



Weather based modeling and pre-harvest forecasting of crop yield using statistical and hybrid models

AKHILESH KUMAR GUPTA^{1*}, KADER ALI SARKAR² and DEBASIS BHATTACHARYA²

¹Department of Agricultural Statistics, College of Agriculture, OUAT, Bhubaneswar, Odisha

²Department of Agricultural Statistics, Institute of Agriculture, Visva-Bharati, Sriniketan (W.B.)

(Received 23 November 2024, Accepted 07 August 2025)

*Corresponding author's email: akgupta@ouat.ac.in

सार – भारत में कृषि और आर्थिक नीतिगत निर्णयों को निर्धारित करने में फसल उपज मॉडलिंग और पूर्वानुमान एक आवश्यक कदम रहा है। वितरण, मूल्य, निर्यात-आयात, भंडारण और अन्य मुद्दों पर योजना और नीतिगत निर्णय इस पर महत्वपूर्ण रूप से निर्भर करते हैं। इस अध्ययन का उद्देश्य फसल की पैदावार और कई मौसम संबंधी कारकों के बीच संबंध का विश्लेषण करके चावल की उपज के लिए एक विश्वसनीय फसल उपज पूर्वानुमान प्रणाली विकसित करना था। मौसम सूचकांक, फसल की पैदावार पर मौसम के साप्ताहिक प्रभावों को समेकित रूप का उपयोग धान की पैदावार पर मौसम कारकों के प्रभाव का अध्ययन करने के लिए किया गया। आंकड़ों की रैखिकता और गैर-रैखिकता के आधार पर कई सांख्यिकीय और न्यूरल नेटवर्क मॉडल विकसित किए गए हैं। सांख्यिकीय मॉडलों के परिणामों से पता चला कि आंकड़ों में रैखिक और गैर-रैखिक दोनों पैटर्न मौजूद थे और वर्षा, न्यूनतम तापमान और टाइम वेरिबल – टी का प्रभाव महत्वपूर्ण था। सटीकता के मामले में न्यूरल नेटवर्क मॉडल, सांख्यिकीय मॉडलों से बेहतर साबित हुए और अध्ययन में पाया गया कि बीरभूम और बर्दवान जिलों में क्रमशः तीन और चार छिपे हुए नोड्स वाले हाईब्रिड मॉडल सबसे उपयुक्त मॉडल थे। विभिन्न नीतिगत निर्णयों के लिए सबसे उपयुक्त मॉडलों का उपयोग करके कटाई से छह से आठ सप्ताह पहले धान की पैदावार का विश्वसनीय अनुमान प्राप्त किया जा सकता है।

ABSTRACT. Crop yield modeling and forecasting have been an essential step in determining agricultural and economic policy decisions in India. Planning and policy decisions on distribution, price, export-import, storage, and other issues are critically dependent on it. This study aimed to develop a trustworthy crop yield prediction system for rice yield by analyzing the relationship between crop yield and several weather variables. Weather indices, an assimilation of the weekly weather effects on crop yield, were used to study the impact of weather factors on rice yield. Several statistical and neural network models have been developed based on the linearity and non-linearity pattern of the data. The results of statistical models demonstrated that both linear and non-linear patterns were present in the data and the effects of rainfall, minimum temperature, and time variable t were significant. The neural network models outperformed statistical models in terms of accuracy, and the study found hybrid models with three and four hidden nodes were the best-fit models in the districts of Birbhum and Burdwan, respectively. A reliable rice yield estimate can be obtained six to eight weeks before harvest by using the best-fit models for various policy decisions.

Key words – Hybrid models, Neural networks, Rice, Weather indices, Yield forecasting.

1. Introduction

The modeling of crop yield and the development of forecasting models for crop yield have been a crucial research area in recent years. It is crucial to have accurate crop yield forecasts for various planning and policy decisions relating to storage, procurement, distribution, pricing, marketing, and export-import (Setiya *et al.*, 2023). There have been various approaches for crop yield modeling and forecasting that utilize data on crop biometrical characters, farmers' eye estimates, the input of crop production, agrometeorological variables, remotely

sensed crop reflectance observations, regression based weather variables, and time series analyses. The traditional methods of crop yield prediction may not always provide accurate results and the timeliness and accuracy of conventional agricultural yield forecasting techniques, such as field surveys and statistical regressions, are limited, especially in light of the more unpredictable weather patterns (Singh *et al.*, 2021). Moreover, climate change has made it even more challenging to predict crop yields with a high degree of certainty. Therefore, there is a need for more robust and data-driven approaches to crop yield prediction

considering various factors contributing to the crop yield. This study utilizes the enhanced capabilities of statistical methods like ARIMAX and machine learning models in combination with regression models and time series analyses to develop a pre-harvest crop yield forecast model.

The use of weather-based models for crop yield forecasting has recently been studied, and these models can offer more accurate and timely information to enhance decision-making throughout the agricultural value chain (Stone and Meinke, 2005; Schauburger *et al.*, 2020; Singh *et al.*, 2021). Although growing season-averaged weather indices and crop yields have been the subject of prior research, a more thorough understanding of the effects of weather extremes on yields and the application of advanced modeling techniques to enhance predictive capabilities are still needed (Lacasa *et al.*, 2023). Data-driven techniques that integrate high-resolution meteorological data at daily or weekly timeframes can improve the precision of in-season crop production estimates (Joshi *et al.*, 2020). These methods are more effective at capturing the nonlinear effects of meteorological factors, such as temperature and precipitation, on crop development and growth.

In recent years with the evolution of advanced computing technology, techniques such as Machine Learning and Artificial Intelligence are being extensively used in different research areas. Artificial neural networks (ANN) fall under the umbrella of machine learning techniques. Jha and Sinha (2013) defined ANN as a multivariate, non-linear, non-parametric data-driven self-adaptive algorithmic method. The main advantage of neural networks is their flexible functional form as fitting a given data set doesn't require specifying a particular model. Inspired by biological systems, ANN can learn from experience and update solutions at different steps. Currently, ANNs are used for a wide variety of problems arising in many different fields of study such as business, industry, science, agriculture, and others.

