

E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Forecasting State – Wise Crop Yield in India

Swaroop Kumar Aduru, Prof. Aniruddha Dasgupta

Abstract

Agriculture is the backbone of India's economy, with millions of farmers dependent on timely insights into crop productivity for sustainable livelihoods and national food security. However, crop yields are highly sensitive to multiple interacting factors, including climatic variability, soil health, and input management. This project focuses on developing an intelligent, data-driven machine learning system for **predicting crop yield across Indian states**, integrating historical agricultural, environmental, and meteorological data into a unified predictive framework.

The primary objective of this project is to build an accurate, interpretable, and scalable model that can forecast the yield of diverse crops based on factors such as area of cultivation, fertilizer and pesticide use, temperature, rainfall, humidity, soil nutrient composition (Nitrogen, Phosphorus, Potassium), soil pH, and seasonal patterns. The system leverages **supervised machine learning algorithms** - including **XGBoost**, **Random Forest**, and **Histogram Gradient Boosting** - to identify the complex, nonlinear dependencies between these input features and observed crop yields. Through model comparison and hyperparameter tuning, the XGBoost model was found to deliver superior predictive accuracy and generalization, making it the foundation of the final deployed pipeline.

The dataset used for training and evaluation was compiled from authentic government and open data sources, encompassing multiple years and covering a wide variety of crops, states, and seasons. Preprocessing steps included data standardization, feature scaling, label encoding for categorical attributes (such as crop, season, and state), and normalization of temporal data using a **YearTransformer** to capture long-term yield trends. The final pipeline integrates all preprocessing and transformation steps with the trained model to ensure seamless, real-time prediction on new data.

To enable accessibility and practical usability, the model has been deployed as an **interactive Streamlit web application**. This app allows users - including farmers, agronomists, and policymakers—to input real-world parameters and instantly obtain yield predictions (expressed in quintals per hectare). The interface automatically handles encoding, scaling, and transformation, ensuring consistent predictions aligned with the trained pipeline. The application's intuitive design, coupled with visual data previews, provides transparency and ease of interpretation for end-users.

This project's broader impact lies in its potential to empower stakeholders with **data-driven decision support**. For farmers, it offers predictive insights to optimize resource allocation - such as fertilizer usage or crop selection - under varying climatic conditions. For policymakers, it provides a scientific foundation for yield forecasting, subsidy distribution, and early warning systems for agricultural risks. Moreover, the system contributes toward achieving **precision agriculture** goals by integrating environmental, soil, and meteorological parameters into a unified predictive framework.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

In conclusion, the proposed **Crop Yield Prediction System** demonstrates how advanced machine learning techniques can transform raw agricultural data into actionable intelligence. By combining domain knowledge with computational analytics, this project enhances the reliability of yield estimation, supports sustainable farming practices, and strengthens resilience against climate variability. Future enhancements may include incorporating satellite imagery, real-time IoT sensor data, and deep learning architectures to further improve prediction granularity and adaptiveness at the regional level.

Keywords: Crop Yield Prediction, Machine Learning, XGBoost, Indian Agriculture, Data Preprocessing, Streamlit Application, Precision Agriculture, Feature Engineering, Climate Variability, Sustainable Farming, Random Forest, Gradient Boosting, Soil Health, Weather Data, Agricultural Forecasting

Problem Statement

Agriculture, a cornerstone of India's economy and a primary livelihood for millions, faces significant volatility due to its high dependence on the unpredictable monsoon and growing climate variability, including irregular rainfall, temperature fluctuations, and extreme weather events. Traditional crop yield forecasting methods, often subjective and manual, fail to capture the complex interplay of factors that influence agricultural output, such as localized weather patterns, soil fertility, fertilizer and pesticide usage, and pest or disease outbreaks. This imprecision leads to financial losses for farmers, inefficiencies in the market, and challenges for policymakers in planning for food security, resource allocation, and disaster preparedness.

There is a pressing need for a data-driven, real-time, and localized predictive system that integrates these multiple factors to generate accurate crop yield forecasts. Such a system can help farmers make informed decisions regarding sowing, irrigation, and input management, assist policymakers in designing effective agricultural policies and subsidies, and provide the market with reliable supply-demand insights. Ultimately, this approach aims to improve agricultural productivity, economic stability, and resilience against climate uncertainty, supporting both sustainable farming practices and national food security.

Objective of the Project

The primary objective of this project is to develop a robust machine learning-based system capable of predicting future crop yield (kg/ha) for major crops such as rice, wheat, sugarcane, and maize across different Indian states. The system aims to provide accurate, data-driven insights to support farmers, policymakers, and other stakeholders in making informed decisions. The specific objectives include:

1. Data Aggregation and Preprocessing

- o Collect and integrate crop-related data from multiple sources, including government databases, agricultural surveys, and weather reports.
- o Clean, standardize, and preprocess the data to handle missing values, inconsistencies, and outliers, ensuring high-quality inputs for modeling.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

2. Feature Engineering and Exploratory Analysis

- Derive meaningful features that capture influential factors such as soil properties, climatic conditions, fertilizer and pesticide usage, and seasonal variations.
- o Conduct Exploratory Data Analysis (EDA) to identify patterns, trends, and correlations in crop yield across states, crops, and seasons.

3. Model Development and Evaluation

- Build and evaluate multiple machine learning models, including Regression, Random Forest, and XGBoost, to accurately predict crop yield.
- o Optimize model performance through feature selection, hyperparameter tuning, and cross-validation techniques.

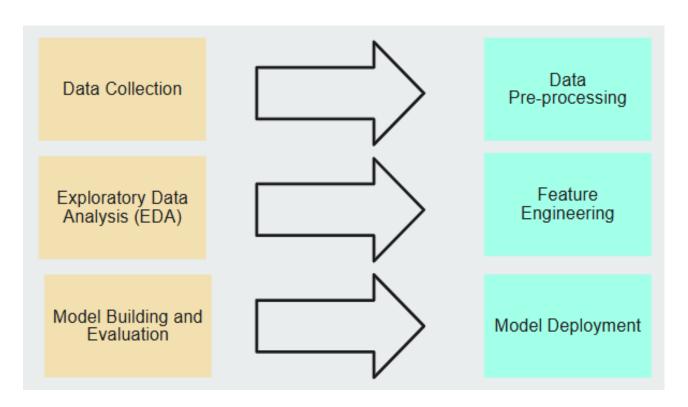
4. Prediction and Insights Generation

- o Generate reliable yield predictions to help farmers plan sowing, irrigation, and input application strategies.
- Provide actionable insights for policymakers to allocate resources efficiently and formulate effective agricultural policies.
- o Enable identification of trends and potential risks, supporting strategic decision-making across the agricultural value chain.

5. Towards Practical Deployment

- Ensure that the predictive system is interpretable and can be integrated into dashboards or decision-support tools for real-world application.
- o Facilitate short-term and long-term planning to enhance agricultural productivity, economic stability, and food security.

Machine Learning Process Flow (Architecture)





E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Resources or Prerequisite

Data Resources

The project uses a combination of crop, soil, and weather datasets to build an accurate predictive model. The datasets were collected from **India.gov.in** (https://india.gov.in/category/agriculture-rural-environment/subcategory/agricultural-produce) and include:

- 1. crop_yield.csv Contains historical yield data for major crops such as rice, wheat, sugarcane, and maize across Indian states, including attributes like year, crop, season, area, production, and yield.
- 2. state_soil_data.csv Provides soil characteristics for different states, including NPK content (Nitrogen, Phosphorus, Potassium), pH value, and other soil fertility parameters.
- 3. state_weather_data_1997_2020.csv Contains historical weather information such as average temperature, rainfall, and humidity from 1997 to 2020, enabling correlation analysis with crop yields.

These datasets are integrated and preprocessed in Python to create a unified dataset suitable for machine learning modeling. Preprocessing steps include handling missing values, encoding categorical variables, scaling numerical features, and feature engineering.

Libraries & Frameworks

The project leverages Python libraries for data manipulation, visualization, modeling, and deployment:

- Data Manipulation: pandas, numpy for cleaning, aggregating, and transforming data.
- Visualization: matplotlib, seaborn for exploratory data analysis (EDA), trend visualization, and outlier detection.
- Machine Learning: scikit-learn, xgboost for model building, training, hyperparameter tuning, and evaluation.
- Deployment: Streamlit to build interactive dashboards for predictions and visualization of crop yield insights.
- Others: pickle for saving trained models, joblib for pipeline serialization.

Software & Environment

- Programming Language: **Python** (3.x)
- Version Control: Git for tracking project changes, collaborating, and maintaining reproducible workflows.
- Development Environment: Jupyter Notebook / VS Code for coding, testing, and iterative model development.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Workflow Highlights from Python Notebook

- Data Integration: Merged crop, soil, and weather datasets into a single modeling dataset.
- Exploratory Data Analysis: Identified key trends, seasonal variations, and state-wise yield differences.
- Feature Engineering: Created additional features like lag yield, crop-area ratio, seasonal dummies, and normalized soil parameters.
- Modeling: Trained multiple regression and tree-based models (Random Forest, XGBoost, Gradient Boosting) and evaluated performance using RMSE, MAE, and R².
- Deployment: Developed a Streamlit app for interactive state-wise and crop-wise yield predictions.
- Model Persistence: Saved trained models using pickle for reusability in dashboards.

Potential Data Challenges & Risks in the project

☐ Manual Data Collection

- Data was gathered manually from multiple sources like Kaggle due to lack of centralized datasets.
- Historical state- or district-level crop yield data is often incomplete or inconsistent, especially for minor crops.
- Limited availability of high-quality, long-term weather and climate data affects prediction accuracy.

☐ Input Variable Challenges

- Weather data at coarse resolution (state-level) may not capture microclimate variations critical for yield prediction.
- Soil health and fertilizer usage data are inconsistent and lack long-term coverage, impacting model reliability.
- Variability in crop practices across regions introduces additional uncertainty.

☐ Data Quality and Completeness

- Multiple data sources are siloed and often incompatible, limiting integration potential.
- Significant gaps exist for newer crops or specialized agricultural activities (floriculture, horticulture, mushrooms).
- Time-lags in published datasets reduce their relevance for real-time predictions.

☐ Scalability Issues

• Manual data collection and labelling are not feasible for nationwide or multi-year expansion.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Detailed Plan of Work

| Week | Activities |
|-------------|---|
| Week 1-2 | Data Exploration and Preprocessing – Collected datasets from Kaggle, examined structure and quality, handled missing values, standardized data, and encoded categorical features. |
| Week 3-4 | Feature Engineering and Exploratory Data Analysis (EDA) – Created derived features (e.g., state ranking, total rainfall), visualized patterns across states, crops, and seasons to identify key factors influencing yield. |
| Week 5-6 | Model Training and Hyperparameter Tuning – Built multiple machine learning models such as Random Forest, XGBoost, and Gradient Boosting; optimized hyperparameters and performed cross-validation for robust performance. |
| Week 7-8 | Evaluation, Deployment, and Documentation – Evaluated models using metrics like RMSE, R ² , and MAE; deployed final model via Streamlit/HTML interface; prepared project documentation. |

Steps involved in Project Implementation

☐ Problem Definition

- Identified the need for a predictive system to forecast crop yield (kg/ha) using features such as rainfall, crop area, production, soil type, fertilizer usage, and weather data.
- Objective: support farmers, policymakers, and market stakeholders in operational planning and decision-making.

