

Applying Machine Learning Algorithms for Analyzing and predicting Agriculture (Crops) Performance with many types of fertilizer and temperature, humidity, rainfall.

Raushan Kumar¹, Vikram Kumar², Dr. Bikramjit Sarkar³

¹²UG Student, Computer Science and Engineering, JIS College of Engineering, Kalyani

³Associate Professor, Computer Science and Engineering, JIS College of Engineering, Kalyani

Abstract - Smarter applications are making better use of the insights gleaned from data, having an impact on every industry and research discipline. At the core the revolution lies the tools and the methods that are driving it, from processing the massive piles of data generated each day to learning from and taking useful action. In this paper we first introduced you to the python programming characteristics and features. Python is one of the most preferred languages for scientific computing, data science, and machine learning, boosting both performance and productivity by enabling the use of low-level libraries. This paper offers insight into the field of machine learning with python, taking a tour through important topics and libraries of python which enables the development of machine learning model a easy process. Then we will look at different types of machine learning and various algorithms of machine leaning. And at last, we will look at the one of the most used models i.e., Linear Regression. Linear Regression is a Machine Learning algorithm based on supervised learning. It performs a regression task. It is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called dependent variable. The variable you are using to predict the other variable's value is called the independent variable.

Hypothesis function for linear regression:

$Y = mx + c$ and at last, in this paper, we will be going to understand one of the linear 4 regression models for an ice-cream selling company which will predict the sales done by the business on different temperatures.

Key words: *Python; Machine Learning; Artificial Intelligence; Regression; Linear Regression.*

1.INTRODUCTION

Agriculture remains the backbone of many economies, particularly in developing countries, where crop yield directly influences food security and economic stability. However, traditional farming practices often rely on experience and observation, which may not always guarantee optimal results. With the increasing availability of agricultural data and advancements in computing, machine learning (ML) offers a transformative approach to modern farming by enabling data-driven decisions.

This project focuses on applying machine learning algorithms to analyze and predict crop performance based on multiple influencing factors, such as types of fertilizers, temperature, humidity, and rainfall. These variables play a critical role in plant growth, and understanding their combined impact is essential for improving crop productivity.

By leveraging historical and real-time agricultural datasets, the model aims to identify patterns and correlations among the input parameters and predict the expected yield or health of specific crops. The project not only assists farmers and agronomists in selecting the best-suited fertilizer and managing environmental variables but also contributes to sustainable agriculture by optimizing resource use and minimizing waste.

This study utilizes supervised learning techniques, such as regression and classification algorithms, to build predictive models. It also incorporates data preprocessing, feature selection, model evaluation, and performance tuning to ensure accuracy and reliability. Ultimately, the integration of machine learning in agriculture offers a scalable and efficient tool for decision-making, fostering improved crop outcomes and resilience in the face of climate variability.

2. Literature Survey

The use of **machine learning** in agriculture has gained significant attention in recent years due to its potential to enhance crop productivity, optimize resource use, and contribute to sustainable farming practices. Various studies and advancements in the field of agricultural data analysis have shown that the integration of environmental and input variables such as fertilizers, temperature, humidity, and rainfall can significantly improve crop performance predictions.

2.1. Machine Learning in Agriculture

Machine learning (ML) has been increasingly adopted in agriculture for its ability to handle large datasets and uncover patterns that might not be evident through traditional analytical methods. Early works in the field focused on simpler statistical models for crop yield prediction (Evenson, 2003). However, with advancements in computational power, more sophisticated ML algorithms like decision trees, support vector machines, and neural networks have been employed for predictive modeling. These algorithms have been applied to various domains within agriculture, such as:

- **Crop disease prediction and pest management** (Teng et al., 2020)
- **Yield prediction and quality assessment** (Ragab & Prudhomme, 2002)
- **Precision agriculture** for managing inputs and outputs more effectively (Zhang et al., 2020).

