# Machine-learning approach to crop yield prediction with the spatial extent of drought

Vitali Diaz[1,2], Ahmed A.A. Osman[3], Gerald A. Corzo Perez[1,2], Henny A.J. Van Lanen[4],
Shreedhar Maskey[2], Dimitri Solomatine[1,2,5]

[1]IHE Delft Institute for Water Education, Hydroinformatics Chair group, Delft, 2601 DA, the Netherlands

[2]Delft University of Technology, Delft, the Netherlands

[3]Arcadis, Wales, United Kingdom

[4]Hydrology and Quantitative Water Management Group, Wageningen University, Wageningen, the Netherlands

[5]Water Problems Institute of the Russian Academy of Sciences, Moscow, Russia

**Corresponding author**: Vitali Diaz; v.diazmercado@tudelft.nl; vitalidime@gmail.com

## Abstract

Crop yield is one of the variables used to assess the impact of droughts on agriculture. Crop growth models calculate yield and variables related to plant development and become more suitable for crop yield estimation. However, these models are limited in that specific data are needed for computation. Given this limitation, machine learning (ML) models are often widely utilised instead, but their use with the spatial characteristics of droughts as input data is limited. This research explored the spatial extent of drought (area) as input data for building an approach to predict seasonal crop yield. This ML approach is made up of two components. The first includes polynomial regression (PR) models, and the second considers artificial neural network (ANN) models. In this approach, the purpose was to evaluate both types of ML models (PR and ANN) and integrate them into one operational tool. The logic is as follows: ANN models determine the most accurate predictions, but in practice, issues regarding data retrieval and processing can make the use of equations, i.e. PR, preferable. The proposed approach provides these PR equations to perform such calculations with early and preliminary input. The estimates can be further improved when the ANN models are run with the final input data. The results indicated that the empirical equations (PR) produced good predictions when using drought area as the input. ANN provides better estimates, in general. This research will improve drought monitoring systems for assessing drought effects. Since it is currently possible to calculate drought areas within these systems, the direct application of the prediction of drought effects is possible to integrate by following approaches such as the one presented or similar.

## Keywords

Spatio-temporal analysis, crop yield, drought impact, machine learning, agricultural drought

## 1 Introduction

Drought continually hits many regions across the world. It negatively affects various human activities such as agriculture, which not only generates economic losses but can also trigger

35  famine, causing millions of deaths (Below et al., 2007; Food and Agriculture Organization of

36  the United Nations (FAO), 2017; Kim et al., 2019; Sheffield and Wood, 2011; World

37  Meteorological Organization (WMO), 2006). Hence, methods that help to improve strategies

38  for drought mitigation are necessary. Within these methods are those that allow predicting the

39  impacts of drought.

40  Assessments of drought impacts confirm that the presence of drought on human activities can

41  be devastating. For instance, the Food and Agriculture Organization of the United Nations

42  (FAO) calculated the damage and losses in the agricultural sector caused by five types of

43  hazards, including drought. FAO estimates that drought causes damages and losses to this

44  sector by up to 80% (FAO, 2017). Although multiple factors are involved in agriculture

45  affectation, drought often plays the primary role, as literature confirms (Dai, 2011; FAO, 2017;

46  Kim et al., 2019).

47  The assessment of drought impacts on agriculture can be performed in terms of crop yield.

48  FAO defines crop yield as the measure of the yield of a crop per unit area of land cultivation

49  (in kg/ha or ton/ha) (FAO and DWFI, 2015). For assessing crop yield under drought affectation,

50  physical models based on crop properties turn out to be more comprehensive and descriptive

51  (Huang et al., 2019; Reynolds et al., 2000; White et al., 1997; Wu et al., 2016). However, an

52  important barrier to such models' realisation is the lack of detailed crop data and the difficulty

53  representing all the processes involved in all stages of crop development (Huang et al., 2019;

54  Reynolds et al., 2000; Wu et al., 2016).

55  Statistical and machine-learning (ML) models, which involve mathematical equations to

56  calculate the output of a model with suitable input(s), can be used to assess crop yield impact

57  by drought without considering any biological or physical process of the crop but the analysis

58  of the input and output data (Chlingaryan et al., 2018; Rahmati et al., 2020; Udmale et al.,

59  2020; van Klompenburg et al., 2020). There have been studies where various inputs, ML

60  techniques and architectures (configurations) have been tested for crop yield prediction (e.g.,

61  Chlingaryan et al., 2018; van Klompenburg et al., 2020). However, the spatial extent of drought

62  (area) is an input that has not been fully explored previously to crop yield prediction. The

63  prediction refers to the calculation of crop yield at the end of the growing season (harvesting)

64  with information available before or during the crop development season (pre-harvesting).

65  This research aims to develop an ML approach to calculate seasonal crop yield (CY) with the

66  monthly drought areas (DAs) as input. The ML approach comprises two components. Each

67  component includes a set of the following types of ML models: polynomial regression (PR)

68  and artificial neural network (ANN). The goal is to build both types of ML models (ANN and

69    MR) and use them as an integrated tool to support the decisions made based on crop yield

70    prediction. The logic is as follows. PR provides the prediction where the crop yield calculation

71    is "clear" to the performer (the end-user) because she/he has access to the equations that have

72    a straightforward interpretation and calculations can be done with early and preliminary input

73    data. For its part, ANN is used as the most accurate model, although the output calculation is

74    not as "clear" as in the case of PR due to the difficulty of interpreting the structure of the

75    resulting ANN. ANN are expected to be used with the final input data.

76    Three East Indian regions where agriculture plays an important role were chosen as a case

77    study. ML models were built for the period 1967-2015. ML models aim to predict rice crop

78    yield since rice is the most cultivated crop in these regions. The ML approach was applied

79    separately to the three regions.

80    **Crop yield prediction in India**

81    In India, as in many other countries, the official crop yield prediction is mainly based on

82    conventional data collections techniques such as ground-field visits (Reynolds et al., 2000;

83    Sawasawa, 2003). The crop yield is measured through crop cutting experiments carried out

84    over sample crop areas. In this country, principal crops' calculations of area and yield are

85    released through the Directorate of Economics and Statistics, Ministry of Agriculture

86    (DESMOA). The production (in kg or ton) of a specific crop is calculated by multiplying the

87    whole field area by its crop yield. The crop production is needed for the decision-makers to

88    take various policy decisions relating to pricing, marketing, distribution, exportation and

89    importation.

90    The Kharif season, as it is locally known, represents about 80% of the annual rainfall (Naresh

91    Kumar et al., 2012). This monsoon season generally goes from June to October. In this season,

92    the highest agricultural production is obtained. Estimation of Kharif crop yield and production

93    is released four times during the year with different levels of sophistication and precision,

94    where the last one is considered the most accurate. The first calculation is presented in

95    September, the second one in January, the third one in March/April, and the fourth, and the last

96    one in June/July. It should be noted that the last two calculations of crop yield and production

97    become available much after the crops have already been harvested in December/January.

98    From the four calculations, the first two can be considered as predictions. These two first

99    predictions serve as primary estimations about how much the final yield and production will

100   be.

101 The existing ground-field visits-based agricultural forecasting system provides reliable

102 information; however, it lacks pre-harvesting forecasting. This limitation motivated the

103 creation of a satellite-based forecasting system to have information at the early stages of crop

104 growth. This system is called the National Crop Forecasting Centre (NCFC) (Sawasawa, 2003).

105 NCFC is continuously verified and continuously updated. Although NCFC advances the one

106 based on ground-field visits, data needed for its execution could be not always available.

107 Therefore, it is necessary to explore other solutions. In this study, it is not intended to replace

108 the previous and new forecasting systems, but to provide a complement to corroborate both

109 estimates, and in a broader sense, to provide the scientific community with an approach to crop

110 yield prediction with information on the spatial extent of drought.

## 2 ML modelling methodology

111 **2 ML modelling methodology**

112 The experiment was carried out with the following methodology that involves the ML

113 construction. The next paragraphs show each step in detail. These steps are (1) data preparation,

114 (2) input variable selection, (3) polynomial regression models calculation, (4) artificial neural

115 network models calculation, and (5) models application and combination.

### 2.1 Step 1. Data preparation

116 **2.1 Step 1. Data preparation**

117 Two types of data were prepared, the crop yield and the percentage of drought areas. For data

118 preparation, three tasks were carried out (1) data retrieving, (2) drought areas calculation, and

119 (3) data de-trending.

#### 2.1.1 Data retrieving

120 **2.1.1 Data retrieving**

121 Section 3 shows what corresponds to data retrieving for crop yield (CY) and the drought

122 indicator. CY data correspond to the largest growing season. CY time series has a value for

123 each year for the period 1966-2015 (49 years). CY was available for each region. On the other

124 hand, drought indicator data is on a monthly basis for the period 1901-2015. The spatial

125 resolution is half a degree.

#### 2.1.2 Drought areas calculation

126 **2.1.2 Drought areas calculation**

127 The drought areas were calculated following the methodology presented below. These areas

128 were calculated for the three regions. Drought areas were calculated from the drought indicator

129 data that is in a grid format, i.e., each cell has associated a geographic location and a time step.

