

Document Version

Final published version

Licence

CC BY

Citation (APA)

Diaz, V., Osman, A. A., Corzo Perez, G. A., Maskey, S., & Solomatine, D. P. (2026). Spatiotemporal changes of drought area as input for a machine-learning approach for crop yield prediction. *Hydrology Research*, 57(2), 125-149. <https://doi.org/10.2166/nh.2026.150>

Important note

To cite this publication, please use the final published version (if applicable). Please check the document version above.

Copyright

In case the licence states "Dutch Copyright Act (Article 25fa)", this publication was made available Green Open Access via the TU Delft Institutional Repository pursuant to Dutch Copyright Act (Article 25fa, the Taverne amendment). This provision does not affect copyright ownership. Unless copyright is transferred by contract or statute, it remains with the copyright holder.

Sharing and reuse

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

Takedown policy

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

Spatiotemporal changes of drought area as input for a machine-learning approach for crop yield prediction

Vitali Diaz ^{a,*}, Ahmed A. Osman^b, Gerald A. Corzo Perez ^{a,c}, Shreedhar Maskey ^c and Dimitri Solomatine ^{a,c,d}

^a Delft University of Technology, Delft, The Netherlands

^b Arcadis, Cardiff, Wales, United Kingdom

^c IHE Delft Institute for Water Education, Westvest 7, 2611 AX, Delft, The Netherlands

^d Water Problems Institute of the Russian Academy of Sciences, Moscow, Russia

*Corresponding author. E-mail: v.diazmercado@tudelft.nl; vitalidime@gmail.com

 VD, 0000-0002-5502-4099; GAC, 0000-0002-2773-7817; SM, 0000-0002-3259-5374; DS, 0000-0003-2031-9871

ABSTRACT

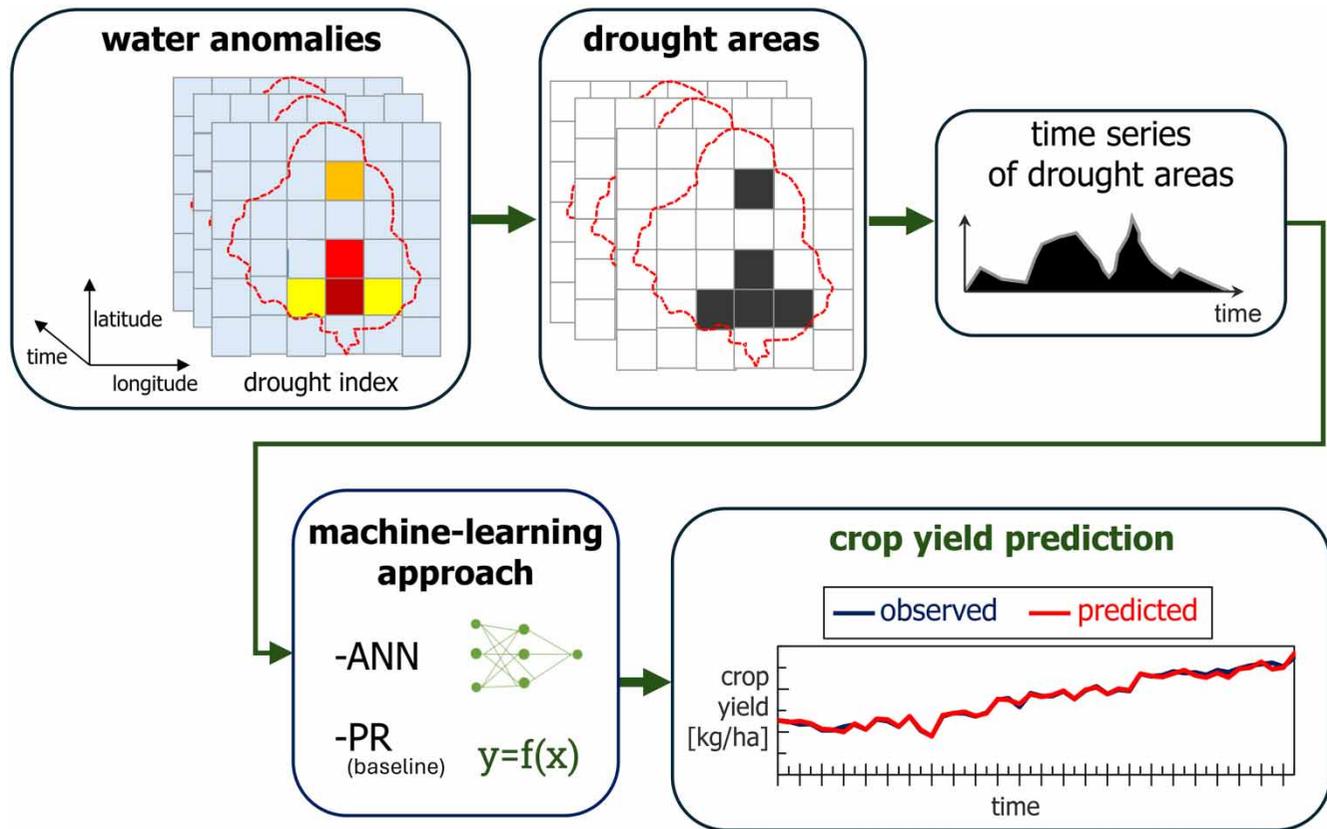
More severe and prolonged droughts observed in recent decades require improved methods to predict impacts on agriculture. Crop-growth models estimate yield and plant development variables and are widely used to assess drought impacts; however, they are not explicit forecasting tools, as their accuracy is constrained by physical assumptions, data availability, and multiple sources of uncertainty. To address these limitations, machine learning (ML) models have been increasingly applied for crop yield prediction, typically using drought indices as input, while spatial drought characteristics remain underexplored. This research develops an ML framework that incorporates the spatial extent of drought to predict seasonal crop yield. The framework combines artificial neural network (ANN) and polynomial regression (PR) models, with PR providing baseline estimates and ANN delivering refined predictions. The approach was tested using 50 years of historical crop yield data and drought areas derived from the Standardised Precipitation Evapotranspiration Index at multiple aggregation periods (1–12 months). Results show ANN models consistently outperform PR models, achieving lower prediction errors, with root mean square error values as low as 48.1 kg/ha in best-performing cases. The results demonstrate that spatiotemporal drought area dynamics and their temporal aggregation provide an effective preprocessing strategy for ML-based drought impact prediction.

Key words: agricultural drought, crop yield, drought area, drought impact, machine learning, spatiotemporal analysis

HIGHLIGHTS

- A step-by-step crop yield prediction framework is developed using ANN and polynomial regression models.
- Spatial extent of meteorological drought is shown to be an effective input for crop yield prediction.
- Different drought aggregation periods enable yield prediction at varying accuracy levels during the crop season.
- Polynomial regression equations provide practical baseline models for early crop yield estimation.

GRAPHICAL ABSTRACT



1. INTRODUCTION

Drought negatively impacts a wide range of human activities, particularly agriculture. These impacts can lead to substantial financial losses and, in severe cases, contribute to food insecurity and loss of life (WMO 2006; Below *et al.* 2007; Sheffield & Wood 2011; FAO 2017; Kim *et al.* 2019). Therefore, there is a critical need for methods that support effective drought mitigation, especially those capable of predicting the impacts of drought (WMO 2006; FAO 2017).

Traditionally, physical crop models have been used to assess the impact of environmental stressors, such as drought. These models simulate plant development by incorporating biological and environmental processes, offering detailed insights into crop performance under drought stress (White *et al.* 1997; Reynolds *et al.* 2000; Wu *et al.* 2016; Huang *et al.* 2019). However, they typically require extensive and high-resolution input data, such as detailed soil, crop, and climate parameters, which are often unavailable or difficult to obtain in many regions (Chlingaryan *et al.* 2018; van Klompenburg *et al.* 2020). This limits their practicality for near-real-time applications.

To address these limitations, statistical and machine learning (ML) models have been used as viable alternatives for predicting the impacts of drought on agriculture. These data-driven approaches rely on historical input–output relationships rather than modelling internal biological mechanisms (Chlingaryan *et al.* 2018; Rahmati *et al.* 2020; Udmale *et al.* 2020; van Klompenburg *et al.* 2020; Araneda-Cabrera *et al.* 2021). One commonly used metric for evaluating the impacts of drought on agriculture is crop yield, defined by the Food and Agriculture Organization of the United Nations (FAO) as the amount of crop produced per unit of land area (kg/ha or ton/ha) (FAO & DWFI 2015). Numerous studies have demonstrated the utility of various ML techniques in crop yield forecasting, typically employing multivariable input data that include weather, soil, and management-related variables (Chlingaryan *et al.* 2018; van Klompenburg *et al.* 2020).

Despite the success of ML models, a critical gap remains in how drought is addressed. Most ML crop yield prediction studies primarily use drought indices as input variables (Chlingaryan *et al.* 2018; van Klompenburg *et al.* 2020). However, spatiotemporal characteristics, particularly drought extent (i.e., the spatial coverage of drought-affected areas), remain

underexplored. Several studies have shown that drought extent correlates strongly with crop yield variability (Diaz *et al.* 2016; Osman 2018; Osman *et al.* 2018; Araneda-Cabrera *et al.* 2021), but the full implementation for crop yield prediction is lacking.

The primary objective of this research is to develop a machine learning approach to predict seasonal crop yields using drought area. Our ML approach includes an artificial neural network (ANN) for accurate predictions and polynomial regression (PR) for baseline calculations. Together, the ANN and PR models can be used as an integrated tool for crop yield prediction.

To demonstrate the utility of the proposed ML approach, three regions in eastern India, where rice cultivation is prevalent and agriculture plays a vital economic role (Ghosh *et al.* 2014), were selected as case studies. Historical data from 1967 to 2015 was used to train and evaluate the models, with a focus on regional rice yield prediction. The methodology was applied independently to each region, allowing for tailored calibration and comparative performance assessment.

2. DATA

2.1. Crop yield

Rice is the most widely cultivated crop in East India, contributing approximately 85% of the country's total rice production (Ghosh *et al.* 2014). State-level rice yield data from 1966 to 2015 was obtained from the Directorate of Economics and Statistics, Department of Agriculture and Cooperation (DAC). This dataset was used to develop machine learning (ML) models for three regions in East India, as shown in Figure 1.

India has three crop seasons: Rabi, Kharif, and Zaid. The Kharif season was selected because it produces the largest crop yield. Kharif crops are sown in June and harvested in November/December. Seasonal crop yield data was obtained from the DAC website and organised into a time series per region. Each year of harvest in the Kharif season was assigned a single value (Figure 1). No data filling was performed in the time series, as data exists for each year in all three regions.

Figure 1 shows the locations of these regions, which are defined as follows: Region 1 includes the states of Bihar and Jharkhand; Region 2 corresponds to West Bengal; and Region 3 comprises Odisha. Two important clarifications regarding the retrieval of crop yield data for these regions are necessary. First, in late 2000, Bihar was split into two states: Bihar and Jharkhand. Subsequently, rice data was reported separately for each new state. In this study, both states' data is aggregated as Region 1, so the yearly crop yield data from 2000 to 2015 is the sum of the current Bihar and Jharkhand figures. Second, in 2011, Orissa was renamed Odisha, although the territory remained unchanged. In this case, Odisha's crop yield data is the union of the former Orissa and the current Odisha time series (Region 3).

2.2. Drought indicator

For the calculation of drought areas, we used Standardised Precipitation-Evapotranspiration Index (SPEI) data. SPEI was proposed by Vicente-Serrano *et al.* (2010) and has proven a useful proxy for assessing agricultural drought (Osman 2018; Osman *et al.* 2018; Araneda-Cabrera *et al.* 2021; Khoshnazar *et al.* 2022). The SPEI follows a similar methodology to the widely used Standardised Precipitation Index (SPI) (Mckee *et al.* 1993), but with added consideration for the difference between precipitation and evapotranspiration.

SPEI data was retrieved from the SPEI Global Drought Monitor (<https://spei.csic.es>) covering the period from 1901 to 2015. The spatial resolution of the drought indicator data is 0.5 degrees. The SPEI data was available at various aggregation periods; for this study, it was retrieved for the periods of one, three, six, nine, and 12 months, designated as DI1, DI3, DI6, DI9, and DI12, respectively. No further temporal aggregation was performed. Drought areas were computed using this drought indicator data following the procedure described in Section 3.1.1.

3. ML MODELLING METHODOLOGY

The ML approach was developed by following the methodology of ML model building. The following paragraphs detail each step. These steps are (1) data preparation, (2) input variable selection, (3) PR model training, (4) ANN model training, and (5) application of the models for crop yield prediction.

3.1. Step 1. Data preparation

Two types of data were prepared: the crop yield (CY) and the drought areas (DA) time series. In this step the following tasks were performed: (1) drought areas calculation, and (2) data de-trending.