□ Data Collection

• Gathered crop, soil, and weather datasets from India.gov.in.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

• Extracted relevant attributes and performed preliminary visual analysis, including min, max, and mode estimation for yield, area, and production.

☐ Data Preprocessing

- Cleaned and standardized datasets to remove inconsistencies.
- Handled missing values and encoded categorical variables such as state, crop type, and season.

☐ Exploratory Data Analysis (EDA)

- Visualized trends and patterns across features like state, crop, area, production, rainfall, and fertilizer usage.
- Identified key variables influencing crop yield, seasonal effects, and state-specific trends.

☐ Feature Engineering

- Created new features such as total rainfall per season, crop-state ranking, and normalized soil parameters.
- Enhanced dataset to improve predictive capability of models.

☐ Model Building and Evaluation

- Trained multiple models including Random Forest, XGBoost, and Gradient Boosting.
- Evaluated performance using metrics such as RMSE, R², and MAE; selected the best-performing model for deployment.

☐ Model Deployment

- Developed an interactive interface using Streamlit/HTML.
- Allowed users to input state, crop, and season details to obtain real-time yield predictions.
- Facilitated actionable insights for farmers and policymakers.

Pre-Processing and Feature Engineering

1. Overview

The dataset used in this project is an integrated collection of crop yield, soil characteristics, and state-wise weather data across multiple years and seasons in India. Each record captures key agricultural parameters such as production, cultivated area, fertilizer and pesticide usage, N-P-K soil nutrients, pH value, average temperature, rainfall, and humidity levels, alongside categorical identifiers like crop type, state, and season.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Preprocessing was performed in two main phases:

- 1. Model Training Stage to clean, transform, and merge raw CSV data into a structured dataset (final_merged_statewise_crop_yield_dataset.csv) suitable for model development.
- 2. App Runtime Stage to preprocess user inputs through a unified function preprocess_input(df) ensuring that live app predictions are handled identically to training data.

2. Data Type Correction

Data from multiple government sources (agricultural statistics, IMD weather records, and soil surveys) were found to contain inconsistent formats.

Steps taken include:

- Renaming Ambiguous Columns: Columns such as "Unnamed: 0" or "Crop_Year" were standardized to clear names like year, crop, and state.
- Data Type Alignment: Numerical features were explicitly cast to float64 or int64 to ensure compatibility during scaling and model training.
- Date Parsing: Temporal columns were converted to Python datetime.date objects for proper chronological sorting and analysis.
- In the Streamlit app, inputs captured through date_input() and time_input() widgets are formatted into ISO date strings before processing.

This ensures that both static (CSV) and dynamic (user-input) data maintain uniform formatting.

3. Handling Missing Values

Agricultural datasets often contain gaps due to measurement errors, unrecorded data, or regional inconsistencies. To maintain data integrity:

- Numerical Columns (e.g., avg_temp_c, total_rainfall_mm, avg_humidity_percent, fertilizer, pesticide) were treated using the Interquartile Range (IQR)-based capping method:
 - Lower bound = $Q1 1.5 \times IQR$
 - \circ Upper bound = Q3 + 1.5 × IQR
 - o Values outside this range were capped to the boundary values rather than dropped.
- Categorical Columns with missing entries (such as season or crop) were imputed with the most frequent category.
- Zero-value anomalies in continuous features (e.g., rainfall or production) were replaced using domain logic or median substitution, depending on state-season context.

This approach avoided loss of records and preserved data variance critical for training generalizable models.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

4. Outlier Detection and Treatment

Outlier detection was conducted both statistically and visually:

- Boxplots and Violin Plots were generated for each numerical attribute to visualize distribution skewness and spot extreme deviations.
- Z-Score and IQR methods were compared; however, due to dataset size constraints, extreme values were not aggressively removed.
- Outliers were instead capped or retained if found to represent realistic agricultural variability (e.g., unusually high sugarcane yield in irrigated states).

Visualization techniques used:

- Boxplots for rainfall, humidity, and yield.
- Scatter plots to assess correlation between fertilizer, area, and production.
- Pair plots and correlation heatmaps for feature interdependencies.

This analysis revealed that while certain variables (like rainfall and production) had high variance, they also carried valuable information reflecting real-world fluctuations across states and years.

5. Encoding of Categorical Variables

Machine learning algorithms such as XGBoost and Random Forest require numerical inputs. Therefore:

- Label Encoding was applied to ordinal features such as state and crop.
- One-Hot Encoding was applied to nominal features like season, using drop_first=True to avoid dummy variable traps.
- The same encoding scheme was integrated into the ColumnTransformer within the model pipeline, ensuring that encoding transformations are replicated automatically during inference in the Streamlit app.

This consistency between training and runtime prevents discrepancies between model expectations and user-input data formats.

6. Feature Scaling

Feature scaling was crucial to normalize numeric attributes with varying units (e.g., rainfall in millimeters vs. temperature in Celsius):

- StandardScaler was applied to normally distributed features (avg_temp_c, avg_humidity_percent, ph).
- MinMaxScaler was applied to skewed or non-linear features (area, fertilizer, production).



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

This hybrid scaling improved model stability and reduced training bias, especially for gradient-boosting algorithms.

7. Data Visualization for Summarization

To better understand feature relationships and guide model design, several visualization steps were performed:

- Correlation Heatmap highlighted strong dependencies for instance, fertilizer and area correlated positively with yield.
- PCA (Principal Component Analysis) helped in visualizing feature contribution and dimensionality reduction.
- Scatter Matrix offered a multi-dimensional perspective on interactions among climatic and soil variables.
- State-wise and Season-wise Yield Distribution charts were plotted to identify regional trends and climate influences.

These visual insights directly informed feature engineering and model selection decisions.

8. Summary

The preprocessing pipeline ensured data integrity, consistency, and reproducibility across both training and deployment stages.

By systematically handling missing values, correcting types, encoding categorical variables, and normalizing numerical features, the dataset was transformed into a robust foundation for machine learning.

The resulting final_merged_statewise_crop_yield_dataset reflects a harmonized, clean, and model-ready data representation, enabling accurate and scalable yield predictions across India.

Exploratory Data Analysis (EDA)

1. Objective

The goal of the Exploratory Data Analysis (EDA) was to understand the underlying structure and relationships within the crop yield dataset, detect outliers, identify correlations, and prepare the data for machine learning modelling. The dataset combined crop yield statistics, weather data, and soil characteristics from multiple Indian states across several years.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

2. Initial Data Understanding

The dataset consisted of attributes such as:

- Crop-related features: crop type, season, area, production, and yield
- Weather variables: average temperature, total rainfall, and average humidity percentage
- Soil attributes: nitrogen (N), phosphorus (P), potassium (K), and pH levels
- Geographic and temporal information: state and year

Data types were verified to ensure numerical and categorical consistency. Ambiguous or unnamed columns were renamed, and the year column was standardized for numerical analysis.

3. EDA Visualizations Summary

| No. | Visualization | Variables / | Purpose | Key Insights / Observations |
|-----|---|---|---|--|
| | Type | Features Used | | |
| 1 | Histogram | yield, fertilizer, rainfall, humidity, temperature | To observe distribution, skewness, and spread of numerical variables. | Most variables showed right- skewed distributions; yield and rainfall had few extreme high values indicating climatic or regional differences. |
| 2 | Boxplot | Each numerical column (e.g., yield, production, area, fertilizer) | To detect outliers and understand value range and quartiles. | Outliers were detected primarily in rainfall and fertilizer; capped instead of removed to preserve natural variability. |
| 3 | Correlation Heatmap | All numerical features | To measure linear relationships between variables. | Strong positive correlation between production and area; moderate correlation between yield and fertilizer; negative relation with excessive humidity. |
| 4 | Pair Plot / Scatter Plot | yield vs rainfall, fertilizer, temperature, humidity | To visualize pairwise relationships and possible non-linear patterns. | Yield increases with fertilizer and rainfall up to a threshold; very high rainfall or humidity negatively impacts yield. |
| 5 | Count Plot (Categorical Analysis) | crop, state, season | To visualize frequency and dominance of different categories. | Kharif and Rabi seasons dominate; rice, wheat, and sugarcane are most represented crops across states. |
| 6 | Group-wise Aggregation Plot | state vs mean(yield) and | To compare yield performance across states and crops. | Maharashtra and Punjab showed higher mean yields; crops like |

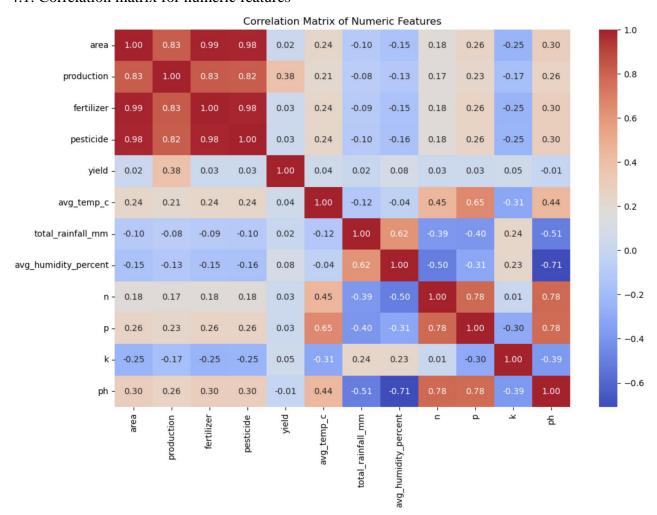


E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

| | | crop vs mean(yield) | | sugarcane and rice perform consistently well. |
|---|---------------|------------------------|----------------------|---|
| 7 | Year-wise | year vs | To examine | Gradual increase in yield till |
| | Trend Plot | mean(yield) | temporal changes in | 2015, slight dip around 2016– |
| | | | yield over years. | 2018 possibly due to monsoon |
| | | | | irregularities. |
| 8 | Soil Nutrient | N, P, K, and pH vs | To explore soil | Balanced NPK ratios correspond |
| | Comparison | yield | parameter effects on | to higher yields; extreme pH |
| | Plot | | productivity. | values show reduced |
| | | | | performance. |

4. EDA Visualizations

4.1. Correlation matrix for numeric features



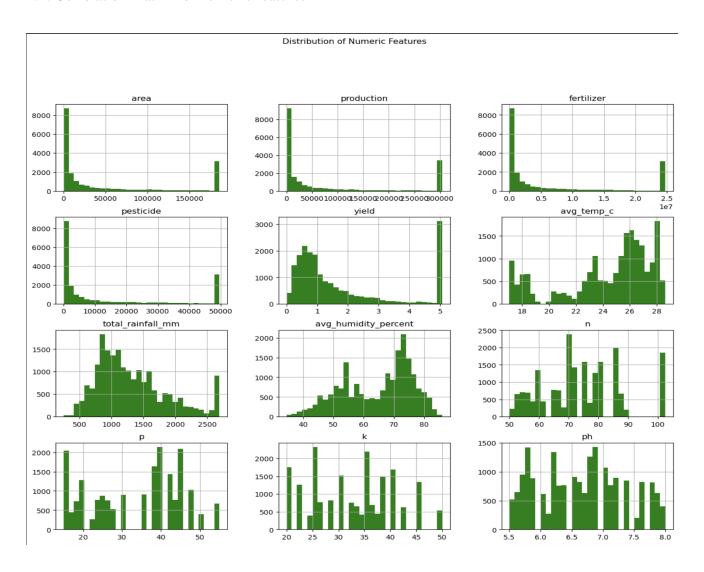
• Values close to +1 (red) → strong positive correlation: as one variable increases, the other tends to increase.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

- Values close to -1 (blue) → strong negative correlation: as one variable increases, the other decreases.
- Values near 0 (white/gray) \rightarrow little or no linear relationship.

