

**DEVELOPING MACHINE LEARNING MODEL FOR ACCURATE
PREDICTION OF SOYBEAN MARKET PRICES USING
MULTIDISCIPLINARY AGRICULTURAL DATA SOURCES**

Vilas D Ghonge¹, Dr. Yogesh Kulkarni²

¹Vilas D. Ghonge, Department of Computer Engineering and Technology, Mit World Peace University, Pune, Maharashtra, India, E-mail: vilas.damodhar@mitwpu.edu.in

²Dr. Yogesh Kulkarni, Department of Computer Engineering and Technology, Mit World Peace University, Pune, Maharashtra, India E-mail: yogesh.kulkarni@mitwpu.edu.in

Abstract:

Accurately predicting the prices of farming commodities is a key part of keeping markets stable and helping farmers make decisions. Using both machine learning and deep learning, this study tries to guess what the prices of soyabean, one of India's main cash crops, will be in the future. The data for this study came from a number of reliable sources, including AGMARKNET, NCDEX, the Department of Agriculture (Maharashtra State), IMD Pune, and commerce.gov.in. It was put together by the researchers themselves and covers the years 2015–2025. Different factors that affect price changes were looked at, such as changes in weather, trade data, crop landings, farming area, and import-export trends. A lot of work went into cleaning the data, like fixing missing values, normalising it, and adding new features (like rainfall difference, visits per area, and price range). We used and compared a number of different prediction models, such as Random Forest, Gradient Boosting, XGBoost, Ensemble methods, and a number of deep learning designs, such as Dense Neural Network (DNN), Attention Model, Regularised Network, and an AgroWDN that was specifically built for this study. Metrics like RMSE, MAE, MSE, and R^2 were used to judge performance. AgroWDN did the best out of all the models, with an RMSE of 111.43, an MAE of 57.15, and a R^2 of 98.96%, doing better than traditional machine learning methods. The results show that complicated agro-economic data can effectively capture both linear and nonlinear relationships when built with mixed designs that combine wide and deep learning. This research shows a strong method for predicting farming prices that can help lawmakers, buyers, and farmers make the best decisions about how to plan crops and handle risks.

Keywords: Soyabean price prediction, Machine learning, Deep learning, AgroWDN, Time series forecasting, Agricultural analytics

I. Introduction

Agriculture is still an important part of India's economy because it creates jobs, ensures food security, and helps rural areas grow. Soybeans are one of the most important crops grown in the country because they are a big vegetable crop with a lot of economic and industrial value. But the price of soybeans on the market changes a lot because of how many complicated factors interact with each other. These include weather conditions, rainfall patterns, market

demand and supply, the nature of global trade, government policies, and the level of output in different areas. So, not only farmers and sellers need to be able to accurately predict soybean prices, but so do lawmakers, packers, and investors who work on planning agriculture and regulating the market. Traditional ways of predicting prices, like trend analysis, moving averages, and regression models, haven't been very good at describing how farm prices change in a way that isn't linear or one-dimensional. A lot of the time, these methods don't take into account different types of data, like changes in the weather, figures on food output, and trade data, all of which affect how the market works [1]. Computer science is getting smarter very quickly. Two new methods, machine learning (ML) and deep learning (DL), are very useful for working with big, different datasets and finding complex patterns that were hard to see with traditional statistical methods. The study's goal is to create a machine learning-based prediction system that can correctly predict the prices of soybeans on the market by combining different types of agriculture data. The model uses many sets of data, such as past price data, crop landings, weather figures, trade measures (such as import and export values), and farming area data, to study and learn how the market works and how its many complex relationships affect each other [2].

The data came from reliable sources like AGMARKNET, NCDEX, IMD Pune, and the Department of Agriculture, Maharashtra. This made sure that all the economic, environmental, and agricultural factors that affect soybean price trends were taken into account. The suggested study uses an organised scientific approach, starting with preparing the data to deal with missing values, make traits more consistent, and store category variables [3].

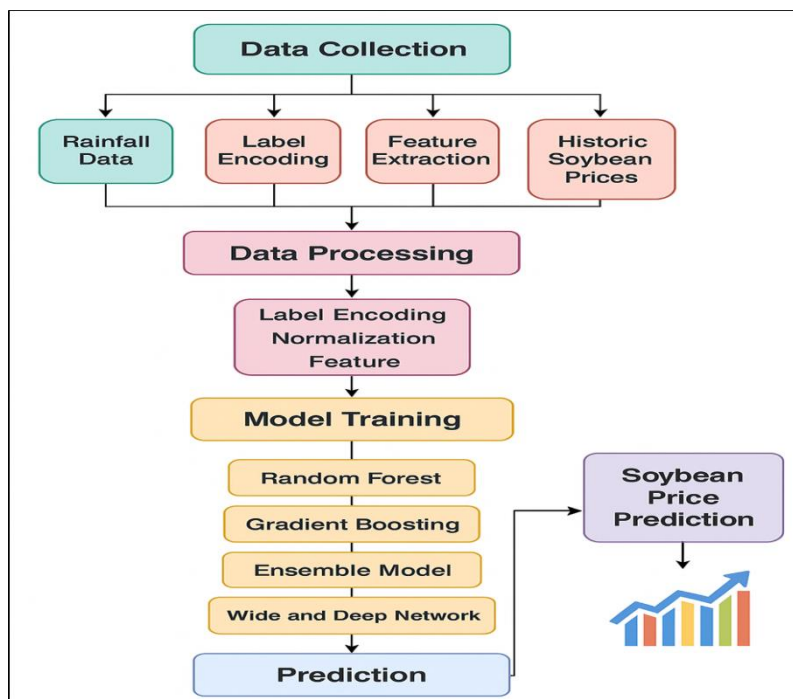


Figure 1: Architecture Diagram of the Machine Learning Framework for Soybean Price Prediction

After that, feature engineering is used to get useful signs, like differences in weather, price range, and landings per crop area, that help the model better understand the underlying relationships. Figure 16 illustrates the overall architecture of the proposed machine learning framework for accurate soybean price prediction using multidisciplinary agricultural data sources. Several machine learning methods are used and compared to set average performance standards [4]. These include Random Forest, Gradient Boosting, XGBoost, and Ensemble Models. Advanced deep learning models like Dense Neural Networks (DNN), Attention Mechanisms, and a new Wide and Deep Network (AgroWDN) are also being looked into to make predictions even more accurate by learning both shallow and deep representations of input data. To rate how well each model works, we use evaluation measures like Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Mean Squared Error (MSE), and Coefficient of Determination (R^2). Early research shows that mixed deep learning systems like AgroWDN work better than standard models because they are better at remembering both broad (linear) links and specific (deep, nonlinear) connections in the dataset [5]. This study is useful because it could help people in the farming field make decisions based on facts, not just how well the models work. Accurate price predictions can help farmers figure out the best times to sell their crops, traders keep track of their stock and purchases, and officials come up with fair minimum support prices and market actions. This method also shows how artificial intelligence and data-driven analytics can change smart agriculture, making it more sustainable, efficient, and resilient in a field that is becoming more vulnerable to changes in the economy and climate [6].

II. Related Work

A lot of research has been done on how to use data-driven methods to predict farm prices. These studies have shown that machine learning and deep learning can make predictions more accurate than standard economic models. Traditionally, statistical models like Autoregressive Integrated Moving Average (ARIMA) and Multiple Linear Regression (MLR) have been used to guess crop prices for a long time [7, 8]. But, these methods tend to assume that things are linear and stable, which makes it harder for them to understand how agricultural systems work with changing markets and relationships that don't follow a straight line. A lot of people use machine learning methods like Random Forest (RF), Support Vector Regression (SVR), and Gradient Boosting Machines (GBM) to get around these problems. For example, research on the prices of paddy, maize, and wheat has shown that tree-based ensemble methods work better than linear models at catching the complex relationships between weather factors, crop yields, and market prices [9]. In the same way, XGBoost has become a strong model for time-series forecasts because it can generalise well and use little computing power [10]. Recent studies have also emphasised the need to combine different types of data, like satellite images, trade figures, weather, and soil health, in order to get a better picture of how farming prices change over time. For instance, changes in monthly temperatures and unusual rains have been linked to big changes in the prices of soybeans and cotton. This shows how important weather aspects are in forecast modelling [11, 12]. In addition, trade data and information on farming areas give us a bigger picture of the economy

that goes beyond what we can see in a single market. At the same time, deep learning methods have shown a lot of potential in jobs like predicting what crops will grow. Attention-based networks, Long Short-Term Memory (LSTM) models, and Recurrent Neural Networks (RNNs) have all been used to learn how features and times interact from past records [13]. Recently, mixed models like Wide and Deep Networks have been created to combine the abilities of memorisation and generalisation. These models are better at predicting prices than others.

