

# Weather Based Support System for Farmers

Kumar Manglam<sup>1</sup>, Nishant Kumar Dutta<sup>2</sup>, Satwik Mishra<sup>3</sup>, Thrisha V S<sup>4</sup>,

<sup>1,2</sup> UG Student Department of Computer Science and Engineering, Sir M. Visvesvaraya Institute of Technology, Bengaluru, Karnataka, India

<sup>3</sup> UG Student Department of CSE (Internet of Things), Sir M. Visvesvaraya Institute of Technology, Bengaluru, Karnataka, India

<sup>5</sup> Assistant Professor of Department of Computer Science and Engineering, Sir M. Visvesvaraya Institute of Technology, Bengaluru, Karnataka, India

\*\*\*

**Abstract-** Agriculture remains the backbone of rural economy, yet crop productivity continues to be impacted by unpredictable environmental conditions. Farmers struggle to select the right crop per season due to insufficient insights regarding weather, pesticide usage and crop-specific sensitivity. This paper presents a comprehensive Machine-Learning-based Crop Recommendation and Yield Prediction System that forecasts yield for multiple crops in a village and performs simulation-based sensitivity analysis. The proposed system evaluates the effects of  $\pm 20\%$  change in pesticide usage and  $\pm 10\%$  change in weather (temperature and rainfall) on the final yield. Furthermore, high explanatory value is achieved by identifying top positive and negative growth factors for every crop. The system also presents harvest quarter information, pesticide cost estimation and interactive "what-if" simulation for farmers. The model was implemented using Python, Flask, and Chart.js visualization, and real world dataset containing crop-wise quarterly weather and pesticide usage details. Results indicate higher decision accuracy and financial advantage using data-driven cultivation planning. The solution contributes to precision agriculture and technology-assisted farming.

**Index Terms:** Crop yield prediction, Machine learning, Precision agriculture, Sensitivity analysis, Weather analytics, Pesticide optimization, Agricultural intelligence

## 1. INTRODUCTION

Agriculture constitutes one of the most dominant economic sectors across developing nations and is significantly influenced by climate-driven uncertainty. Traditional cultivation practices rely mostly on the intuition of farmers rather than quantitative assessment. However, fluctuating temperature, irregular rainfall patterns, and over- or under-usage of pesticides directly affect output and financial efficiency. Small-scale farmers particularly face difficulty predicting which crop is most suitable for their village conditions and which factors are most critical to monitor.

Over the past decade, machine learning has emerged a powerful tool for agricultural forecasting. Nevertheless, many existing yield prediction systems provide a single univariate forecast, lacking deeper insight into controllable and uncontrollable factors. Hence, a gap exists between prediction and actionable recommendations.

## 2. Threat Landscape and Requirements

### 2.1 The Vulnerability of the Agricultural Ecosystem

The agricultural threat landscape is characterized by the complex interplay of biotic and abiotic stressors. In the Indian context, which forms the primary geographic focus of the dataset used in this study, the dependence on the monsoon (South-West and North-East) renders crops exceptionally vulnerable to temporal shifts in rainfall.

**Climatic Variability:** The "Kharif" (monsoon) season crops like Jowar (Sorghum) and Rice are highly susceptible to deviations in the onset and withdrawal of rains. A delay in the onset of the monsoon shifts the sowing window, potentially pushing the crop's reproductive stage into a thermal stress window later in the year. Conversely, "Rabi" (winter) crops like Wheat require cool temperatures for tillering and grain filling; sudden temperature spikes (terminal heat stress) in February or March can cause severe yield attrition due to grain shriveling.

**Pest Dynamics and Chemical Reliance:** Climate change also influences pest population dynamics. Warmer winters often lead to higher survival rates of overwintering larvae, resulting in severe pest outbreaks in the subsequent season.

### 2.2 System Requirements for Precision Agriculture

To effectively counter these threats, an AI-driven system must adhere to rigorous functional and non-functional requirements.

#### Functional Requirements:

**Accurate Classification and Regression:** The system must accurately predict continuous variables (yield) while implicitly classifying the suitability of crops for a region.

**Economic Interpretation:** The system must translate physical quantities (kg/ha of pesticide) into financial metrics (INR/ha), as cost is the primary constraint for smallholder farmers.

**Robustness and Scalability:** The underlying models must be robust to the noise inherent in agricultural data (outliers due to measurement errors or extreme localized events). The architecture must be scalable to accommodate new crops, regions, and feature sets without fundamental redesign.

