

<sup>1</sup>Shilpa Choudhary<sup>2</sup>Sandeep Kumar<sup>3</sup>Monali Gulhane,<sup>4</sup>Nitin Rakesh

## Empowering Agriculture: Leveraging Intelligent Systems for Sustainable Farming



**Abstract:** - Climate change, market volatility, water scarcity, and pest control issues are just a few of the challenges that farmers face. Limited access to resources and education, as well as the slow adoption of new technologies, aggravate these problems. Creating intelligent systems is necessary to handle these pressing issues. This method makes use of key components such as crop type, weather, and soil properties to forecast agricultural yields. This system offers significant yield prediction capabilities for 53 crops by utilizing complex algorithms including XG-Boost, random forest, and decision tree. It takes into account crucial factors to accurately assess crop production potential, such as temperature, rainfall, nutrient levels (N, P, and K), soil pH, and temperature data. Cutting-edge machine learning techniques look at past data and trends to provide farmers with crucial information they need to make wise decisions. In order to solve the problems farmers, encounter in the modern agricultural sector, this intelligent system aims to enhance resource efficiency, farming systems' resilience and production, and the use of resources.

**Keywords:** Crop Yield Prediction, Data Preprocessing, Decision Tree, Random Forest, XG Boost.

### I. INTRODUCTION

Numerous problems confronting modern farmers put livelihoods, food security, and environmental sustainability in agriculture under jeopardy. Global agricultural systems are confronted with complex issues such as pest and disease outbreaks, water scarcity, market volatility, and climate change [1]. We must employ innovative strategies that make use of state-of-the-art technologies to address these issues and raise agricultural output, resilience, and sustainability [2–3]. Our research work presents a complex model for predicting crop yield that incorporates multiple variables, such as meteorological conditions, soil characteristics, crop features, and historical yield data, with cutting-edge machine learning techniques [4]. The objective of this model is to provide accurate and timely crop yield projections to help farmers, policymakers, and other agricultural supply chain stakeholders make educated decisions and lower crop production risks. In order to comprehend the complex relationships between input factors and crop yields, the model integrates a number of machine learning techniques [5–6].

Our method makes use of algorithms' predictive power to estimate crop yields—both specialty and staple crops—in various agroclimatic zones. We evaluate the model's performance using two main assessment metrics: accuracy and R2 score [7-8]. While accuracy shows the percentage of cases that are correctly predicted, R2 score shows how much of the variance in the target variable is explained by the model. We conduct extensive testing and validation of the model in different crops and regions to assess its accuracy, stability, and longevity [9]. Our study contrasts various machine learning techniques to identify the best models for agricultural yield prediction. We evaluate each algorithm's strengths and weaknesses and provide insights into how they might be applied in different agricultural contexts [10]. The practical insights our study offers for improved resource allocation, agricultural management techniques, and resistance to climatic variability and market uncertainty have a substantial impact on agricultural stakeholders [11]. Our crop production forecast model is useful for enhancing agricultural decision-making and advancing sustainable growth in the agricultural sector. It makes use of contemporary machine learning techniques and combines a variety of datasets.

<sup>1</sup>Department of Computer Science and Engineering, Neil Gogte Institute of Technology, Hyderabad, India

<sup>2</sup>Department of Computer Science and Engineering, Koneru Lakshmaiah Education Foundation, Vijayawada, A.P., India

<sup>3</sup>Department of Computer of Computer Science and Engineering, Symbiosis Institute of Technology, Nagpur Campus, Symbiosis International (Deemed University), Pune, India

<sup>4</sup>Department of Computer of Computer Science and Engineering, Symbiosis Institute of Technology, Nagpur Campus, Symbiosis International (Deemed University), Pune, India

Email: <sup>1</sup>shilpachoudhary2020@gmail.com, <sup>2</sup>er.sandeepsahratia@gmail.com, <sup>3</sup>monali.gulhane4@gmail.com, <sup>4</sup>nitin.rakesh@gmail.com

Corresponding Mail: <sup>3</sup>monali.gulhane4@gmail.com, <sup>4</sup>nitin.rakesh@gmail.com

## II. LITERATURE SURVEY

Through the analysis of several datasets and algorithms, machine learning forecasts agricultural production by taking farming practices, weather patterns, and soil characteristics into account. In order to produce accurate forecasts, machine learning algorithms analyze historical data and relevant attributes to find patterns and connections. In order to assist farmers and other stakeholders in making educated decisions on planting, harvesting, and resource distribution, regression, ensemble learning, and deep learning approaches are commonly employed to estimate crop yield. This predictive ability optimizes agricultural operations, enhancing productivity, resource efficiency, and economic results in the farming sector.

