

1 **Season-Ahead Forecasting of Water Storage and Irrigation**
2 **Requirements**

3 *An Application to the Southwest Monsoon in India*

4
5 **Arun Ravindranath^{1*}, Naresh Devineni¹, Upmanu Lall², Paulina Concha Larrauri³**

6 ¹Department of Civil Engineering, Center for Water Resources and Environmental Research (City Water Center),
7 NOAA Center for Earth System Sciences and Remote Sensing Technologies, City University of New York (City
8 College), New York, NY 10031

9 ²Department of Earth and Environmental Engineering, Columbia Water Center, The Earth Institute, Columbia
10 University, New York, NY 10027, USA

11 ³Columbia Water Center, The Earth Institute, Columbia University, New York, NY 10027, USA

12
13 **Short title:** Forecasting seasonal crop water stress

14
15
16
17
18
19 * - Contact Author's email: aravind000@citymail.cuny.edu

20

21

22

23

24 **Abstract**

25 Water risk management is a ubiquitous challenge faced by stakeholders in the water or
26 agricultural sector. We present a methodological framework for forecasting water storage
27 requirements and present an application of this methodology to risk assessment in India. The
28 application focused on forecasting crop water stress for potatoes grown during the monsoon
29 season in the Satara district of Maharashtra. Pre-season large-scale climate predictors used to
30 forecast water stress were selected based on an exhaustive search method that evaluates for
31 highest Rank Probability Skill Score and lowest Root Mean Squared Error in a leave-one-out
32 cross validation mode. Adaptive forecasts were made over the years 2001 through 2013 using the
33 identified predictors and a non-parametric k-nearest neighbors approach. The accuracy of the
34 adaptive forecasts (2001-2013) was judged based on directional concordance and contingency
35 metrics such as hit/miss rate and false alarms. Based on these criteria, our forecasts were correct
36 nine out of thirteen times, with two misses and two false alarms. The results of these drought
37 forecasts were compared with precipitation forecasts from the Indian Meteorological Department
38 (IMD). We assert that it is necessary to couple informative water stress indices with an effective
39 forecasting methodology to maximize the utility of such indices, thereby optimizing water
40 management decisions.

41

42 **Keywords:** Crop stress, water risk, seasonal forecasts, climate-information, deficit, monsoon
43 prediction, contract farming, agricultural drought risk

44

45

46

47

48

49

50

51

52

53

54

55

56

57 **1. Introduction**

58 Monitoring and forecasting systems can aid in pinpointing mitigation tactics for water
59 security and water resources management. There is a continued interest in forecasting and
60 monitoring systems that can inform planners and decision-makers in various water-dependent
61 sectors at sufficient lead times and with increasingly higher levels of accuracy and reliability.
62 The agricultural sector is perhaps the greatest example of this, being a heavily water-dependent
63 sector that serves as the economic backbone of a country. The agricultural sector consumes
64 more freshwater than any other economic sector, with an estimated 1,300 m³/cap/yr needed to
65 maintain an adequate diet (Rockstrom et al., 2009). Significant increases of water will be
66 required to produce food by 2050, ranging from 8,500 to 11,000 km³/yr, depending on to what
67 extent rainfed and irrigated agricultural systems improve (Rockstrom et al., 2009). Additionally,
68 to maintain high yields, irrigation will continue to be an important buffer to climate shocks. This
69 is especially true when one considers that almost all of the world's major agricultural lands are
70 located in the most drought-prone areas of the world (Mishra and Desai, 2006). Hence,
71 developing forecasting techniques to improve how we address irrigation requirements, water
72 storage requirements and crop water stress is a major step in dealing with the larger issue of
73 water resources management at local, regional and global scales. The present study focuses on
74 forecasting water storage and irrigation requirements in the agricultural sector as one important
75 dimension to the larger issue of drought forecasting and water resources management, with an
76 application of such forecasting to the monsoonal climate of India.

