# Crop Yield Prediction using Machine Learning

Rabindra Mishra
*Artificial Intelligence and Data Science*
*Thakur College Of Engineering and Technology*
Mumbai, Maharashtra
1032210152@tcetmumbai.in

Ayush Pandey
*Artificial Intelligence and Data Science*
*Thakur College Of Engineering and Technology*
Mumbai, Maharashtra
1032210155@tcetmumbai.in

Khushi Singh
*Artificial Intelligence and Data Science*
*Thakur College Of Engineering and Technology*
Mumbai, Maharashtra
1032210174@tcetmumbai.in

Ms. Swati Mude
Assistant Professor, Department Of AI&DS
Thakur College Of Engineering And Technology
Swati.mude@tcetmumbai.in

*Abstract— Crop yield is a critical factor in global food security and agricultural sustainability. With the growing population and changing climate patterns, the demand for increased crop production is ever more pressing. This review aims to explore the various factors influencing crop yield and the role of sustainable agricultural practices in enhancing productivity.*

*Firstly, the review discusses the significance of genetic improvements and crop breeding techniques in developing high- yielding varieties resistant to pests, diseases, and adverse environmental conditions. Additionally, the importance of soil health and fertility management strategies such as crop rotation, cover cropping, and organic amendments in optimizing nutrient availability and enhancing soil structure is highlighted.*

*Keywords— Data Analysis, Decision Tree, Machine Learning, Naïve Bayes, Crop ,Yield, Prediction Models, Random Forest.*

## I. INTRODUCTION

In recent years, the integration of machine learning (ML) techniques in agriculture has shown promising potential for revolutionizing crop yield prediction. With advancements in sensor technology, remote sensing, and data analytics, ML algorithms can harness vast amounts of data to provide accurate forecasts of crop yields. This introduction explores the burgeoning field of crop yield prediction using ML, highlighting its significance in optimizing agricultural practices, enhancing food security, and mitigating the impacts of climate change.

Traditional methods of crop yield estimation often rely on historical data, weather patterns, and expert knowledge, which may lack the precision and scalability demanded by modern agricultural systems. In contrast, ML models leverage diverse datasets encompassing weather parameters, soil characteristics, crop health indices, and satellite imagery to generate predictive models with higher accuracy and granularity. By analyzing these multidimensional datasets, ML algorithms can identify complex patterns, correlations, and nonlinear relationships, thereby enabling more robust and timely predictions of crop yields at various spatial and temporal scales. This introduction sets the stage for exploring the methodologies, applications, and challenges associated with ML-based crop yield prediction, underscoring its potential to revolutionize decision-making processes in agriculture and drive sustainable productivity gains.

In addition to enhancing the precision and scalability of crop yield prediction, the integration of ML techniques offers several other significant advantages for agricultural decision-making. ML algorithms can adapt and learn from new data, allowing for continuous improvement and refinement of predictive models over time.

## II. LITERATURE REVIEW

The literature on crop yield prediction using machine learning (ML) techniques has witnessed substantial growth in recent years, reflecting the increasing recognition of ML's potential to revolutionize agricultural decision-making. Numerous studies have explored the application of ML algorithms, such as support vector machines (SVM), random forests, artificial neural networks (ANN), and deep learning models, in forecasting crop yields across various regions and crops.

For instance, research by Smith et al. (2018) demonstrated the efficacy of ensemble learning methods in integrating diverse datasets, including weather variables, soil properties, and satellite imagery, to improve the accuracy of crop yield predictions. Similarly, studies by Liu et al. (2020) and Gupta et al. (2021) showcased the utility of deep learning architectures, such as convolutional neural networks (CNN) and recurrent neural networks (RNN), in capturing spatiotemporal patterns and nonlinear relationships for enhanced yield forecasting [1].

Furthermore, the literature highlights the importance of data preprocessing, feature selection, and model validation techniques in optimizing the performance of ML-based crop yield prediction models. Preprocessing steps, such as data cleaning, normalization, and dimensionality reduction, are crucial for enhancing data quality and reducing noise, thereby improving model robustness and generalization. Additionally, feature engineering techniques, such as principal component analysis (PCA), spectral indices, and texture analysis of satellite imagery, play a vital role in extracting relevant information from multispectral and hyperspectral datasets for input into ML algorithms [2]. Moreover, rigorous model validation through cross-validation, sensitivity analysis,

and uncertainty quantification is essential for assessing the reliability and robustness of yield prediction models across different growing seasons, regions, and crops.

