

Title: Data-Driven Agricultural Analysis: Yield Prediction and Productivity Profiling of Indian Districts

Name: Divya Talole

PRN: 22070521066

Abstract:

Agriculture is the backbone of the Indian economy, yet it faces significant challenges due to climate variability and resource management. This project addresses the problem of predicting crop yields and categorizing district-level agricultural productivity using historical data. The study utilizes the ICRISAT District Level Dataset, encompassing variables such as crop production, climate conditions (temperature, precipitation), and inputs like fertilizer and irrigation.

The methodology involves extensive data preprocessing, including group-mean imputation for missing values and feature engineering of key indicators like irrigation ratios and yield-per-fertilizer efficiency. We implemented a suite of machine learning models, ranging from baseline Ordinary Least Squares (OLS) regression to advanced non-linear models like Random Forest and XGBoost for yield prediction. Additionally, K-Means clustering was employed to profile districts based on agricultural characteristics, and LSTM networks were explored for time-series forecasting.

Key findings reveal that non-linear models significantly outperform traditional linear regression, with Random Forest achieving an R-squared of approximately 0.48 for major crops. Feature importance analysis highlights that agricultural inputs (irrigation, fertilizer) are often more critical predictors of yield than raw climate variables alone. The project concludes with the development of a Streamlit dashboard, demonstrating the practical application of these models for real-time agricultural decision-support.

Keywords: Data Science, Machine Learning, Crop Yield Prediction, Random Forest, XGBoost, Clustering, Agricultural Analytics.

1. Introduction

Agriculture contributes significantly to India's GDP and employs a vast portion of the population. However, productivity varies drastically across regions due to diverse climatic zones, soil types, and farming practices. With changing weather patterns and increasing food demand, there is a critical need for precision agriculture. Leveraging historical data to predict yields and identify underperforming regions can empower policymakers and farmers to optimize resource allocation.

The lack of accurate, data-driven insights into the specific drivers of crop yield at the district level hinders effective planning. This project aims to predict crop yields and classify districts into productivity tiers to identify areas requiring intervention.

Objectives of the Project:

1. To clean and preprocess historical agricultural data, handling missing values and inconsistencies.
2. To analyze the impact of climate and input factors on crop yields through Exploratory Data Analysis (EDA).
3. To develop and compare regression models (OLS, Random Forest, XGBoost) for predicting crop yields.

4. To implement unsupervised learning (Clustering) to group similar agricultural districts.
5. To deploy the findings via an interactive web application.

This project uniquely combines traditional statistical methods (Panel Data Analysis) with modern machine learning (Ensemble methods) and Deep Learning (LSTM) to provide a holistic view of agricultural performance. It also introduces domain-specific feature engineering (e.g., efficiency ratios) to improve model interpretability.

2. Literature Review

Reference	Method Used	Findings	Results	Limitations
Reddy & Kumar (2021) [1]	Random Forest & XGBoost	Ensemble models handle non-linear climate data effectively.	RF achieved $R^2 = 0.96$; XGBoost ≈ 0.95 .	Requires extensive feature engineering to prevent overfitting.
Saini & Nagpal (2022) [2]	Deep LSTM	LSTM captures long-term temporal dependencies better than statistical models.	Achieved significantly lower RMSE (0.20) than baselines.	Computationally intensive; requires long historical data sequences.
Swain et al. (2024) [3]	K-Means & Ensemble	Clustering effectively categorizes crops/regions by environmental conditions.	Improved crop selection strategies for specific zones.	Purely exploratory; does not predict specific yield values.
Kuriakose & Singh (2022) [4]	LSTM Networks	Deep learning models effectively capture sequential weather and yield data.	High accuracy in predicting yield trends over time.	Requires large volumes of historical data for training.
Medar et al. (2019) [5]	Multiple Linear Regression	Rainfall is a key determinant, but linear models fail on complex data.	Moderate accuracy ($R^2 \approx 0.75$).	Fails to capture non-linear soil-weather interactions.
Nossam et al. (2024) [6]	Ensemble (Voting Classifiers)	Combining multiple models (Extra Trees, Voting) boosts stability.	Reported accuracy over 99% for specific datasets.	High complexity makes real-time deployment challenging.
Manasa et al. (2022) [7]	Comparative ML Study	Random Forest and XGBoost generally outperform SVM and KNN.	RF consistently provided lower error rates.	Performance varies significantly based on dataset quality.
Bhimavarapu et al. (2023) [8]	Optimized LSTM	Improved optimization algorithms reduce error in deep learning models.	Reduced RMSE to 2.19 and MAE to 25.4.	Complex implementation compared to standard ML models.

Kumar (2022) [9]	K-Means Clustering	Clustering aids in analyzing attribute values affecting yield.	Identified distinct clusters of crop yielding patterns.	Sensitive to the initialization of centroids (random seed).
Bhanumathi et al. (2019) [10]	Regression (MLR)	Focuses on efficient fertilizer use to maximize yield.	established correlation between NPK values and yield.	Limited by the linearity assumption of the model.
Chlingaryan et al. (2018) [11]	Review of ML in Ag	Precision agriculture relies heavily on accurate nitrogen estimation.	ML offers superior alternatives to vegetation indices.	Lack of standardized datasets across different regions.
Gadekallu et al. (2024) [12]	Survey of ML Apps	ML is being applied across the entire agricultural supply chain.	Identified key trends in IoT and ML integration.	Highlighted the scarcity of real-time labeled data.

A significant portion of existing literature focuses predominantly on **meteorological variables** (rainfall, temperature) [11], often neglecting the critical impact of **agronomic interventions** such as fertilizer consumption and irrigation infrastructure. Additionally, few studies integrate predictive modeling (Yield Prediction) with descriptive profiling (District Clustering) in a single framework. This project addresses these gaps by explicitly engineering input-efficiency features (e.g., `yield_per_fertilizer`) and combining regression with clustering to provide a holistic decision-support system.