130 The calculation of drought areas started with the reclassification of all the cells of the drought

131 indicator data by non-drought and drought cells. The drought indicator data was evaluated cell

132 by cell to determine those that are in drought, i.e. drought condition. To determine drought and

133 non-drought condition ($D_S$), the Eq. 1 was applied (Corzo Perez et al., 2011; Diaz et al., 2019,

134    2020; Herrera-Estrada et al., 2017). Eq. 1 represents the following. When the drought indicator

135    is below to the selected threshold $T$, the value of 1 is used to indicate drought in the cell and

136    non-drought is represented by the value of 0. This classification is performed for all the cells

137    of the grid data in each time step ($t$).

138    $$D_S(t) = \begin{cases} 1 \text{ if } \text{DI}(t) \leq T \\ 0 \text{ if } \text{DI}(t) > T \end{cases}$$    (Eq. 1)

139    Once the ones-and-zeros data was obtained, the drought areas (DAs) were calculated for each

140    region with Eq. 2. DA was computed as the ratio between the cells in drought and the total

141    number of cells of the region ($N$). In Eq. 2, the number of cell is denoted by $c$.

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

143    The number of cells ($N$) of the mask is 63, 31 and 54 for region 1, 2 and 3. The masks in raster

144    format were built for each region. The mask is an array of ones and zeros, where the value of

145    1 indicates the land. We used the threshold $T = -1$ to calculate cells in droughts. This threshold

146    is widely used to identify a cell in drought when working with standardised indices such as the

147    used in this research (Sect. 3.2). Usually, drought indicator data is calculated at different

148    aggregations periods. We retrieved this data for 1, 3, 6, 9, and 12 months of aggregation period

149    (Sect. 3.2). DAs' time series were calculated for each aggregation period and are indicated as

150    DA1, DA3, DA6, DA9, and DA12.

151    **2.1.3 Data de-trending**

152    Data stationarity is typically assumed when modelling. However, the present study uses crop

153    yield, which is non-stationary in nature. The crop yield depends on factors that affect its trend,

154    such as drought, flood, cultivars and its own management. Therefore, it is advisable to remove

155    short-term fluctuations in crop yield before constructing the model (Montesino Pouzols and

156    Lendasse, 2010).

157    Among the methods available to de-trend data, the 'first difference' method is popular due to

158    its simplicity. In this method, the trend is removed from the time series by subtracting the

159    previous value $x^*(t-1)$ from the current one $x^*(t)$, as shown in Eq. 3. The de-trended value for

160    the first time step ($t = 1$) is not calculated. The length of the de-trended time series is $n = m - 1$,

161    where $m$ is the length of the original time series. The de-trended data $x(t)$ has the same units as

162    the original data $x^*(t)$.

163    $$x(t) = x^*(t) - x^*(t-1)$$    (Eq. 3)

164    The trended of CY and DA time series was removed with Eq. 3. For the case of CY, the de-

165    trended time series retained one value per season, i.e. one per year. As noted, the method for

166    removing the trend does not generate the value for the first time step; therefore, the de-trended

167    CY data corresponds to the period 1967-2015 (49 years).

168    In the case of DA, Eq. 3 was applied as follows. Because the DA data is monthly, i.e. 12 values

169    per year, and CY data is seasonal, i.e. one value per year, first DA time series were extracted

170    for each month. The monthly values for January were extracted for each year and so on until

171    December. These twelve DA time series were compiled for each of the five DA1, 3, 6, 9 and

172    12 time series. A total of 60 DA time series ($12 \times 5$) were obtained. To refer to these time

173    series, a number (suffix) was added to indicate the month. In this way, for example, the time

174    series DA3_7 indicates the drought areas for July calculated from the drought indicator with

175    3-month aggregation period. Eq. 3 for the removal of the trend was applied to each of the 60

176    DA time series. The DA time series run from 1901-2015. For the construction of the ML

177    models, the common period 1967-2015 (49 years) was chosen.

178    **2.2 Step 2. Input variable selection**

179    In an ML model, the input, known as the predictor, is generally made up of independent

180    variables. Often these variables are arranged in different ways to determine the best model

181    input representation. An example arrangement is the selection of the independent variable using

182    different previous time steps, such as $t−1$ (the previous time), $t−2$ and so on. When using

183    drought indicators as the predictors, these arrangements include the different aggregation

184    periods (i.e. different aggregation periods are tested). The idea is not to include all the variables

185    and all their different possible arrangements but rather to find the best ones and discard those

186    that do not contribute significantly to the model's results. Other arrangements of the input

187    variable include the average, or other statistics, over a period.

188    There are different methods for selecting input variables. Based on the procedure, these

189    methods are classified into model-based and filter types (May et al., 2011). The first includes

190    all those where the model runs, and based on its performance, a specific variable is chosen or

191    discarded. The latter include methods where the variable is chosen *a priori* through a generally

192    statistical process and does not require the model to be run. Correlation analysis, which falls

193    under the second category, is often chosen for its simplicity and wide application. Correlation

194    is calculated between the time series of the output variable (CY in this case) and the different

195    input variables, including their various arrangements.

196    In this study, for the selection of the relevant input variables, the correlation analysis was done.

197    The correlation was calculated between the de-trended time series of the seasonal CY and the

198    60 DAs. As mentioned before, due to DAs are monthly and CY is seasonal, 12 time series of

199    DAs were prepared, one per month, for each aggregation period. The DAs were then correlated

200    with the CY. Another option could be to use the yearly average value of the DAs, such as the

201    average of the DAs of the months of the cultivation period, or something similar. However, we

202    opted to identify the DAs of the months that have the highest correlation with the seasonal CY

203    and use them as inputs.

204    The approach of the selection of the most correlated DAs was chosen for two main reasons.

205    On the one hand, rice responds to the climate variations differently from one growth stage to

206    another over the year, which could be better captured with the information of some months

207    than others. On the other hand, different types of drought (i.e. meteorological, agricultural, and

208    hydrological) are expected to affect (impact) the crop yield to different degrees. This level of

209    affectation could be taken into account either by using different hydro-meteorological variables

210    or selecting different aggregation periods of the meteorological variables, as in this case. An

211    average of DAs could "hide" a significant drought area that could contribute more (or less) to

212    the final crop yield. In addition, in this research, ML models were built to be used at different

213    stages of crop cultivation, i.e. models to be applied in June, July, and so on, each of them with

214    a different expected degree of accuracy. Therefore, the use of time series for each month

215    extracted from the DAs for all the different aggregation periods is more appropriate.

216    Based on the correlation coefficient, the input variables were selected. In total, 15 sets of input

217    variables (Table 2) were selected for each month from January to December. Each set is made

218    up of different DA time series, out of the 60 de-trended DAs. The number of variables is

219    different in each set. These sets of input variables are presented in the results section. All sets

220    include the de-trended CY from the previous year ($CY_{t-1}$). $CY_{t-1}$ was used because, in the

221    particular case of the study area, CY of the current year is planned to be reached based on data

222    of the previous year.

223    **2.3 Step 3. Polynomial regression models calculation**

224    For the case of PR, four types of models were tested (Table 1). All the PR models were built

225    for each month from January to December following Eq. 5 to 8. A total of 15 sets of

226    combinations of input variables were tested in each PR model. The best PR model was

227    identified for each month following the RMSE criterion (Eq. 9). MATLAB software was used

228    for implementation.

229    PR is an extension of linear regression that allows the use of more than one input variable to

230    calculate the output variable (Eq. 4).

231 $$y = b_0 + \sum_{i=1}^{n} b_i x_i + e \qquad \text{(Eq. 4)}$$

232 In Eq. 4, $y$ is the output variable, also known as the response, which in this case is the crop

233 yield. The term $x_i$ is the $i$-th input variable (predictor) from a total of $n$ variables. The regression

234 coefficients vector is represented by $b$. From the coefficients vector, $b_0$ is known as the

235 intercept. The vector of errors is indicated by $e$.

236 Table 1 shows four formulations of PR. The PR models are indicated as linear, pure-quadratic,

237 quadratic and interactions. Descriptions of each and their equations are presented in Table 1

238 (Eq. 5 to 8). The input variable ($x_i$) was selected based on the correlation analysis (Sect. 2.2).

239 **Table 1** Polynomial regression (PR) types followed in this study.

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

240

241 The best PR model was identified from four types using the root mean square error (RMSE)

242 criterion. The RMSE is calculated between the observations ($o$) and the predictions ($p$), as

243 shown in Eq. 9. RMSE is one of the most widely used criteria in the comparison of observations

244 and model calculations.

245 $$\text{RMSE} = \sqrt{\frac{\sum_{i=1}^{n} (o_i - p_i)^2}{n}} \qquad \text{(Eq. 9)}$$

246 **2.4 Step 4. Artificial neural network models calculation**

247 ANN is a method loosely based on imitating the basic functionality of neurons (i.e. the working

248 units of the human brain) (Govindaraju, 2000; Maier and Dandy, 2000). The input variables

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

249  (predictors) are connected to each other through mathematical formulations that allow complex

250  non-linear relationships to be represented. These connexions are symbolised as nodes

251  interconnected within a network aimed at calculating the output variable (response).

252  Of the different proposed ANN architectures (network designs), one of the most widely used

253  is the feedforward neural network (FFNN). The FFNN is schematised by a series of nodes

254  located in one of three layers: input, hidden or output. The number of input nodes is equal to

255  the number of input variables in the input layer (Elshorbagy et al., 2010). This first layer is in

256  turn connected to the hidden layer, which receives this name because the connections made

257  there may not be immediately evident to the model performer. In this hidden layer, the number

258  of nodes is not defined by default; rather, the greater the number of nodes, the more complex

259  the model. Finally, the nodes of the hidden layer are connected to those of the output layer. In

260  a single-output variable problem, there is only one node. ANNs are typically trained by non-

261  linear optimisation gradient-based algorithms, e.g. the Levenberg-Marquardt algorithm.