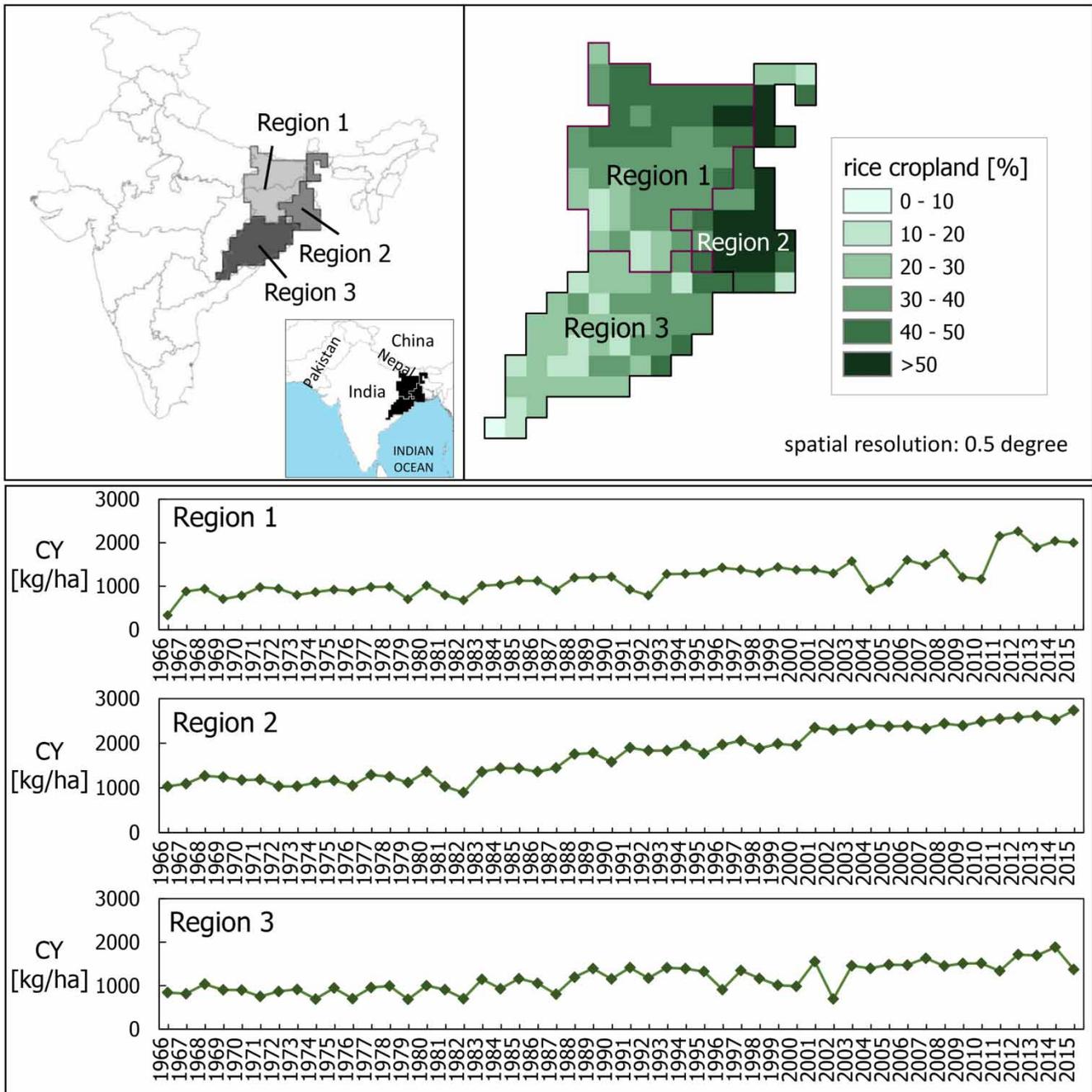


Figure 1 | Case study location (top) and crop yield (CY) data for Kharif season (bottom). The case study includes region 1 (Bihar and Jharkhand), region 2 (West Bengal), and region 3 (Odisha). The percentage of rice cropland is indicated (top). Source of rice cropland: Monfreda *et al.* (2008).

3.1.1. Drought areas calculation

Drought area (DA) calculation begins with reclassifying all the cells of the drought indicator data as non-drought (0 s) and drought (1 s) cells. To indicate drought and non-drought conditions (DS), Equation (1) was applied (Corzo Perez *et al.* 2011; Herrera-Estrada *et al.* 2017; Diaz *et al.* 2019, 2020). Equation (1) works as follows: when the drought indicator falls below the chosen threshold τ , the value of 1 is used to indicate drought in the cell; otherwise, 0 is used to indicate

non-drought. This classification is performed for all grid data cells at each time step (t):

$$D_S(t) = \begin{cases} 1 & \text{if } DI(t) \leq \tau \\ 0 & \text{if } DI(t) > \tau \end{cases} \quad (1)$$

After, DAs were computed as the ratio between the cells in drought and the total number of cells of the region (N) (Equation (2)). In Equation (2), the cell's id is denoted by c :

$$DA(t) = 100/N \cdot \sum_{c=1}^N D_S(t) \quad (2)$$

Masks were also created for each region. The mask is an array of ones and zeros, where 1 signifies land. The number of cells (N) in the mask is 63, 31, and 54 for regions 1, 2, and 3, respectively. We used the threshold $\tau = -1$ to identify drought cells. This threshold is commonly used to detect drought conditions in standardised indices (Diaz *et al.* 2019). DAs were computed for each aggregation period of the drought indicator data and are denoted as DA1, DA3, DA6, DA9, and DA12.

3.1.2. Data de-trending

Short-term fluctuations were removed from the time series before building the ML models. We used the 'first difference' method for its straightforward implementation (Montesino Pouzols & Lendasse 2010). In this method, the trend is removed from the time series of both drought areas and crop yield by subtracting the previous value $x^*(t-1)$ from the current value $x^*(t)$, as shown in Equation (3). The de-trended value for the first time step ($t=1$) is not calculated. The length of the de-trended time series is $n = m - 1$, where m is the length of the original time series. The de-trended data $x(t)$ has the same units as the original data $x^*(t)$:

$$x(t) = x^*(t) - x^*(t-1) \quad (3)$$

Once the time series is de-trended, all the ML model construction steps are carried out using the de-trended time series. After the ML models are built, the de-trending process must be reversed to obtain predictions in the original scale. The reverse de-trending process can be performed with Equation (4), which is the inverse of Equation (3) for the de-trended prediction $x(t+1)$. Practically, the prediction $x^*(t+1)$ in the original magnitude is calculated by adding the de-trended prediction $x(t+1)$ to the last value of the original time series, i.e. $x^*(t)$:

$$x^*(t+1) = x^*(t) + x(t+1) \quad (4)$$

As observed in Equation (3), the 'first difference' method does not produce a value for the initial time step; therefore, the CY data, which initially covers the period from 1966 to 2015, now runs from 1967 to 2015 after de-trending.

Regarding DA, Equation (3) was used as follows. First, the monthly DAs were rearranged to organise the drought areas for each month from January to December (Figure 2). This rearrangement process was performed for each of the five aggregation periods, DA1, DA3, DA6, DA9, and DA12 months, resulting in a total of 60 rearranged DA time series (12×5). A suffix was added to identify these series by month; for example, the time series DA3_7 indicates drought areas for all Julys of the period, calculated from the drought indicator with a 3-month aggregation period. Equation (3) was applied to each of the 60 DA time series (Figure 2). The de-trended DA series span from 1902 to 2015. For building the ML models, the common period with de-trended CY of 1967–2015 (49 years) was selected.

3.2. Step 2. Input variable selection

The correlation was calculated between the de-trended time series of CY and the de-trended, rearranged time series of DA, based on the CY–DA pairings described in Section 3.1.2 (see Figure 2).

We selected the time series of DAs with the highest correlation as input variables. This selection is based on two main considerations. First, rice responds differently to climate variations at various growth stages throughout the year. Therefore, certain months may provide more relevant information than others for capturing these stage-specific impacts (Osman 2018; Osman *et al.* 2018). Second, different types of droughts (meteorological, agricultural, and hydrological) affect crop

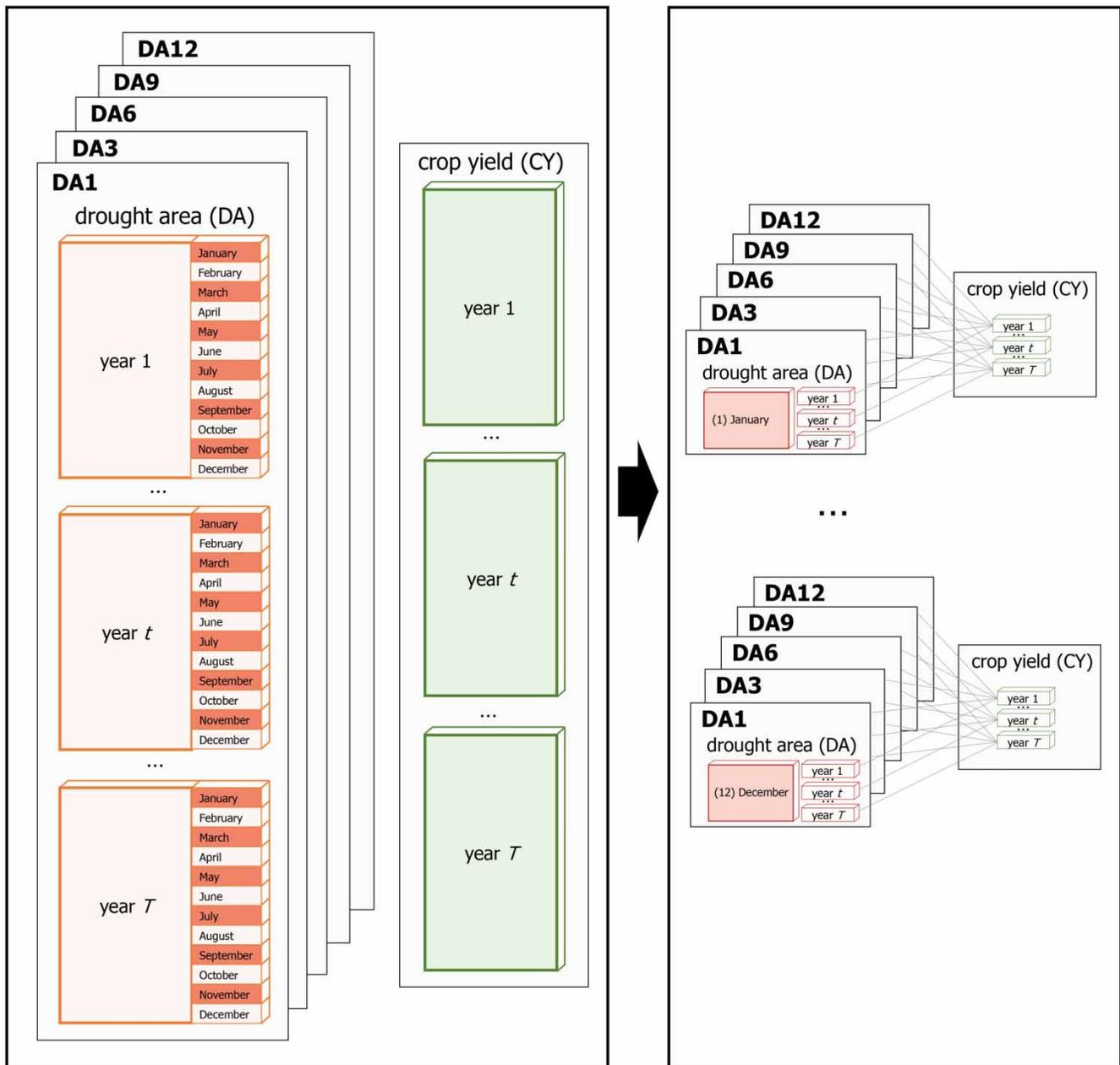


Figure 2 | The diagram illustrates how the monthly DAs time series were reorganised to align with seasonal CY data. (Left) For each year, there are 12 DA values and one CY value. DAs were calculated for the aggregation periods of 1, 3, 6, 9, and 12 months (DA1 to DA12). (Right) DAs were reorganised by month from January to December. The procedures of data de-trending, correlation analysis, input variable selection, and ML models construction were carried out for each month. The entire procedure was carried out independently for each of the three case studies.

yield in distinct ways depending on the crop's developmental stage (Araneda-Cabrera *et al.* 2021). These varying impacts can be accounted for by using different hydro-meteorological variables or by applying different temporal aggregation periods to the meteorological data (Araneda-Cabrera *et al.* 2021; Khoshnazar *et al.* 2022), as done in this study.

A simple average of DAs across the season could obscure significant drought events in specific months that may have a stronger or weaker influence on final crop yield (Osman 2018; Osman *et al.* 2018). Additionally, the ML models in this research were designed to be applicable at different stages of crop cultivation – for example, models tailored for June, July, and other months – each with its own expected level of accuracy. Therefore, using monthly time series of DAs for all aggregation periods one, three, six, nine, and 12 months) is more appropriate than relying on seasonal averages (Figure 2).

Table 1 | PR types followed in this study

PR type	Equation	Description
Linear	$y = b_0 + \sum_{i=1}^n b_i x_i$ Eq. (5)	It has an intercept and linear terms of predictors
Pure-quadratic	$y = b_0 + \sum_{i=1}^n b_i x_i + \sum_{i=1}^n b_{n+i} x_i^2$ Eq. (6)	It has an intercept, as well as linear and squared terms of predictors
Quadratic	$y = b_0 + \sum_{i=1}^n b_i x_i + \sum_{i=1}^n b_{n+i} x_i^2 + \sum_{i=1}^{n-1} \sum_{j=i+1}^n b_{2n+(i-1)n-\frac{(i-1)j}{2}+(j-i)} x_i x_j$ Eq. (7)	It has an intercept, linear and squared terms and all products of pairs of distinct predictors
Interactions	$y = b_0 + \sum_{i=1}^n b_i x_i + \sum_{i=1}^{n-1} \sum_{j=i+1}^n b_{n+(i-1)n-\frac{(i-1)j}{2}+(j-i)} x_i x_j$ Eq. (8)	It has an intercept, linear terms of predictors, all products of pairs of distinct predictors and no squared terms

3.3. Step 3. PR models training

For the case of PR, four types of models were tested (Table 1). All the PR models were constructed for each month from January to December following Equations (5)–(8). A total of 15 sets of input variable combinations were examined. MATLAB software was used for implementing the PR models. PR is an extension of linear regression that allows multiple input variables to calculate the output variable (Equation (9)):

$$y = b_0 + \sum_{i=1}^n b_i x_i + e \quad (9)$$

In Equation (10), y is the output variable, also known as the response, which in this case is the crop yield. The term x_i represents the i -th input variable (predictor) out of a total of n variables. The regression coefficients vector is denoted by b . From this vector, b_0 is recognised as the intercept. The vector of errors is indicated by e . Table 1 displays the four formulations of PR models: linear, pure-quadratic, quadratic, and interactions (Equations (5)–(8)). The input variable (x_i) was chosen based on the correlation analysis (Section 2.2).

Of the four types of PR models, the best was identified using the root mean square error (RMSE) (Equation (9)). The RMSE is calculated as the difference between the observations (o) and the predictions (p):

$$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^n (o_i - p_i)^2}{n}} \quad (10)$$

3.4. Step 4. ANN models training

ANN is a method loosely based on mimicking the basic functionality of neurons (i.e., the working units of the human brain) (Govindaraju 2000; Maier & Dandy 2000). The input variables (predictors) are connected to each other through mathematical formulations that enable the representation of complex, non-linear relationships. These connections are symbolised as nodes interconnected within a network to calculate the output variable (response).

In the ANN setup, the number of nodes in the input layer equalled the number of variables in the respective combination. The number of nodes in the output layer was one and represented the seasonal CY. An iterative optimisation process was performed for the hidden layer, by varying the number of nodes from 1 to 10. For each number of nodes, 100 iterations were conducted, totalling 1,000. To ensure reproducibility, random values were reset to default at the start of each variation in node number. ANNs were constructed for each month, from January to December. MATLAB software was used to implement the ANNs, employing the Levenberg–Marquardt algorithm for training. In each ANN, 85% of the data was used for training and 15% for testing (verification), similar to Elshorbagy *et al.* (2010). For each month, the best model was identified for each number of hidden nodes, resulting in 10 candidate models. From these, the overall optimal model for that month was selected. RMSE was used to select the best models. RMSE was calculated for (1) the training dataset

(RMSE_cal), (2) the testing dataset (RMSE_test), and (3) the entire dataset (RMSE). In all cases, the final (best) model was chosen based on RMSE over the full period. The iteration optimisation process, including RMSE calculation, was conducted for each set of input variables and each month (Section 4.2).