With minimal ground truth data, Y. Alebele et al.[1] suggest Gaussian kernel regression for rice yield estimate from optical and SAR imaging. The approach performs better than Bayesian linear inference and probabilistic Gaussian regression. Combining RDVI1 with interferometric coherence at the heading stage yields the best prediction accuracy. The study makes the case for the benefits of combining optical indices with satellite interferometric coherence for mapping agricultural yield using Gaussian kernel regression. In order to estimate crop yields on the Canadian Prairies, J. Liu et al. [2] compared crop metrics derived from Terra/MODIS to known yields for spring wheat, canola, and barley. According to the results, vegetation indices at the height of growth were superior to GPP or NPP as yield predictors, while EVI2 outperformed NDVI. The models demonstrated stability across a range of years, however there were interannual variations. Annual agricultural yields were mapped at the polygon level of Soil Landscapes of Canada using the best-performing models [2]. For the purpose of predicting crop yield, Jhajharia, Kavita, et al.[3] employed decision trees, random forests, XGBoost regression, CNN, and extended short-term memory networks. CNN and random forests have the lowest loss and maximum accuracy, respectively. Compared to existing algorithms, a model that accurately predicts crop yield has been devised.

The long-term sustainability of agriculture is threatened by P. Sharma et al.'s attempt to anticipate crop yield utilizing factors including rainfall, crop, climatic conditions, area, production, and yield [4]. Machine learning and deep learning are used in crop yield prediction to determine crop production and growing season. Regression methods such as XGBoost, random forest, and decision trees are applied; random forest and convolutional neural networks exhibit superior accuracy. To comprehend errors, the model is compared to alternative methods and examined using the root mean square error measure. An Automated Rice Crop Yield Prediction utilizing a Sine Cosine Algorithm with a Weighted Regularized Extreme Learning Machine (SCA-WRELM) is introduced by S. Thirumal and R. Latha [5]. It produces better prediction results by using the WRELM model for yield prediction along with a min-max data normalization technique. A dataset of rice yield is used to evaluate the method.

In their discussion of machine learning's application in agriculture, V. K. G. Kalaiselvi et al.[6] point out both potential and problems. Crop growth estimates can be improved by incorporating real-time data from IoT sensors. Fifteen algorithms are evaluated in the study; a newly feature-enhanced algorithm achieves 99.59% classification accuracy. This could lead to increased production rates, reduced costs, and more resilient infrastructure. The results may also aid in the early detection of illnesses by farmers, improve crop productivity, and lower food costs. A. k. Gajula et al.[7] focuses on soil quality detection to predict crop suitability, cultivation requirements, and yield, aiding farmers in better planning and increasing production. It also provides fertilizer requirements. However, the model is limited by data and does not consider climatic disasters. Future improvements may include geospatial analysis.

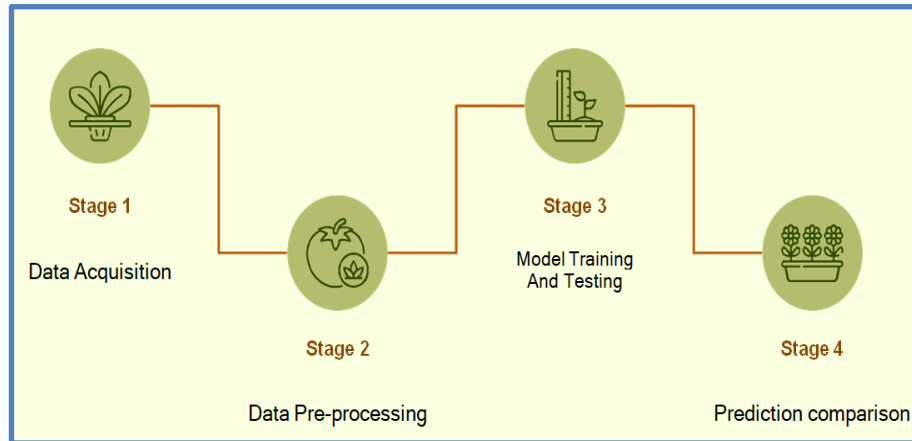
## III. PROPOSED WORK

The proposed method begins with preparing the dataset according to specific requirements, followed by rigorous data preprocessing steps, including merging disparate datasets and handling missing values. Subsequently, a comprehensive comparative analysis of various machine learning algorithms is conducted to evaluate their performance. This analysis involves training and testing multiple models using the prepared dataset to assess their predictive capabilities. A detailed description of the proposed methodology is elaborated in Figure 1.