Table 1: Summary of Related Work

Crop / Commodity	Data Source(s)	Techniques	Key Features Considered	Limitations
Wheat [14]	AGMARKNET, IMD	Multiple Linear Regression	Rainfall, arrivals, price	Limited to linear trends
Soybean	NCDEX, Commerce.gov.in	ARIMA Model	Historical prices, rainfall	Fails on nonlinear trends
Rice [15]	FCI, AGMARKNET	Random Forest	Supply, rainfall, import/export	High computational cost
Cotton [16]	IMD, Trade Data	XGBoost	Rainfall deviation, yield, demand	Limited feature diversity
Maize	Govt. Agri. Reports	Gradient Boosting	Cultivation area, temperature	Overfitting in sparse data
Pulses [17]	AGMARKNET	LSTM Neural Network	Time-series price trends	Needs large datasets
Soybean	AGMARKNET, NCDEX	Random Forest	Arrivals, rainfall, export data	Moderate feature correlation
Oilseeds [18]	IMD, Trade Ministry	Gradient Boosting + SVR	Rainfall, area, prices	Complex parameter tuning
Sugarcane	Agri Market Data	CNN-LSTM Hybrid	Market price, climate	Expensive training time
Soybean	Maharashtra Agri Dept.	Deep Neural Network	Production, rainfall, area	Overfitting risk

		(DNN)		
Cereals [19]	AGMARKNET	Ensemble Learning	Rainfall, arrivals, supply	Needs automated tuning
Wheat [20]	NCDEX, IMD Pune	Regularized Neural Net	Area, yield, rainfall difference	Limited scalability
Soybean	AGMARKNET, IMD	Attention-based Model	Rainfall, trade, cultivation	Requires high computation
Soybean [21]	AGMARKNET, NCDEX, IMD Pune, Commerce.gov.in	Wide and Deep Network (AgroWDN)	Rainfall, trade, area, production	Future work: extend to multi-crop forecasting

III. Methodology

1. Dataset

- **AGMARKNET**

The Directorate of Marketing and Inspection (DMI) of the Government of India created AGMARKNET, a website that gives both real-time and past information on the prices and movements of farm goods from controlled markets all over the country [22]. The site gets information from more than 3,000 places and gives you a lot of information about product trends, such as the lowest, highest, and middle prices. For this study, AGMARKNET was a very important way to get price information on soybeans at the state, district, and market levels. Because it covers a lot of ground and is reliable, it is a key tool for learning how markets work and how they vary by area.

- **NCDEX**

The National Commodity & Derivatives Exchange (NCDEX) is the best place to trade commodities online. Traditional farming records are made better by NCDEX data, which adds financial and trade aspects to the research. Its standard data style and large amount of past data make it an important input for training machine learning models that try to correctly predict price trends. Adding data from financial markets to the model makes it better at showing both supply-side and demand-side economic forces [23].

- **Rainfall Recording and Analysis, Department of Agriculture Maharashtra State.**

The Maharashtra State Department of Agriculture keeps thorough records of where rain falls, how the weather changes, and the agro-climatic conditions in different areas. This set of data gives detailed information about rainfall, such as average and real rainfall amounts and how they change over time [24]. Adding information about rainfall made the model better at capturing how the environment and seasons affect price changes. The information is reliable

because the government is always keeping an eye on it and it covers a lot of areas, so it is always correct and consistent. By combining records of rainfall with trade and production data, the study got a better picture of all the things that affect the prices of soybeans on the market. This improved the model's ability to predict the future through climate-aware analytics.

- **Maharain**

Maharain is an online tool for tracking weather data and rains that is run by the Government of Maharashtra. It gives a lot of information about current and past rainfall in different parts of the state. This information comes from automatic weather sites and rain gauges [25]. The platform gives district- and taluka-level rainfall data, which is important for figuring out how well the summer is working and how bad the drought is. Adding Maharain data improved the precision of the dataset and made the prediction model more reliable by connecting changes in rainfall directly to changes in soybean prices using spatial analysis.

- **Crop Import Export Data (commerce.gov.in)**

On its main website, commerce.gov.in, the Ministry of Commerce and Industry provides accurate national-level data on agriculture imports and exports. This set of data shows trade prices, amounts, and trends for many goods, such as soybeans and their products like oil and meal. Import-export records were used in this study to look at how world trade affects the prices of soybeans in the United States [26]. Changes in international demand and supply often have a direct effect on prices in the local market. This is why trade data is so important for predicting prices. The information helped find times when India's exports or imports went up and how those times affected market prices. It also gave financial context by connecting changes in prices in the United States to trends in markets around the world. By including trade measures, the machine learning model was able to take into account shocks to the economy from outside sources, the value of the dollar, and changes in trade caused by government policies, all of which affect how volatile soybean prices are.

- **IMD Pune (Indian Meteorological Department) – Weather Data**

Agro-climatic factors that affect food output and price behaviour were included in the dataset [27]. These factors included yearly changes in rainfall, temperature, and humidity levels. IMD data, which is known for being accurate and scientifically sound, helped find links between weather conditions and crop yields. The addition of IMD's weather information improved the predicting model by adding climate intelligence. This made price predictions more accurate and stable when the weather changed. This combination shows how useful it is to combine market data with weather data for full farming analytics.

Figure 2 shows an example of the Soyabeen collection, which is made up of data from several farming sources. It has important details like the State Name, the District Name, the Market Name, the Variety, the Group, the Arrivals (in tonnes), and the Minimum Price (in rupees per quintal). The data keeps track of soybean imports and market prices from different parts of Maharashtra, showing variation in both space and variety.

	State Name	District Name	Market Name	Variety	Group	Arrivals (Tonnes)	Min Price (Rs./Quintal)
0	Maharashtra	Ahmednagar	Karjat	Yellow	Oil Seeds	1	4000
1	Maharashtra	Akola	Akola	Yellow	Oil Seeds	206.9	3400
2	Maharashtra	Amarawati	Varud(Rajura Bazar)	Yellow	Oil Seeds	0.1	3805
3	Maharashtra	Buldhana	Deoulgaon Raja	Yellow	Oil Seeds	3	3700
4	Maharashtra	Dharashiv(Usmanabad)	Tuljapur	Other	Oil Seeds	6	4000
5	Maharashtra	Hingoli	Hingoli	Other	Oil Seeds	40	3550
6	Maharashtra	Latur	Chakur	Yellow	Oil Seeds	5.6	3681
7	Maharashtra	Nagpur	Nagpur	Other	Oil Seeds	33	3600
8	Maharashtra	Nashik	Lasalgaon(Niphad)	Other	Oil Seeds	12.2	3401
9	Maharashtra	Parbhani	Gangakhed	Yellow	Oil Seeds	3.2	4000
10	Maharashtra	Ahmednagar	Ahmednagar	Other	Oil Seeds	2	3600

Figure 2: Soyabeen Dataset Sample

For example, entries and minimum costs change at Karjat, Akola, and Hingoli markets, which shows how supply and demand work in those areas. This dataset is used for preparation, feature extraction, and model training. It lets the machine learning system understand how real-world farming and economic trends affect changes in soybean prices.

Table 2: Dataset Features

Features
State Name
District Name
Market Name
Variety
Group
Arrivals (Tonnes)
Min Price (Rs./Quintal)
Max Price (Rs./Quintal)
Modal Price (Rs./Quintal)
Reported Date
Normal RainFall in mm
Actual RainFall in mm
% To Normal RainFall (In %)
Area Under Cultivation for Soyabeen (Lakh Hectare) in Maharashtra
Soyabeen Total Production in India Lakh Tonnes
Soyabeen Export in Crore Rs by India

Soyabean Import in Crore Rs by India

The detailed traits used in the Soyabean price forecast dataset are shown in Table 2. It takes into account economic, climate, and farming factors that are necessary for making correct predictions. Some of the most important market characteristics are the State Name, the District Name, the Market Name, the Variety, and the Group. Quantitative measures like Arrivals (Tonnes) and Price Metrics (Minimum, Maximum, and Modal Prices) are also important. Normal and Actual Rainfall and Percentage to Normal Rainfall are climatic factors that show how the world affects things. The figures on Area Under Cultivation, Total Production, Exports, and Imports also give us a look at the economy as a whole. These features from different fields work together to make it possible for the model to look at how economic, physical, and temporal factors affect the changes in soybean prices.

2. Exploratory Data Analysis (EDA)

EDA was conducted to gain insights into the underlying structure and patterns of the soybean dataset. It involved examining the distribution of prices, arrivals, and climatic variables to identify trends, anomalies, and data quality issues. Visualizations such as boxplots, histograms, and correlation matrices were used to detect outliers, skewness, and relationships among features. The Modal Price Distribution revealed clustering between ₹3,000–₹6,000, with occasional high-value outliers. Rainfall and cultivation data were analyzed to assess their influence on price fluctuations.

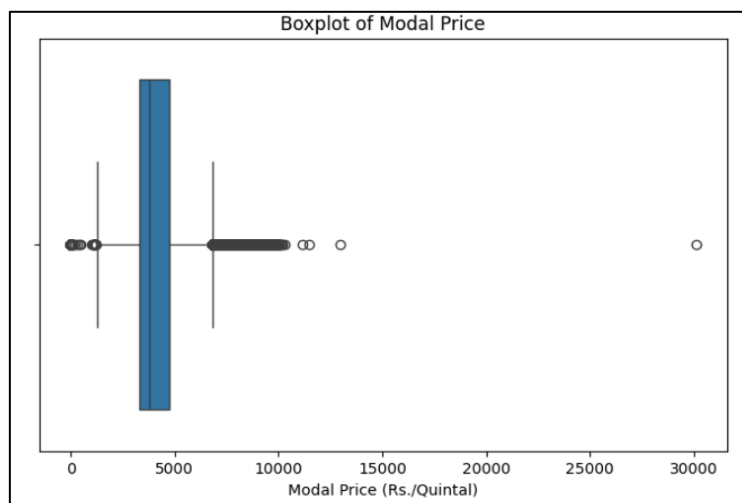


Figure 3: Modal Price Distribution

Through EDA, key patterns and dependencies were uncovered, guiding feature engineering and model selection, ultimately enhancing the predictive accuracy of soybean price forecasting. Figure 3 illustrates the boxplot of the modal price of soyabean, highlighting its overall distribution, spread, and presence of outliers across the dataset. Most modal prices are concentrated between ₹3,000 and ₹6,000 per quintal, indicating the common market price range during the observed period. The plot also reveals a few extreme values extending up to

₹30,000, suggesting occasional price spikes possibly due to market shortages or climatic disruptions.