3.5. Step 5. Application of the ANN/PR models for crop yield prediction

Once the most effective ML models were built, a set of models was selected for each month based on their performance. The ANN models can be used independently or in combination with the PR models as a baseline, as proposed in our approach, enabling a comparison of the outputs of ANN models. Moreover, when spatial input data required to calculate drought area inputs for the ANN models are not yet available, the PR models can be applied using early estimations of drought areas, allowing for continuous early-season predictions.

4. RESULTS AND DISCUSSION

4.1. Data preparation: drought areas and crop yield

Figure 3 illustrates the drought areas computed for the three regions. Region 1 (Figure 3, the upper panel) displays the highest values among the three regions. The 1990s show greater drought area than most other years, which aligns with Guha-Sapir (2019); during this decade, there were three droughts: 1993, 1996, and 2000. At the start of the period, large drought areas were also observed in all three regions; these findings correspond with Bhalme & Mooley (1980).

Figure 3 reveals a pattern in the distribution of drought areas across all aggregation periods, from DA1 to DA12. In DA1, most areas are concentrated in the early months of the year; even December appears nearly white (without drought). Later, in DA3, large areas are situated from April to November. Subsequently, for DA6 and DA9, the largest areas are centred in the second half of the year. Some droughts extend into the following year; these are indicated by the reddish lines visible during the first half of the year (initial columns). Finally, in DA12, consecutive large areas are marked by reddish lines; droughts typically start in the second semester and persist into the subsequent year. These findings highlight the importance of considering multiple aggregation periods when applying drought indicators based on meteorological variables; each period serves as a proxy for capturing different types of drought.

Figure 4 illustrates the time series of de-trended CY and DA for the three regions. For DA (shown in red), the values are presented in reverse order to make interpretation easier. Generally, when the drought area increases, crop yield decreases; conversely, crop yield tends to increase when the drought area decreases. Across all three regions, decreases in CY often align with increases in DA.

In each region, the fluctuations in de-trended CY are more frequent during specific periods. For example, in region 1, fluctuations are more common from 2003 to 2015; in region 2, from 1967 to 2001; and in region 3, also notably from 1967 to 2001. Region 1 is the most northerly of the three.

The overall pattern of DA variation is as follows: values fluctuate throughout the year for the one and three-month aggregation periods (DA1 and DA3). For DA6 to DA12, values tend to concentrate in the second half of the year. These findings also highlight the usefulness of different aggregation periods in detecting various types of drought. Increasing DA does not always correspond to reduced CY across all aggregation periods. For instance, in region 1 (Figure 4, upper panel), the decline in CY in 2004 does not coincide with increases in DA9 and DA12; however, it does align with DA1, DA3, and DA6. These outcomes further support the importance of using multiple aggregation periods to assess drought assessments.

4.2. Input variable selection (correlation analysis)

Figure 5 summarises the correlation between de-trended CY and DAs, while Figure 6 displays the correlation for each monthly DA time series. Both figures demonstrate that the correlation varies throughout the year across the three regions. In all cases, the correlation coefficient rises to a peak and then declines. The month when this maximum occurs differs for each region but is within the crop season (i.e., June to November/December). In region 1, the maximum occurs in July. Region 2 exhibits this pattern in four months: June, July, October, and November. Finally, region 3 shows it in October, November, and December.

These correlation results are valuable for monitoring agricultural drought. For instance, in region 1, the drought-affected areas (DA6) reach their highest correlation in July. This suggests that the accumulated effect of the previous six months is vital for the crop yield of the Kharif season, which roughly spans from June to November/December.

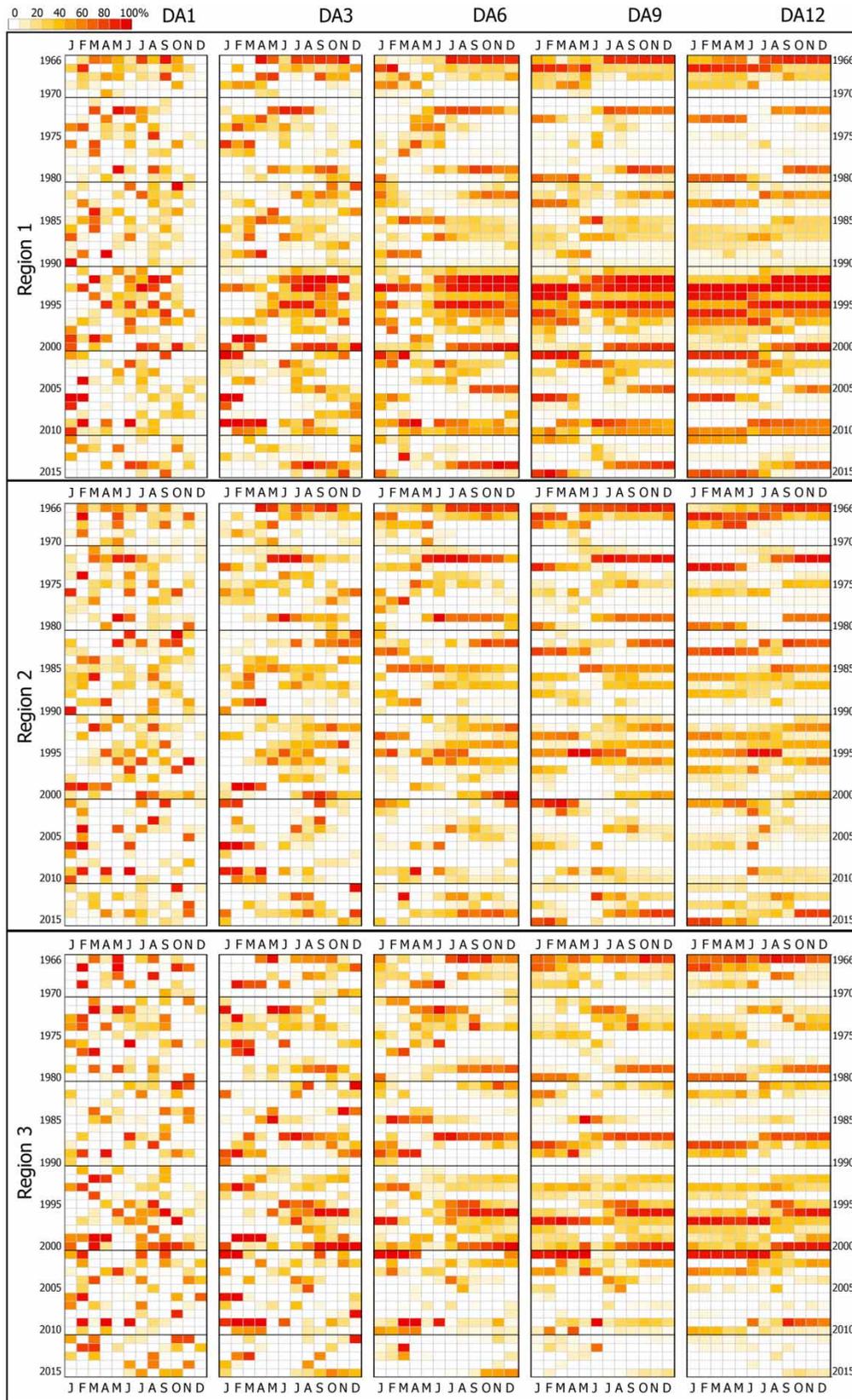


Figure 3 | DAs for each aggregation period (one, three, six, nine, and 12 months) and region. The top, middle, and bottom panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal) and region 3 (Odisha).

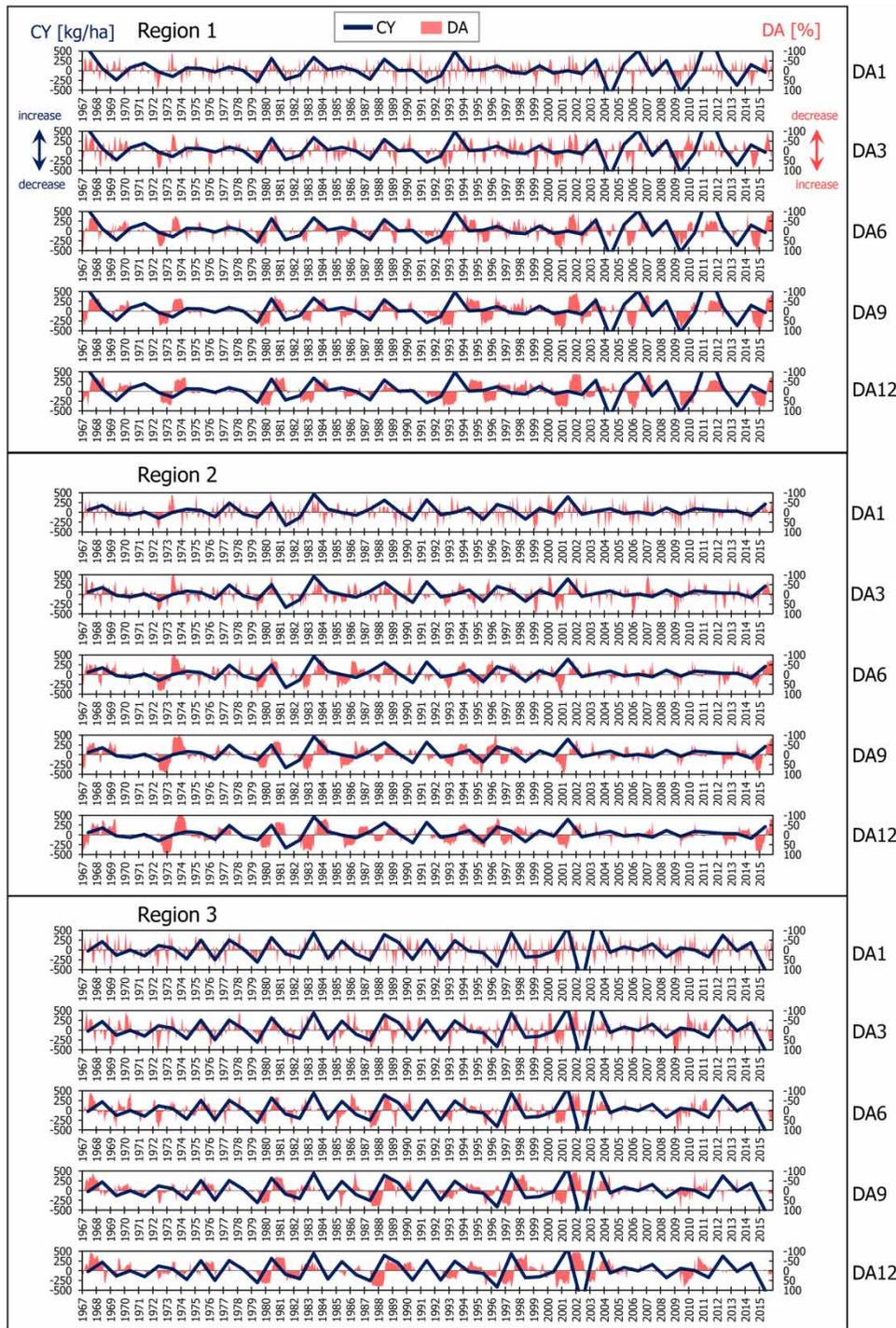


Figure 4 | Time series of the de-trended CY and DAS for each aggregation period (one, three, six, nine, and 12 months) and region. The top, middle, and bottom panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal) and region 3 (Odisha).

Figure 5 reveals a general pattern in the correlation coefficient (R). In region 1 (Figure 5(a)), DA6, DA9, and DA12 display increasing R values from June onwards. In region 2 (Figure 5(b)), a similar trend occurs, with two peaks: one for DA1 and DA3, and another for DA6, DA9, and DA12. The initial peak for DA1 and DA3 indicates the importance of monitoring the conditions in the immediate one to three months. The second peak suggests that medium- and long-term conditions, over six

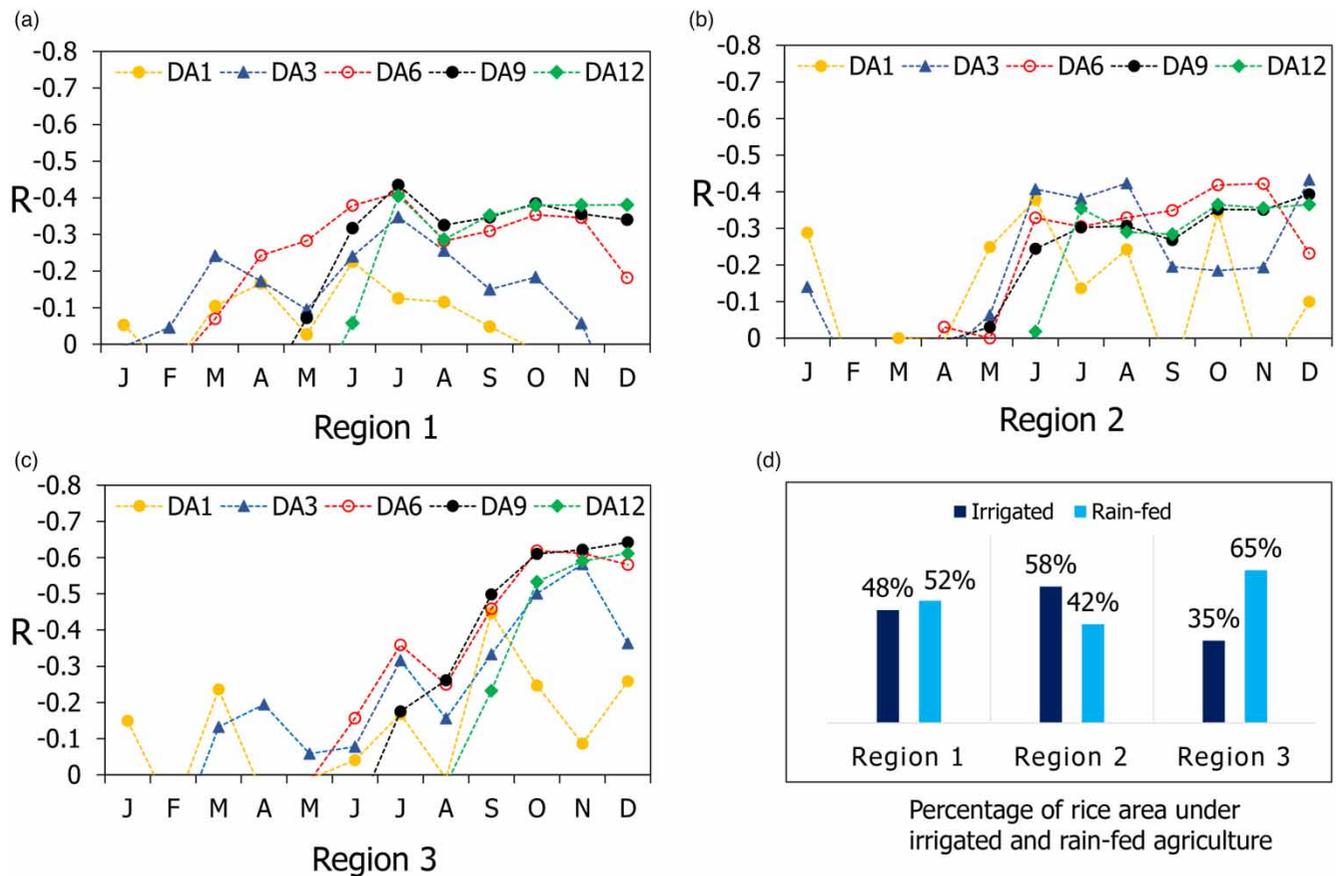


Figure 5 | Summary of correlation between de-trended CY and DAs for each aggregation period (one, three, six, nine, and 12 months) and region: (a) region 1 (Bihar and Jharkhand), (b) region 2 (West Bengal) and (c) region 3 (Odisha). Negative R indicates the correlation between the increase in DA and the decrease in CY. (d) Percentage of rice area under irrigated and rain-fed agriculture. Source of irrigated and rain-fed agriculture data: Directorate of Rice Development (DRD) (2014).

to 12 months, are crucial for the harvest period. In region 3 (Figure 5(c)), most peaks in R happen at the end of the growing season, highlighting that pre-harvest conditions are decisive for crop yield.