### 3.1 Data Acquisition

We had to use several datasets in order to predict crop yields with any degree of accuracy. For yield estimates to be accurate, a single dataset needs to contain all the features and factors. So, in order to precisely match the requirements of our research, we started a thorough process of creating the dataset. We merged numerous datasets,

each of which included distinct and essential data required to predict crop yields with any degree of accuracy. Data on rainfall, temperatures, crop features, and crop productivity in various states and seasons were all included in the databases related to agriculture. Table I contains the datasets comprehensive descriptions.



**Figure. 1: Architectural Diagram of Proposed Work**

**Table I: Various Dataset & Descriptions**

Sr. No.	Dataset Name	Description	Features
1	States with Season Rainfalls	Season-wise rainfall data is categorized by state.	State, Season, Rainfall
2	States with Seasonal Temperatures	Season-wise temperature data categorized by states	State, Season, Temperature
3	Crop Parameters (n, p, k, pH values)	Parameters crucial for a specific crop's growth	Crop, n, p, k, pH
4	State Season Crop Area Production	Crop production data across states and seasons	State, Season, Crop, Area, Production

### 3.2 Data Preprocessing

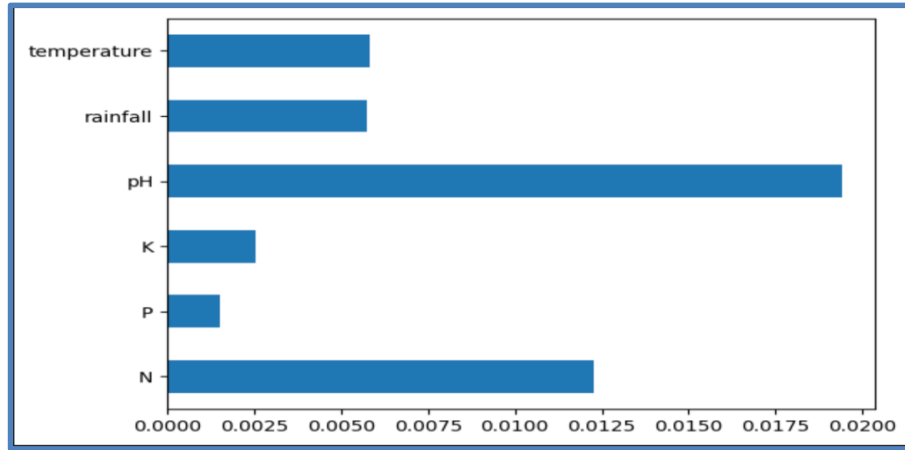
#### *Merging of different Datasets*

Standardizing the data format across all datasets comes next after data preprocessing. In order to guarantee interoperability and seamless dataset integration, standardization is crucial. Following confirmation of data correctness, the merging process starts. The datasets comprising rainfall and temperature information are combined by matching the shared attributes of state and season. This integration produces a unified dataset containing extensive weather-related data. The combined dataset is enhanced by combining it with the crop parameter dataset. The merging process is made more accessible by the shared characteristics of the crop. The dataset gains strength by including crop factors like nitrogen, phosphorus, potassium, and pH levels, offering essential insights into the growth needs of specific crops. The comprehensive dataset is combined with the dataset comprising crop area production information. Completing this last merging process, all pertinent datasets are integrated to create a complete dataset ready for analysis to estimate crop yield. The integrated features of different datasets are shown in Table II, and their significance is mentioned in the Figure 2.

