# CROP RECOMMENDATION SYSTEM USING GRADIENT BOOSTING ALGORITHM

**[1]Dr. G ROSLINE NESA KUMARI, [2]E SAI SATHVIK, [3]M SARATH CHANDRA,
[4]S PRADEEP REDDY, [5]S KALYAN RAJU**

Department of Computer Science and Engineering
Bharath Institute of Higher Education and Research

*Abstract*- **Crop Recommendation System (CRS) leveraging the Gradient Boosting algorithm to enhance precision in agricultural decision-making. With the evolving challenges in agriculture, traditional methods of crop selection often struggle to account for the dynamic interactions among environmental variables. Gradient Boosting Algorithm which is known of its strong feature for solving diversified and complex problems wins over as a reliable and adaptable method for developing exact and dynamic crop recommendations. Introduction of the study includes a review on relevance of literature that shows us the evolutionary part of machine learning in agriculture and its significance in application of crop recommendation systems through accomplishment with the help of Gradient Boosting algorithms. Where the practical part of the scientific method is given more attention, the paper dwells on the critical stages of data acquisition and preprocessing including the consideration of plant properties, the history of inputs methods, clay data and weather information. After this, the Gradient Boosting algorithm is expounded on, exposing which enhancing mechanisms are implemented, basic principles of ensemble learning techniques and the best model configuration. CRS's efficacy is revealed by many different sets of case studies from agrarian areas, in which it shows its ability to generate particular solutions depending on local departments. In the side-wise comparison of the algorithm with other machine-learning models, the former stands out in terms of predictive prowess. The paper winds up discussing the scalability issues by recognizing constraints, challenges, and possible options for future research, indicating that the mentioned CRS based on Gradient Boosting is a potential instrument for systematic crop decision, output maximization, and ecologically friendly agriculture.**

*Keywords:* **decision tree, CRS, Gradient Boosting Algorithm, machine learning.**

## I. INTRODUCTION

Agriculture though grapples with a horde of problems has to cater not only to the increasing population but also to the demanding factors of climate variability, soil complexity and dynamic market situations. The key farming activity of crop selection plays a central position in relation to the farm's success, with its effect on yields, resources use, and agricultural sustainability as a whole. Nevertheless, standard way of plants recommendation fails to absorb the complex and fast-changing ecosystem of current world agriculture. We summon here an innovative method of plants recommendation based on the application of the complex Gradient Boosting algorithm as the driving force. The main target is to enhance recommendations for crops efficiency and responsiveness using advanced data analytics which are updated in real-time knowledge-based process for informed decision-making and allocation of vital resources. As we live in a world where weather systems become unpredictable and where natural environment doesn't follow expected dynamics, the necessity of these robust and innovative agricultural technologies gets clearer. Gradient Boosting algorithm, with proven capability of dealing across diverse datasets and managing overfitting issues, becomes one of the top choices for the complexities of current day agriculture technology. The scientific basis is made possible by its ability for miniature data capturing, learning from daily failures and enhancing predictive accuracy; it is therefore likely to bring a revolution in the manner grain and seed is recommended.
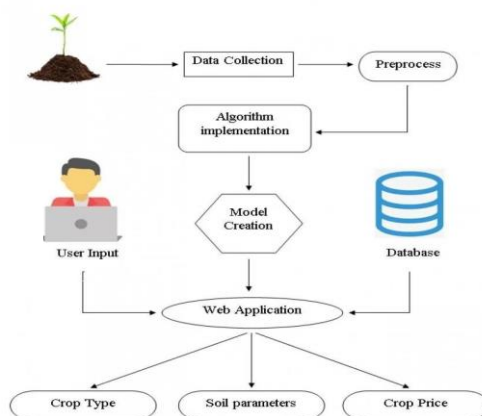
Creating CRS, implementation, and evaluation as of the project goals. Through examination of the details of data collection and preparation the paper strives to ascertain the quality and scope of input data that includes the same as soil rate, climate features, historical yield, and crop properties.

A detailed exposition of the Gradient Boosting algorithm follows, elucidating its boosting mechanisms, ensemble learning principles, and optimal model configurations.