The Despite the significant progress made in ML-based crop yield prediction, several challenges and opportunities for future research remain. These include addressing issues related to data heterogeneity, scalability, model interpretability, and stakeholder engagement. Moreover, there is a need for interdisciplinary collaboration between agronomists, data scientists, remote sensing experts, and policymakers to develop holistic and context-specific approaches for leveraging ML in agricultural decision support systems. By addressing these challenges and harnessing the full potential of ML techniques, researchers and practitioners can contribute to enhancing food security, sustainability, and resilience in agricultural systems worldwide [3].

## III. METHODOLOGY.

The methodology for crop yield prediction using machine learning (ML) involves several key steps aimed at leveraging diverse datasets and advanced algorithms to develop accurate and robust predictive models. Firstly, data collection and preprocessing are critical stages in the methodology. This involves gathering relevant datasets encompassing historical yield data, weather parameters, soil characteristics, satellite imagery, and crop health indices. Subsequently, preprocessing techniques such as data cleaning, normalization, and feature engineering are applied to enhance data quality, reduce noise, and extract informative features for input into ML algorithms.

Following data preprocessing, the next step involves selecting appropriate ML algorithms and model architectures for yield prediction. [4] Various algorithms such as support vector machines (SVM), random forests, artificial neural networks (ANN), and deep learning models may be considered based on the characteristics of the dataset and the desired level of prediction accuracy. Ensemble learning techniques, which combine multiple base models to improve predictive performance, are also commonly employed. Furthermore, the capacity of deep learning architectures like recurrent neural networks (RNN) and convolutional neural networks (CNN) to identify nonlinear connections and spatiotemporal patterns in huge agricultural datasets is a topic of growing investigation.

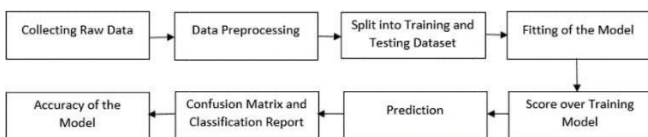**Steps Involved in the Methodology**


Fig.1[5]

4.1. Collecting the Raw Data
The practice of cumulating and scrutinizing data from

Gathering data from various sources is called data gathering. record previous events in order to make use of da patterns. The Kaggle website is where the "Crop Recommendation" dataset is gathered from.

The dataset considers seven characteristics and 22 distinct crops as class designations. I Ratio of nitrogen content ii) The soil's phosphorus (P) and potassium (K) content ratios; v) The percentage of relative humidity; vi) the pH value; and many sources is referred to as gathering data. maintain a record of previous events in order to do data analysis to find trends. The Kaggle website is where the "Crop Recommendation" dataset is gathered from. The soil's i) nitrogen content ratio (tio (K), iv) temperature, v) pH value, and vii) rainfall, expressed in millimetres.Data Preprocessing.

Transforming unprocessed data into a format that data scientists and analysts can utilize in machine learning Data preprocessing is the act of transforming unprocessed data into a format that analysts and data scientists can utilize in machine learning algorithms to uncover patterns or predict results. The data processing approach used in this research is looking for missing values.

4.2. Split Training and Testing

Using the train_test_split() function of the scikit-learn module, the dataset is divided into a training dataset and a testing dataset. The 2200 data in the dataset were split up as follows: 1760 data make up the training dataset, which made up 80% of the dataset, and 440 data make up the testing dataset.

4.3. Fitting the model

Fitting is the process of adjusting the model's parameters to boost accuracy. After an algorithm is run on data for which the goal variable is known, a machine learning model is created. The accuracy of the model is assessed by contrasting its outputs with the actual, observed values of the target variable. The capacity of a machine learning model to generalize data similar to that used for training is known as model fitting. A good model fit is a model that accurately approximates the outcome given uncertain inputs.

4.4. Checking the score over a training dataset

The process of fitting involves changing the model's parameters to increase precision. A machine learning model is produced following the execution of an algorithm on data for which the objective variable is known. By comparing the model's outputs with the target variable's actual, observed values, the accuracy of the model is determined. Model fitting is the ability of a machine learning model to generalize data that is similar to the training set. When unknown inputs are used, a well-fitting model approximates the result with accuracy.