Figure 4 shows how the average price of soybeans is spread out across India's top 10 districts. This shows how prices change in different areas. Prices are about the same in districts like Ahmednagar, Amaravati, Buldhana, and Latur; the usual average price is around ₹3,500 to ₹5,000 per quintal.

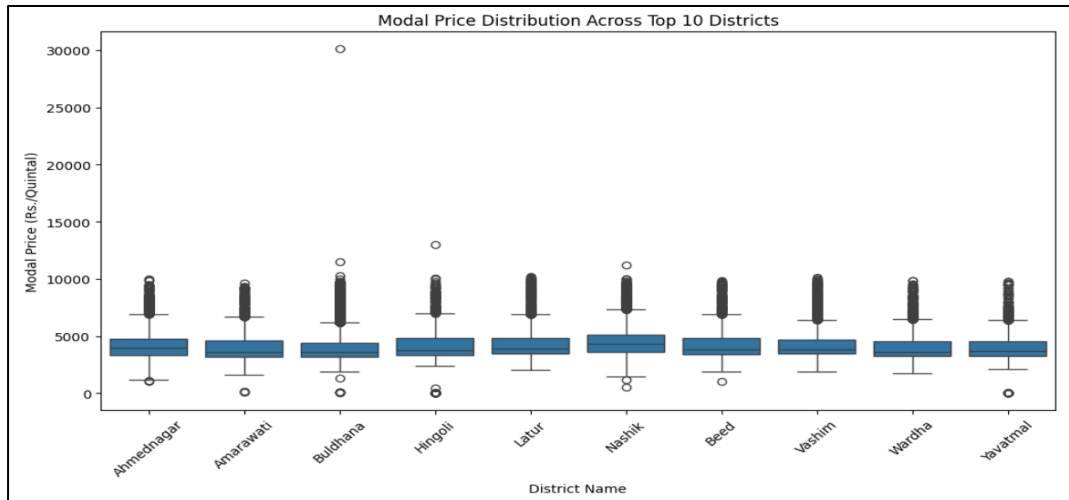


Figure 4: Modal Price Distribution Across Top 10 Districts of India

The boxplots show modest diversity, but some areas have big peaks that show prices sometimes go up because of differences in supply and demand, bad weather, or problems in the market. The fact that prices are mostly the same across districts says that prices are fixed, with only small changes happening in local stores. This comparison helps us understand how area markets work, which is important for making accurate models that predict prices.

Figure 5 shows how the average price (Rs./Quintal) for soyabeans in India is related to the total amount grown (Lakh Tonnes). The scatter plot shows a negative connection, which means that as production goes up, market prices tend to go down.

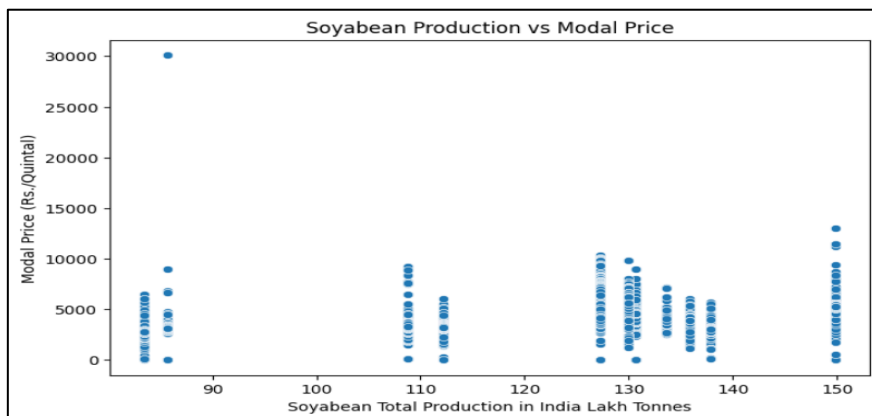


Figure 5: Soyabean Production vs Model Price

This is the basic supply-and-demand relationship. When production goes down, prices tend to go up because there are fewer goods on the market. Some "outliers" have unusual price jumps that may be caused by trade with other countries or weather conditions.

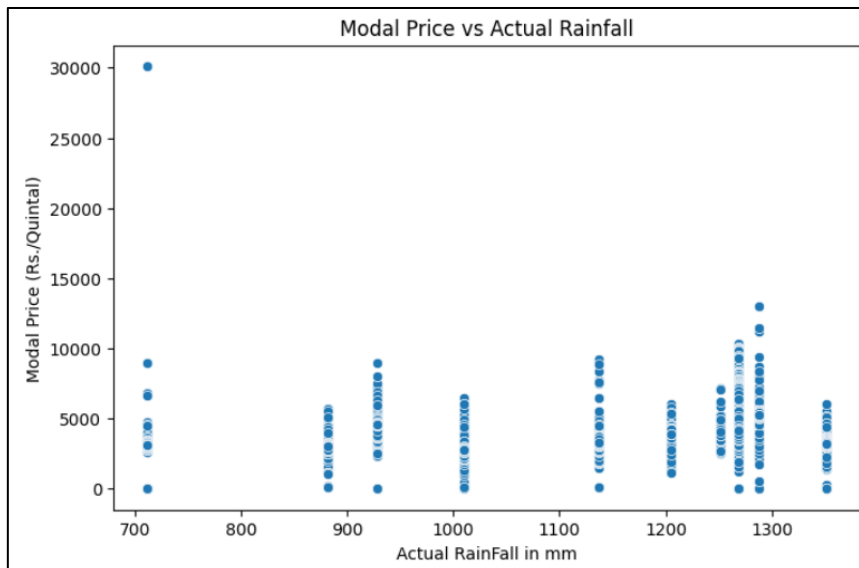


Figure 6: Model Price vs Actual Rainfall

Figure 6 shows how the average price (Rs./Quintal) and real rainfall (mm) change over time and in different places and years. The scatter plot shows that moderate rainfall levels (around 1000–1200 mm) are usually linked to stable soybean prices. On the other hand, prices tend to change a lot when rainfall levels are too low or too high.

3. DATA PRE-PROCESSING

a. Check Missing Values

For categorical variables, the most common group was used to fill in the blanks. When a lot of records were missing, the affected records were thrown out to keep the data's purity. This step made sure that the information was uniform and clean, and that it correctly reflected farming and market situations in the real world so that the model could be trained well.

b. Check Duplicate Values

Duplicate values can make data quality much worse by adding extra work, skewed results, and putting too much emphasis on certain trends during model training. Since the soybean dataset is made up of data from many different farming and trade systems, it was likely that there would be duplicate records, like entries for the same market and date. To fix this, data sorting and key feature comparison were used to find duplicates in characteristics like State Name, District Name, Market Name, and Reported Date. To find examples that overlap, the pandas copied() method and record groups were used. Once they were found, the duplicate rows were taken out, but one correct record was kept to keep important data safe. This process trimmed down the dataset, got rid of duplicates, and stopped overfitting that can happen when data trends are repeated too many times. Making sure that each sample was

unique was important for fair model learning. This way, the prediction algorithms could see real trends in changes in soybean prices instead of just artefacts caused by duplicate data.

c. Label Encoding

Label encoding was used to turn category characteristics into numbers that machine learning systems could understand. There were many non-numeric entries in the information, like State Name, District Name, Market Name, Variety, and Group, that had to be transformed before they could be processed by computers. LabelEncoder from the scikit-learn library was used to give each unique category value its own unique numeric code. This kept any ordinal connections that were present. This change made it easier for algorithms like Random Forest, Gradient Boosting, and XGBoost to understand category data. It was carefully checked to make sure that encoding did not add any unwanted ranking among traits that were not in any particular order. Encoding helped optimise memory and speed up model calculations for factors with a high cardinality. The encoding process was also written down to make sure that it was always the same during the model rollout and forecast stages. Label encoding improved feature compatibility and the general learning ability of the price forecast models by turning category data into a structure that computers can read.

d. Data Normalization

Data normalisation was used to make the range of number features more regular. This made the model more stable and helped it converge during training. Since the soybean dataset had factors with different units, like Price (Rs./Quintal), Area under Cultivation (Lakh Hectare), and Rainfall (mm), unscaled values could have a big effect on learning methods. Depending on the needs of the model, methods like Min-Max Normalisation and Standardisation (Z-score scaling) were used. Normalisation changed all continuous features to similar scales, which are usually between 0 and 1. This made sure that no single feature controlled the optimisation process. This step before processing also sped up training for gradient-based algorithms and made deep learning models more stable numerically. Normalisation also lessened the distortion caused by very high or very low numbers and made it easier to see patterns in both weather and market data. Overall, data normalisation was a key part of making the soybean price predicting system better at making predictions and being able to apply those predictions to other situations.