Machine learning techniques have been substantial during the past few decades. ML is conceptually different from much of traditional statistics since its primary focus is outcome prediction rather than inference into the mechanical processes underlying those outcomes. With the aid of computer-based tools called Crop Recommendation Systems (CRS), farmers may make well-informed judgments about which crops to sow depending on variables like soil composition, weather patterns, and past crop yields.

Furthermore, the paper demonstrates the practical application of the CRS through insightful case studies in diverse agricultural regions, showcasing its efficacy in delivering tailored recommendations based on local conditions. Comparative analyses with other machine learning models underscore the algorithm's superiority in predictive capabilities. The introduction sets the stage for a comprehensive exploration of the proposed CRS, highlighting its potential to revolutionize crop selection processes and contribute to the sustainability and efficiency of global agriculture. Machine learning techniques have been substantial during the past few decades. ML is conceptually different from much of traditional statistics since its primary focus is outcome prediction rather than inference into the mechanical processes underlying those outcomes.

## II.  SYSTEM ARCHITECTURE:



## III.  DATA CLEANING AND PREPROCESSING:
### 1.Identification of Relevant Variables:
The first step in data collection involves identifying the key variables that influence crop growth and yield. These may include soil properties (such as pH, moisture, and nutrient levels), climate data (temperature, precipitation, and humidity), historical yields, and specific crop characteristics. This step requires collaboration with agricultural experts to ensure comprehensive coverage of factors affecting crop performance.

### 2. Data Sources and Acquisition:
Once the variables are identified, the next step is to gather relevant data from diverse sources. Soil data can be obtained through soil testing laboratories, climate data from meteorological agencies, and historical yields from agricultural records. Additionally, satellite imagery and remote sensing technologies can provide real-time data on crop conditions, adding a dynamic element to the dataset. Ensuring data accuracy and consistency is crucial at this stage.
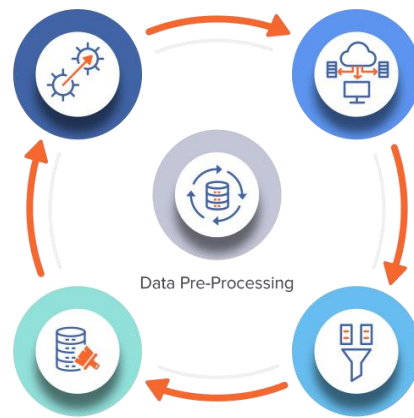
### 3. Data Cleaning:
Data collected from different sources may contain errors, outliers, or missing values. Cleaning the dataset involves identifying and rectifying such discrepancies to ensure the accuracy and reliability of the information. Techniques such as imputation, removal of outliers, and consistency checks are applied to enhance the quality of the dataset.

### 4. Data Normalization:
To bring uniformity to the dataset and ensure that variables are on a similar scale, normalization is employed. This step prevents certain features from disproportionately influencing the model due to differences in their scales. Common normalization techniques include Min-Max scaling or Z-score normalization.

### 5. Feature Engineering:
Feature engineering involves creating new features or modifying existing ones to enhance the predictive power of the model. For crop recommendation, this might include deriving variables such as growing degree days, which reflect the accumulation of heat over time, or creating categorical features to represent specific soil types.

Data Pre-Processing

**6.Handling Categorical Data:**
Some variables, such as crop types or soil classifications, may be categorical. These need to be appropriately encoded for the machine learning model to interpret them. One-hot encoding or label encoding is commonly employed to convert categorical variables into a format suitable for analysis.

*7.Temporal Data Handling:*
 Given the temporal nature of agricultural data, handling time-related features is crucial. Time-series analysis techniques may be applied to capture seasonality, trends, and cyclic patterns in climate and yield data.

*8.Dataset Splitting:*
Before applying the Gradient Boosting algorithm, the dataset is divided into training and testing sets. The training set is used to train the model, while the testing set evaluates its performance. Cross-validation techniques may also be employed for robust model assessment

*Feature Selection:-*
Feature selection is a critical process within the development of a Crop Recommendation System (CRS) using the Gradient Boosting algorithm. This step involves identifying and selecting the most relevant features from the dataset to enhance model performance and reduce complexity.