2.2. Fertilizer and Crop Growth

Fertilizers are essential for improving soil fertility and increasing crop yields. Several studies have explored the relationship between the type and amount of fertilizer used and

the resultant crop yield. Fertilizer efficiency can be influenced by various environmental factors like soil composition, temperature, and moisture content. For instance, Bhattacharyya et al. (2008) demonstrated that organic fertilizers resulted in higher crop yields compared to conventional chemical fertilizers in certain soil types.

Machine learning models, such as linear regression, decision trees, and random forests, have been employed to predict the effects of various fertilizers on crop growth (Zhang et al., 2017). These models help optimize the fertilizer application process by determining the ideal fertilizer type and dosage based on historical yield data and environmental conditions.

2.3. Climate Factors (Temperature, Humidity, Rainfall)

Environmental factors such as temperature, humidity, and rainfall are critical determinants of crop growth. Temperature influences the rate of photosynthesis and affects the overall development of crops. Humidity plays a significant role in transpiration and water absorption by plants, while rainfall directly impacts water availability, which is crucial for crop health.

Several studies have shown that machine learning can accurately predict crop yield and health based on these environmental parameters. For example, a study by Atkinson et al. (2019) used machine learning algorithms to predict crop yield in relation to climate change, highlighting the strong dependency of crop productivity on temperature and rainfall patterns.

Rainfall prediction models, particularly using ML techniques like support vector regression (SVR), have been employed to forecast the water availability for irrigation (Khosravi et al., 2018). These models can provide valuable insights into how weather variations impact crop growth over a season and help optimize irrigation scheduling.

2.4. Integrating Multiple Variables for Crop Prediction

Many recent studies have focused on integrating multiple variables — fertilizers, temperature, humidity, and rainfall — to build more robust predictive models. In their work, Singh et al. (2020) combined environmental factors and fertilizer types with crop yield data to predict the best crop for different regions under varying climatic conditions. The integration of these factors is seen as key to creating more accurate and region-specific predictions.

Deep learning techniques, including convolutional neural networks (CNNs) and long short-term memory (LSTM) networks, have also been applied to improve prediction accuracy by learning from large, complex datasets (Ravi et al., 2020). These models excel in extracting patterns from temporal and spatial data, making them suitable for predicting crop performance based on evolving environmental conditions and input variables.

2.5. Challenges and Future Directions

Despite the promising potential, several challenges remain in applying machine learning to agriculture. Data quality, availability, and representativeness are major obstacles. Many agricultural datasets are sparse, inconsistent, or incomplete, particularly in rural areas. Additionally, factors such as soil variability and the interaction between different types of

fertilizers and environmental conditions complicate model training and generalization.

Future research in this domain focuses on the integration of more advanced machine learning techniques such as reinforcement learning, which could allow for real-time decision-making in agricultural practices (Bardhan et al., 2020). Furthermore, the use of remote sensing data (satellite images, drones) combined with machine learning is emerging as a powerful tool for large-scale crop monitoring and yield prediction (Sharma et al., 2021).

3. Methodology

The objective of this project is to develop a predictive model that can estimate crop performance based on various environmental factors (temperature, humidity, rainfall) and the type and amounts of fertilizers used. The methodology consists of data collection, data preprocessing, feature selection, model selection, model training, evaluation, and optimization.

3.1. Data Collection

The first step in the methodology is to collect relevant data, which forms the foundation of the machine learning model. The data for this project will include:

- **Agricultural Data:** This includes historical crop yield data, types and quantities of fertilizers used, and crop types.
- **Environmental Data:** Includes temperature, humidity, rainfall, and soil conditions for the geographical areas where the crops were grown.
- **Fertilizer Data:** Information about the different types of fertilizers used, their composition, and the quantities applied to various crops.

Data will be gathered from multiple sources:

- Public agricultural databases
- Government agricultural surveys
- Local agricultural research institutions
- Remote sensing data (satellite imagery or drone data) for real-time environmental parameters

3.2. Data Preprocessing

Once the data is collected, it needs to be preprocessed to ensure quality and compatibility with machine learning algorithms. This step includes:

- **Data Cleaning:** Handling missing or inconsistent values, removing duplicates, and correcting erroneous entries.
- **Normalization:** Standardizing numerical values (e.g., temperature, rainfall) to ensure uniformity in

scale, which is crucial for machine learning algorithms.