4. FEATURE EXTRACTION

Feature extraction is a key part of making the dataset more accurate by finding new characteristics that better show the trends that lie beneath the data. Price Range is the first derived feature that measures the difference between the highest and lowest prices of soybeans that have been recorded on the market. This helps show how prices change in certain areas and over time, showing instability and the relationship between supply and demand. The second feature is Rainfall Difference, which shows how much real rainfall is different from regular rainfall. It shows how changes in the weather affect food output, the quality of the yield, and, in the end, market costs. The third factor, Arrivals per Area, shows

how the total amount of soybeans that come into markets is related to the land that is being farmed. It shows how productive different areas are and how much stock there is on the market.

a. Price Range

$$X[Price'_{Range} = X[Max Price \left(\frac{Rs}{Quintal}\right)' - X[Min Price \left(\frac{Rs.}{Quintal}\right)']$$

b. Rainfall_difference

$$X[Rainfall'_{Difference} = X[Actual RainFall in mm' - X[Normal RainFall in mm']$$

c. Arrivals_per_area

$$X[Arrivals'_{per'Area} = \frac{X[Arrivals (Tonnes)']}{X[Area Under Cultivation for Spyabean (Lakh Hectare)in maharashtra']}$$

d. Handle_infinity_values

$$X[Arrivals'_{per'Area} = X[Arrivals'_{per'Area} .replace([np.inf, -np.inf], 0)$$

5. TRAIN TEST SPLIT

To make sure the model could be tested and applied to new situations, the soybean dataset was split into training, validation, and test groups with an 80–20 split ratio. The Modal Price was used as the goal variable (class name). The 108,893 samples in the training set were used to teach the machine learning and deep learning models how to find complex connections between things like rainfall, landings, and farming area. The validation set, which has 36,298 samples, was used to fine-tune hyperparameters, keep an eye on model performance, and stop training from becoming too perfect.

6. MACHINE LEARNING MODELS

a. Random Forest

Random Forest is used to predict soybean prices because it can handle the complicated relationships between weather, market, and output factors. It can handle noise and missing data, which makes it a good starting point for looking at big farming datasets and making stable, accurate price predictions.

- Step 1: Bootstrap Sampling

From the training dataset $D = \{(x_i, y_i)\}$ create B random subsets using sampling with replacement.

- Step 2: Tree Construction

For each subset, grow a decision tree by selecting a random subset of features $F \subseteq \{f_1, f_2, \dots, f_m\}$.

- Step 3: Node Splitting

At each node, choose the best split feature f_j that minimizes impurity (e.g., Gini Index):

$$G = 1 - \sum (p_k)^2$$

- Step 4: Final Prediction

For regression:

$$\hat{y} = \left(\frac{1}{B}\right) \sum T b(x)$$

For classification:

$$\hat{y} = \text{mode}(T1(x), T2(x), \dots, TB(x))$$

b. Ensemble Model

The Ensemble Model takes guesses from several algorithms and puts them all together to make the general performance better and lower the flaws of each model. It takes advantage of the benefits of methods like Random Forest, Gradient Boosting, and XGBoost by combining them. This mixed method makes generalisation, stability, and accuracy better. The ensemble model combines different learning patterns to make predictions about soybean prices that are fair and take into account both short-term market trends and long-term climate or production changes that affect price changes.

- Step 1: Model Selection

Choose base learners (Random Forest, Gradient Boosting,

XGBoost): $M1, M2, \dots, Mn$.

- Step 2: Training

Train each model on the same dataset $D = \{(x_i, y_i)\}$.

- Step 3: Prediction Generation

Each model produces its prediction $\hat{y}_j = M_j(x)$.

c. Gradient Boosting

Iteratively building models, where each new tree fixes the mistakes of the ones that came before it, is what gradient boosting does. It uses gradient descent to improve a loss function, which makes predictions very accurate. Gradient Boosting is a good way to model how weather, market arrivals, and supply data are not linearly dependent on each other when trying to predict soybean prices. Because it can be fine-tuned and doesn't get too comfortable, it's perfect for picking up on small changes in how farming markets work.

Initialize Model

Start with a base prediction using the mean of target values:

$$F_0(x) = \operatorname{argmin}_{\gamma} \sum L(y_i, \gamma)$$

Compute Residuals

For each iteration m , compute pseudo-residuals:

$$r_{im} = -\frac{\partial L(y_i, F_{\{m-1\}}(x_i))}{\partial F_{\{m-1\}}(x_i)}$$

Fit Weak Learner

Train a regression tree $h_m(x)$ to predict the residuals r_{im} .

Update Model

Compute optimal multiplier γ_m :

$$\gamma_m = \operatorname{argmin}_{\gamma} \sum L(y_i, F_{\{m-1\}}(x_i) + \gamma h_m(x_i))$$

d. XGBoost

Extreme Gradient Boosting, or XGBoost, is a more advanced version of Gradient Boosting that focusses on regularisation and making computations more efficient. It makes training go faster, improves accuracy, and handles lost or uneven data better than anything else. XGBoost takes into account complex relationships between weather, economy, and farming factors when predicting the price of soybeans. It works well with big datasets because it can be scaled up and down easily. It also makes very accurate and reliable predictions that are important for making decisions in market research and farming economics.

Step 1: Model Initialization

Start with initial prediction $\hat{y}_i^0 = \text{average}(y)$.

Step 2: Additive Learning

At each iteration t , add a new tree $f_t(x)$ to minimize:

$$Obj^t = \sum l(y_i, \hat{y}_i^{t-1} + f_t(x_i)) + \Omega(f_t)$$

$$\text{where } \Omega(f_t) = \gamma T + 0.5\lambda ||w||^2$$

Step 3: Second-Order Approximation

Compute gradient (g_i) and hessian (h_i):

$$g_i = \partial l(y_i, \hat{y}_i)$$

$$h_i = \partial^2 l(y_i, \hat{y}_i)$$

Step 4: Leaf Weight Calculation

Compute optimal leaf weights:

$$w_j * = - \frac{(\sum g_i)}{(\sum h_i + \lambda)}$$

Step 5: Prediction Update

Update final prediction:

$$\hat{y}_i^t = \hat{y}_i^{t-1} + \eta * f_t(x_i)$$

where η is the learning rate.

7. CALLBACKS

a. Early Stopping

Early Stopping is a strong callback method used to stop overfitting and improve generalisation during model training. During training, it checks a certain performance measure, usually accuracy or validation loss, and stops when there is no more growth after a certain number of epochs, which is also known as patient. The model doesn't learn noise or trends that aren't important in the information because it stops training early. Early Stopping was used during the deep learning training stages of the soybean price forecast model to improve performance while lowering the cost of computing. This method made sure that the model kept its best weights, which were the ones with the lowest validation loss. This made the predictions more accurate and stable when it was introduced to new weather and farming data.

b. reduce Learning rate

The Reduce Learning Rate on Plateau callback changes the learning rate on the fly while the model is being trained to speed up convergence and keep it from stopping. If a watched measure like validation loss stops getting better after a certain number of epochs, the learning rate is slowed down by a set amount. This lets the optimiser make smaller, more accurate changes. This method keeps the model from going too far past the best minimum in the loss environment. This reply improved the performance of the deep learning model in the system that predicted the price of soybeans by making optimisation smoother and improving accuracy. Lowering the learning rate at the right time stabilised training, kept loss from going up and down too much, and improved the general efficiency of convergence for difficult tasks like predicting farm prices.

8. DEEP LEARNING MODELS

a. Dense Neural Network (DNN)

A Dense Neural Network (DNN) is a type of deep learning design in which every neurone in one layer is linked to every neurone in the next layer. It shows complicated, nonlinear connections between many input factors, like prices, landings, and weather. In the model for predicting the price of soybeans, the DNN works with normalised numerical data by using several hidden layers and activation functions like ReLU. Through backpropagation and optimisation, the model learns to describe features in a structured way while minimising loss.

Because it can generalise trends, it can be used to make accurate predictions about crop prices.

Input Transformation

Each input feature (e.g., rainfall, arrivals, production) is passed through the first layer.

$$z(1) = W(1) * x + b(1)$$

Hidden Layer Activation

Non-linear activation (e.g., ReLU) is applied to introduce complexity.

$$a(1) = ReLU(z(1))$$

Forward Propagation

The output of each layer becomes input to the next layer.

