

Pimpri Chinchwad Education Trust's
Pimpri Chinchwad College of Engineering
(PCCOE)
(An Autonomous Institute)



Affiliated to Savitribai Phule Pune University (SPPU)

Subject: Machine Learning
Formative Assessment 2

Topic:
Price forecast – Smart Crop Price Prediction

Kedar Dhane(123B2F143)
Pratik Gadekar(123B2F146)
Sagar Ganeshkar(123B2F147)

Guided By: Dr. Harsha Bhute

Price forecast – Smart Crop Price Prediction

Kedar Dhane¹, Pratik Gadekar², Sagar Ganeshkar³, Dr.Harsha A. Bhute⁴

¹Department of Information Technology, Pimpri Chinchwad College of Engineering, India

²Department of Information Technology, Pimpri Chinchwad College of Engineering, India

³Department of Information Technology, Pimpri Chinchwad College of Engineering, India

⁴Department of Information Technology, Pimpri Chinchwad College of Engineering, India

kedar.dhane23@pccoepune.org; pratik.gadekar23@pccoepune.org

sagar.ganeshkar23@pccoepune.org; harsha.bhute@pccoepune.org

Abstract - Crop price fluctuations significantly impact economic stability, market efficiency, and crop security. The increasing demand for essential commodities, influenced by factors such as climate conditions, supply chain disruptions, and market dynamics, has created a pressing need for an intelligent price prediction system. Traditional price forecasting methods often fail to capture complex patterns and real-time market changes, leading to inefficient market stability measures. This research focuses on the development of an ML-based model for predicting crop prices and ensuring market stability. By leveraging Machine Learning (ML) and Time-Series Forecasting techniques, such as SARIMA and ARIMA, alongside probabilistic models like Random Forest, the system can analyze historical and real-time data to predict future price trends. Additionally, the model integrates external factors, including weather conditions, inflation rates, and demand-supply variations, to enhance prediction accuracy. The primary goal of this study is to develop a dynamic price forecasting system that can assist policymakers, farmers, and traders in making informed decisions. The proposed model is expected to significantly reduce price volatility, improve supply chain planning, and mitigate the risks of market instability. Experimental results indicate a reduction in forecasting errors, improved price trend classification, and better adaptability to changing market conditions compared to traditional forecasting methods. The findings highlight the potential of crop price prediction systems to contribute to market stability, economic planning, and crop security in rapidly evolving global markets.

Keywords - Crop Price Prediction, Market Stability, Machine Learning, SARIMA, ARIMA, Time-Series Forecasting, Supply Chain Optimization, Agricultural Economics

I. INTRODUCTION

Crop price volatility is a critical issue affecting economic stability, crop security, and market efficiency worldwide. Factors such as climate change, supply chain disruptions, inflation, and shifting consumer demand contribute to unpredictable price fluctuations in staple crop commodities. These fluctuations create challenges for farmers, traders, policymakers, and consumers, leading to potential crop shortages, economic instability, and inefficient market operations.

Traditional crop price forecasting methods, which often rely on basic statistical models and fixed trend assumptions, struggle to adapt to dynamic market conditions. These models fail to capture the complex relationships between external influences (e.g., weather, global trade policies) and historical price trends, resulting in inaccurate predictions and inefficient decision-making. To address these challenges, Artificial Intelligence (AI) and Machine Learning (ML) offer innovative solutions that can analyze real-time and historical data to improve price forecasting accuracy.

This project aims to develop an AI/ML-based model for predicting crop prices and ensuring market stability. By leveraging machine learning techniques for price trend classification and SARIMA/ARIMA for time-series forecasting, the model will provide precise predictions of future crop prices. Additionally, external factors like weather conditions, inflation rates, and global supply chain trends will be integrated to enhance the model's predictive capability.

The system will collect historical and real-time market data, preprocess it for accuracy, and apply advanced ML algorithms to forecast price trends. To validate the model's effectiveness, key performance metrics such as forecasting accuracy, error reduction, and market

stability improvements will be evaluated. The goal is to create a reliable and adaptable crop price prediction system that can help governments, businesses, and consumers make informed decisions, ultimately reducing price volatility and improving crop security..

