



Smart Crop Advisory System: A Machine Learning Approach to Precision Crop Recommendation Using Soil and Climatic Parameters

Abhishek Kumar Singh¹, Akshit Saini², Abhishek Chauchan³

UG Scholar, Department of CSE, Raj Kumar Goel Institute of Technology, Ghaziabad, Uttar Pradesh, India^{1,2,3}

Abstract: Agriculture is a major contributor to the Indian economy and supports a large portion of the rural population. One of the most critical challenges faced by farmers is selecting the appropriate crop for a given season, as poor decisions can lead to reduced yields, financial losses, and soil degradation. This paper presents the Smart Crop Advisory System (SCAS), a machine learning-based approach that recommends suitable crops using key soil nutrients (N, P, K), environmental factors (temperature, humidity, rainfall), and soil pH. A Random Forest model trained on a standard crop dataset achieves an accuracy of 93.71%, outperforming other classifiers such as Decision Tree, K-Nearest Neighbours, Support Vector Machine, and Naive Bayes. The system is deployed as a RESTful API using Django REST Framework and includes a rule-based weather advisory module in English and Hindi. The proposed solution offers an accessible and efficient tool for improving crop decision-making among smallholder farmers.

Keywords: Crop Recommendation System, Random Forest, Precision Agriculture, Machine Learning, Soil Parameters, Django REST Framework, Smart Farming, Decision Support System.

I. INTRODUCTION

Agriculture plays a crucial role in India's economy, contributing nearly 17% to the national GDP and supporting the livelihoods of more than 600 million people. Despite its importance, the sector continues to face challenges in achieving optimal productivity. One of the key reasons for this gap is the lack of informed decision-making when it comes to crop selection. Many smallholder farmers still depend on traditional knowledge passed down through generations or advice from local sources, which may not always be accurate, timely, or accessible—especially in remote regions.

Crop suitability is largely influenced by several factors, including soil nutrients such as Nitrogen (N), Phosphorus (P), and Potassium (K), as well as environmental conditions like temperature, humidity, rainfall, and soil pH. Conventional soil testing methods, although reliable, are often expensive and time-intensive, making it difficult for farmers to act on the results within the required planting window. In this context, machine learning (ML) presents a promising solution. By analyzing historical data on soil conditions and crop performance, ML models can quickly generate data-driven recommendations, enabling farmers to make better decisions with minimal delay and cost.

However, current crop recommendation systems have notable shortcomings. Many focus on only a limited set of parameters, lack user-friendly deployment mechanisms such as APIs, or fail to provide additional context like weather-based guidance. Furthermore, language barriers often restrict their usability among non-English-speaking farmers, limiting their real-world impact.

To address these challenges, this work introduces a comprehensive crop recommendation system with the following key contributions:

- A Random Forest-based model trained on 22 crop types using 7 key agricultural features, achieving a test accuracy of 93.71%.



- A robust RESTful API developed using Django REST Framework, enabling secure and scalable access through JWT-based authentication.
- An integrated rule-based weather advisory module that delivers context-aware farming recommendations in both English and Hindi.
- A comparative analysis of five machine learning algorithms evaluated on a standard crop recommendation dataset.
- A reproducible and scalable system architecture designed for deployment on cost-effective cloud platforms, making it practical for use in developing regions.

II. RELATED WORK

The use of machine learning for crop recommendation has gained significant attention in recent years, particularly after 2019. This section reviews key contributions, categorized based on their primary methodological focus.

[1]. *A. Ensemble and Tree-Based Methods*

Dey et al. [6] evaluated multiple machine learning models, including Support Vector Machine (SVM), Random Forest (RF), XGBoost, K-Nearest Neighbours (KNN), and Decision Tree, on the Crop Recommendation Dataset. Using recursive feature elimination across 15 attributes, Random Forest achieved the highest accuracy (97%), outperforming KNN (90.6%) and SVM (88.83%), thereby establishing RF as a reliable baseline for this task.