The use of ANN and other machine-learning techniques has increased immensely in agricultural research in recent years. Klomplenburg *et al.* (2020) did a systematic literature review of the research work done on the application of machine-learning techniques in crop yield prediction and found that the important independent variables for crop yield prediction are temperature, rainfall, and soil type, and the most used technique for crop yield modeling is Artificial Neural Networks. Jain *et al.* (1980) developed weather indices to develop a crop yield forecast model for rice, which was subsequently used by Kheda, Bhulsar (Chauhan *et al.*, 2009), Budwan & Birbhum (Gupta *et al.*, 2023), Rapeseed & Mustard in

Faizabad (Azfar *et al.*, 2021), and Potato in West Bengal (Gupta *et al.*, 2022b). Laxmi and Kumar (2011) developed neural network-based forecast models for crop yield at the district level. Similar efforts have been made to develop crop yield prediction models for potato (Gupta *et al.*, 2022), rice (Das *et al.*, 2018), coconut (Das *et al.*, 2020), and cashew (Das *et al.*, 2022).

2. Data and methodology

2.1. Data

Rice (*Oryza sativa* L.) is the most dominant crop in India, and West Bengal is the largest producer of rice in India (Agricultural Statistics at a Glance, 2020). In this paper, the Kharif rice yield data in the two districts viz. Burdwan and Birbhum with the highest yield rate of rice in West Bengal have been used. District-wise yearly yield data of the Kharif rice were collected from the released issues of yield estimates of the Bureau of Applied Economics and Statistics (BAES), Department of Statistics and Programme Implementation, Government of West Bengal for 43 years from 1977-78 to 2019-20.

Data on the different weather variables for the districts mentioned here were collected from the NASA Power Data Access Viewer (Das *et al.*, 2018, 2022; Setiya *et al.*, 2023) from 1977 to 2020. For this study, weather data on the minimum temperature, maximum temperature, rainfall, and relative humidity were considered. Daily data were first converted into weekly data per standard meteorological weeks. The weather during the pre-sowing period affects the field preparation, sowing, and germination of the crop. Heavy rainfall or drought during this period can delay sowing or severely damage the sown crop, adversely affecting the crop establishment and, ultimately, crop yield. The weather data from the pre-sowing period, that is, the 23rd standard meteorological week through to the harvesting period, that is, the 41st standard meteorological week, were considered.

2.2. Methodology

2.2.1. Weather indices & multiple linear regression

Weather variables have varying effects on crop growth and development at different stages. As a result, the volume and distribution pattern of weather variables throughout the crop season determines the extent to which they influence crop yield. Furthermore, the influence of a specific meteorological variable on agricultural productivity varies dramatically among weeks. Using the weekly data for each weather variable, a weighted variable is generated concerning each weather variable for each year under consideration. The weights assigned to weekly

weather data for a certain weather variable are simple correlation coefficients calculated over time between crop yield (adjusted for time trend) and weekly weather data for the i^{th} weather variable. The weight assigned to a meteorological variable's m^{th} week is determined by the association between yield and m^{th} week weather data across time. These weighted variables are referred to as weather indices. The weather indices of a particular weather variable in a specific year which is the representative of the weather variable value for that year obtained by combining the weekly influence of weather variable in the following manner

$$Z_{t(i)} = \sum_{w=1}^m r_{\{y_t, X_{tw(i)}\}} X_{tw(i)} / \sum_{w=1}^m r_{\{y_t, X_{tw(i)}\}}, \quad (1)$$

where m denotes the number of weeks in any particular crop season, and p denotes the number of weather variables taken in the study. Let X_{iw} ($i=1, 2, \dots, p; w=1, 2, \dots, m$) denote the value of i^{th} weather variable in w^{th} week and y_t denote the crop yield for the year t ($t=1, 2, \dots, 43$). Gupta (2023) has developed an R package to obtain weather indices from weekly weather data and yearly yield data.

The multiple linear regression model was fitted by taking weather indices and time variable t as independent variables, and the final model was obtained through a stepwise variable selection technique by selecting significant variables. The regression model takes the form as follows:

$$y_t = a + \sum_{i=1}^p b_i Z_{t(i)} + ct + \varepsilon \quad (2)$$

2.2.2. ARIMAX model

The ARIMAX models are similar to ARIMA models with an added layer of complexity, which can incorporate exogenous variables in the expression of ARIMA (p, d, q). The crop yield is a real-world time series data that possesses auto-correlation properties; therefore, an ARIMAX model utilizes the auto-correlation of yield and effects of exogenous variables. In this study, a version of ARIMAX (Hyndman, 2021) was used, which can be termed regression with ARIMA errors. Pandit *et al.*, (2023) used a similar model to develop yield forecast models for rabi crops, and Alam *et al.*, (2018) for rice in India. In this case, the β_i coefficients can be interpreted as typical regression coefficients. It is defined below as

$$y_t = \sum_{i=1}^k \beta_i x_{t(i)} + \eta_t, \quad (3)$$

where

$$\Phi(B)(1 - B)^d \eta_t = \theta(B)\varepsilon_t$$

Therefore, the ARIMAX model becomes

$$y_t = \sum_{i=1}^k \beta_i x_{t(i)} + \frac{\theta(B)\varepsilon_t}{\Phi(B)(1 - B)^d}$$

$$\Phi(B)(1 - B)^d y_t = \Phi(B)(1 - B)^d \sum_{i=1}^k \beta_i x_{t(i)} + \theta(B)\varepsilon_t \quad (4)$$

If the drift term μ is included in the model, the above expression is modified to

$$\Phi(B)(1 - B)^d \left(y_t - \frac{\mu t^d}{d!} \right)$$

$$= \Phi(B)(1 - B)^d \sum_{i=1}^k \beta_i x_{t(i)} + \theta(B)\varepsilon_t \quad (5)$$

The value of μ is close to the actual value mean of the time series when $d=0$ and it is equal to the mean of $(1 - B)^d y_t$ for $d>0$.