$$a(l) = \sigma(W(l) * a(l - 1) + b(l))$$

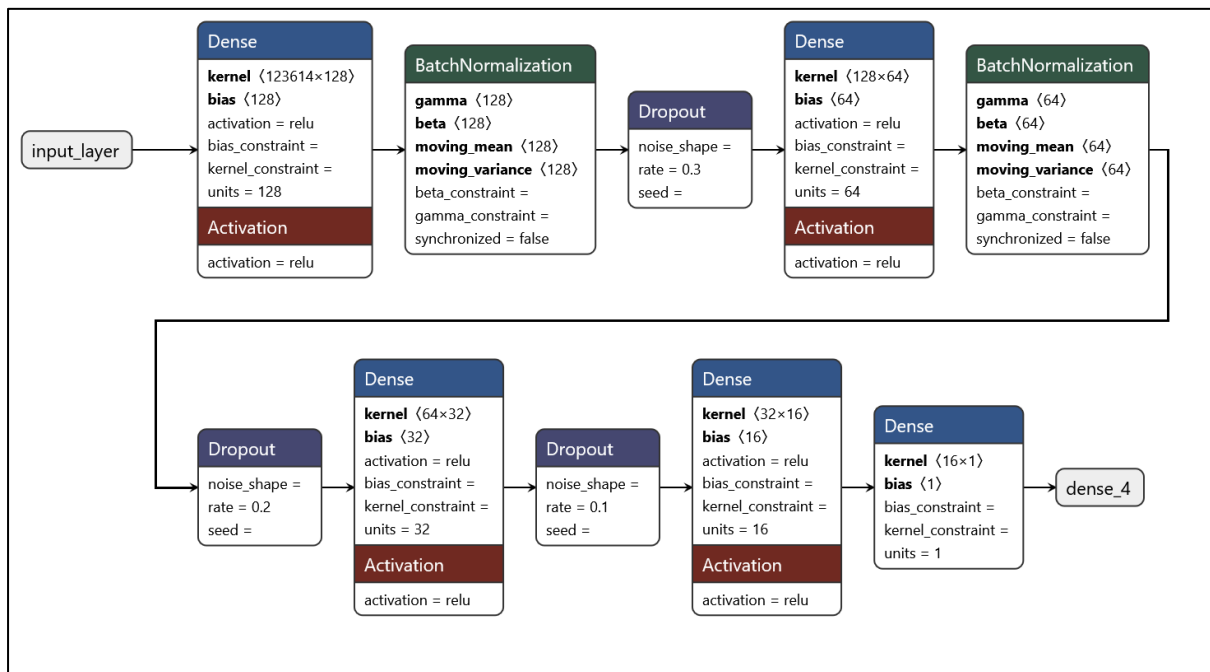


Figure 7: Architecture of the Dense Neural Network (DNN) Model

b. Attention Model

The Attention Model improves deep learning by letting the network focus on the most important input traits while making predictions about what will happen. Instead of giving all inputs the same amount of weight, it figures out attention weights that put important factors like weather variation, landings, or output levels at the top of the list.

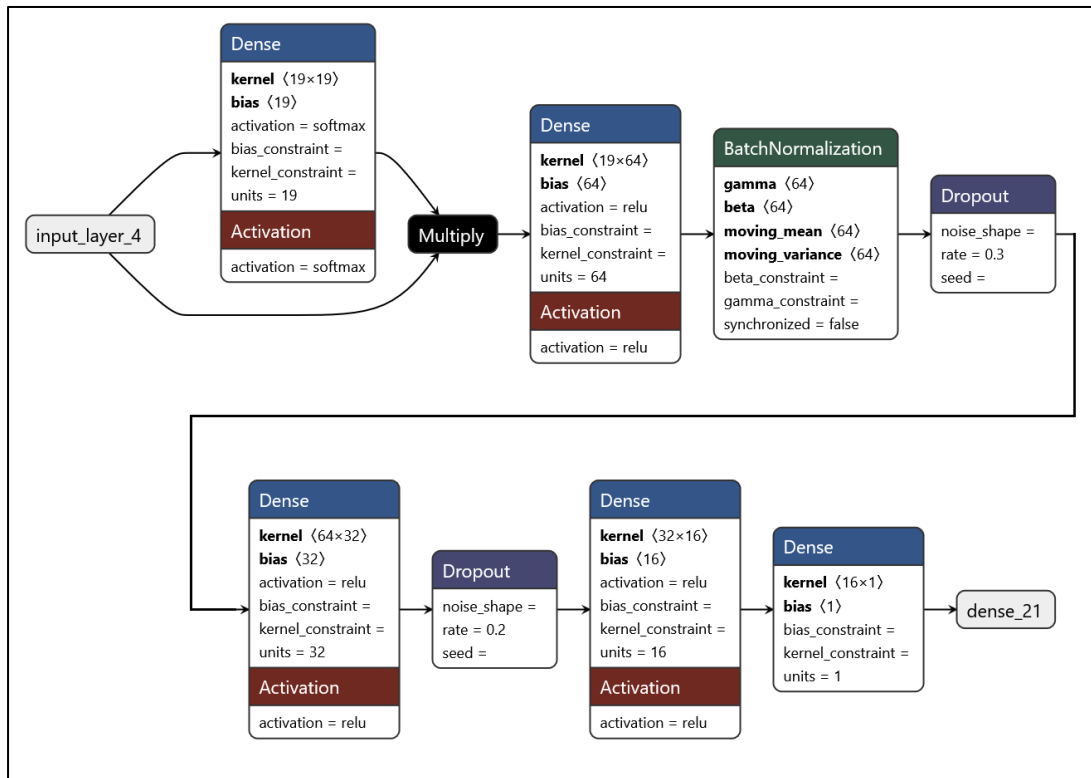


Figure 8: Attention-Based Neural Network Model

The attention mechanism makes soybean price predictions easier to understand and more accurate by drawing attention to important factors that affect market changes. Figure 8 shows the model focusing on key features, enhancing soybean price prediction accuracy. This focused learning method helps the model easily record complex relationships, which leads to more accurate and situation-aware price predictions.

Step 1: Compute Query, Key, and Value Vectors

Each input vector x_i is converted into three representations.

$$Q = W_Q * x$$

$$K = W_K * x$$

$$V = W_V * x$$

Step 2: Calculate Attention Scores

Compute similarity between Query and Key.

$$score(Q, K) = \frac{(Q * K^T)}{\sqrt{d_k}}$$

Step 3: Softmax Normalization

Convert scores into attention weights.

$$\alpha_i = \frac{\exp(\text{score}(Q, K_i))}{\sum_j \exp(\text{score}(Q, K_j))}$$

Step 4: Weighted Feature Aggregation

Generate context vector as a weighted sum of Value vectors.

$$c = \sum_i \alpha_i * V_i$$

c. Regularized Network

To stop overfitting and improve generalisation, the Regularised Network uses methods like L1 (Lasso), L2 (Ridge), or Dropout regularisation. By adding a punishment term to the loss function, it stops models from being too complicated and weight values from being too high. Figure 9 shows regularization improving model stability and prediction accuracy. This network makes sure that when predicting the price of soybeans, the model looks at important links instead of noise or patterns that are used more than once.

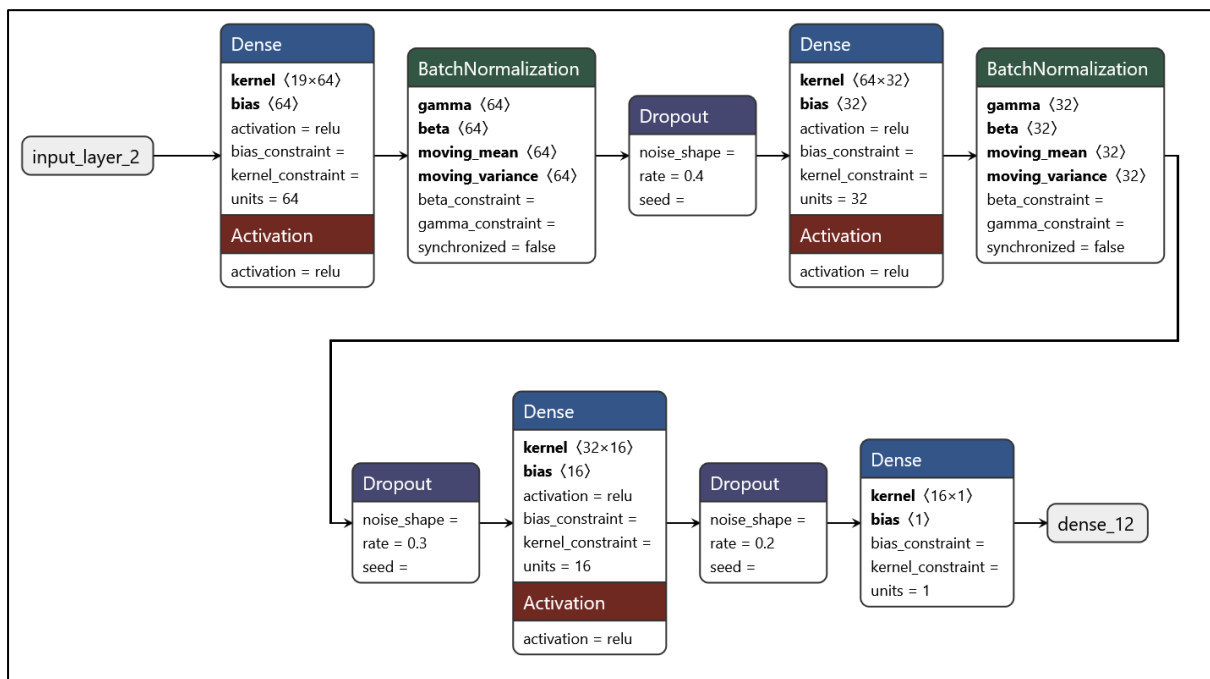


Figure 9: Architecture of the Regularized Neural Network Model

Regularisation keeps predictions accurate, makes learning more stable, and makes it more resistant to changes in the data. This makes it a good method for looking at multidimensional farming and climate datasets with a lot of feature overlap.