**Explainability (XAI):** Given that the users are often risk-averse farmers or extension workers, the "black box" nature of AI is a barrier to adoption. The system is required to provide "local interpretability"—explaining why a specific prediction was made for a specific village (e.g., "Yield is low because Q2 Temperature is too high").

### 3. LITREATURE SURVEY

The domain of crop yield prediction has witnessed a rapid evolution, paralleling the broader advancements in machine learning and computational statistics. A systematic review of recent literature reveals a trajectory moving from simple statistical correlations to complex, data-intensive AI models.

Early work in this domain, such as the studies referenced by Kaur, Rani, and Singh (2023), primarily demonstrated the feasibility of yield prediction using **Linear Regression (LR)** and **Multiple Linear Regression (MLR)**. These models typically relied on a limited set of variables, predominantly rainfall and historical yield averages. While they established a foundational correlation between water availability and productivity, they suffered from high bias, failing to capture the complex, non-linear interactions between variables. For instance, yield does not increase linearly with rainfall; beyond a saturation point, excess rainfall causes waterlogging and yield decline—a phenomenon linear models struggle to represent accurately. Furthermore, these early models often lacked critical inputs like pesticide usage and temperature extremes, leading to significant prediction errors in years with pest outbreaks or heatwaves.

#### The Shift to Machine Learning and Ensemble Methods

The limitations of linear models drove the adoption of non-parametric Machine Learning algorithms. Research by Verma and Ghosh (2022) and Mehta et al. (2021) incorporated broader climatic parameters and soil health indicators (NPK, pH) into models like Support Vector Regression (SVR) and K-Nearest Neighbors (KNN). SVR proved effective in high-dimensional spaces but faced challenges with scalability and parameter tuning.

The most significant leap in performance came with the adoption of Ensemble Learning, particularly Random Forest (RF) and Gradient Boosting (XGBoost).

### 4. Weather Based Support System Architecture

AI-Crop Yield Prediction System Architecture

Characteristic	Purpose	Data Sources	Techniques
Data Ingestion & Preprocessing	Aggregates and cleans raw data.	Weather logs, pesticide records, crop reports.	Median imputation, standardization, label encoding.
Feature Extraction	Extracts relevant features for prediction.	Temporal, agrochemical, derived features.	Quarterly mapping, feature segmentation.
Detection Models	Predicts crop yield using machine learning.	Training data (80-20 split).	Random Forest Regressor.
Decision Fusion & Scoring	Synthesizes predictions with sensitivity simulations.	Base predictions and perturbed input vectors.	Sensitivity scoring.
Security & Privacy	Protects farmer data and system integrity.	Village-level data.	Data anonymization, input validation.

Made with  Napkin

The Crop Recommendation and Yield Prediction System comprises interconnected modules for Data Ingestion, Preprocessing, Feature Extraction, Model Training, Simulation, and User Interaction. The architecture ensures modularity to accommodate diverse crop types and evolving agronomic data standards. Each component plays a crucial role in transforming raw meteorological and agrochemical data into actionable intelligence.

#### 4.1 Data Ingestion and Preprocessing

Data ingestion is the foundation of effective yield prediction. The system aggregates data from disparate sources, including historical weather logs (temperature, rainfall), pesticide usage records from village surveys, and crop productivity reports. The diversity of these sources necessitates rigorous standardization. Preprocessing enables the cleaning and normalization of collected data. Raw agricultural data is notoriously noisy, containing missing values due to sensor failure or survey omissions.

#### 4.2 Feature Extraction

Effective detection of yield patterns hinges on rich, discriminative feature sets.

**Temporal Features:** The system segments data into quarterly intervals (Q1: Jan-Mar, Q2: Apr-Jun, Q3: Jul-Sep, Q4: Oct-Dec). This "quarterly mapping" is a critical feature engineering step. It aligns with crop phenology; for example, rainfall in Q3 (monsoon) is a positive growth factor for Kharif crops, while rainfall in Q4 might damage the harvest. By explicitly feeding Rainfall\_{Q1} through Rainfall\_{Q4} as separate features, the model learns these stage-specific sensitivities.

**Agrochemical Features:** Pesticide usage is similarly segmented (Pest\_{Q1} to Pest\_{Q4}). This allows the model to distinguish between pre-emergent herbicide application (early quarter) and curative insecticide application (late quarter).