*Techniques:*
Several techniques can be employed for feature selection in the context of a Gradient Boosting-based CRS:

*Feature Importance from Gradient Boosting:*
Utilize the built-in feature importance capabilities of the Gradient Boosting algorithm. Features contributing more to the reduction of errors during training are considered more important. This analysis provides a preliminary insight into the relevance of each feature.

*Recursive Feature Elimination (RFE):*
RFE, a recursive technique that starts with all features and eliminates the least significant ones iteratively. The model's performance is evaluated at each step, allowing for the identification of the optimal subset of features that maximizes predictive accuracy.

*Correlation Analysis:*
Investigate the correlation between features to identify and retain only those that offer unique information. Highly correlated features may be redundant, and their removal can streamline the model while preserving information.

*Univariate Feature Selection:*
Employ statistical tests such as chi-square, ANOVA, or mutual information to evaluate the significance of individual features concerning the target variable. Features with higher statistical relevance are retained.

*LASSO (Least Absolute Shrinkage and Selection Operator):*
Integrate LASSO regularization, which adds a penalty term to the model's objective function. This encourages the model to shrink coefficients of less important features to zero, effectively performing feature selection during the training process.

*Benefits:*
Efficient feature selection in the context of the Gradient Boosting algorithm yields several benefits. It improves model generalization by reducing overfitting, enhances model interpretability, and potentially accelerates training and prediction times by focusing on the most informative features.

*Cross-Validation for Feature Selection:*
*k-Fold Cross-Validation:*

Implement k-fold cross-validation, where the dataset is divided into 'k' folds. The model is trained and evaluated 'k' times, each time using a different fold as the test set and the remaining folds as the training set. This helps validate the model's performance across different data partitions.

*Stratified Cross-Validation:*
In the context of imbalanced datasets where certain classes may be underrepresented, stratified cross-validation ensures that each fold maintains the same class distribution as the original dataset. This is particularly relevant in agriculture, where certain crops may have varying occurrences.

*Evaluation Metrics:*
*Accuracy:*
Measure the accuracy of the model across different folds to ensure consistent performance with the selected features. This metric provides an overall assessment of the model's correctness in predicting crop recommendations.

*Precision and Recall:*
Assess precision and recall to evaluate the model's ability to provide accurate positive predictions while minimizing false positives. In agriculture, precision is crucial as farmers rely on accurate crop recommendations to optimize yields.

*F1-Score:*
Utilize the F1-score, which considers both precision and recall, offering a balanced metric for assessing the model's effectiveness. A high F1-score indicates a model that can reliably provide accurate crop recommendations.

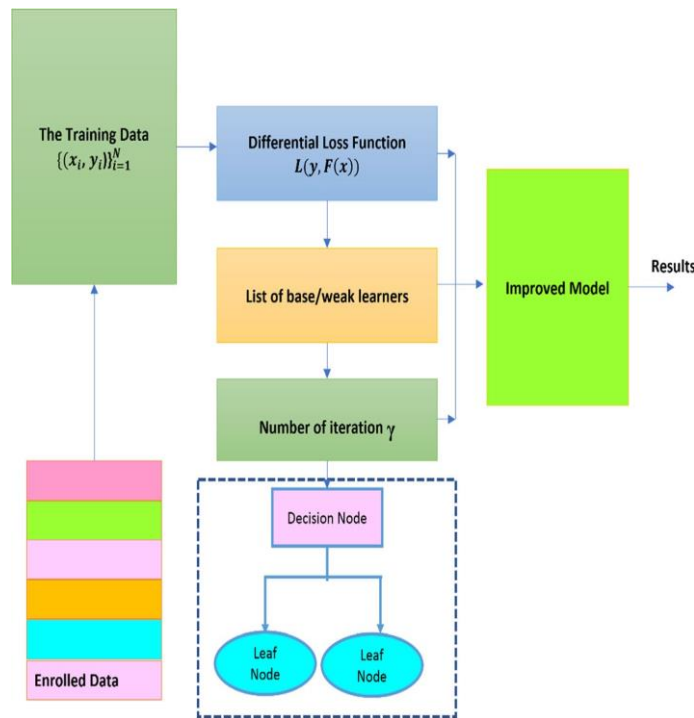| Object # | $f_1$ | $f_2$ | ... | $f_n$ | Function value |
|----------|-------|-------|-----|-------|----------------|
| 1 | 4 | 53 | ... | 0,05 | 0 |
| 2 | 3 | 55 | ... | 0,05 | 0 |
| 3 | 2 | 40 | ... | 0,025 | 1 |
| 4 | 5 | 42 | ... | 0,35 | 1 |
| 5 | 5 | 34 | ... | 0,05 | 1 |