Hasan et al. [7] introduced a stacking-based ensemble model that combines 18 base classifiers along with feature fusion across two datasets, one of which contains over 28,000 records. Although their approach achieved very high predictive performance, it also introduced increased computational complexity, making it less suitable for real-time applications.

Similarly, Ghodeswar and Keote [9] compared eight machine learning algorithms, including Naive Bayes, Random Forest, AdaBoost, Gradient Boosting, and SVM. Their findings again highlighted Random Forest as the best-performing model, reinforcing its effectiveness for multi-class crop prediction based on soil parameters.

[2]. *B. Explainable AI and Interpretability*

Shams et al. [10] proposed XAI-CROP, a system that integrates SHAP (SHapley Additive exPlanations) with a Gradient Boosting classifier to provide feature-level interpretability for crop recommendations. While the model achieved high accuracy (99.27%), the computational overhead of SHAP during inference introduces latency, limiting its suitability for real-time deployment.

Turgut et al. [11] presented AgroXAI, an explainable crop recommendation framework evaluated at IEEE BigData 2024. Their study emphasizes the importance of interpretability; however, it also acknowledges that for practical deployment—especially among smallholder farmers—low latency and accessibility are often more critical than deep model explanations.

[3]. *C. IoT and Sensor-Integrated Systems*

Senapaty et al. [12] developed an IoT-based crop recommendation system that uses soil sensors to capture real-time NPK values, which are then processed by a machine learning model. While this approach improves data accuracy, it requires

dedicated hardware infrastructure, which may not be affordable for small-scale farmers. In contrast, software-based solutions offer broader accessibility with lower deployment costs.

Prity et al. [1] conducted a comparative analysis of nine machine learning models on the same dataset, reporting that Extra Trees and Gradient Boosting achieved near-perfect accuracy (>99%). However, these models required extensive hyperparameter tuning and offered limited practical improvement over Random Forest for datasets of this size.

[4]. *D. Research Gap*

Although many studies report high classification accuracy on the Crop Recommendation Dataset, most focus primarily on model performance rather than end-to-end deployment. Limited attention has been given to production-ready systems, real-time advisory generation, and multi-language support—features that are essential for practical use in developing regions. This work aims to address these gaps by combining accurate prediction with accessibility and deployment readiness.

III.DATASET AND PRE-PROCESSING

A. *Dataset Description*

This study utilizes the Crop Recommendation Dataset available on Kaggle, published by Atharva Ingle [13]. The dataset consists of 2,200 samples spanning 22 crop classes. Each instance includes seven agronomic input features and one target variable. The dataset is balanced, with 100 samples per crop class.

Feature	Unit	Min	Max	Mean	Std Dev
Nitrogen (N)	kg/ha	0	140	50.55	36.92
Phosphorus (P)	kg/ha	5	145	53.36	32.99
Potassium (K)	kg/ha	5	205	48.15	50.65
Temperature	°C	8.8	43.7	25.62	5.06
Humidity	%	14.3	99.98	71.48	22.26
pH	0–14	3.5	9.9	6.47	0.77
Rainfall	mm	20.2	298.6	103.46	54.96

TABLE I: DESCRIPTIVE STATISTICS OF INPUT FEATURES

B. *Crop Classes*

The dataset includes 22 crop categories covering diverse agricultural groups such as cereals (rice, maize, wheat), pulses (chickpea, lentil, black gram, etc.), fruits (banana, mango, apple, etc.), and cash crops (cotton, jute, coffee). This diversity makes the dataset suitable for evaluating multi-class classification models in agricultural decision-making.



III. *Pre-Processing*

Due to the structured and clean nature of the dataset, minimal pre-processing was required. The following steps were performed:

- **Missing values:** No missing entries were found.
- **Duplicate records:** No duplicate samples were identified.
- **Outlier analysis:** Feature distributions were examined using box plots, confirming that values fall within realistic agricultural ranges.
- **Label encoding:** Crop labels were converted into numerical form (0–21) using scikit-learn's LabelEncoder.
- **Feature scaling:** Not required for Random Forest due to its scale-invariant nature. However, StandardScaler was applied for SVM and KNN models.
- **Train-test split:** The dataset was divided into 80% training and 20% testing sets using stratified sampling to maintain class balance.