Step 1: Model Prediction

The network predicts price based on weights and biases.

$$\hat{y} = f(W * x + b)$$

Step 2: Define Regularized Loss Function

Add a penalty term to control large weights.

$$L = \left(\frac{1}{N}\right) * \Sigma(y_i - \hat{y}_i)^2 + \lambda * R(W)$$

Step 3: Regularization Types

$$L1 (Lasso): R(W) = \Sigma|W_i|$$

$$L2 (Ridge): R(W) = \Sigma W_i^2$$

d. Wide and Deep Network (AgroWDN)

The AgroWDN is a mixed design that takes the best parts of linear models (the wide component) and deep neural networks (the deep component). The wide part remembers past interactions and interactions between features, and the deep part generalises complicated, nonlinear connections. AgroWDN uses farming, environmental, and economic factors to predict soybean prices. It does this by taking into account both short-term trends and long-term relationships.

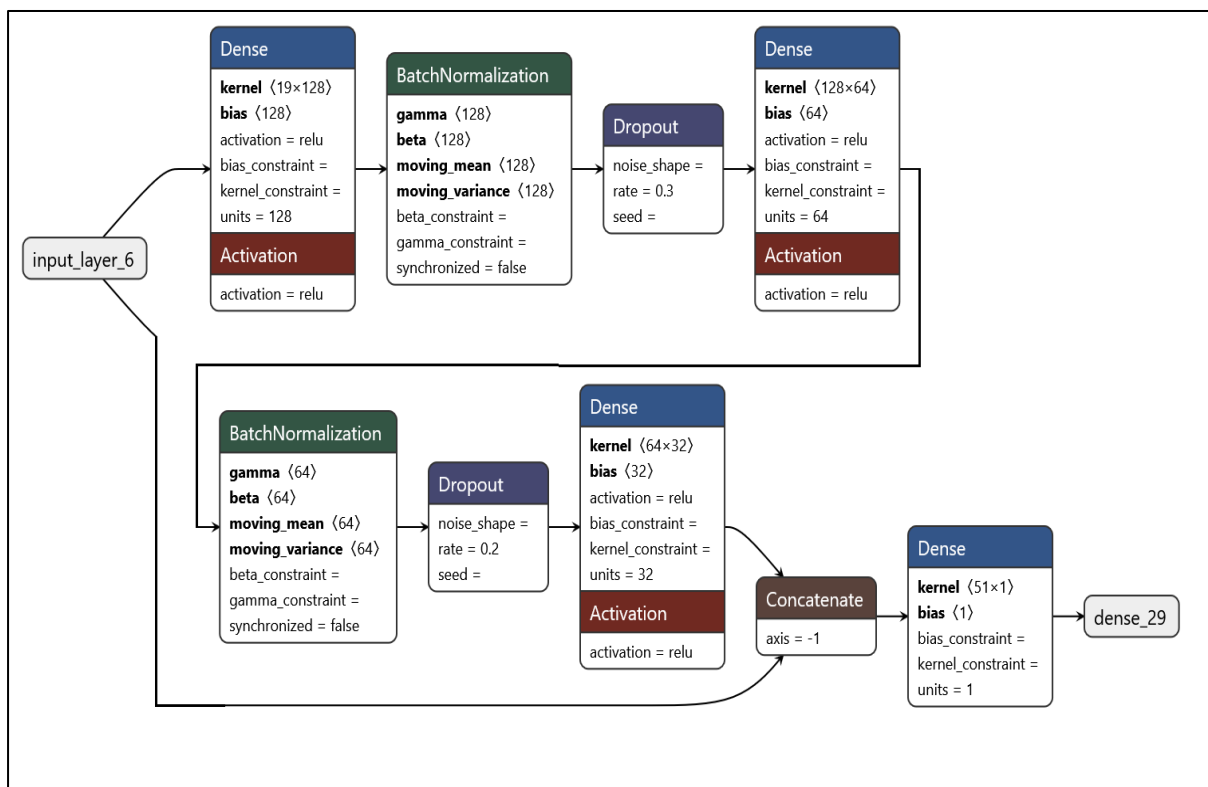


Figure 10: Wide and Deep Network (AgroWDN) Model

This balanced learning method makes the model easier to understand, more scalable, and more accurate. Figure 10 shows AgroWDN combining linear and nonlinear feature learning. It also does a better job of predicting dynamic and multidimensional farm price data than traditional machine learning models.

Step 1: Wide Component (Linear Part)

Memorizes interactions and cross-features.

$$y_{wide} = w^T * x + b$$

Step 2: Deep Component (Nonlinear Part)

Learns complex feature representations.

$$y_{deep} = f(W2 * f(W1 * x + b1) + b2)$$

Step 3: Combined Prediction

Merge wide and deep outputs for final prediction.

$$y_{pred} = \sigma(y_{wide} + y_{deep})$$

Step 4: Optimization Objective

Minimize loss (e.g., MSE) during training.

$$L = \left(\frac{1}{N}\right) * \sum (y_i - y_{pred_i})^2$$

IV. RESULT ANALYSIS

A. Accuracy and Loss Comparison Graphs

Loss and Mean Absolute Error (MAE) over time are shown in Figure 11 for the Wide and Deep Network (AgroWDN) during training and evaluation. Both measures drop sharply in the first few rounds but level off after a few epochs, which means that learning and convergence are going well.

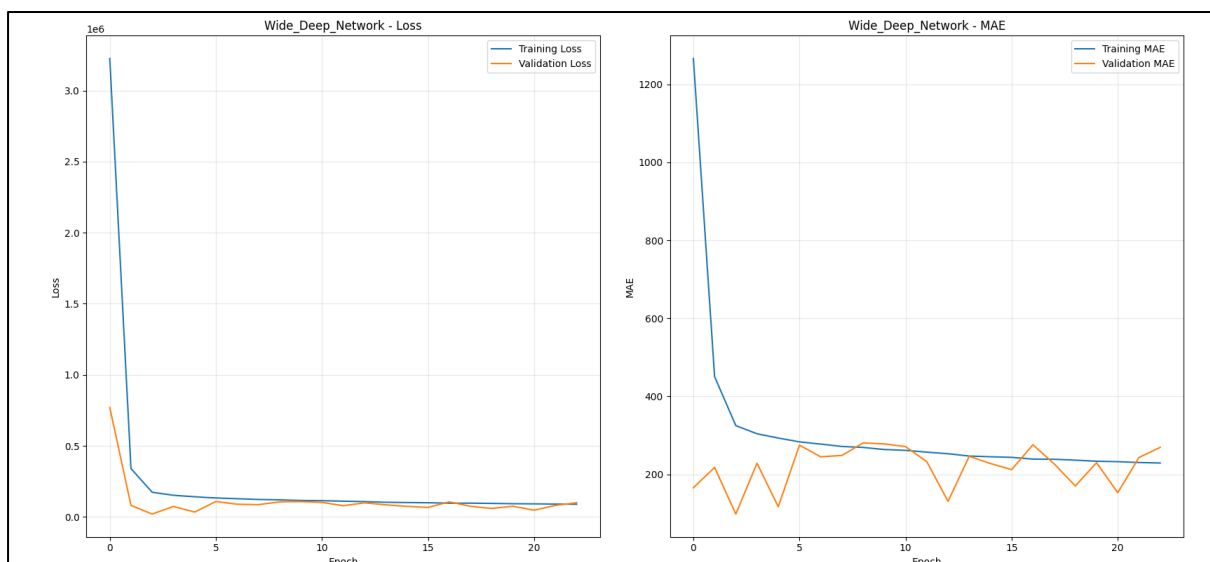


Figure 11: Wide and Deep Network (AgroWDN) Loss and MAE Comparison Graph

The fact that the training and validation curves are very close to each other shows that there isn't much overfitting and that the model can generalise well. This shows that the AgroWDN model is reliable and accurate at predicting soybean prices from large farming datasets.

B. Price Prediction Comparison Graph

Figure 12 shows a comparison of the price distributions of real and projected soybean prices using the AgroWDN model's Wide and Deep Network (WDN). Prices are mostly between ₹3,000 and ₹6,000, which is the main market frequency. The overlapped ranges show that the model can make accurate and reliable predictions. This proves that it is good at learning how crop prices change over time and predicting soybean market trends.

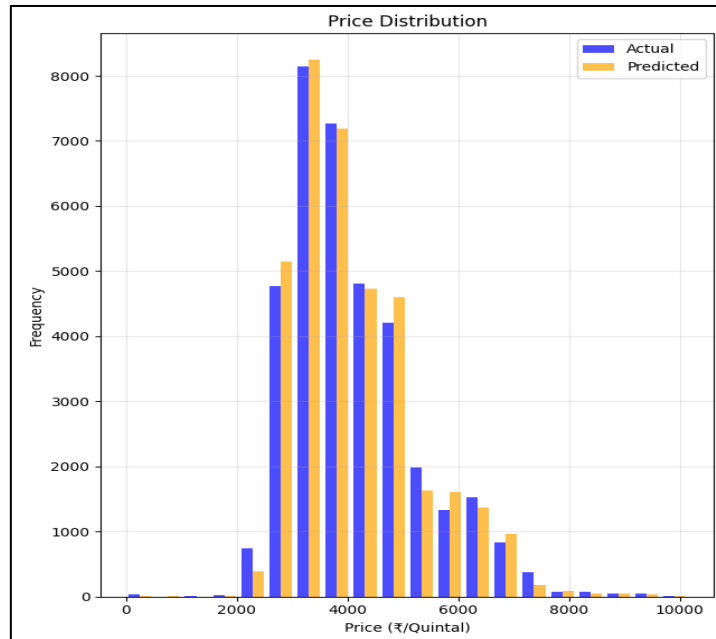


Figure 12: Price Prediction Using Wide and Deep Network (AgroWDN) Graph

C. Actual and predicted prices

The link between real and projected soybean prices made by the Wide and Deep Network (AgroWDN) model is shown in Figure 13. Each data point is a recording of a price, and the red straight line shows the ideal 1:1 connection.