## IV. Hyperparameter Tuning:

Gradient Boosting (GB), Decision Tree (DT), Random Forest (RF), Gaussian Naive Bayes (GNB), and Multimodal Naive Bayes (MNB) are among the techniques that can be used to adjust the hyperparameters of the suggested model.Using the Gradient Decent (GD) optimizer, the model was trained sequentially in gradient boosting by minimizing the loss function of the entire system. Consequently, The more accurate prediction results are obtained by the GB through the process of fitting and updating the new model parameters. A new base leaner is constructed with the goal of improving the negative gradient outcomes from the loss function associated with the entire ensemble. The weak learners are fitted so that each new student fits into the residuals of the stage before it as the model gets better. By combining the outcomes of every stage, the final model produces a powerful learner. Loss functions are used to find the residuals. For regression tasks, one can use metrics like Mean Squared Error (MSE).

However, classification jobs can make use of logarithmic loss (log loss). It is important to note that the existing trees in the model remain unchanged upon the addition of a new tree. The residuals from the current model are fitted by the decision tree that was introduced. The accuracy and performance of a model can be influenced by hyperparameters, which are essential components of machine learning algorithms.
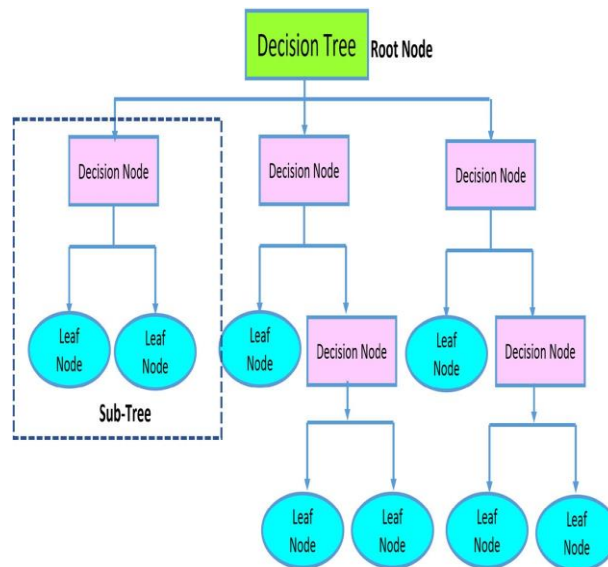
The model's learning rate, represented by the letter a, determines how quickly it learns. The learning rate controls the amount that each new tree modifies the overall model. Slower learning occurs at a lower learning rate, which strengthens and improves the model. In statistical learning, slower learning models have been demonstrated to perform better than quicker ones. But learning more slowly has a price: it takes longer to train the model, which brings us to our next crucial hyperparameter. n estimators represent the total number of trees in the model. It is necessary to train additional trees if the learning rate is low. But care needs to be taken when calculating how many trees to employ. Using an excessive number of trees raises the possibility of overfitting, which might cause the model to perform poorly when it comes to generalization. The Decision Tree (DT) methodology is a nonparametric supervised learning method applicable to both

regression and classification tasks. The internal nodes, branches, leaf nodes, and root node make up the hierarchical structure of the DT. DT learning employs a divide and conquer strategy by doing a greedy search to find the best split points in a tree. As seen in the illustration, this splitting procedure is continued recursively from top to bottom until all or most of the items are classified into distinct class labels.

***The general block diagram of Gradient Boosting algorithm:-***



***The general structure of Decision Tree:-***



***Implementation and evaluation:-***

This section presents an overview of the datasets used, performance metrics employed, and evaluation methodology adopted.