77

78 Existing forecasts either deal directly with basic hydrologic or meteorological variables,
79 such as precipitation, temperature and soil moisture, or they work with proxies of drought, often
80 in the form of indices such as the Standardized Precipitation Index, or SPI (McKee et al, 1993),
81 the Palmer Drought Severity Index, or PDSI (Palmer, 1965), the Standardized Precipitation
82 Evapotranspiration Index, or SPEI (Serrano et al, 2010), and the Normalized Difference
83 Vegetation Index, or NDVI, among others. A comprehensive list of indices used in drought
84 forecasting can be found in Heim (2002), Mishra and Singh (2010) and Liu and Pan (2016). The
85 forecast of basic variables requires subsequently integrating these forecasts into a product that
86 can estimate water storage or irrigation requirements, as these variables do not immediately
87 divulge such information. This represents a challenge by itself. In light of this limitation, in this
88 paper, we present a crop water stress index that is defined and constructed based on the work by
89 Devineni et al (2013). The advantage of this particular index, hereby known as the cumulative
90 deficit index (CDI), is that it accounts for the variability in water supply and demand while
91 incorporating information specific to a particular crop of interest. CDI is derived by
92 accumulating differences in supply (rainfall) and demand (crop water requirement), and with
93 very few crop input parameters. The CDI is a determinant of water stress faced by the crop and
94 hence of the dependence of the crop yield on water availability. It can be interpreted as the water
95 that is required from external storage beyond rainfall to meet demand (Devineni et al, 2013;
96 Devineni et al, 2015). Therefore, the index directly informs water storage and irrigation
97 requirements.

98

99 The primary focus of this paper will be on exploring the possibility of providing forecasts
100 for CDI by investigating the sources of predictability and developing statistically verifiable
101 models for the season-ahead probabilistic forecasts. Significant crop water deficits can adversely
102 impact the crop production or water reserves and lead to high-energy costs for pumping
103 groundwater for irrigation to maintain yield. The seasonal forecasting of CDI provides a way for
104 institutional planning and action in this context to reduce the climate-related water risks in
105 agriculture, which is one of the largest consumers of water. An application of CDI forecasting is
106 presented for the state of Maharashtra in India to verify whether advance reliable forecasts for
107 potato-based CDI can be developed. A non-parametric k-nearest neighbor (kNN) bootstrapping
108 algorithm as described in Lall and Sharma (1996) is employed for forecasting CDI using pre-
109 season large-scale climate indices. This is a simple probabilistic forecasting procedure that
110 captures uncertainty. We examine these forecasts and suggest ways of interpreting them in a
111 manner that can aid stakeholders in the agricultural water resources sector in addressing the
112 fundamental questions about irrigation and water storage requirements. These forecasts will then
113 be compared to precipitation forecasts for the same season in the same area of India as given by
114 the Indian Meteorological Department (IMD).

115

116 In section 2, we present a survey of the existing forecasting systems in monsoonal
117 climates and their skill and limitations. In section 3, we discuss the background and scientific
118 basis of CDI, including its explicit formulation and governing equations. In section 4, we get
119 into a thorough description of the case study and all steps involved, including background
120 information relating to the case study and location, data collection and processing, a complete
121 description of the forecasting model, methods and predictor selection scheme. Section 5 presents
122 the results of the forecast, a discussion of these results and their implications, and a comparison
123 of our results with those of IMD. Finally, section 6 summarizes and concludes the paper.