262  In the ANN setup, the number of nodes of the input layer was equal to the number of variables

263  of the respective combination. The number of nodes in the output layer was one and

264  corresponded to the seasonal crop production (CY). An iteration optimisation procedure was

265  carried out regarding the hidden layer, varying the number of nodes from 1 to 10. For each

266  number of nodes, 100 iterations were done, being 1,000 in total. For reproducibility of the

267  results, the random values were set to default at the beginning of the number of nodes change.

268  For each month, from January to December, the ANNs were built. MATLAB software was

269  used to implement the ANNs with the Levenberg-Marquardt algorithm for training. In each of

270  the ANNs, 85 % of the data was used for training-validation, and the rest for testing

271  (verification). The best model corresponding to each number of hidden nodes was identified,

272  i.e. ten models per month and the best model for each month. RMSE was used to identify the

273  best models. RMSE was calculated for (1) the training-validation dataset (RMSE_cal), (2) the

274  testing dataset (RMSE_test), and (3) the entire period (RMSE). In all the cases, the final (best)

275  model was chosen based on RMSE for the entire period.

276  **2.5 Step 5. Models application and combination**

277  Once the best ML models, PR and ANN, were known, the pair of models were selected for

278  each month. Depending on the performance of these models (and experience of their use), they

279  can be used either separately or combined, e.g. being run in parallel so that a modeller could

280  see the cases when models produce different results. An alternative is to use a dynamic

281    weighting of the models' outputs (e.g. with the weights being proportional to the historical

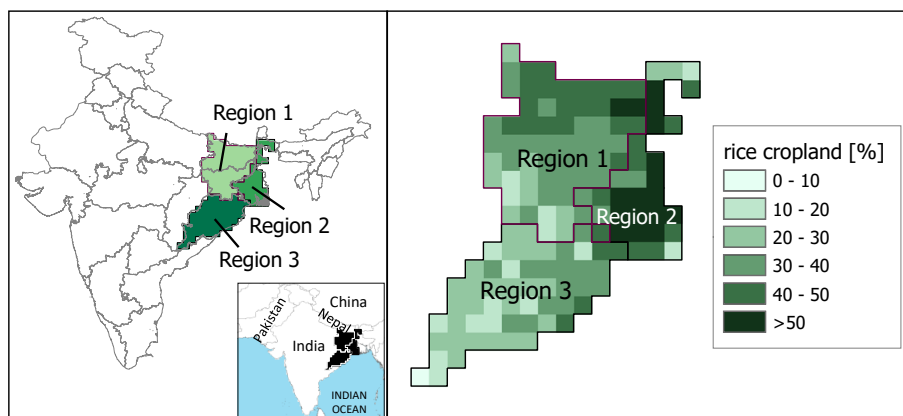282    performance) to form a "model committee".

283    **3 Data**

284    **3.1 Crop yield**

285    Rice is the most important food grain in East India, so it was selected to assess our ML-oriented

286    crop-yield predictions. Rice from this region accounts for roughly 85 percent of the total rice

287    production in India (Ghosh et al., 2014). As mentioned, ML models were constructed for three

288    regions of the eastern Indian (Figure 1). State-wise crop-yield data was retrieved from 1966 to

289    2015 (49 years) through the Indian Directorate of Economic and Statistics from the Department

290    of Agriculture (DAC) (http://eands.dacnet.nic.in/).

291    There are three crop seasons in India: Rabi, Kharif and Zaid. Of these, the Kharif season was

292    chosen for study because it is the largest in terms of crop production. Kharif crops are sown in

293    June and harvested in November/December. Seasonal crop-yield data was obtained from the

294    DAC website and arranged into time series per region. One value was assigned to each year of

295    crops harvested in the Kharif season.

296    Figure 1 shows the location of the three regions. These are made up as follows. Region 1

297    includes the current states of Bihar and Jharkhand; region 2 corresponds to the state of West

298    Bengal; and region 3 makes up the state of Odisha. Two important clarifications have to be

299    made regarding crop yield data retrieving for these regions. First, in late 2000, Bihar was

300    divided into two states: Bihar and Jharkhand. Thereafter, rice data was reported separately. In

301    this study, both states are marked as region 1; the crop-yield data from 2000 to 2015 is the

302    reported sum of current Bihar and Jharkhand. Second, in 2011, Orissa was renamed Odisha

303    (region 3), but the territory remains the same. In this case, crop yield data for Odisha is that

304    reported for the former Orissa and the current Odisha.

**Figure 1** Case study location and rice cropland (in percentage). Case study comprises region 1 (Bihar and Jharkhand), region 2 (West Bengal) and region 3 (Odisha). Source of rice cropland: Monfreda et al. (2008).

**3.2 Drought indicator**

Soil moisture is the preferred variable for calculating agricultural drought indicators. However, another widely disseminated way to indirectly infer this type of drought indicator is to use meteorological drought indicators as proxies. Among these, the Standardised Precipitation Evaporation Index (SPEI) proposed by Vicente-Serrano et al. (2010) has shown to be useful in assessing agricultural drought. The SPEI follows a similar methodology as that of the widely used Standardized 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) between 1901 and 2015. The spatial resolution of the drought indicator data is 0.5 degrees. The SPEI data was available at different aggregation periods; for this study, it was retrieved for the aggregation periods of 1, 3, 6, 9 and 12 months, indicated as DI1, DI3, DI6, DI9 and DI12, respectively.
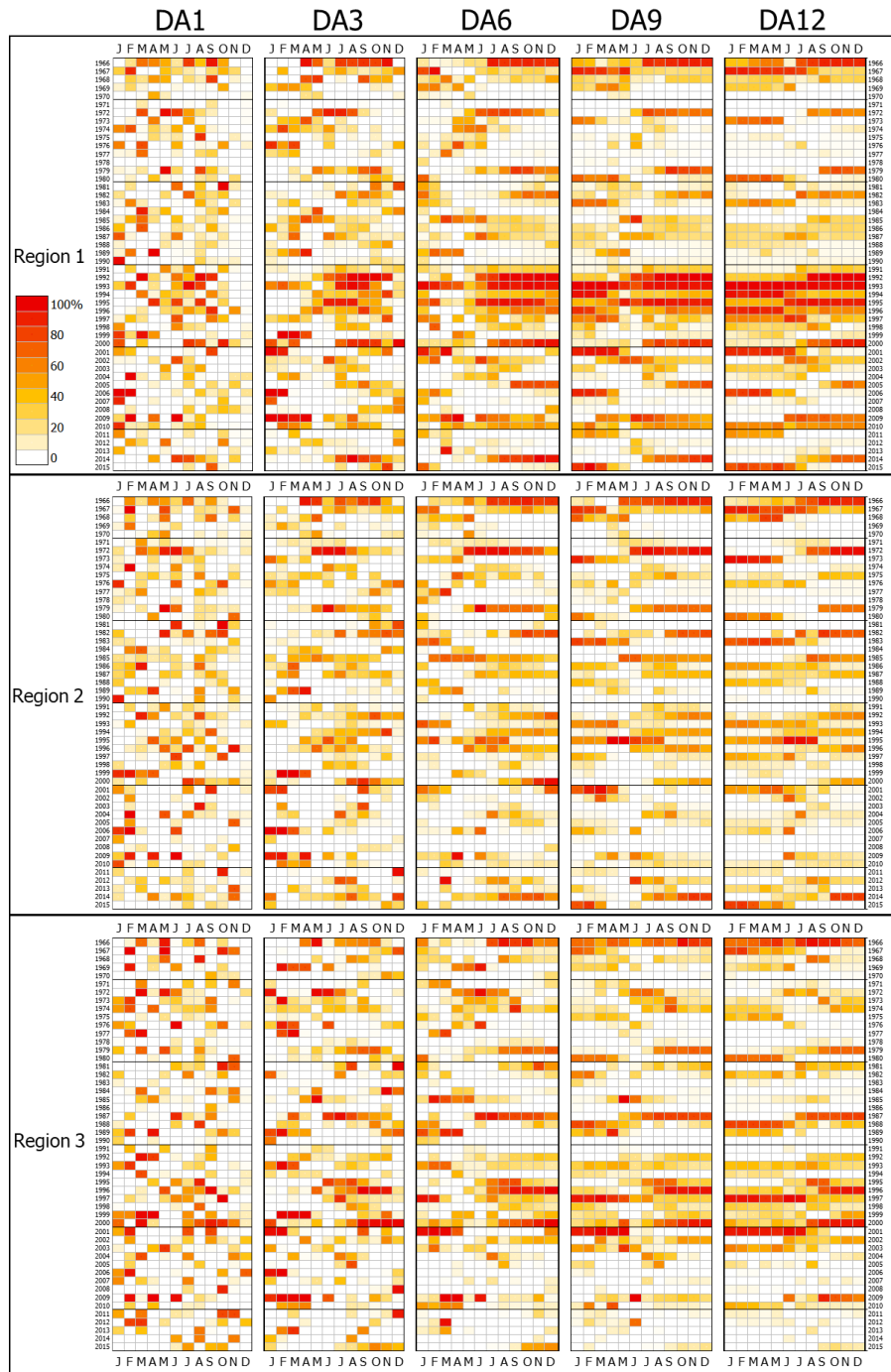
**4 Results and discussion**

**4.1 Data preparation: drought areas and crop yield**

Figure 2 show the drought areas calculated for the three regions. In this heat map, columns indicate the months and rows point out the years. The redder the colour, the larger the drought area. In general, region 1 (Figure 2, the upper panel) presents the highest values concerning the other two regions. In general, the 1990s show higher values of areas with respect to the rest of the period, which agrees with Guha-Sapir (2019); in this decade, there were three droughts, 1993, 1996 and 2000. At the beginning of the period, large areas are also observed in the three regions; these results align with Bhalme and Mooley (1980).