4.2. Correlation matrix for numeric features

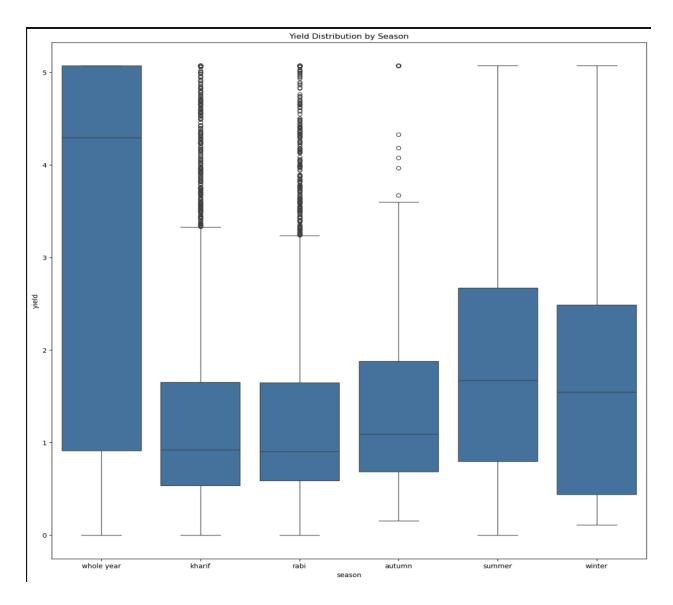


This code visualizes the distribution of key numerical features in the dataset using histograms. It helps identify the spread, skewness, and presence of outliers in variables like area, production, rainfall, and yield during EDA.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

4.3. Boxplots to check outliers grouped by season

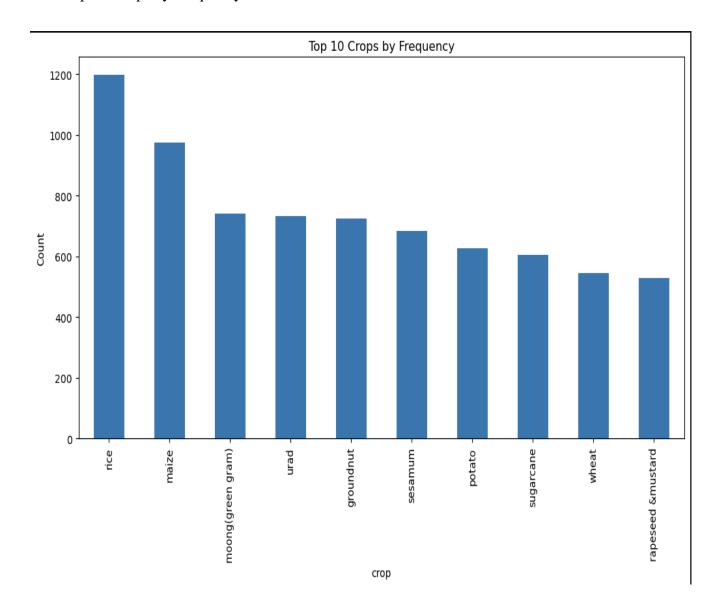


This code creates a **boxplot of crop yield across different seasons** to visualize how yield varies seasonally. It helps in **detecting outliers and comparing yield distribution** among Kharif, Rabi, and other seasons.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

4.4. Top 10 Crops by Frequency

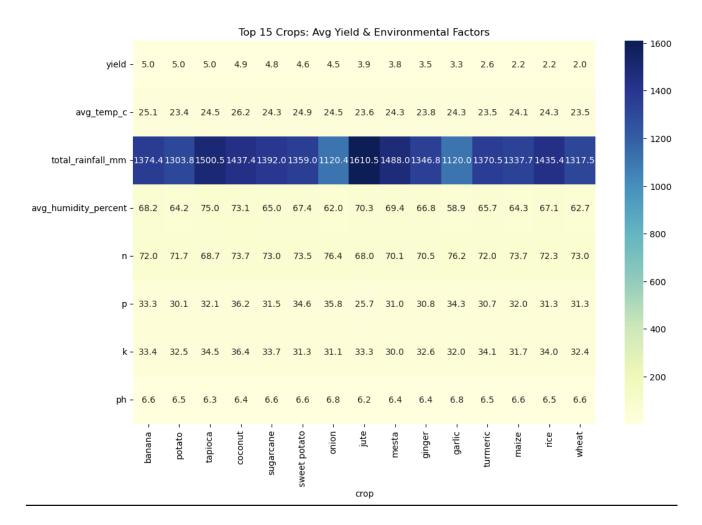


This code plots a bar chart showing the 10 most frequently occurring crops in the dataset. It helps quickly identify which crops are most common.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

4.5. Top 15 Crops: Average Yield & Environmental Factors

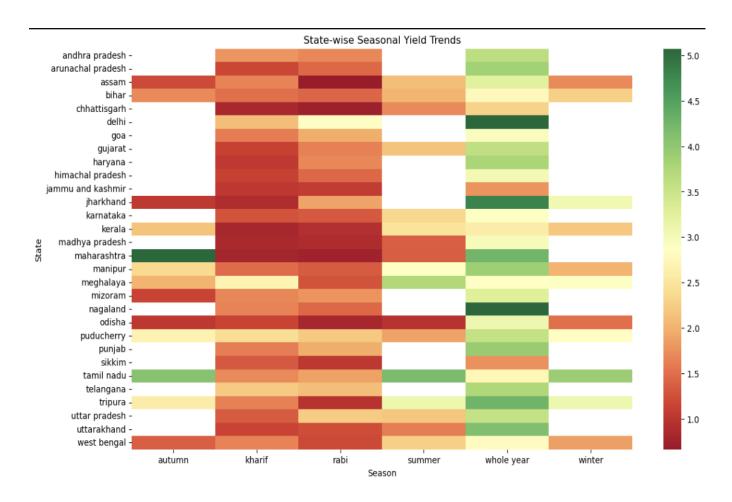


This analysis groups crops by their average yield and associated environmental/soil factors like temperature, rainfall, humidity, and nutrients (N, P, K, pH). A heatmap visualizes the top 15 crops, making it easy to compare how different factors relate to yield. Banana, potato, and tapioca show the highest average yields, typically growing under moderate temperatures (23–26°C) and high rainfall (1300–1500 mm). Nutrient levels and soil pH also vary slightly across crops, indicating specific soil requirements for optimal growth. Overall, the visualization highlights key crops and the environmental conditions that favor higher yields.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

4.6. State-wise Seasonal Yield Trends

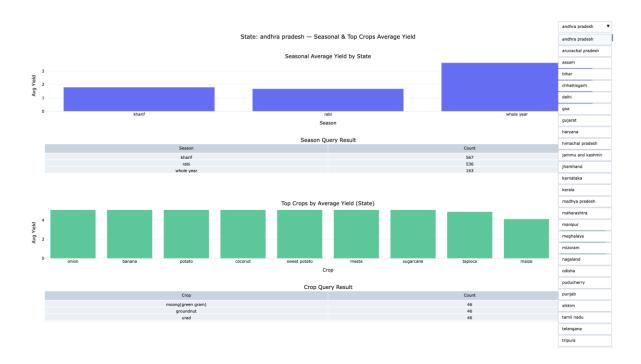


This analysis calculates the average crop yield for each state across different seasons. The resulting heatmap visualizes state-wise seasonal yield trends, highlighting how productivity varies throughout the year. States like Delhi and Andhra Pradesh show high yields during the "kharif" and "whole year" seasons, while regions like Assam and Bihar have more consistent yields across multiple seasons. This helps identify which states perform best in specific seasons and informs targeted agricultural planning.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

5. Interactive state-level dashboard showing



An interactive dashboard was created using Plotly to visualize crop yield trends across Indian states. The dashboard presents both seasonal and crop-based yield patterns for each state, helping to identify key crops and high-performing seasons.

The visualization includes:

- Seasonal Average Yield Chart Displays the mean yield for each agricultural season (Kharif, Rabi, etc.) within a selected state.
- Season Summary Table Shows the number of records available per season for that state.
- Top Crops by Average Yield Chart Highlights the top 10 crops based on their average yield for the selected state.
- Crop Summary Table Lists the most frequently cultivated crops and their occurrence count.
- A dropdown menu enables dynamic selection of states, updating all charts and tables interactively.

Machine Learning Modelling & Techniques Applied

1. Problem Statement

The goal is to predict crop yield using historical agricultural data, including crop type, state, seasonal factors, weather conditions, and soil nutrient levels. Machine learning models were applied to capture complex relationships between these features and crop yield.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

2. Data Preparation

- The dataset was copied to a modeling dataframe df_model to avoid altering the original data.
- The target variable y was defined as yield.
- Columns like production were dropped since they are dependent on area and yield, preventing data leakage.
- Features were categorized for appropriate preprocessing:
 - o Categorical Label Encoding: crop, state
 - o One-Hot Encoding: season
 - Numerical transformations:
 - Log transformation for skewed features: area, fertilizer, pesticide, total_rainfall_mm
 - Standard scaling for temperature and humidity: avg_temp_c, avg_humidity_percent
 - Min-Max scaling for soil nutrients: n, p, k, ph
 - Year Transformation: Scaled to 0–1 range to normalize chronological effects.

3. Preprocessing Pipeline

- A ColumnTransformer was created to apply specific preprocessing steps to each feature type, ensuring reproducible and consistent transformations.
- Custom YearTransformer scales the year to a 0–1 range, preserving temporal information without biasing models due to raw year values.

4. Train/Test Split

- Data was split chronologically:
 - o Training set: years \leq 2017
 - \circ Test set: years > 2017
- This approach mimics real-world prediction scenarios where future data is unknown during training.

5. Machine Learning Models Applied

Several regression models were trained to predict crop yield:

- 1. Linear Regression baseline linear model for simple relationships.
- 2. Random Forest Regressor ensemble tree-based model capturing non-linear patterns.
- 3. HistGradientBoosting Regressor advanced gradient boosting with histogram optimization.
- 4. Support Vector Regressor (SVR) kernel-based model for capturing complex non-linear relationships.
- 5. MLP Regressor feedforward neural network with hidden layers (64, 32).
- 6. XGBoost Regressor gradient boosting algorithm with hyperparameters:



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

- o n_estimators=1000, max_depth=8, learning_rate=0.05
- o subsample=0.8, colsample_bytree=0.9

Each model was integrated into a Pipeline with the preprocessor for streamlined training and evaluation.