124

125 **2. A Brief Review of the Current Forecasting Systems for Water Management in** 126 **Monsoonal Climates**

127 A number of forecasting methodologies have been proposed or developed for water
128 management and agricultural planning. Shah and Mishra (2016) investigated the goodness of the
129 Global Ensemble Forecast System (GEFS) for generating medium-range (~7 day) drought
130 forecasts in India, and found that the GEFS has higher forecasting skill during the non-monsoon
131 season than monsoon season for both temperature and precipitation, largely due to inability to
132 represent the intraseasonal variability during the monsoon season. This forecasting system tends
133 to forecast temperature variables with higher skill than precipitation and has variable skill
134 according to region. Hence, there is sensitivity to intraseasonal variation, which monsoon
135 climates are notorious for, and regional variation as well. Mishra and Desai (2005) used well-
136 chosen linear stochastic models (ARIMA) to forecast SPI- 3, 6, 9, 12, and 24 as a drought proxy
137 in the Kansabati River Basin, an important source of water for irrigation and an area in which
138 crops are grown, in the Purulia district of West Bengal, India at lead times of 1, 2, 3, 4, 5 and 6
139 months. Highest skill, as measured by the correlation coefficient between observed and model-
140 predicted SPI series, occurred at shorter lead times, with correlation values between 0.80 and
141 0.93 depending on which SPI series was forecasted.

142 Asoka and Mishra (2015) forecast vegetation anomalies (as NDVI) at the regional scale
143 as a proxy of vegetation health, and thus moisture availability. The model used NDVI, root-zone
144 soil moisture, and sea surface temperature (SST) at one to three months lead time to develop the
145 vegetation anomaly forecast, and skill was highest at one month lead time and much lower for
146 two and three months lead time as measured in a validation phase by examining the R^2 statistic
147 and by plotting the observed NDVI against the model-interpolated series for one-, two-, and
148 three- month lead times. Skill also varied based on location in space and, in particular, was
149 lower during the monsoon season (JJAS) likely due to the effect of intraseasonal variability of
150 the monsoon system on agricultural practices. Belayneh and Adamowski (2012), in the interest
151 of drought forecasting, forecasted SPI 3 and SPI 12 over lead times of one and six months in the
152 Awash River Basin in Ethiopia using Artificial Neural Network, Wavelet Neural Network and
153 Support Vector Regression models and similarly found that forecast skill was higher at the
154 shorter lead time.

155 Kar et al (2012) considered Multi-Model Ensemble (MME) methods in both a
156 deterministic and probabilistic context. It was found that the individual member models showed
157 poor skill in simulating monsoon interannual variability and that on average spatially, a MME
158 scheme that uses the member models as predictors in a point-by-point multiple regression as a
159 means of averaging the member model forecasts outperforms the other schemes mentioned in the
160 paper in forecasting precipitation. However, it was found that even here, none of the three MME
161 schemes had any usable skill in a certain region of India, and it was concluded that a
162 probabilistic system would work better. When probabilistic forecasts were generated
163 (probabilistic MME) and evaluated for skill, Rank Probability Skill Score (RPSS) was positive
164 for the best scheme, in only the northern most parts of India and a few scattered points in north
165 and central India.

166 Finally, Shah et al (2017) examined how different forecast products can be used
167 operationally to provide hydrologic forecasts (e.g. for precipitation, temperature) for India at a 7
168 – 45 day accumulation period, which is critical for agricultural and water resource planning.
169 Forecast skill was evaluated on the basis of correlation with observations, median absolute error
170 (MAE) and the critical success index (CSI). Four forecast products from Indian Institute of
171 Tropical Meteorology (IITM) were compared with Climate Forecast System version 2 (CFSv2)
172 and Global Ensemble Forecast System version 2 (GEFSv2) forecast products, and it was found
173 that the meteorological variables predicted from the IITM products showed superior skill for all
174 accumulation periods. The key point here is that the IITM ensemble is postulated to capture
175 intraseasonal variability of rainfall during the monsoon season.

176 A variety of forecasts for seasonal rainfall are available at different lead times and with
177 different skills depending on method, location and measure of skill as demonstrated in the review
178 above. However, none of these directly inform irrigation water requirements for a specific crop
179 or of the potential reduction in yield due to a water deficit that occurs depending on the actual
180 sequence of daily rainfall amounts. Ours is the first paper to directly address forecasting a
181 measure that can be tuned to a specific crop using historical observations and crop models or
182 crop performance data.