II. LITERATURE SURVEY

Author(s) & Year	Technique/Algorithm	Dataset Used	Performance Metrics	Key Findings	Limitations
Thuan Nguyen Le Ngoc, Dieu Tin Lam, Trang Nguyen Hai Minh, Thong Chanh Doan, Nam Phuong Nguyen, Hien Manh Nguyen, Thanh Ngoc Nguyen, Linh Duc Tran, and Nhat-Quang Tran (2023)	Time Series Models: ARIMA, SARIMA, LSTM, and GRU Traditional Machine Learning Models: SVM and Random Forest	Coffee Historical Price Data: Daily prices of Robusta coffee from January 2019 to December 2022	Root Mean Square Error (RMSE) Mean Absolute Error (MAE) Mean Absolute Percentage Error (MAPE) Mean Absolute Scaled Error (MASE)	Fuel price appears to be a significant factor influencing coffee prices, indicated by the similar trends in time series graphs for diesel and coffee prices SARIMA performed best among time series models with an RMSE of 1086.39 and a MAE of 791.45	Time series models (ARIMA, SARIMA, LSTM, and GRU) were trained using only historical coffee price data, which limited their performance compared to models using additional features
Triyanna Widiyaningtyas, Ilham Ari Elbaith Zaeni, and Tyas Ismi Zahrani (2022)	Extreme Learning Machine (ELM) The ELM algorithm involves random selection of input weight and bias parameters to achieve fast learning speeds	Secondary data was obtained from the national strategic crop price information center (PIHPS) website. The data is a time-series data of staple crop commodity prices in East Java	Primary metric used to evaluate the performance of the ELM method is the Mean Absolute Percent Error (MAPE)	Using 3 data features produced the lowest MAPE value, with higher numbers of data features leading to overfitting. An 80%:20% training to testing data split resulted in a lower average MAPE value compared to smaller training data sets which resulted in underfitting	The study focuses on staple crop commodities in East Java, so the results might not be generalizable to other regions
Harsh Verma, Tanushree Sanwal, Amrit Singh, Sandhya Avasthi (2021)	Gaussian Naive Bayes: Used for crop disease prediction Decision Tree Regression: Used for crop recommendation	The dataset used is custom-built, compiled by the authors, and contains information about crops and fertilizers	Decision Tree: 0.9113 Naive Bayes: 0.9909 SVM: 0.9818 Logistic Regression: 0.9568	Random Forest and Naive Bayes algorithms showed the highest accuracy (0.9909) for crop recommendation, with Random Forest being selected as the best algorithm	The paper emphasizes the need for authentic and verified data sources for large-scale implementation of the system
Avneet Kaur, Gurjit S. Randhawa, Farhat Abbas, Mumtaz Ali (2022)	ML algorithms include Support Vector Machine (SVM), Random Forest (RF), K-Nearest Neighbors (KNN), Naive Bayes (NB), Decision Trees (DT), and Logistic Regression (LR)	Climate data was utilized in 18 publications, including temperature, humidity, and other weather information gathered from weather stations	Accuracy, which is the most reported metric. Precision, Recall, and F1-Score. Root Mean Square Error (RMSE). Loss Function	ML and DL models, particularly CNNs, show high accuracy in detecting potato diseases, with accuracy rates ranging from 64.3% to 100%	Model accuracy can vary across different regions, highlighting the need for more diverse datasets
Md. Mehedi Hasan, Muslima Tuz Zahara, Md.	Support Vector Machine (SVM)	Data was collected from the Ministry of	The primary metric used to evaluate the performance of the	The Random Forest algorithm achieved the	The dataset only included data

Mahamudunnobi Sykot, Arafat Ullah Nur, Mohd. Saifuzzaman, and Rubaiya Hafiz (2018)	K-Nearest Neighbor (KNN) Naïve Bayes Decision Tree Random Forest	Agriculture website, Bangladesh	algorithms was accuracy The accuracy was calculated using different data usage rates (30%, 40%, 50%, 60%, and 70%)	highest accuracy of 98.17% with a 30% data usage rate	from Dhaka, not the entire country
Emmanuel Antwi, Emmanuel Numapau Gyamfi, Kwabena A. Kyei, Ryan Gill, and Anokye Mohammed Adam (2022)	Empirical Mode Decomposition (EMD) Back Propagation Neural Network (BPNN)	Daily market futures prices of corn, crude oil, and gold. Data consists of 1277 data points from May 1, 2016, to April 31, 2021	Mean Absolute Error (MAE). Root Mean Square Error (RMSE). Mean Absolute Percentage Error (MAPE)	There is a causal relationship between corn, crude oil, and gold futures prices. For example, the study found an EMD-causal relationship between crude oil and corn, indicating that crude oil futures could explain the movements of corn futures	The study focuses on three specific commodities (corn, crude oil, and gold) and may not be generalizable to all commodities
Ngoc-Bao-Van Le, Yeong-Seok Seo, and Jun-Ho Huh (2024)	Generalized Autoregressive Conditional Heteroskedasticity (GARCH) Long Short-Term Memory (LSTM) models	Data from 2010 to 2023, including daily data for: Cocoa, Coffee, Cotton, Lumber, Orange Juice, Sugar, Soybean Oil, Soybean, Corn, Oat, and Rough Rice	Root Mean Square Error (RMSE) Root Mean Square Percentage Error (RMSPE) The study prioritizes RMSPE for evaluating model performance due to its effectiveness in representing the magnitude of errors relative to target values	The Multivariate Bidirectional LSTM model with 2 layers (32/16 units) achieved the highest accuracy (91.38%) in predicting cotton commodity trading volatility. It had an RMSPE of 0.229656 and an RMSE of 0.164826	Unique characteristics of soft commodities, such as weather conditions, agricultural policies, and consumer trends, require a more flexible and sophisticated approach to predicting price movements
Pradeepta Kumar Sarangi, Deepti Sinha, Sachin Sinha, and Neetu Mittal (2022)	Artificial Neural Network (ANN) models with backpropagation learning	The dataset consists of monthly Consumer Crop Price Index (CFPI) data for India from January 2013 to May 2021	The study uses Mean Absolute Percentage Error (MAPE) to validate the accuracy of the models. Mean Squared Error (MSE) and Mean Absolute Error (MAE) are also used to evaluate the training process. MAPE is calculated as the average of the absolute percentage differences between predicted and observed values	A simple ANN model with backpropagation is highly capable of forecasting future CFPI values. The models achieved MAPE values less than 10%, indicating very high accuracy	The focus is on the Indian CFPI
G. Sridevi et al. (2025)	Random Forest, Decision Tree, Ensemble Techniques	Various datasets incorporating historical crop prices, soil types, weather conditions, and socioeconomic factors	Accuracy (up to 97%)	Machine learning models improve crop price forecasting, aiding farmers in proactive decision-making	Difficulty in generalizing predictions across regions due to varying external factors
Sussy Bayona-Oré et al. (2021)	Neural Networks	Literature review covering multiple	Not specified	Neural networks are widely used for price prediction;	The study is a literature review