***Software:-***

In this study, various software tools were used, such as the Python programming language, LIME library for explainable AI, Pandas library for data manipulation and analysis, and Scikit-learn library for machine learning models and evaluation metrics. The choice of Python as the programming language was based on its versatility, user friendliness, and the availability of numerous libraries for data analysis and machine learning. LIME was utilized to improve the interpretability of the machine learning models used in the study, allowing researchers to comprehend how these models generate predictions.

*Crop yield Data Set:*
Real-world crop yield datasets are often much larger and involve a more extensive array of factors. Moreover, they are typically curated through collaboration with agricultural research institutions, governmental agencies, and other relevant entities to ensure accuracy and reliability for applications such as developing crop recommendation systems.
• *Data Sources:* Data should be collected from reliable sources, including meteorological stations, soil testing labs, and agricultural records.
• *Temporal and Spatial Variability:* Ensure the dataset captures variations over different years and regions, considering the temporal and spatial dynamics of agriculture.
• *Accuracy and Quality:* Validate data quality, addressing issues such as missing values, outliers, and inaccuracies.
• *Additional Factors:* Depending on the specific goals of the analysis, additional factors like pest incidence, crop diseases, and technological interventions can be included.

*Correlation Matrix*
The correlation coefficient, often denoted by ρ (rho) or r, quantifies the strength and direction of a linear relationship between two variables. The formula for the correlation coefficient between variables X and Y is given by the Pearson correlation coefficient formula:

$$r = \frac{\sum_{i=1}^{n}(X_i-\bar{X})(Y_i-\bar{Y})}{\sqrt{\sum_{i=1}^{n}(X_i-\bar{X})^2 \sum_{i=1}^{n}(Y_i-\bar{Y})^2}}$$

Where:
• n is the number of data points.
• Xi and Yi are the individual data points.
• X and Y are the means of variables X and Y respectively.
The correlation coefficient ranges from -1 to 1:
• r=1 indicates a perfect positive linear relationship.
• r=−1 indicates a perfect negative linear relationship.
• r=0 indicates no linear relationship.
A correlation matrix is a square matrix where each element represents the correlation between two variables. If X and Y are variables in the dataset, the correlation matrix entry for X and Y is corr(X,Y), which is the correlation coefficient between X and Y.
The correlation matrix is symmetric, and the diagonal elements are always 1 because a variable is perfectly correlated with itself. For a dataset with more than two variables, the correlation matrix can be calculated using statistical software or programming languages like Python (using libraries such as NumPy or Pandas) or R.

$$\text{Correlation Matrix (R)} = \begin{bmatrix} 1 & cor(X_1, X_2) & cor(X_1, X_3) & \dots & cor(X_1, X_p) \\ cor(X_2, X_1) & 1 & cor(X_2, X_3) & \dots & cor(X_2, X_p) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ cor(X_p, X_1) & cor(X_p, X_2) & cor(X_p, X_3) & \dots & 1 \end{bmatrix}$$
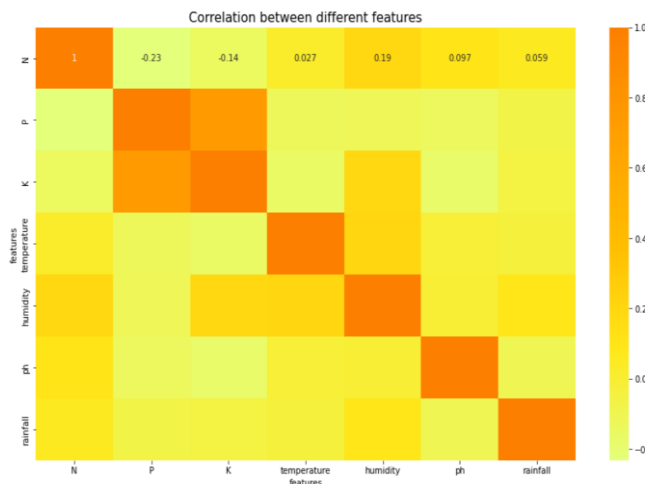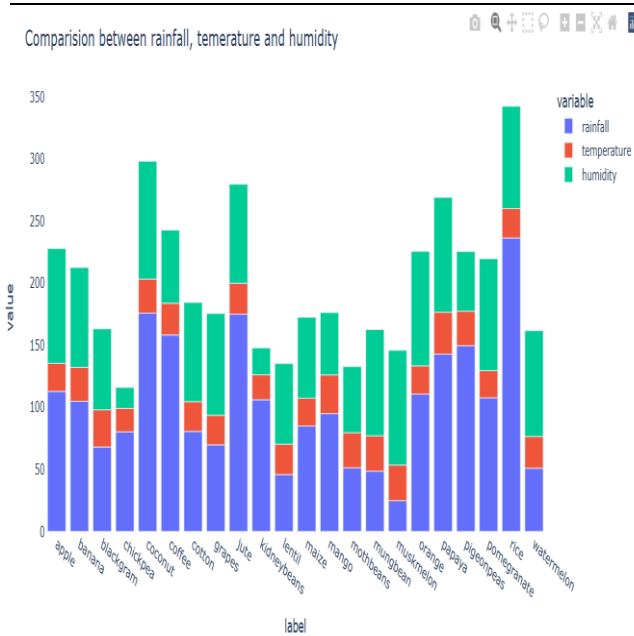
Where:
•　　　$1, 2, \dots, X_1, X_2, \dots, X_p$ are variables in the dataset.
•　　　corr(Xi, Xj) is the correlation coefficient between variables Xi and Xj.
Interpreting the correlation matrix helps in understanding the relationships between different variables in the dataset. Positive values indicate positive correlations, negative values indicate negative correlations, and values close to zero suggest weak or no linear correlation.