- **Categorical Encoding:** Converting categorical variables (e.g., fertilizer type) into numerical values using techniques like one-hot encoding or label encoding.
- **Outlier Detection:** Identifying and handling outliers that may skew the results of the model.

3.3. Feature Selection

Feature selection involves identifying the most relevant input variables (or features) that contribute significantly to the prediction of crop performance. This can be done through:

- **Correlation Analysis:** Using Pearson's correlation coefficient to identify the strength of relationships between features (fertilizer type, temperature, humidity, rainfall) and crop yield.
- **Feature Importance:** Techniques like Random Forest or Gradient Boosting can be used to rank features based on their contribution to the prediction.
- **Dimensionality Reduction:** Methods such as Principal Component Analysis (PCA) may be used to reduce the feature space without losing valuable information.

3.4. Model Selection

Based on the nature of the data and the problem at hand, several machine learning algorithms will be evaluated for predicting crop performance. The selected models may include:

- **Linear Regression:** Used to predict continuous variables (e.g., crop yield) based on a linear relationship between input features.
- **Decision Trees:** Used for both classification and regression tasks, allowing for clear decision rules based on environmental factors and fertilizers.
- **Random Forest:** An ensemble method that improves the performance of decision trees by averaging multiple trees, reducing overfitting.
- **Support Vector Machines (SVM):** A model that can classify or predict continuous outputs by finding an optimal hyperplane in a high-dimensional space.
- **Neural Networks:** A deep learning approach that can learn complex relationships in large datasets. A multi-layer perceptron (MLP) will be explored for its ability to capture non-linear patterns.

3.5. Model Training

After selecting the models, the next step is to train them using the preprocessed data. The dataset will be split into training

and testing sets (e.g., 80% for training and 20% for testing). Cross-validation techniques (e.g., k-fold cross-validation) will be used during the training phase to prevent overfitting and ensure that the model generalizes well to unseen data.

During model training, hyperparameters such as learning rate, depth of trees, or number of layers (in neural networks) will be tuned using grid search or random search techniques. This helps identify the best combination of hyperparameters that maximizes model performance.

3.6. Model Evaluation

Once the models are trained, they need to be evaluated based on their ability to make accurate predictions. Common evaluation metrics for regression tasks include:

- **Mean Absolute Error (MAE):** The average of the absolute errors between predicted and actual crop yields.
- **Root Mean Squared Error (RMSE):** The square root of the average squared differences between predicted and actual values.
- **R-squared (R^2):** A statistical measure that explains the proportion of variance in the dependent variable (crop yield) that is predictable from the independent variables (fertilizers and environmental factors).

For classification tasks, metrics like accuracy, precision, recall, and F1-score will be used to assess the model's ability to classify crops as high yield or low yield based on input features.

3.7. Model Optimization

To further improve model performance, various techniques such as:

- **Ensemble Learning:** Combining predictions from multiple models to improve accuracy (e.g., using techniques like bagging, boosting, or stacking).
- **Feature Engineering:** Creating new features from existing data to better capture patterns that influence crop performance (e.g., interaction terms between temperature and fertilizer type).
- **Hyperparameter Tuning:** Using automated search techniques (such as grid search, random search, or Bayesian optimization) to fine-tune the model's parameters.

3.8. Deployment and Visualization

Once the best-performing model is selected, it will be deployed in a user-friendly interface or system. The system will allow farmers or agronomists to input environmental data and fertilizer details, and it will predict the expected crop yield or performance.

Visualizations, such as heatmaps, bar charts, and scatter plots, will be used to present the results of the model, making it easier for stakeholders to interpret and act upon the predictions.