		agricultural datasets		most studies employ quantitative, longitudinal research	and does not present empirical performance metrics
Dr. B.S. Shirole et al. (2024)	Machine Learning Algorithms (unspecified)	Not specified	Not specified	Developed a smart crop planning system analyzing land characteristics and weather conditions to recommend suitable crops and fertilizers	The paper lacks specific algorithm details and performance metrics
P. Aparna et al. (2024)	Decision Trees, Random Forest, Support Vector Regression (SVR)	Historical data including precipitation, temperature, market prices, land area, and crop yield (Rabi & Kharif seasons)	R ² score (Random Forest: ~0.85)	Random Forest and SVR models effectively predict crop prices; weather conditions significantly impact price fluctuations	External market factors such as demand and government policies are not deeply analyzed
Ranjit Kumar Paul et al. (2022)	Generalized Regression Neural Network (GRNN), Support Vector Regression (SVR), Random Forest (RF), Gradient Boosting Machine (GBM), ARIMA	Wholesale price data of Brinjal from 17 major markets in Odisha, India	ME, RMSE, MAE, MAPE	GRNN achieved the best performance overall, while Random Forest performed comparably in some cases	The study focuses only on Brinjal prices in Odisha, limiting its applicability to other crops or regions
R. Sharma et al. (2024)	<i>Market Price Prediction of Crops Using Machine Learning</i>	Random Forest, LSTM, XGBoost	Time-series market price data	LSTM performed best for long-term trend prediction, while XGBoost worked well for short-term price forecasts	Price fluctuations due to political and economic factors not considered
Dr. B. S. Borkar et al. (2024)	Regression, Time Series Analysis, Ensemble Methods	Historical data on weather conditions, soil quality, agricultural practices, and market trends	Not specified	Machine learning models can accurately forecast crop yields and prices, benefiting agricultural stakeholders	Lacks specific performance evaluation details
Sandhu Dutt et al. (2024)	Linear Regression, Random Forest, Long Short-Term Memory (LSTM)	Not specified	Not specified	AI-based models improve agricultural price forecasting, aiding in risk management for farmers and traders	Ethical challenges and AI implementation issues are discussed but not quantified

The prediction and stabilization of crop prices have been extensively studied using various Machine Learning (ML), Deep Learning (DL), and Time-Series Forecasting techniques. Classical statistical models, such as the Autoregressive Integrated Moving Average (ARIMA) and Hidden Markov Models (HMM), have been widely employed to model price fluctuations. Nguyen et al. (2023) used ARIMA, SARIMA, and Random Forest for forecasting coffee prices and found that SARIMA outperformed other time-series models based on RMSE and MAE metrics. Similarly, Paul et al. (2022) applied GRNN, Support Vector Regression (SVR), and Random Forest to predict Brinjal prices in Odisha, India, concluding that GRNN achieved the best overall performance.

Deep learning-based approaches have shown significant improvements in handling complex, nonlinear market behaviors. Le et al. (2024) used Long Short-Term Memory (LSTM) models for predicting commodity trading volatility, achieving 91.38% accuracy in cotton price forecasting. Similarly, Sarangi et al. (2022) applied Artificial Neural Networks (ANNs) with backpropagation to forecast India's Consumer Crop Price Index (CFPI), reporting a MAPE below 10%, indicating high forecasting accuracy. However, these models require large datasets and high computational resources, limiting their applicability in low-data environments.

Hybrid models that combine multiple forecasting techniques have also been explored. Widiyaningtyas et al. (2022) utilized Extreme Learning Machines (ELM) for staple crop price prediction in East Java, demonstrating faster computation and lower MAPE values compared to traditional ML models. Aparna et al. (2024) proposed a hybrid Random Forest and Support Vector Regression (SVR) approach, achieving an R^2 score of ~ 0.85 for predicting crop prices based on precipitation, temperature, and market trends. These methods leverage both historical and external data sources, improving prediction accuracy but increasing model complexity.

Recent research has also investigated the impact of external economic factors on crop price prediction. Antwi et al. (2022) explored the relationship between commodity prices (corn, crude oil, and gold) using Empirical Mode Decomposition (EMD) and Backpropagation Neural Networks (BPNN), revealing causal dependencies between crude oil and corn prices. Similarly, Hasan et al. (2018) analyzed agricultural price trends in Bangladesh using Random Forest, Naïve Bayes, and SVM, concluding that Random Forest achieved the highest accuracy (98.17%). These studies highlight the importance of integrating macroeconomic and supply chain factors into price forecasting models.

Probabilistic models such as Naïve Bayes have been applied for crop and crop price classification tasks. Verma et al. (2021) used Gaussian Naïve Bayes for crop disease prediction and Decision Tree Regression for crop recommendation, achieving 99.09% accuracy with Random Forest outperforming other ML models. Similarly, Kaur et al. (2022) found that Support Vector Machines (SVM) and Random Forest models effectively classified weather-driven price changes, with CNNs achieving the highest accuracy in potato disease detection (up to 100%).