**Table II: Features Integrated from Various Datasets**

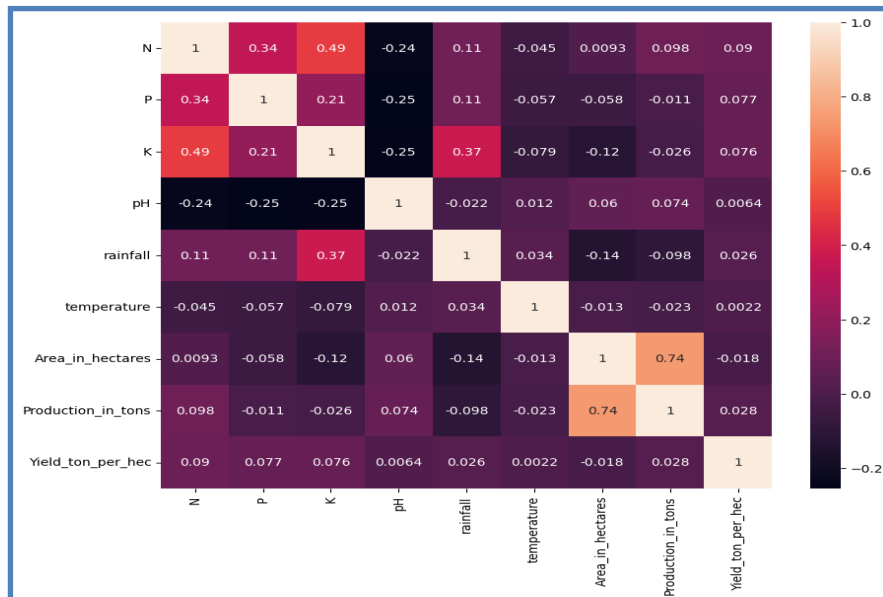
Sr. No.	Features	Description
1	State	Represents the state where the data corresponds.
2	Season	Represents the specific season (e.g., summer, winter).
3	Rainfall	Seasonal rainfall values for each state.
4	Temperature	Seasonal temperature values for each state.

5	Crop	Represents the specific crop.
6	n, p, k, pH	The parameter value for the crop.
7	Area	Crop area for a specific state and season.
8	Production	Crop production for a specific state and season.



**Figure. 2: Significant Features and Their Importance**

Modern farmers face numerous challenges in agriculture. Within crop yield prediction, a correlation heat map is crucial for researchers to investigate complex correlations between agricultural parameters and crop yields. Figure 3 illustrates the correlation coefficients between variables, including meteorological conditions, soil qualities, and crop traits. Researchers can analyze the heat map to find essential relationships and detect highly correlated variables, which can help select features and detect multicollinearity. The heat map helped us to explore data, allowing us to discover patterns and trends in the information. The insights obtained from the heat map improve model interpretation, enabling researchers to create more precise and dependable predictive models for agricultural decision-making and resource allocation.



**Figure. 3: Correlation Heat map**

**Handling Missing Values**

Addressing any missing data that resulted from the merger is essential before proceeding. Methods for managing missing data and maintaining data quality and integrity include statistical imputation and the removal of partial information. Once any issues with missing data have been resolved, the complete dataset is ready for analysis.

### 3.3 Prediction using various Classifiers

Our methodology used a combination of classifiers and thoroughly assessed their performance to determine the most efficient models. Our talk will concentrate on the Random Forest Regressor (RF Regressor) and XGBoost Regressor, the top-performing models due to limited space. We optimized these classifiers by tweaking hyperparameters and achieved higher accuracy than other models in our evaluation.

#### *Random Forest Regressor (RF Regressor)*

Random Forest Regressor (RF Regressor) is an ensemble learning technique that relies on decision trees. It creates numerous decision trees during training and combines their predictions to generate a final result. The technique incorporates randomness by training each tree on a random portion of the training data and evaluating a random subset of characteristics for splitting at each node. Introducing randomization aids in diminishing overfitting and enhancing generalization performance.

#### Algorithm-1

**Step 1:** Selecting N samples from the training set with replacement can be represented as.

$$D_t = \{(X_i, y_i)\}_{i=1}^N$$

**Step 2:**

- Choose at random a portion of the training data for each tree in the forest.
- Splitting nodes based on a feature j and split value s to minimize some impurity measure (e.g., mean squared error):

$$Q(D, j, S) = \frac{|D_{\text{left}}|}{|D|} \text{MSE}(D_{\text{left}}) + \frac{|D_{\text{right}}|}{|D|} \text{MSE}(D_{\text{right}})$$

**Step 3:**

- Predicting the target value for a new sample test  $X_{\text{test}}$  by averaging predictions of all trees:

$$Y_{\text{test}} = \frac{1}{T} \sum_{t=1}^T \text{Tree}_t(X_{\text{test}})$$

**Step 4:** Display the ultimate forecast and minimize the overall loss function.

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - Y_i)^2$$

```
RandomizedSearchCV(cv=5, estimator=RandomForestRegressor(), n_jobs=1,
                  param_distributions={'max_depth': [5, 10, 15, 20, 25, 30],
                                      'max_features': ['auto', 'sqrt'],
                                      'min_samples_leaf': [1, 2, 5, 10],
                                      'min_samples_split': [2, 5, 10, 15,
                                                          100],
                                      'n_estimators': [40, 50, 60, 70, 80, 90,
                                                      100, 110]},
                  random_state=42, scoring='neg_mean_squared_error',
                  verbose=2)
```