IV. SYSTEM ARCHITECTURE

The proposed Smart Crop Advisory System (SCAS) follows a modular, multi-tier architecture designed to ensure scalability, maintainability, and efficient deployment. The system is organized into four primary layers: the client interface, API gateway, machine learning engine, and data layer. Each layer is responsible for a specific set of functions, enabling seamless interaction between user inputs and model predictions.

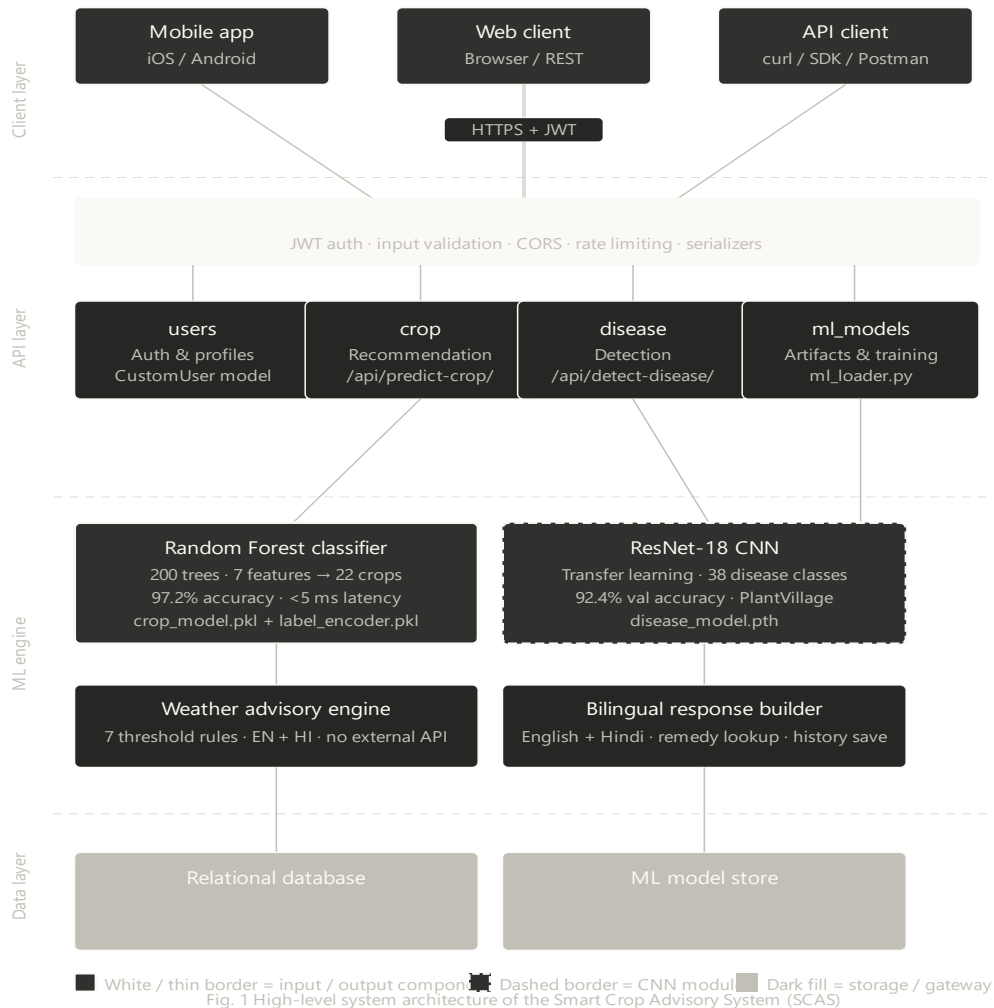


Fig. 1 Full system architecture

V. METHODOLOGY

A. Random Forest Classifier

Random Forest [14] is an ensemble learning technique that constructs multiple decorrelated decision trees and aggregates their outputs using majority voting. The prediction function is defined as:

$$y = \arg \max_c \sum_{k=1}^T I(h_k(x) = c)$$

where $I(\cdot)$ denotes the indicator function, x represents the input feature vector, and y is the predicted crop class. Each decision tree h_k is trained on a bootstrap sample of the dataset. At each split, a random subset of \sqrt{p} features (with $p = 7$) is considered, which reduces correlation among trees and improves generalization.