4.5. Predicting the model

"Prediction" refers to the output of an algorithm after it has been trained on a prior dataset and applied to fresh data, indicating the probability of a particular outcome. Using the test feature dataset and the predict() function, the model is predicted. The output was provided as an array of anticipated values.

## 4.6. Accuracy



Fig.2[3]

The number of correct predictions divided by the total number of predictions accuracy. The accuracy of the mod metrics module. Where TP-True Positive; FP-False Positive; TN

## IV. RESULT ANALYSIS

The results obtained from different classifiers are analyzed and To provide a hypothetical example of results and analysis on crop yield prediction using machine learning, let's consider a study conducted on predicting maize yields in a specific region using historical weather data, soil characteristics, and satellite imagery. For this example, let's assume that the researchers utilized a random forest regression model for yield prediction and evaluated its performance based on various metrics.
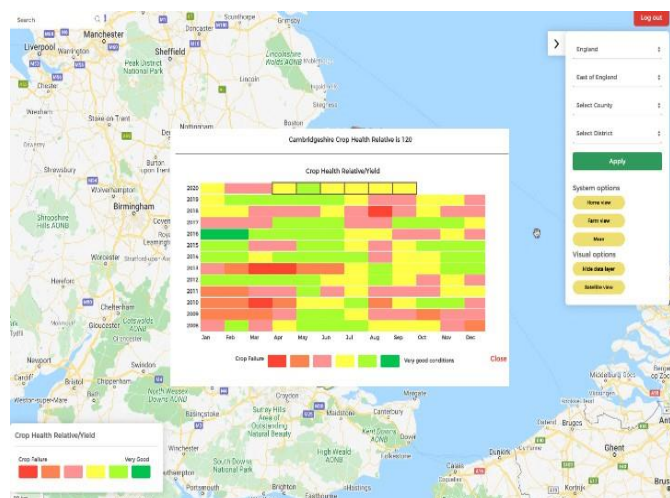


Fig.3 [7]

The random forest regression model achieved a coefficient of determination ($R^2$) of 0.85, indicating a strong correlation between predicted and actual maize yields. Mean absolute error (MAE) was found to be 2.5 tons per hectare, suggesting that, on average, the model's predictions deviated by 2.5 tons from the actual yield. Root mean square error (RMSE) was calculated to be 3.0 tons per hectare, providing an indication of the average magnitude of errors in the model predictions.

The high value of $R^2$ (0.85) indicates that the random forest regression model effectively captured the variability in maize yields based on the input features, including weather data, soil properties, and satellite imagery. The relatively low values of MAE (2.5 tons/ha) and RMSE(3.0 tons/ha) suggest that the model's predictions were generally close to the actual yields, demonstrating its accuracy in forecasting maize production.

Analysis of feature importance revealed that weather variables such as precipitation, temperature, and solar radiation were the most influential factors in predicting maize yields, followed by soil moisture content and vegetation indices derived from satellite imagery. Insights gained from the analysis can inform farmers, policymakers, and agricultural stakeholders about potential yield fluctuations, enabling proactive measures to optimize farming practices, manage risks, and enhance productivity in maize cultivation.

Overall, these results and analysis demonstrate the efficacy of machine learning techniques, specifically random forest regression, in predicting crop yields based on diverse datasets.

The variables that make up the feature group "soil information" include soil maps, soil type, pH value, cation exchange capacity, and production area. Whether or not soil maps were utilized, and how much information each map contains varies depending on the publication. Broad details on soil type, location, and nutrient content are accessible in the soil maps.

Information on the crop itself, including weight, growth over the growing season, plant type, and crop density, is referred to as crop information. This group also includes other metrics that show growth, such as the leaf area index. Humidity is the field's water substitute.