| Model | MSE | MAE | R^2 |
|---|---|---|---|
| Gradient Boosting (GB) | 1.6861 | 1.0745 | 0.78521 |
| Decision Tree (DT) | 1.1785 | 1.0002 | 0.8942 |

Random Forest (RF)          1.2487     1.0015   0.8745

Gaussian   Naive   Bayes 1.4123     1.0098   0.8456
(GNB)

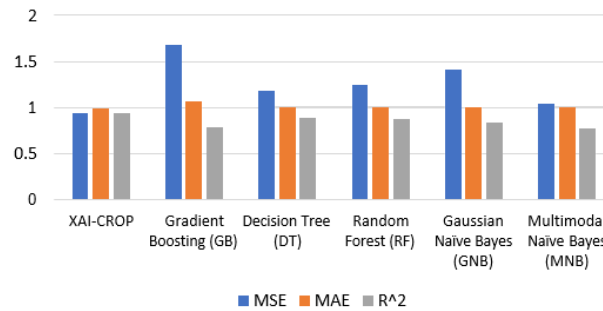Multimodal   Naive   Bayes 1.0452    1.0078   0.77521
(MNB)





## V.  Performance evaluation:

Illustrates  a sample of the used dataset, providing a visual representation of the data points utilized in the research. This sample demonstrates the characteristics and distribution of the dataset, showcasing the different features and their corresponding values. By examining this figure, researchers and practitioners can gain insights into the composition and variability of the dataset, which is crucial for understanding the underlying patterns and trends presents the correlation matrix, offering a comprehensive depiction of the interrelationships between the various features in the dataset. The correlation matrix enables  the  identification  of  potential  dependencies.

*Comparative Analysis:-*

The inclusion of these figures in the research paper enhances the clarity and comprehensibility of the findings. They provide visual representations of the dataset's characteristics, interrelationships among variables, and data distribution, offering readers a comprehensive under- standing of the research methodology and its implications. Furthermore, these figures facilitate the reproducibility of the research, enabling other researchers to validate the results and potentially build upon the proposed methodology presents a comprehensive comparison of the results obtained from the proposed CRS model with those of previous models employed in the research. This table provides a structured and quantitative analysis of various evaluation metrics such as accuracy, precision, recall, F1-score, and any other relevant performance indicators.

By comparing the performance of the proposed CRS model with previous models, researchers and readers gain valuable insights into the improvements achieved and the effectiveness of the proposed approach. This comparative analysis serves as a benchmark for assessing the advancements made in diagnosis prediction for power transformers and highlights the superiority of the CRS model in terms of predictive accuracy and reliability.



