



**India's BIGGEST**  
**Analytics Case Competition!**  
**LTF Challenge- Farmer Income Prediction**

Team Name : marisettiharshini2005

Team members: M.Harshini & D.Nikitha

# Problem Statement

## Objective:

- Predict the Total Income of farmers using demographic, financial, and resource-related features.

## Context & Importance:

- Many farmers lack predictable income, which affects subsidy planning, creditworthiness evaluation, and agricultural policies.
- Accurately predicting income can:
  - Identify at-risk farmers
  - Improve financial inclusion
  - Support targeted interventions from the government and NGOs

## Challenge:

- The dataset includes diverse and noisy real-world features.
- Balancing accuracy, interpretability, and scalability is crucial for deploying such models in practice.

## Impact

- Predicts farmer income for smarter subsidies and financial aid.
- Helps identify vulnerable farmers for early support.
- Supports data-driven rural policy and economic planning.

## Feasibility

- Uses efficient LightGBM model with clean preprocessing.
- Fully automated pipeline from data to predictions.
- Easily deployable with minimal adjustments.

## Innovation

- SHAP and MAPE-based insights improve model transparency.
- Auto-generated diagnostics (residuals, feature impact).
- Advanced explainability uncommon in student-level projects.

## Scalability

- Adaptable across regions and crops.
- Can integrate external data (e.g., weather, satellite).
- Suitable for APIs, dashboards, or mobile deployment.

# Our Approach



## Data Loading & Understanding

- Loaded training and test data from Excel sheets
- Explored feature types, distribution, and target variable



## Preprocessing

- Cleaned column names, handled missing values
- Encoded categorical features using `pd.factorize()`



## Exploratory Data Analysis (EDA)

- Visualized income distribution, correlations, and feature trends
- Identified top influencing features



## Model Selection & Training

- Chose LightGBM for performance and interpretability
- Tuned hyperparameters:  
`n_estimators=1000, learning_rate=0.05,`



## Model Evaluation

- Used MAPEto assess performance
- Visualized residuals, prediction errors, actual vs. predicted



## Prediction & Submission

- Predicted income for test data
- Formatted output with FarmerID and Predicted\_Income

# Model Selection & Justification

## ✓ Why LightGBM?

- Chosen for its efficiency, high performance on tabular data, and support for categorical variables
- Easily handles large datasets and missing values without complex precesssing
- Outperforms many models in terms of speed, accuracy-and interpretability – making it ideal for real-world deployment

## ⚙️ Advantages Over Other Models

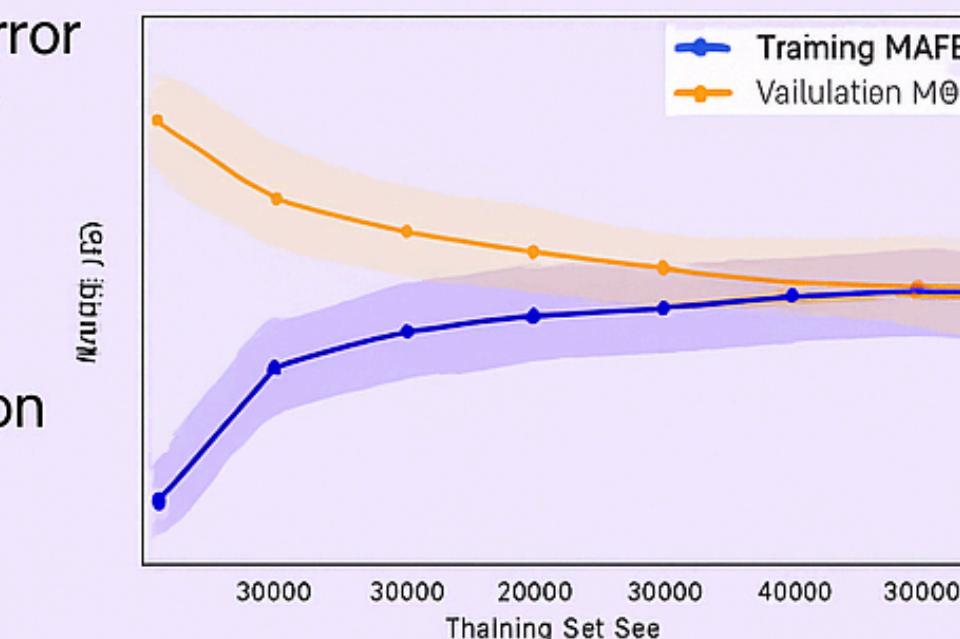
Compared To	Advantage of LightGBM
XGBoost	Faster training with leaf-wise growth
Random Forest	Better tuning flexibility predictive accuracy
Linear Models	Captures non-linear interaction between features
Neural Nets	Requires less data around 2%, proving model reliability