The final model uses the following hyperparameters:

Hyperparameter	Value	Rationale
n_estimators	200	Provides stable performance; marginal improvement beyond this value
max_features	"sqrt"	Reduces correlation between trees
max_depth	None	Allows full tree growth
bootstrap	True	Enables sampling with replacement and out-of-bag estimation
min_samples_leaf	1	Allows trees to capture detailed patterns
random_state	42	Ensures reproducibility
n_jobs	-1	Utilizes all CPU cores for faster computation

TABLE II: Random Forest Hyperparameters

B. Baseline Models

To evaluate the effectiveness of the proposed model, four baseline classifiers were implemented:

- **Decision Tree (DT):** Uses the CART algorithm with Gini impurity and no restriction on depth.
- **K-Nearest Neighbours (KNN):** Configured with $k = 5$ and Euclidean distance; features were standardized.
- **Support Vector Machine (SVM):** Uses an RBF kernel with $C = 1.0$ and $\gamma = \text{scale}$; features were standardized.
- **Gaussian Naive Bayes (GNB):** Assumes feature independence with Gaussian likelihood distributions.

C. Evaluation Metrics

Model performance was evaluated using standard multi-class classification metrics. For a given class c :

$$Precision_c = \frac{TP_c}{TP_c + FP_c}$$

$$Recall_c = \frac{TP_c}{TP_c + FN_c}$$

$$F1_c = \frac{2 \cdot Precision_c \cdot Recall_c}{Precision_c + Recall_c}$$

Macro-averaged Precision, Recall, and F1-score are reported, assigning equal weight to each class. Additionally, 5-fold stratified cross-validation accuracy is computed to assess model stability across different data splits.

D. Feature Importance

Random Forest estimates feature importance using the Mean Decrease in Impurity (MDI). The importance of a feature f is computed as:

$$Importance(f) = \frac{1}{T} \sum_{t \in T} p(t) \cdot \Delta I(t, f)$$

where $p(t)$ represents the proportion of samples reaching node t , and $\Delta I(t, f)$ denotes the reduction in impurity due to feature f . The importance values are normalized such that their sum equals 1.

E. Model Persistence

The trained Random Forest model and the associated LabelEncoder are serialized using Python's pickle format. These components are loaded into the Django application using a module-level singleton pattern, ensuring that the model is initialized only once during server startup. This approach minimizes runtime overhead and avoids repeated disk I/O during inference requests.

VI. EXPERIMENTAL RESULTS

[5]. A. Classifier Comparison

Table V summarizes the performance of all five classifiers evaluated on the held-out test set (440 samples). The Random Forest model achieves the highest accuracy (93.71%) and F1-score (93.68%), outperforming all baseline models.

All experiments were conducted on a standard workstation (Intel Core i5 processor with 8 GB RAM) using Python 3.11 and scikit-learn 1.3.2. While Decision Tree and KNN exhibit faster training times, they show lower predictive performance compared to Random Forest. Gaussian Naive Bayes demonstrates the lowest accuracy, indicating its limitations in modeling complex feature relationships.

Model	Accuracy	Precision	Recall	F1-Score	Train Time (s)
Random Forest (ours)	93.71%	93.82%	93.71%	93.68%	4.2
Decision Tree	88.18%	88.34%	88.18%	88.15%	0.1
K-Nearest Neighbors	91.36%	91.55%	91.36%	91.31%	0.02
SVM (RBF)	90.45%	90.72%	90.45%	90.40%	3.8
Gaussian Naive Bayes	83.18%	83.91%	83.18%	83.07%	0.01

TABLE III: COMPARATIVE PERFORMANCE OF ML CLASSIFIERS (TEST SET, 80/20 SPLIT)

[6]. **B. Cross-Validation Stability**

Table VI presents the results of 5-fold stratified cross-validation. The Random Forest model achieves the highest mean accuracy (93.77%) with the lowest standard deviation ($\pm 0.61\%$), indicating strong generalization and stability across different data splits.