11

329    In Figure 2, a pattern is observed in the drought areas distribution for all the aggregation
330    periods, i.e. from DA1 to DA12. In DA1, the areas mainly concentrate in the first months; even
331    the December column is almost white (without drought). Later, for DA3, the large areas are
332    located from April to November. Successively, for DA6 and DA9, the largest areas are
333    concentrated in the second half of the year. There are even droughts that end in the following
334    year; they are the reddish lines that are observed in the first semester (first columns). Finally,
335    in DA12, there are consecutive large areas indicated by the reddish lines; droughts usually
336    begin in the second semester and extend until the following year. These results show the
337    importance of considering more than one period of aggregation when using indicators based
338    on meteorological variables; each aggregation period can be a proxy for analysing different
339    types of drought and its effects.

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU



**Figure 2** Drought areas (DAs) for each aggregation period (1, 3, 6, 9 and 12 months) and region. Top, middle, and bottom panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal) and region 3 (Odisha).

343 Figure 3 shows the time series of de-trended CY and DA for the three regions. In the case of

344 DA (indicated in red), the values are displayed in inverse order to facilitate interpretation. In

345 general, when drought areas increase, this is expected to affect crop yield (decreasing).

346 Otherwise, when the drought area decreases, this effect favours an increase in crop yield. In

347 general, for the three regions, the decreases in CY coincide with the increases in DA. The

348 general pattern regarding DA variations is as follows. The values fluctuate throughout the year

349 for the aggregation periods of one and three months (DA1 and DA3). Subsequently, for DA6

350 to DA12, the values are concentrated in the second half of the year. These results also show

351 the usefulness of the different aggregation periods to capture different types of drought. The

352 effect of increasing DA seems not to be observed in decreasing CY for all cases of DAs. For

353 example, in region 1 (Figure 3, the upper panel), the decrease in 2004, one of the maximums,

354 does not coincide with increases in DA9 and DA12, but it does for DA1, DA3 and DA6. These

355 results also support the use of the different aggregation periods on drought assessments.

356 **4.2 Input variable selection (correlation analysis)**

357 Figure 4 summarises the correlation between the de-trended CY and the DAs, and Figure 5

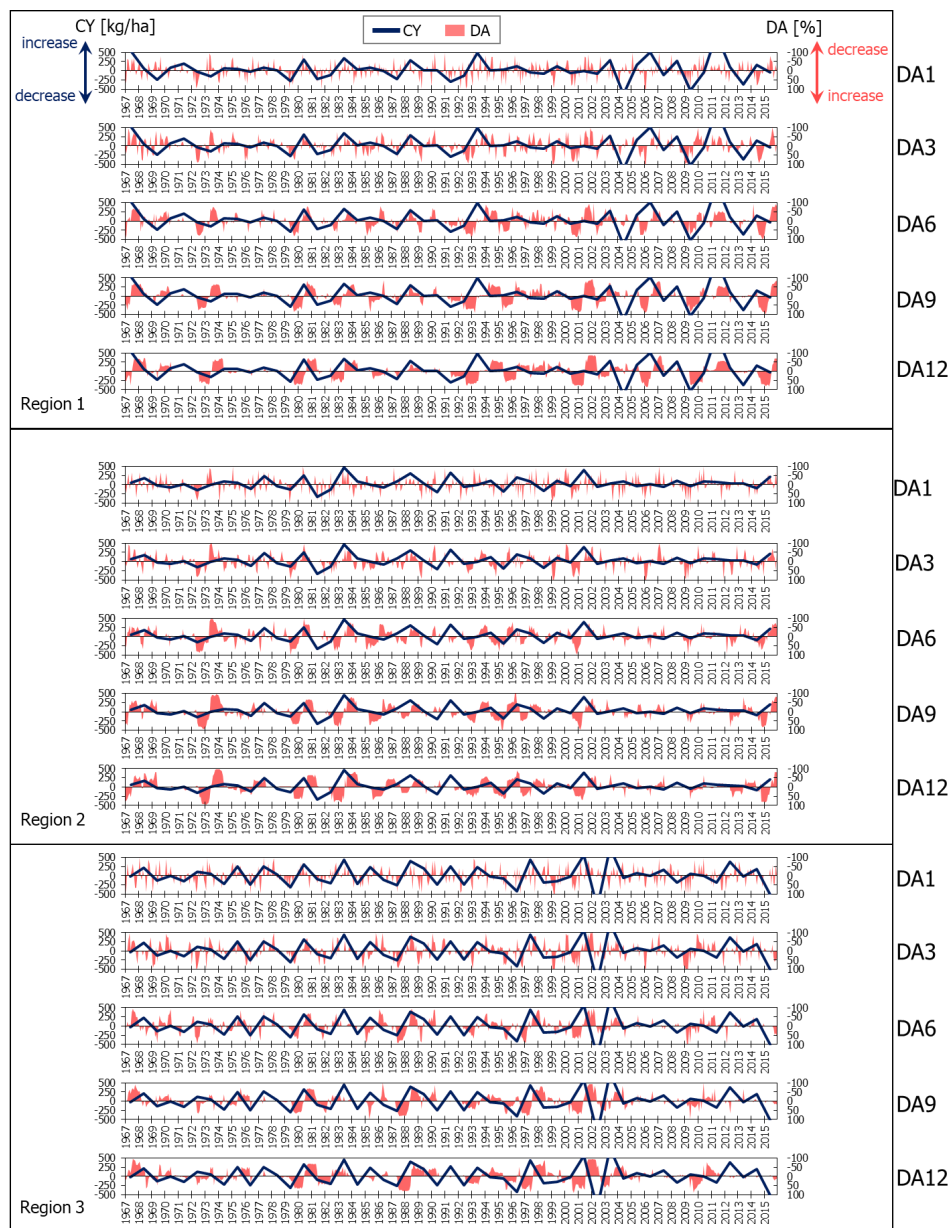358 presents the correlation for each monthly DA time series.

359 Figure 4 and 5 shows that the correlation is different over the year in the three regions. In all

360 cases, the correlation coefficient increases until a maximum and then decreases. The month in

361 which the maximum value is reached is different for each region but falls within the crop season

362 (i.e. June to November/December). For region 1, it is in July. For region 2, there are four

363 months with this pattern, June, July, October, and November. Finally, for region 3, it is

364 October, November and December.

365 These results of correlation can be useful for monitoring agricultural drought. For example, in

366 region 1, the drought areas calculated from SPEI6 (i.e. DA6) show a maximum correlation in

367 July. This correlation value means that the previous six months' accumulated effect is crucial

368 for the crop yield of the Kharif season, which covers more or less from June to

369 November/December.

370 Figure 4 shows the following pattern. In general, for region 1, results similar to DA6 are

371 observed for DA3, 9 and 12. For region 2, a similar pattern happens in the peaks, but in this

372 case two, one corresponding to DA1 and 3, and the other to DA6, 9, and 12. The first peak of

373 DA1 and DA3 may indicate that it is crucial to pay attention to the immediate period conditions

374 of one to three months. In the case of the second peak, the medium and long-term conditions,

375 6 to 12 months, are more important to monitor for the harvest month. For region 3, the peak

376    occurs at the end of the growing season, in almost all cases. Hence, the condition before the

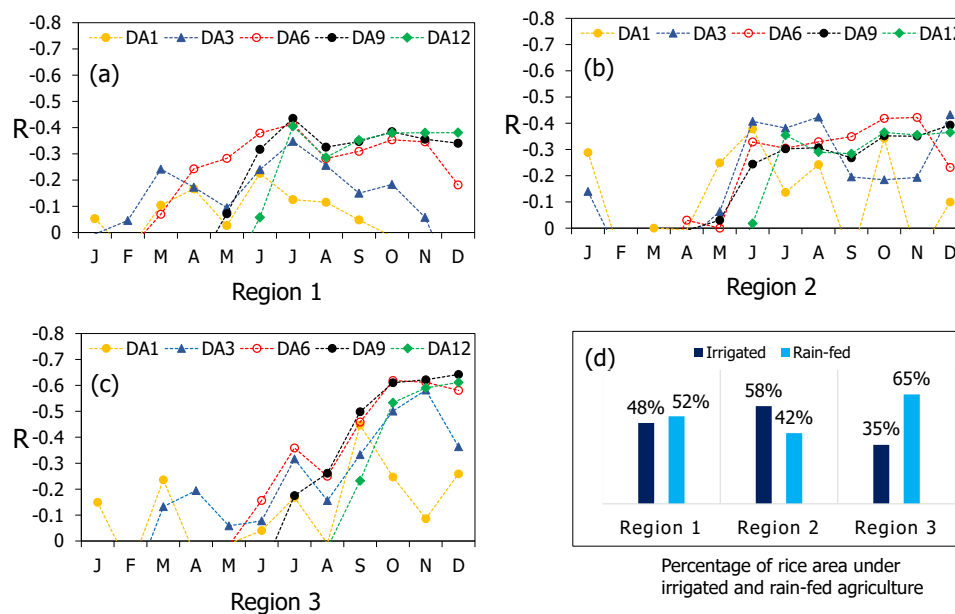377    growing season is decisive for the crop yield.