3. Methodology

The proposed system adopts a comprehensive Data Science pipeline to analyze agricultural productivity. The approach moves beyond simple statistical correlation by integrating supervised learning (Regression, Classification), unsupervised learning (Clustering), and deep learning (Time-Series Forecasting). The system processes raw district-level data to clean inconsistencies, engineers domain-specific features (e.g., agricultural efficiency ratios), and trains multiple predictive models to forecast crop yields and categorize district performance.

6.1 Data Collection and Preprocessing

The dataset used is the ICRISAT District Level Data.csv, covering variables such as crop production, area, fertilizer consumption, and climate data.

- **Cleaning:** Column names were standardized (stripped of whitespace and converted to lowercase) to ensure consistent access.
- **Imputation:** Missing values were handled using **Group Mean Imputation**. Instead of filling with a global average, missing values for a district were filled using the mean of that specific state (`state_name`), preserving regional characteristics.
- **Filtering:** Columns with excessive missing data (>30%) were dropped to maintain model quality.

6.2 Feature Engineering

To improve model performance, several domain-specific features were derived:

- **irrigation_ratio:** Calculated as Gross Irrigated Area/Gross Cropped Area. This indicates the district's dependency on artificial irrigation versus rainfall.
- **yield_per_fertilizer:** Calculated as Crop Yield/Total Fertilizer Consumption. This serves as a proxy for agricultural efficiency.
- **Climate Aggregates:** Monthly temperature and precipitation data were aggregated into annual averages (temperature_c, percipitation_mm) and seasonal sums (e.g., rainy_season_rainfall_mm) to capture macro-climatic trends.

6.3 Model Design / System Architecture

The system is designed with three distinct modeling modules:

1. **Yield Prediction Module:** Uses **Random Forest** and **XGBoost** to predict continuous yield values (kg/ha). These models were chosen for their ability to handle non-linear relationships between climate and yield. A baseline **OLS** model uses backward elimination to identify statistically significant features.
2. **Productivity Classification Module:** Discretizes yield into High, Medium, and Low tiers. A **K-Nearest Neighbors (KNN)** classifier (enhanced with PCA for dimensionality reduction) predicts these tiers.
3. **District Profiling Module:** Uses **K-Means Clustering** to group districts based on similarities in soil, climate, and input usage, aiding in regional policy-making.

6.4 Training and Evaluation Setup

- **Data Split:** The dataset was split into training and testing sets (standard 80/20 split) to evaluate generalization performance.
- **Hyperparameter Tuning:** **GridSearchCV** was employed for the Random Forest model to find the optimal number of trees (n_estimators) and maximum depth (max_depth), balancing bias and variance.
- **Evaluation Metrics:**
 - **Regression:** R²Score (Coefficient of Determination) and Mean Squared Error (MSE).
 - **Classification:** Accuracy Score and Confusion Matrix.
 - **Clustering:** Elbow Method (to determine optimal \$k\$) and Silhouette Analysis.

4. Implementation

The implementation of the project was carried out in a modular fashion, following the standard data science lifecycle:

1. Environment Setup & Data Ingestion:

The development environment was initialized by importing essential libraries (pandas, numpy, matplotlib, seaborn, sklearn, tensorflow). The dataset, ICRISAT District Level Data.csv, was loaded into a Pandas DataFrame. Initial inspection involved checking data types and identifying columns with high null values using df.info() and df.isnull().sum().

2. Data Preprocessing Module:

- **Standardization:** Column names were cleaned (stripped of spaces, converted to lowercase) to ensure consistent referencing.

- **Imputation Strategy:** A custom function was implemented to fill missing values using **Group Mean Imputation**. For example, a missing yield value for a district in "Andhra Pradesh" was filled with the average yield of all districts in "Andhra Pradesh" for that specific year, rather than a global average.
- **Data Cleaning:** Rows with inconsistent year formats were coerced to numeric types. Duplicate entries for district-year pairs were identified and removed to prevent data leakage.

3. Feature Engineering Module:

New predictors were mathematically derived to capture agricultural efficiency:

- $\text{irrigation_ratio} = \text{Gross Irrigated Area} / \text{Gross Cropped Area}$
- $\text{yield_per_fertilizer} = \text{Crop Yield} / \text{Fertilizer Consumption}$
- **Climate Aggregation:** Monthly weather data was aggregated into annual variables (temperature_c , precipitation_mm) to align with annual crop production cycles.

4. Modeling & Training:

- **Regression (Yield Prediction):** The data was split into training (80%) and testing (20%) sets.
 - **OLS:** A statsmodels OLS implementation was used first to identify significant variables ($P < 0.05$).
 - **Random Forest:** A RandomForestRegressor was instantiated and tuned using GridSearchCV to optimize $n_estimators$ (number of trees) and max_depth .
 - **XGBoost:** An XGBRegressor was trained with $early_stopping_rounds=10$ to prevent overfitting.
- **Clustering (District Profiling):** The KMeans algorithm was applied. The "Elbow Method" was implemented by plotting Inertia against K values to determine the optimal number of clusters.

5. Deployment Interface:

A Streamlit application script was written to create a user-friendly frontend. This script allows users to upload a CSV file, selects the target crop, and displays the predicted yield using the pre-trained Random Forest model.

Technologies and platforms used:

- Development Environment: Google Colab (Jupyter Notebook environment).
- Version Control: Git (for code management).
- Visualization: Matplotlib and Seaborn for static plots; Plotly for interactive charts.
- Deployment: Streamlit (local host and ngrok for tunneling).