## 🔧 Model Configuration

- Model: LGBM Regressor
- n\_estimators = 1000
- learning\_rate = 0.05
- num\_leaves = 31
- random\_state = 42

## 📈 Learning Curve Analysis (MAPE)

- Plotted Training vs, Validation MAPE over increasing training set sizes
- Validation error consistently decreases; showing improved generalization
- Final MAPE stabilizes around 24% proving reliability



# Evaluation Metrics



## Evaluation Metrics Used

- MAPE (Mean Absolute Percentage Error)
- MAE (Mean Absolute Error)



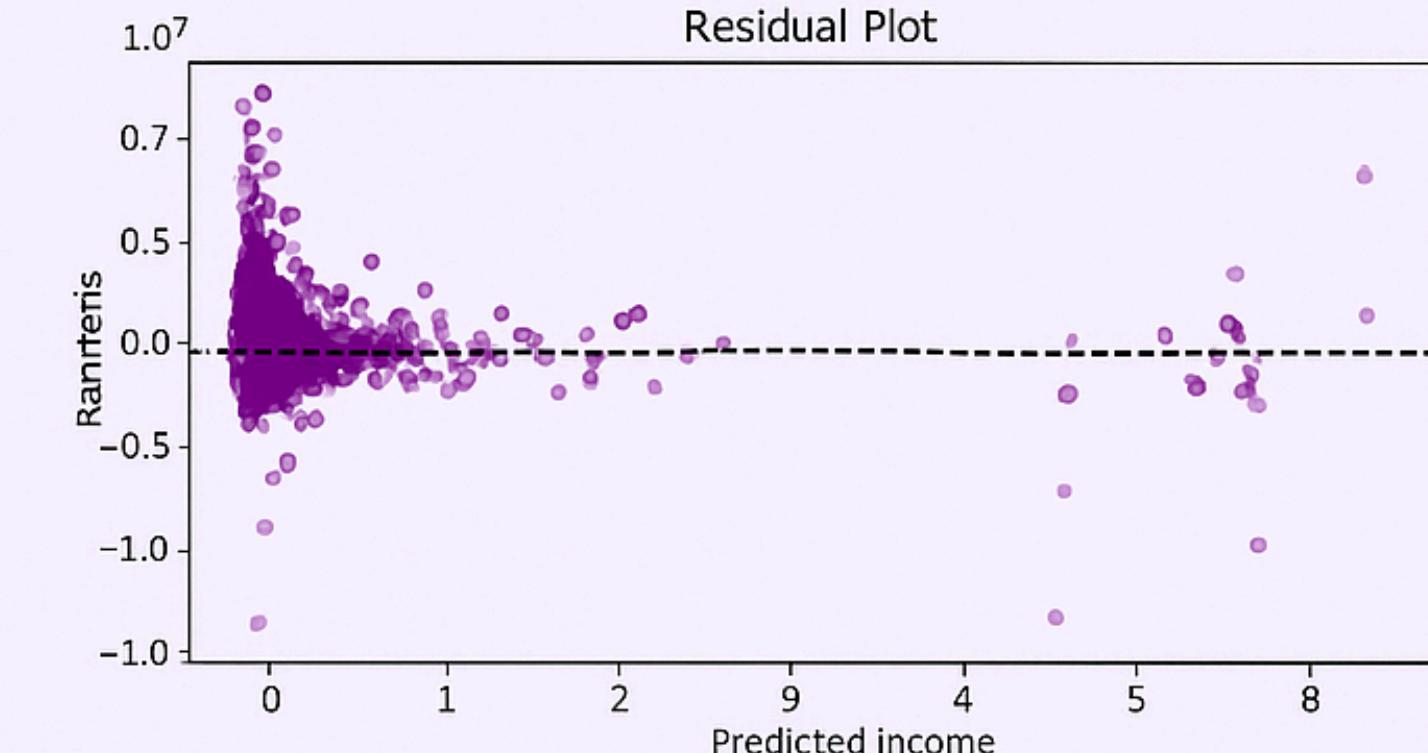
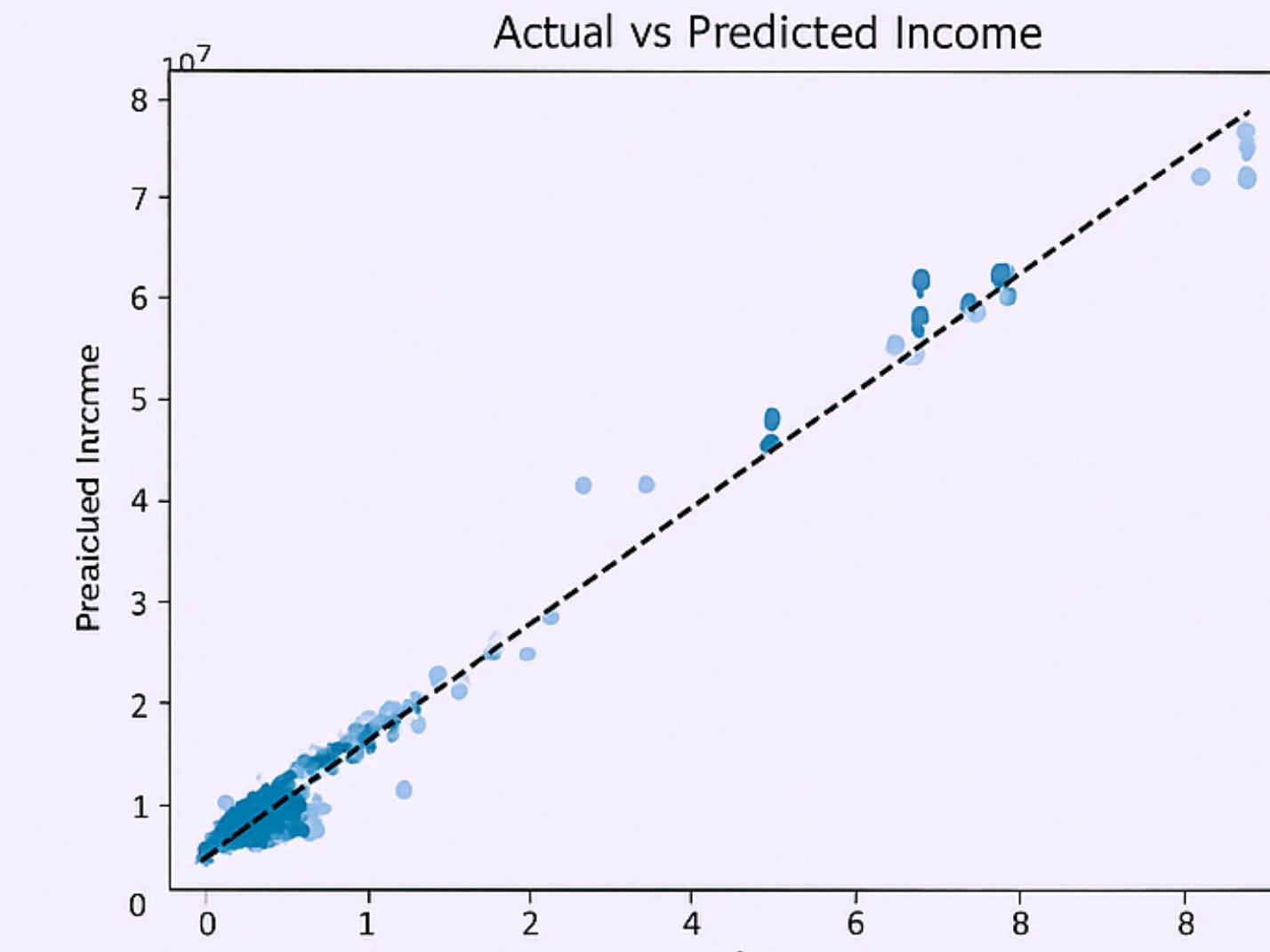
## Why MAPE?