183

184 **3. The Cumulative Deficit Index: Background and Scientific Basis**

185 Our interest in this study is to provide one-season-ahead forecasts of irrigation and water
 186 storage requirements for water resources management in the agricultural sector, and
 187 subsequently compare the outcomes of these forecasts with the forecasts issued by IMD. We
 188 begin by developing an index for crop water stress as a means of gauging irrigation
 189 requirements. The index developed and used in this study computes the maximum cumulative
 190 deficit over a growing season between daily water requirement for optimal crop growth and daily
 191 effective rainfall. Variants of this method have been presented in our previous studies for
 192 quantifying the water stress globally (Devineni et al, 2013; Devineni et al, 2015; Chen et al,
 193 2014), and drought indexing for the United States (Etienne et al, 2016; Ho et al, 2016). Given an
 194 n -year record of daily data, our water stress index calculates the day-by-day accumulation of
 195 deficit in rainfall in each of the n growing seasons. The maximum of these seasonal daily deficit
 196 values is taken to be the value of the index for the season. Hence, we give this index the name
 197 *cumulative deficit index*, abbreviated CDI. On a practical level, such an index gives a worst-
 198 case scenario in terms of the seasonal water stress on the crop, and can therefore be interpreted as
 199 the amount of water that should be drawn from external storage to meet water demand. This
 200 may include irrigation, ground water pumping, interbasin transfers, and/or withdrawing water
 201 from a storage or water-harvesting facility.

202 Deficit is estimated as the difference between the seasonal crop water requirement and
 203 effective rainfall for each crop in a given location in the season. Effective rainfall is given as

204

$$205 \quad S_{j,d} = \alpha_j * P_{j,d} \dots (1)$$

206 In Eq. (1), $P_{j,d}$ is the rainfall for a day d in any given year at a location j . α_j is the parameter that
 207 determines the fraction of rainfall that can be utilized by the crops for location j . It accounts for
 208 losses to direct runoff, evaporation and groundwater infiltration. In our study, we set $\alpha_j = 0.7$
 209 (Devineni et al, 2013).

210 The water use for a given crop is estimated based on the expected growth stage and daily
 211 evapotranspiration as

212

$$213 \quad D_{j,d} = k_{c,d}^{(j)} * ET_{0j,d} \dots (2)$$

214

215 In Eq. (2), $k_{c,d}^{(j)}$ is the crop coefficient, which is the ratio of actual evapotranspiration (ET_d) of
 216 a given crop under non-stressed conditions to reference crop evaporation (ET_0). It represents
 217 crop-specific water use at various growth stages of the crop and is typically derived empirically
 218 based on local climatic conditions (Doorenbos and Pruitt, 1977). The accumulated deficit over
 219 a season is then given as

220

$$221 \quad deficit_{j,d} = \max(deficit_{j,d-1} + D_{j,d} - S_{j,d}, 0) \text{ where } deficit_{j,d=0} = 0 \dots (3)$$

222

223
$$CDI_{j,t} = \max(\text{deficit}_{j,d(y)}; d = 1:n_s; t = 1:n); \text{ where } \text{deficit}_{j,d(0)}=0, y=1,\dots,n \dots (4)$$