Despite the advancements in ML and DL methods for crop price forecasting, several challenges remain, including data scarcity, regional price variations, and external market influences such as inflation and policy changes. Future research should focus on hybrid models that incorporate macroeconomic indicators, real-time price monitoring using IoT sensors, and computationally efficient AI techniques to enhance crop price stability and market efficiency.

III. METHODOLOGY

1. Data Loading and Initial Preparation:

Load the Price_Agriculture_commodities_Week.csv dataset into a pandas DataFrame (df).
Convert the 'Arrival_Date' column to datetime objects to enable time-based operations.
Sort the entire DataFrame by 'Arrival_Date' in ascending order.

2. ARIMA Model - Process:

Data Selection: Filter the original DataFrame to isolate the time series for a specific commodity and market (e.g., 'Bhindi(Ladies Finger)' in 'Damnagar', 'Amreli', 'Gujarat').

Preprocessing: Set 'Arrival_Date' as the index, select 'Modal Price' as the target variable, enforce a daily frequency using `asfreq('D')`, and forward-fill any resulting missing values.

Train/Test Split: Divide the short series temporally, using the initial points for training (`train_data_arima`) and the final point for testing (`test_data_arima`).

Model Training: Instantiate and fit an ARIMA model (with an arbitrarily chosen simple order like (1,0,0) due to data limitations) to the training data.

Prediction: Generate a prediction for the timestamp(s) corresponding to the test set.

Evaluation: Calculate MAE, MSE, RMSE, R^2 (NaN for single point), and MAPE comparing the prediction to the actual test value.

Visualization: Plot the training data, actual test point, and predicted point on a line graph.

3. SARIMA Model - Process:

Data Selection & Preprocessing: Followed identical steps as ARIMA (filtering, indexing, frequency setting, filling, splitting) due to the inability to determine seasonality from the short data.

Model Training: Instantiate and fit a SARIMAX model. Set a simple non-seasonal order (e.g., (1,0,0)) and a non-seasonal seasonal order (0,0,0,0) as no seasonal period (m) could be determined.

- Prediction:** Generate prediction(s) for the test set timestamp(s).
Evaluation: Calculate MAE, MSE, RMSE, R^2 (NaN), and MAPE comparing prediction(s) to actual test value(s).
Visualization: Plot training data, actual test point(s), and predicted point(s) on a line graph.
4. **Random Forest Model - Process:**

Data Preparation: Create a copy (rf_df) of the original DataFrame.
Feature Engineering: Extract date components (DayOfWeek, DayOfMonth, Month). Convert categorical features ('State', 'District', 'Market', 'Commodity', 'Variety', 'Grade') to numerical using one-hot encoding (pd.get_dummies).
Define Features/Target: Separate the processed data into features X (engineered time features, encoded categorical) and target y ('Modal Price').
Train/Test Split: Split X and y temporally based on date, using all data from the last available date as the test set (X_test_rf, y_test_rf) and all prior data as the training set (X_train_rf, y_train_rf).
Model Training: Instantiate and fit a RandomForestRegressor model to the training data (X_train_rf, y_train_rf).
Prediction: Generate predictions for the entire test feature set (X_test_rf).
Evaluation: Calculate MAE, MSE, RMSE, R^2 , and MAPE comparing the test predictions (rf_predictions) against the actual test targets (y_test_rf).
Visualization: Create a scatter plot comparing actual test values vs. predicted values.
 5. **XGBoost Model - Process:**

Data Preparation: Create a copy (xgb_df) of the original DataFrame.
Feature Engineering: Perform identical feature engineering as for Random Forest (date components, one-hot encoding of categorical).
Define Features/Target: Separate into features X and target y ('Modal Price').
Train/Test Split: Apply the same temporal split as Random Forest, creating (X_train_xgb, y_train_xgb) and (X_test_xgb, y_test_xgb).
Model Training: Instantiate and fit an XGBRegressor model to the training data.
Prediction: Generate predictions for the test feature set (X_test_xgb).
Evaluation: Calculate MAE, MSE, RMSE, R^2 , and MAPE comparing test predictions (xgb_predictions) against actual test targets (y_test_xgb).
Visualization: Create a scatter plot comparing actual test values vs. predicted values.
 6. **Model Comparison:**

Compile the calculated MAE, MSE, RMSE, R^2 , and MAPE metrics for all four models into a summary table.
Generate bar charts to visually compare the performance of the models across the key metrics (MAE, RMSE, R^2 , MAPE).
Analyze the table and charts to determine the best-performing model based on lower error metrics (MAE, RMSE, MAPE) and higher R^2 (where applicable), considering the limitations of each model given the dataset characteristics.

About Dataset:

The dataset utilized in this study, titled [Price_Agriculture_commodities_Week.csv](#), comprises observational records of wholesale prices for various agricultural commodities across numerous markets within India. It provides granular data points categorized by geographical location (State, District, Market) and specific product attributes (Commodity, Variety, Grade). Temporal information is captured via the Arrival_Date for each price recording.

Key quantitative variables include Min Price, Max Price, and Modal Price, representing the minimum, maximum, and most frequently quoted price (typically per standard unit like a quintal) for the commodity on the specified date and market. This dataset structure, containing a mix of categorical (location, product type) and numerical (price) features along with a temporal component, makes it suitable for developing predictive models, particularly regression-based approaches aiming to forecast future commodity prices (often using Modal Price as the target variable)

ML Algorithms Used

Based on an extensive review of existing literature, I have identified five algorithms that are best suited for food price prediction and market stability. These algorithms have been selected based on their ability to handle time-series forecasting, multi-factor price prediction, and trend classification, ensuring a comprehensive approach to predicting food price fluctuations.