2.2.3. Artificial neural network

A completely connected network composed of several fully connected layers that link every neuron in one layer to every other layer's neuron is known as a multilayer feed-forward neural network (MLP) architecture. There are no loops in a feed-forward network, and it has a one-way flow. For a wide range of applications, including forecasting, the multilayered perceptron (MLP) architecture is the most often used neural network architecture (Zhang *et al.*, 1998). An input layer, hidden layers, and output layer constitute this network. Neurons, which are also known as nodes, are nonlinear components that make up each layer. Neurons of the succeeding hidden layers receive input from the preceding layers, and the nodes in the input layer accept the values of the input (predictor) variables. Neurons in every layer are completely coupled to neurons in every other layer. As a result, the neurons in the following layer receive input from the outputs of the neurons in each layer. After all the processing in the hidden levels, the final layer, known as the output layer, provides the output value. Using an activation function or transfer function, each neuron processes the input value locally before transforming the output and sending it to neighboring neurons. The number of hidden layers in the MLP architecture can vary, as can the number of neurons in each layer. In a causal connection problem, predictor variables are sent to neurons in the ANN input layer and

the predicted response variable value is received from the neurons in the output layer of ANN.

The functional relationship estimated by the ANN can be written as $y=f(x_1, x_2, \dots, x_p)$, where x_1, x_2, \dots, x_p are predictor variables and y is the response. The output of j^{th} node in the neural network is given by

$$y_j = g \left(\theta_j + \sum_{i=1}^p w_{ij} x_i \right), \quad (6)$$

where g is a transfer or activation function, θ_j is the bias of the node j , w_{1j}, \dots, w_{pj} are weights of node j and x_i ($i=1, 2, \dots, p$) are the input variables. The most important aspect to consider when creating a multilayer feedforward neural network architecture for a prediction problem is the number of hidden layers and nodes in each hidden layer. Neural networks can perform complex nonlinear mapping between input and output variables by detecting and capturing features and patterns in the data through the use of hidden layers and nodes within them. As Zhang *et al.*, (1998) pointed out, there is currently no theoretical foundation for selecting these parameters. Determining hidden layers and nodes is most commonly accomplished through trial and error.

The activation function and learning algorithm are two crucial components of an ANN. The activation function and learning algorithm set the neuron's weight in the hidden layer. The robust backpropagation (Riedmiller and Braun, 1993) technique, which is a quicker and better version of the widely used backpropagation training process, was utilized in this study together with a sigmoid (logistic) activation function ($f(x) = 1/(1 + e^{-x})$). The sigmoid (logistic) function as an activation function has been used only for hidden layers, whereas for output nodes, the linear activation function has been used, as suggested by Rumelhart *et al.*, (1995) for forecasting problems. Before training the ANNs and hybrid models, the data were scaled on a scale of [0, 1].

2.2.4. Regression-ANN hybrid approach

The variable studied in this paper, rice yield, is real-world time series data, which are not necessarily fully linear or nonlinear. These data often contain both linear and nonlinear components. Then, data y_t can be expressed as:

$$y_t = \hat{y}_t + \varepsilon_t, \quad (7)$$

where \hat{y}_t represents the linear component of the data, is the predicted value obtained from the linear regression models, and ε_t is the residual term from the linear regression model, which represents the nonlinear component of the data.

If the data are truly a combination of linear and nonlinear patterns, neither linear regression nor ANNs can adequately capture and model the overall pattern in the data because linear regression cannot deal with the nonlinear relationship between a dependent variable and independent variables, whereas ANNs alone may be unable to adequately model both linear and nonlinear components of the data. To better model complicated data, it is recommended that linear and nonlinear models be combined. Zhang (2003) proposed this approach for the first time and Ray *et al.*, (2016) applied it for wheat yield forecasting and recently, Gupta *et al.*, (2022 a) applied the same for potato yield modeling.

The hybrid modeling approach is essentially a two-step procedure in which the regression model deals with linear patterns of data and the ANN model deals with nonlinear patterns of data. The findings of both models are then integrated to generate the hybrid model. More precisely, in the first phase, a linear regression model, is applied to the data under consideration, and the residuals of this model, which contain the nonlinear pattern of the data, are modeled using the ANN method. The fitted values from the linear regression model estimate the linear pattern, whereas the fitted residuals from the ANN model will estimate the nonlinear pattern of the data. The final combined estimated model can be given as

$$\hat{y}'_t = \hat{y}_t + \hat{\varepsilon}_t, \quad (8)$$

where \hat{y}_t are the fitted values from the linear regression model, and $\hat{\varepsilon}_t$ is the fitted values of the residuals from the ANN model.

The *Modus Operandi* of the study was to fit a multiple linear regression at first and ultimately obtain a final regression model after variable selection. Two approaches were considered for the ANN models: first, only the significant weather variables were used as inputs, and second, all the independent variables (four weather indices and time variable t) were used as inputs. The number of hidden nodes was varied from one to five. The hybrid MLR-ANN model was fitted with the same procedure as ANN, except that the output variable in this approach was the residuals obtained from the regression model. All analyses were performed in the R statistical package (R Core Team, 2024).

3. Results and discussion

3.1. Multiple linear regression and ARIMAX models

The multiple regression model with four weather indices and time variable t was fitted for rice yield and the variable selection was employed to get the final models which are represented in Table 1.