224

225 In equation (3), $\text{deficit}_{j,d}$ refers to the accumulated daily deficit for any given year with a crop
226 growth period of n_s days in the year, $D_{j,d}$ to total daily water demand, $S_{j,d}$ to the total daily
227 effective rainfall, for geographical location j , and day d ; t refers to a calendar or cropping year;
228 and n is the total number of years in the analysis. For an n -year record, seasonal water stress is
229 evaluated as the maximum cumulative deficit each season and defined here as $CDI_{j,t}$. CDI
230 focuses on the rainfall distribution within the season relative to the crop water demand. It
231 therefore accounts for the timing of planting, different stages of crop growth, and the timing and
232 distribution of rainfall in the season. The index may also be treated as a hydrologic index and
233 forecasted exactly as one would forecast precipitation or temperature variables, or any other
234 water stress or drought index. Depending on the lead time of such forecasts, this can give
235 farmers and other agricultural stakeholders a sufficient amount of planning and preparation time,
236 thus providing them a critical edge in hedging agricultural water risk. This is critical for
237 irrigation and water storage planning. The computation of CDI is illustrated in Fig.1. This
238 figure provides insights on the time-evolving vulnerability to stress arising from deficient rainfall
239 and changes in crop demand.

240

241 **4. Case Study: Forecasting Irrigation Requirements for Potatoes in Maharashtra, India**

242 We provide an application of our general approach to forecast CDI for potatoes grown in
243 the Satara district in Maharashtra, India as an application. The Satara district in Maharashtra is
244 one of the primary regions for sourcing potatoes during the monsoon season (June - September).
245 Satara supplies the majority of the potatoes processed by the Frito-Lay manufacturing plant in
246 Pune, Maharashtra (Economic Times, 2013). Potato is a major cash crop in Maharashtra and
247 accounts for at least 75% of total production (Nikam, *et al.*, 2008). The average annual rainfall in
248 this arid to semi-arid region is around 350 mm with high inter-annual variability. The region has
249 experienced four droughts (seasonal rainfall below long-term average) since 2001. The ability to
250 predict such droughts with a reasonable accuracy at lead times of three to six months could
251 suggest ways to adapt existing agricultural operations to the anticipated conditions and minimize
252 the impacts of droughts on the agricultural supply chain. Hence, we develop, present and
253 evaluate the results from retrospective forecasts of CDI for the monsoon season over the period
254 2001-2013. The June-July-August-September (JJAS) season is the growing season for potatoes
255 in the Satara district. It is also the core monsoon season for the Indian sub-continent. The
256 forecasts use climate data from three to six months prior to the beginning of the monsoon season
257 as predictors, and forecasts are to be issued in May, one month prior to monsoon onset. This
258 section discusses the full forecasting procedure used to predict CDI for potatoes grown during
259 the JJAS monsoon season in Satara, India. This discussion covers all data used, the data
260 processing steps, prediction selection routine and its results, and the forecasting model itself.
261 Figure 2 presents a flowchart summarizing the entire process.

262

263 4.1: Data Collection and Processing

264 4.1.1: Precipitation and Temperature Data and the CDI

265 Gridded daily rainfall data from 1901 – 2004 available at $1^0 \times 1^0$ spatial resolution from
266 the India Meteorological Department (Rajeevan *et al.*, 2006), and gridded daily temperature data
267 from 1969 – 2005, available at the same spatial resolution from India Meteorological Department
268 are used in this study. Since the daily temperature data is available only for 37 years, we used
269 the daily climatology, i.e. the mean daily temperature, for the remaining 77 years (Devineni *et*
270 *al.*, 2013). The daily climate time series grids were spatially averaged over the Satara district.
271 This process resulted in a time series of daily precipitation and temperature estimates for 104
272 years. The daily Reference Crop Evapotranspiration (ET_0) was developed based on the daily time
273 series of minimum, mean and maximum temperature data, and extraterrestrial solar radiation
274 (Hargreaves and Samani, 1982). The Hargreaves method is used globally to predict ET_0 in
275 regions where data availability is limited to air temperature data (Allen, *et al.*, 1998). Seasonal
276 daily rainfall data from 2005 to 2013 for the Satara district were collected separately from a
277 website maintained by the Agricultural Department of Maharashtra State and used to augment
278 the 104 years of rainfall and temperature data. The CDI was computed for each of these 113
279 seasons using the daily rainfall data and reference crop evapotranspiration. This will serve as the
280 predictand for our forecast model. We remind the reader that Figure 1 illustrates the
281 computation of CDI.