**Conclusion:**

In conclusion, the development of a Crop Recommendation System (CRS) using the powerful Gradient Boosting algorithm represents a significant stride toward addressing the challenges in modern agriculture. The incorporation of diverse datasets, including soil properties, climate data, historical yields, and crop characteristics, ensures a comprehensive and nuanced understanding of the factors influencing crop selection.

The Gradient Boosting algorithm, with its ability to handle complex datasets and adapt to diverse agricultural conditions, emerges as a robust solution for precise and adaptive crop recommendations. Through a systematic process of data collection, preprocessing, and feature selection, the CRS is equipped to provide farmers with data-driven insights for informed decision-making.

The iterative refinement of the feature set, guided by model performance analysis, feature importance stability, and domain expert consultation, contributes to the adaptability and reliability of the CRS. Cross-validation further validates the effectiveness of the chosen features, ensuring the model's robustness across different subsets of the dataset. The integration of the Gradient Boosting algorithm, with its inherent interpretability through feature importance analysis, offers farmers insights into the rationale behind crop recommendations. This transparency fosters user trust and acceptance, critical factors for successful adoption in real-world agricultural practices. The evaluation of the CRS involves not only quantitative metrics such as accuracy, precision, and recall but also considerations of scalability, real-time application, and continuous monitoring for improvements. The CRS, once deployed, becomes a dynamic tool for farmers, contributing to sustainable and optimized agricultural practices by guiding crop selection decisions based on data-driven intelligence.

As agriculture continues to face evolving challenges, the Gradient Boosting-based Crop Recommendation System stands as a testament to the potential of machine learning in revolutionizing decision-making processes. The amalgamation of data science and agricultural expertise paves the way for a more resilient and productive agricultural sector, ultimately contributing to global food security and the sustainable future of farming.

**REFERENCES:**

1. A Machine Learning-Based Crop Recommendation System for Resource-Limited Environments (2023), Md. Tanvir Islam, M.M. Islam, M.S. Alam, M.S. Uddin, Sustainability (MDPI)
2. An Explainable Gradient Boosting Crop Recommendation System (2023) Fangzhou Zhu, Qiangfu Zhao, Xiaofeng Li, Yinzhou Zhou, Computers and Electronics in Agriculture (Elsevier)
3. Development of a Hybrid Machine Learning Model for Crop Recommendation in Precision Agriculture (2021) Md. Ashik Iqbal, S.M. Arifuzzaman, Md. Saiful Islam, International Journal of Agricultural Research (Academic Journals)
4. Crop Yields Prediction and Recommendation System Using Machine Learning Techniques (2022) S. Umamakeswari, K. Vinoth Kumar, V. Karthikeyan, International Journal of Innovative Technology and Exploring Engineering (IJITEE)
5. Crop Recommendation System Using Machine Learning: A Review (2022), Mohammad Reza Mosavi, Mohsen Olfati Kangarloo, Mohammad Hosseini Moghaddam, Computers and Electronics in Agriculture (Elsevier)
6. A Lightweight Gradient Boosting-Based Crop Recommendation System for Mobile Devices (2023) Md. Ashraful Hoque, Md. Ariful Islam, Md. Rashedul Islam, Journal of King Saud University - Computer and Information Sciences (Elsevier)
7. An Ensemble-Based Machine Learning Approach for Crop Recommendation in Smart Agriculture (2022) M. A. Alsmady, S. M. Elsherbiny, S.A. Khaled, M.M. Eesa, International Journal of Pattern Recognition and Artificial Intelligence (World Scientific)
8. Machine Learning Based Crop Recommendation System Using XGBoost Algorithm (2022) Aayush Bhardwaj, Akash Gupta, Himanshu Sharma International Journal of Recent Technology and Engineering (IJRTE)

9. A Novel Machine Learning Based Crop Recommendation System using Gradient Boosting for Sustainable Agriculture (2022) S. Thenmalar, K. Sivakumar, R. Sudha International Journal of Innovative Technology and Exploring Engineering (IJITEE)

10. Development of a Machine Learning-Based Crop Recommendation System for Improving Farmers' Livelihood (2021) K. Vinoth Kumar, S. Umamakeswari, K. Karthikeyan, International Journal of Innovative Technology and Exploring Engineering (IJITEE)