Figure 6 illustrates that the correlation between CY and DAs is positive outside the growing season and negative during it. However, this pattern is less clear for DA1 and DA3. These correlation patterns support the idea that drought significantly influences crop yield, as months with lower drought areas are more associated with increases in CY. Conversely, months with higher drought areas often see declines in CY.

Figure 5(d) presents the percentage of irrigated versus rain-fed agriculture. Regions 1 and 2 rely on irrigation for approximately half of their agriculture, whereas region 3 depends on irrigation for only about 35%. This lower percentage of irrigation in region 3 may help explain why its correlation coefficients are higher compared to the other regions (Figures 5 and 6(c)). As region 3 is more dependent on rainfall, this condition is better captured when drought is calculated using precipitation, as in this case (Section 3.2).

Figure 5(a)–5(c) illustrates the following pattern across the three regions. The correlation coefficients between CY and DAs increase with longer aggregation periods and as the year progresses. DA1 and DA3 show better correlations in the early months of the year. DA6 exhibits stronger correlations in subsequent months, from May/June onwards. Lastly, DA9 and DA12 do so in the second half of the year, particularly towards the end of the growing season.

Based on the correlation coefficients, the input variables were selected. A total of 15 input variable sets (Table 2) were chosen. Each set comprises different DA time series, namely DA1, DA3, DA6, DA9, and DA12. The CY of the previous year was included in all combinations and is indicated as CY_{t-1} . Combinations 1 to 5 include only a single DA time series. Combinations 6 to 9 consist of pairs of DAs calculated with the drought indicator at successive aggregation times.

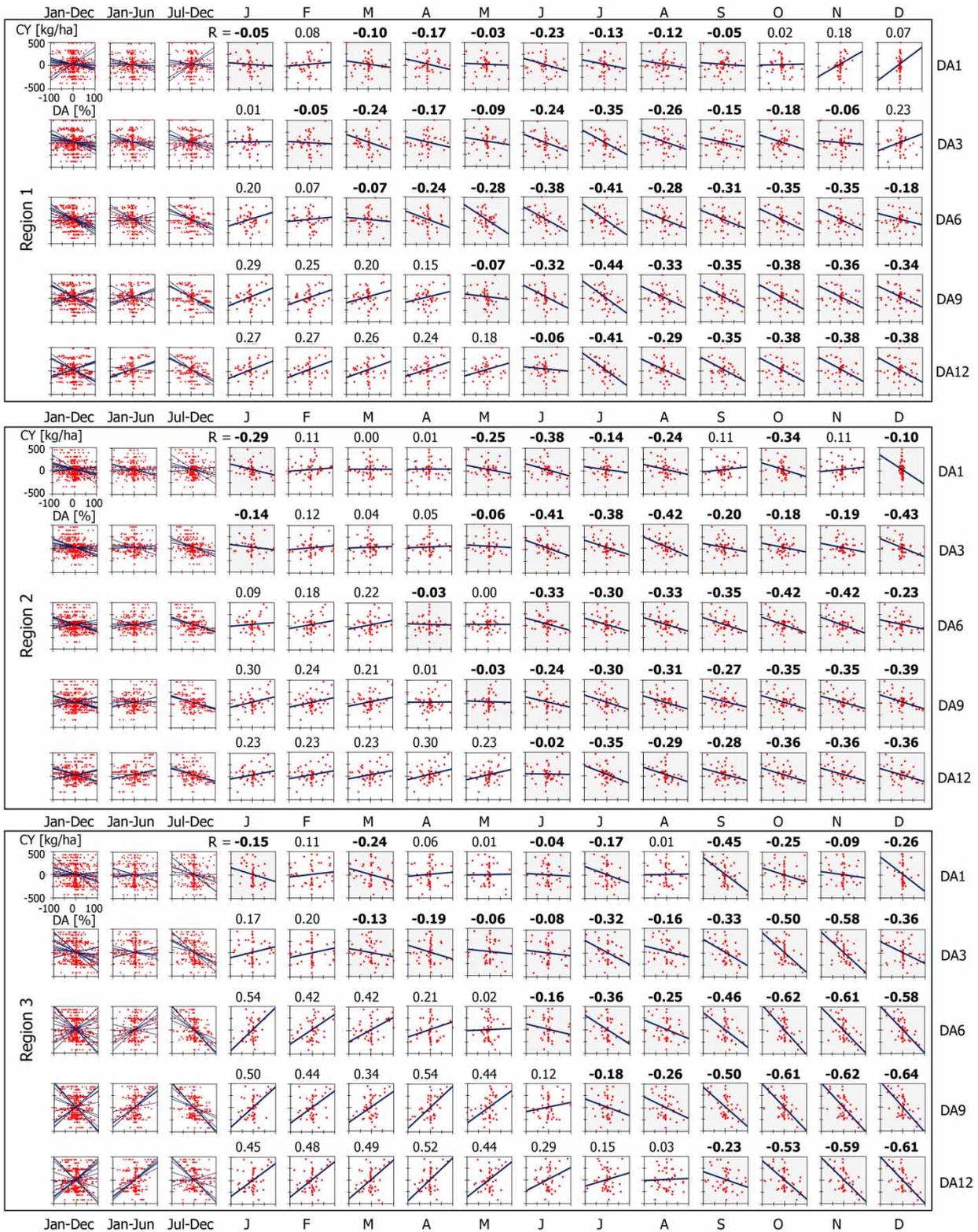


Figure 6 | Correlation (R) between de-trended CY and DAs for each aggregation period (one, three, six, nine, and 12 months) and region. DA is on the x-axis, and CY is on the y-axis. Results are presented for annual (January–December), seasonal (January–June and July–December), and individual monthly DA time series (January to December). The top, middle, and bottom panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal), and region 3 (Odisha). Negative R indicates the correlation between the increase in DA and the decrease in CY.

For example, combination 6 includes DA1 and DA3; combination 7 includes DA3 and DA6, and so on. Similarly, combinations 10 to 12 involve triples, while combinations 13 and 14 involve four DA series, and combination 15 includes all DA time series.

The models were developed for each month (January to December), using the 15 combinations listed in Table 2. For instance, in January, the series of DAs extracted for that month, such as DA1_1, DA3_1, DA6_1, DA9_1, and DA12_1, were applied. The suffix indicates the month. These DAs were then used following the 15 combinations in Table 2 to build the January models (ANN and PR). This process was similarly applied from February onwards to December.

4.3. ANN and PR models training

The results show varying magnitudes of error between observed and predicted CY. The models with the lowest errors are presented in Figures 7, 8, and 9 for each of the three regions. For each month, the best-performing ANN and PR models are shown, with the RMSE indicated in each case. Conversely, Figure 10 displays the error for each input set (combination); the lowest error achieved in each month is presented for each ANN and PR. As expected, ANN exhibits the least errors (Figure 10). However, PR results are not significantly worse; for example, in some cases, the errors shown by linear PR are very close to those of ANN (e.g., Figure 10, region 2). Generally, models with the lowest errors correspond to region 2, followed by regions 3 and 1 (Figure 10). This is attributed to varying levels of crop irrigation and differences in irrigation water sources (surface water and groundwater), which influence the accuracy of the models across regions. Another factor affecting model performance is the sudden changes in CY data, with regions 1 and 3 exhibiting the most pronounced fluctuations, and to a lesser extent, region 2. Figure 10 also indicates that different types of PR yield better results across the three regions. Overall, linear and pure-quadratic models tend to produce more stable results (with fewer abrupt changes among realisations), but not necessarily better performing than quadratic and interaction models. Quadratic and interaction models generally yield superior results, often very close to those of ANN, for example, in PR interactions (Figure 10, region 1).

4.4. Application of ANN/PR models for crop yield prediction

The best-performing models were selected for each month. Table 3 summarises these models, including the input set (combination), number of nodes, and errors for ANN, as well as the input set, model type, and errors for PR. The number of nodes indicates the model complexity in each model. In this way, the more nodes, the more complex the model in the case of ANN. Conversely, quadratic and interaction types showed the best performance in PR models. In all cases, no single DA time series

Table 2 | Input sets (combinations) used to build the ML models. CY and DA stand for crop yield and drought area

Input set (combination)	Input variables
1	CY _{t-1} , DA1
2	CY _{t-1} , DA3
3	CY _{t-1} , DA6
4	CY _{t-1} , DA9
5	CY _{t-1} , DA12
6	CY _{t-1} , DA1, DA3
7	CY _{t-1} , DA3, DA6
8	CY _{t-1} , DA6, DA9
9	CY _{t-1} , DA9, DA12
10	CY _{t-1} , DA1, DA3, DA6
11	CY _{t-1} , DA3, DA6, DA9
12	CY _{t-1} , DA6, DA9, DA12
13	CY _{t-1} , DA1, DA3, DA6, DA9
14	CY _{t-1} , DA3, DA6, DA9, DA12
15	CY _{t-1} , DA1, DA3, DA6, DA9, DA12

DAs are calculated with the drought indicator at aggregation periods of one, three, six, nine, and 12 months (details in Section 4.2).

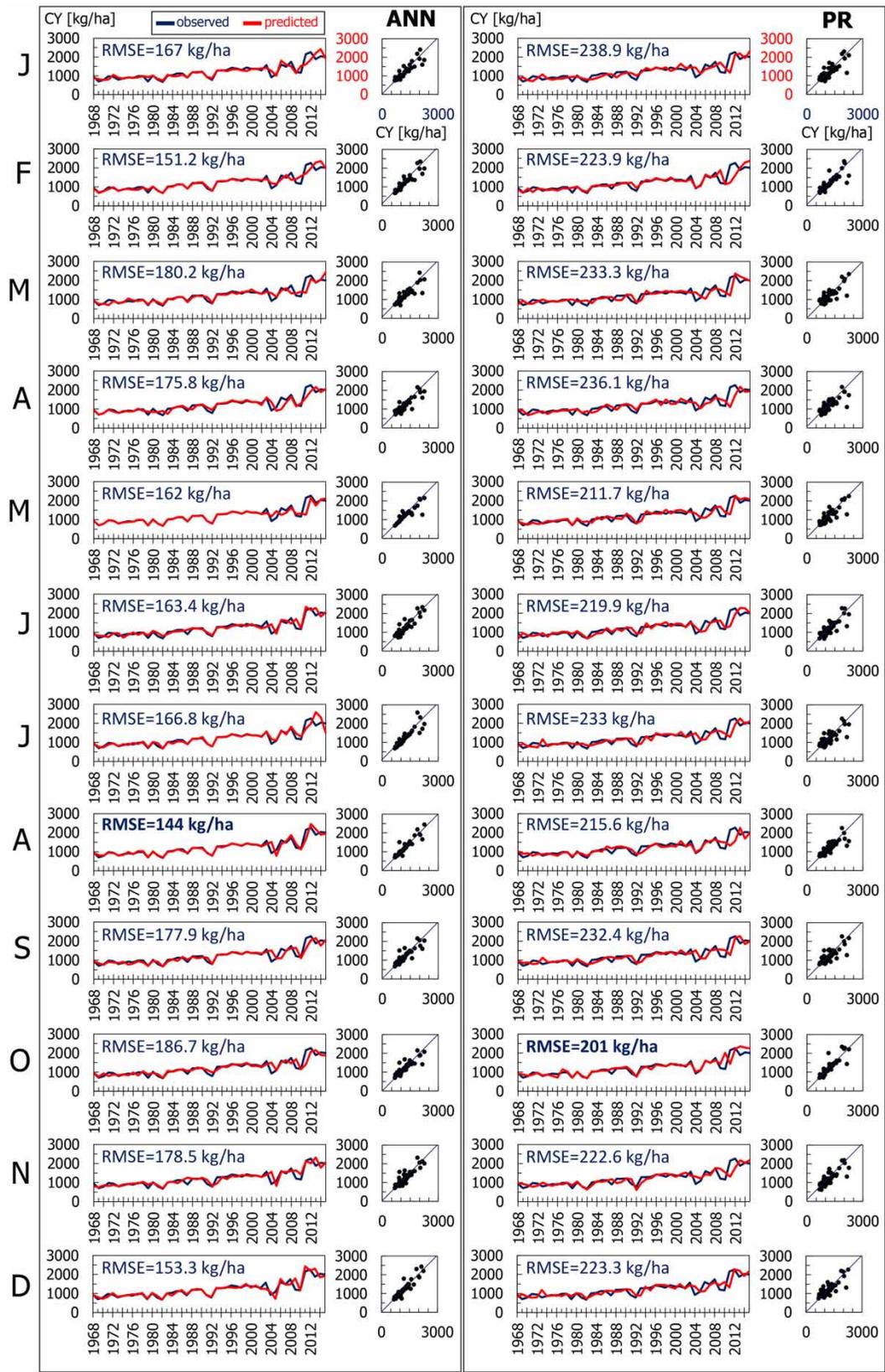


Figure 7 | ANN and PR models for predicting seasonal CY built using each time series of monthly DAS: region 1 (Bihar and Jharkhand). The model with the lowest error (RMSE) is presented for each month, from January to December (J to D).