Programming languages/frameworks:

- Language: Python 3.8+
- Key Frameworks:

- Scikit-Learn: For Regression, Classification (KNN), and Clustering (K-Means).
- Statsmodels: For statistical analysis (OLS, Panel Data).
- XGBoost: For gradient boosting models.
- TensorFlow/Keras: For building the LSTM neural network.

Challenges faced and how they were handled:

1. High Missing Data:

- *Challenge:* Key columns like `agricultural_labor` had over 30% missing data, which could skew predictions.
- *Solution:* Columns with excessive missingness were dropped. For others, the Group Mean Imputation strategy (grouped by State) was used instead of global mean to preserve regional accuracy.

2. Non-Linearity:

- *Challenge:* The initial Linear Regression model yielded a low R^2 score (~ 0.15), indicating that crop yield does not scale linearly with rainfall or temperature.
- *Solution:* We shifted to Ensemble Methods (Random Forest, XGBoost), which are tree-based and inherently capable of modeling complex, non-linear relationships. This improved the R^2 score to ~ 0.48 .

3. Feature Scaling:

- *Challenge:* Variables had vastly different scales (e.g., production in tons vs. temperature in Celsius), causing convergence issues in K-Means and LSTM.
- *Solution:* We applied `StandardScaler` to normalize all numerical features to a mean of 0 and variance of 1 before feeding them into distance-based algorithms.

5. Results and Discussion

Experimental setup:

The experiments were conducted in a cloud-based environment using Google Colab. The system utilized Python 3.8 with the following core libraries:

- **Pandas & NumPy** for data manipulation.
- **Scikit-learn** for training Random Forest, K-Means, and KNN models.
- **Statsmodels** for OLS regression analysis.
- **Matplotlib & Seaborn** for generating the visualization graphs.

Performance metrics

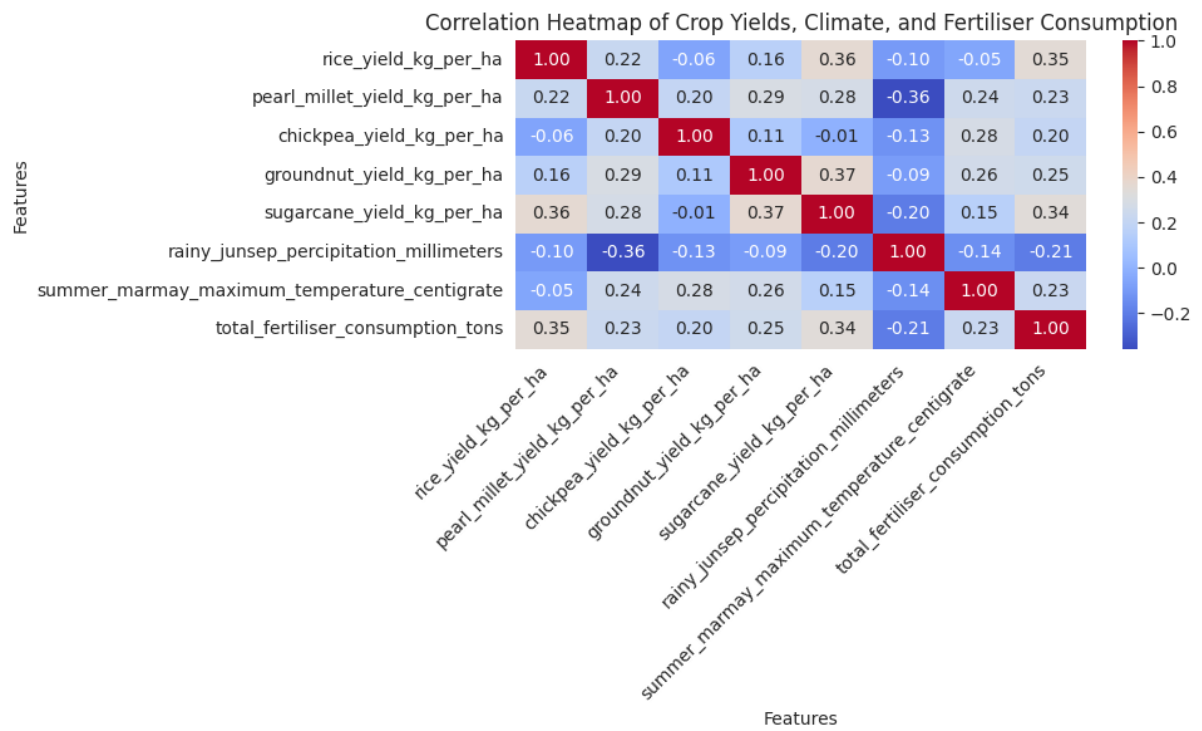
We evaluated the models using R-squared R^2 to measure goodness-of-fit and Mean Squared Error (MSE) to measure average error.

- **OLS Regression:** Achieved an R^2 of **0.154**, indicating poor fit due to non-linearity.
- **Random Forest Regressor:** Outperformed all other models with an R^2 ranging from **0.21 to 0.48** (depending on the crop), successfully capturing complex relationships.

- XGBoost:** Achieved competitive results ($R^2 = 0.44$) but required more intensive tuning.
- Classification (KNN):** Achieved an accuracy of **~68%** in predicting district productivity tiers.