**Figure. 4: Hyperparameter specifications of Random Forest Regressor**

#### *XGBoost Regressor*

XGBoost Regressor is a gradient-boosting implementation specifically created for regression purposes. The algorithm constructs a series of weak learners, usually decision trees, one after the other, with each new learner correcting the errors made by the previous ones. XGBoost utilizes a gradient-based optimization approach to

minimize a loss function and enhance the model's predictive accuracy. The XGBoost Regressor's hyper-parameters consist of the learning rate (eta), maximum tree depth (max\_depth), minimal child weight (min\_child\_weight), training instance subsample ratio (subsample), and regularization parameters lambda and alpha.

### Algorithm-2

**Step 1:** Set up model parameters such as learning rate (eta), maximum tree depth (max\_depth), minimal child weight (min\_child\_weight), subsample ratio for training instances (subsample), and regularization parameters like lambda and alpha.

**Step 2:** Initialize the model prediction as the mean of the target values:

$$Y_i = \text{mean}(y_i) \text{ for all } i = 1, 2, 3, \dots, N.$$

**Step 3:** Calculate the initial prediction residuals: ( $r_i = y_i - y_i$ ).

For each Boosting round  $k = 1, 2, 3, \dots, K$ .

- Compute the gradient of the loss function with respect to the residuals:

$$g_{ik} = \left. \frac{\partial L(y_i, Y_i)}{\partial Y_i} \right|_{Y_i = Y_i^{(k-1)}}$$

$Y_i^{(k-1)}$  represents the prediction of the ensemble up to round (k-1)

- Compute the second-order gradient (Hessian) of the loss function with respect to the residuals:

$$h_{ik} = \left. \frac{\partial^2 L(y_i, Y_i)}{\partial Y_i^2} \right|_{Y_i = Y_i^{(k-1)}}$$

- Fit a regression tree to the negative gradients  $-g_{ik}$  using the training data

$$\{(x_i, g_{ik})\}_{i=1}^N \text{ as target}$$

- Compute the leaf weights for each leaf node in the tree using the following formula:

$$w_j = - \frac{\sum_{i \in I_j} g_{i,k}}{\sum_{i \in I_j} h_{i,k} + \lambda}$$

where  $I_j$  represents the set of training samples falling into leaf node  $j$ , and  $\lambda$  is the regularization parameter.

- Update the ensemble model prediction for each sample  $i$ :

$$Y_{\cdot i}^{(k)} = Y_{\cdot i}^{(k-1)} + \eta \sum_{j=1}^J w_j I(x_i \in R_j)$$

- Update the residuals for each sample:

$$r_{\cdot i}^{(k)} = r_{\cdot i}^{(k-1)} - \eta \sum_{j=1}^J w_j I(x_i \in R_j)$$

**Step 4:** Calculate the final prediction by summing the predictions from all trees in the ensemble.

$$Y_{test} = \sum_{k=1}^K Y_{test}^{(k)}$$

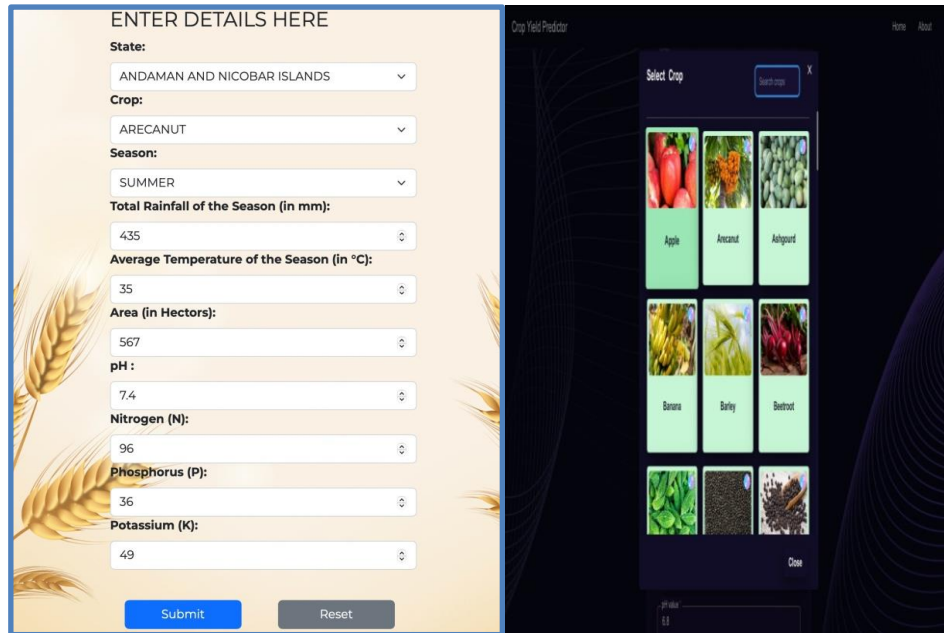
```
param_grid={
    'max_depth': range(4,10),
    'learning_rate': np.logspace(-3, 0, 50),
    'subsample': np.linspace(0.5, 1, 50),
    'n_estimators': range(40, 100, 10)
}
xgb_cv=RandomizedSearchCV(estimator=xgb_random,param_distributions=param_grid,
scoring='neg_mean_squared_error',n_iter=50,cv=5,n_jobs=-1,random_state=42)
```