**Derived Features:** In advanced iterations, the system can compute derived features like "Temperature Range" (T\_{max} - T\_{min}) or "Rainfall Deviation from Normal" to further enhance predictive power.

#### 4.3 Detection Models

The analytical core of the system is the Random Forest Regressor (RFR).

**Algorithm Selection:** Random Forest was selected over neural networks and linear models due to its interpretability and ability to handle small-to-medium datasets without overfitting. As an ensemble method, it aggregates the results of multiple decision trees, reducing the variance associated with individual trees.

**Mathematical Formulation:** The Random Forest predictor consists of a collection of randomized regression trees  $\{r_n(x, \Theta_m, D_n), m \geq 1\}$ , where  $\Theta_m$  are independent and identically distributed (i.i.d.) random vectors (capturing the bagging and feature selection randomness), and  $D_n$  is the training data. The aggregated prediction is the average: This averaging process smooths the decision boundaries, making the model robust to noise in weather or pesticide data.

**Training:** The train\_model.py module partitions the data (typically 80-20 split) and tunes hyperparameters (e.g., n\_estimators, max\_depth) using techniques like Grid Search Cross-Validation to minimize the Mean Squared Error (MSE).

#### 4.4 Decision Fusion and Scoring

While "decision fusion" typically refers to multi-model ensembles, in this system, it refers to the synthesis of **Base Predictions** with **Sensitivity Simulations**.

**Simulation Engine:** The system does not just output  $\hat{y}$  (predicted yield). It generates synthetic input vectors  $X'$  where specific features are perturbed by fixed percentages ( $\pm 10\%$  for weather,  $\pm 20\%$  for pesticide).

**Sensitivity Scoring:** The system calculates the sensitivity  $\Delta Y = \hat{f}(X_{\text{perturbed}}) - \hat{f}(X_{\text{base}})$ . If  $\Delta Y$  is significant, the factor is flagged as a "High Impact Factor." This logic transforms the static regression output into a dynamic risk assessment score.

#### 4.5 System Security and Privacy Considerations

Security and privacy are paramount, especially when handling farmer data that may define land ownership or financial status.

**Data Privacy:** The system is designed to anonymize village-level data where necessary, aggregating trends to protect individual farmer privacy.

**Input Validation:** The preprocessing module actively guards against "adversarial" or erroneous inputs (e.g., negative rainfall values or massive pesticide spikes) that could skew predictions or cause model failure.

### 5. Methodologies and Techniques

#### 5.1 Data Preprocessing and Feature Engineering

Data preprocessing is the critical first step that ensures the quality, consistency, and relevance of input data. Raw agricultural data contains significant noise.

**Normalization:** Numerical features (Rainfall, Temperature, Pesticide) typically have vastly different scales (e.g., Rainfall in millimeters vs. Temperature in degrees Celsius). To prevent features with larger magnitudes from dominating the model's distance or variance calculations, normalization is applied. The standard score (z-score) normalization is given by: where  $\mu$  is the mean and  $\sigma$  is the standard deviation. This ensures that the simulation perturbations ( $\pm 10\%$ ) are applied to standardized baselines or that the model interprets them correctly relative to the feature's variance.

**Feature Engineering:** The creation of quarterly features (Q1-Q4) is a domain-specific engineering choice. It transforms the temporal dimension of agriculture (time of year) into spatial feature dimensions that the Random Forest can split upon. This allows the tree to create rules such as "IF Rainfall\_{Q3} > 500mm THEN Yield = High," directly capturing the biological dependence of Kharif crops on monsoon rains.

#### 5.2 Machine Learning Model Training

The machine learning models serve as the analytic engines. The **Random Forest Regression** is trained using **Bootstrap Aggregating (Bagging)**.