- Training MAPE ~23.5%
- Validation MAPE ~24% (from learning curve)



## Visual Analysis Included

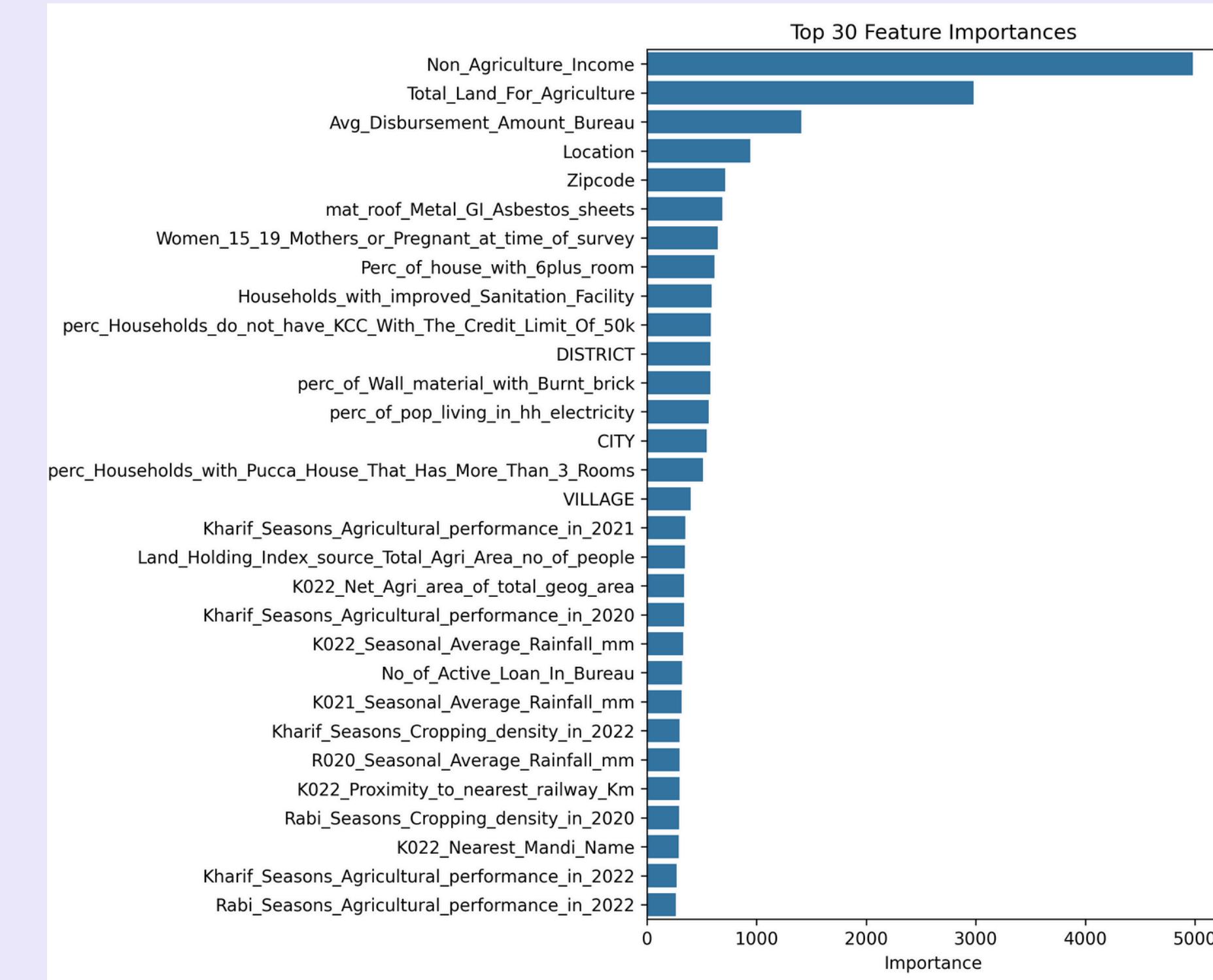
- Prediction Error Distribution
- Actual vs. Predicted Plot
- Residuals vs. Predicted Values