Figure. 5: Hyperparameter specifications of XGBoost Regressor

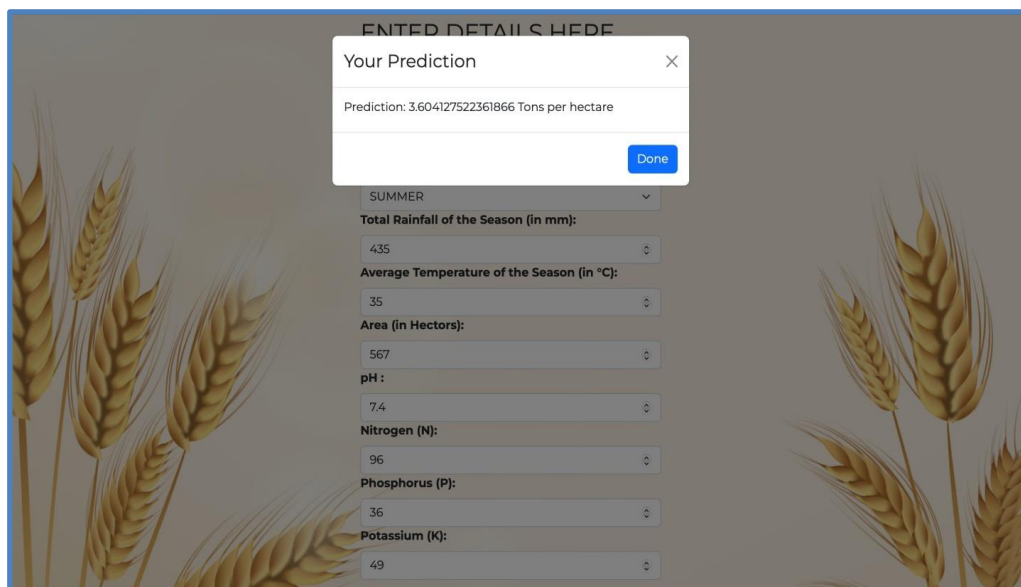
We improved the RF Regressor's and XGBoost Regressor's prediction accuracy by carefully adjusting their hyper-parameters (as shown in Figure 4 & 5), allowing us to use these models effectively for regression tasks. We verified their supremacy in identifying the fundamental patterns in the data and providing dependable predictions for our application through thorough assessment and comparison.

#### IV. RESULTS & DISCUSSIONS

A GUI has been created (as shown in Figure 6&7) to improve user engagement with the model's prediction algorithm. Users can enter precise information, including state, crops, seasons, total rainfall, average temperature, area, pH, nitrogen, phosphorus, and potassium. The graphical user interface (GUI) offers a user-friendly platform for smooth interaction with the prediction system, enhancing the overall user experience and usability.



**Figure 6: GUI for Crop Yield Prediction**



**Figure 7: Predicted Values from Machine Learning Model**

We assessed many models in our crop yield prediction study using two primary assessment metrics: R2 score and accuracy. The R2 score quantifies the amount of variance in the target variable that the independent variables can explain, whereas accuracy quantifies the percentage of correctly categorized cases in a classification task. Regarding evaluation measures, the tuned XG Boost model with Randomized Search and KNN\_CV consistently outperfor