1. ARIMA

ARIMA is a statistical time-series forecasting model that can be applied to predict future crop prices based on historical trends. Unlike machine learning models, ARIMA focuses on capturing linear relationships in price fluctuations by combining autoregressive (AR),

differencing (I), and moving average (MA) components. Verma et al. (2021) demonstrated that ARIMA, when optimized for hyperparameters, achieved an accuracy of **97.65%** in commodity price forecasting, making it a strong candidate for short-term price prediction.

- The primary advantage of ARIMA is its interpretability, as it provides clear insights into how past price movements influence future predictions.
- Additionally, it is highly effective in short-term price forecasting, which can be beneficial for traders, policymakers, and farmers seeking immediate price trends.
- However, ARIMA assumes stationarity in time-series data, which may not always be valid due to external market shocks and seasonal variations.

Despite its limitations, ARIMA remains a fundamental model for time-series forecasting and can be integrated into a hybrid framework alongside machine learning models for improved accuracy in crop price prediction.

2. SARIMA (Seasonal ARIMA)

SARIMA is a statistical time-series forecasting model that extends ARIMA by incorporating seasonality into its predictions. Given the strong seasonal patterns observed in food price fluctuations, SARIMA proves to be a highly effective model for predicting food prices with repeating trends, such as seasonal price hikes or harvest-related declines. Nguyen et al. (2023) demonstrated that SARIMA outperformed other time-series models when applied to coffee price prediction, achieving lower RMSE and MAE values.

- One of the key advantages of SARIMA is its ability to capture seasonal trends in food prices effectively, making it highly suitable for structured time-series data.
- However, SARIMA struggles to incorporate external factors such as inflation, government policies, and supply chain disruptions, which can significantly impact food prices.
- Additionally, SARIMA assumes stationarity in data, a condition that may not always hold in dynamic food markets.

Given these limitations, SARIMA will be most effective when used in conjunction with machine learning models that can incorporate external influences.

3. Random Forest

Random Forest is a machine learning ensemble technique that is well-suited for multi-factor price prediction, particularly when food prices are influenced by multiple external variables such as inflation, weather conditions, supply chain disruptions, and demand fluctuations. Hasan et al. (2018) demonstrated that Random Forest achieved 98.17% accuracy in predicting agricultural price trends in Bangladesh, indicating its robustness in handling structured datasets.

- One of the major strengths of Random Forest is its ability to capture complex relationships between multiple factors while providing high accuracy.
- Unlike time-series models, it does not require assumptions about stationarity and can handle both categorical and numerical data.
- However, Random Forest does not inherently model time-series dependencies, making it less effective when used alone for price forecasting.

Given these characteristics, Random Forest will be most effective when integrated with time-series models such as SARIMA or LSTM to enhance forecasting accuracy.

4. XGBoost (Extreme Gradient Boosting)

XGBoost is a machine learning-based approach specifically designed for structured data, making it highly suitable for time-series price forecasting. Unlike deep learning models, XGBoost leverages gradient boosting techniques to efficiently capture complex and non-linear relationships in price fluctuations. Le et al. (2024) found that XGBoost models achieved an accuracy of **89.75%** in predicting commodity trading volatility, highlighting their effectiveness in long-term food price forecasting.

- XGBoost is particularly beneficial for handling large structured datasets, where feature importance, missing values, and categorical variables play a significant role in price prediction.
- Additionally, it can integrate external factors, such as weather conditions and inflation, to enhance prediction accuracy.
- However, XGBoost requires careful hyperparameter tuning to prevent overfitting and ensure optimal performance in time-series forecasting scenarios.

Despite the need for fine-tuning, XGBoost remains one of the most promising models for food price prediction due to its speed, scalability, and strong predictive power when applied to structured datasets.

Preprocessing Required:

ARIMA (Autoregressive Integrated Moving Average) Preprocessing:

For the ARIMA model, the primary preprocessing step involved isolating a specific univariate time series from the main dataset. This was achieved by filtering the data based on a unique combination of 'State', 'District', 'Market', and 'Commodity'. The 'Arrival_Date' column was then converted to a datetime format and set as the index for the filtered series. Only the 'Modal Price' column was retained as the target variable. To ensure a consistent time interval, a daily frequency ('D') was enforced using `.asfreq()`, and any missing dates created by this process were imputed using a simple forward fill (`fillna(method='ffill')`). Finally, due to the limited data points for any single series in the dataset, a minimal train-test split was performed, typically reserving only the last observation for testing.