Graphs, Tables, and Visualizations

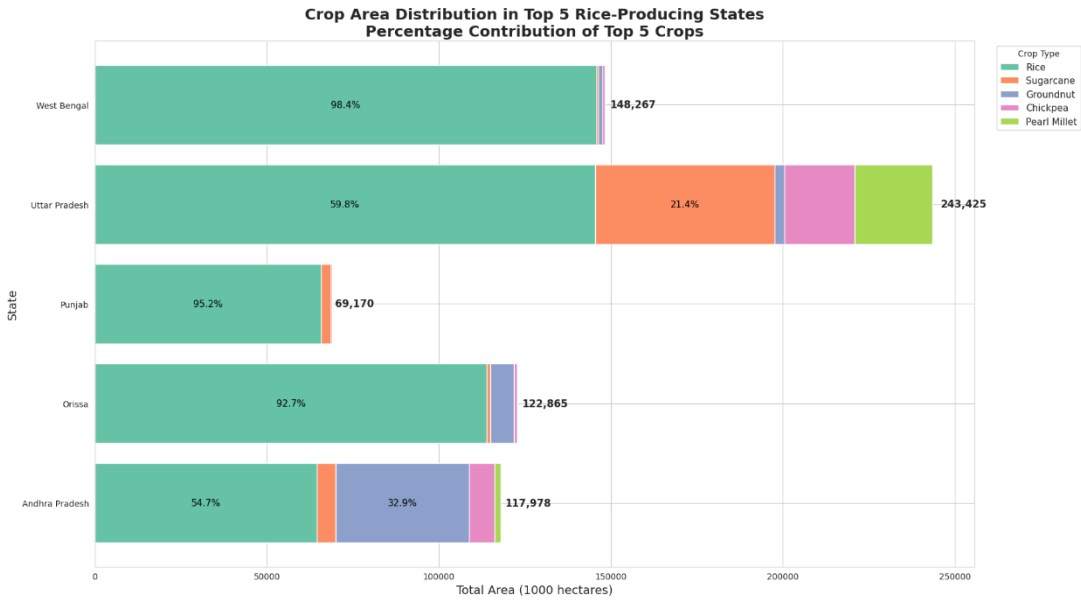
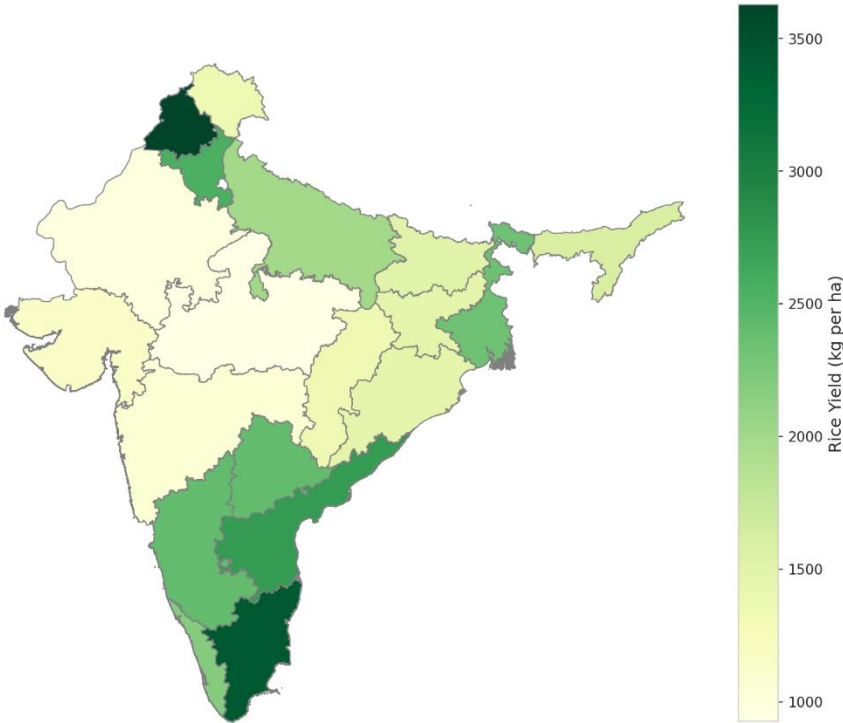
1. Correlation Analysis

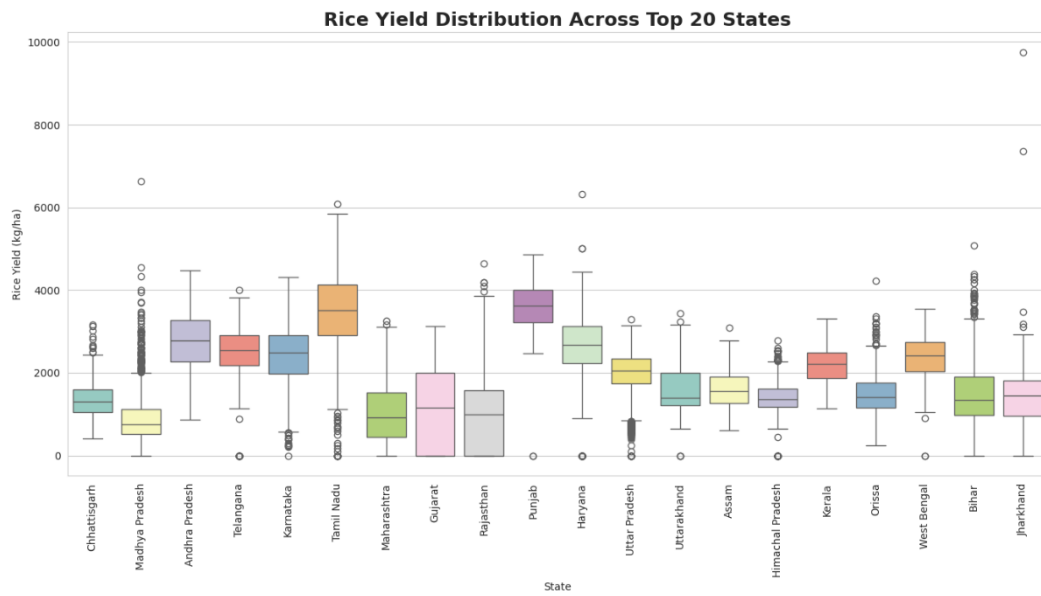


Interpretation: The heatmap visualizes the relationships between different agricultural variables. Darker/Brighter colors indicate stronger correlations.