282 CDI as a water stress measure is a proxy of not only crop water stress but also irrigation
283 and water storage requirements. Consider Fig. 1. When daily seasonal rainfall is low or when
284 rainfall enters an inactive phase for a considerable period of time, as displayed by the vertical
285 cyan bars, the amount of daily accumulated water deficit increases to reflect the disparity
286 between water supplied as rainfall and the water required by the crop to sustain itself, as
287 displayed by the red curve in Fig. 1. The highest point, or peak, on the black deficit time series
288 in Fig. 1 is the value of CDI, and it prepares us for the worst-case scenario of deficient water
289 supply for the crop. This can be calculated for multiple crops, each CDI value depending on the
290 specific crop's water demand and the location and time of planting. This gives the stakeholder a
291 conservative estimate of how much additional water is needed beyond what Nature is willing to
292 supply in order to maintain critical yields while apportioning water resources intelligently. Since
293 agriculture tends to be one of the largest consumers of water --- about seventy-percent of all the
294 world's freshwater withdrawals go towards irrigation use (USGS, 2017), and this is in addition to
295 what is rainfed --- this is an integral part of water resources management.

296 The annual time series of the CDI computed for the JJAS season (referred to as Kharif
297 season in India sub-continent) in Satara is presented in Fig. 3. We have standardized the CDI
298 values as the percentage difference each year from the 113-year average of CDI. The long-term
299 average CDI for growing potatoes in Satara is 241 mm. This is equivalent to approximately
300 257,644 gallons of water used for irrigating a one-acre farm of potatoes on average throughout
301 the season. The percent differences in Fig. 3 refer to percentages of this number, i.e. a 10%
302 increase in CDI indicates an additional requirement of 25,764 gallons. From Fig. 3, it is clear that
303 (a) Satara experiences recurrent droughts with intermediate wet periods and (b) there is year-to-
304 year persistence in the incidence of these droughts. Such variations and epochal changes are
305 typically modulated through large-scale global climate patterns. Investigating the relationship
306 between monsoon deficit and the large-scale climate teleconnections could enable the

307 development of models that can be used to understand and predict the variability in the CDI in
308 the region.

309

310 4.1.2: Climate Precursors and Climate Data

311 Our goal was to develop a simple statistical model for predicting CDI for potatoes grown in
312 Satara. The generalized climate forecast models available at low spatial resolution are not
313 specific enough for this task. Consequently, the first objective was to identify appropriate climate
314 predictors before the monsoon starts in June. There is an extensive history of developing long-
315 range predictions of monsoon rainfall that are based on various regional to large-scale climate
316 predictors (Walker, 1924; Thapliyal, 1987). A variety of seasonal forecasts of the all India
317 Summer Monsoon Rainfall (ISMR) are documented and available for reference (Gadgil et al.,
318 | 2007; Kumar et al., 1995).

319

320 It is well established that inter-annual climate modes such as ENSO associated with
321 anomalous Sea Surface Temperature (SST) conditions in the tropical Pacific Ocean influence the
322 inter-annual variability of ISMR (Parthasarathy and Pant, 1985; Shukla and Paolino, 1983).
323 Anomalously warm tropical eastern Pacific SSTs (El Niño) are associated with a drier-than-
324 normal ISMR, whereas anomalously cool tropical eastern Pacific SSTs (La Niña) are associated
325 with a wetter-than-normal ISMR (Sikka, 1980; Parthasarathy and Panth, 1985; Rasmusson and
326 Carpenter, 1983). Ihara, *et al.* (2007) have suggested that the ENSO warm (cool) phases shift the
327 location of the tropical Walker circulation and cause deficient (excessive) rainfall by suppressing
328 (enhancing) the convection over India. Hence, ENSO indices were chosen to be among the
329 candidate predictors for the forecast model. Raw monthly SST data for the Niño 3, Niño 4, Niño
330 12 and Niño 34 indices were taken from the KNMI climate explorer database (KNMI, 2016).