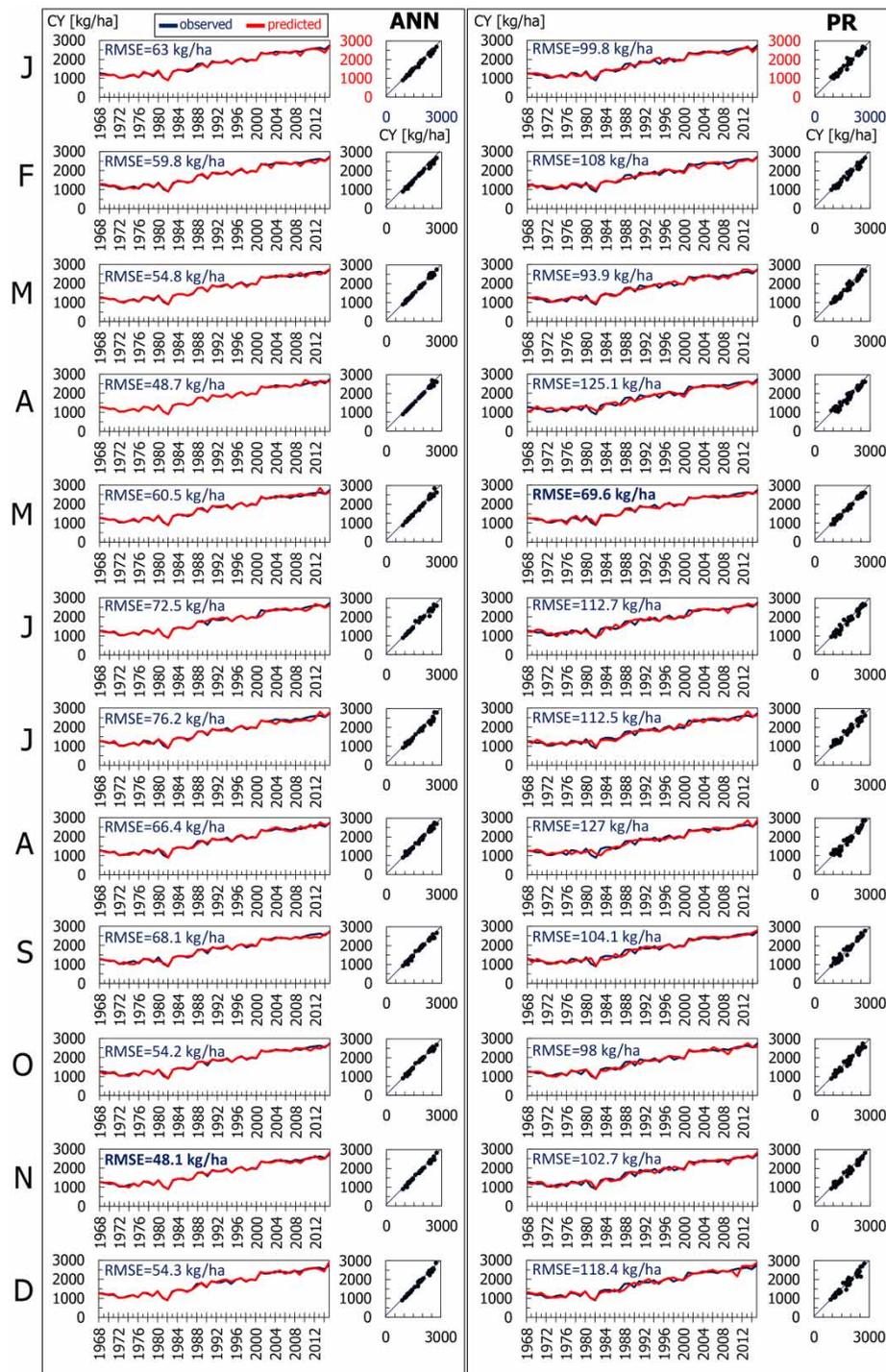


Figure 8 | ANN and PR models for predicting seasonal CY built using each time series of monthly DAs: region 2 (West Bengal). The model with the lowest error (RMSE) is presented for each month, from January to December (J to D).

corresponding to one of the various aggregation periods (DA1, DA3, DA6, DA9 or DA12) that produced good results was identified within the combinations of input variables. The input sets are comprised of multiple DAs corresponding to the various aggregation periods. Therefore, using multiple aggregation periods of drought indicator results in a better model performance.

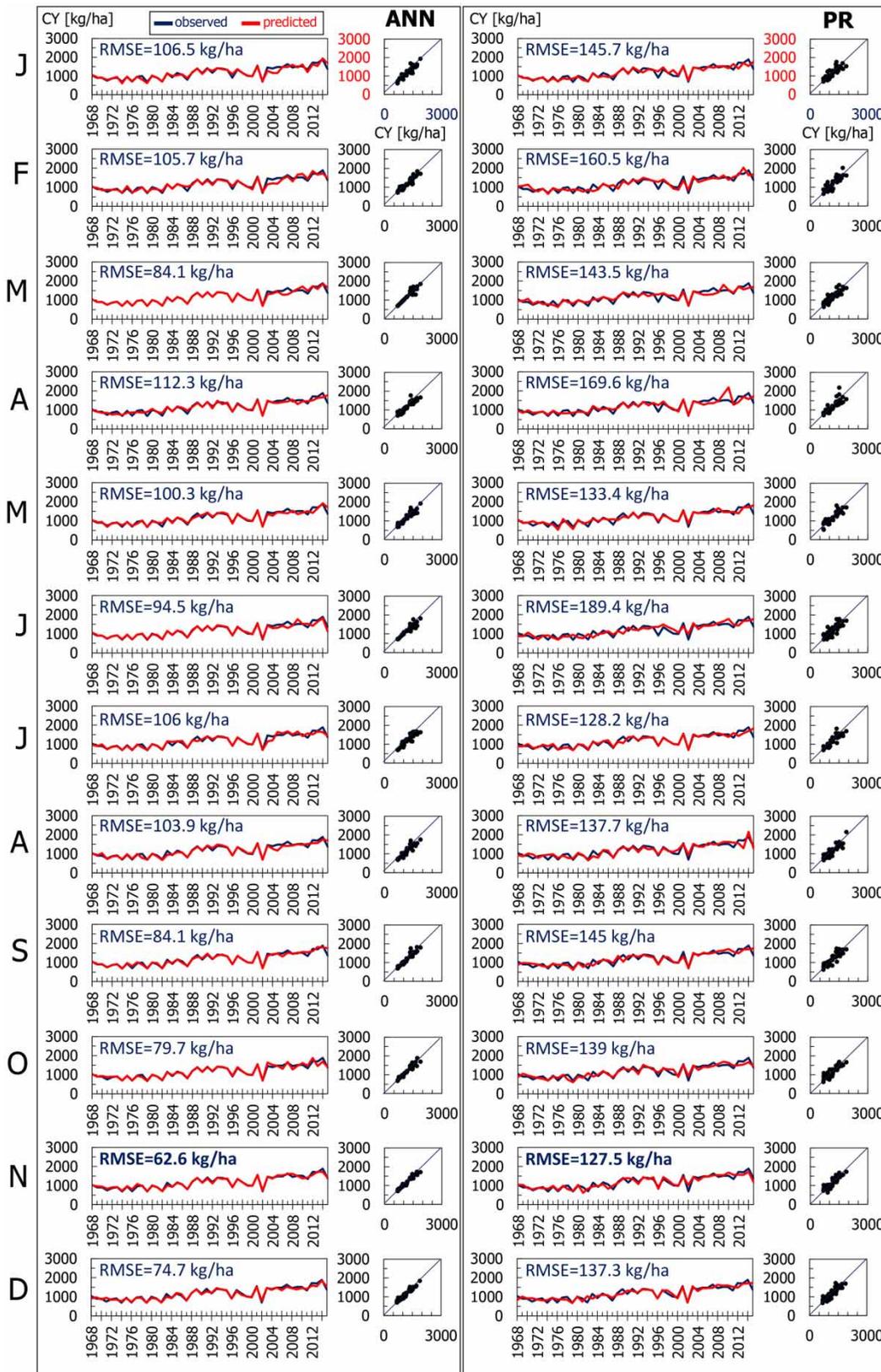


Figure 9 | ANN and PR models for predicting seasonal CY built using each time series of monthly DAs: region 3 (Odisha). The model with the lowest error (RMSE) is presented for each month, from January to December (J to D).

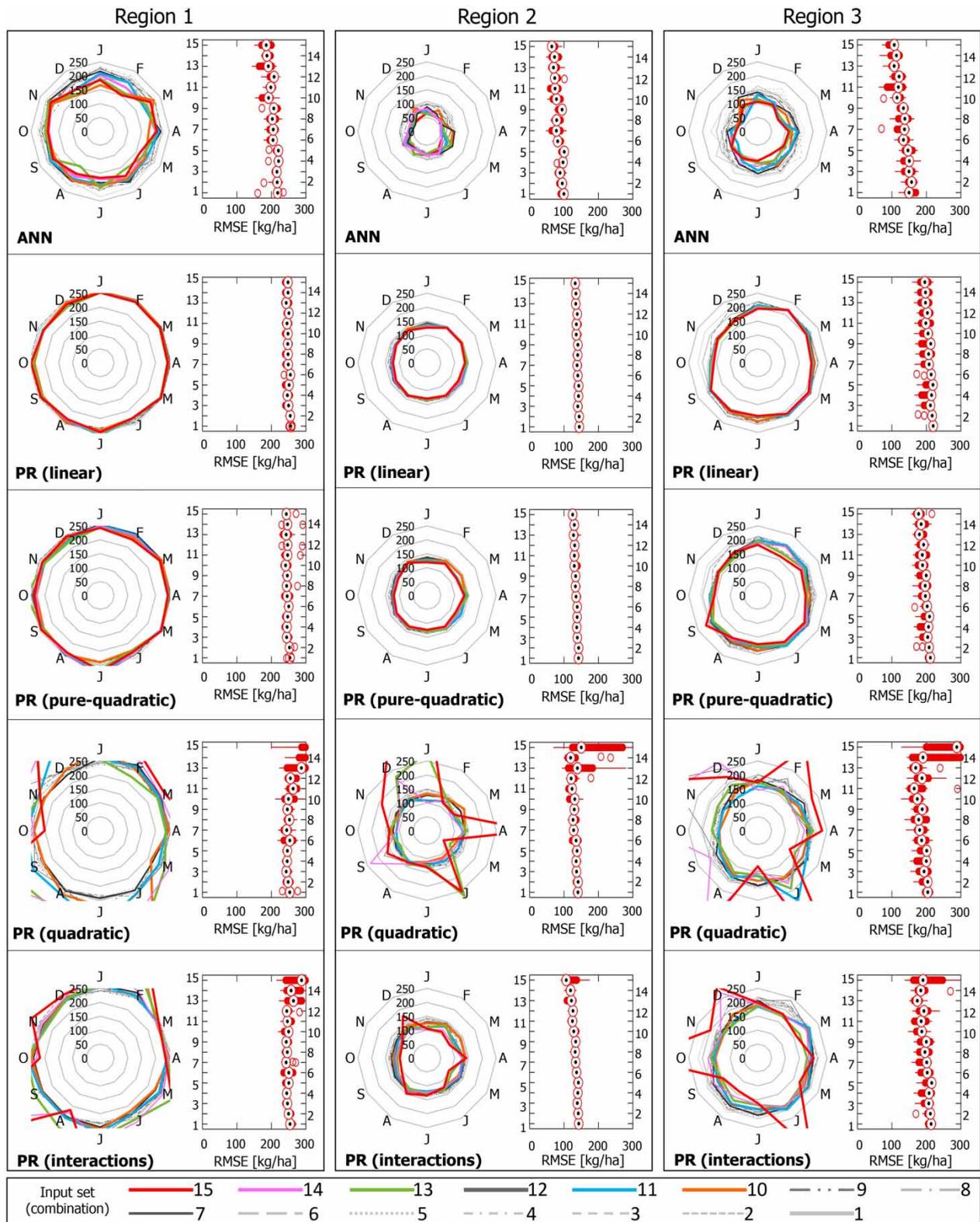


Figure 10 | RMSE [kg/ha] for each of the 15 input sets (combinations) of the ANN and PR models built for each region. For each set of inputs (from one to 15), the lowest RMSE values are presented for each month (January to December). The results of each input set are shown with lines to aid visual comparison. The left, middle, and right panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal), and region 3 (Odisha), respectively.