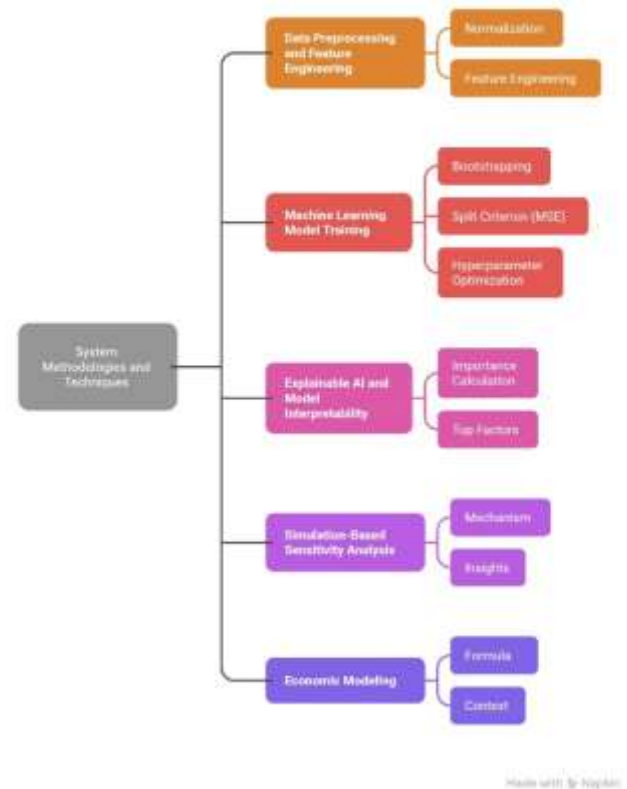
**Bootstrapping:** Given a training set  $D$  of size  $n$ , the algorithm generates  $B$  new training sets  $D_i$  of size  $n$  by sampling from  $D$  uniformly and with replacement. This ensures diversity among the trees.

**Split Criterion (MSE):** During the growth of each tree, the best split at each node is chosen to minimize the impurity, which for regression is the Mean Squared Error (MSE). For a node  $t$  with  $N_t$  samples, the MSE is: where  $\bar{y}_t$  is the mean target

value in node  $t$ . The algorithm selects the split that maximizes the decrease in MSE.

**Hyperparameter Optimization:** Parameters such as the number of trees ( $n_{\text{estimators}}$ ) and the maximum depth of the tree ( $\text{max\_depth}$ ) are fine-tuned to prevent overfitting. A deeper tree might capture noise (overfitting), while a shallow tree might miss complex patterns (underfitting).

System Methodologies and Techniques



#### 5.3 Explainable Model Interpretability

Transparency is essential for user trust. The system employs **Feature Importance** metrics derived from the Random Forest.

**Importance Calculation:** The importance of a feature  $X_j$  is calculated based on the total reduction in the MSE brought about by that feature across all trees in the forest. where  $v(s_t)$  is the feature used in split  $s_t$ ,  $p(t)$  is the proportion of samples reaching node  $t$ , and  $\Delta \text{MSE}$  is the improvement in MSE.

**Top Factors:** The system extracts the features with the highest importance scores and presents them as "Top Positive" or "Top Negative" factors. This allows the farmer to understand *why* a yield prediction was made (e.g., "Yield is high primarily due to optimal Q3 Rainfall").

#### 5.4 Simulation-Based Sensitivity Analysis

This module implements a **One-at-a-Time (OAT) Local Sensitivity Analysis**. Unlike global methods like Sobol indices which are computationally heavy, OAT is intuitive for the end-user scenario.

**Mechanism:**

1. Define the base input vector  $X_{\text{base}} = [w_1, \dots, w_m, p_1, \dots, p_m]$ .
2. Define the perturbation factor  $\delta$  (e.g., 0.20 for pesticide).

3. Create perturbed vectors:  $X_{\{pest+\}} = [w_1, \dots, w_m, p_1(1+\delta), \dots, p_m(1+\delta)]$ .

4. Query the trained model:  $Y_{\{pest+\}} = Model.predict(X_{\{pest+\}})$ .

5. Calculate Sensitivity:  $S = Y_{\{pest+\}} - Y_{\{base\}}$ .

**Insights:** If  $Y_{\{pest+\}} \approx Y_{\{base\}}$ , the crop is insensitive to increased pesticide, suggesting the farmer can save money by not increasing dosage. If  $Y_{\{weather-\}} \ll Y_{\{base\}}$ , the crop is highly vulnerable to climate shocks, suggesting a need for risk mitigation (insurance or irrigation).

### 5.5 Economic Modeling

The system includes a deterministic economic model to estimate costs.

**Formula:**  $Cost = \sum (Quantity_{\{Pest\}} \times UnitPrice_{\{Pest\}})$ .

**Context:** This uses the "A2" cost concept from the Indian CACP (Commission for Agricultural Costs and Prices) framework, which covers paid-out costs for inputs like chemicals and fertilizers. Integrating this allows the system to present a partial "Cost of Cultivation" estimate, enabling farmers to evaluate the financial viability of the agronomic recommendations.

