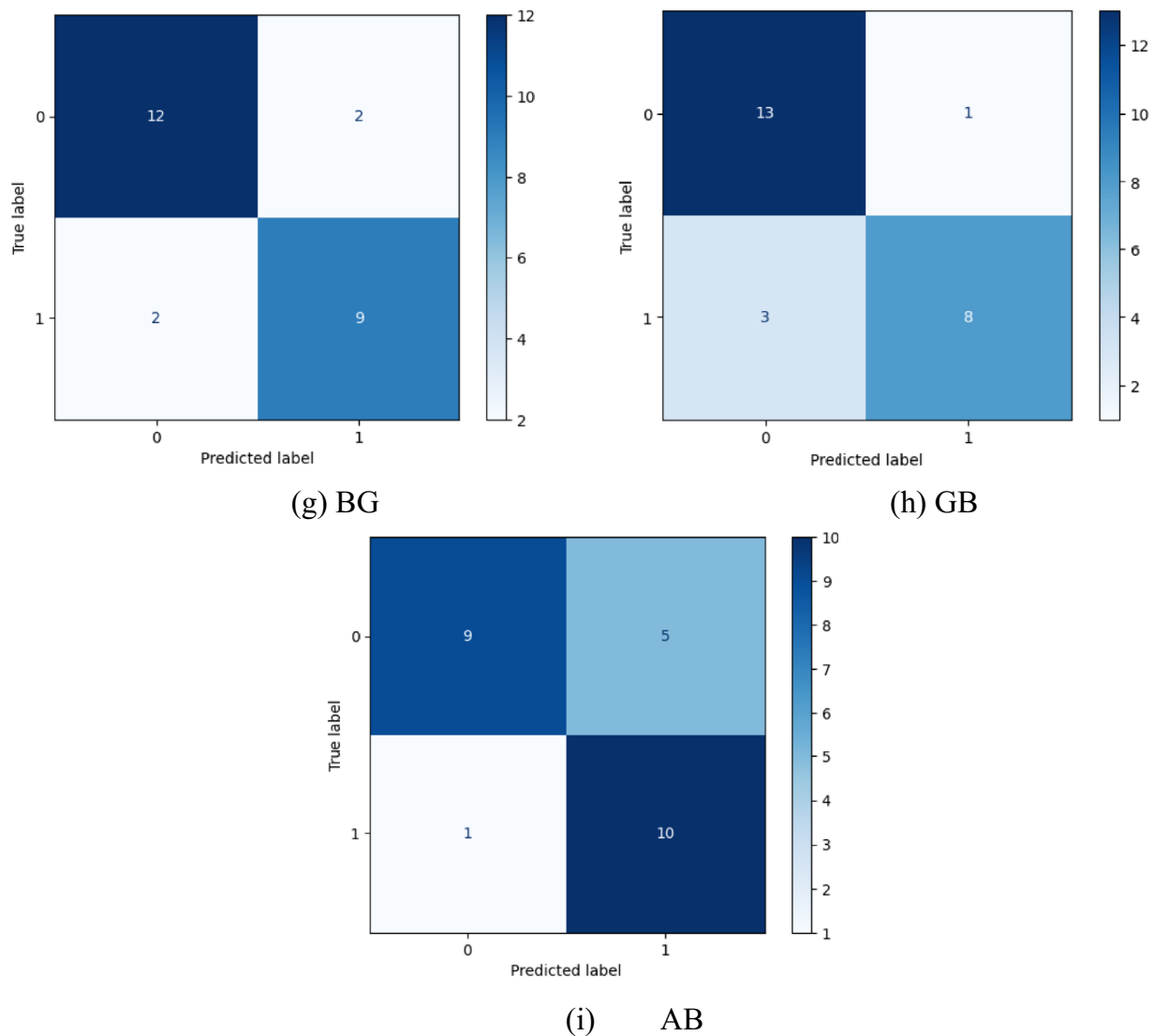
(g) BG

(h) GB

(i)    AB

**Fig. 4** (continued)

AdaBoost displayed the lowest test accuracy of 14.09% for crops. Conversely, the performance of models for horticultural crops showed nearly identical levels of prediction accuracy, precision, F1 score, and recall.
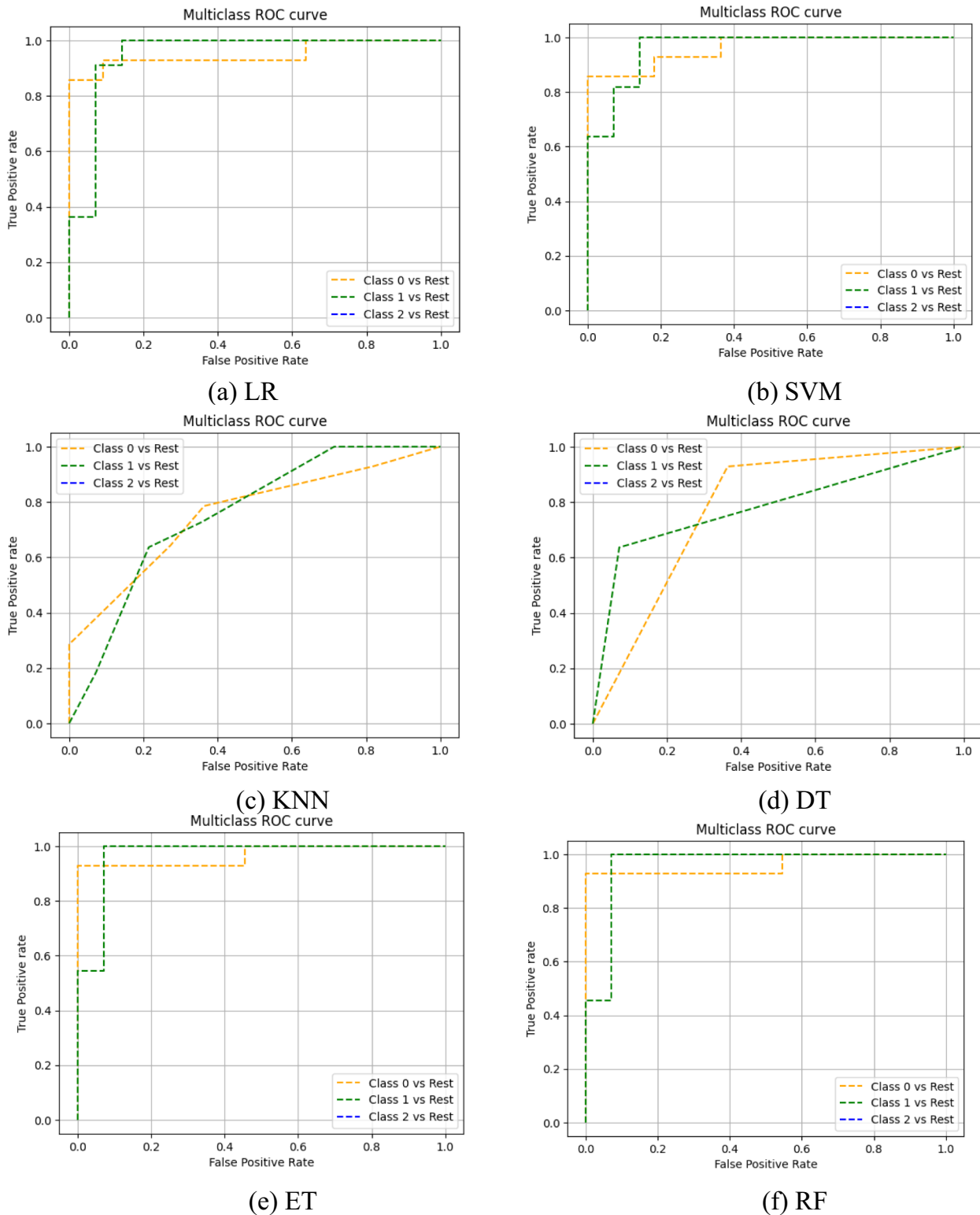
### 3.4 Discussion

The proposed crop recommendation system is a pivotal application of ML in agriculture, aiming to empower farmers with data-driven insights for strategic crop selection. By leveraging various factors like soil properties, climate conditions, historical data, and market demand, this system aids farmers in making well-informed decisions regarding crop cultivation. Identifying pertinent features that influence crop growth and yield, such as temperature, rainfall, soil pH, and nutrient levels, forms the cornerstone of this endeavour.