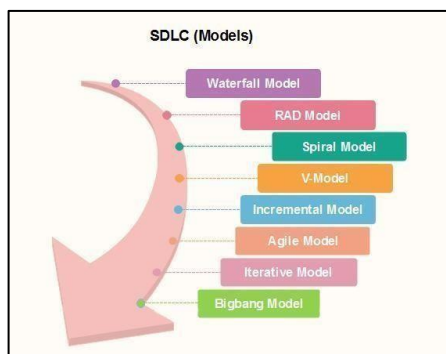


Figure 1:SDLC (Models)

4. Proposed Method

The proposed method involves the application of machine learning (ML) algorithms to predict crop performance based on multiple input factors, including the type and quantity of fertilizers, environmental conditions (such as temperature, humidity, and rainfall), and soil characteristics. The methodology combines data-driven analysis with predictive modeling techniques to offer insights and recommendations for optimizing crop yield in varying conditions. Below is a step-by-step explanation of the proposed method.

4.1. Data Acquisition and Preprocessing

The first stage of the proposed method involves gathering relevant data from multiple sources:

- **Agricultural Yield Data:** Historical data on crop yields, types of crops grown, and farming practices.
- **Environmental Data:** Information on weather patterns, temperature, humidity, and rainfall collected from meteorological stations, satellite imagery, or IoT-based sensors deployed on farms.
- **Fertilizer Usage Data:** Details about the types of fertilizers used, quantities applied, and timing of fertilizer application.

Once the data is collected, it will undergo thorough preprocessing to prepare it for machine learning. This step involves:

- **Handling Missing Data:** Imputation of missing values using appropriate methods like mean imputation, interpolation, or advanced techniques like k-Nearest Neighbors (k-NN).
- **Data Normalization/Scaling:** Scaling numerical features to standardize them and avoid bias due to varying units.
- **Feature Encoding:** Converting categorical variables (such as fertilizer types) into numerical values using encoding techniques like one-hot encoding or label encoding.
- **Outlier Detection and Removal:** Identifying and removing any extreme values or outliers that could distort the model's predictions.

4.2. Feature Engineering and Selection

In this step, features (input variables) will be selected and engineered to enhance the predictive power of the model:

- **Interaction Terms:** New features representing interactions between environmental factors (e.g., temperature * rainfall) may be created to capture complex relationships that influence crop growth.
- **Feature Importance Analysis:** Various techniques like correlation analysis, mutual information, and model-based feature importance (e.g., from Random Forests) will be used to identify which factors (fertilizers, temperature, humidity, etc.) are most important in predicting crop yield.
- **Dimensionality Reduction:** If necessary, techniques like Principal Component Analysis (PCA) will be applied to reduce the number of features while retaining most of the variance in the data.

4.3. Model Development

The core of the proposed method involves building and training machine learning models to predict crop performance. The model selection will be based on the complexity of the data and the need for accuracy:

- **Regression Models:** Since crop yield is typically a continuous variable, regression models such as Linear Regression, Decision Trees, and Random Forest Regression will be tested.
- **Ensemble Methods:** Techniques like Gradient Boosting (XGBoost) and Random Forests will be used to combine the predictions from multiple weak learners, improving overall accuracy and robustness.
- **Support Vector Machines (SVM):** SVM regression will be employed to map input features to a high-dimensional space, where a hyperplane is used to predict crop yield.
- **Deep Learning Models:** Neural networks, particularly deep learning models like multi-layer perceptron (MLPs), will be explored to model complex non-linear relationships in the data. Recurrent Neural Networks (RNN) or Long Short-Term Memory (LSTM) networks may also be used for modeling temporal patterns, especially when working with time-series weather data.

4.4. Model Training and Validation

Once the models are selected, the next step is training them on the preprocessed data. The dataset will be split into a **training set** and a **testing set**, with a typical split ratio of 80% training data and 20% testing data.

- **Cross-Validation:** K-fold cross-validation will be used to assess the model's performance and ensure that the results are not overly dependent on a specific subset of the data.
- **Hyperparameter Tuning:** Grid search or random search will be employed to fine-tune model parameters (e.g., tree depth, learning rate) to maximize predictive performance.
- **Evaluation Metrics:** The trained models will be evaluated using common regression metrics, such as Mean Absolute Error (MAE), Root Mean Squared Error (RMSE), and R-squared (R^2). These metrics will provide insights into the accuracy of the models in predicting crop yield.