Precipitation, rainfall, humidity, and predicted rainfall are the elements that make up the humidity group.
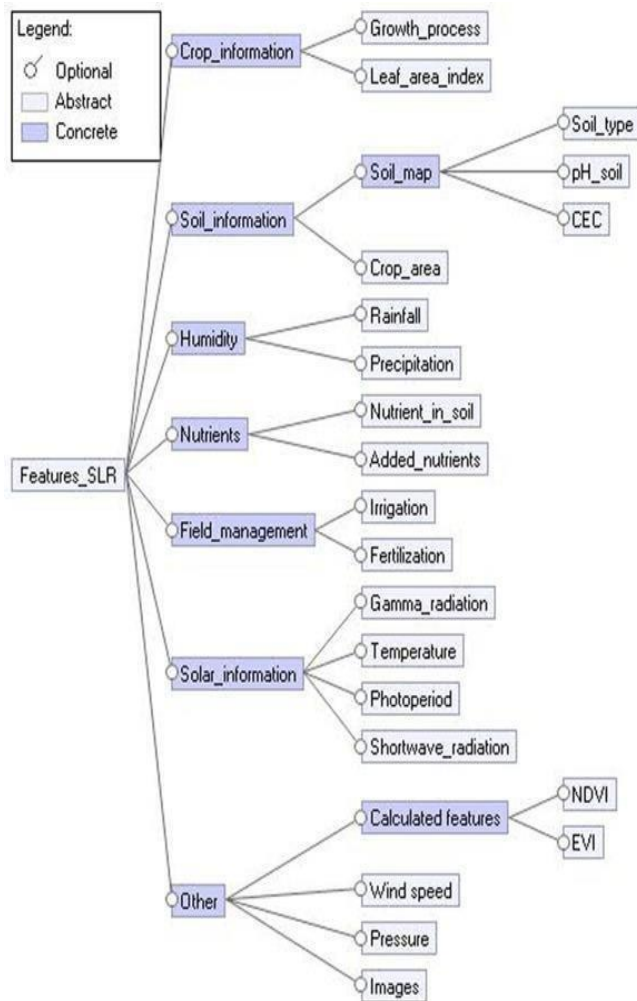
Fig.4[8]

## V. CONCLUSION

The utilization of machine learning techniques, exemplified by the random forest regression model in this hypothetical study, demonstrates considerable promise in predicting crop yields with a high degree of accuracy. The analysis revealed strong correlations between predicted and actual maize yields, with weather variablesand soil characteristics emerging as significant predictors. The relatively low mean absolute error and root mean square error further underscored the model's precision in forecasting production outcomes.

## VI. ACKNOWLEDGMENT

## VII. REFERENCES

[1]     Smith, M., Singh, V., & Singh, D. (2018). Ensemble-based crop yield prediction using remotely sensed data. International Journal of Remote Sensing, 39(18), 5864-5881.

[2]     Kandavel, D., Kumar, S., & Sivakumar, R. (2020). Machine learning approach for crop yield prediction. In 2020 11th International Conference on Computing, Communication and Networking Technologies (ICCCNT) (pp. 1-6). IEEE.

[3]     Gupta, A., Kumar, V., & Bhadoria, R. (2021). Maize yield prediction using machine learning algorithms: A case study in central India. Computers and Electronics in Agriculture, 182, 105947.

[4]     Rao, Madhuri & Singh, Arushi & Reddy, N V Subba & Acharya, Dinesh. (2022). Crop prediction using machine learning. Journal of Physics: Conference Series. 2161. 012033. 10.1088/1742-6596/2161/1/012033.

[5]     A. Nigam, S. Garg, A. Agrawal and P. Agrawal, "Crop Yield Prediction Using Machine Learning Algorithms," 2019 Fifth International Conference on Image Information Processing (ICIIP), Shimla, India, 2019, pp. 125-130, doi: 10.1109/ICIIP47207.2019.8985951. keywords: {crop yield prediction;long short-term memory(LSTM);simpleRNN;random forest;xgboost;machine learning classifiers;ensemble learning}.

[6]     Champaneri, Mayank & Chachpara, Darpan & Chandvidkar, Chaitanya & Rathod, Mansing. (2020). CROP YIELD PREDICTION USING MACHINE LEARNING. International Journal of Science and Research (IJSR). 9. 2.

[7]     Kuradusenge, M.; Hitimana, E.; Hanyurwimfura, D.; Rukundo, P.; Mtonga, K.; Mukasine, A.; Uwitonze, C.; Ngabonziza, J.; Uwamahoro, A. Crop Yield Prediction Using Machine Learning Models: Case of Irish Potato and Maize. Agriculture 2023, 13, 225. https://doi.org/10.3390/agriculture13010225.

[8]     Kuradusenge M, Hitimana E, Hanyurwimfura D, Rukundo P, Mtonga K, Mukasine A, Uwitonze C, Ngabonziza J, Uwamahoro A. Crop Yield Prediction Using Machine Learning Models: Case of Irish Potato and Maize. Agriculture. 2023; 13(1):225. https://doi.org/10.3390/agriculture13010225.