331

332 For each given raw ENSO index (3, 4, 12 and 34), we considered three different types of
333 derived ENSO indices: a December-January-February (DJF) seasonal average, a March-April-
334 May (MAM) seasonal average, and a MAM minus DJF (MAM-DJF) differenced time series.
335 Among the Niño indices calculated, the change in the tropical Pacific SSTs from December to
336 May (MAM-DJF trend) was found to be of significance by previous investigators. Shukla and
337 Paolino (1983) found the correlation coefficient between the MAM-DJF trend pressure
338 anomalies and the ISMR to be a significant -0.42. Their investigation showed that the Darwin
339 pressure anomalies decrease from DJF to MAM before the occurrence of heavy monsoon rainfall
340 and increase prior to the occurrence of deficit monsoon rainfall. Parthasarathy et al. (1988) found
341 the correlation coefficient between this winter-to-spring trend and ISMR over the period 1951-
342 1980 to be between 0.40 and 0.52 in magnitude, depending on the specific region within the
343 tropical Pacific. Hence, MAM-DJF trends from Niño 3, Niño 4, Niño 12 and Niño 34 were
344 considered to be potential model predictors. Parthasarathy et al. (1988) found that the MAM-
345 averaged tropical Pacific SSTs over the box 14 N to 20 N, 176 E to 160 W had a correlation of -
346 0.40 with ISMR, convincing us to consider this average as well. In addition to the MAM and
347 MAM-DJF averages, we computed the winter season (DJF) average, although DJF-averaged
348 tropical Pacific SSTs were not found to be significant in the literature. However, it is worth
349 noting that Parthasarathy et al. (1988) found that the correlation coefficient between the Darwin
350 | SLP during the DJF season and ISMR was +0.39.

351 As the concurrent season (JJAS) state of ENSO has an important, well-documented impact
352 on ISMR, we also elected to include the Niño 3.4 JJAS average. As mentioned earlier, an El
353 Niño event during the JJAS season is strongly associated with an anomalously dry JJAS rainfall
354 season in India, while a La Niña event during the JJAS season is strongly associated with an
355 anomalously wet JJAS rainfall season in India, prompting our choice. We coupled the JJAS
356 seasonal average for the Niño 3.4 index with forecasts of the JJA and JAS seasonal averages for
357 the Niño 3.4 index. These forecasts were obtained from the International Research Institute for
358 Climate and Society (IRI) ENSO forecast page and covered the period 2002-2013. These
359 forecasts can be used to forecast JJAS monsoon CDI in place of the observed Niño 3.4 JJAS
360 values on a real-time basis. These forecasted values were averages of the projections from at
361 least six distinct statistical/dynamical models, with one average for the JJA season and one
362 | average for the JAS season. Together, we start with a total of thirteen ENSO-based indices.