Integrating domain knowledge, including agronomic expertise and local farming practices, enriches the accuracy and relevance of the system's recommendations. Moreover, establishing farmer feedback loops plays a vital role in refining and enhancing the recommendation system iteratively, ensuring its adaptability to evolving agricultural landscapes.

Recognizing the inherent uncertainty in predictions stemming from factors like variability in weather patterns and market conditions underscores the importance of incorporating robust risk assessment techniques. By accounting for associated risks alongside potential yields, farmers can make more informed decisions that balance productivity and risk mitigation.

In essence, the convergence of ML techniques with agricultural domain knowledge fosters a dynamic ecosystem wherein farmers can navigate complex agricultural

**Fig. 5** ROC curve of the ML models

challenges with confidence and resilience. Table 3 compares the proposed work with existing studies, showcasing the performance metrics and outcomes. The data presented

in Table 3 demonstrate that our proposed approach outperforms existing works, highlighting superior performance across various evaluation criteria.
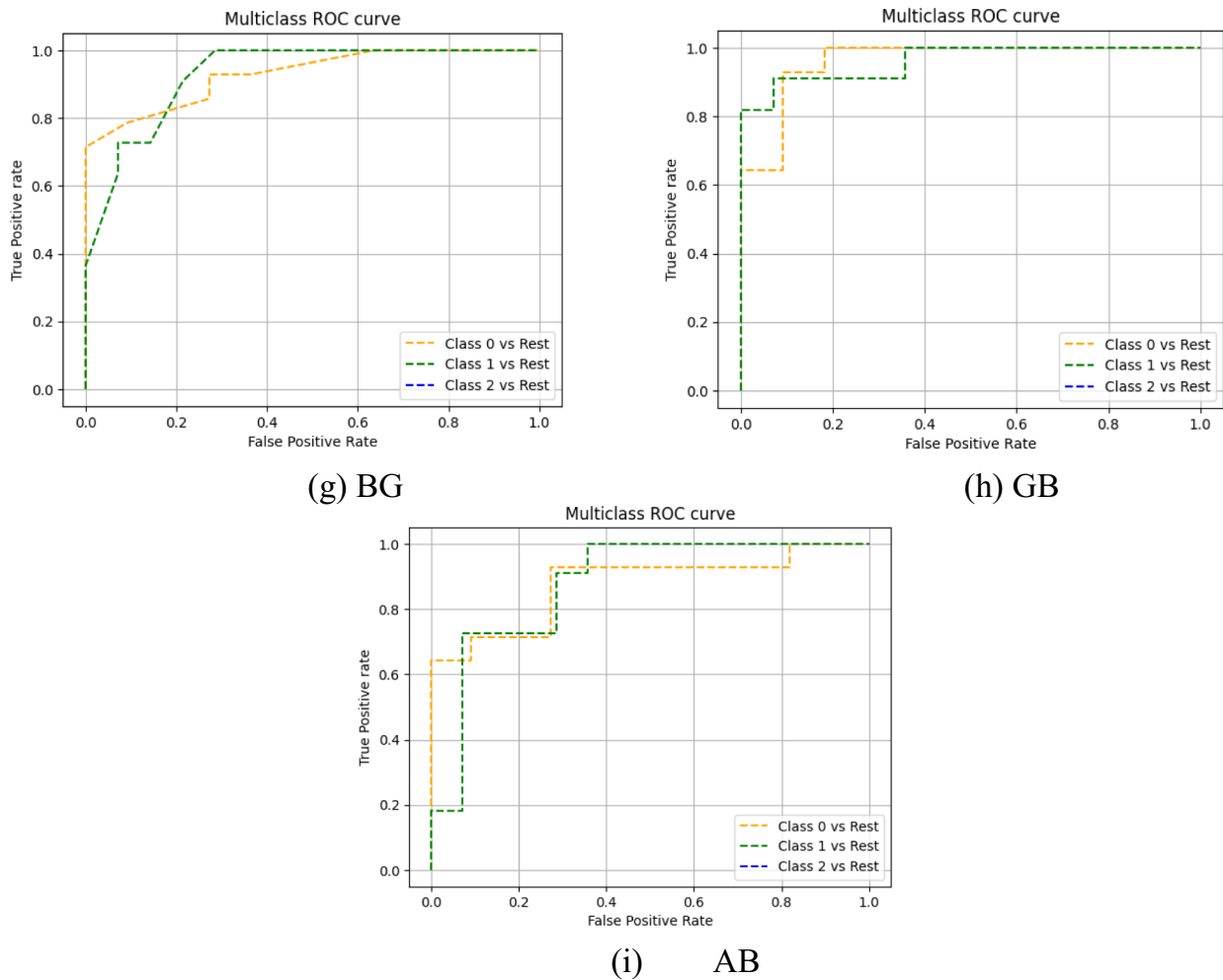
(g) BG