Table 3 | Summary of the ANN and PR models for predicting CY built for each month and region: (1) Bihar and Jharkhand, (2) West Bengal, and (3) Odisha

Region	ANN				PR						
	Month	Input set (combination)	No. nodes	RMSE [kg/ha]	Month	Input set (combination)	Type	RMSE [kg/ha]			
Region 1	Jan	10	CY _{t-1} , DA1,3,6	4	167.0	Jan	8	CY _{t-1} , DA6,9	quadratic	238.9	
	Feb	15	CY _{t-1} , DA1,3,6,9,12	6	151.2	Feb	13	CY _{t-1} , DA1,3,6,9	quadratic	223.9	
	Mar	11	CY _{t-1} , DA3,6,9	7	180.2	Mar	6	CY _{t-1} , DA1,3	quadratic	233.3	
	Apr	10	CY _{t-1} , DA1,3,6	9	175.8	Apr	15	CY _{t-1} , DA1,3,6,9,12	interactions	236.1	
	May	15	CY _{t-1} , DA1,3,6,9,12	5	162.0	May	10	CY _{t-1} , DA1,3,6	quadratic	211.7	
	Jun	13	CY _{t-1} , DA1,3,6,9	2	163.4	Jun	10	CY _{t-1} , DA1,3,6	interactions	219.9	
	Jul	15	CY _{t-1} , DA1,3,6,9,12	10	166.8	Jul	6	CY _{t-1} , DA1,3	quadratic	233.0	
	Aug	13	CY _{t-1} , DA1,3,6,9	5	144.0	Aug	15	CY _{t-1} , DA1,3,6,9,12	interactions	215.6	
	Sep	6	CY _{t-1} , DA1,3	5	177.9	Sep	7	CY _{t-1} , DA3,6	quadratic	232.4	
	Oct	14	CY _{t-1} , DA3,6,9,12	6	186.7	Oct	15	CY _{t-1} , DA1,3,6,9,12	quadratic	201.0	
	Nov	8	CY _{t-1} , DA6,9	4	178.5	Nov	13	CY _{t-1} , DA1,3,6,9	interactions	222.6	
	Dec	10	CY _{t-1} , DA1,3,6	4	153.3	Dec	13	CY _{t-1} , DA1,3,6,9	pure-quadratic	223.3	
Region 2	Jan	13	CY _{t-1} , DA1,3,6,9	8	63.0	Jan	14	CY _{t-1} , DA3,6,9,12	quadratic	99.8	
	Feb	11	CY _{t-1} , DA3,6,9	10	59.8	Feb	15	CY _{t-1} , DA1,3,6,9,12	interactions	108.0	
	Mar	7	CY _{t-1} , DA3,6	8	54.8	Mar	15	CY _{t-1} , DA1,3,6,9,12	interactions	93.9	
	Apr	14	CY _{t-1} , DA3,6,9,12	7	48.7	Apr	14	CY _{t-1} , DA3,6,9,12	interactions	125.1	
	May	15	CY _{t-1} , DA1,3,6,9,12	10	60.5	May	15	CY _{t-1} , DA1,3,6,9,12	quadratic	69.6	
	Jun	13	CY _{t-1} , DA1,3,6,9	7	72.5	Jun	10	CY _{t-1} , DA1,3,6	quadratic	112.7	
	Jul	6	CY _{t-1} , DA1,3	6	76.2	Jul	10	CY _{t-1} , DA1,3,6	quadratic	112.5	
	Aug	6	CY _{t-1} , DA1,3	9	66.4	Aug	13	CY _{t-1} , DA1,3,6,9	interactions	127.0	
	Sep	6	CY _{t-1} , DA1,3	10	68.1	Sep	15	CY _{t-1} , DA1,3,6,9,12	interactions	104.1	
	Oct	7	CY _{t-1} , DA3,6	10	54.2	Oct	15	CY _{t-1} , DA1,3,6,9,12	interactions	98.0	
	Nov	7	CY _{t-1} , DA3,6	10	48.1	Nov	15	CY _{t-1} , DA1,3,6,9,12	interactions	102.7	
	Dec	15	CY _{t-1} , DA1,3,6,9,12	8	54.3	Dec	14	CY _{t-1} , DA3,6,9,12	interactions	118.4	
Region 3	Jan	15	CY _{t-1} , DA1,3,6,9,12	7	106.5	Jan	14	CY _{t-1} , DA3,6,9,12	quadratic	145.7	
	Feb	13	CY _{t-1} , DA1,3,6,9	10	105.7	Feb	10	CY _{t-1} , DA1,3,6	quadratic	160.5	
	Mar	15	CY _{t-1} , DA1,3,6,9,12	9	84.1	Mar	12	CY _{t-1} , DA6,9,12	quadratic	143.5	
	Apr	15	CY _{t-1} , DA1,3,6,9,12	4	112.3	Apr	14	CY _{t-1} , DA3,6,9,12	quadratic	169.6	
	May	12	CY _{t-1} , DA6,9,12	10	100.3	May	15	CY _{t-1} , DA1,3,6,9,12	quadratic	133.4	
	Jun	15	CY _{t-1} , DA1,3,6,9,12	9	94.5	Jun	12	CY _{t-1} , DA6,9,12	quadratic	189.4	
	Jul	15	CY _{t-1} , DA1,3,6,9,12	7	106.0	Jul	15	CY _{t-1} , DA1,3,6,9,12	quadratic	128.2	
	Aug	12	CY _{t-1} , DA6,9,12	7	103.9	Aug	15	CY _{t-1} , DA1,3,6,9,12	interactions	137.7	
	Sep	11	CY _{t-1} , DA3,6,9	9	84.1	Sep	13	CY _{t-1} , DA1,3,6,9	quadratic	145.0	
	Oct	15	CY _{t-1} , DA1,3,6,9,12	10	79.7	Oct	10	CY _{t-1} , DA1,3,6	quadratic	139.0	
	Nov	11	CY _{t-1} , DA3,6,9	10	62.6	Nov	10	CY _{t-1} , DA1,3,6	quadratic	127.5	
	Dec	11	CY _{t-1} , DA3,6,9	9	74.7	Dec	8	CY _{t-1} , DA6,9	quadratic	137.3	

The table shows the models built with the lowest error (RMSE). DA stands for drought area.

Tables 4, 5 and 6 were created from Table 3. These three tables present the PR formulas for regions 1, 2, and 3. Each table shows the PR formula and the inputs used. These formulas are also designed to work independently for the CY prediction for each region.

The process of applying PR models begins by choosing the appropriate formula from Table 4, 5, or 6. For example, in region 1, if drought indicator data is available up to and including March, the formula for March is selected from Table 4. Afterwards, DAs are calculated (Section 3.1.2), and the DA time series are updated. In this case, Table 4 indicates that DA1 and DA3 are needed. From these DA time series, values for March are extracted, i.e., DA1_3 and DA3_3 (see Section 3.2 and 4.2). Then, the de-trending process is applied to each time series (Section 3.1.3). Next, CY is computed. Finally, the reverse de-trending is performed to ensure the predicted CY has the same order of magnitude as the original CY data (Section 3.1.3). Simultaneously, or when it becomes possible, the ANN model for the month under analysis is applied.

Table 4 | PR models for predicting CY built for each month: region 1 (Bihar and Jharkhand)

Month	Input						PR model
	x_1	x_2	x_3	x_4	x_5	x_6	
Jan	CY_{t-1}	DA6	DA9				$-60.7111 - 0.1944x_1 - 0.2201x_2 + 1.2033x_3 - 0.0023x_1x_2 + 0.0043x_1x_3 - 0.0372x_2x_3 + 0.0003x_1^2 + 0.0504x_2^2 + 0.0308x_3^2$
Feb	CY_{t-1}	DA1	DA3	DA6	DA9		$-27.4716 - 0.4688x_1 + 1.8718x_2 - 1.3313x_3 - 0.2611x_4 + 1.3878x_5 - 0.0137x_1x_2 + 0.0135x_1x_3 + 0.0032x_1x_4 + 0.0064x_1x_5 + 0.0823x_2x_3 + 0.0574x_2x_4 + 0.0935x_2x_5 - 0.0544x_3x_4 - 0.0746x_3x_5 - 0.0241x_4x_5 + 0.0014x_1^2 - 0.0496x_2^2 - 0.0202x_3^2 - 0.0016x_4^2 + 0.0227x_5^2$
Mar	CY_{t-1}	DA1	DA3				$28.1213 - 0.5204x_1 - 0.4908x_2 + 0.0545x_3 + 0.0051x_1x_2 - 0.0093x_1x_3 + 0.0033x_2x_3 + 0.0003x_1^2 - 0.0107x_2^2 + 0.0086x_3^2$
Apr	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$-24.3419 - 0.4785x_1 - 0.1965x_2 - 0.1356x_3 + 0.0848x_4 - 0.4774x_5 + 0.8029x_6 + 0.0066x_1x_2 + 0.0031x_1x_3 - 0.0128x_1x_4 + 0.0081x_1x_5 - 0.0003x_1x_6 + 0.0067x_2x_3 - 0.0604x_2x_4 + 0.1495x_2x_5 - 0.0169x_2x_6 + 0.0248x_3x_4 - 0.1295x_3x_5 - 0.0306x_3x_6 + 0.0458x_4x_5 + 0.0516x_4x_6 + 0.0595x_5x_6$
May	CY_{t-1}	DA1	DA3	DA6			$113.2521 - 0.5132x_1 + 1.0101x_2 - 1.4019x_3 - 1.1130x_4 + 0.0100x_1x_2 + 0.0150x_1x_3 - 0.0027x_1x_4 + 0.0250x_2x_3 - 0.0655x_2x_4 + 0.0596x_3x_4 - 0.0006x_1^2 - 0.0358x_2^2 - 0.0380x_3^2 - 0.0495x_4^2$
Jun	CY_{t-1}	DA1	DA3	DA6			$54.3 - 0.3715x_1 + 1.4832x_2 + 0.1432x_3 - 3.0648x_4 - 0.0106x_1x_2 + 0.0256x_1x_3 - 0.0111x_1x_4 - 0.0556x_2x_3 + 0.0648x_2x_4 - 0.0172x_3x_4$
Jul	CY_{t-1}	DA1	DA3				$18.7237 - 0.3166x_1 + 1.3310x_2 - 3.0099x_3 - 0.0030x_1x_2 + 0.0024x_1x_3 + 0.0054x_2x_3 + 0.0001x_1^2 + 0.0065x_2^2 - 0.0065x_3^2$
Aug	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$59.2373 - 0.6972x_1 + 0.1791x_2 + 5.1900x_3 - 1.3783x_4 - 6.9753x_5 + 1.5471x_6 - 0.0142x_1x_2 + 0.0072x_1x_3 + 0.1163x_1x_4 - 0.1285x_1x_5 + 0.0294x_1x_6 - 0.3670x_2x_3 + 0.0897x_2x_4 + 0.2332x_2x_5 + 0.0922x_2x_6 + 0.3014x_3x_4 + 0.3444x_3x_5 - 0.4160x_3x_6 - 0.5819x_4x_5 - 0.0450x_4x_6 + 0.3299x_5x_6$
Sep	CY_{t-1}	DA3	DA6				$44.8563 - 0.4565x_1 + 0.6884x_2 - 1.9466x_3 + 0.0053x_1x_2 - 0.0005x_1x_3 + 0.0012x_2x_3 + 0.0004x_1^2 - 0.0172x_2^2 - 0.0002x_3^2$
Oct	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$76.1546 + 0.0046x_1 - 2.2220x_2 + 1.0816x_3 + 19.1690x_4 - 53.2338x_5 + 29.1398x_6 + 0.0048x_1x_2 + 0.0155x_1x_3 - 0.0383x_1x_4 - 0.0868x_1x_5 + 0.1254x_1x_6 - 0.0444x_2x_3 + 0.0448x_2x_4 + 0.0175x_2x_5 - 0.0552x_2x_6 + 0.2154x_3x_4 - 1.0260x_3x_5 + 0.7776x_3x_6 + 3.2060x_4x_5 - 3.3267x_4x_6 + 11.6655x_5x_6 + 0.0002x_1^2 - 0.0547x_2^2 + 0.1171x_3^2 + 0.2874x_4^2 - 7.7995x_5^2 - 4.0845x_6^2$
Nov	CY_{t-1}	DA1	DA3	DA6	DA9		$30.0286 - 0.4536x_1 - 0.6721x_2 - 0.8270x_3 - 7.0981x_4 + 5.3007x_5 - 0.0339x_1x_2 + 0.0086x_1x_3 + 0.0107x_1x_4 - 0.0084x_1x_5 + 0.1347x_2x_3 + 0.1123x_2x_4 - 0.0596x_2x_5 + 0.2355x_3x_4 - 0.2262x_3x_5 - 0.0117x_4x_5$
Dec	CY_{t-1}	DA1	DA3	DA6	DA9		$29.2005 - 0.3816x_1 - 0.6953x_2 + 0.8469x_3 + 1.2024x_4 - 3.2563x_5 + 0.0005x_1^2 - 0.5339x_2^2 - 0.0047x_3^2 - 0.0119x_4^2 + 0.0083x_5^2$

For each month, the input (x_1 to x_6) and the PR formula are indicated. DA stands for drought area.

Table 5 | PR models for predicting CY built for each month: region 2 (West Bengal)