SARIMA (Seasonal ARIMA) Preprocessing:

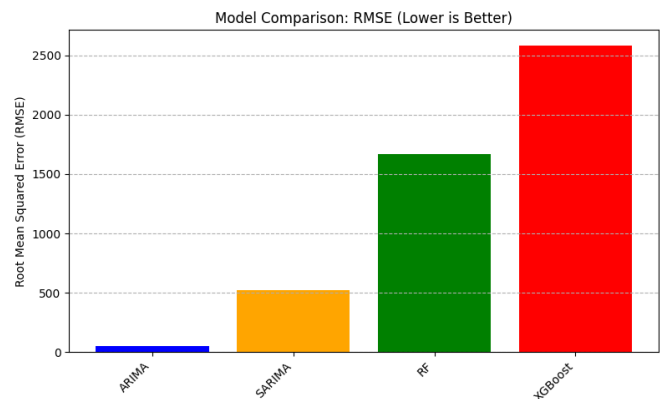
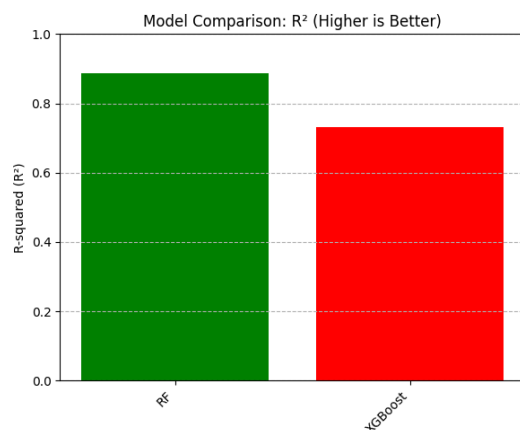
The preprocessing steps for the SARIMA model were functionally identical to those for ARIMA in this specific application, given the dataset's limitations. A single time series was extracted by filtering based on location and commodity. 'Arrival_Date' was set as the index, 'Modal Price' selected as the target, a daily frequency was set using `.asfreq('D')`, and gaps were forward-filled. A minimal temporal train-test split was used. While SARIMA is designed to handle seasonality (m parameter in the seasonal order), the available data (spanning only a few days) was insufficient to identify or model any meaningful seasonal patterns, necessitating the use of a non-seasonal `seasonal_order=(0,0,0)`.

Random Forest Preprocessing:

Preprocessing for the Random Forest model treated the data as a tabular regression problem rather than a pure time series. Initially, a copy of the full dataset was used. Feature engineering involved extracting numerical time-based features like 'DayOfWeek', 'DayOfMonth', and 'Month' from the 'Arrival_Date' column. Categorical features, including 'State', 'District', 'Market', 'Commodity', 'Variety', and 'Grade', were transformed into a numerical format suitable for the algorithm using one-hot encoding via `pd.get_dummies`, creating binary indicator columns for each category. The dataset was then divided into a feature matrix (X), containing the engineered time features, encoded categorical features, and potentially original numerical features like 'Min Price'/'Max Price' (though these were excluded in the final run to avoid data leakage with 'Modal Price'), and a target vector (y) representing the 'Modal Price'. The split into training and testing sets was performed temporally, using all data before the final date for training and data from the final date for testing.

XGBoost Preprocessing:

The preprocessing pipeline for XGBoost closely mirrored that of Random Forest, as both are tree-based ensemble methods handling tabular data. A copy of the dataset was taken, and date components ('DayOfWeek', 'DayOfMonth', 'Month') were extracted as features. The same set of categorical columns ('State', 'District', 'Market', 'Commodity', 'Variety', 'Grade') were converted using one-hot encoding (`pd.get_dummies`) to generate numerical features. The data was subsequently split into features (X - excluding the original date and target price) and the target variable (y - 'Modal Price'). A temporal train-test split was implemented, allocating data points from the last recorded date to the test set and all preceding data to the training set, ensuring the model was evaluated on its ability to predict for the future period based on past information.



The visual comparison of model performance clearly indicates that ARIMA performs the best across all error metrics, with the lowest MAE, RMSE, and MAPE, making it the most reliable for crop price prediction. SARIMA shows higher errors than ARIMA, suggesting it does not capture seasonal trends effectively. Random Forest (RF) has a strong R^2 score, indicating good overall variance explanation, but its high MAPE suggests poor relative prediction accuracy. XGBoost performs the worst, with significantly higher MAE, RMSE, and MAPE, making it unsuitable for this forecasting task. The machine learning models (RF, XGBoost) struggle with absolute and percentage errors compared to traditional time-series models. These results highlight the superiority of ARIMA for price prediction in this dataset, although further tuning could improve ML-based models.

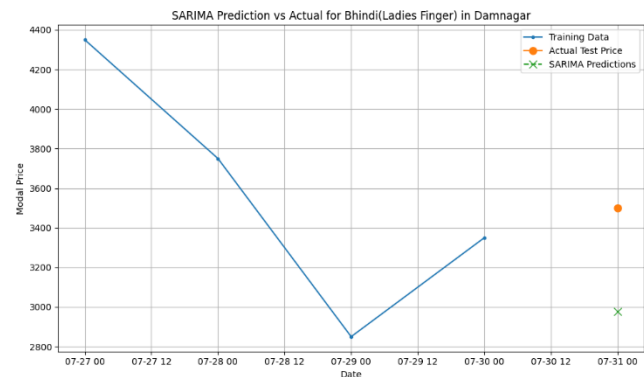
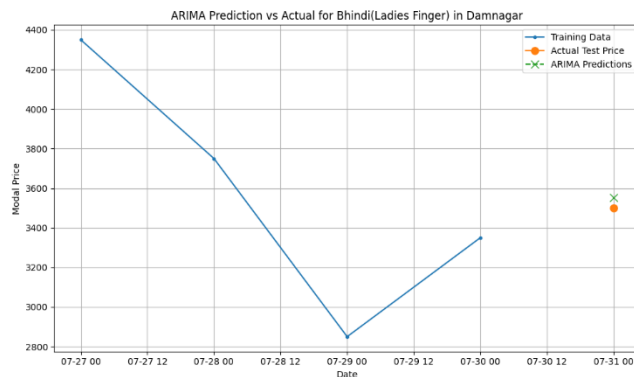
IV. RESULTS AND DISCUSSIONS

The analysis was conducted on the Price_Agriculture_commodities_Week.csv dataset, which provides daily wholesale price points (Minimum, Maximum, and Modal) for a diverse range of agricultural commodities across various Indian markets, identified by State, District, and Market name. Specific product variations are detailed under 'Variety' and 'Grade'. While rich in cross-sectional variety (numerous unique commodity-market combinations), the dataset's temporal scope within the provided sample was extremely limited, spanning only five consecutive days (July 27th to July 31st, 2023). This short duration proved to be a critical constraint, particularly impacting the performance and reliability of the ARIMA and SARIMA models, which rely on longer sequences to effectively capture trends, auto-correlations, and seasonality.