(h) GB



(i)     AB

**Fig. 5** (continued)

### 3.4.1 Possible Causes of Miss-Recommendation

Crop suitability is subject to considerable temporal variability, influenced by seasonal fluctuations, climate dynamics, and evolving agricultural practices. Failing to account for these temporal dynamics in the model adequately may result in outdated or irrelevant recommendations that do not align with current conditions. Moreover, the suitability of crops can vary significantly across different geographic regions due to variations in soil types, climatic conditions, and other factors. Models trained on data from one region may struggle to generalize effectively to others, potentially leading to inaccurate recommendations when applied in diverse geographic areas.

Inaccurate model evaluation metrics or insufficient validation processes can also pose challenges. If the evaluation metrics fail to reflect real-world performance accurately or the validation process lacks rigor, it may result in overly optimistic estimates of model performance. Consequently,

this increases the likelihood of erroneous recommendations when the model is deployed in production settings. Additionally, certain factors influencing crop adaptation may not be well understood or adequately modelled, leading to unexpected phenomena that can influence recommendations in unforeseen ways. Therefore, addressing these challenges and enhancing the robustness of the model's temporal adaptability and geographic generalizability is critical for ensuring the reliability and efficacy of crop recommendation systems.

### 3.4.2 Implication of the Study

The proposed systems use machine learning techniques to analyze data and provide farmers with personalized advice on choosing which crops to grow in a specific location and conditions. This system can help maximize crop productivity by suggesting the most suitable crops for a given location and season. By considering factors such as soil type, moisture levels, temperature, and sunlight, machine learning
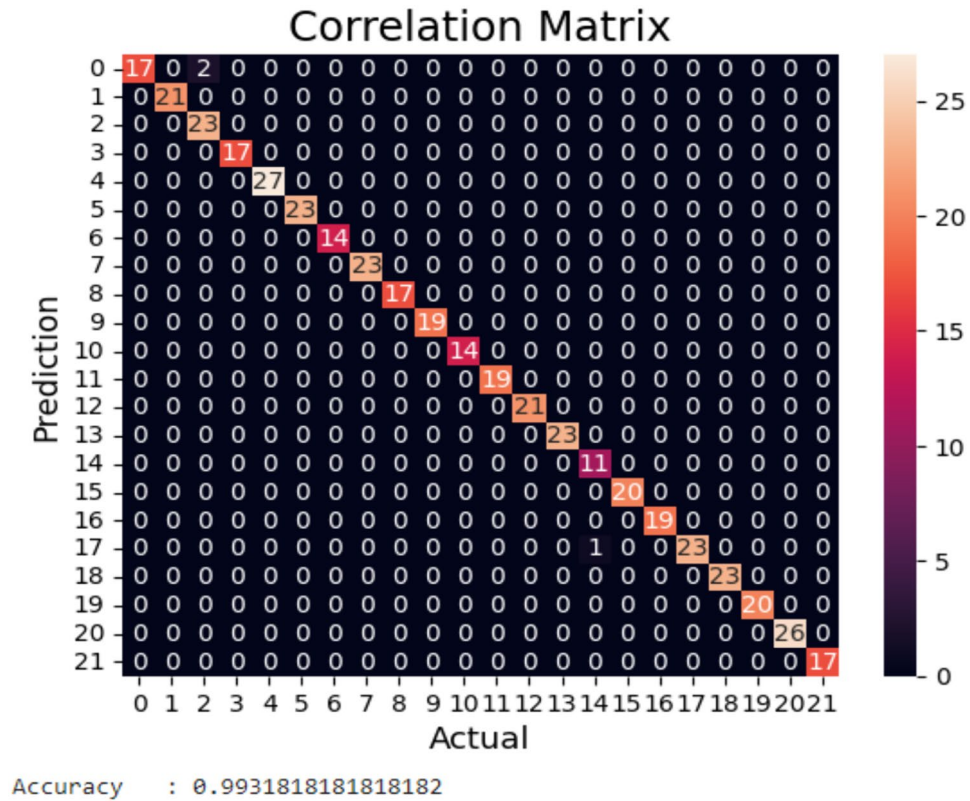
**Fig. 6** Correlation of Crops