378



379

**Figure 3** Time series of the de-trended crop yield (CY) and drought areas (DAs) for each aggregation period (1,

3, 6, 9 and 12 months) and region. Top, middle, and bottom panels indicate region 1 (Bihar and Jharkhand), region

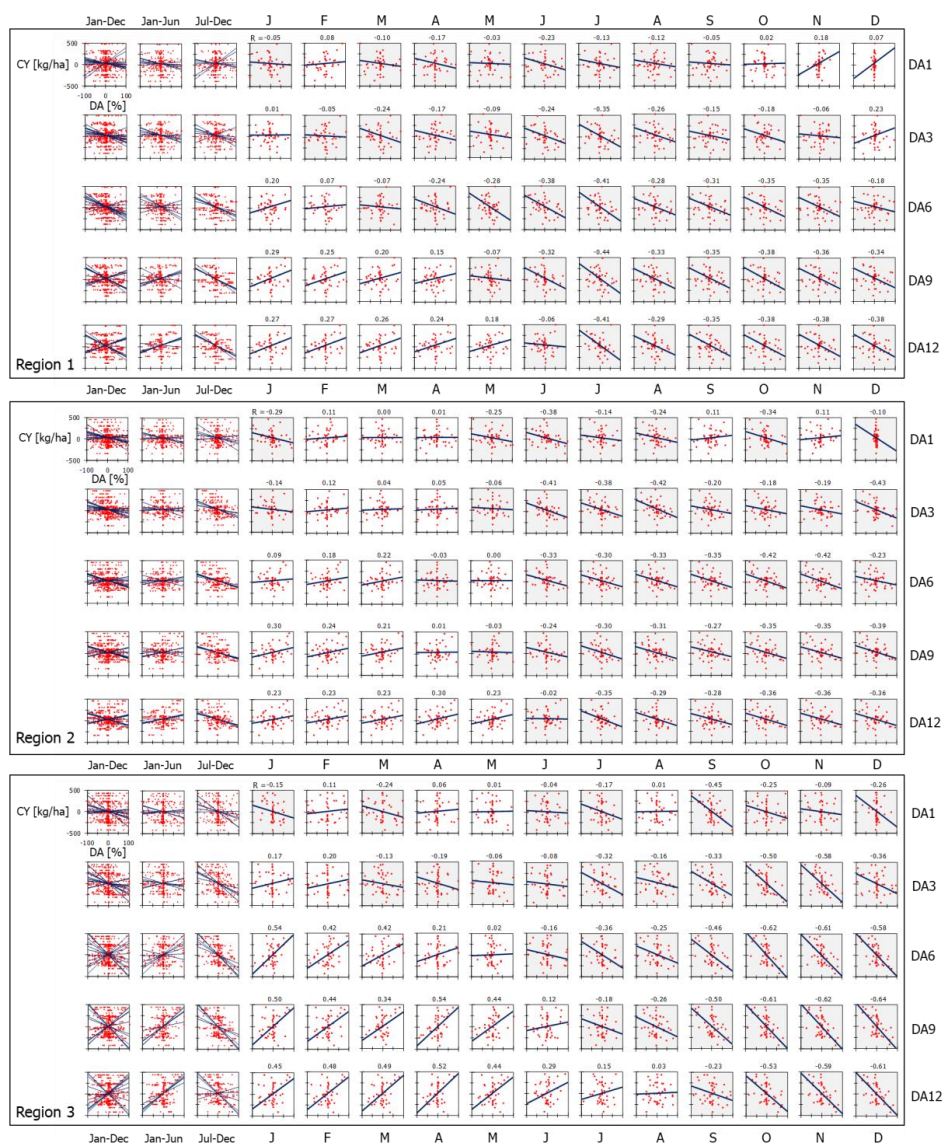2 (West Bengal) and region 3 (Odisha).

15

383    Figure 5 shows how the correlation coefficients between CY and DA are positive outside the

384    growing season and negative within that season. However, this pattern is less evident for DA1

385    and DA3. The pattern shown by the correlation coefficients in Figure 5 supports the idea that

386    drought is an important factor in crop yield since the months with less drought are more

387    correlated with the increase in CY, and the months with more drought do so with decrease in

388    CY.

389    Figure 4 (d) shows the percentage of irrigated and rain-fed agriculture. For regions 1 and 2,

390    about half is by irrigation, while in region 3, only 35%. Perhaps this percentage of irrigation

391    for region 3 explains why the correlation coefficients for this region are higher than for the

392    other two (Figure 4, and 5 (c)). Region 3 is more dependent on rain for agriculture; therefore,

393    this condition is best captured when calculating drought with the precipitation, as in this case

394    (Sect. 3.2).



395

**Figure 4** Summary of correlation between de-trended crop yield (CY) and drought areas (DAs) for each

aggregation period (1, 3, 6, 9 and 12 months) and region: (a) region 1 (Bihar and Jharkhand), (b) region 2 (West

Bengal) and (c) region 3 (Odisha). Percentage of rice area under irrigated and rein-fed agriculture (d). Source of

irrigated and rein-fed agriculture data: Directorate of Rice Development (DRD), (2014).

Hydrology and
Earth System
Sciences
Discussions
Open Access
EGU

**Figure 5** Correlation between de-trended crop yield (CY) and drought areas (DAs) for each aggregation period (1, 3, 6, 9 and 12 months) and region. Results are shown for each monthly DA time series from June to December (J to D). Top, middle, and bottom panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal) and region 3 (Odisha).

17

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

408  Figure 4 (a, b, and c) shows the following pattern in the three regions. The correlation

409  coefficients between CY and DAs increase according to the aggregation periods and the month

410  of analysis. DA1 and DA3 have a better correlation in the first months of the year. DA6 has a

411  better correlation in the subsequent months, between May and June. Finally, DA9 and 12 do

412  so within the second half of the year.

413  Each respective DA time series reaches a maximum (or maximums) of correlation, and then

414  correlation decreases. According to this pattern, the 15 combinations of input variables shown

415  in Table 2 were selected. As earlier mentioned, the CY of the previous year was included in all

416  combinations and is indicated as $CY_{t-1}$. Combinations 1 to 5 only include a DA time series.

417  Combinations 6 to 9 are DA pairs that were calculated with the drought indicator of successive

418  aggregation times. For example, combination 6 forms DA1 and 3, combination 7 includes DA3

419  and 6, and so on. Similarly, combinations 10 to 13 are proposed, but for triples. Combinations

420  13 and 14 are fourfold. Finally, the last combination (15th) is made up of all the DA series.

421  As mentioned, the models were built for each DA time series using the 15 combinations

422  corresponding to each month. For example, the monthly series of DAs extracted for January

423  were used for the case of January. These DAs are DA1_1, DA3_1, DA6_1, DA9_1 and

424  DA12_1. The suffix indicates the month. Then, the different DA1_1 to DA12_1 were used

425  following the 15 combinations shown in Table 2 to build the ML models (ANN and PR) for