Month	Input						PR model
	x_1	x_2	x_3	x_4	x_5	x_6	
Jan	CY_{t-1}	DA3	DA6	DA9	DA12		$8.5606 - 0.2404x_1 - 1.1236x_2 - 0.7606x_3 + 6.6535x_4 - 5.3772x_5 + 0.0087x_1x_2 - 0.0044x_1x_3 - 0.0182x_1x_4 + 0.0234x_1x_5 + 0.0080x_2x_3 + 0.0234x_2x_4 - 0.0037x_2x_5 - 0.0402x_3x_4 + 0.1648x_3x_5 + 0.0200x_4x_5 + 0.0001x_1^2 - 0.0145x_2^2 - 0.0657x_3^2 + 0.0544x_4^2 - 0.0952x_5^2$
Feb	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$-24.8742 - 0.5460x_1 - 0.1190x_2 + 0.2175x_3 + 0.7776x_4 - 8.6335x_5 + 6.4022x_6 - 0.0164x_1x_2 + 0.0095x_1x_3 - 0.0251x_1x_4 + 0.0262x_1x_5 - 0.0057x_1x_6 - 0.0179x_2x_3 - 0.0241x_2x_4 - 0.1705x_2x_5 + 0.1579x_2x_6 + 0.0064x_3x_4 + 0.2383x_3x_5 - 0.2779x_3x_6 - 0.0117x_4x_5 + 0.0266x_4x_6 + 0.0614x_5x_6$
Mar	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$35.6904 - 0.3835x_1 - 0.9286x_2 + 0.1960x_3 - 0.3445x_4 - 0.3559x_5 + 0.6370x_6 - 0.0025x_1x_2 - 0.0009x_1x_3 + 0.0111x_1x_4 - 0.0252x_1x_5 + 0.0144x_1x_6 - 0.0059x_2x_3 + 0.0426x_2x_4 + 0.0063x_2x_5 + 0.0012x_2x_6 - 0.0362x_3x_4 - 0.1287x_3x_5 - 0.0038x_3x_6 + 0.0242x_4x_5 - 0.0355x_4x_6 + 0.0394x_5x_6$
Apr	CY_{t-1}	DA3	DA6	DA9	DA12		$8.5856 - 0.1865x_1 + 1.5824x_2 - 1.0816x_3 - 1.0256x_4 + 1.7846x_5 - 0.0164x_1x_2 + 0.0242x_1x_3 - 0.0013x_1x_4 + 0.0009x_1x_5 - 0.0084x_2x_3 + 0.0073x_2x_4 - 0.0710x_2x_5 - 0.0430x_3x_4 + 0.0659x_3x_5 + 0.0317x_4x_5$
May	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$-25.0101 - 0.8233x_1 - 1.8073x_2 + 1.1145x_3 + 1.6217x_4 + 0.9651x_5 + 0.5729x_6 + 0.0254x_1x_2 - 0.1198x_1x_3 + 0.0959x_1x_4 - 0.0112x_1x_5 + 0.0311x_1x_6 - 0.2178x_2x_3 + 0.3465x_2x_4 - 0.3214x_2x_5 + 0.0602x_2x_6 - 0.9192x_3x_4 + 1.2301x_3x_5 - 0.2167x_3x_6 - 0.8955x_4x_5 + 0.1015x_4x_6 + 0.0662x_5x_6 + 0.0048x_1^2 - 0.0096x_2^2 + 0.3527x_3^2 + 0.4308x_4^2 - 0.0492x_5^2 + 0.0639x_6^2$
Jun	CY_{t-1}	DA1	DA3	DA6			$90.7623 - 0.5785x_1 + 0.1582x_2 - 2.7914x_3 + 0.8655x_4 - 0.0176x_1x_2 + 0.0093x_1x_3 - 0.0108x_1x_4 + 0.0533x_2x_3 - 0.0521x_2x_4 + 0.1589x_3x_4 + 0.0012x_1^2 + 0.0072x_2^2 - 0.0974x_3^2 - 0.0714x_4^2$
Jul	CY_{t-1}	DA1	DA3	DA6			$26.1164 - 0.6892x_1 - 0.6723x_2 - 5.5280x_3 + 4.6922x_4 + 0.0070x_1x_2 + 0.0111x_1x_3 - 0.0148x_1x_4 - 0.1301x_2x_3 + 0.0838x_2x_4 + 0.5157x_3x_4 + 0.0014x_1^2 + 0.0679x_2^2 - 0.1671x_3^2 - 0.3540x_4^2$
Aug	CY_{t-1}	DA1	DA3	DA6	DA9		$55.6167 - 0.2284x_1 - 0.0182x_2 - 1.7996x_3 - 4.0674x_4 + 3.7965x_5 + 0.0117x_1x_2 - 0.0259x_1x_3 + 0.0556x_1x_4 - 0.0484x_1x_5 - 0.0176x_2x_3 - 0.1459x_2x_4 + 0.1017x_2x_5 - 0.0487x_3x_4 + 0.2346x_3x_5 - 0.1273x_4x_5$
Sep	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$35.6058 - 0.3263x_1 + 1.9755x_2 - 0.4197x_3 - 3.5963x_4 + 2.7383x_5 - 1.2234x_6 + 0.0013x_1x_2 - 0.0057x_1x_3 - 0.0470x_1x_4 + 0.0042x_1x_5 + 0.0475x_1x_6 + 0.0033x_2x_3 - 0.1889x_2x_4 + 0.0749x_2x_5 + 0.1060x_2x_6 + 0.0179x_3x_4 - 0.0003x_3x_5 + 0.0412x_3x_6 + 0.0291x_4x_5 - 0.0312x_4x_6 - 0.0379x_5x_6$
Oct	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$7.7675 - 0.1875x_1 - 0.1476x_2 - 0.8333x_3 - 5.1327x_4 + 15.3857x_5 - 10.6323x_6 - 0.0012x_1x_2 - 0.0011x_1x_3 + 0.0588x_1x_4 + 0.0365x_1x_5 - 0.0886x_1x_6 - 0.1339x_2x_3 + 0.1763x_2x_4 - 0.5955x_2x_5 + 0.4854x_2x_6 - 0.4231x_3x_4 - 0.2159x_3x_5 + 0.6868x_3x_6 + 0.3521x_4x_5 + 0.0666x_4x_6 - 0.4145x_5x_6$
Nov	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$38.3601 - 0.2443x_1 + 1.7236x_2 - 0.6584x_3 - 6.7484x_4 + 13.3609x_5 - 9.4895x_6 + 0.0114x_1x_2 + 0.0162x_1x_3 + 0.0331x_1x_4 - 0.0817x_1x_5 + 0.0478x_1x_6 + 0.0370x_2x_3 - 0.1350x_2x_4 - 0.0212x_2x_5 + 0.1631x_2x_6 - 0.1562x_3x_4 - 0.0082x_3x_5 + 0.1229x_3x_6 + 0.2672x_4x_5 - 0.0938x_4x_6 - 0.1335x_5x_6$
Dec	CY_{t-1}	DA3	DA6	DA9	DA12		$24.769 - 0.1091x_1 - 2.9747x_2 + 2.9990x_3 - 5.4144x_4 + 3.3374x_5 + 0.0083x_1x_2 - 0.0069x_1x_3 + 0.0596x_1x_4 - 0.0630x_1x_5 + 0.0755x_2x_3 + 0.0127x_2x_4 + 0.0094x_2x_5 - 0.0052x_3x_4 - 0.0884x_3x_5 + 0.0361x_4x_5$

For each month, the input (x_1 to x_6) and the PR formula are indicated. DA stands for drought area.

Table 6 | PR models for predicting CY built for each month: region 3 (Odisha)

Month	Input						PR model
	x_1	x_2	x_3	x_4	x_5	x_6	
Jan	CY_{t-1}	DA3	DA6	DA9	DA12		$-149.3429 - 0.4867x_1 - 1.5749x_2 + 2.0827x_3 + 5.9761x_4 - 6.0586x_5 - 0.0022x_1x_2 + 0.0100x_1x_3 + 0.0200x_1x_4 + 0.0045x_1x_5 - 0.0142x_2x_3 - 0.2414x_2x_4 + 0.1392x_2x_5 - 0.1332x_3x_4 + 0.1123x_3x_5 + 0.2083x_4x_5 + 0.0022x_1^2 + 0.0262x_2^2 + 0.0771x_3^2 + 0.0431x_4^2 - 0.1405x_5^2$
Feb	CY_{t-1}	DA1	DA3	DA6			$-90.6767 - 0.6674x_1 + 0.1283x_2 + 0.2580x_3 + 0.4540x_4 - 0.0041x_1x_2 + 0.0141x_1x_3 - 0.0009x_1x_4 + 0.0055x_2x_3 - 0.0195x_2x_4 + 0.0771x_3x_4 + 0.0006x_1^2 + 0.0313x_2^2 - 0.0207x_3^2 + 0.0129x_4^2$
Mar	CY_{t-1}	DA6	DA9	DA12			$-168.6741 - 0.7249x_1 + 0.2079x_2 - 2.2594x_3 + 2.2421x_4 + 0.0074x_1x_2 - 0.0102x_1x_3 + 0.0347x_1x_4 - 0.0159x_2x_3 + 0.0009x_2x_4 + 0.1147x_3x_4 + 0.0025x_1^2 + 0.0454x_2^2 - 0.0197x_3^2 + 0.0318x_4^2$
Apr	CY_{t-1}	DA3	DA6	DA9	DA12		$-116.7973 - 0.6789x_1 - 0.4066x_2 - 0.5459x_3 + 3.4428x_4 - 3.2126x_5 + 0.0008x_1x_2 - 0.0110x_1x_3 + 0.0063x_1x_4 + 0.0337x_1x_5 + 0.0647x_2x_3 - 0.1280x_2x_4 + 0.0847x_2x_5 - 0.0041x_3x_4 - 0.1576x_3x_5 - 0.0357x_4x_5 + 0.0025x_1^2 - 0.0386x_2^2 + 0.0180x_3^2 + 0.0968x_4^2 + 0.1431x_5^2$
May	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$-56.0895 - 0.8435x_1 - 1.5688x_2 + 5.5848x_3 - 5.6556x_4 - 0.0876x_5 - 0.4449x_6 + 0.0396x_1x_2 - 0.0552x_1x_3 + 0.0130x_1x_4 + 0.0414x_1x_5 - 0.0155x_1x_6 + 0.0691x_2x_3 - 0.1386x_2x_4 + 0.4106x_2x_5 + 0.0874x_2x_6 + 0.2997x_3x_4 - 0.2552x_3x_5 - 0.4282x_3x_6 - 0.0482x_4x_5 + 0.2264x_4x_6 - 0.2702x_5x_6 + 0.0040x_1^2 - 0.0721x_2^2 - 0.0198x_3^2 - 0.2076x_4^2 + 0.2160x_5^2 - 0.0223x_6^2$
Jun	CY_{t-1}	DA6	DA9	DA12			$-23.8562 - 0.3639x_1 - 1.8924x_2 - 0.0052x_3 + 1.3074x_4 - 0.0060x_1x_2 - 0.0057x_1x_3 + 0.0205x_1x_4 - 0.0135x_2x_3 - 0.0965x_2x_4 + 0.1034x_3x_4 + 0.0004x_1^2 + 0.0110x_2^2 - 0.0171x_3^2 + 0.0913x_4^2$
Jul	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$-18.8884 - 0.7725x_1 + 2.8997x_2 - 1.9129x_3 - 0.9194x_4 - 0.5636x_5 - 0.6886x_6 - 0.0070x_1x_2 + 0.0320x_1x_3 - 0.0220x_1x_4 - 0.0221x_1x_5 - 0.0042x_1x_6 + 0.3776x_2x_3 - 0.0748x_2x_4 - 0.1803x_2x_5 - 0.2590x_2x_6 - 0.5984x_3x_4 + 0.6811x_3x_5 - 0.0178x_3x_6 + 0.8957x_4x_5 + 0.0173x_4x_6 - 0.1524x_5x_6 + 0.0012x_1^2 - 0.1151x_2^2 - 0.1006x_3^2 - 0.0306x_4^2 - 0.7603x_5^2 + 0.1200x_6^2$
Aug	CY_{t-1}	DA1	DA3	DA6	DA9	DA12	$4.8997 - 0.7900x_1 - 0.9225x_2 + 3.8372x_3 - 0.0832x_4 - 9.7835x_5 + 4.0199x_6 - 0.0065x_1x_2 + 0.0352x_1x_3 + 0.0005x_1x_4 - 0.0461x_1x_5 - 0.0019x_1x_6 - 0.0759x_2x_3 - 0.1196x_2x_4 + 0.1775x_2x_5 + 0.0748x_2x_6 + 0.0694x_3x_4 + 0.2503x_3x_5 - 0.3715x_3x_6 - 0.2022x_4x_5 + 0.4167x_4x_6 - 0.2192x_5x_6$
Sep	CY_{t-1}	DA1	DA3	DA6	DA9		$41.4745 - 0.5431x_1 - 0.0366x_2 + 0.9681x_3 + 3.6023x_4 - 4.3272x_5 - 0.0002x_1x_2 + 0.0115x_1x_3 - 0.0191x_1x_4 + 0.0139x_1x_5 - 0.0809x_2x_3 + 0.0508x_2x_4 + 0.0205x_2x_5 + 0.4602x_3x_4 - 0.5016x_3x_5 + 0.3000x_4x_5 + 0.0002x_1^2 + 0.0172x_2^2 - 0.0339x_3^2 - 0.3409x_4^2 + 0.0831x_5^2$
Oct	CY_{t-1}	DA1	DA3	DA6			$-48.806 - 0.6966x_1 - 0.4241x_2 - 1.7664x_3 - 3.0097x_4 + 0.0040x_1x_2 + 0.0053x_1x_3 - 0.0175x_1x_4 - 0.0038x_2x_3 + 0.0111x_2x_4 - 0.1443x_3x_4 + 0.0008x_1^2 + 0.0073x_2^2 + 0.0861x_3^2 + 0.0558x_4^2$
Nov	CY_{t-1}	DA1	DA3	DA6			$47.8316 - 0.6925x_1 + 0.7765x_2 - 2.3671x_3 - 2.9813x_4 + 0.0043x_1x_2 + 0.0011x_1x_3 - 0.0066x_1x_4 + 0.0797x_2x_3 - 0.0306x_2x_4 - 0.0144x_3x_4 + 0.0004x_1^2 - 0.0064x_2^2 - 0.0407x_3^2 + 0.0200x_4^2$
Dec	CY_{t-1}	DA6	DA9				$13.0378 - 0.5111x_1 + 0.5765x_2 - 3.4820x_3 + 0.0177x_1x_2 - 0.0158x_1x_3 + 0.0155x_2x_3 + 0.0004x_1^2 - 0.0691x_2^2 + 0.0343x_3^2$

For each month, the input (x_1 to x_6) and the PR formula are indicated. DA stands for drought area.

4.5. Modelling limitations

The modelling limitations of the presented approach include the following.

- (1) To identify drought-affected areas, a threshold value of $SPEI \leq -1$ was employed. Relying on a single threshold might lead to overestimating the actual impacts of drought on crop yield.
- (2) The gridded SPEI data at a spatial resolution of $0.5^\circ \times 0.5^\circ$ was utilised for each region individually. Such a coarse spatial resolution across different region sizes might not accurately delineate drought areas, resulting in over- or underestimations of their extent.
- (3) Due to limitations in the yield data – specifically, the availability of only time series information – irrigated and rain-fed areas could not be modelled separately. Using spatially explicit crop data that distinguish between irrigated and rain-fed areas, along with finer-resolution drought areas data, could improve the results.
- (4) This study assumes drought as the sole causative factor; however, floods also negatively affect crop yield and overall regional production. Flood impacts are not incorporated into the models.
- (5) Additional factors such as market conditions, technological advancements, and management practices may also influence rice yields. Nonetheless, this study presumes drought to be the most significant factor.
- (6) Limited crop yield data was available for model development, with only the Kharif season data used, resulting in the crop yield time series containing a single value for each year.

4.6. Consideration of other factors, drought types, and indices

While soil moisture is widely recognised as a key variable for assessing agricultural drought (Carrão *et al.* 2016; Zhu *et al.* 2020), its use is often limited by data availability and spatial resolution, especially in local or regional contexts. In such cases, proxy indicators like the SPEI can offer a practical alternative for estimating drought impacts (Osman 2018; Osman *et al.* 2018; Araneda-Cabrera *et al.* 2021; Khoshnazar *et al.* 2022).

Our study shows that DA derived from SPEI can effectively act as a predictor of crop yield variability. The ML approach presented here can be further explored by examining other drought indicators (e.g., Hao & AghaKouchak 2013; Khoshnazar *et al.* 2022) to investigate any additional improvement in crop yield prediction.

Incorporating additional factors, such as irrigation practices, soil properties, and farmer practices, would further improve prediction models. However, such data are often difficult to gather or standardise. As an alternative, we suggest modifying the drought area input by applying spatial or temporal weights derived from these auxiliary factors. For example, drought areas could be adjusted – expanded or contracted – based on access to irrigation, land management practices, or socioeconomic vulnerability. These weighted drought areas could serve as more refined inputs for ML models.

Moreover, future implementations could divide the study region into spatial units (e.g., grid cells) to enable spatially resolved ML models. Advances in remote sensing now make it possible to access and process such high-resolution data. Combining anthropogenic and natural variables – either separately or together – through feature selection and classification could also help identify the most influential drivers of crop performance.

Lastly, our modelling framework can be adapted to simulate climate change scenarios, providing insights into how future drought patterns may impact crop productivity and guiding the development of more resilient agricultural strategies.

5. SUMMARY AND CONCLUSIONS

This research introduces an ML approach for predicting CY using DA data as input. The approach was applied to three regions in India to predict rice yield. The framework is based on ANN models, with PR models serving as a baseline. Together, ANN and PR models form an integrated tool for crop yield prediction.

ANN models can be used independently or in combination with PR models, as proposed in this study, allowing for performance comparison and validation. Moreover, when the spatial input data required to calculate drought areas inputs for ANN models are not yet available, the PR models can be applied using early estimates of drought area data, enabling timely and continuous early-season crop yield predictions.