**Table 2** Comparative Analysis of ML Algorithms

| Name of classifiers | Precision (%) | Recall (%) | F1-score (%) | Accuracy (%) |
|---|---|---|---|---|
| LR | 88 | 87 | 88 | 96.36 |
| SVM | 88 | 88 | 88 | 96.82 |
| KNN | 68 | 69 | 68 | 96 |
| DT | 72 | 72 | 72 | 98.2 |
| RF | 92 | 93 | 92 | 99.3 |
| BG | 92 | 92 | 92 | 98.9 |
| AB | 78 | 78 | 76 | 14.1 |
| GB | 85 | 83 | 83 | 98.2 |
| ET | 96 | 96 | 96 | 91 |

algorithms can identify crops likely to thrive in specific environmental conditions, thereby increasing yields and profitability for farmers [26]. This study can help reduce risks associated with climate fluctuations, pest infestations, and market fluctuations. By diversifying crop selection and considering factors such as disease resistance and market demand, farmers can reduce risk exposure and improve their resilience to external shocks. This system can promote sustainable agriculture practices by encouraging the cultivation of environmentally friendly and economically viable crops. Promoting crop rotation, soil conservation [27], and biodiversity can help farmers maintain their land's long-term health and productivity [28, 29]. The importance of this type

**Table 3** Comparative analysis of the proposed work with existing studies

| Author | Precision (%) | Recall (%) | F1-Score (%) | Accuracy (%) |
|---|---|---|---|---|
| [19] Islam et al | – | – | – | 90% |
| [12] Doshi et al | – | – | – | 91% |
| [13] Chougule et al | 65% | – | – | – |
| [30] Jain et al | – | – | – | 88% |
| [17] Liu et al | 97.4% | – | 97.8% | 98.7% |
| Our Proposed Model (Random Forest) | 99% | 99% | 99% | 99% |

of model lies in informed decision-making when practising precision agriculture, as a variety of external factors influence crop performance. Attempting to understand the combined effects of these factors without the aid of models represents a daunting challenge. Agricultural productivity is profoundly influenced by a range of biotic and abiotic stressors. The projected global GDP from agriculture is offset by annual economic losses driven by factors such as climate change, market instability, and resource limitations. This loss affects not only agricultural producers but also extends to related industries and economies worldwide, posing significant challenges to food security, livelihoods, and sustainable development efforts on a global scale.

## 4 Conclusion

Agriculture is the cornerstone of the world's economy, employing a significant workforce and contributing substantially to nations' GDPs. However, agricultural productivity has many challenges, ranging from unpredictable climatic conditions to varying soil types and limited access to modern agricultural practices. Developing intelligent agricultural systems is imperative to overcome these hurdles and bolster agricultural productivity.

This paper addresses this need by presenting a Machine Learning-based Crop Recommendation System tailored to the agricultural landscape. The proposed system aims to deliver personalized crop recommendations by leveraging historical data on climatic conditions, soil properties, crop yields, and farmer preferences. Through the evaluation of nine distinct ML models, including LR, SVM, KNN, DT, RF, BG, AB, GB, and ET, this study seeks to generate practical recommendations for crop selection. Employing various preprocessing techniques to clean and normalize the data ensures its suitability for model training. The machine learning models are trained on historical datasets encompassing temperature, rainfall, humidity, soil pH, and nutrient levels, correlating crop yields with environmental and agronomic factors. Fine-tuning of the models, employing techniques such as cross-validation, optimizes their performance and ensures robustness. Among these models, Random Forest emerges as the top performer, achieving the highest accuracy rate of 99.31%. The proposed Machine Learning-based Crop Recommendation system offers a promising avenue for addressing the challenges faced by farmers. By harnessing advanced data analytics and artificial intelligence techniques, the system empowers farmers with timely and personalized recommendations, ultimately leading to improved agricultural productivity, food security, and economic prosperity.