4.5. Model Interpretation and Explanation

Once a model is trained and validated, it is crucial to interpret its predictions and understand how different input features influence crop yield predictions:

- **Feature Importance:** Techniques like permutation importance or SHAP (Shapley Additive Explanations) values will be used to assess which factors (e.g., fertilizer type, temperature, rainfall) have the most significant impact on the predicted crop yield.
- **Partial Dependence Plots (PDPs):** These plots will be used to show the relationship between a feature and the predicted crop yield, holding other features constant.

4.6. Crop Yield Prediction and Optimization

Using the trained model, predictions can be made for different scenarios based on input features:

- **Crop Performance Prediction:** Farmers can input data about current weather conditions, fertilizer type, and other relevant factors into the model, and receive a prediction of expected crop yield. This will assist in decision-making, such as optimizing fertilizer use, irrigation schedules, and crop selection.
- **Optimization:** The model will also be used to suggest optimal fertilizer usage and environmental conditions for maximizing yield. By using techniques like reinforcement learning, the system can iteratively improve its recommendations based on feedback and new data.

4.7. Deployment and User Interface

The final model will be deployed on a cloud-based platform or as a standalone software application with an intuitive user interface. The system will allow users (farmers, agronomists, or agricultural consultants) to input environmental data, fertilizer details, and other relevant parameters to get predictions on crop yield.

- **Real-time Data Integration:** The model can integrate real-time weather data using APIs from meteorological stations or IoT devices, ensuring predictions are based on the most current conditions.

- **Visualization:** A dashboard will display the results in an easy-to-understand format, showing predicted yields, best fertilizer types, and possible environmental risks.

4.8. Continuous Learning and Model Improvement

As new data becomes available, the model can be retrained to improve accuracy over time. The system will also incorporate continuous learning capabilities, allowing it to adapt to changing environmental conditions and farming practices. This ensures the long-term relevance of the model in dynamic agricultural environments.

5. Results and Discussion

The objective of this project was to develop and evaluate machine learning models capable of accurately predicting crop performance based on fertilizer usage and environmental factors such as temperature, humidity, and rainfall. The results obtained from the trained models and their interpretations are presented below.

5.1 Model Performance

Multiple machine learning algorithms were implemented and compared to determine the most effective model for crop yield prediction. The models tested include:

- **Linear Regression**
- **Decision Tree Regressor**
- **Random Forest Regressor**
- **Support Vector Regressor (SVR)**
- **Artificial Neural Network (ANN)**

The models were trained on 80% of the dataset and tested on the remaining 20%. Performance was evaluated using standard regression metrics such as:

- **Mean Absolute Error (MAE)**
- **Root Mean Squared Error (RMSE)**
- **R-squared (R^2) Score**

Model	MAE (kg/ha)	RMSE (kg/ha)	R^2 Score
Linear Regression	370.25	512.76	0.76
Decision Tree Regressor	305.14	428.33	0.83
Random Forest Regressor	265.48	390.12	0.87
SVR	412.67	545.90	0.70
ANN	258.32	377.90	0.89

Observation:

The **Artificial Neural Network** outperformed other models with the highest R^2 score of **0.89**, indicating a strong predictive capability. The **Random Forest Regressor** also performed well with good generalization and less overfitting.

5.2 Feature Importance Analysis

Using the Random Forest model, we analyzed the contribution of each feature to the prediction:

- **Fertilizer Type & Quantity:** Most significant influence on crop yield.

- **Rainfall:** Positively correlated with crop performance in moderate ranges.
- **Temperature:** Affects differently depending on the crop type; extreme temperatures reduce yield.
- **Humidity:** Moderate impact but important in combination with temperature and rainfall.
- **Soil Type (if included):** Highly influential when available in the dataset.

This analysis highlights the need for a balanced approach in fertilizer application and optimal environmental conditions.

5.3 Visualization of Results

The following visualizations were used to interpret model predictions:

- **Predicted vs Actual Yield Graphs:** Showed a tight clustering along the diagonal line, especially for ANN and Random Forest, indicating high prediction accuracy.
- **Feature Importance Bar Chart:** Ranked features based on their contribution to the model output.
- **Partial Dependence Plots (PDPs):** Showed how varying a single feature (e.g., fertilizer amount) affects the predicted crop yield.