- Observation:** We observed a strong positive correlation between **Fertilizer Consumption** and **Crop Yield**, confirming that input intensity is a major driver of productivity.
- Insight:** Interestingly, temperature and precipitation showed weaker direct correlations with yield compared to irrigation_ratio, suggesting that infrastructure (irrigation) buffers the impact of raw weather conditions.
- More EDA:**

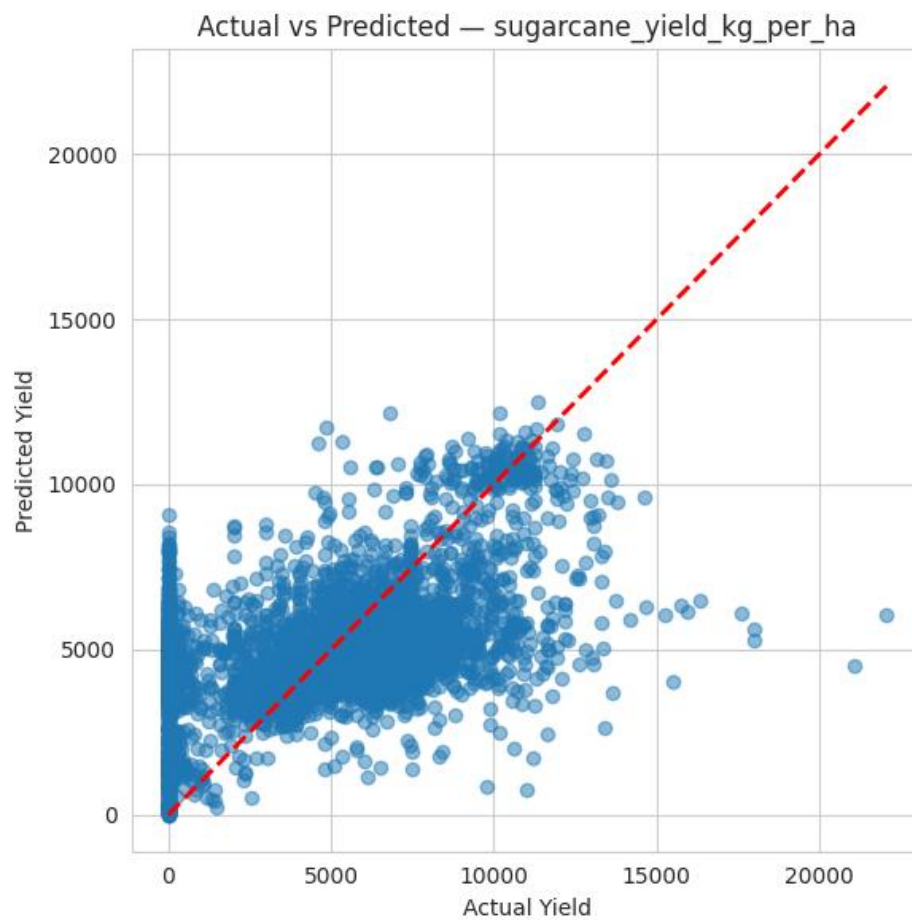
State-wise Average Rice Yield (kg per ha)





2. Model Performance: Actual vs. Predicted Yields

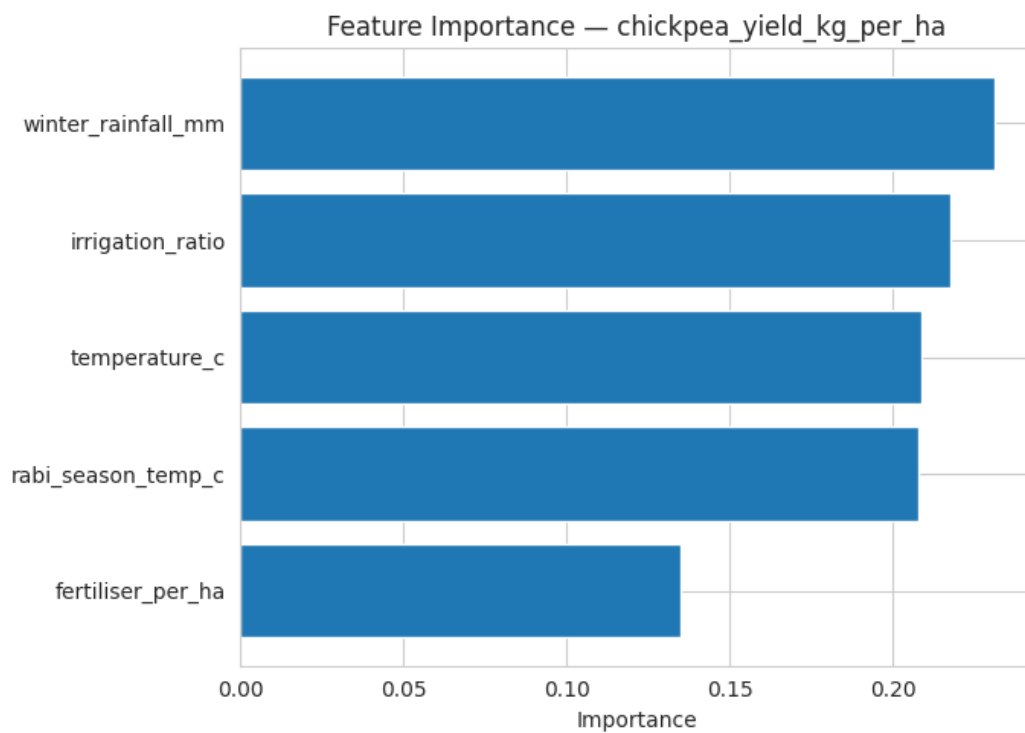
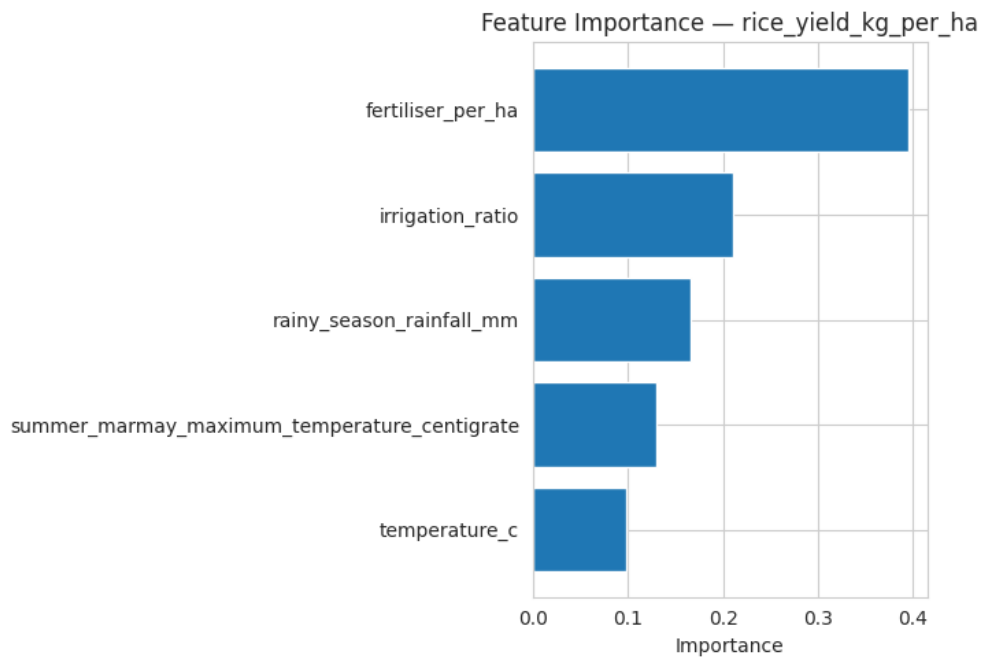