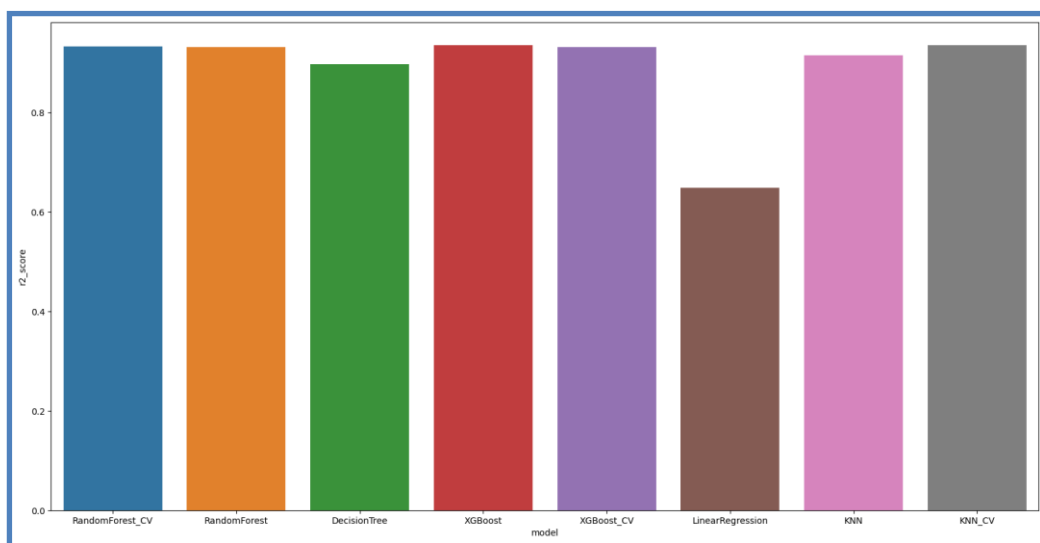


med other models. With an accuracy of 0.937 and a noteworthy R2 score of 0.937, the modified XG Boost model with Randomized Search demonstrated its high performance in accurately predicting crop yields. Findings are displayed in Figures 8 and 9.

With an accuracy of 0.934 and an R2 score of 0.934, the KNN model with Cross-Validation performed similarly, demonstrating its dependability in crop yield predicting under varied conditions. The findings highlight the importance of optimizing and fine-tuning models to improve forecast accuracy. By utilizing advanced techniques like Randomized Search and Cross-Validation and modifying hyper-parameters, we were able to increase the precision and accuracy of our crop yield prediction models. The best-performing models, according to our analysis, were the optimized XG Boost model with Randomized Search and the KNN with Cross-Validation model. These models have the potential to increase agricultural productivity and enhance crop management decision-making. The outcomes are displayed in Table III.

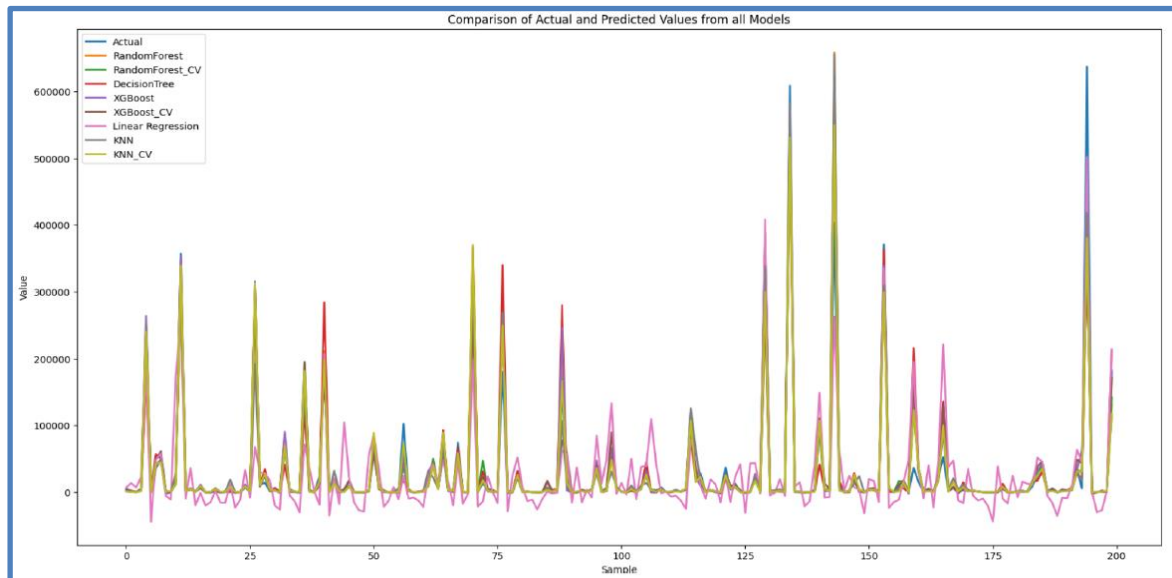
**Table III: Comparative Analysis of Different Models**