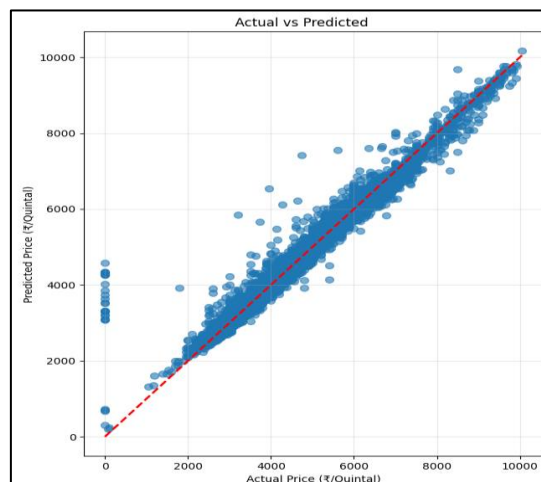


Figure 13: Actual vs Predicted Soybean Prices Using AgroWDN Model

Most of the points are close to this line, which shows that the projected values are very close to the real market prices. The close grouping along the vertical means that the model is very accurate and there is little variation. This strong linear relationship shows how well the AgroWDN model can learn complicated market behaviour and make accurate soybean price predictions when farming conditions change.

D. Comparative Analysis Graph

Table 3 shows a summary of the performance rating measures for different machine learning and deep learning models created to predict the price of soybeans. XGBoost, Gradient Boosting, Ensemble, Random Forest, Attention Model, Regularised Model, Dense Neural Network, and the Wide and Deep Network (AgroWDN) are some of the models that were compared. Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), Coefficient of Determination (R^2), and Mean Squared Error (MSE) are all ways to measure performance. The AgroWDN model has the best accuracy and usefulness in predicting soybean prices across a wide range of farming and weather factors, as shown by its low RMSE and MAE values and high R^2 value (98.96%).

Table 3: Comparative Performance Analysis of Machine Learning and Deep Learning Models

Model	RMSE	MAE	R² (%)	MSE
XGBoost	122.933687	59.565613	98.8958	15112.691343
Gradient_Boosting	125.627610	61.579617	98.8469	15782.296317
Ensemble	149.639837	70.939589	98.3640	22392.080671
Random_Forest	153.509256	67.810029	98.2783	23565.091527
Attention_Model	177.101545	82.201901	97.7084	31364.957118
Regularized_Model	191.561418	100.036928	97.3189	36695.776888
Dense_Neural_Network	194.133548	100.075968	97.2464	37687.834524
Wide and Deep Network (AgroWDN)	111.434823	57.154470	98.9651	14879.954517

This picture (Figure 14) shows how well different machine learning and deep learning models used to predict soybean prices did using three different performance metrics: Root Mean Squared Error (RMSE), Mean Absolute Error (MAE), and Coefficient of Determination (R^2). The Wide and Deep Network (AgroWDN) has the lowest RMSE and MAE numbers and the best R^2 score, which means it can make more accurate predictions and apply its model to more situations. Other models, such as XGBoost and Gradient Boosting,

also do well, but not as well as AgroWDN. This shows that AgroWDN is good at finding nonlinear patterns in large farming datasets.

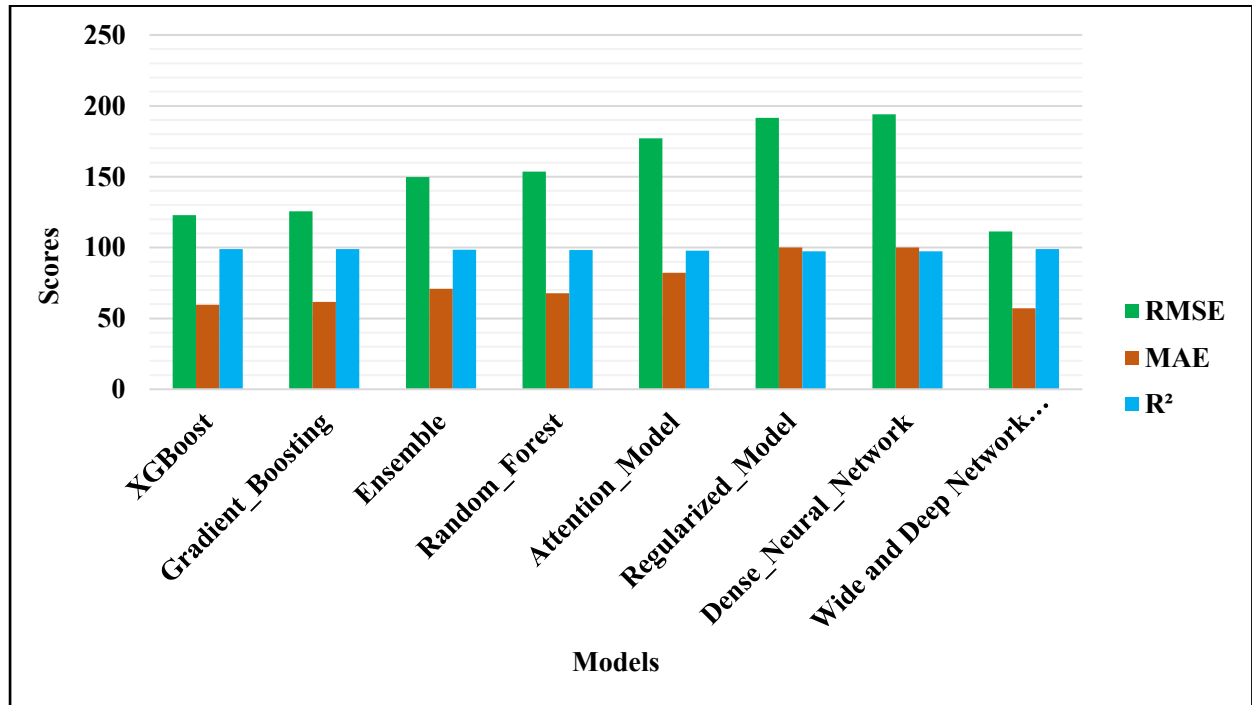


Figure 14: Comparison of Model Performance Parameters (RMSE, MAE, and R²)

Figure 15 shows the Mean Squared Error (MSE) comparison between different machine learning and deep learning models used to guess the price of soybeans. When MSE numbers are low, it means that the model is more accurate. It was the Wide and Deep Network (AgroWDN) and XGBoost models that got the lowest MSE, which shows how well they can predict and reduce errors.

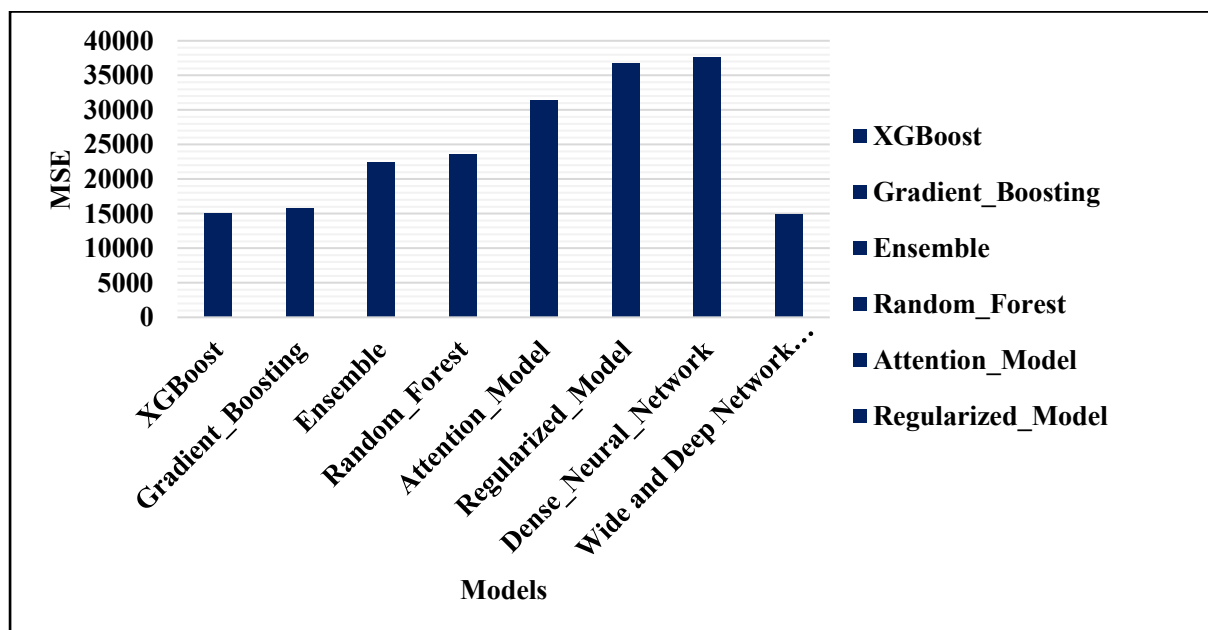


Figure 15: Mean Squared Error (MSE) Comparison Graph of Different Models

Models like the Regularised and Dense Neural Networks, on the other hand, have higher MSE, which means that the differences are bigger. This study proves that AgroWDN is good at predicting prices for farm markets in a way that is accurate, consistent, and based on data.