TABLE 1(a)

Parameter estimates of multiple linear regression models for rice yield

District	Model structure
Burdwan	$Y_t = -35.17 - 0.07^{***} Z_{prep} + 2.14^{*} Z_{mint} + 0.41^{***} t$ (26.36) (0.02) (1.02) (0.04)
Birbhum	$Y_t = -30.07^{*} - 0.02^{***} Z_{prep} + 1.89^{**} Z_{mint} + 0.42^{***} t$ (14.41) (0.004) (0.59) (0.03)

TABLE 1(b)

Parameter estimates of ARIMAX models for rice yield

District	Model structure	drift	ar ₁	ma ₁	Z _{mint}	Z _{prep}
Burdwan	ARIMA (1,1,1)-X	0.41 (0.28)	-0.92 ^{***} (0.13)	0.71 ^{**} (0.243)	1.55 (0.90)	-0.05 ^{***} (0.014)
Birbhum	ARIMA (1,1,0)-X	0.38 (0.20)	-0.44 ^{**} (0.17)		1.88 ^{***} (0.54)	-0.02 ^{***} (0.004)

Figures in the parentheses represent the standard errors and the symbols *, **, and *** represent the significance of estimates at 5%, 1%, and 0.1%, respectively.

The final multiple regression models for Burwan and Birbhum were found to be highly significant with $F=78.88^{***}$ and $F=170.1^{***}$ respectively. Furthermore, it was revealed that the time variable (t), minimum temperature, and precipitation have a significant effect on the yield of rice in both districts and therefore were retained in the model after variable selection. These significant weather variables were then used as exogenous variables in ARIMAX models.

The ARIMAX models were identified using Box and Jenkins (2015) methodology for the rice yield data in both districts. Various tentative ARIMAX models of several order combinations were fitted, and finally, selected models are presented in Table 1(b). The training and testing accuracy measures of ARIMAX models are presented in Table 3 for comparison purposes.

3.2. Non-linearity of Data

The residuals obtained from the regression model in Table 1 were subjected to the BDS test and the results in Table 2 show most of the p-values were less than 0.05 but not all. These results indicate the presence of both linear and non-linear components in the rice yield of the Burdwan and Birbhum districts. Pandit *et al.* (2023) and Ray *et al.* (2016) also used BDS test to check the non-linearity in the data and reported presence of non-linearity in the yield data.

3.3. ANN & MLR-ANN hybrid models

The BDS test results of MLR residuals in Table 2 indicated the presence of both linear and non-linear patterns in the rice yield. Therefore, several neural

TABLE 2

Non-linearity test results of rice yield

Embedding dimension	Burdwan				Birbhum			
	$\epsilon=0.5\sigma$	$\epsilon=0.5\sigma$	$\epsilon=\sigma$	$\epsilon=1.5\sigma$	$\epsilon=2\sigma$	$\epsilon=\sigma$	$\epsilon=1.5\sigma$	$\epsilon=2\sigma$
m=2	2.36 (0.018)	87.58 (0.000)	1.14 (0.255)	0.56 (0.573)	2.40 (0.016)	2.96 (0.003)	-1.26 (0.208)	-2.46 (0.014)
m=3	3.55 (0.000)	89.12 (0.000)	2.39 (0.017)	2.37 (0.018)	3.24 (0.001)	6.73 (0.000)	2.1 (0.036)	0.57 (0.565)

Figures in the parenthesis represent the p-values of the test statistics

network models and hybrid models with varying numbers of inputs and hidden nodes were fitted. The Z_{prep} , Z_{mint} , and time variable t were taken as inputs and yield as output. The model structures and accuracy measures of various ANNs and hybrid models are given in Appendices 1 and 2. These rice yield prediction models were subjected to scrutiny for better accuracy among different models within their approach and the models with lower errors from each approach are presented in Table 3.

The models in Table 3 were compared for accuracy to find out the best-fit model. Among all the approaches, the hybrid MLR-ANN model with three and four hidden nodes was found to be the best-fit model having the lowest training and testing errors among all models. These models combine the linear result of MLR and the non-linear result of artificial neural networks.

Khan *et al.*, (2024) in their study also reported that it can be inferred that the hybrid model, PCA-ANN, outperformed the individual models. Saravanan & Bhagavathiappan (2024) also found hybrid deep learning-based crop yield prediction model were successful in modeling the crop yield data. A structural presentation of fitted neural fitted with weights at each neuron is given in Fig. 1.

Furthermore, actual values, predicted values, and their corresponding percentage errors for the period of testing data *i.e.* 2015-16 to 2019-20 have been presented in Table 4. Fig. 2 represents the actual values, the training set fitted values and the testing set predicted values of the best-fit models. In addition, Table 5 represents the 6-8 weeks before pre-harvest forecasts by best-fit models. These forecasts are very close to the actual yield values in both districts.

3.4. Model diagnostics

The assumption of no autocorrelation among residuals was tested using the “Ljung-Box” test, which yielded a test statistic value of $\chi^2=16.75$ (p-value=0.669) for Burdwan and $\chi^2=23.54$ (p-value=0.263) for Birbhum, indicating no autocorrelation among residuals. In addition,