426  January. Similarly, it was carried out from February to December.

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

428  DAs are calculated with the drought indicator for the aggregate period of 1, 3, 6, 9 and 12 months.

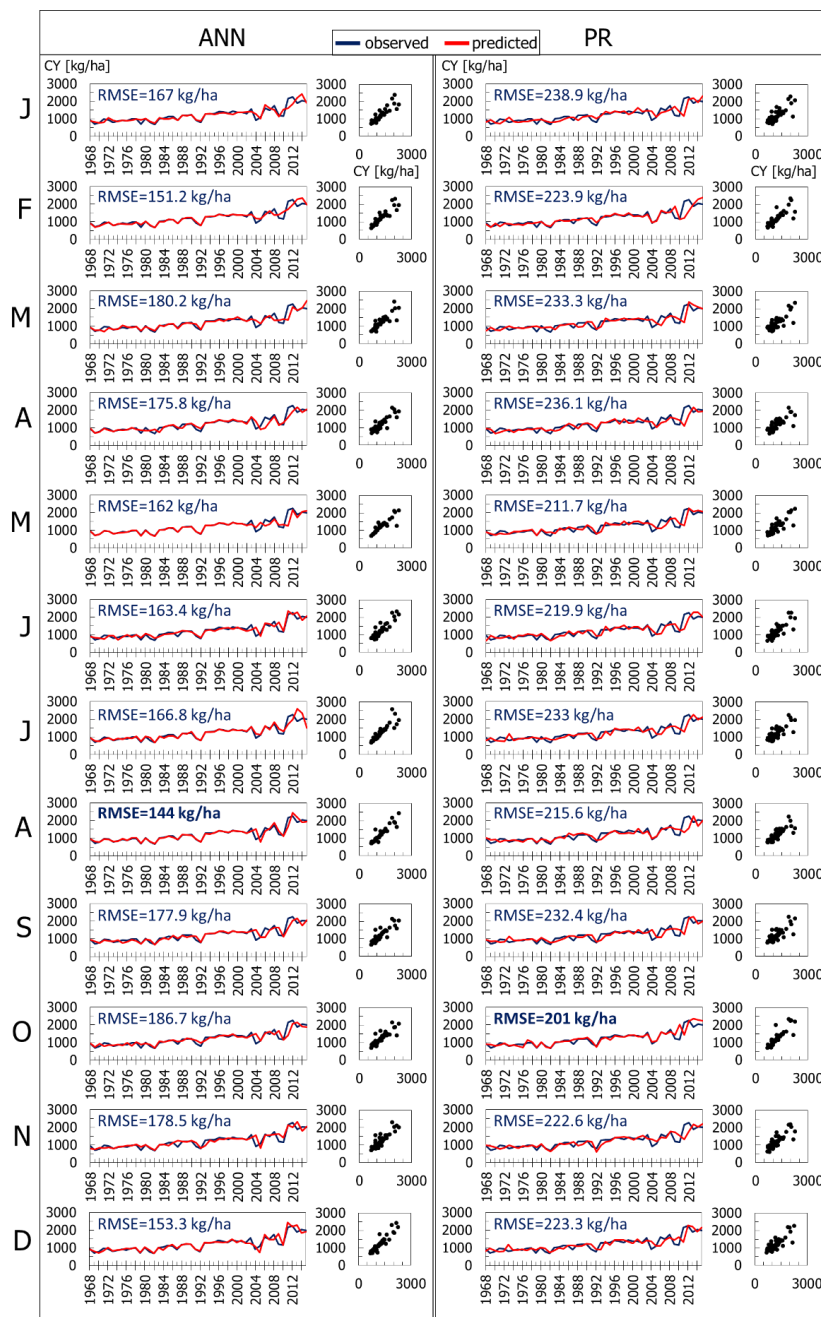| Input set (combination) | Input variables |
| --- | --- |
| 1 | $CY_{t-1}$, DA1 |
| 2 | $CY_{t-1}$, DA3 |
| 3 | $CY_{t-1}$, DA6 |
| 4 | $CY_{t-1}$, DA9 |
| 5 | $CY_{t-1}$, DA12 |
| 6 | $CY_{t-1}$, DA1,3 |
| 7 | $CY_{t-1}$, DA3,6 |
| 8 | $CY_{t-1}$, DA6,9 |
| 9 | $CY_{t-1}$, DA9,12 |
| 10 | $CY_{t-1}$,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

### 4.3 ANN and PR models

430    The results show different magnitudes of error between the observed and predicted CY. The

431    models with the lowest error are presented in Figures 6, 7 and 8, for each of the three regions.

432    The pair of ANN and PR that best predicts CY is shown for each month. The RMSE is also

433    indicated in each case. On the other hand, Figure 9 shows the error for each input set

434    (combination); the lowest error achieved in each month is presented in each case both for each
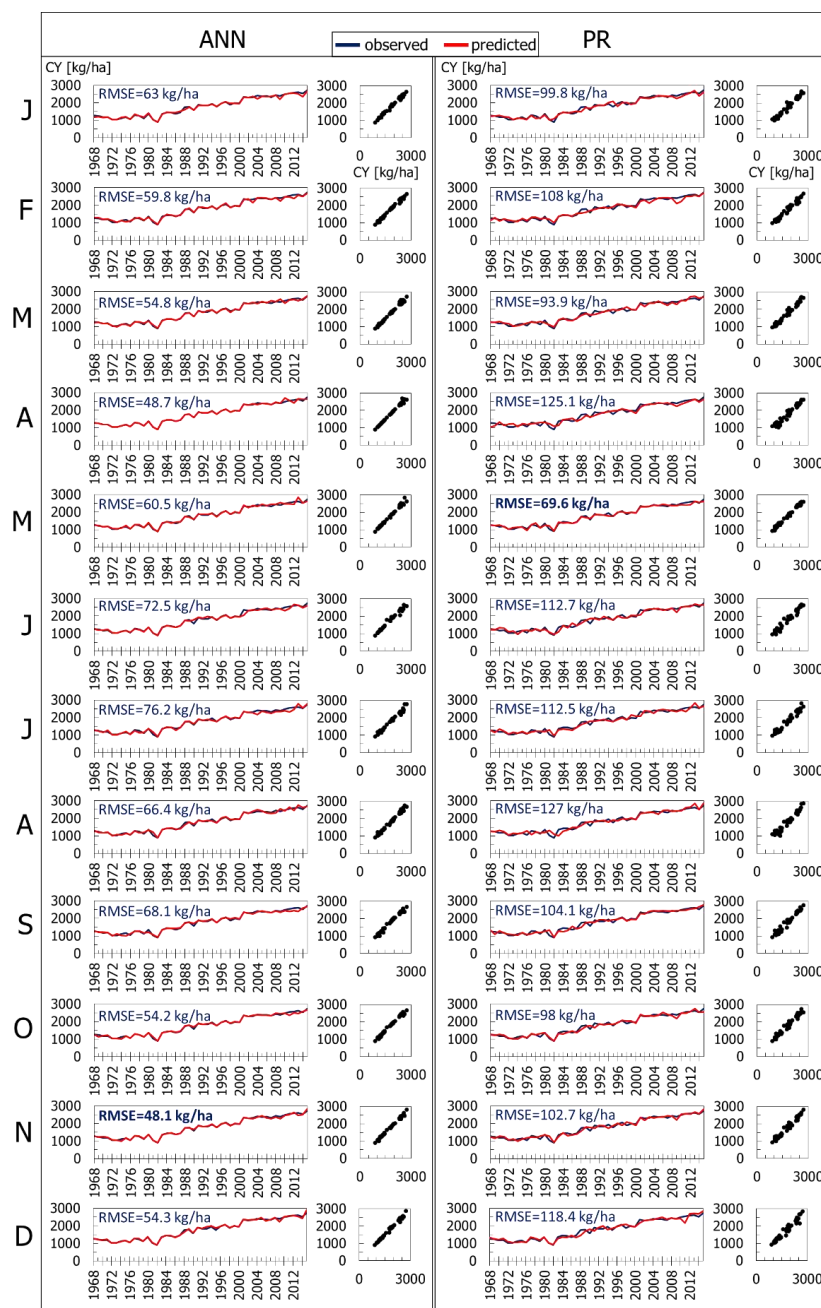
435    ANN and PR.

436    In general, ANN shows the least errors, as expected (Figure 9). However, the results of PR are

437    not much worse compared to those of ANN; for example, in some cases, the errors shown by

438    linear PR are very close to those of ANN (e.g. Figure 9, region 2). In general, it is observed

439    that the models with the lowest errors correspond to region 2, followed by region 3 and region

440    1 (Figure 9). It is attributed to the different degrees of crop irrigation with surface and mostly

441    groundwater, which determines the accuracy of the modelling in the different regions. Another

442    factor contributing to the models' performance is the drastic changes in the CY data, where

443    regions 1 and 3 are the ones that presented the most, and to a much lesser extent, region 2.

444    Figure 9 shows that in the three regions, different types of PR showed better results. In general,

445    linear and pure-quadratic indicate more stable results (no sudden changes among the different

446    realisations) but not better than quadratic and interactions. In general quadratic and interactions

447    present better results, being in some cases very close to those shown by ANN, e.g. PR

448    interactions (Figure 9, region 1).

**Figure 6** ANN and PR models for predicting seasonal crop yield (CY) built for each time series of monthly drought areas (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).
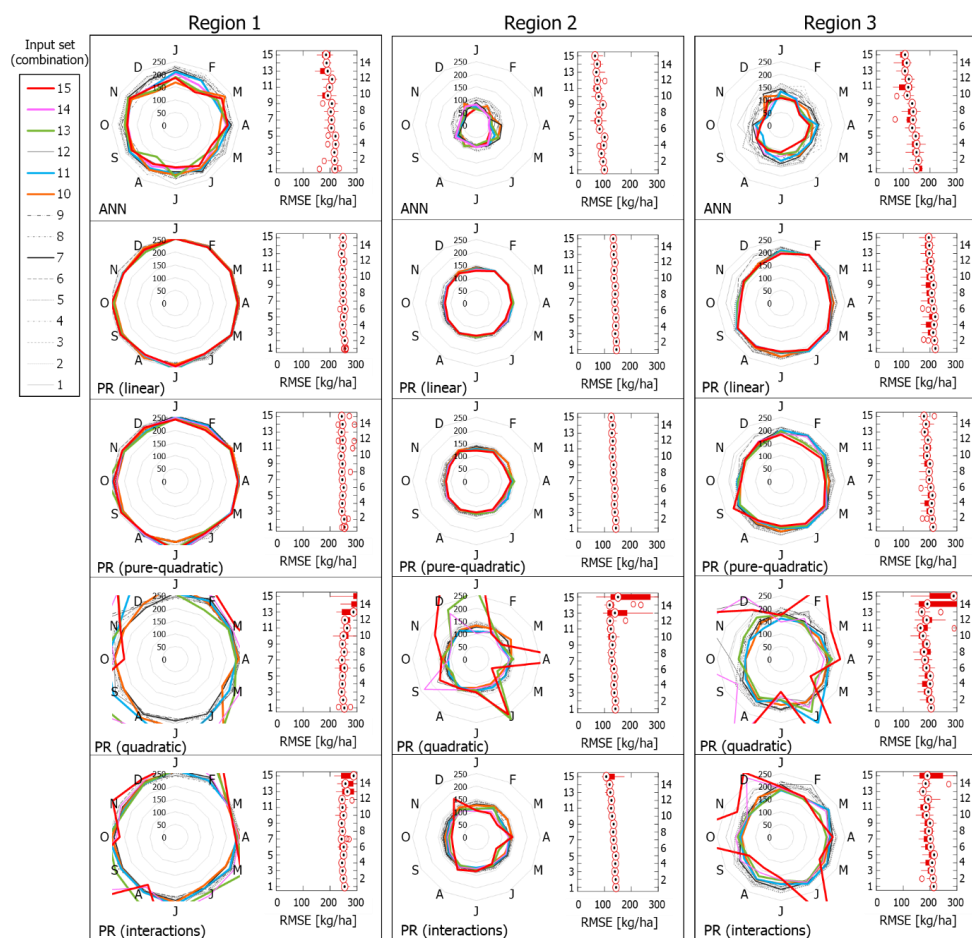
**Figure 7** ANN and PR models for predicting seasonal crop yield (CY) built for each time series of monthly drought areas (DAs): region 2 (West Bengal). The model with the lowest error (RMSE) is presented for each month, from January to December (J to D).