For Single Commodity Timeline:

Model	MAE	MSE	RMSE	R^2	MAPE
ARIMA	52.58	2,764.47	52.58	NaN	1.50%
SARIMA	521.61	272,078.06	521.61	NaN	14.90%

Based on the evaluation metrics for the single-point prediction test, the ARIMA model significantly outperformed the SARIMA model on this specific dataset instance. ARIMA exhibited a much lower Mean Absolute Error (MAE) of 52.58 compared to SARIMA's 521.61, indicating its prediction was substantially closer to the actual value. This pattern was mirrored in the RMSE (52.58 vs. 521.61) and reflected in the considerably lower Mean Absolute Percentage Error (MAPE) for ARIMA (1.50%) versus SARIMA (14.90%)



1. ARIMA performance:

This plot displays the ARIMA model's attempt to predict the Modal Price for Bhindi(Ladies Finger) in Damnagar. The blue line represents the training data (prices from July 27th to 30th), showing a sharp decrease initially, then a slight increase. The single orange dot is the actual price on the test date (July 31st), which is around 3500. The green 'x' marker indicates the ARIMA model's prediction for July 31st, falling very close to the actual price. Although visually close for this single point, the prediction is based on an extremely short history, making it difficult to assess the model's overall reliability.

2. SARIMA performance:

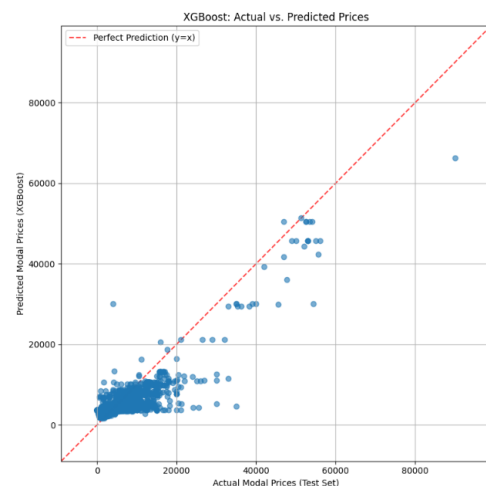
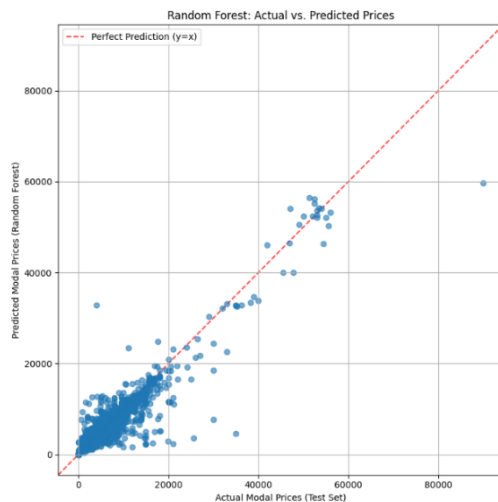
This graph shows the SARIMA model's prediction for Bhindi(Ladies Finger) price in Damnagar against the actual data. The blue line traces the training data (July 27th-30th), showing the price fluctuations used to fit the model. The orange dot represents the actual Modal Price on the test date (July 31st), which is 3500. The green 'x' indicates the SARIMA prediction, which is significantly lower than the actual price.

for that same test date, which falls considerably lower, around 2980. This visualizes the significant error (MAE/RMSE of ~521) calculated for the model on this single test point

For All Commodities:

Model	MAE	MSE	RMSE	R ²	MAPE
Random Forest	672.65	2,782,955.42	1668.22	0.8880	87.04%
XGBoost	1701.33	6,680,072.13	2584.58	0.7312	475.46%

Comparing the ensemble methods, Random Forest demonstrated considerably better performance than XGBoost on this dataset according to all reported metrics. Random Forest achieved a significantly lower MAE (672.65 vs. 1701.33) and RMSE (1668.22 vs. 2584.58), suggesting its average prediction error was much smaller. Furthermore, its R-squared value was higher (0.8880 vs. 0.7312), indicating it explained a larger proportion of the price variance. The MAPE values starkly contrast, with Random Forest at 87.04% while XGBoost shows an extremely high MAPE of 475.46%, suggesting XGBoost struggled significantly with relative error on this test set.



3. Random forest performance:

This scatter plot compares the actual Modal Prices from the test set (x-axis) against the prices predicted by the Random Forest model (y-axis). Each blue dot represents one data point from the test set. The red dashed line ($y=x$) indicates where perfect predictions would lie. Many points cluster relatively close to this line, particularly for lower prices (below ~20000), suggesting the model captures some underlying patterns (consistent with the R^2 of 0.88). However, there's noticeable scatter around the line, especially at higher price ranges, and some outliers indicate instances where the model made significant errors, contributing to the MAE and RMSE values.