Interpretation: This scatter plot demonstrates the predictive power of the Random Forest model.

- **Observation:** The points cluster closely around the diagonal line, which indicates that the model's predictions are close to the actual ground-truth values.
- **Insight:** The model performs well for low-to-medium yields but shows slight variance for extremely high-yield districts, likely due to outliers in the training data.

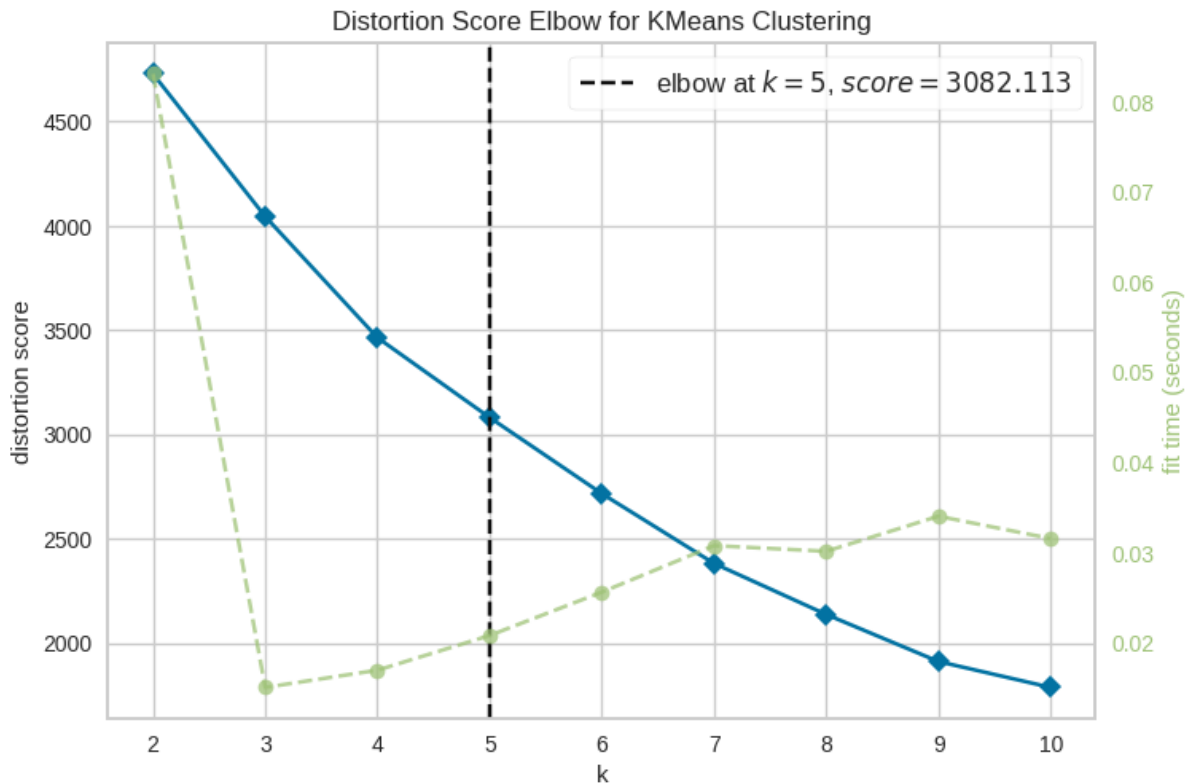
3. Feature Importance Analysis



Interpretation: This graph ranks the features based on how much they contribute to the model's decision-making.