457

**Figure 8** ANN and PR models for predicting seasonal crop yield (CY) built for each time series of monthly
drought areas (DAs): region 3 (Odisha). The model with the lowest error (RMSE) is presented for each month,
from January to December (J to D).

**Figure 9** Root mean square error (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 input (from one to 15), the lowest errors are presented for each month (January to December). Results of each input set are shown with lines to facilitate the analysis. Left, middle, and right panels indicate region 1 (Bihar and Jharkhand), region 2 (West Bengal) and region 3 (Odisha).

## 4.4 Models application and combination

The best performing models were selected for each month. Table 3 shows the summary of these models, which includes the input set (combination), number of nodes, and errors for ANN, and input set, type and errors for PR. The number of nodes indicates the degree of non-linearity presented in each model. In this way, the more nodes, the more complex the model is in the case of ANN. On the other hand, quadratic and interactions are the types that showed the best performance in PR models. In all cases, within the combinations of input variables, a single DA time series corresponding to one of the various aggregation periods (D1, D3, D6, D9 or

23

475   D12) that by itself produced good results was not found. The input sets are made up of two and

476   up to six different DAs corresponding to the various aggregation periods. Thus, using more

477   than one aggregation period of drought indicator results in better model performance.

478   Tables 4, 5 and 6 are derived from Table 3. These three tables show the PR formulas for region

479   1, 2 and 3, respectively. In each table, the PR formula and the inputs are indicated. These

480   formulas are also intended to be a stand-alone tool in the CY prediction for each region.

481   The proposed procedure for applying the ML models is as follows.

482   The calculation begins by selecting the formula of the PR model for each month. Then the CY

483   is calculated with the chosen formula and the corresponding input variables. At the same time,

484   or when it can be computed, the ANN model of the month under analysis is applied. An

485   alternative is to use an approach based on the dynamic weighting of the models' outputs to form

486   a model committee.

487   **Table 3** Summary of the ANN and PR models for predicting crop yield (CY) built for each month and region: (1)

488   Bihar and Jharkhand, (2) West Bengal and (3) Odisha. The table shows the models built with the lowest error

489   (RMSE). DA stands for drought area.

| | | ANN | | | | | PR | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Region | 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 |

490

491    **Table 4** PR models for predicting crop yield (CY) built for each month: region 1 (Bihar and Jharkhand). For each

492    moth, it is indicated the input (x1 to x6) and the PR formula. DA stands for drought area.

| 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 | | $29.2005 -0.3816x_1 -0.6953x_2 +0.8469x_3 +1.2024x_4 -3.2563x_5 +0.0005x_1^2 -0.5339x_2^2 -0.0047x_3^2 -0.0119x_4^2 +0.0083x_5^2$ |

493

494

495

496

497

498

499

500 **Table 5** PR models for predicting crop yield (CY) built for each month: region 2 (West Bengal). For each moth,

501 it is indicated the input (x1 to x6) and the PR formula. DA stands for drought area.

| 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 | | $24.769 -0.1091x_1 -2.9747x_2 +2.9990x_3 -5.4144x_4 +3.3374x_5 +0.0083x_1x_2 -0.0069x_1x_3$ $+0.0596x_1x_4 -0.0630x_1x_5 +0.0755x_2x_3 +0.0127x_2x_4 +0.0094x_2x_5 -0.0052x_3x_4$ $-0.0884x_3x_5 +0.0361x_4x_5$ |

502

503

504

505

506

507

508

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

509   **Table 6** PR models for predicting crop yield (CY) built for each month: region 3 (Odisha). For each moth, it is

510   indicated the input (x1 to x6) and the PR formula. DA stands for drought area.

| 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 | | | | $13.0378 - 0.5111x_1 + 0.5765x_2 - 3.4820x_3 + 0.0177x_1x_2 - 0.0158x_1x_3 + 0.0155x_2x_3 + 0.0004x_1^2 - 0.0691x_2^2 + 0.0343x_3^2$ |

511

## 4.5 ML modelling limitations

513   The limitations of the presented approach are the following.

514   (1) To determine drought areas, a threshold value of the Standardised Precipitation

515   Evapotranspiration Index (SPEI) drought index (SPEI ≤ -1) was used. Using just one threshold

516   might lead to over or underestimation of the actual drought impacts over crop yield.

517 (2) Gridded data of SPEI at spatial resolution (0.5°x0.5°) was used in this study over each

518 region individually. Using such a coarse spatial resolution on different region sizes might not

519 capture the drought area correctly, leading to over or underestimating its magnitude.

520 (3) The study area has a diverse ecosystem of irrigated and rain-fed land, which may influence

521 the correlation between DA and crop yield more or less.

522 (4) This study assumes that drought is the only causative factor; however, floods negatively

523 impact crop yield in the region, thus in the total production in the regions. Flood impacts are

524 not considered in the models.

525 (5) Many other factors might influence rice yield, such as market, technologies, management,

526 etc. In this study, it was assumed that drought plays the prominent role.

527 (6) Insufficient crop yield data for the ML model building was an issue because the CY time

528 series only had one value for each year.

### 5 Summary and conclusions

530 This research introduced a step-by-step ML approach for predicting crop yield (CY) with

531 drought areas (DAs) as input. The ML approach comprises two components. Each component

532 employs two types of ML models: polynomial regression (PR) and artificial neural network

533 (ANN). The goal was to build the ML models (ANN and PR) and use them as an integrated

534 tool to crop yield prediction. The formulas of the PR models were also provided. The ML

535 approach was applied in three East India regions.

536 The following conclusions are drawn from this research.

537 • Based on the performance of PR and ANN models, results show drought area to be a
538   suitable variable to predict crop yield.

539 • The correlation analysis between DA and CY showed high negative correlations in
540   Odisha (region 3). The correlation gradually decreases in Bihar and Jharkhand (region
541   1) and West Bengal (region 2). These correlation values can be because West Bengal
542   has better access to irrigation facilities than Odisha and Bihar & Jharkhand.

543 • On comparing ANN models and PR models, the ANN were more accurate than PR
544   models to predict crop yield for all regions. This could have been expected since the
545   drought–crop relationship is a highly non-linear problem.

546 • It can be concluded that ANN has a high capability to predict CY in the pre-harvesting
547   stage with good accuracy, considering the drought indicator used (SPEI), which uses
548   climate variables such as precipitation and temperature (for evapotranspiration
549   calculation).

550 From the analysis and findings of this research, the following recommendations can be
551 provided for further improvement.

552 • Sensitivity analysis should be performed to identify the parameters that can impact the
553 model results. For instance, different spatial resolutions of drought indicator and
554 different thresholds should be investigated.

555 • Wet extreme events should be considered, especially in the flood-prone regions such as
556 the coastal areas of West Bengal (region 2) and Odisha (region 3) and North Bihar
557 (region 1), where floods also influence crop yield.

558 • Non-climatic factors such as econometric, fertilisers, and management practices might
559 be considered because they influence crop yield.

560 • In order to improve the model accuracy, more input data should be used in further
561 studies. For CY, this can be estimated by remote sensing techniques on a monthly basis
562 so that the ML models can be built for this temporal resolution and the spatial coverage
563 can be better addressed.

564 • The performance of other ML models has to be investigated, especially committee
565 (ensemble) methods like random forests or boosting methods. In the case of data at
566 scales less than monthly, the use of deep learning algorithms (e.g. LSTM networks)
567 could be recommended to explore.

568 We envision that this research will improve drought monitoring systems for assessing drought
569 effects. Since it is currently possible to calculate drought areas within these systems, the direct
570 application of the prediction of drought effects is possible to integrate by following approaches
571 such as the one presented or similar.

## Coda and data availability

## Competing interests

576 
577 An author is member of the editorial board of journal HESS. The peer-review process was guided by an
578 independent editor, and the authors have also no other competing interests to declare.

## Acknowledgements

589 **References**

590 Below, R., Grover-Kopec, E., and Dilley, M. (2007). Documenting Drought-Related Disasters:

591 A Global Reassessment. *J. Environ. Dev.*, *16*(3), 328–344.

592 https://doi.org/10.1177/1070496507306222

593 Bhalme, H.N. and Mooley, D. a. (1980). Large-Scale Droughts/Floods and Monsoon

594 Circulation. *Monthly Weather Review*, *108*(8), 1197–1211. https://doi.org/10.1175/1520-

595 0493(1980)108<1197:LSDAMC>2.0.CO;2

596 Chlingaryan, A., Sukkarieh, S., and Whelan, B. (2018). Machine learning approaches for crop

597 yield prediction and nitrogen status estimation in precision agriculture: A review.