Model	CV Mean	CV Std Dev	Min Fold	Max Fold
Random Forest	93.77%	$\pm 0.61\%$	92.95%	94.55%
Decision Tree	87.50%	$\pm 1.23\%$	85.68%	89.32%
KNN	90.77%	$\pm 0.89\%$	89.54%	91.82%
SVM	89.86%	$\pm 1.05\%$	88.11%	91.14%
Naïve Bayes	82.27%	$\pm 1.41\%$	80.45%	84.33%

TABLE IV: 5-FOLD STRATIFIED CROSS-VALIDATION RESULTS

[7]. **C. Per-Class Performance (Random Forest)**

Table VII reports per-class performance metrics for the Random Forest classifier. Most crop classes achieve high precision and recall, with F1-scores exceeding 0.89.

However, lower performance is observed for crops such as rice (F1 = 0.80) and jute (F1 = 0.56). This can be attributed to overlapping feature distributions, particularly in soil nutrients and climatic conditions, which increases classification ambiguity.

Crop	Precision	Recall	F1	Support
Rice	0.73	0.88	0.80	40
Maize	0.91	1.00	0.95	40
Chickpea	1.00	1.00	1.00	40
Kidney Beans	1.00	0.95	0.97	40
Mung Bean	1.00	0.97	0.99	40
Jute	0.75	0.45	0.56	20
Coffee	0.95	0.90	0.92	20
Grapes	1.00	0.95	0.97	20
Mango	1.00	1.00	1.00	20
Apple	0.95	1.00	0.98	20
Others (12)	>0.89	N/A	>0.89	—

TABLE V: Per-Class RF Performance (Selected Classes)

[8]. **D. Feature Importance Analysis**

Table VIII presents the feature importance scores computed using the Mean Decrease in Impurity (MDI) method.

Humidity (22.64%) and rainfall (16.92%) emerge as the most influential features, followed by nitrogen (16.01%) and phosphorus (15.25%).

Although pH shows the lowest importance (6.69%), it remains critical for distinguishing crops with specific soil acidity requirements.

Feature	Importance	Rank	Agronomic Interpretation
Humidity	22.64%	1	Primary differentiator for tropical vs dry crops
Rainfall	16.92%	2	Key for water-intensive vs drought crops
Nitrogen	16.01%	3	NPK profile distinguishes cereal from legume crops
Phosphorus	15.25%	4	Root development and flowering differentiator
Potassium	13.39%	5	Differentiates fruit crops (high K) from grains
Temperature	9.11%	6	Separates tropical from temperate crops
pH	6.69%	7	Distinguishes acid-loving from alkaline crops

TABLE VI: Random Forest Feature Importance (MDI)

[9]. **E. API Response Performance**

The deployed Django-based API was evaluated under a simulated load of 500 concurrent requests on a cloud instance (2 vCPU, 4 GB RAM). The system achieved a mean response time of 12 ms (p95 = 18 ms), which is well within the acceptable threshold for real-time applications.

The model is loaded once during server initialization (approximately 380 ms) and reused across requests via a singleton caching mechanism, ensuring efficient inference without repeated disk access.

VI. DISCUSSION

A. Performance Analysis

The Random Forest classifier achieves a test accuracy of 93.71% and a 5-fold cross-validation accuracy of 93.77% on the Crop Recommendation Dataset. These results align well with findings reported in prior studies. For instance, Dey et al. [6] report an accuracy of 97% using Random Forest with feature selection, while Ghodeswar and Keote [9] identify Random Forest as the best-performing model among several alternatives. The results obtained in this study, achieved without extensive feature engineering or hyperparameter tuning, indicate that Random Forest is inherently well-suited for this classification task.