While our current study effectively utilizes the available dataset to develop a crop recommendation system, we acknowledge certain limitations. Specifically, our dataset does not fully capture variations in land quality, climate changes, or historical crop planting data, which are crucial for accurate recommendations. Future research will focus on incorporating more detailed and temporal soil and climate data, as well as historical crop planting information. We also plan to integrate farmers' professional knowledge and explore advanced Machine Learning techniques and hybrid models to address the complexities of agricultural decision-making better. These enhancements will improve the system's adaptability and accuracy, ultimately providing more personalized and effective recommendations for farmers.

## Declarations

## References

1. Wang SW, Lee WK, Son Y. An assessment of climate change impacts and adaptation in South Asian agriculture. Int J Clim Chang Strateg Manag. 2017;9:517–34. https://doi.org/10.1108/IJCCSM-05-2016-0069.
2. Bouguettaya A, Zarzour H, Kechida A, Taberkit AM. Deep learning techniques to classify agricultural crops through UAV imagery: a review. Neural Comput Appl. 2022;34:9511–36. https://doi.org/10.1007/S00521-022-07104-9.20223412.
3. Gathala MK, Timsina J, Islam MS, Rahman MM, Hossain MI, Harun-Ar-Rashid M, McDonald A. Conservation agriculture based tillage and crop establishment options can maintain farmers'

yields and increase profits in South Asia's rice-maize systems: evidence from Bangladesh. Field Crops Res. 2015;172:85–98.

4. Fei S, Hassan MA, Xiao Y, Su X, Chen Z, Cheng Q, Duan F, Chen R, Ma Y. UAV-based multi-sensor data fusion and machine learning algorithm for yield prediction in wheat. Precis Agric. 2022. https://doi.org/10.1007/s11119-022-09938-8.

5. Paudel D, Boogaard H, de Wit A, van der Velde M, Claverie M, Nisini L, Janssen S, Osinga S, Athanasiadis IN. Machine learning for regional crop yield forecasting in Europe. Field Crops Res. 2022. https://doi.org/10.1016/J.FCR.2021.108377.

6. Burdett H, Wellen C. Statistical and machine learning methods for crop yield prediction in the context of precision agriculture. Precis Agric. 2022;23:1553–74. https://doi.org/10.1007/s11119-022-09897-0.

7. Duke OP, Alabi T, Neeti N, Adewopo J. Comparison of UAV and SAR performance for Crop type classification using machine learning algorithms: a case study of humid forest ecology experimental research site of West Africa. Int J Remote Sens. 2022;43:4259–86. https://doi.org/10.1080/01431161.2022.2109444.

8. Too EC, Yujian L, Njuki S, Yingchun L. A comparative study of fine-tuning deep learning models for plant disease identification. Comput Electron Agric. 2019;161:272–9. https://doi.org/10.1016/j.compag.2018.03.032.

9. Chu H, Zhang C, Wang M, Gouda M, Wei X, He Y, Liu Y. Hyperspectral imaging with shallow convolutional neural networks (SCNN) predicts the early herbicide stress in wheat cultivars. J Hazard Mater. 2022. https://doi.org/10.1016/j.jhazmat.2021.126706.

10. Dey B, Masum Ul Haque M, Khatun R, Ahmed R. Comparative performance of four CNN-based deep learning variants in detecting Hispa pest, two fungal diseases, and NPK deficiency symptoms of rice (Oryza sativa) Comput. Electron Agric. 2022. https://doi.org/10.1016/j.compag.2022.107340.

11. Sai Sankar PR, Ramakrishna SDPS, Venkata Rakesh MM, Raja P, Hoang VT, Szczepanski C. Intelligent health assessment system for paddy crop using CNN, 2021 3rd. Int Conf Signal Process Commun ICPSC. 2021;2021:382–7. https://doi.org/10.1109/ICSPC51351.2021.9451644.

12. Z. Doshi S, Nadkarni R, Agrawal N. Shah, Agro-consultant: intelligent crop recommendation system using machine learning algorithms. In: Fourth International Conference on Computing Communication Control and Automation (ICCUBEA), IEEE, 2018, pp. 1–6. https://doi.org/10.1109/ICCUBEA.2018.8697349.