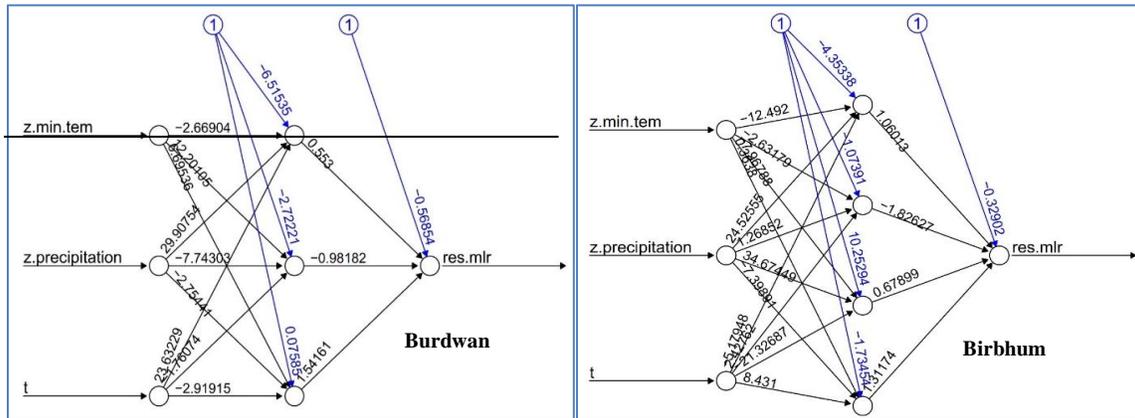


Fig. 1. NN structure of Hybrid models for rice yield

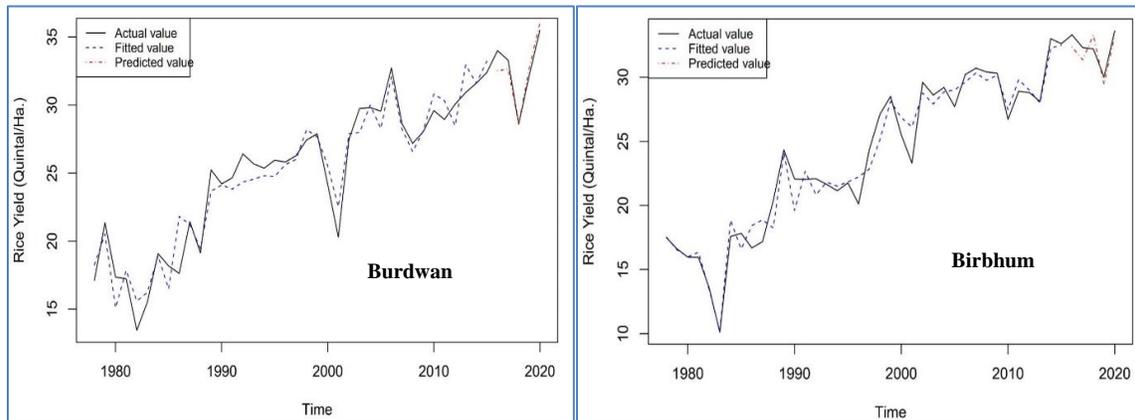


Fig. 2. Actual and model fitted values of rice yield

TABLE 3

Comparison of prediction models for rice yield

District	Model	Structure	Training			Testing		
			RMSE	MAE	MAPE	RMSE	MAE	MAPE
Burdwan	MLR	$Y_t = -35.17 - 0.07Z_{prep} + 2.14Z_{min} + 0.41t$	1.82	1.50	6.81	1.56	1.03	3.06
	ARIMAX	ARIMA (1,1,1)-X	1.86	1.45	6.32	1.42	1.23	3.87
	ANN.5	3:5:1	1.82	1.50	6.81	1.44	0.98	2.96
	ANN.9	5:4:1	1.47	1.27	5.71	1.50	1.13	3.34
	HANN.3	3:3:1	1.34	1.03	4.65	0.80	0.65	1.92
	HANN.6	5:1:1	1.72	1.42	6.43	1.46	1.00	3.02
Birbhum	MLR	$Y_t = -30.07 - 0.02Z_{prep} + 1.89Z_{min} + 0.42t$	1.45	1.23	5.90	1.30	0.97	3.02
	ARIMAX	ARIMA (1,1,0)-X	1.72	1.36	6.23	1.06	0.84	2.64
	ANN.3	3:3:1	1.31	1.11	5.22	1.06	0.79	2.43
	ANN.9	5:4:1	1.33	1.07	4.96	1.18	0.90	2.78
	HANN.4	3:4:1	1.13	0.84	3.75	0.84	0.80	2.47
	HANN.8	5:3:1	1.27	1.03	4.73	1.01	0.79	2.45

the assumption of normality of residuals was tested using the Shapiro-Wilk test, which yielded a test statistic value of $W=0.96$ (p -value=0.163) for Burdwan and $W=0.96$ (p -value=0.724) for Birbhum.

These results indicated that the residuals did not deviate significantly from normality, and the same was reflected in the QQ plot in Fig. 3, where most of the points are on the normal QQ line.