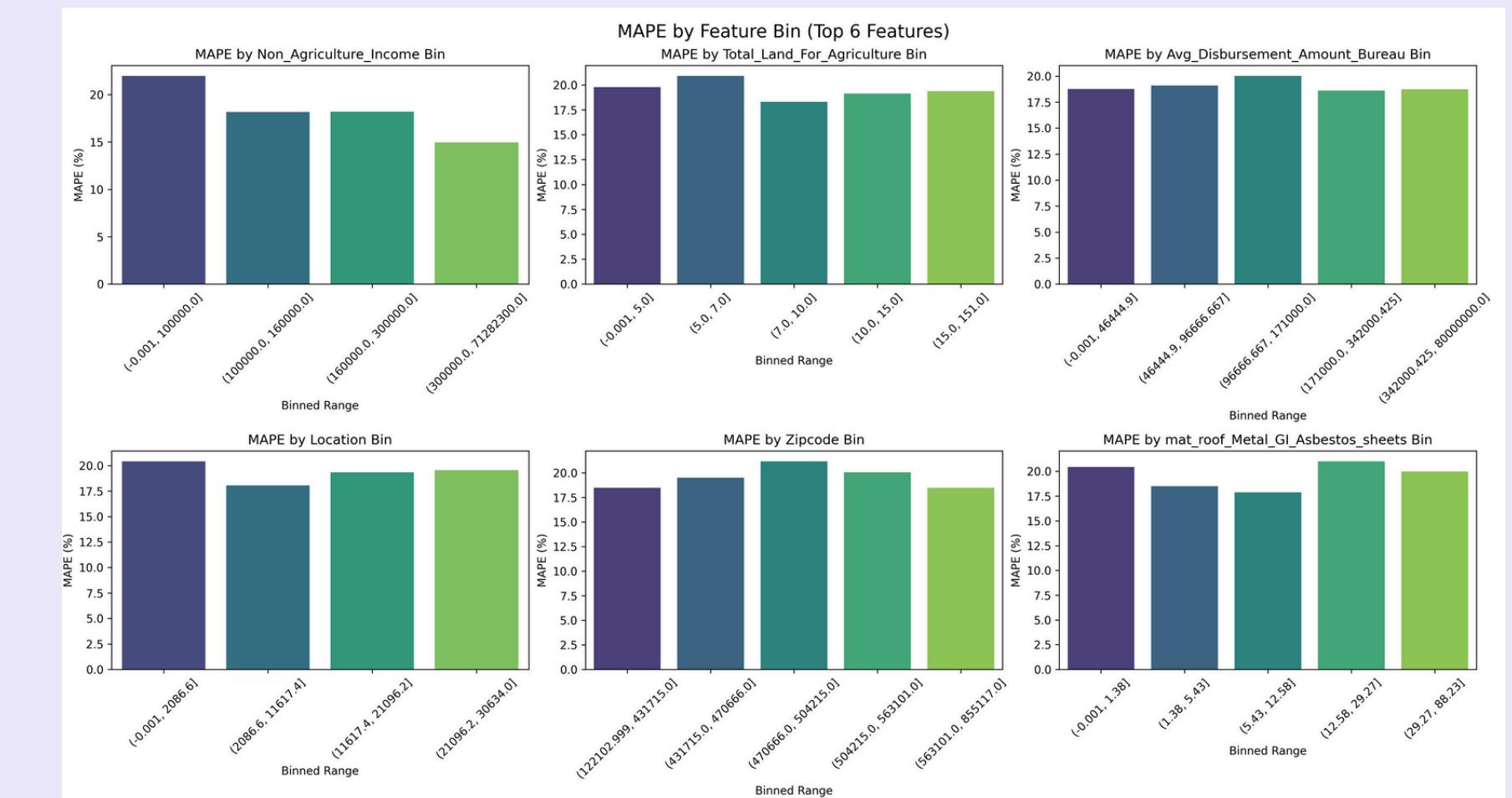
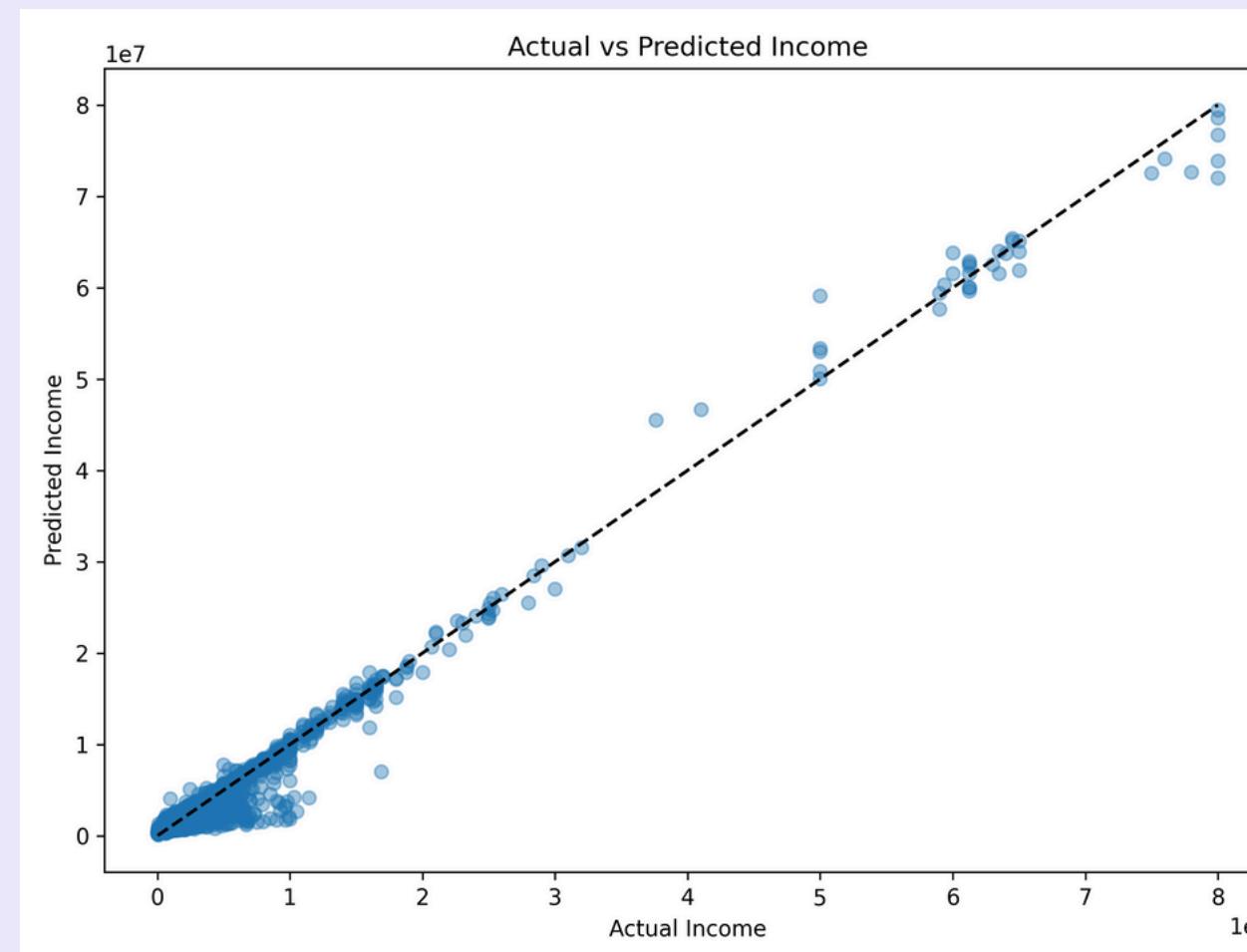
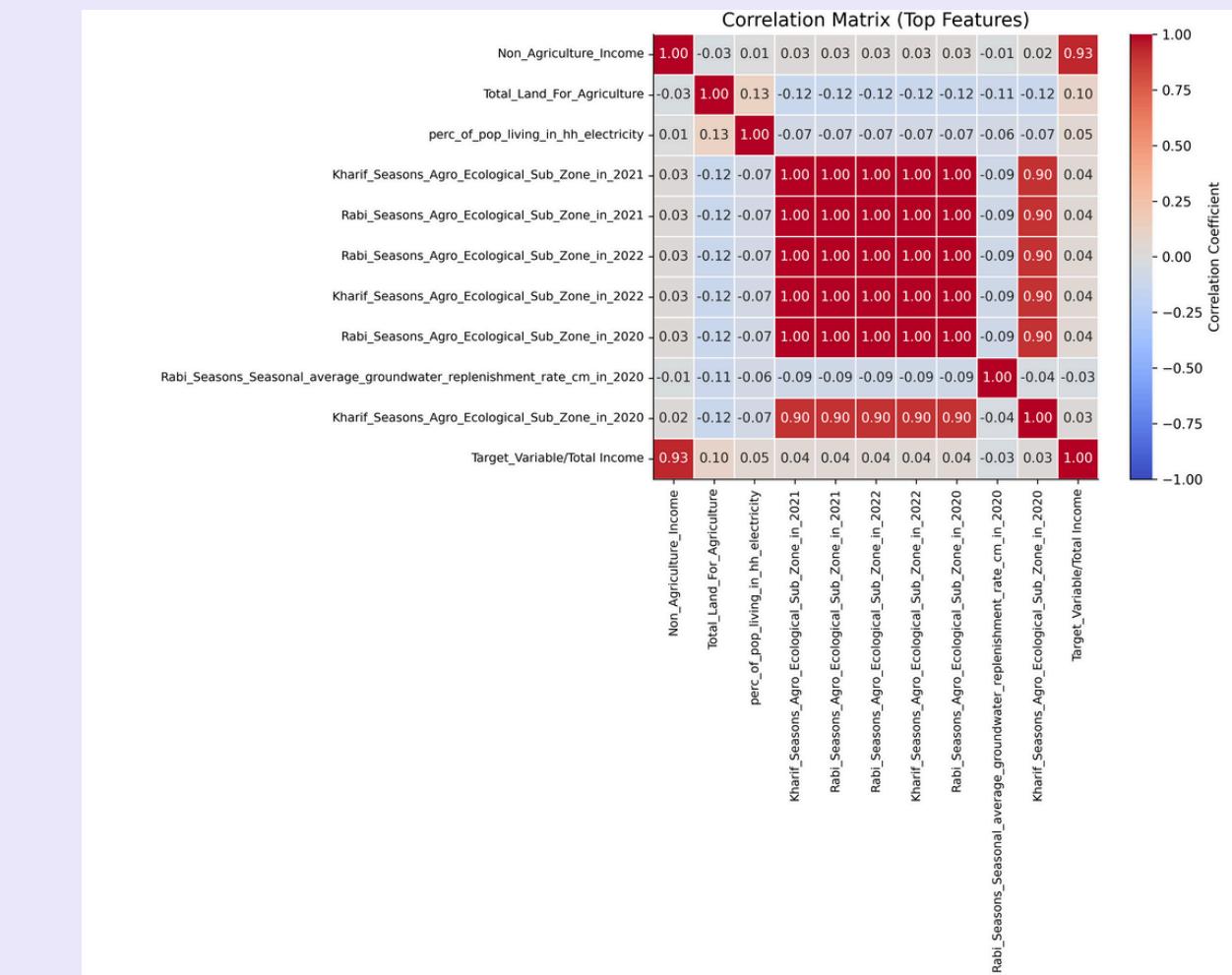
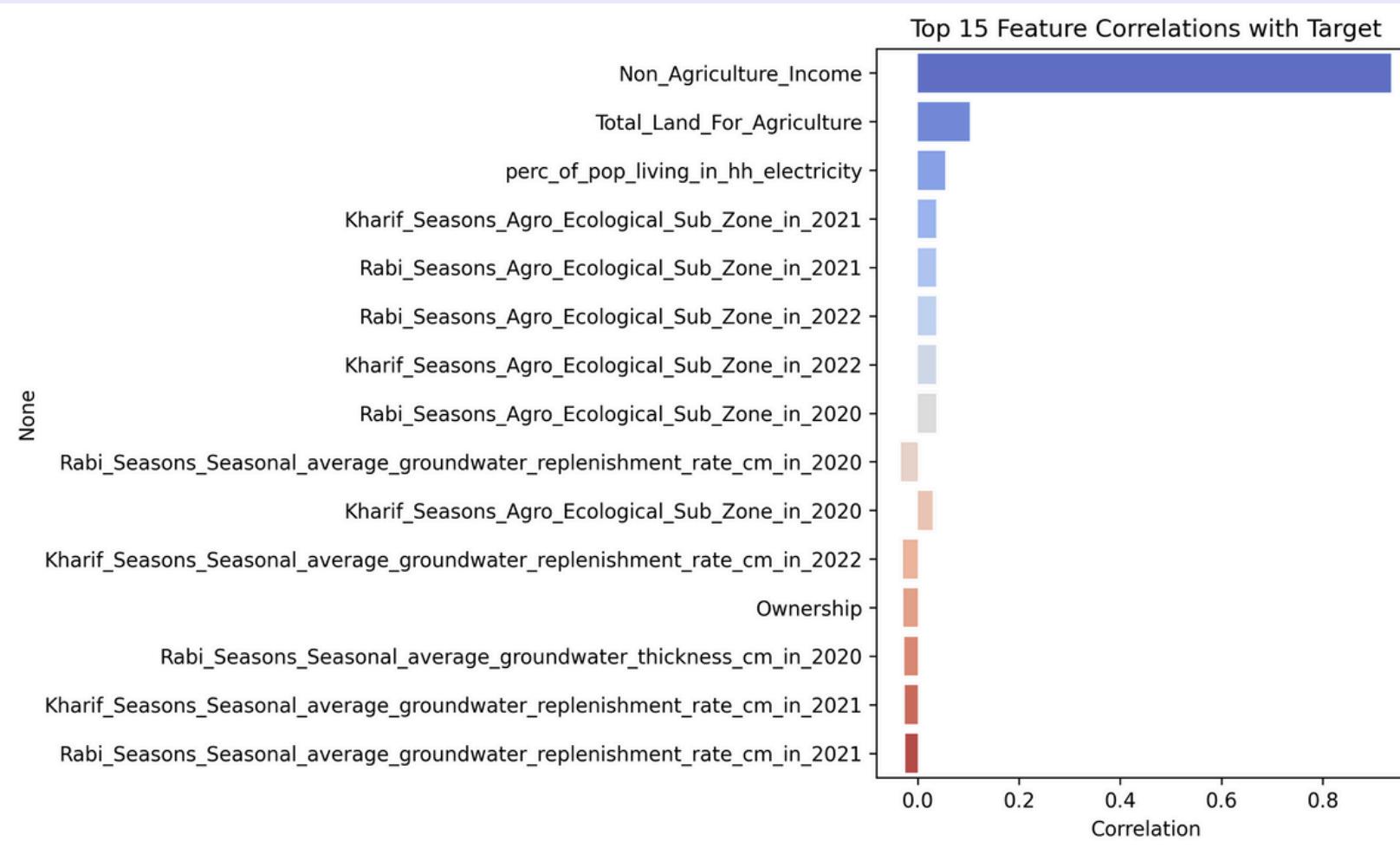




# Data Analysis

- Generated a feature importance plot to highlight the most impactful variables influencing farmer income predictions.
- Plotted Top 15 feature correlations to identify strong predictors of income.
- Created a correlation heatmap to visualize feature-target relationships.
- Used pairplots on top features to explore interactions and clustering.
- Validated key features for modeling through visual and statistical insights.





# KEY BUSINESS INSIGHTS



## Non-Agricultural Income

Farmers with higher secondary income are more financially resilient.



## Agricultural Performance

Recent seasonal yields (2020–2022) directly influence income predictions.



## Landholding & Amenities

Larger land size and access to facilities (e.g., sanitation, roofs, electricity) drive higher incomes



## Geographic Disparities

Location-specific variables (ZIP code, district, city) reveal income gaps and intervention

## REAL-WORLD IMPACT



### Targeted Farmer Support

Helps identify financially vulnerable households for early intervention.



### Informed Loan/Policy Decisions

Supports accurate subsidy distribution and loan disbursement strategies.



### Data-Driven Agricultural Planning

Enables local governments and stakeholders to design evidence-based programs

# Thank You