6. Model Evaluation

- Models were evaluated using R² score on both training and test sets.
- Results Summary:

| Model | R ² Train | R ² Test |
|------------------|----------------------|---------------------|
| XGBoost | 0.997 | 0.922 |
| RandomForest | 0.990 | 0.903 |
| HistGradient | 0.899 | 0.850 |
| MLPRegressor | 0.506 | 0.403 |
| LinearRegression | 0.202 | 0.151 |
| SVR | 0.070 | -0.027 |

- XGBoost outperformed all models, achieving the highest R² on the test set (0.922), indicating strong predictive performance with minimal overfitting.
- 7. Feature Importance (XGBoost): The top 10 predictive features were identified:

| Feature | Importance | | | | |
|-------------------|------------|--|--|--|--|
| season_whole year | 0.470 | | | | |
| crop | 0.105 | | | | |
| k | 0.061 | | | | |
| n | 0.059 | | | | |
| p | 0.056 | | | | |
| ph | 0.051 | | | | |
| season_winter | 0.033 | | | | |
| season_summer | 0.024 | | | | |
| season_rabi | 0.018 | | | | |
| fertilizer | 0.018 | | | | |

- Seasonal factors (season_whole year, season_winter) and crop type are the most influential predictors.
- Soil nutrients (n, p, k, ph) also contribute significantly, highlighting their role in yield determination.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

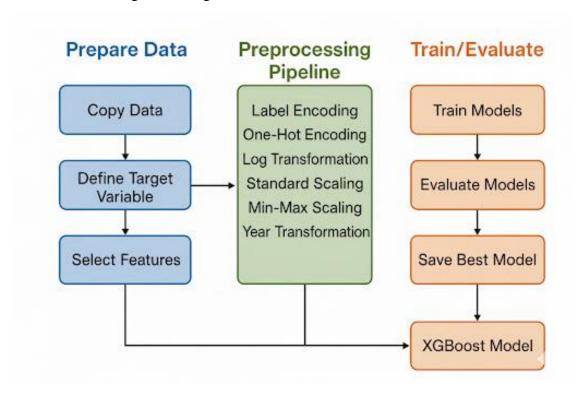
8. Model Deployment

• The best-performing XGBoost model, along with the preprocessing pipeline, was saved as a pickle file (best_model_XGBoost.pkl) for future predictions and integration into applications.

9. Summary

- A comprehensive machine learning pipeline was built, from data preprocessing to model evaluation.
- XGBoost was identified as the most effective model for predicting crop yield.
- The pipeline captures non-linear relationships, handles categorical and numerical transformations, and provides interpretable feature importance.
- This approach ensures a reproducible, scalable, and accurate yield prediction system suitable for practical agricultural applications.

10. Model Building Flow Diagram



11. Code Screenshots

(Cited

https://www.researchgate.net/publication/364028532 Crop Yield Prediction using XG Boost Algorit hm)



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com



(Cited: https://ijsret.com/wp-content/uploads/IJSRET_V11_issue3_1159.pdf)



| | Understanding and Preprocessing the data | | | | | | | | | | | | |
|-------|--|------------------------------------|--------|------|--------|------|-------|---------|------------|-----------------|------------|-----------|-------------|
| [1835 | <pre># Load all three datasets crop_yield = pd.read_csv("crop_yield.csv") state_soil = pd.read_csv("state_soil_data.csv") state_weather = pd.read_csv("state_weather_data_1997_2020.csv") # Inspect first few rows crop_yield_head = crop_yield.head(10)</pre> | | | | | | | | | | | | |
| | | ate_soil_head = ate_weather_hea | | | | | | 10) | | | | | |
| [1837 | | int("\nCrop Yei op_yield_head | | прle | | set | |) | | | | | |
| | Crop Yeild Sample Dataset : | | | | | | | | | | | | |
| [1837 | | Crop Crop | p_Year | | Seas | on | State | Area | Production | Annual_Rainfall | Fertilizer | Pesticide | Yield |
| | 0 | Arecanut | 1997 | W | hole Y | ear | Assam | 73814.0 | 56708 | 2051.4 | 7024878.38 | 22882.34 | 0.796087 |
| | | Arhar/Tur | 1997 | | Kh | arif | Assam | 6637.0 | 4685 | 2051.4 | 631643.29 | 2057.47 | 0.710435 |
| | 2 | Castor seed | 1997 | | Kh | arif | Assam | 796.0 | 22 | 2051.4 | 75755.32 | 246.76 | 0.238333 |
| | 3 | Coconut | 1997 | W | hole Y | ear | Assam | 19656.0 | 126905000 | 2051.4 | 1870661.52 | 6093.36 | 5238.051739 |
| | 4 | Cotton(lint) | 1997 | | Kh | arif | Assam | 1739.0 | 794 | 2051.4 | 165500.63 | 539.09 | 0.420909 |
| | 5 | Dry chillies | 1997 | W | hole Y | ear | Assam | 13587.0 | 9073 | 2051.4 | 1293074.79 | 4211.97 | 0.643636 |
| | 6 | Gram | 1997 | | R | abi | Assam | 2979.0 | 1507 | 2051.4 | 283511.43 | 923.49 | 0.465455 |
| | 7 | Jute | 1997 | | Kh | arif | Assam | 94520.0 | 904095 | 2051.4 | 8995468.40 | 29301.20 | 9.919565 |
| | 8 | Linseed | 1997 | | R | abi | Assam | 10098.0 | 5158 | 2051.4 | 961026.66 | 3130.38 | 0.461364 |
| | 9 | Maize | 1997 | | Kh | arif | Assam | 19216.0 | 14721 | 2051.4 | 1828786.72 | 5956.96 | 0.615652 |
| [1839 | | int("\nState So ate_soil_head | il Sa | nple | | set | |) | | | | | |
| | St | ate Soil Sample | Data | set | | | | | | | | | |
| [1839 | | state | N | P | ĸ | pН | | | | | | | |
| | 0 | Andhra Pradesh | 78 | 45 | 22 | 6.8 | | | | | | | |
| | | Arunachal Pradesh | 55 | 15 | 35 | 5.5 | | | | | | | |
| | 2 | Assam | 60 | 18 | 38 | 5.8 | | | | | | | |
| | 3 | Bihar | r 85 | 30 | 25 | 7.2 | | | | | | | |
| | 4 | Chhattisgarh | 70 | 35 | 20 | 6.5 | | | | | | | |
| | 5 | Delh | i 90 | 40 | 30 | 7.5 | | | | | | | |
| | 6 | Goa | 65 | 25 | 45 | 6.2 | | | | | | | |



```
[1841... print("\nState Weather Sample Dataset : \n") state_weather_head
         State Weather Sample Dataset :
                   state year avg_temp_c total_rainfall_mm avg_humidity_percent
        0 Andhra Pradesh 1997 28.21 1191.08
         1 Andhra Pradesh 1998

        2
        Andhra Pradesh
        1999
        28.03
        603.67

        3
        Andhra Pradesh
        2000
        27.74
        1070.25

                                                                            66.91

        4
        Andhra Pradesh
        2001
        28.08
        910.13
        68.69

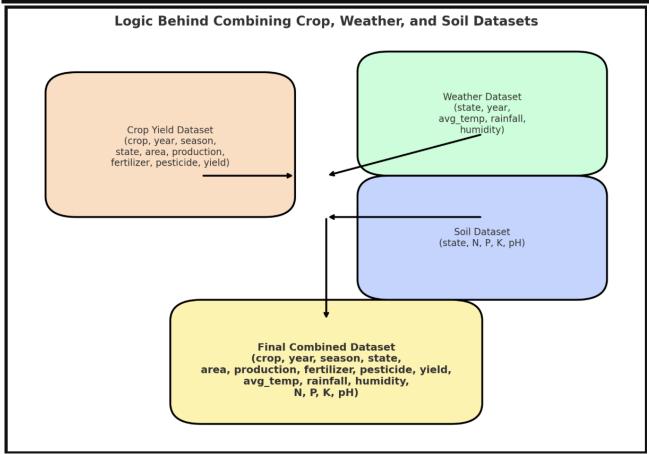
        5
        Andhra Pradesh
        2002
        28.54
        768.22
        66.52

                                                                   68.83
        6 Andhra Pradesh 2003 28.31 857.23
        7 Andhra Pradesh 2004
                                              1192.26
        8 Andhra Pradesh 2005
9 Andhra Pradesh 2006
                                      27.95
                                                      1343.62
[1843_ # Unique States from all three datasets
crop_yield_states = set(crop_yield['State'].unique())
state_soil_states = set(state_soil['state'].unique())
state_weather_states = set(state_weather['state'].unique())
        print("\n------ Unique States from Crop Yeil Dataset : -----\n", crop_yield_states)
print("\n----- Unique States from State Soil Dataset : ----\n", state_soil_states)
print("\n----- Unique States from State Weather Dataset : ----\n", state_weather_states)
```