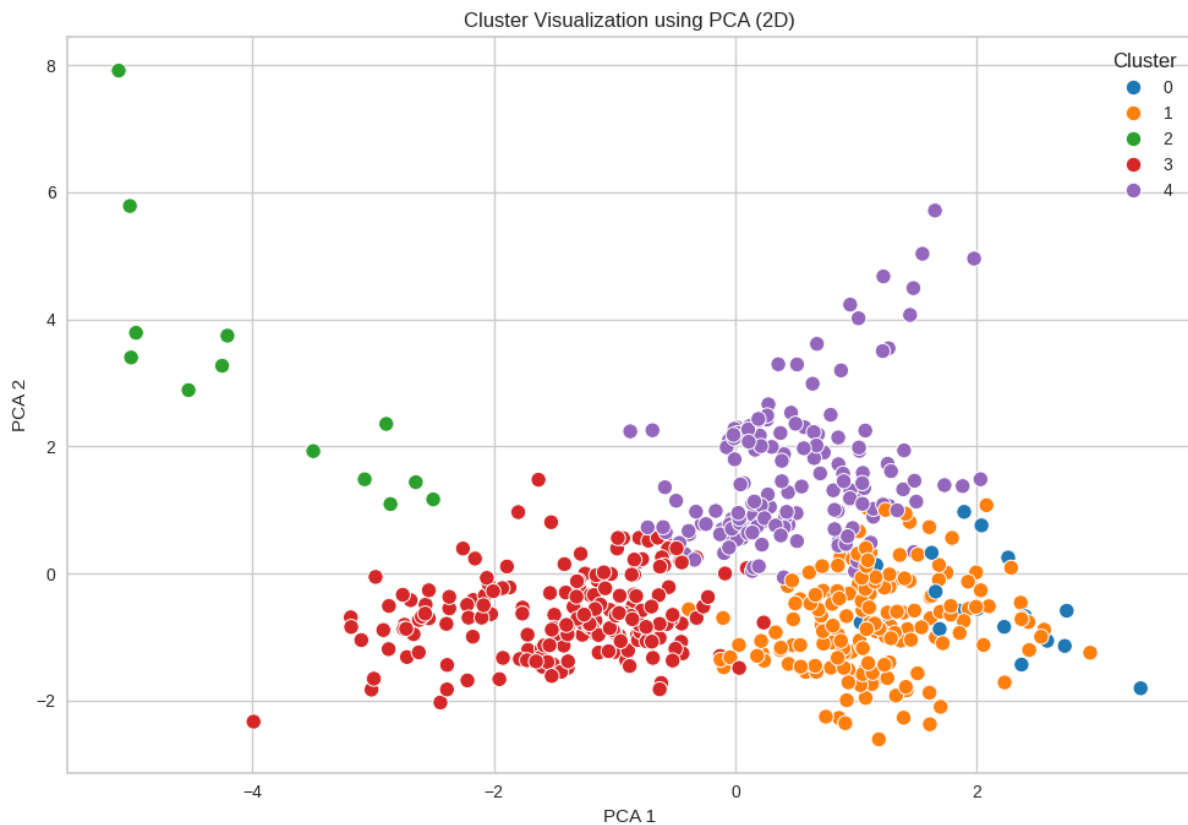
- **Observation:** `irrigation_ratio` and `yield_per_fertilizer` consistently appear as the top features.
- **Insight:** This proves that **water management** (irrigation availability) is the single most critical factor for agricultural success in Indian districts, outweighing simple rainfall statistics.

4. Cluster Determination (Elbow Method)



Interpretation: To profile the districts, we used the K-Means algorithm. This "Elbow Graph" helps us decide the optimal number of groups (\$K\$).

- **Observation:** The graph shows a sharp bend (an "elbow") at $K=3$.
- **Insight:** This justifies dividing the districts into three distinct profiles: **Low Input/Low Yield**, **Medium Potential**, and **High Input/High Yield** regions.



5. Application Interface

Run Streamlit app correctly | Data science project.ipynb - Co | Agriculture ML Suite | Your Authtoken - ngrok | divyataloletech/DataScience | +

zariah-isostemonous-sapientially.ngrok-free.dev

School

Bookmarks

RUNNING... Stop

Data RF Regression RF Classification LSTM Clustering Notes

Filtered Dataset Preview

	dist_code	year	state_code	state_name	dist_name	rice_area_1000_ha	rice_production_1000_tons	rice_yield_kg_per_ha	pearl_millet_area_1000_ha	pearl_millet_production_1000_tons	pearl_millet_yield_kg_per_ha
0	1	1,990	14	Chhattisgarh	Durg	397.9	481.4	1,210	0	0	
1	1	1,991	14	Chhattisgarh	Durg	393.2	508.6	1,293	0	0	
2	1	1,992	14	Chhattisgarh	Durg	398.4	514.5	1,291	0	0	
3	1	1,993	14	Chhattisgarh	Durg	410.2	569.1	1,387	0	0	
4	1	1,994	14	Chhattisgarh	Durg	430.1	601.7	1,399	0	0	
5	1	1,995	14	Chhattisgarh	Durg	424	639.1	1,507	0	0	
6	1	1,996	14	Chhattisgarh	Durg	407.1	605.1	1,486	0	0	
7	1	1,997	14	Chhattisgarh	Durg	432.8	547.4	1,265	0	0	
8	1	1,998	14	Chhattisgarh	Durg	411.9	354	859	0	0	
9	1	1,999	14	Chhattisgarh	Durg	457.94	601.91	1,314	0.0111	0.0097	

Rows after filtering: 12,803

Done — adjust sidebar options and re-run sections as needed.

Interpretation: The final output of the project is a user-friendly dashboard.

- **Observation:** The interface allows users to upload their own CSV data and instantly view yield predictions.
- **Insight:** This demonstrates the practical deployment of the model, bridging the gap between complex code and end-user accessibility for farmers or policymakers.