E. Soyabean Price Prediction

In Table 4, demonstrate that by comparison between the real soybean prices and the prices that the Wide and Deep Network (AgroWDN) model said they would be. The index number, real market prices, model-predicted prices, and forecast mistakes are all shown in the table. The model's high prediction accuracy and stability are shown by the small gaps between real and expected numbers.

Table 4: Comparison of Actual and Predicted Soybean Prices Using AgroWDN Model

Index	Actual	Predicted	Error
1	₹2670.00	₹2642.71	₹27.29
2	₹2300.00	₹2318.10	₹-18.10
3	₹4100.00	₹4124.81	₹-24.81
4	₹4970.00	₹4920.54	₹49.46
5	₹4381.00	₹4339.16	₹41.84
6	₹4700.00	₹4786.99	₹-86.99
7	₹4550.00	₹4488.30	₹61.70
8	₹5120.00	₹4958.61	₹161.39
9	₹3700.00	₹3977.17	₹2.83
10	₹3650.00	₹3546.74	₹103.26

The results show that the AgroWDN model can correctly predict soybean prices and understand how complicated markets work. This shows that it is a useful tool for predicting crop prices and helping people make decisions.

V. Conclusion

The study called “Developing Machine Learning Model for Accurate Prediction of Soybean Market Prices Using Multidisciplinary Agricultural Data Sources” shows that advanced machine learning and deep learning techniques can be used to accurately predict soybean prices. By combining different sets of data, such as those about climate, agriculture, and the economy, the study created a complete framework for making predictions that accurately depicts the complicated, nonlinear connections that cause price changes in the market. Several algorithms were tested, such as Random Forest, Gradient Boosting, XGBoost,

Ensemble models, and deep learning architectures like the Dense Neural Network (DNN), Attention Model, Regularised Network, and the hybrid Wide and Deep Network (AgroWDN). A lot of data preprocessing, feature engineering, and model training were done. The AgroWDN model did the best, with a R^2 of 98.96%, RMSE of 111.43, and MAE of 57.15, showing that it was more accurate and could be used in more situations. The best thing about the model is that it can learn from both linear and nonlinear correlations between many farming factors, such as weather patterns, crop area, trade data, and market landings. It can do this by combining wide (memorisation) and deep (generalisation) learning. The model's stability and reliability were also confirmed by visual studies like loss-MAE convergence and real versus expected comparisons.

References

- [1] Qin, P.; Wang, T.; Luo, Y. A review on plant-based proteins from soybean: Health benefits and soy product development. *J. Agric. Food Res.* 2022, 7, 100265.
- [2] Bazzana, D.; Foltz, J.; Zhang, Y. Impact of climate smart agriculture on food security: An agent-based analysis. *Food Policy* 2022, 111, 102304.
- [3] Wadas, W.; Kondraciuk, T. The Role of Foliar-Applied Silicon in Improving the Growth and Productivity of Early Potatoes. *Agriculture* 2025, 15, 556.
- [4] Osinga, S.A.; Paudel, D.; Mouzakitis, S.A.; Athanasiadis, I.N. Big data in agriculture: Between opportunity and solution. *Agric. Syst.* 2022, 195, 103298.
- [5] Feng, H.; Fan, Y.; Yue, J.; Bian, M.; Liu, Y.; Chen, R.; Ma, Y.; Fan, J.; Yang, G.; Zhao, C. Estimation of potato above-ground biomass based on the VGC-AGB model and deep learning. *Comput. Electron. Agric.* 2025, 232, 110122.
- [6] L. Thakre, M. Daware, A. Mohite, A. Sakhare and P. Khobragade, "Smart Farming: Enhancing Crop Recommendation and Price Prediction with Advanced Machine Learning," 2025 6th International Conference for Emerging Technology (INCET), BELGAUM, India, 2025, pp. 1-6, doi: 10.1109/INCET64471.2025.11140186.
- [7] Triantakostas, D.; Karakostas, A. Soil Organic Carbon Monitoring and Modelling via Machine Learning Methods Using Soil and Remote Sensing Data. *Agriculture* 2025, 15, 910.
- [8] Bregaglio, S.; Ginaldi, F.; Raparelli, E.; Fila, G.; Bajocco, S. Improving crop yield prediction accuracy by embedding phenological heterogeneity into model parameter sets. *Agric. Syst.* 2023, 209, 103666.
- [9] Arshad, S.; Kazmi, J.H.; Javed, M.G.; Mohammed, S. Applicability of machine learning techniques in predicting wheat yield based on remote sensing and climate data in Pakistan, South Asia. *Eur. J. Agron.* 2023, 147, 126837.
- [10] Lu, C.; Leng, G.; Liao, X.; Tu, H.; Qiu, J.; Li, J.; Huang, S.; Peng, J. In-season maize yield prediction in Northeast China: The phase-dependent benefits of assimilating climate forecast and satellite observations. *Agric. For. Meteorol.* 2024, 358, 110242.
- [11] Li, Y.; Liu, X.; Zhang, X.; Gu, X.; Yu, L.; Cai, H.; Peng, X. Using solar-induced chlorophyll fluorescence to predict winter wheat actual evapotranspiration through

- machine learning and deep learning methods. *Agric. Water Manag.* 2025, 309, 109322.
- [12] P. Khobragade, P. K. Dhankar, A. Titarmare, M. Dhone, S. Thakur and P. Saraf, "Quantum-Enhanced AI Robotics for Sustainable Agriculture: Pioneering Autonomous Systems in Precision Farming," 2024 International Conference on Artificial Intelligence and Quantum Computation-Based Sensor Application (ICAIQSA), Nagpur, India, 2024, pp. 1-7, doi:10.1109/ICAIQSA64000.2024.10882412.
- [13] Lu, J.; Li, J.; Fu, H.; Tang, X.; Liu, Z.; Chen, H.; Sun, Y.; Ning, X. Deep Learning for Multi-Source Data-Driven Crop Yield Prediction in Northeast China. *Agriculture* 2024, 14, 794.
- [14] Khan, S.N.; Li, D.; Maimaitijiang, M. Using gross primary production data and deep transfer learning for crop yield prediction in the US Corn Belt. *Int. J. Appl. Earth Obs. Geoinf.* 2024, 131, 103965.
- [15] Du, J.; Zhang, Y.; Wang, P.; Tansey, K.; Liu, J.; Zhang, S. Enhancing Winter Wheat Yield Estimation With a CNN-Transformer Hybrid Framework Utilizing Multiple Remotely Sensed Parameters. *IEEE Trans. Geosci. Remote Sens.* 2025, 63, 4405213.
- [16] Lu, J.; Li, J.; Fu, H.; Zou, W.; Kang, J.; Yu, H.; Lin, X. Estimation of rice yield using multi-source remote sensing data combined with crop growth model and deep learning algorithm. *Agric. For. Meteorol.* 2025, 370, 110600.
- [17] Wang, W.; Deng, X.; Yue, H. Black soil conservation will boost China's grain supply and reduce agricultural greenhouse gas emissions in the future. *Environ. Impact Assess. Rev.* 2024, 106, 107482.
- [18] Xin, M.; Zhang, Z.; Han, Y.; Feng, L.; Lei, Y.; Li, X.; Wu, F.; Wang, J.; Wang, Z.; Li, Y. Soybean phenological changes in response to climate warming in three northeastern provinces of China. *Field Crops Res.* 2023, 302, 109082.
- [19] Gueymard, C.A. A review of validation methodologies and statistical performance indicators for modeled solar radiation data: Towards a better bankability of solar projects. *Renew. Sustain. Energy Rev.* 2014, 39, 1024–1034.
- [20] A. Chavan, A. Desai, and K. S. Oza, "Sugarcane Crop Disease Detection", *Int Journal Adv Comp Theory Engg*, vol. 14, no. 1, pp. 20–27, Apr. 2025.
- [21] Ramandanis, D.; Xinogalos, S. Investigating the Support Provided by Chatbots to Educational Institutions and Their Students: A Systematic Literature Review. *Multimodal Technol. Interact.* 2023, 7, 103.
- [22] Karger, E.; Kureljusić, M. Using Artificial Intelligence for Drug Discovery: A Bibliometric Study and Future Research Agenda. *Pharmaceuticals* 2022, 15, 1492.
- [23] Lundberg, L.; Boldt, M.; Borg, A.; Grahn, H. Bibliometric Mining of Research Trends in Machine Learning. *AI* 2024, 5, 208–236.
- [24] Yadav, A. A Comparative Study of Time Series, Machine Learning, and Deep Learning Models for Forecasting Global Price of Wheat. *Oper. Res. Forum* 2024, 5, 113.

- [25] Kumar, R. Predicting Wheat Futures Prices in India. *Asia-Pac. Financ. Mark.* 2021, 28, 121–140.
- [26] M. A. Y. B. Prof. Sharif Shaikh, “The structure Investigation of Rectangular Underneath and Ground-level Water Containers Using Advanced Simulation Tools”, *IJRAET*, vol. 14, no. 2, pp. 19–28, Aug. 2025.
- [27] Dewi, C.; Prasatya, G.S.K.; Christanto, H.J.; Widiarto, S.O.B.; Dai, G. Modified Random Forest Regression Model for Predicting Wholesale Rice Prices. *J. Theor. Appl. Inf. Technol.* 2023, 101, 7749–7759.