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

| Final Attribute | Source | How It's Combined | Reasoning |
|----------------------|---|---|--|
| crop | crop_yield.Crop | Taken directly. | Crop name is unique to the crop yield dataset — no in soil/weather. |
| year | crop_yield.Crop_Year | Align with weather.year. | Crop yield is for a specific year; matched with state's weather for the same year. |
| season | crop_yield.Season | Taken directly. | Only available in crop data — defines growing period. |
| state | crop_yield.State | Join key with soil + weather. | All three datasets share state-level info. |
| area | crop_yield.Area | Taken directly. | Crop-specific land area (hectares). |
| production | crop_yield.Production | Taken directly. | Crop-specific production (tonnes). |
| fertilizer | crop_yield.Fertilizer | Taken directly. | Represents total fertilizer use in that crop's area. |
| pesticide | crop_yield.Pesticide | Taken directly. | Crop-specific pesticide usage. |
| yield | crop_yield.Yield | Taken directly (but cross-check with Production/Area). | This is the target variable of interest (output). |
| avg_temp_c | state_weather_data.avg_temp_c | Join on (State, Year) . | State's yearly average temperature affects yield. |
| total_rainfall_mm | state_weather_data.total_rainfall_mm | Join on (State, Year) . | State's yearly rainfall drives water availability. |
| avg_humidity_percent | state_weather_data.avg_humidity_percent | Join on (State, Year) . | Humidity affects evapotranspiration and disease risk. |
| N | state_soil_data.N | Join on State . | Soil nutrient baseline, doesn't change year-to-yea |
| P | state_soil_data.P | Join on State . | Same reasoning as N. |
| К | state_soil_data.K | Join on State . | Same reasoning as N. |
| pН | state_soil_data.pH | Join on State . | Soil acidity/alkalinity constraint; constant for state. |





```
# Standardize column names
crop_yield.columns = crop_yield.columns.str.strip().str.lower()
state_weather.columns = state_weather.columns.str.strip().str.lower()
state_soil.columns = state_soil.columns.str.strip().str.lower()
         # After Standardizing column names
print("\n1. crop_yield.columns: \n", crop_yield.columns)
print("\n1. state_weather.columns: \n", state_weather.columns)
print("\n2. state_soil.columns: \n", state_soil.columns)
                                     -- Initial column names --

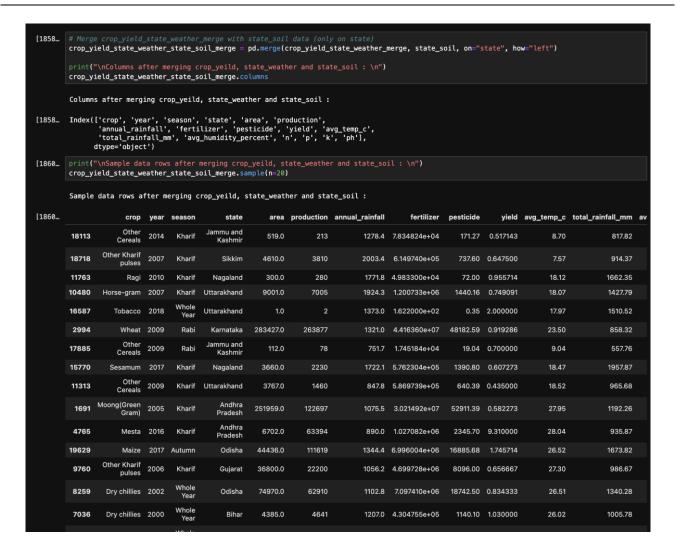
    crop_yield.columns:
    Index(['Crop', 'Crop_Year', 'Season', 'State', 'Area', 'Production', 'Annual_Rainfall', 'Fertilizer', 'Pesticide', 'Yield'], dtype='object')

         3. state_soil.columns :
   Index(['state', 'N', 'P', 'K', 'pH'], dtype='object')
                                     -- After Standardizing column names -

    crop_yield.columns:
    Index(['crop', 'crop_year', 'season', 'state', 'area', 'production', 'annual_rainfall', 'fertilizer', 'pesticide', 'yield'], dtype='object')

          3. state_soil.columns :
   Index(['state', 'n', 'p', 'k', 'ph'], dtype='object')
             "crop_year" in crop_yield.columns:
    crop_yield.rename(columns={"crop_year": "year"}, inplace=True)
         print("\ncrop_yield.columns : \n", crop_yield.columns)
         # Merge crop_yield + state_weather on (state, year)
crop_yield_state_weather_merge = pd.merge( crop_yield, state_weather, on=["state", "year"], how="left" )
[1854...
         \label{linear_print} $$ print("\ncolumns after merging crop_yeild and state_weather : \n") $$ crop\_yield_state_weather_merge.columns $$
         Columns after merging crop_yeild and state_weather :
Sample data rows after merging crop_yeild and state_weather :
[1856...
                         crop year season
                                                             area production annual_rainfall
                                                                                                                             yield avg_temp_c total_rainfall_mm avg_hur
                    Coriander 2000 Whole
Year
                                              Madhya 94705.0
Pradesh
          7202
                                                                       27475
                                                                                        723.3 9297189.85 24623.30 0.307674
                                                                                                                                          25.33
                                                                                                                                                            867.75
         18740
                         Rice 2009 Kharif
                                                Sikkim 12270.0
                                                                                       2176.6 1911911.40 2085.90 1.595000
                                                                                                                                                            924.82
                     Turmeric 2002 Whole Arunachal
Year Pradesh
                                                                                      2832.8 48660.38 128.50 4.459091
                                                                                                                                                           1782.25
                                                Jammu
and
Kashmir
                      other 2005 Kharif
         17736
                                                            20.0
                                                                                       1721.8
                                                                                                 2398.40
                                                                                                                4.20 0.415000
                                                                                                                                           8.63
                                                                                                                                                            575.81
          1315
                         Rice 2003 Kharif
                                                  Goa 35471.0
                                                                     115996
                                                                                       3011.6 3510919.58 8513.04 3.270000
                                                                                                                                          27.88
                                                                                                                                                           1413.96
                    Other Rabi
pulses 2010
                                                 West
Bengal
          11972
                                       Rabi
                                                            255.0
                                                                                       1096.0
                                                                                                  42358.05
                                                                                                                61.20 0.446667
                                                                                                                                          26.24
                                                                                                                                                           1356.15
         17894 Sesamum 2009 Kharif
                                                          4624.0
                                                                     2031
                                                                                    751.7 720511.68 786.08 0.443333
                                                                                                                                          9.04
                                                                                                                                                           557.76
```







```
[1862... # Select final required columns (dropping off annual_rainfall, since we have total_rainfall_mm column)
           crop_yield_state_weather_state_soil_merge = crop_yield_state_weather_state_soil_merge[
                   crop, 'year', 'season', 'state', 'area', 'production', 'fertilizer', 'pesticide',
'yield', 'avg_temp_c', 'total_rainfall_mm', 'avg_humidity_percent', 'n', 'p', 'k', 'ph']
           crop_yield_state_weather_state_soil_merge.to_csv("final_merged_statewise_crop_yeild_dataset.csv", index=False)
           print("\nFinal merged dataset created with shape :\n", crop_yield_state_weather_state_soil_merge.shape)
           Final merged dataset created with shape : (19689, 16)
[1864... df = pd.read_csv("final_merged_statewise_crop_yeild_dataset.csv")
[1866... df.columns.tolist()
             'year',
'season',
'state',
'area',
'production',
'fertilizer',
             'pesticide',
'pesticide',
'yield',
'avg_temp_c',
'total_rainfall_mm',
'avg_humidity_percent',
'n',
              'ph']
[1868... df.info()
           <class 'pandas.core.frame.DataFrame'>
RangeIndex: 19689 entries, 0 to 19688
Data columns (total 16 columns):
# Column Non-Null Co
                                                   Non-Null Count Dtype
                 0 crop
1 year
2 season
3 state
4 area
5 production
6 fertilizer
7 pesticide
8 yield
9 avg_temp_c
10 total_rainfa
11 avg_humidity
12 n
                                                                           float64
                                                   19689 non-null
            dtypes: float64(8), int64(5), object(3) memory usage: 2.4+ MB
             # Change year column to object (categorical)
df['year'] = df['year'].astype(object)  # or .astype('object')
 [1872... df.info()
             <class 'pandas.core.frame.DataFrame'>
             RangeIndex: 19689 entries, 0 to 19688
Data columns (total 16 columns):
# Column Non-Null Count Dtype
           object
object
                                                                           object
object
float64
                                                                          float64
int64
float64
float64
float64
float64
float64
                                                                           int64
int64
                                                                           int64
```



```
[1878…  # Check missing values
[1878... crop
year
season
                        state
                      area
production
fertilizer
pesticide
yield
                       avg_temp_c
total_rainfall_mm
avg_humidity_percent
                        dtype: int64
 [1880... # Check duplicate
                      df.duplicated().sum()
[1880... 0
 [1882... # Check for Numerical columns
df.select_dtypes(include=['number']).columns.tolist()
['area',
'production',
'fertilizer',
'pesticide',
'yield',
'avg_temp_c',
'total_rainfall_mm',
'avg_humidity_percent',
'n',
'n',
                           'ph']
                      # Check for Categorical columns
categorical_columns = df.select_dtypes(include=['object']).columns.tolist()
                        categorical_columns
[1884... ['crop', 'year', 'season', 'state']
                    # Check categorical unique values
print("\nTotal Unique years:", df['year'].nunique())
print("\n", df['year'].unique())
print("\n", df['rcop'].unique())
print("\nTotal Unique states:", df['state'].unique())
print("\n" off['state'].unique())
print("\n" off['state'].unique())
print("\n", df['season'].unique())
                       Total Unique years: 24
                          [1997 1998 1999 2000 2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020]
                        Total Unique crops: 55
                         ['Arecanut' 'Arhar/Tur' 'Castor seed' 'Coconut' 'Cotton(lint)'
'Dry chillies' 'Gram' 'Jute' 'Linseed' 'Maize' 'Mesta' 'Niger seed'
'Onion' 'Other Rabi pulses' 'Potato' 'Rapeseed &Mustard' 'Rice'
'Sesamum' 'Small millets' 'Sugarcane' 'Sweet potato' 'Tapioca' 'Tobacco'
'Turmeric' 'Wheat' 'Bajra' 'Black pepper' 'Cardamom' 'Coriander' 'Garlic'
'Ginger' 'Groundnut' 'Horse-gram' 'Jowar' 'Ragi' 'Cashewnut' 'Banana'
'Soyabean' 'Barley' 'Khesari' 'Masoor' 'Moong(Green Gram)'
'Other Kharif pulses' 'Safflower' 'Sannhamp' 'Sunflower' 'Urad'
'Peas & beans (Pulses)' 'other oilseeds' 'Other Cereals' 'Cowpea(Lobia)'
'Oilseeds total' 'Guar seed' 'Other Summer Pulses' 'Moth']
                        Total Unique states: 30
                          ['Assam' 'Karnataka' 'Kerala' 'Meghalaya' 'West Bengal' 'Puducherry' 'Goa' 'Andhra Pradesh' 'Tamil Nadu' 'Odisha' 'Bihar' 'Gujarat' 'Madhya Pradesh' 'Maharashtra' 'Mizoram' 'Punjab' 'Uttar Pradesh' 'Haryana' 'Himachal Pradesh' 'Tripura' 'Nagaland' 'Chhattisgarh' 'Uttarakhand' 'Jharkhand' 'Delhi' 'Manipur' 'Jammu and Kashmir' 'Telangana' 'Arunachal Pradesh' 'Sikkim']
                        Total Unique season: 6
                          ['Whole Year ' 'Kharif 'Winter ']
                                                                                                                                         ''Autumn ''Summer
                                                                                            ' 'Rabi
```



```
1. Dataset Overview
           • Rows: 19,689
           • Columns: 16
           • Types:
                • Categorical: crop, season, state, year
               • Numerical: area, production, fertilizer, pesticide, yield, avg_temp_c, total_rainfall_mm, avg_humidity_percent, N, P, K, pH
          2. Missing / Null Values
           1. No missing values detected — all columns are complete.
          3. Categorical Variables
           1. crop → crop name (string)
           2. season → cropping season (e.g., Kharif, Rabi, Whole Year)
           3. state → Indian states/regions
           4. year → yeild year (ranging from ~1997 onward)
           5. Cleaning to apply:
               A. Standardize text (remove trailing spaces, unify case, handle alternate spellings).
               B. Encode categories if needed (Label Encoding / One-Hot Encoding for ML models).
          4. Numerical Variables
           1. area, production, fertilizer, pesticide, yield → Continuous values
           2. avg_temp_c, total_rainfall_mm, avg_humidity_percent → Climate variables
           3. N, P, K, pH → Soil attributes
           4. Cleaning to apply:
               A. Check for outliers (e.g., abnormally high production values or negative values).
               B. Convert year into datetime index if doing time-series analysis.
               C. Normalize/scale variables for ML (e.g., MinMaxScaler or StandardScaler).
          5. Data Cleaning Steps Applied / Recommended
           1. Verified no null/missing values.
           2. Identify outliers in area, production, fertilizer, pesticide, vield,
           3. Standardize categorical strings (crop, season, state, year).
[1889... # Function to detect outliers using IQR (Interquartile Range)
```

```
[1889. # Function to detect outliers using IQR (Interquartile Range)
#IQR = 03 - 01
#01 (25th percentile) = value below which 25% of data lies.
#03 (75th percentile) = value below which 55% of data lies.
#IQR measures the spread of the middle 50% of the data.

def detect_outliers(df, col):
    01 = df(col).quantile(0.75)
    108 = 03 - 01
    lower = 01 - 1.5 * IOR
    upper = 03 + 1.5 * IOR
    upper = 03 + 1.5 * IOR
    outliers = df([df(col] < lower) | (df(col] > upper)]
    return outliers

# Get numerica columns
numerical_columns = df.select_dtypes(include=[np.number]).columns

# Loop through numeric columns
for col in numerical_columns:
    outliers = detect_outliers(df, col)
    if not outliers.empty:
        print(f"Nos" Column: (col)")
        print("Sample outliers:")
        print("Sample outliers:")
        print("Sample outliers:")
        print(sample outliers:")
        print(siguze=[a,8))
        sns.boxplot(x=df(col))
        plt.figure(figsize=[a,8))
        sns.boxplot(x=df(col))
        plt.show()
    else:
        print("Nos outliers detected in Column: {col}")
```



```
# Apply outlier handling to all numeric columns
numerical_columns = df.select_dtypes(include=[np.number]).columns
for col in numerical_columns:
    data = cap_outliers(df, col)
[1893... #Verify again if outliers exists
           # Function to detect outliers using IOR
def detect_outliers(df, col):
01 = df[col].quantile(0.25)
03 = df[col].quantile(0.75)
IOR = 03 - 01
                 IQR = Q3 - Q1
lower = Q1 - 1.5 * IQR
upper = Q3 + 1.5 * IQR
outliers = df[(df[col] < lower) | (df[col] > upper)]
return outliers
            numerical_columns = df.select_dtypes(include=[np.number]).columns
            for col in numerical_columns:
                 cot a numeriza_cottamic
outliers = detect_outliers(df, col)
if not outliers.empty:
    print(f"\n\formall column: {col}")
    print(f"\n\formall column: {col}")
    print("Sample outliers:")
    print(outliers[[col]].head(20))
                       plt.figure(figsize=(8,8))
sns.boxplot(x=df[col])
                       plt.title(f"Outliers in {col}")
                       plt.show()
                 else:
                      print(f"\nNo outliers detected in Column : {col}")
            No outliers detected in Column : area
            No outliers detected in Column : production
            No outliers detected in Column : fertilizer
            No outliers detected in Column : pesticide
            No outliers detected in Column : yield
            No outliers detected in Column : avg temp c
            No outliers detected in Column : total rainfall mm
            No outliers detected in Column : avg_humidity_percent
            No outliers detected in Column : n
            No outliers detected in Column : p
            No outliers detected in Column : k
            No outliers detected in Column : ph
[1895...  # Check the data size after handling outliers print("\nData Shape: \n", df.shape)
            Data Shape : (19689, 16)
```