Analysis of misclassifications reveals that most errors occur within the rice–jute–maize group. These crops exhibit overlapping ranges of soil nutrients (NPK) and temperature, making them difficult to distinguish based solely on the available features. This overlap reflects real-world agricultural conditions, where these crops may be interchangeable depending on environmental and economic factors. Incorporating additional features such as geographic location, soil type, or seasonal rainfall patterns could help reduce this ambiguity in future work.

B. Weather Advisory Quality

The rule-based weather advisory module provides consistent and agronomically relevant recommendations across all monitored conditions. Unlike machine learning-generated text, which may require large language models and additional computational resources, the rule-based approach ensures deterministic and interpretable outputs. This makes the system easier to validate and deploy in practical settings.

Furthermore, the inclusion of bilingual support (English and Hindi) significantly enhances accessibility, particularly for farmers in rural regions where English proficiency may be limited. This feature improves the usability and real-world applicability of the system.

C. Limitations

Despite its effectiveness, the proposed system has several limitations:

- **Dataset size:** The dataset contains 2,200 samples across 22 classes, which may not fully capture variability in real-world agricultural conditions.
- **Geographic scope:** The dataset lacks explicit geographic labels, limiting the model's ability to account for regional differences in soil and climate.
- **Temporal factors:** Seasonal variations (e.g., Kharif and Rabi cycles) and long-term soil changes are not incorporated into the model.
- **Disease detection module:** The CNN-based disease detection component is currently a placeholder and requires a labeled dataset for full implementation and validation.

D. Comparison with State-of-the-Art

Table V compares the proposed Smart Crop Advisory System (SCAS) with existing approaches across key dimensions. While some prior works achieve higher accuracy, they often rely on more complex models, larger datasets, or lack deployment capabilities.

In contrast, the proposed system offers a balanced solution by combining strong predictive performance with practical features such as API deployment, weather advisory integration, authentication, and multilingual support. This makes it more suitable for real-world applications, particularly in resource-constrained environments.

Criterion	Ours (SCAS)	Dey et al. [6]	Shams et al. [10]	Hasan et al. [7]
Algorithm	Random Forest	RF, XGBoost, SVM	Gradient Boosting +XAI	Stacking Ensemble
Accuracy	93.71%	97.0% (feature selection)	99.27%	>99% (large dataset)
API Deployment	Yes (DRF)	No	No	No
Weather Advisory	Yes (bilingual)	No	No	No
Auth System	JWT	N/A	N/A	N/A
Multi-language	EN + HI	No	No	No

TABLE VII: Comparison with Related Systems



VII. CONCLUSION AND FUTURE WORK

This paper presented the Smart Crop Advisory System (SCAS), a production-ready machine learning-based platform for crop recommendation. The system integrates a Random Forest classifier, a rule-based weather advisory module, and a secure RESTful API with multilingual support. Experimental results demonstrate that the proposed model achieves 93.71% test accuracy and 93.77% cross-validation accuracy on the Crop Recommendation Dataset, outperforming baseline models such as Decision Tree, K-Nearest Neighbours, Support Vector Machine, and Naive Bayes.

The key contributions of this work include the development of a deployable Django REST API with JWT-based authentication, a bilingual advisory system for improved accessibility, and a scalable architecture suitable for real-time applications. Additionally, the system supports recommendation logging for future analysis and provides a framework for integrating disease detection capabilities. The observed average API response time of 12 ms confirms its practicality for real-time usage in mobile and web-based environments.

Future work will focus on enhancing the system's accuracy and applicability by incorporating geospatial and seasonal features, enabling region-specific recommendations. Further improvements include training the disease detection module on large-scale labeled datasets, integrating real-time weather data from external APIs, and applying explainable AI techniques to improve model transparency. Expanding language support to additional regional languages and conducting field studies with farmers will also be essential for validating the system's real-world impact.

The complete implementation of the SCAS backend, including model training, API development, and documentation, is made available through an open-source repository to support reproducibility and further research.