The following conclusions have been drawn from this research.

- Based on the performance of PR and ANN models, results indicate that drought area is a suitable variable for predicting crop yield outcomes.

- The correlation analysis between DA and CY showed high negative correlations in Odisha (region 3). The correlation gradually decreases in Bihar and Jharkhand (region 1) and West Bengal (region 2). These correlation values may be because West Bengal has better access to irrigation facilities than Odisha and Bihar and Jharkhand.
- When comparing ANN models and PR models, the ANN was more accurate than PR models in predicting crop yield across all regions. This is expected since the drought–crop relationship is a highly non-linear problem.
- It can be concluded that ANN can predict CY in the pre-harvesting stage with good accuracy, given the drought indicator (SPEI), which uses climate variables such as precipitation and temperature (for evapotranspiration calculation).

Based on the analysis and findings of this research, the following recommendations can be made for further improvement of the study.

- Sensitivity analysis should be carried out to identify the parameters that can influence the model results. For instance, varying the spatial resolution of the drought indicator and using different drought indicator thresholds should be investigated.
- Extreme wet events should be considered, especially in flood-prone regions such as the coastal areas of West Bengal (region 2), Odisha (region 3), and North Bihar (region 1), where floods also affect crop yield.
- Non-climatic factors such as economic conditions, fertiliser use, and management practices, may also be considered, as they influence crop yield.
- To enhance the accuracy of the model, more input data should be incorporated in future research. For CY, this can be assessed through remote sensing techniques on a monthly basis so that the machine learning models can be constructed for this temporal resolution and the spatial coverage can be improved.
- Other ML models should be examined further, especially committee (ensemble) methods like random forests or boosting techniques. For data at scales smaller than monthly, deep learning algorithms (e.g., LSTM networks) may also be recommended for exploration.

We envision that this research will enhance drought monitoring systems for evaluating the impacts of drought. Although these systems can currently calculate drought areas, the direct application of drought impact predictions can be integrated using approaches such as the one presented or similar ones.

ACKNOWLEDGEMENTS

VD thanks the Mexican National Council for Science and Technology (CONACYT) and Alianza FiiDEM for the study grant 217776/382365. AAAO was supported by the Orange Knowledge Programme (former NFP) and the World Meteorological Organization (WMO). GACP and VD acknowledge the grant No. 2579 of the Prince Albert II of Monaco Foundation. DS acknowledges the grant No. 17-77-30006 of the Russian Science Foundation, and the Hydroinformatics Research Fund of IHE Delft, in whose framework some research ideas and components were developed.

LANGUAGE EDITING STATEMENT

This manuscript benefited from grammar and language improvements using Grammarly and OpenAI's ChatGPT. These tools were used exclusively to assist with editing; they were not used for generating content, conducting analyses, or interpreting results.

DATA AVAILABILITY STATEMENT

All relevant data are available from an online repository or repositories: <https://doi.org/10.5281/zenodo.18404801>.

CONFLICT OF INTEREST

The authors declare there is no conflict.

REFERENCES

- Araneda-Cabrera, R. J., Bermúdez, M. & Puertas, J. (2021) Assessment of the performance of drought indices for explaining crop yield variability at the national scale: methodological framework and application to Mozambique, *Agricultural Water Management*, 246 (September 2020), 106692. <https://doi.org/10.1016/j.agwat.2020.106692>.

- Below, R., Grover-Kopec, E. & Dilley, M. (2007) Documenting drought-related disasters: a global reassessment, *Journal of Environment and Development*, **16** (3), 328–344. <https://doi.org/10.1177/1070496507306222>.
- Bhalme, H. N. & Mooley, D. a. (1980) Large-Scale droughts/Floods and monsoon circulation, *Monthly Weather Review*, **108** (8), 1197–1211. [https://doi.org/10.1175/1520-0493\(1980\)108<1197:LSDAMC>2.0.CO;2](https://doi.org/10.1175/1520-0493(1980)108<1197:LSDAMC>2.0.CO;2).
- Carrão, H., Russo, S., Sepulcre-Canto, G. & Barbosa, P. (2016) An empirical standardized soil moisture index for agricultural drought assessment from remotely sensed data, *International Journal of Applied Earth Observation and Geoinformation*, **48**, 74–84. <https://doi.org/10.1016/j.jag.2015.06.011>.
- Chlingaryan, A., Sukkarieh, S. & Whelan, B. (2018) Machine learning approaches for crop yield prediction and nitrogen status estimation in precision agriculture: a review, *Computers and Electronics in Agriculture*, **151** (May), 61–69. <https://doi.org/10.1016/j.compag.2018.05.012>.
- Corzo Perez, G. A., van Huijgevoort, M. H. J., Voß, F. & van Lanen, H. A. J. (2011) On the spatio-temporal analysis of hydrological droughts from global hydrological models, *Hydrology and Earth System Sciences*, **15** (9), 2963–2978. <https://doi.org/10.5194/hess-15-2963-2011>.
- Diaz, V., Corzo Perez, G. A., Van Lanen, H. A. J. & Solomatine, D. (2016) Spatio-temporal analysis of large-scale meteorological drought: helping to achieve the SDGs 6.A and 11.5, *12th Kovacs Colloquium*. Paris, France: UNESCO-IHP. <https://doi.org/10.13140/RG.2.1.2595.2888>.
- Diaz, V., Corzo, G., Van Lanen, H. A. J. & Solomatine, D. P. (2019) Spatiotemporal drought analysis at country scale through the application of the STAND toolbox. In: Corzo, G. & Varouchakis, E. A. (eds). *Spatiotemporal Analysis of Extreme Hydrological Events*. Amsterdam: Elsevier, pp. 77–93. <https://doi.org/10.1016/B978-0-12-811689-0.00004-5>.
- Diaz, V., Corzo Perez, G. A., Van Lanen, H. A. J., Solomatine, D. & Varouchakis, E. A. (2020) An approach to characterise spatio-temporal drought dynamics, *Advances in Water Resources*, **137**, 103512. <https://doi.org/10.1016/j.advwatres.2020.103512>.
- Directorate of Rice Development (DRD) (2014) *Rice – A Status Paper*. Patna, India: Government of India, Ministry of Agriculture (Department of Agriculture & Cooperation), Directorate of Rice Development, 144 pp.
- Elshorbagy, A., Corzo, G., Srinivasulu, S. & Solomatine, D. P. (2010) Experimental investigation of the predictive capabilities of data driven modeling techniques in hydrology – part 2: application, *Hydrology and Earth System Sciences*, **14** (10), 1943–1961. <https://doi.org/10.5194/hess-14-1943-2010>.
- Food and Agriculture Organization of the United Nations (FAO) (2017) *The Impact of Disasters and Crises on Agriculture and Food Security*. Rome: FAO. Available at: www.fao.org/publications.
- Food and Agriculture Organization of the United Nations (FAO) and Robert B Daugherty Water for Food Institute at the University of Nebraska (2015) *Yield gap analysis of field crops, Methods and case studies*. In: Sadras, V. O., Cassman, K. G. G., Grassini, P., Hall, A. J., Bastiaanssen, W. G. M., Laborte, A. G. & Steduto, P. (eds.) *FAO Water Reports*, Vol. 41, Rome: FAO.
- Ghosh, K., Balasubramanian, R., Bandopadhyay, S., Chattopadhyay, N., Singh, K. K. & Rathore, L. S. (2014) Development of crop yield forecast models under FASAL-a case study of kharif rice in West Bengal, *Journal of Agrometeorology*, **16** (1), 1–8. doi:10.54386/jam.v16i1.1479.
- Govindaraju, R. S. (2000) Artificial neural networks in hydrology. I: preliminary concepts. *Journal of Hydrologic Engineering*, **5** (2), 115–123. [https://doi.org/10.1061/\(ASCE\)1084-0699\(2000\)5:2\(115\)](https://doi.org/10.1061/(ASCE)1084-0699(2000)5:2(115)).
- Guha-Sapir, D. (2019) EM-DAT: The Emergency Events Database – Université catholique de Louvain (UCL) – CRED. Available at: www.emdat.be.
- Hao, Z. & AghaKouchak, A. (2013) Multivariate standardized drought index: a parametric multi-index model, *Advances in Water Resources*, **57**, 12–18. <https://doi.org/10.1016/j.advwatres.2013.03.009>.
- Herrera-Estrada, J. E., Satoh, Y. & Sheffield, J. (2017) Spatio-temporal dynamics of global drought, *Geophysical Research Letters*, **44**, 2254–2263. <https://doi.org/10.1002/2016GL071768>.
- Huang, J., Gómez-Dans, J. L., Huang, H., Ma, H., Wu, Q., Lewis, P. E., Liang, S., Chen, Z., Xue, J. H., Wu, Y., Zhao, F., Wang, J. & Xie, X. (2019) Assimilation of remote sensing into crop growth models: current status and perspectives, *Agricultural and Forest Meteorology*, **276–277** (July), 107609. <https://doi.org/10.1016/j.agrformet.2019.06.008>.
- Khoshnazar, A., Corzo Perez, G. A., Diaz, V., Aminzadeh, M. & Cerón Pineda, R. A. (2022) Wet-environment evapotranspiration and precipitation standardized index (WEPSI) for drought assessment and monitoring, *Hydrology Research*, **53** (11), 1393–1413. <https://doi.org/10.2166/nh.2022.062>.
- Kim, W., Iizumi, T. & Nishimori, M. (2019) Global patterns of crop production losses associated with droughts from 1983 to 2009, *Journal of Applied Meteorology and Climatology*, **58** (6), 1233–1244. <https://doi.org/10.1175/JAMC-D-18-0174.1>.
- Maier, H. R. & Dandy, G. C. (2000) Neural networks for the prediction and forecasting of water resources variables: a review of modelling issues and applications, *Environmental Modelling & Software*, **15** (1), 101–124. [https://doi.org/10.1016/S1364-8152\(99\)00007-9](https://doi.org/10.1016/S1364-8152(99)00007-9).
- Mckee, T. B., Doesken, N. J. & Kleist, J. (1993) The relationship of drought frequency and duration to time scales, *AMS 8th Conf. Appl. Climatol.* January. Anaheim, CA: American Meteorological Society, pp. 179–184. <https://doi.org/citeulike-article-id:10490403>.
- Monfreda, C., Ramankutty, N. & Foley, J. A. (2008) Farming the planet: 2. geographic distribution of crop areas, yields, physiological types, and net primary production in the year 2000, *Global Biogeochemical Cycles*, **22** (1), 1–19. <https://doi.org/10.1029/2007GB002947>.
- Montesino Pouzols, F. & Lendasse, A. (2010) Effect of different detrending approaches on computational intelligence models of time series, *The 2010 International Joint Conference on Neural Networks (IJCNN)*. Barcelona: IEEE, pp. 1–8. <https://doi.org/10.1109/IJCNN.2010.5596314>.

- Osman, A. A. A. (2018) *Spatiotemporal Analysis and Prediction of Crop Yield Using Data-Driven Models and Drought Areas: Case Study of India*. Delft, The Netherlands. <https://ihedelftrepository.contentdm.oclc.org/digital/collection/masters1/id/309029/rec/4>. Available at: <https://ihedelftrepository.contentdm.oclc.org/digital/collection/masters1/id/309029/rec/4>.
- Osman, A. A. A., Diaz, V., Corzo Perez, G. A., Van Lanen, H. A. J. & Solomatine, D. (2018) Finding negative response of crop yield to drought: a spatiotemporal approach over East India, *International Conference on Water, Environment, Energy and Society (ICWEES)*. Djerba Island, Tunisia: UN Office for Disaster Risk Reduction.
- Rahmati, O., Falah, F., Dayal, K. S., Deo, R. C., Mohammadi, F., Biggs, T., Moghaddam, D. D., Naghibi, S. A. & Bui, D. T. (2020) Machine learning approaches for spatial modeling of agricultural droughts in the south-east region of Queensland Australia, *Science of the Total Environment*, **699**, 134230. <https://doi.org/10.1016/j.scitotenv.2019.134230>.
- Reynolds, C. A., Yitayew, M., Slack, D. C., Hutchinson, C. F., Huete, A. & Petersen, M. S. (2000) Estimating crop yields and production by integrating the FAO crop specific water balance model with real-time satellite data and ground-based ancillary data, *International Journal of Remote Sensing*, **21** (18), 3487–3508. <https://doi.org/10.1080/014311600750037516>.
- Sheffield, J. & Wood, E. F. (2011) *Drought: Past Problems and Future Scenarios*. London: Earthscan.
- Udmale, P., Ichikawa, Y., Ning, S., Shrestha, S. & Pal, I. (2020) A statistical approach towards defining national-scale meteorological droughts in India using crop data, *Environmental Research Letters*, **15** (9), 094090. <https://doi.org/10.1088/1748-9326/abacfa>.
- van Klompenburg, T., Kassahun, A. & Catal, C. (2020) Crop yield prediction using machine learning: a systematic literature review, *Computers and Electronics in Agriculture*, **177** (July), 105709. <https://doi.org/10.1016/j.compag.2020.105709>.
- Vicente-Serrano, S. M., Beguería, S. & López-Moreno, J. I. (2010) A multiscalar drought index sensitive to global warming: the standardized precipitation evapotranspiration index, *Journal of Climate*, **23** (7), 1696–1718. <https://doi.org/10.1175/2009JCLI2909.1>.
- White, M. A., Thornton, P. E. & Running, S. W. (1997) A continental phenology model for monitoring vegetation responses to interannual climatic variability, *Global Biogeochemical Cycles*, **11** (2), 217–234. <https://doi.org/10.1029/97GB00330>.
- World Meteorological Organization (WMO) (2006) *Drought Monitoring and Early Warning: Concepts, Progress and Future Challenges*. WMO-No. 1006. Geneva, Switzerland: WMO. Available at: http://www.droughtmanagement.info/literature/WMO_drought_monitoring_early_warning_2006.pdf.
- Wu, X., Vuichard, N., Ciais, P., Viovy, N., Wang, X., Magliulo, V. & Wattenbach, M. (2016) ORCHIDEE-CROP (v0), a new process-based agro-land surface model: model description and evaluation over Europe, *Geoscientific Model Development*, **9** (2), 857–873. <https://doi.org/10.5194/gmd-9-857-2016>.
- Zhu, X., Hou, C., Xu, K. & Liu, Y. (2020) Establishment of agricultural drought loss models: a comparison of statistical methods, *Ecological Indicators*, **112** (July), 106084. <https://doi.org/10.1016/j.ecolind.2020.106084>.

First received 23 September 2024; accepted in revised form 23 December 2025. Available online 21 January 2026