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

(Cited: Nigam A, Garg S, Agrawal A, Agrawal P. Crop yield prediction using machine learning algorithms. In: Proceedings of the Fifth International Conference on Image Information Processing (ICIIP); 2019; Shimla, India. IEEE; 2019. https://doi.org/10.1109/ICIIP47207.2019.898595)

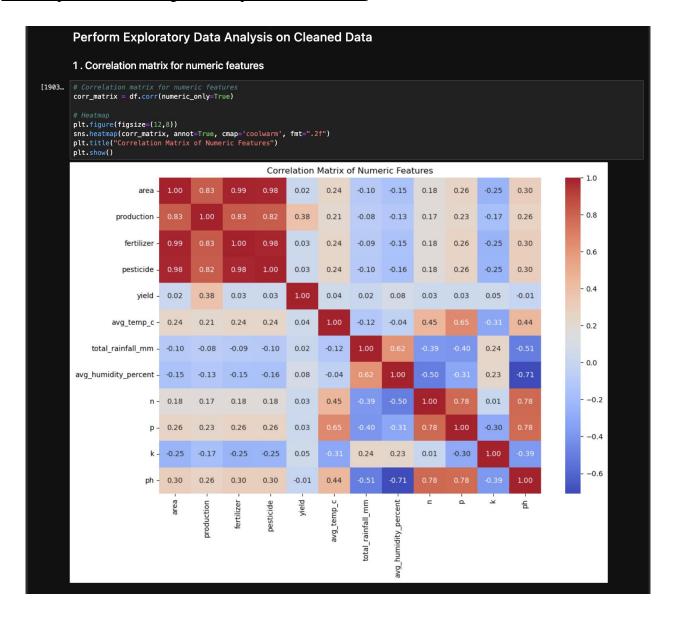
```
[1899...  # Fix text formatting - making all values from categorical columns in lowercase
                                print("\nCrop Values Before making lowercase :\n", df['crop'].unique())
print("\nSeason Values Before making lowercase :\n", df['season'].unique())
print("\nState Values Before making lowercase :\n", df['state'].unique())
                                for col in ['crop', 'season', 'state']:
    data[col] = df[col].astype(str).str.strip().str.lower()
                               print("\nCrop Values After making lowercase :\n", df['crop'].unique())
print("\nSeason Values After making lowercase :\n", df['season'].unique())
print("\nState Values After making lowercase :\n", df['state'].unique())
                               Crop Values Before making lowercase:
['Arecanut' 'Arhar/Tur' 'Castor seed' 'Coconut' 'Cotton(lint)'
'Dry chillies' 'Gram' 'Jute' 'Linseed' 'Maize' 'Mesta' 'Niger seed'
'Onion' 'Other Rabi pulses' 'Potato' 'Rapeseed &Mustard' 'Rice'
'Sesamum' 'Small millets' 'Sugarcane' 'Sweet potato' 'Tapioca' 'Tobacco'
'Turmeric' 'Wheat' 'Bajra' 'Black pepper' 'Cardamom' 'Coriander' 'Garlic'
'Ginger' 'Groundnut' 'Horse-gram' 'Jowar' 'Ragi' 'Cashewnut' 'Banana'
'Soyabean' 'Barley' 'Khesari' 'Masoor' 'Moong(Green Gram)'
'Other Kharif pulses' 'Safflower' 'Sannhamp' 'Sunflower' 'Urad'
'Peas & beans (Pulses)' 'other oilseeds' 'Other Cereals' 'Cowpea(Lobia)'
'Oilseeds total' 'Guar seed' 'Other Summer Pulses' 'Moth']
                                  Season Values Before making lowercase :
['Whole Year ' 'Kharif ' 'Rabi ' 'Autumn
                                                                                                                                                                                                                                                                       ' 'Summer
                                        ['Whole Year '
'Winter ']
                                State Values Before making lowercase:
['Assam' 'Karnataka' 'Kerala' 'Meghalaya' 'West Bengal' 'Puducherry' 'Goa'
'Andhra Pradesh' 'Tamil Nadu' 'Odisha' 'Bihar' 'Gujarat' 'Madhya Pradesh'
'Maharashtra' 'Mizoram' 'Punjab' 'Uttar Pradesh' 'Haryana'
'Himachal Pradesh' 'Tripura' 'Nagaland' 'Chhattisgarh' 'Uttarakhand'
'Jharkhand' 'Delhi' 'Manipur' 'Jammu and Kashmir' 'Telangana'
'Arunachal Pradesh' 'Sikkim']
                               Crop Values After making lowercase:
['arecanut' 'arhar/tur' 'castor seed' 'coconut' 'cotton(lint)'
'dry chillies' 'gram' 'jute' 'linseed' 'maize' 'mesta' 'niger seed'
'onion' 'other rabi pulses' 'potato' 'rapeseed &mustard' 'rice'
'sesamum' 'small millets' 'sugarcane' 'sweet potato' 'tapioca' 'tobacco'
'turmeric' 'wheat' 'bajra' 'black pepper' 'cardamom' 'coriander' 'garlic'
'ginger' 'groundnut' 'horse-gram' 'jowar' 'ragi' 'cashewnut' 'banana'
'soyabean' 'barley' 'khesari' 'masoor' 'moong(green gram)'
'other kharif pulses' 'safflower' 'sannhamp' 'sunflower' 'urad'
'peas & beans (pulses)' 'other oilseeds' 'other cereals' 'cowpea(lobia)'
'oilseeds total' 'guar seed' 'other summer pulses' 'moth']
                                 Season Values After making lowercase :
['whole year' 'kharif' 'rabi' 'autumn' 'summer' 'winter']
                                State Values After making lowercase:
['assam' 'karnataka' 'kerala' 'meghalaya' 'west bengal' 'puducherry' 'goa'
'andhra pradesh' 'tamil nadu' 'odisha' 'bihar' 'gujarat' 'madhya pradesh'
'maharashtra' 'mizoram' 'punjab' 'uttar pradesh' 'haryana'
'himachal pradesh' 'tripura' 'nagaland' 'chhattisgarh' 'uttarakhand'
'jharkhand' 'delhi' 'manipur' 'jammu and kashmir' 'telangana'
'arunachal pradesh' 'sikkim']
```



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

(Cited :

https://www.researchgate.net/publication/393426488 Multivariate Regressive Gradient Spiral Optimi zed_Deep_Belief_Learning_For_Crop_Yield_Prediction)