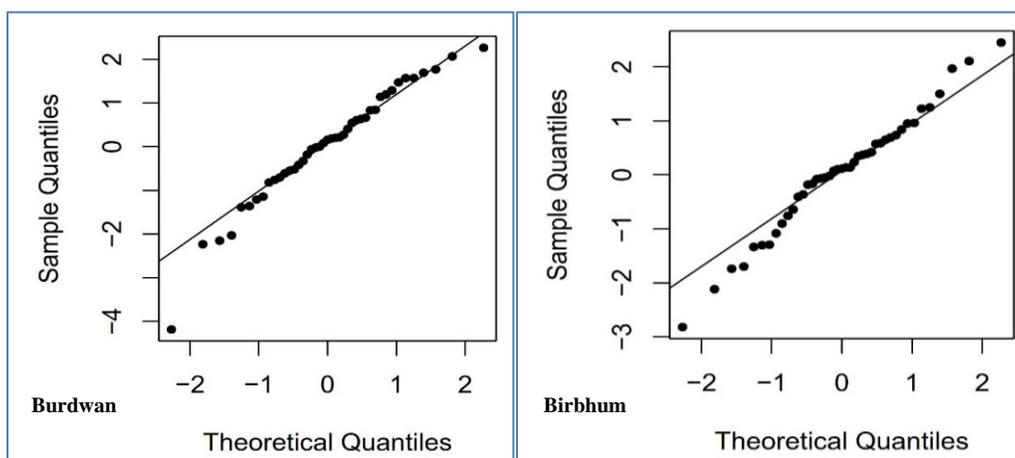


Fig. 3. Normal QQ plot of best-fit model residuals

TABLE 4

Predicted values and percentage prediction errors for rice yield

District	Model	2015-16	2016-17	2017-18	2018-19	2019-20
Burdwan	Actual yield	34.00	33.30	28.60	32.20	35.50
	MLR	30.67(10.85)	32.41(2.74)	28.73(-0.44)	32.76(-1.70)	35.74(-0.68)
	ARIMAX	31.85(6.76)	32.44(2.65)	30.6(-6.53)	32.97(-2.35)	35.88(-1.05)
	ANN.5	30.96(9.81)	32.55(2.31)	28.96(-1.23)	32.91(-2.16)	35.45(0.15)
	ANN.9	31.09(9.37)	32.6(2.13)	28.64(-0.15)	31.55(2.07)	34.13(4.00)
	HANN.3	32.53(4.52)	32.64(2.03)	28.62(-0.08)	32.73(-1.6)	36.05(-1.53)
	HANN.6	30.94(9.89)	32.63(2.06)	29.02(-1.45)	33.03(-2.52)	35.51(-0.04)
Birbhum	Actual yield	33.30	32.30	32.20	30.00	33.60
	MLR	32.06(3.87)	29.81(8.36)	31.91(0.92)	30.74(-2.42)	33.51(0.26)
	ARIMAX	32.64(2.02)	30.36(6.39)	32.43(-0.71)	31.18(-3.77)	33.79(-0.58)
	ANN.13	32.44(2.65)	30.19(6.99)	32.44(-0.75)	29.86(0.47)	33.00(1.83)
	ANN.19	31.75(4.89)	30.25(6.79)	32.00(0.61)	30.5(-1.65)	33.4(0.60)
	HANN.4	32.34(2.97)	31.35(3.04)	33.28(-3.26)	29.43(1.95)	33.18(1.26)
	HANN.8	32.64(2.02)	30.3(6.60)	32.9(-2.11)	30.45(-1.47)	33.46(0.43)

Figures in the parenthesis represent the % error of predicted values.

TABLE 5

Pre-harvest forecasts by best-fit models in both districts

District	Model	2015-16	2016-17	2017-18	2018-19	2019-20
Burdwan	Actual yield	34.00	33.30	28.60	32.20	35.50
	6 Weeks before	28.82(17.98)	29.98(11.09)	24.63(16.09)	32.17(0.09)	35.54(-0.11)
	8 weeks before	28.87(17.75)	30.62(8.74)	25.97(10.14)	32.09(0.33)	35.72(-0.61)
Birbhum	Actual yield	33.30	32.30	32.20	30.00	33.60
	6 Weeks before	32.63(1.44)	32.63(1.44)	32.63(1.44)	32.63(1.44)	32.63(1.44)
	8 weeks before	32.48(2.52)	32.48(2.52)	32.48(2.52)	32.48(2.52)	32.48(2.52)

Figures in the parenthesis represent the % error of predicted values

TABLE 6

Prediction accuracy test of models for rice yield

Burdwan			Birbhum		
Model pairs	DM statistic	P-value	Model pair	DM statistic	P-value
HANN.3-MLR	4.66	0.000	HANN.4-MLR	2.80	0.004
HANN.3-ANN.5	4.74	0.000	HANN.4-ANN.3	2.45	0.009
HANN.3-ANN.9	2.08	0.022	HANN.4-ANN.9	1.79	0.041
HANN.3-HANN.6	3.61	0.000	HANN.4-HANN.8	0.64	0.264

3.5. DM test among models

The model performance results in Table 3 established that the HANN.3 and HANN.4 models have the lowest error measures among other approaches. Also, Fig. 3 indicates that the fitted values of the hybrid model are very close to the actual values. However, to further validate the difference in accuracy of the models, the statistical significance of the difference in the accuracy of the models was tested by the DM test. The prediction accuracy of the best fit hybrid model has been compared with all other models and the DM test results are given in Table 6.

The results in Table 6 revealed that the hybrid model has significantly different prediction accuracy as compared to other approaches except for HANN.8 model in Birbhum. It indicates that the hybrid models HANN.3 in Burdwan and hybrid model HANN.9 in Birbhum have significantly better prediction accuracy than other models. Hence, on combining the results of Table 3, Table 4, and Table 6 the hybrid MLR-ANN models with three and four hidden nodes are the best-fit prediction model for rice yield in the Burdwan and Birbhum districts respectively.

Weather variability, driven by climate change, significantly impacts rice production. Numerous studies have demonstrated the profound influence of weather variables on rice yield. Temperature, rainfall, and relative humidity are among the most critical factors. Precipitation, especially rainfall during the growing season, is a critical factor affecting the growth and yield of crops, particularly Kharif monsoon rice (Bowden *et al.*, 2023; Maiti *et al.*, 2024). We found that there was a highly significant negative effect of rainfall on rice yield. This may be caused by the inundation following the flood situation in the Ganges region and rainfall received during the rice harvesting period which has adverse effects on the rice yield. Asada and Matsumoto (2009) reported a similar result of a negative correlation between rice production and rainfall in districts of the lower Ganges region. Tack *et al.*, (2015) highlighted how extreme temperatures, both

hot and cold, can negatively impact crop yields. Optimal temperatures are essential for rice growth and development (Peng *et al.*, 2004). Excessive heat can lead to sterility and reduced grain filling (Xu *et al.*, 2021), while cold temperatures can delay germination and growth (Hussain *et al.*, 2019). The effect of minimum temperature was found to be positive in this study. Furthermore, Nemoto *et al.*, (2016) highlight the need for forecasting models to consider not just the magnitude of weather events but also their timing within the crop's growth cycle which is accounted for by weather indices used in this research.

Accurately forecasting rice yield is crucial for global food security, especially in regions where rice is a staple crop. Ghosh *et al.* (2015) showcased the development of a rice yield prediction system in India, combining extended-range forecasts with crop models. This approach underscores the potential of integrating meteorological data with crop modeling to enhance forecasting accuracy. However, Mosleh *et al.* (2015) remind us that many existing forecasting methods are empirical and require further calibration and validation before being applied to different geographical locations. This highlights the need for robust and adaptable forecasting models that can account for regional variations in climate and rice varieties.

4. Conclusions

The minimum temperature and rainfall are the most important weather variables for rice yield in West Bengal. Interestingly, the effect of rainfall on rice yield came out to be negative while the effect of minimum temperature was positive on rice yield. It is concluded that neural network models with proper tuning of hyperparameters perform better than linear models for forecasting purposes. However, the importance of statistical models cannot be marginalized as they can contribute in the sense of choosing the right inputs for the neural network models. The selected best-fit models can be used to obtain a reliable forecast of crop yield 6-8 weeks before harvest for various policy decisions at various levels.