13. ChouguleVKA, Mukhopadhyay D. Crop suitability and fertilizers recommendation using data mining techniques, in: Advances in Intelligent Systems and Computing, Vol. 714, Springer Verlag, 2019, pp. 205–213. https://doi.org/10.1007/978-981-13-0224-419.

14. Kulkarni NH, Srinivasan GN, Sagar BM, Cauvery NK. Improving crop productivity through A crop recommendation system using ensembling technique. Proc 2018 3rd Int Conf Comput Syst Inf Technol Sustain Solut CSITSS. 2018. https://doi.org/10.1109/CSITSS.2018.8768790.

15. Modi D, Sutagundar AV, Yalavigi V, Aravatagimath A. Crop recommendation using machine learning algorithm. 2021 5th Int Conf Inf Syst Comput Networks ISCON. 2021.

16. Prabhu S, Revandekar P, Shirdhankar S, Paygude S. Soil analysis and crop prediction. Int J Sci Res Sci Technol. 2020;7(4):117–23.

17. Gosai D, Raval C, Nayak R, Jayswal H, Patel A. Crop recommendation system using machine learning. Int J Sci Res Comput Sci Eng Inf Technol. 2021;7(3):558–69.

18. Viviliya B, Vaidhehi V. The design of hybrid crop recommendation system using machine learning algorithms. Int J Innov Technol Explor Eng (IJITEE). 2019;9:4305–11.

19. Abrougui K, Gabsi K, Mercatoris B, Khemis C, Amami R, Chehaibi S. Prediction of organic potato yield using tillage systems and soil properties by artificial neural network (ANN) and multiple linear regressions (MLR). Soil Tillage Res. 2019;190:202–8.

20. Villanueva MB, Salenga MLM. Bitter melon crop yield prediction using machine learning algorithm. Int J Adv Comput Sci Appl. 2018;9:1–6.

21. Kaggle. https://www.kaggle.com/datasets/atharvaingle/crop-recommendation-dataset. Accessed date: 12 January 2024

22. Suresh G, Kumar DAS, Lekashri DS, Manikandan DR, Head C-O. Efficient crop yield recommendation system using machine learning for digital farming. Int J Mod Agric. 2021;10:906–14.

23. Hua Y, Li F, Yang S. Application of support vector machine model based on machine learning in art teaching. Wireless Commun Mobile Comput. 2022. https://doi.org/10.1155/2022/7954589.

24. Dey B, Abir KAM, Ahmed R, Salam MA, Redowan M, Miah MD, Iqbal MA. Monitoring groundwater potential dynamics of north-eastern Bengal Basin in Bangladesh using AHP-Machine learning approaches. Ecol Indicat. 2023. https://doi.org/10.1016/j.ecolind.2023.110886.

25. Charoen-Ung P, Mittrapiyanuruk P. Sugarcane yield grade prediction using random forest with forward feature selection and hyper-parameter tuning. Adv Intell Syst Comput. 2019;769:33–42. https://doi.org/10.1007/978-3-319-93692-5_4.

26. Sridevi V, Chellamuthu V. Impact of weather on rice—a review. Int J Appl Res. 2015;1:825–31.

27. Swathi T, Sudha S. Crop classification and prediction based on soil nutrition using machine learning methods. Int J Inf Technol. 2023;15:2951–60. https://doi.org/10.1007/S41870-023-01345-0.

28. Singh Jatav Sri Karan H, Jena J, Maitra S, Hossain A, Pramanick B, Gitari HI, Praharaj S, Shankar T, Bharati Palai J, Rathore A, Kumar Mandal T, Singh Jatav H. Role of legumes in cropping system for soil ecosystem improvement improvement of cropping system view project precision agriculture view project role of legumes in cropping system for soil ecosystem improvement. Nova Science Publishers, Inc. 2022.

29. Thilakarathne NN, Bakar MSA, Abas PE, Yassin H. A cloud enabled crop recommendation platform for machine learning-driven precision farming. Sensors. 2022. https://doi.org/10.3390/s22166299.

30. JainS, Ramesh D. Machine Learning convergence for weather-based crop selection, IEEE International Students' Conference on Electrical, Electronics and Computer Science (2020). https://doi.org/10.1109/SCEECS48394.2020.75.