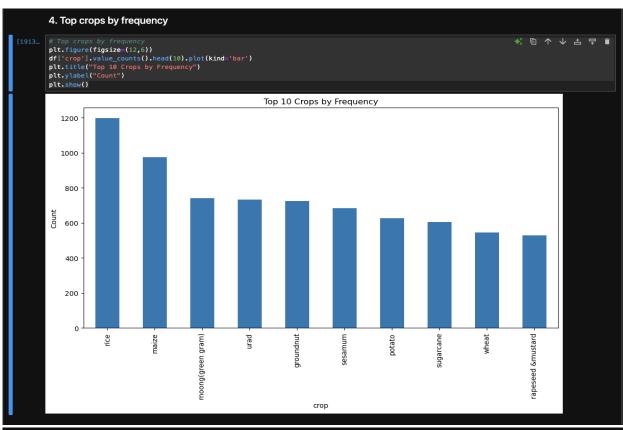




E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com



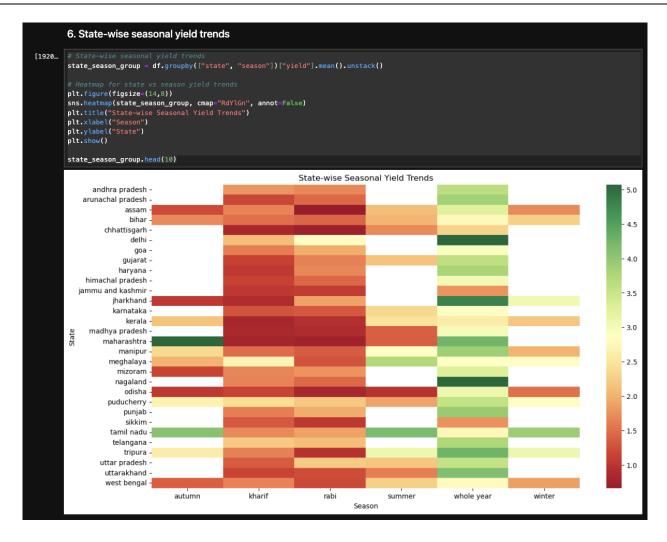








E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com



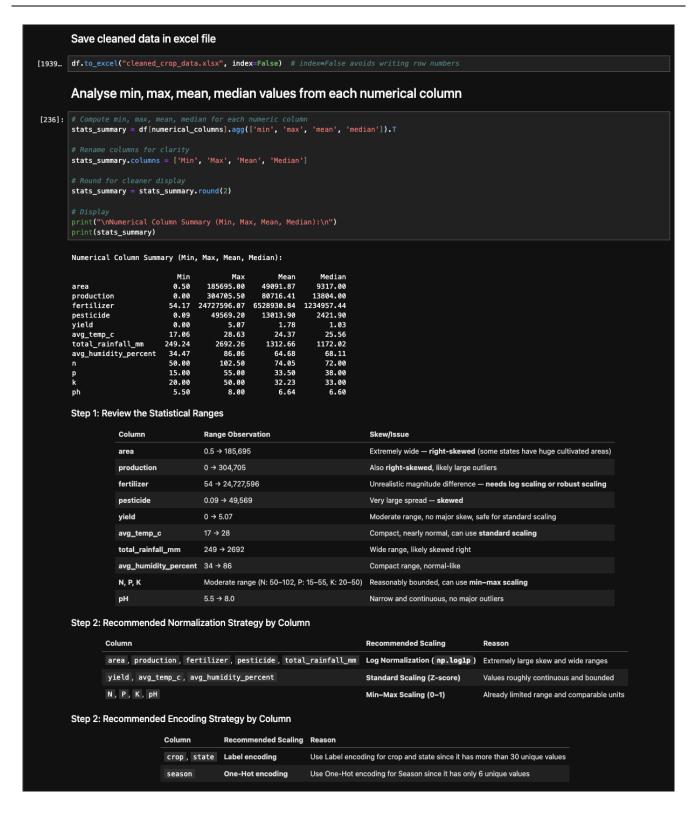


E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

```
7. Plot yield vs environmental factors for each crop
# Select major crops for analysis
major_crops = ["rice", "wheat", "sugarcane"]
 # Plot yield vs environmental factors for each crop
factors = ["total_rainfall_mm", "avg_temp_c", "avg_humidity_percent", "n", "p", "k", "ph"]
plt.figure(figsize=(16,12))
 for i, crop in enumerate(major_crops):
    subset = df[df["crop"] == crop]
    for j, factor in enumerate(factors):
        plt.subplot(len(major_crops), len(factors), i*len(factors) + j + 1)
        sns.scatterplot(data=subset, x=factor, y="yield", alpha=0.5)
        corr = subset["yield"].corr(subset[factor])
        plt.title(f"{crop}\n{factor}\nCorr={corr:.2f}")
        plt.xlabet("")
        plt.ylabel("")
 plt.tight_layout()
plt.show()
                                                                                                         avg_humidity_percent
Corr=-0.54
                 sugarcane
al_rainfall_m
Corr=-0.08
                                                                                                                                                                                                                       sugarcane
                                                                                                                                                                                                                                                                                                                        sugarcane
ph
Corr=0.17
                                                                                                                                                                     Corr=0.24
                                                                                                                                                                                                                      Corr=0.21
                                                                                                                                                                                                                                                                        Corr=0.05
                               2000
```



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com





E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Build Regression Model for yield predictions

1. Why Dropping production Makes Sense

Why it's correct

production is usually derived from [\text{production} = \text{yield} \times \text{area}] So if you keep production, you create data leakage — because the target (yield) can be indirectly reconstructed.

Impact

Dropping production prevents overfitting and gives realistic accuracy. If you included it before, the model likely had artificially high R² (0.9+) — now you're seeing the true generalization performance (~0.4–0.5).

2. Why Dropping year Is Tricky

Why dropping year can make sense

- If year is just a sequence number (1997–2020) and doesn't directly influence yield, it can confuse the model.
- Helps avoid overfitting to time order.

Why dropping year might hurt accuracy

- Yield often improves over time due to technological advancement, fertilizer use, or climate trends.
- Without year , the model can't learn this temporal trend.

Fix:

You don't have to reintroduce raw year directly. Instead, use derived time-based features that carry trend information but don't cause leakage.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

```
# Example: train on years <= 2017, test on years > 2017
train_mask = X['year'].astype(int) <= 2017
X_train = X[train_mask]
X_test = X[~train_mask]
X_test = X[~train_mask]
y_test = y[~train_mask]
 xgb_params = dict(
    n_estimators=1000,
            max_depth=8,
learning_rate=0.05,
subsample=0.8,
colsample_bytree=0.9,
            random_state=42
  models = {
                LinearRegression": LinearRegression(),
                                        rest": RandomForestRegressor(random_state=42),
ient": HistGradientBoostingRegressor(random_state=42),
              'Historausent , misser | historausent , misser |
'SVR': SVR(kernel='rbf'),
'MLPRegressor': MLPRegressor(hidden_layer_sizes=(64,32), max_iter=1000, random_state=42),
'XGBoost": XGBRegressor(***xgb_params)
best_r2 = -np.inf
best_model_name = None
best_pipeline = None
 for name, model in models.items():
    pipe = Pipeline([
                       ('preprocessor', preprocessor), ('model', model)
            pipe.fit(X_train, y_train)
y_train_pred = pipe.predict(X_train)
y_test_pred = pipe.predict(X_test)
            r2_train = r2_score(y_train, y_train_pred)
r2_test = r2_score(y_test, y_test_pred)
                        'Model': name,
'R2_Train': r2_train,
'R2_Test': r2_test
           if r2_test > best_r2:
  best_r2 = r2_test
  best_model_name = nam
  best_pipeline = pipe
✓ -
                                                                             ---- Model Comparison (R2 Train/Test): ---
       | Model | R2_Train | R2_Test | XGBoost | 0.997092 | 0.921520 | 0.980678 | 0.985978 | 0.985973 | 0.850320 | 0.505554 | 0.403017 | 0.151068 | 0.69679 | 0.027287 | 0.151068 | 0.696679 | 0.0270858 | 0.069679 | 0.0270858 | 0.069679 | 0.0270858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.069679 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0086858 | 0.0
 1. Model Performance (R<sup>2</sup> Train/Test)
                                                                             R<sup>2</sup>_Train R<sup>2</sup>_Test
                                                                             0.997
                                                                                                 0.922 Excellent fit; test R<sup>2</sup> slightly lower than train, minor overfitting but reasonable for real-world agricultural data.
                                   XGBoost
                                                                                                    0.903 Slightly worse than XGBoost; consistent.
                                   RandomForest
                                                                            0.990
                                   HistGradient
                                                                             0.899
                                                                                                   0.850 Good performance; weaker than ensemble trees (RF/XGB).
                                                                                                    0.403 Poor; likely due to not enough tuning or small feature scaling issues.
                                   MLPRegressor
                                                                             0.506
                                   LinearRegression 0.202 0.151 Expected; crop yield relationships are highly non-linear.
                                                                            0.070 -0.027 Fails here; SVR with default kernel cannot handle this high-dimensional, heterogeneous data.

▼ Observation:

    1. Tree-based models dominate (XGBoost > RF > HistGradient).
    2.\ Linear/MLP/SVR\ perform\ poorly\ -- \ logical\ given\ dataset\ complexity\ and\ mixed\ numeric/categorical\ features.
    3. So these numbers are logically consistent with the dataset and typical crop yield prediction tasks.
```



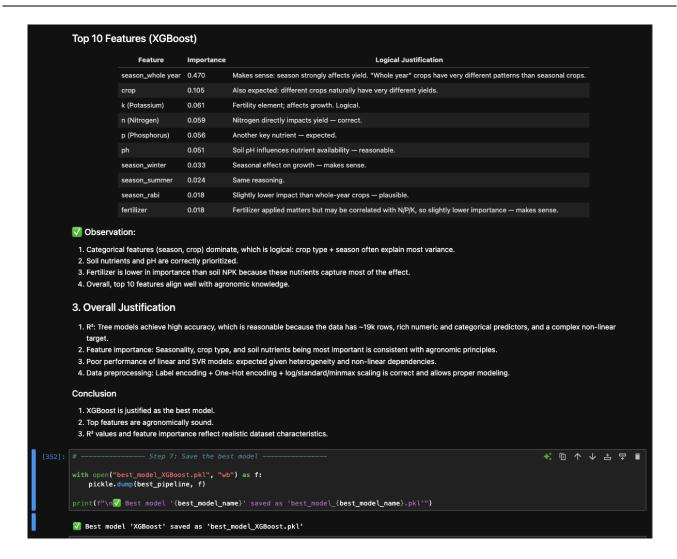
E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

(Cited: https://www.mdpi.com/2076-3417/13/16/9288)

```
[348]:
              # ----- Step 6: Top 10
if best_model_name == 'XGBoost':
                             __modet_mame == Xboost .
ktract preprocessed feature names
get_feature_names(preprocessor, X):
feature_names = []
                              Teature_names = []
for name, transformer, cols in preprocessor.transformers:
   if transformer == 'passthrough':
        feature_names.extend(cols)
   elif hasatr(transformer, 'get_feature_names_out'):
        feature_names.extend(transformer.get_feature_names_out(cols))
}
                             feature_names.extend(cols)
return feature_names
                      feature_names = get_feature_names(best_pipeline.named_steps['preprocessor'], X_train)
xgb_model = best_pipeline.named_steps['model']
importances = xgb_model.feature_importances_
                      min_len = min(len(feature_names), len(importances))
feat_imp = pd.DataFrame{{
    'Feature': feature_names|:min_len],
    'Importance': importances[:min_len]
                      })
top_10_features = feat_imp.sort_values(by='Importance', ascending=False).head(10)
                                                                                         --- 🔽 Top 10 Features (XGBoost): -
                      print(top_10_features)
print("\n")
                     # PLOT
plt.figure(figsize=(8,6))
sns.barplot(x='Importance', y='Feature', data=top_10_features, palette='viridis')
plt.title("Top 10 XGBoost Feature Importances")
plt.show()
                                                                   – V Top 10 Features (XGBoost): --
                                                                                            Top 10 XGBoost Feature Importances
                   season_whole year
                            season_winter
                        season_summer
                               season_rabi
                                      fertilizer
                                                                                                                        Importance
```



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com





E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Model Deployment

Model Deployment

After training and evaluating the machine learning models, the best-performing XGBoost pipeline was selected and deployed as an interactive application using Streamlit, a Python framework for building data apps and dashboards. Deployment makes the model accessible for end-users who can input real-world parameters and get predictions without needing to run Python scripts or understand the underlying code.

1. Loading the Model

The trained pipeline (best_model_XGBoost.pkl) is loaded using Python's pickle module. This ensures that all preprocessing steps (encoding, scaling, log transformation, year normalization) and the trained XGBoost model are included together for consistent predictions.

```
with open("best_model_XGBoost.pkl", "wb") as f:
    pickle.dump(best_pipeline, f)
```

2. User Interface with Streamlit

Streamlit provides an easy-to-use interface to create web apps directly from Python code. In this app:

- Dropdowns (st.selectbox) allow selection of crop, state, and season.
- Numeric input fields (st.number_input) let the user enter area, fertilizer, pesticide, weather, and soil nutrient values.
- Input values are combined into a single DataFrame matching the format used during training.

3. Real-Time Prediction

Once the user clicks the "Predict Crop Yield" button:

- The app feeds the input data through the pre-loaded pipeline.
- Preprocessing steps automatically transform the inputs (e.g., encoding categorical variables, scaling numeric values, log-transforming skewed features).
- The XGBoost model predicts the crop yield in q/ha, and the result is displayed dynamically in a styled output box.

(Cited: https://docs.streamlit.io/develop/api-reference)



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

4. Custom Styling

CSS styles are embedded in the Streamlit app for improved aesthetics:

- 1. Centered, bold title
- 2. Highlighted input and result boxes
- 3. Button hover effects

This enhances user experience, making the application more professional and visually appealing.

5. Benefits of Using Streamlit

(Cited: https://docs.streamlit.io/develop/concepts)

- Interactive Deployment: Users can quickly input parameters and get predictions in real-time.
- No Backend Setup: Streamlit handles the web server and rendering; no additional web development is required.
- Integration of ML Pipeline: Preprocessing and modeling are fully integrated, ensuring predictions are consistent with training.
- Shareable and Scalable: The app can be hosted on cloud platforms (Streamlit Cloud, Heroku, AWS) to make it accessible to farmers, agronomists, or researchers.