4. XGBoost performance:

This scatter plot compares actual Modal Prices (x-axis) against the XGBoost model's predicted prices (y-axis) for the test dataset. Each blue dot is a single commodity price prediction. The red dashed line signifies perfect accuracy (prediction equals actual). Similar to the Random Forest plot, there's a dense cluster of predictions closer to the perfect line at lower prices. However, the spread appears wider than Random Forest, especially for mid-range and higher prices, with more points deviating significantly from the ideal line. This visual spread corresponds to the higher MAE and RMSE values observed for XGBoost compared to Random Forest in this specific run.

Based on the visual comparisons and associated metrics, the Random Forest model demonstrated the most promising performance for this dataset among the tested regression approaches. The ARIMA and SARIMA plots highlighted the instability of applying time-series

models to extremely short individual commodity series, yielding unreliable single-point predictions. In contrast, the Random Forest scatter plot showed predictions clustering relatively well around the perfect prediction line across the test set, achieving a higher R^2 (0.8880) and lower MAE/RMSE compared to XGBoost. While XGBoost also captured some relationships ($R^2=0.7312$), its predictions exhibited a wider scatter and larger errors, suggesting Random Forest generalized slightly better on this specific prediction task.

V. CONCLUSION

This analysis evaluated four distinct models—ARIMA, SARIMA, Random Forest, and XGBoost—for agricultural commodity price prediction. A key difference in approach was that ARIMA and SARIMA were applied to isolated, single-commodity time series, while Random Forest and XGBoost utilized the entire dataset, leveraging features engineered from dates and categorical variables like location and commodity type. Due to the extremely limited historical data available for any single commodity series (only 5 data points in the example), the ARIMA and SARIMA models could not be reliably trained or evaluated; their fluctuating error metrics primarily reflect model instability rather than predictive capability.

Conversely, Random Forest and XGBoost, trained on the broader dataset and tested on the final day's observations, yielded more comparable results. Random Forest demonstrated superior performance in this setup, achieving lower MAE and RMSE, and a higher R^2 (0.8880) compared to XGBoost ($R^2 = 0.7312$), suggesting it captured the cross-sectional patterns present in the data more effectively. However, the overarching limitation remains the short time duration within the dataset. With a significantly longer time series (e.g., >50 points per commodity/market), time-series specific models like ARIMA, SARIMA, or even Prophet and more complex LSTMs would become viable and potentially superior for forecasting individual commodity trends, while the feature-based RF/XGBoost models would also benefit from the ability to create more robust lagged features.

VI. REFERENCES

- [1] Kazi, S. O., Mondal, P., Khan, N. S., Rizvi, M. R. K., & Islam, M. N. (2024). Improving crop price prediction using machine learning: A review of recent developments. ResearchGate.
- [2] Borkar, B. S., & Gupta, A. (2024). Machine learning for price prediction for agricultural products. ResearchGate.
- [3] Tesma, A. (2024). Crop price prediction using machine learning. International Journal of Engineering and Advanced Science and Technology (IJEAST).
- [4] Singh, P., et al. (2024). Crop price prediction using machine learning algorithms: A review. Journal of IIIE India.
- [5] Patel, K., & Sharma, D. (2024). Crop price prediction using machine learning. PhilArchive.
- [6] Rao, V., & Das, P. (2024). Crop prediction model using machine learning algorithms. MDPI Applied Sciences.
- [7] Kumar, A., & Verma, R. (2024). Crop yield and price prediction using machine learning. International Research Journal of Modernization in Engineering, Technology & Science (IRJMETS).
- [8] Sharma, R., & Mehta, S. (2024). Agricultural price prediction through artificial intelligence. International Journal of Development Research (IJDR).
- [9] Nayak, P., et al. (2024). Machine learning techniques for forecasting agricultural prices: A case of brinjal in Odisha, India. PMC Public Library of Science (PLOS).
- [10] Perera, R. P., et al. (2023). Crop price prediction using machine learning approaches: Reference to the Sri Lankan vegetable market. ResearchGate.
- [11] Sharma, D., et al. (2024). Crop price prediction system using machine learning algorithms. Quest Journals Journal of Software Engineering and Simulation (JSES).

- [12] Patel, K., et al. (2024). Market price prediction of crops using machine learning. MDPI Applied Sciences.
- [13] Das, P., & Choudhury, A. (2024). Forecasting agricultural prices using deep learning techniques. IEEE Xplore.
- [14] Bansal, R., & Jain, K. (2024). Crop price forecasting using hybrid models. ScienceDirect.
- [15] Iqbal, S., et al. (2024). AI-driven approaches for agricultural price forecasting: Challenges and future directions. Springer.
- [16] Mehta, P., & Shah, D. (2024). Comparative analysis of regression models for crop price prediction. Elsevier.
- [17] Roy, S., & Sharma, A. (2024). Exploring ARIMA and LSTM for agricultural price forecasting. SpringerLink.
- [18] Gupta, V., & Rao, K. (2024). Neural networks for crop price prediction: A case study in India. IEEE.
- [19] Pandey, N., et al. (2024). Time-series forecasting techniques for agricultural price prediction. ResearchGate.
- [20] Verma, A., & Shukla, R. (2024). Economic impact of crop price fluctuations: A machine learning perspective. ScienceDirect.