598 *Computers and Electronics in Agriculture*, *151*(May), 61–69.

599 https://doi.org/10.1016/j.compag.2018.05.012

600 Corzo Perez, G.A., van Huijgevoort, M.H.J., Voß, F., and van Lanen, H.A.J. (2011). On the

601 spatio-temporal analysis of hydrological droughts from global hydrological models.

602 *Hydrology and Earth System Sciences*, *15*(9), 2963–2978. https://doi.org/10.5194/hess-

603 15-2963-2011

604 Dai, A. (2011). Characteristics and trends in various forms of the Palmer Drought Severity

605 Index during 1900 – 2008. *Journal of Geophysical Research*, *116*(March), 1–26.

606 https://doi.org/10.1029/2010JD015541

607 Diaz, V., Corzo, G., Van Lanen, H.A.J., and Solomatine, D.P. (2019). Spatiotemporal Drought

608 Analysis at Country Scale Through the Application of the STAND Toolbox.

609 *Spatiotemporal Analysis of Extreme Hydrological Events*, 77–93.

610 https://doi.org/10.1016/B978-0-12-811689-0.00004-5

611 Diaz, V., Corzo Perez, G.A., Van Lanen, H.A.J., Solomatine, D., and Varouchakis, E.A.

612 (2020). An approach to characterise spatio-temporal drought dynamics. *Advances in*

613 *Water Resources*, *137*, 103512.

614 https://doi.org/https://doi.org/10.1016/j.advwatres.2020.103512

615 Elshorbagy, A., Corzo, G., Srinivasulu, S., and Solomatine, D.P. (2010). Experimental

616 investigation of the predictive capabilities of data driven modeling techniques in

617 hydrology - Part 2: Application. *Hydrology and Earth System Sciences*, *14*(10), 1943–

618 1961. https://doi.org/10.5194/hess-14-1943-2010

619 Food and Agriculture Organization of the United Nations (FAO). (2017). *The Impact of*

620 *disasters and crises on agriculture and Food Security*. Retrieved from

621 www.fao.org/publications

Hydrology and
Earth System
Sciences
Discussions

622    Food and Agriculture Organization of the United Nations (FAO) and Robert B Daugherty
623        Water for Food Institute at the University of Nebraska. (2015). *Yield gap analysis of field*
624        *crops, Methods and case studies*. (V. O. Sadras, K. G. G. Cassman, P. Grassini, A. J. Hall,
625        W. G. M. Bastiaanssen, A. G. Laborte, … P. Steduto, Eds.), *FAO Water Reports* (Vol.
626        41). Rome, Italy.

627    Ghosh, K., Balasubramanian, R., Bandopadhyay, S., Chattopadhyay, N., Singh, K.K., and
628        Rathore, L.S. (2014). Development of crop yield forecast models under FASAL-a case
629        study of kharif rice in West Bengal. *Journal of Agrometeorology*, *16*(1), 1–8.

630    Govindaraju, R.S. (2000). Artificial Neural Networks in Hydrology. I: Preliminary Concepts.
631        *Journal of Hydrologic Engineering*, *5*(2), 115–123. https://doi.org/10.1061/(ASCE)1084-
632        0699(2000)5:2(115)

633    Guha-Sapir, D. (2019). EM-DAT: The Emergency Events Database - Université catholique de
634        Louvain (UCL) - CRED. Retrieved from www.emdat.be

635    Herrera-Estrada, J.E., Satoh, Y., and Sheffield, J. (2017). Spatio-Temporal Dynamics of Global
636        Drought.      *Geophysical      Research      Letters*,      2254–2263.
637        https://doi.org/10.1002/2016GL071768

638    Huang, J., Gómez-Dans, J.L., Huang, H., Ma, H., Wu, Q., Lewis, P.E., Liang, S., Chen, Z.,
639        Xue, J.H., Wu, Y., Zhao, F., Wang, J., and Xie, X. (2019). Assimilation of remote sensing
640        into crop growth models: Current status and perspectives. *Agricultural and Forest*
641        *Meteorology*, *276–277*(July), 107609. https://doi.org/10.1016/j.agrformet.2019.06.008

642    Kim, W., Iizumi, T., and Nishimori, M. (2019). Global Patterns of Crop Production Losses
643        Associated with Droughts from 1983 to 2009. *Journal of Applied Meteorology and*
644        *Climatology*, *58*(6), 1233–1244. https://doi.org/10.1175/JAMC-D-18-0174.1

645    Maier, H.R. and Dandy, G.C. (2000). Neural networks for the prediction and forecasting of
646        water resources variables: a review of modelling issues and applications. *Environmental*
647        *Modelling & Software*, *15*(1), 101–124. https://doi.org/10.1016/S1364-8152(99)00007-9

648    May, R., Dandy, G., and Maier, H. (2011). Review of Input Variable Selection Methods for
649        Artificial Neural Networks. In G. Dandy (Ed.), *Artificial Neural Networks -*
650        *Methodological Advances and Biomedical Applications* (p. Ch. 2). Rijeka: InTech.
651        https://doi.org/10.5772/16004

652    Mckee, T.B., Doesken, N.J., and Kleist, J. (1993). The relationship of drought frequency and
653        duration to time scales. *AMS 8th Conf. Appl. Climatol.*, (January), 179–184.
654        https://doi.org/citeulike-article-id:10490403

655    Monfreda, C., Ramankutty, N., and Foley, J.A. (2008). Farming the planet: 2. Geographic

Hydrology and
Earth System
Sciences
Discussions

EGU

Open Access

656    distribution of crop areas, yields, physiological types, and net primary production in the
657    year    2000.    *Global    Biogeochem.    Cycles*,    *22*(1),    1–19.
658    https://doi.org/10.1029/2007GB002947

659    Montesino Pouzols, F. and Lendasse, A. (2010). Effect of different detrending approaches on
660    computational intelligence models of time series. In *The 2010 International Joint*
661    *Conference    on    Neural    Networks    (IJCNN)*    (pp.    1–8).    IEEE.
662    https://doi.org/10.1109/IJCNN.2010.5596314

663    Naresh Kumar, M., Murthy, C.S., Sesha Sai, M.V.R., and Roy, P.S. (2012). Spatiotemporal
664    analysis of meteorological drought variability in the Indian region using standardized
665    precipitation    index.    *Meteorological    Applications*,    *19*(2),    256–264.
666    https://doi.org/10.1002/met.277

667    Rahmati, O., Falah, F., Dayal, K.S., Deo, R.C., Mohammadi, F., Biggs, T., Moghaddam, D.D.,
668    Naghibi, S.A., and Bui, D.T. (2020). Machine learning approaches for spatial modeling
669    of agricultural droughts in the south-east region of Queensland Australia. *Science of the*
670    *Total Environment*, *699*, 134230. https://doi.org/10.1016/j.scitotenv.2019.134230

671    Reynolds, C.A., Yitayew, M., Slack, D.C., Hutchinson, C.F., Huete, A., and Petersen, M.S.
672    (2000). Estimating crop yields and production by integrating the FAO Crop Specific
673    Water Balance model with real-time satellite data and ground-based ancillary data.
674    *International    Journal    of    Remote    Sensing*,    *21*(18),    3487–3508.
675    https://doi.org/10.1080/014311600750037516

676    Sawasawa, H. (2003). *Crop yield estimation: Integrating RS, GIS and management factors, a*
677    *case study of Birkoor and Kortigiri Mandals. MSc thesis*. International Institute for Geo-
678    information    Science    and    Earth    Observation.    Retrieved    from
679    http://www.itc.nl/library/papers_2003/msc/nrm/sawasawa.pdf

680    Sheffield, J. and Wood, E.F. (2011). *Drought: Past problems and future scenarios*. (P.
681    Earthscan, Ed.). London.

682    Udmale, P., Ichikawa, Y., Ning, S., Shrestha, S., and Pal, I. (2020). A statistical approach
683    towards defining national-scale meteorological droughts in India using crop data.
684    *Environmental Research Letters*, *15*(9). https://doi.org/10.1088/1748-9326/abacfa

685    van Klompenburg, T., Kassahun, A., and Catal, C. (2020). Crop yield prediction using machine
686    learning: A systematic literature review. *Computers and Electronics in Agriculture*,
687    *177*(July), 105709. https://doi.org/10.1016/j.compag.2020.105709

688    White, M.A., Thornton, P.E., and Running, S.W. (1997). A continental phenology model for
689    monitoring    vegetation    responses    to    interannual    climatic    variability.    *Global*

690     *Biogeochemical Cycles*, *11*(2), 217–234. https://doi.org/10.1029/97GB00330

691     World Meteorological Organization (WMO). (2006). *Drought monitoring and early warning:*

692     *concepts, progress and future challenges*. *WMO-No. 1006*. Geneva, Switzerland.

693     Retrieved                                                                      from

694     http://www.droughtmanagement.info/literature/WMO_drought_monitoring_early_warni

695     ng_2006.pdf

696     Wu, X., Vuichard, N., Ciais, P., Viovy, N., Wang, X., Magliulo, V., and Wattenbach, M.

697     (2016). ORCHIDEE-CROP (v0), a new process-based agro-land surface model: model

698     description and evaluation over Europe, 857–873. https://doi.org/10.5194/gmd-9-857-

699     2016

700