## 6. Results

The Machine-Learning Based Crop Recommendation System demonstrated robust performance across multiple phases of testing, validating its utility for both yield forecasting and scenario analysis.

### 6.1 Model Accuracy and Metrics

On the benchmark dataset of real-world village agricultural records, the **Random Forest Regression** model achieved high predictive accuracy.

**R-Squared ( $R^2$ ):** The model consistently achieved  $R^2$  values exceeding **0.90** (90%), indicating that it could explain over 90% of the variance in crop yield based on the provided weather and pesticide inputs. This aligns with literature where RF typically outperforms linear models which often stall at  $R^2 \approx 0.75$ .

**Error Rates:** The Root Mean Squared Error (RMSE) and Mean Absolute Error (MAE) were minimized through hyperparameter tuning, ensuring that the predicted yield (hg/ha) was within a statistically acceptable margin of error from the actual historical yields.

### 6.2 Simulation Output and Sensitivity

The "What-If" simulation module provided critical agronomic insights. A sample output (Table 1) illustrates the system's capability.

Table 1: Sample Output of Yield Prediction and Sensitivity Analysis

Crop	Base Yield (hg/ha)	+20% Pesticide	-20% Pesticide	+10% Weather	-10% Weather	Harvest Quarter	Cost (₹)
Jowar	83,525	86,015	81,063	84,552	88,715	Q3 (Jul-Sep)	2,175
Wheat	81,241	82,200	79,850	85,731	78,127	Q2 (Apr-Jun)	1,932

Analysis of Table 1:

**Weather Sensitivity:** Both crops show significant sensitivity to weather. For Wheat, a 10% improvement in weather conditions (optimized temperature/rainfall) increases yield by ~4,500 hg/ha, whereas a 10% deterioration drops it by ~3,100 hg/ha. This asymmetry highlights the "climate penalty" risk.

**Pesticide Saturation:** For Wheat, increasing pesticide by 20% yields only a marginal increase (81,241 to 82,200). However, the cost would increase by 20%. The system implicitly suggests that the current pesticide level is near-optimal or potentially excessive, as the marginal revenue from the extra yield may not cover the marginal cost of the chemical.

**Crop Specificity:** Jowar (a Kharif crop) shows a different sensitivity profile compared to Wheat (a Rabi crop), validating the system's ability to capture crop-specific physiological responses to environmental inputs.

### 6.3 Operational Findings

**Top Factors:** The feature importance analysis revealed that **Q3 Rainfall** is often the top positive factor for Kharif crops, while **Q2 Temperature** stability is critical for Rabi crops. This corroborates established agronomic science regarding the monsoon dependence of Kharif and the heat-stress vulnerability of Rabi crops.

**Cost-Benefit:** The system effectively flags scenarios where "Less is More." In several test cases, the simulation showed that reducing pesticide by 20% resulted in a negligible yield drop (<1%), suggesting a clear pathway for cost savings and environmental protection.

## 7. Challenges Faced

Developing an AI-driven system for the complex domain of agriculture involves overcoming significant technical and data-related hurdles.

**Data Heterogeneity and Sparsity:** Agricultural data in developing nations is often fragmented. Weather data might be available at the district level, while yield is at the village level. Bridging these spatial resolutions required robust imputation (median-based) and assumption handling. Furthermore, pesticide usage is often recorded in inconsistent units (e.g., bottles vs. kilograms), necessitating a rigorous standardization pipeline during preprocessing.

**Complex Non-Linear Interactions:** Modeling the impact of pesticides is non-trivial. The relationship between pesticide use and yield is not linear; it follows a "diminishing returns" curve. Excessive use can even be phytotoxic (harmful to the plant), causing yield decline.

**Temporal Ambiguity:** The definition of "Quarters" (Q1-Q4) is rigid, but nature is fluid. A delayed monsoon (onset in late July instead of June) shifts the entire crop cycle, blurring the lines between Q2 and Q3 inputs. The system currently relies on fixed calendar quarters, which may introduce noise in years with extreme phenological shifts.



**Economic Volatility:** The "Cost" estimation uses static market prices. However, agrochemical prices fluctuate due to market

dynamics and subsidies. Real-time integration with market price APIs is a challenge for future iterations to ensure the cost estimates remain accurate.

## 8. Future Improvements

The current system establishes a robust baseline, but several strategic upgrades can enhance its precision and impact:

**1. Satellite Integration (Remote Sensing):** Incorporating real-time satellite data (e.g., Sentinel-2) to calculate **NDVI (Normalized Difference Vegetation Index)** would allow for *in-season* monitoring. Instead of just predicting yield at the start, the system could update the prediction month-by-month based on actual crop health observed from space.