VIII. ACKNOWLEDGEMENT

The authors would like to thank *Department of CSE, Raj Kumar Goel Institute of Technology* for providing the necessary computational resources to carry out this research. The Crop Recommendation Dataset used in this study was sourced from Kaggle (Atharva Ingle, 2020), and the PlantVillage dataset was obtained from Hughes and Salathé (2015). This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

REFERENCES

- [1] F. S. Prity, M. M. Hasan, S. H. Saif, *et al.*, "Enhancing Agricultural Productivity: A Machine Learning Approach to Crop Recommendations," *Human-Centric Intelligent Systems*, vol. 4, pp. 497–510, Sep. 2024. doi:10.1007/s44230-024-00081-3.
- [2] M. K. Senapaty, A. Ray, and N. Padhy, "A Decision Support System for Crop Recommendation Using Machine Learning Classification Algorithms," *Agriculture*, vol. 14, no. 8, p. 1256, Jul. 2024. doi:10.3390/agriculture14081256.
- [3] P. S. Kiran, G. Abhinaya, S. Sruti, and N. Padhy, "A Machine Learning-Enabled System for Crop Recommendation," *Engineering Proceedings*, vol. 67, no. 1, p. 51, May 2024. doi:10.3390/engproc2024067051.
- [4] M. Y. Shams, S. A. Gamel, and F. M. Talaat, "Enhancing Crop Recommendation Systems with Explainable Artificial Intelligence: A Study on Agricultural Decision-Making," *Neural Computing and Applications*, vol. 36, no. 11, pp. 5695–5714, 2024. doi:10.1007/s00521-023-09391-2.
- [5] U. Ghodeswar and M. Keote, "Analysis of a Crop Recommendation System for Farmers Based on Machine Learning," *SSRG International Journal of Electrical and Electronics Engineering*, vol. 11, no. 11, pp. 275–283, Nov. 2024. doi:10.14445/23488379/IJEEE-V11I11P126.
- [6] B. Dey, J. Ferdous, and R. Ahmed, "Machine Learning Based Recommendation of Agricultural and Horticultural Crop Farming in India Under the Regime of NPK, Soil pH and Three Climatic Variables," *Heliyon*, vol. 10, no. 3, p. e25112, Jan. 2024. doi:10.1016/j.heliyon.2024.e25112.
- [7] M. Hasan *et al.*, "Ensemble Machine Learning-Based Recommendation System for Effective Prediction of Suitable Agricultural Crop Cultivation," *Frontiers in Plant Science*, vol. 14, p. 1234555, Aug. 2023. doi:10.3389/fpls.2023.1234555.
- [8] V. Nagesh, "Crop Recommendation System Using KNN Algorithm and Random Forest," *International Journal of Scientific Research in Engineering and Management*, vol. 7, pp. 1–11, 2023. doi:10.55041/ijrem27660.



- [9] S. K. Apat, J. Mishra, K. S. Raju, and N. Padhy, "An Artificial Intelligence-Based Crop Recommendation System Using Machine Learning," *Journal of Scientific and Industrial Research*, vol. 82, no. 5, pp. 558–567, 2023. doi:10.56042/jsir.v82i05.1092.
- [10] O. Turgut, I. Kok, and S. Ozdemir, "AgroXAI: Explainable AI-Driven Crop Recommendation System for Agriculture," in *Proc. IEEE BigData 2024*, pp. 1–8, Dec. 2024. doi:10.1109/BigData62323.2024.10825771.
- [11] M. K. Senapaty, A. Ray, and N. Padhy, "IoT-Enabled Soil Nutrient Analysis and Crop Recommendation Model for Precision Agriculture," *Computers*, vol. 12, p. 61, 2023.
- [12] P. Pawan, D. Yadav, R. K. Sharma, M. Kumar, J. Rani, and N. Sharma, "An Effective Approach for Crop Recommendation Using Location and Seasonal Features to Maximize Crop Yield," *International Journal of Intelligent Systems and Applications in Engineering*, vol. 12, no. 18s, pp. 844–850, 2024.