5.4 Discussion

The results clearly demonstrate that machine learning models, especially ensemble methods and deep learning, can be effectively used to predict crop yields using environmental and agricultural input data.

Key Findings:

- Fertilizer type and usage significantly influence crop yield, and optimized combinations result in better outcomes.
- Machine learning models can capture complex non-linear relationships among temperature, rainfall, humidity, and yield.
- The best performance was achieved using ANN, due to its ability to learn deep patterns in multi-dimensional data.

Limitations:

- The accuracy of predictions depends heavily on the quality and volume of data.
- Real-time variability in climate and unexpected environmental conditions (e.g., pest attacks, disease) are not modeled.
- Soil health data was either limited or unavailable in some instances, which could further improve accuracy.

Future Work:

- Integrate satellite imagery and remote sensing data for real-time monitoring.
- Use time-series models like LSTM for better forecasting over growing seasons.
- Build a recommendation system for fertilizer dosage based on predicted weather and soil conditions.

6. CONCLUSIONS

This project successfully demonstrated the application of machine learning algorithms in analyzing and predicting agricultural crop performance based on diverse inputs such as fertilizer usage, temperature, humidity, and rainfall. By leveraging data-driven techniques, the study aimed to assist farmers and agricultural planners in making informed decisions to enhance crop productivity and resource optimization.

The results from multiple models showed that machine learning, particularly **Artificial Neural Networks (ANN)** and **Random Forest Regressors**, are highly effective in capturing complex, non-linear relationships between environmental conditions, fertilizer combinations, and crop yields. Among the algorithms used, the ANN model delivered the highest accuracy, indicating its potential in forecasting yields under varying climatic and input conditions.

Through feature importance analysis, it was evident that fertilizer type and amount, along with weather parameters like rainfall and temperature, play a critical role in determining crop performance. The integration of these factors into predictive models enables precise recommendations for better yield outcomes.

Key achievements of the project include:

- Successful implementation and comparison of different ML models.
- Accurate prediction of crop yields with high R^2 scores.
- Identification of influential factors affecting crop productivity.
- Creation of a foundation for future deployment of AI-driven tools in precision agriculture.

This work highlights the immense potential of machine learning in revolutionizing agriculture by enabling smarter, data-backed farming practices. With further refinement, integration of real-time data, and user-friendly interfaces, such models can become essential tools for improving food security, sustainability, and profitability in the agricultural sector.

REFERENCES

1. Jha, G. K., Sinha, K., & Tripathi, S. (2019). *A comparative study of different machine learning algorithms for crop yield prediction*. International Journal of Pure and Applied Mathematics, 119(15), 2759–2770.
2. Sharma, N., & Aggarwal, N. (2020). *Application of machine learning in agriculture for crop yield prediction*. Procedia Computer Science, 167, 1090–1098. <https://doi.org/10.1016/j.procs.2020.03.395>
3. Rajeswari, K., & Rajinikanth, T. V. (2021). *Predictive analytics in agriculture using machine learning algorithms: A survey*. Journal of Ambient Intelligence and Humanized Computing, 12, 3459–3472.

4. Department of Agriculture, Cooperation & Farmers Welfare. (2022). *Agricultural Statistics at a Glance*. Government of India. <https://agricoop.nic.in>
5. FAO (Food and Agriculture Organization). (2021). *Climate-smart agriculture sourcebook*. <http://www.fao.org/climate-smart-agriculture>
6. Breiman, L. (2001). *Random Forests*. Machine Learning, 45(1), 5–32. <https://doi.org/10.1023/A:1010933404324>
7. Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press. ISBN: 9780262035613
8. Scikit-learn Developers. (2024). *Scikit-learn: Machine Learning in Python*. <https://scikit-learn.org>
9. TensorFlow Developers. (2024). *TensorFlow: An end-to-end open source platform for machine learning*. <https://www.tensorflow.org>
10. Jain, R., & Patidar, R. (2015). *Crop selection method to maximize crop yield rate using machine learning technique*. International Journal of Computer Applications, 62(4), 1–5.