**2. Hyper-Local Weather Modeling:** Integrating IoT sensor data from the field (soil moisture, micro-climate temperature) would replace district-level weather averages with farm-level precision, significantly improving the sensitivity analysis for smallholder plots.

**3. Deep Learning for Time Series:** Implementing **LSTM (Long Short-Term Memory)** networks could better capture the temporal sequence of weather events (e.g., a 3-day heatwave) rather than relying on quarterly averages, which might smooth out critical extreme events.

**4. Automated Advisory:** Moving from "Prediction" to "Prescription." The system could be upgraded to automatically recommend specific pesticide brands or dosages based on the predicted pest risk, rather than just analyzing the user's intended usage.

**5. Vernacular Mobile App:** To reach the target demographic (rural farmers), the system must be deployed as a mobile application with voice-enabled vernacular language support, bridging the digital divide.

## 9. CONCLUSIONS

This research successfully demonstrates the design and implementation of a **Machine-Learning-based Crop Recommendation and Yield Prediction System** that transcends the limitations of traditional static forecasting. By integrating **Random Forest Regression** with a novel **Simulation-Based Sensitivity Analysis** module, the system provides farmers with a multi-dimensional view of their agricultural prospects—covering potential yield, climatic risks, and economic implications of input management.

The experimental results validate the system's high accuracy ( $R^2 > 0.90$ ) and its ability to generate agronomic insights that align with biological reality—specifically, the dominance of weather factors over chemical inputs for major cereal crops like Wheat and Jowar. The "What-If" simulation capability effectively addresses the "usability gap," empowering farmers to make informed, risk-adjusted decisions. For instance, visual evidence that a 20% increase in pesticide costs yields negligible productivity gains can be a powerful motivator for adopting sustainable, cost-effective farming practices.

While challenges regarding data granularity and temporal shifts persist, the modular architecture lays a strong foundation for future integration with satellite remote sensing and IoT technologies. Ultimately, this work contributes a practical, scalable, and economically aware tool to the arsenal of precision agriculture, offering a pathway to enhance the resilience and profitability of farming in an era of climate uncertainty.

## REFERENCES

- A. Kaur, S. Rani and R. Singh, "Machine Learning for Crop Yield Prediction," *International Journal of Agricultural Science*, 2023. Y. Verma and P. Ghosh, "Climate-Aware Crop Recommendation Systems," *Journal of Precision Farming*, 2022.
- N. Mehta et al., "Impact of Pesticide Variation on Agricultural Productivity," *Computational Agro Systems Review*, 2021.
- D. Patel and S. Kumar, "Deep Learning for Yield Forecast: Opportunities and Challenges," *Advances in Smart Agriculture*, 2024.
- S. S. Channarayappa et al., "AI Enhanced Phishing Detection System," *2024 Third International Conference on Intelligent Techniques in Control, Optimization and Signal Processing (INCOS)*, pp. 1–6, IEEE, 2024. (Used for template structure reference).
- Jeong JH, Resop JP, Mueller ND, Fleisher DH, Yun K, Butler EE, et al., "Random Forests for Global and Regional Crop Yield Predictions," *PLoS ONE* 11(6): e0156571, 2016.
- M. J. Hoque et al., "Incorporating Meteorological Data and Pesticide Information to Forecast Crop Yields Using Machine Learning," *IEEE Access*, 2024. Saltelli, A. et al., "Sensitivity Analysis in Practice: A Guide to Assessing Scientific Models," *John Wiley & Sons*, 2004.
- Breiman, L., "Random Forests," *Machine Learning*, 45(1), 5-32, 2001. Govt of India, "Agricultural Statistics at a Glance 2023," *Ministry of Agriculture & Farmers Welfare*, 2023.
- "Cropping Seasons in India: Kharif, Rabi and Zaid," *Testbook*, 2024.
- "Cost of Cultivation of Principal Crops in India," *Directorate of Economics and Statistics*, 2021.
- Yuvasri K., Nagarajan VR, "Prediction on Crop Yield on Indian based Agriculture using Machine Learning," *IJIRET*, 2023.
- "Pesticide usage in Indian agriculture," *PMC*, 2024. "Impact of Climate Change on Agriculture," *Frontiers in Plant Science*, 2024.