6. Practical Use Case

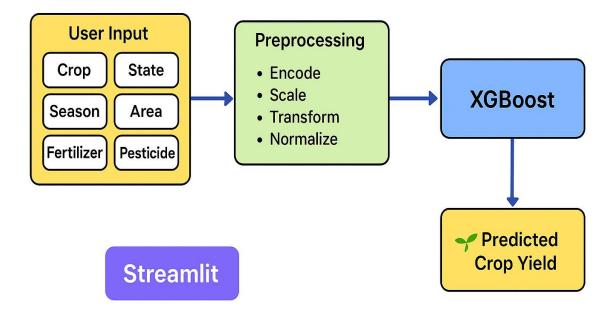
This deployment allows agricultural stakeholders to:

- 1. Estimate crop yield based on planned field practices and environmental conditions.
- 2. Explore "what-if" scenarios by changing fertilizer, pesticide, or weather parameters.
- 3. Make informed decisions on crop selection, resource allocation, and harvesting strategies.
- 4. In short, Streamlit provides a low-code, interactive way to deploy machine learning models for real-world use, transforming your XGBoost model from a script into a practical decision-making tool.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Machine Learning Model Deployment





E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

6. XGBoost_App.py

(Cited by : Available online: https://www.kaggle.com/code/theeyeschico/crop-analysis-and-prediction)



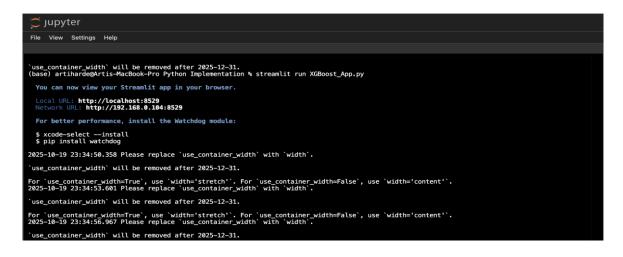
E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

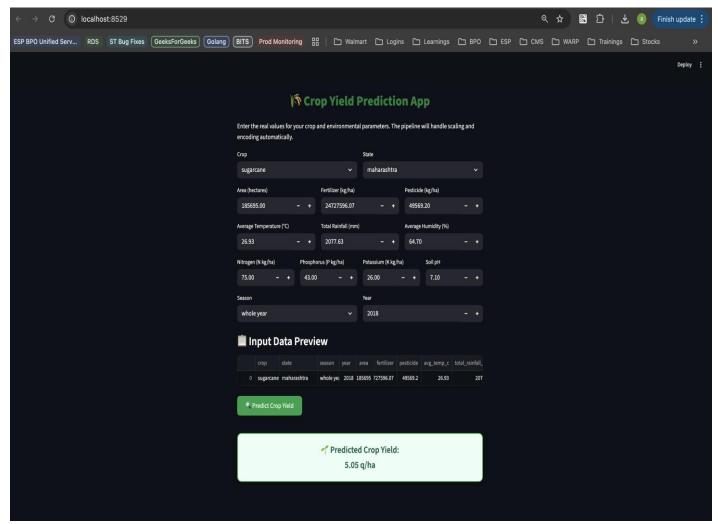
```
Jupyter XGBoost_App.py Last Checkpoint: 3 days ago
File Edit View Settings Help
Ħ
    col1, col2, col3 = st.columns(3)
88 with col1:
        avg_temp_c = st.number_input("Average Temperature (°C)")
    with col2:
        total_rainfall_mm = st.number_input("Total Rainfall (mm)")
92 with col3:
         avg_humidity_percent = st.number_input("Average Humidity (%)")
95 col1, col2, col3, col4 = st.columns(4)
96 with col1:
97     n = st.number_input("Nitrogen (N kg/ha)")
98 with col2:
       p = st.number_input("Phosphorus (P kg/ha)")
100 with col3:
k = st.number_input("Potassium (K kg/ha)")
with col4:
        ph = st.number_input("Soil pH")
105 col1, col2 = st.columns(2)
106 with col1:
        season = st.selectbox("Season", season_options)
108 with col2:
        year = st.number_input("Year", min_value=1997, max_value=2025, value=2020, step=1)
112 input_data = pd.DataFrame([{
        rut_data = pd.DataFrame({{
    'crop': crop_name,
    'state': state_name,
    'season': season,
    'year': year,
    'area': area,
    'fertilizer': fertilizer,
    'pesticide': pesticide,
    'avg_temp_c': avg_temp_c,
    'total_rainfall_mm': total_rainfall_mm,
    'avg_bumidity_execut': avg_bumidity_per
         'avg_humidity_percent': avg_humidity_percent,
         'n': n,
'p': p,
'k': k,
127 }])
130 st.subheader(" Input Data Preview")
131 st.dataframe(input_data, use_container_width=True)
134 if st.button(" Predict Crop Yield"):
               prediction = model_pipeline.predict(input_data)
st.markdown(f"""
                        Predicted Crop Yield: <br> <strong>{prediction[0]:.2f} q/ha</strong>
               """, unsafe_allow_html=True)
          except Exception as e:
               st.error(f"X Error during prediction: {e}")
```



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

7. Local Deployment of XGBoost_App.py Application Using Streamlit







E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

Project Files Structure

Project Github Link:

https://github.com/swaroopkumaraduru/Swaroop-Kumar-Aduru-Forecasting-Crop-Yield-in-India

1. Data Files

- crop_yield.csv: Crop-wise yield, area, production, fertilizer, pesticide, rainfall, season, and state (1997–2020).
- state_soil_data.csv: State-wise soil nutrients (N, P, K) and pH.
- state_weather_data_1997_2020.csv: State-wise annual weather (temperature, rainfall, humidity).
- final_merged_statewise_crop_yeild_dataset.csv: Merged dataset combining crop, soil, and weather data for modelling.
- final_merged_statewise_crop_yeild_dataset.xlsx: Excel version of the merged dataset.
- cleaned_crop_data.xlsx: Cleaned and preprocessed dataset after outlier handling and normalization.

2. Notebooks

- Forecasting State-Wise Crop Yield in India_Final.ipynb: Main analysis notebook. Covers:
 - o Data loading, merging, and cleaning
 - Exploratory Data Analysis (EDA)
 - o Feature engineering and normalization
 - o Model building (Linear Regression, Random Forest, XGBoost, etc.)
 - o Model evaluation and feature importance
 - o Dashboard generation (HTML)

3. Python Apps & Models

- XGBoost_App.py: Streamlit app for interactive crop yield prediction using the trained XGBoost model
- best_model_XGBoost.pkl: Serialized pipeline for XGBoost regression (used by the app).

4. Dashboards & Visualizations

- state_season_crop_dashboard.html: Interactive dashboard (state-wise seasonal yield and top crops).
- bivariate_analysis_dashboard.html: Dashboard for bivariate analysis (yield vs. rainfall, temperature, humidity, soil nutrients, crop diversity).
- Flow_Chart_Data_Merge.png: Visual flowchart showing how datasets are merged.

Key Insights of Project

□ Data Overview & Patterns

- Crop yield varies significantly across states and seasons.
- Certain crops dominate the dataset, showing higher frequency and production.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

☐ Feature Importance

- Features like area, fertilizer, avg_temp_c, and total_rainfall_mm strongly influence yield.
- Seasonal and state-level variations also play a significant role in prediction.

☐ Outlier & Data Quality Analysis

- Outliers exist in yield, area, and production, affecting model accuracy.
- Proper data cleaning and scaling improved model stability.

■ Model Performance

- Tree-based models like Random Forest and XGBoost performed better than linear models.
- Hyperparameter tuning and cross-validation enhanced predictive performance.

■ EDA & Visualization Insights

- Boxplots and distribution plots revealed seasonal trends and crop-wise yield differences.
- Correlation analysis highlighted strong relationships between fertilizer usage and yield.

☐ Shortcomings in Data & Model

- Sparse data for some crops or regions limits prediction accuracy.
- Lack of real-time or future weather features restricts forward-looking predictions.

☐ Actionable Predictions

- Model can help farmers estimate expected yield per crop and region.
- Insights can guide optimal fertilizer application and crop planning.

☐ Potential for Dashboard & Decision Support

- Predictions can be visualized in dashboards for state-wise, season-wise, and crop-specific insights.
- Can serve as a decision support tool for policymakers and farmers.

Future Work & Extension or Scope of improvements

☐ Crop-Specific Yield Predictions

- Provide separate yield predictions for each major crop to help farmers plan more effectively.
- Highlight high-risk crops in certain regions based on past performance.



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

□ Seasonal Forecasting

- Offer short-term predictions for upcoming seasons to guide sowing and harvesting schedules.
- Include alerts for potential low-yield seasons based on historical trends.

☐ Basic Decision Support Dashboards

- Develop lightweight dashboards using Streamlit or Power BI to visualize:
 - Yield trends over years
 - o Top-performing crops and regions
 - o Key factors affecting yield (fertilizer, rainfall, temperature)

☐ Predictive Insights for Resource Planning

- Estimate fertilizer or pesticide needs based on predicted yield and crop type.
- Suggest optimal cropping patterns for maximizing productivity on limited land.

☐ Simple Risk Assessment

- Highlight areas prone to low yield due to adverse weather conditions.
- Provide basic "what-if" analysis using existing data (e.g., impact of rainfall shortage on yield).

□ Data-Driven Recommendations

- Recommend suitable crops for specific states or regions based on historical trends.
- Provide guidelines for crop rotation and seasonal planning.

☐ Interactive Visualization Features

- Add filters for year, state, crop, or season to make dashboards more user-friendly.
- Include charts showing correlations between inputs (fertilizer, rainfall) and yield.

☐ Incremental Model Updates

- 1. Update predictions regularly as new yearly data becomes available.
- 2. Keep dashboards relevant without requiring heavy computational resources.

Bibliography

- 1. India Crop Production. https://www.kaggle.com/datasets/thedevastator/statewise-crop-production-in-india-a-statistical
- 2. Crop Production Trends https://www.kaggle.com/code/abhijitdahatonde/crop-production-trends



E-ISSN: 3048-7641 • Website: www.aijfr.com • Email: editor@aijfr.com

References

- 1. Overview of Regression Models and How to Determine the Best Model for Data : https://journaljsrr.com/index.php/JSRR/article/view/2452/5060
- 2. Technology and Management for Social Innovation (IATMSI), Gwalior, India,21–23 December 2022Ranjan, P.; Garg, R.; Rai, J.K. Artificial Intelligence Applications in Soil & Crop Management. In Proceedings of the IEEE Conference on Interdisciplinary Approaches in
- 3. Gehlot, A.; Sidana, N.; Jawale, D.; Jain, N.; Singh, B.P.; Singh, B. Technical analysis of crop production prediction using MachineLearning and Deep Learning Algorithms. In Proceedings of the International Conference on Innovative Computing, Intelligent.Communication and Smart Electrical Systems (ICSES), Chennai, India, 24–25 September 2022;
- 3. Vivek, S.; Ashish, K.T.; Himanshu, M. Technological revolutions in smart farming: Current trends, challenges & future directions. Compute. Electron. Agric. **2022**.
- 4. Mamatha, J.C.K. Machine learning based crop growth management in greenhouse environment using hydroponics farming techniques. Meas. Sens. 2023, 25, 100665.