Authors' statement

The contents and views expressed in this research paper/article are the views of the authors and do not necessarily reflect the views of the organizations they belong to.

Conflict of interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Authors' contribution

Akhilesh Kumar Gupta: Conceptualization, methodology, data curation, formal analysis, validation, original draft preparation, visualization. (email: akgupta@ouat.ac.in).

Kader Ali Sarkar: methodology, original draft preparation, supervision. (email: kader804@gmail.com).

Debasis Bhattacharya: validation, review and editing, supervision. (email: debasis_us@yahoo.com).

References

- Alam, W. Ray, M. Kumar, R. R. Sinha, K. Rathod, S. and Singh, K. N., 2018, "Improved ARIMAX modal based on ANN and SVM approaches for forecasting rice yield using weather variables", *Indian J. Agric. Sci.*, **88**, 12. 1909-1913.
- Anonymous, 2020, "Agricultural Statistics at a Glance", *Ministry of Agriculture & Farmers Welfare.*, Government of India.
- Azfar, M. Sisodia, B. Rai, V. and Devi, M., 2021, "Pre-harvest forecast models for rapeseed & mustard yield using principal component analysis of weather variables", *Mausam*, **66**, 4. 761-766. <https://doi.org/10.54302/mausam.v66i4.583>
- Asada, H, and Matsumoto, J., 2009, "Effects of rainfall variation on rice production in the Ganges-Brahmaputra Basin", *Clim. Res.*, **38**. 249-260. <https://doi.org/10.3354/cr00785>
- Bowden, C. Foster, T. and Parkes, B., 2023, "Identifying links between monsoon variability and rice production in India through machine learning", *Sci. Rep.*, **13**:2446. <https://doi.org/10.1038/s41598-023-27752-8>
- Box, G.E.P. Jenkins, G.M. Reinsel, G.C. and Ljung, G.M., 2015, "Time Series Analysis: Forecasting and Control", John Wiley & Sons
- Chauhan, V.S. Shekh, A.M. Dixit, S.K. Mishra, A.P. and Kumar, S., 2009, "Yield prediction model of rice in Bulsar district of Gujarat", *J. Agrometeorology.*, **11**, 2.162-168. <https://doi.org/10.54386/jam.v11i2.1245>
- Das, B. Murganekar, D. Navyashree, S. and Kumar, P., 2022, "Novel combination artificial neural network models could not outperform individual models for weather-based cashew yield prediction", *Int. J. Biometeorol.*, **66**. 1627-1638. <https://doi.org/10.1007/s00484-022-02306-1>
- Das, B. Nair, B. Arunachalam, V. Reddy, K. V. Venkatesh, P. Chakraborty, D. and Desai, S., 2020, "Comparative evaluation of linear and nonlinear weather-based models for coconut yield prediction in the west coast of India", *Int. J. Biometeorol.*, **64**. 1111-1123. <https://doi.org/10.1007/s00484-020-01884-2>
- Das, B. Nair, B. Reddy, V. K. and Venkatesh, P., 2018, "Evaluation of multiple linear, neural network and penalised regression models for prediction of rice yield based on weather parameters for west coast of India", *Int. J. Biometeorol.*, **62**, 10. 1809-1822. <https://doi.org/10.1007/s00484-018-1583-6>
- Ghosh, K. Singh, A. Mohanty, U. C. Acharya, N. Pal, R. K. Singh, K. K. and Pasupalak, S., 2015, "Development of a rice yield prediction system over Bhubaneswar, India: combination of extended range forecast and CERES-rice model", *Meteorological Applications.*, **22**. 525-533. <https://doi.org/10.1002/met.1483>
- Gupta, A., 2023, "_weatherindices: Calculate Weather Indices.", *R package version 0.1.0.*, <<https://CRAN.R-project.org/package=weatherindices>>.
- Gupta, A.K. Sarkar, K.A. Dhakre, D.S. and Bhattacharya, D., 2023, "Weather based crop yield prediction using artificial neural networks: A comparative study with other approaches", *Mausam*, **74**, 3. 825-832. <https://doi.org/10.54302/mausam.v74i3.174>
- Gupta, A.K. Sarkar, K.A. Bhattacharya, D. and Dhakre, D.S., 2022a, "Potato yield modeling based on meteorological factors using discriminant analysis and artificial neural networks", *Int. J. Veg. Sci.*, **28**, 5. 465-476. <https://doi.org/10.1080/19315260.2021.2021342>
- Gupta, A.K. Sarkar, K.A. Dhakre, D.S. and Bhattacharya, D., 2022b, "Weather Based Potato Yield Modelling using Statistical and Machine Learning Technique", *Environment and Ecology.*, **40**. 1444-1449.
- Hussain, S. Khaliq, A. Ali, B. Hussain, H. A. Qadir, T. and Hussain, S., 2019, "Temperature Extremes: Impact on Rice Growth and Development", In: Hasanuzzaman M, Hakeem KR, Nahar K, Alharby HF ,eds, "Plant Abiotic Stress Tolerance: Agronomic, Molecular and Biotechnological Approaches", *Springer International Publishing*, Cham. 153-171
- Hyndman, R.J. and Athanasopoulos, G., 2021, "Forecasting: Principles and Practice, 3rd edn", *OTexts.*, Melbourne, Australia
- Jain, R.C. Agrawal, R. and Jha, M.P., 1980, "Effect of climatic variables on rice yield and its forecast", *Mausam*, **31**, 4. 591-596. <https://doi.org/10.54302/mausam.v31i4.3477>
- Jha, G.K. and Sinha, K., 2013, "Agricultural Price Forecasting Using Neural Network Model: An Innovative Information Delivery System", *Agric. Econ. Res. Rev.*, **26**. 229-239
- Joshi, V. R. Kazula, M. J. Coulter, J. A. Naeve, S. L. and Garcia y Garcia, A., 2021, "In-season weather data provide reliable yield estimates of maize and soybean in the US central Corn Belt", *Int. J. Biometeorol.*, **65**. 489-502. <https://doi.org/10.1007/s00484-020-02039-z>
- Lacasa, J. Messina, C.D. and Ciampitti, I.A., 2023, "A probabilistic framework for forecasting maize yield response to agricultural inputs with sub-seasonal climate predictions", *Environ. Res. Lett.*, **18**:074042. <https://doi.org/10.1088/1748-9326/acd8d1>
- Laxmi, R.R. and Kumar, A., 2011, "Weather based forecasting model for crops yield using neural network approach", *Stat. Appl.*, **9**, 1&2. 55-69.
- Maiti, A. Hasan, M. K. Sannigrahi, S. Bar, S. Chakraborti, S. Mahto, S. S. Chatterjee, S. Pramanik, S. Pilla, F. Auerbach, J. Sonnentag, O. Song, C. and Zhang, Q., 2024, "Optimal rainfall threshold