Model	R2_Score	Accuracy
RF Regressor	0.932	0.932
RF	0.930	0.891
Decision Tree	0.896	0.896
XG Boost Regressor	0.931	0.931
<b>XG Boost with Randomized Search</b>	<b>0.937</b>	<b>0.937</b>
Linear Regression	0.648	0.648
KNN	0.913	0.913
<b>KNN_CV</b>	<b>0.934</b>	<b>0.934</b>



**Figure. 8: Results of R2\_Score on Different Models**





**Figure 9: Comparison of Actual and Predicted values of all Models**

## V. CONCLUSION & FUTURE SCOPE

Intelligent technologies for agricultural yield prediction show great potential in helping farmers overcome the various issues in today's agricultural industry. These systems utilize sophisticated machine learning algorithms to analyze elements, including weather conditions, soil characteristics, and crop type, to provide essential insights about crop output potential. Assessing different models such as Random Forest, Decision Tree, XGBoost, and others shows their effectiveness in correctly forecasting agricultural yields for various crops. Furthermore, these models' excellent R2 scores and accuracies highlight their reliability and robustness in predicting yield. Intelligent systems empower farmers by providing actionable insights from historical data and patterns to optimize resource usage, enhance agricultural efficiency, and increase the resilience of farming systems. This helps to reduce the impact of climate change, market fluctuations, and other challenges.

It will be necessary to improve machine learning algorithms and incorporate a variety of data sources, such as satellite imaging and Internet of Things sensors, in order to advance intelligent systems for agricultural production estimates in the future. The use of ensemble methods and deep learning presents a chance to improve forecast accuracy. For deployment to be successful and potentially alter agriculture to improve food security and sustainability, stakeholders must work together.

## REFERENCES

- [1] Y. Alebele et al., "Estimation of Crop Yield From Combined Optical and SAR Imagery Using Gaussian Kernel Regression," in *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 14, pp. 10520-10534, 2021.
- [2] J. Liu et al., "Crop Yield Estimation in the Canadian Prairies Using Terra/MODIS-Derived Crop Metrics," in *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, vol. 13, pp. 2685-2697, 2020.
- [3] Jhajharia, Kavita, Pratistha Mathur, Sanchit Jain, and Sukriti Nijhawan. "Crop yield prediction using machine learning and deep learning techniques." *Procedia Computer Science* 218 (2023): 406-417.
- [4] P. Sharma, P. Dadheech, N. Aneja and S. Aneja, "Predicting Agriculture Yields Based on Machine Learning Using Regression and Deep Learning," in *IEEE Access*, vol. 11, pp. 111255-111264, 2023.
- [5] S. Thirumal and R. Latha, "Automated Rice Crop Yield Prediction using Sine Cosine Algorithm with Weighted Regularized Extreme Learning Machine," 2023 7th International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2023, pp. 35-40.
- [6] R. J, V. K. G. Kalaiselvi, A. Sheela, D. S. D and J. G, "Crop Yield Prediction Using Machine Learning Algorithm," 2021 4th International Conference on Computing and Communications Technologies (ICCCT), Chennai, India, 2021, pp. 611-616.
- [7] A. k. Gajula, J. Singamsetty, V. C. Dodda and L. Kuruguntla, "Prediction of crop and yield in agriculture using machine learning technique," 2021 12th International Conference on Computing Communication and Networking Technologies (ICCCNT), Kharagpur, India, 2021, pp. 1-5

- [8] E. -S. Ibrahim et al., "In vitro imaging of kidney stones using ultra-short echo-time magnetic resonance," 2015 IEEE 12th International Symposium on Biomedical Imaging (ISBI), Brooklyn, NY, USA, 2015, pp. 1466-1469.
- [9] Shilpa Choudhary, Kamlesh Lakhwani and Sandeep Kumar, "Three Dimensional Objects Recognition & Pattern Recognition Technique; Related Challenges: A Review," Multimedia Tool and Application, vol. 23, no. 1, pp. 1-44, 2022.
- [10] Rani, Shilpa, Deepika Ghai, and Sandeep Kumar, "Reconstruction of Simple and Complex Three Dimensional Images Using Pattern Recognition Algorithm," Journal of Information Technology Management, pp.235-247, 2022.
- [11] Shilpa Rani, Deepika Ghai and Sandeep Kumar, "Object Detection and Recognition using Contour based Edge Detection and Fast R-CNN" in Multimedia Tools and Application, vol. 22, no. 2, pp. 1-25, 2022