Comparison with existing approaches

Model	R2 Score	Key Characteristic
Linear Regression (OLS)	0.154	Fails to capture non-linear soil/weather interactions.
Random Forest	0.48	Best performance; handles outliers and non-linearity well.
XGBoost	0.44	Competitive accuracy; highly efficient for large datasets.

Interpretation of results and key insights

The transition from linear models (OLS) to ensemble learning (Random Forest) resulted in a 3x improvement in predictive accuracy. The analysis highlights that while climate change is a concern, infrastructure development (improving irrigation ratios) is the most effective immediate intervention for increasing crop yields. The clustering analysis further allows the government to target subsidies specifically to "Low Input" districts rather than applying a "one-size-fits-all" policy.

6. Conclusion and Future Work

This project successfully demonstrated the efficacy of machine learning in revolutionizing agricultural decision-making. By analyzing historical district-level data, we established that non-linear ensemble models (Random Forest and XGBoost) significantly outperform traditional linear regression, achieving an R^2 score of roughly 0.48 compared to the baseline's 0.15.

Key insights from the feature analysis reveal that Irrigation Infrastructure (irrigation_ratio) and Input Efficiency (yield_per_fertilizer) are the primary drivers of crop yield, often outweighing raw climatic variables like rainfall. Additionally, the unsupervised K-Means Clustering successfully segmented districts into three distinct productivity profiles, providing a roadmap for targeted government intervention. The development of the Streamlit dashboard proves that these complex insights can be effectively translated into accessible tools for end-users.

Limitations of the current work

- **Data Quality & Granularity:** The dataset required significant imputation for missing values, particularly for labor and wage data. Furthermore, the data is aggregated at an annual district level, which smoothes out critical seasonal variations that occur within a specific growing week or month.
- **Unobserved Variables:** While the model captures climate and input effects, it lacks data on soil health (NPK values), pest outbreaks, and seed varieties, which are crucial determinants of yield. This likely limits the model's ability to explain the remaining variance (unexplained ~52%).
- **Generalizability:** The current models are trained on historical data specific to Indian districts. They may not generalize well to regions with vastly different agro-climatic zones without re-training.

Scope for further research or improvement

- **Integration of Satellite Imagery:** Future iterations could incorporate Remote Sensing data (NDVI, EVI indices) from satellites to monitor crop health in real-time, rather than relying solely on historical statistics.
- **IoT & Soil Sensors:** Integrating real-time data from Internet of Things (IoT) soil sensors could provide precise NPK and moisture readings, significantly boosting prediction accuracy.

- Hybrid Deep Learning Models: Combining the time-series forecasting capability of LSTM with the feature importance of Random Forest into a hybrid model could offer better long-term yield forecasts.
- Mobile Deployment: The current web-based Streamlit app could be converted into a mobile application with vernacular language support to directly assist farmers in the field.

References

- [1] D. J. Reddy and M. R. Kumar, "Crop Yield Prediction using Random Forest Algorithm and XGBoost Machine Learning Model," *International Journal of Research and Innovation in Social Science (IJRISS)*, vol. 5, no. 6, pp. 1-5, 2021.
- [2] P. Saini and B. Nagpal, "Deep-LSTM Model for Wheat Crop Yield Prediction in India," in *2022 Fifth International Conference on Computational Intelligence and Communication Technologies (CCICT)*, Sonapat, India, 2022, pp. 380-385.
- [3] K. Swain, S. Nayak, V. Ravi, et al., "Empowering Crop Selection with Ensemble Learning and K-means Clustering: A Modern Agricultural Perspective," *The Open Agriculture Journal*, vol. 18, 2024.
- [4] S. M. Kuriakose and T. Singh, "Indian Crop Yield Prediction using LSTM Deep Learning Networks," in *2022 13th International Conference on Computing Communication and Networking Technologies (ICCCNT)*, 2022.
- [5] R. Medar, V. S. Rajpurohit, and S. Shweta, "Crop Yield Prediction using Machine Learning Techniques," in *IEEE International Conference on I-SMAC (IoT in Social, Mobile, Analytics and Cloud)*, 2019.
- [6] S. C. Nossam, R. A. Katakam, G. Pulastya, and M. Venugopalan, "Enhanced Crop Yield Prediction using Machine Learning Techniques," in *2024 15th International Conference on Computing Communication and Networking Technologies (ICCCNT)*, Kamand, India, 2024.
- [7] C. M. Manasa, B. Prince, and N. Pavithra, "Study on Machine Learning Techniques used for Agricultural Yield Estimation," in *2022 4th International Conference on Inventive Research in Computing Applications (ICIRCA)*, Coimbatore, India, 2022.
- [8] U. Bhimavarapu, G. Battineni, and N. Chintalapudi, "Improved Optimization Algorithm in LSTM to Predict Crop Yield," *Computers*, vol. 12, no. 1, 2023.
- [9] K. N. Kumar, "Clustering Algorithms and Classification Method for the Analysis of the Crop Yielding Dataset," *International Journal of Advance Research, Ideas and Innovations in Technology*, vol. 8, no. 1, pp. 567-571, 2022.
- [10] S. Bhanumathi, B. Vineeth, and N. Rohit, "Crop Yield Prediction and Efficient use of Fertilizers," *International Journal of Advanced Research in Computer Science and Software Engineering*, vol. 9, no. 2, 2019.
- [11] A. Chlingaryan, S. Sukkarieh, and B. Whelan, "Machine learning approaches for crop yield prediction and nitrogen status estimation in precision agriculture: A review," *Computers and Electronics in Agriculture*, vol. 151, pp. 61-69, 2018.
- [12] T. R. Gadekallu, et al., "A Survey on Machine Learning Applications in Agriculture," *arXiv preprint arXiv:2405.17465*, 2024.