- for monsoon rice production in India varies across space and time”, *Commun. Earth Environ.*, **5**, 1. 1–8. <https://doi.org/10.1038/s43247-024-01414-7>
- Mosleh, M.K. Hassan, Q.K. and Chowdhury, E.H., 2015, “Application of Remote Sensors in Mapping Rice Area and Forecasting Its Production: A Review”, *Sensors*, **15**. 769–791. <https://doi.org/10.3390/s150100769>
- Nemoto, M. Hamasaki, T. Matsuba, S. Hayashi, S. and Yanagihara, S., 2016, “Estimation of Rice Yield Components with Meteorological Elements Divided According to Developmental Stages”, *J. Agricul. Met.*, **72**, 3-4. 128–141. <https://doi.org/10.2480/agrmet.D-15-00017>
- Pandit, P. Sagar, A. Ghose, B. Dey, P. Paul, M. Alqadhi, S. Mallick, J. Almohamad, H. and Abdo, H. G., 2023, “Hybrid time series models with exogenous variable for improved yield forecasting of major Rabi crops in India”, *Sci. Rep.*, **13**, 1. <https://doi.org/10.1038/s41598-023-49544-w>
- Peng, S. Huang, J. Sheehy, J. E. Laza, R. C. Visperas, R. M. Zhong, X. Centeno, G. S. Khush, G. S. and Cassman, K. G., 2004, “Rice yields decline with higher night temperature from global warming”, *Proc. Natl. Acad. Sci. USA.*, **101**, 27. 9971–9975. <https://doi.org/10.1073/pnas.0403720101>
- R Core Team, 2024, “R: A language and environment for statistical computing”, *R Foundation for Statistical Computing, Vienna, Austria.*, URL <https://www.R-project.org/>
- Ray, M. Rai, A. Vaidhyanathan, R. and Singh, K., 2016, “ARIMA-WNN Hybrid Model for Forecasting Wheat Yield Time-Series Data”, *J Indian Soc. Agric. Stat.*, **70**, 1. 63–70
- Riedmiller, M. and Braun, H., 1993, “A direct adaptive method for faster backpropagation learning: the RPROP algorithm”, In: *IEEE International Conference on Neural Networks.*, 586–591 vol.1
- Rumelhart, D.E. Durbin, R. Golden, R. and Chauvin, Y., 1995, “Backpropagation: The basic theory”, In: *Backpropagation: Theory, architectures, and applications*, Lawrence Erlbaum Associates., Inc, Hillsdale, NJ, US, 1–34.
- Schauberger, B. Jägermeyr, J. and Gornott, C., 2020, “A systematic review of local to regional yield forecasting approaches and frequently used data resources”, *Eur. J. Agron.*, 120:126153. <https://doi.org/10.1016/j.eja.2020.126153>
- Setiya, P. Satpathi, A. and Nain, A.S., 2023, “Predicting rice yield based on weather variables using multiple linear, neural networks, and penalized regression models”, *Theor. Appl. Climatol.*, **154**. 365–375. <https://doi.org/10.1007/s00704-023-04563-5>
- Singh, A.K. Singh, N. Singh, H. and Kushwaha, H.S., 2021, “Role and Importance of Weather Forecasts in Modern Agriculture”, *Int. J. Curr. Microbiol. App. Sci.*, **10**, 1. 2646–2662. <https://doi.org/10.20546/ijcmas.2021.1001.308>
- Stone, R.C. and Meinke, H., 2005, “Operational seasonal forecasting of crop performance”, *Philosophical Transactions of the Royal Society B: Biological Sciences.*, 360:2109–2124. <https://doi.org/10.1098/rstb.2005.1753>
- Tack, J. Barkley, A. and Nalley, L.L. 2015, “Effect of warming temperatures on US wheat yields”, *Proceedings of the National Academy of Sciences.*, **112**. 6931–6936. <https://doi.org/10.1073/pnas.1415181112>
- Van Klompenburg, T. Kassahun, A. and Catal, C. 2020, “Crop yield prediction using machine learning: A systematic literature review”, *Comp. Electron. Agricul.*, 177:105709. <https://doi.org/10.1016/j.compag.2020.105709>
- Xu, Y. Chu, C. and Yao, S., 2021, “The impact of high-temperature stress on rice: Challenges and solutions”, *Crop J.*, **9**, 5. 963–976. <https://doi.org/10.1016/j.cj.2021.02.011>
- Zhang, G. Eddy Patuwo, B. Y. Hu, M., 1998, “Forecasting with artificial neural networks: The state of the art”, *Int. J. Forecast.*, **14**, 1. 35–62. [https://doi.org/10.1016/S0169-2070\(97\)00044-7](https://doi.org/10.1016/S0169-2070(97)00044-7)
- Zhang, G.P., 2003, “Time series forecasting using a hybrid ARIMA and neural network model”, *Neurocomputing*, **50**. 159–175. [https://doi.org/10.1016/S0925-2312\(01\)00702-0](https://doi.org/10.1016/S0925-2312(01)00702-0)