363
364 Other candidate predictor variables include concurrent season (JJAS) eastern Indian Ocean
365 SSTs known as the Indonesian Throughflow, or ITF. Warm, low-salinity water from the Pacific
366 is introduced into the Indian Ocean via the ITF and is considered to be an integral component in
367 the heat and hydrological budget of the Indian Ocean (Gordon et al., 1997). The ITF waters are
368 also believed to influence SSTs and associated ocean-atmosphere coupling within the Indian
369 Ocean, making it an important aspect of monsoon climate research (Gordon et al., 1997). Thus,
370 the ITF was also selected to be a candidate predictor in the model. During the JJAS monsoon
371 season, the ITF is strengthened considerably, allowing an abundant amount of relatively warm
372 water to be injected into the Indian Ocean. Eastern Indian Ocean SSTs during the JJAS season
373 correspond to enhanced (suppressed) atmospheric convection during the anomalous warming
374 (cooling) of the Indian Ocean waters, which in turn supplies (robs) the developing monsoon of
375 much-needed moisture. We found that the Spearman rank correlation coefficient between CDI in
376 Satara and the average SST anomalies over 20° N and 5° S and 100° E and 130° E (the region
377 representing ITF) during the JJAS season is around -0.35 (statistically significant at the 95%
378 level), suggesting that warm conditions in the ITF region result in below-normal CDI, or low
379 crop water stress. Figure 4 presents the field correlation map of SST anomalies with CDI. For
380 these reasons, we chose concurrent season ITF data to be a candidate predictor. The ITF data was
381 collected from the IRI data library and consists of two components: an observation component
382 and a forecasted component. The observations consist of measured eastern Indian Ocean SST
383 anomalies during the JJAS season from 1901 through 2013. The forecasts consist of JJAS-
384 season ITF values retrospective from the ECHAM4.5 global climate model and cover the period
385 2001-2013. Skillful forecasts for the tropical SSTs based on coupled ocean-atmospheric general
386 circulation models have been in operation from various climate centers since 1998. Hence, in
387 the forecasting scheme, we used the ITF derived from forecasted SST state issued in May from
388 ECHAM4.5 operational forecasting center (available from IRI data library:
389 http://iridl.ldeo.columbia.edu/SOURCES/.IRI/.FD/.ECHAM4p5/.Forecast/.ca_sst/ensemble24/;
390 Li and Goddard, 2005; van den Dool, 2007; Roeckner et al., 1996). The observed JJAS ITF data
391 are used to train the model, while the retrospective JJAS ITF forecasts are used to make forecasts
392 for the years 2001 – 2013.

393

394 **4.2: The Forecasting Procedure**

396 Given a pool of candidate predictors, the next step is to select the best subset of those
397 predictors. The predictors used in the forecasting model were chosen based on an exhaustive
398 search method. In the exhaustive search method, all possible combinations of the candidate
399 predictor variables are used to develop models that are cross-validated on historical data. Skill
400 metrics are then used to compare the predictive accuracy of each combination. In the present
401 study, we began with 113 years of CDI data and fourteen candidates: Niño 3 DJF, Niño 3 MAM,
402 Niño 3 MAM-DJF, Niño 4 DJF, Niño 4 MAM, Niño 4 MAM-DJF, Niño 12 DJF, Niño 12
403 MAM, Niño 12 MAM-DJF, Niño 34 DJF, Niño 34 MAM, Niño 34 MAM-DJF, Niño 34 JJAS
404 and ITF. The exhaustive search method utilized the kNN cross-validation algorithm and forty
405 years of training data (1901-1940) to build forecast distributions for each of the years 1941-2013.
406 At each step, the training data was updated to include data from all of the years up until the year
407 being cross-validated. Thus, we always use only the historical data and update the model each
408 year with the information of the previous year, much as a regular user of the forecast system
409 would have to do. These forecasting distributions, built over a 73-year record (1941 to 2013)
410 were created successively for every unique combination of two variables, every unique
411 combination of three variables, so on and so forth until we reached the entire pool of predictors.

412 For each and every possible unique combination of the predictor variables, we obtain a
413 matrix of seventy-three columns. For each of these seventy-three (73) years, the squared error
414 and rank probability score (Epstein, 1969; Murphy, 1969, 1971; Candille and Talagrand, 2005)
415 were computed, and from this the root mean squared error (RMSE) and rank probability skill
416 score (RPSS) were computed. In this manner, a single RPSS value and RMSE value were
417 calculated for every possible combination of the predictor variables. We chose the following
418 combination of predictors based on the relative optimality of both their RPSS and RMSE scores:
419 Niño 12 MAM-DJF, Niño 34 MAM-DJF, and ITF, and this set of variables had an RMSE of
420 49.25 mm of required (JJAS) seasonal water storage and RPSS of 0.26. We devised a simple but
421 effective decision rule for determining the optimal choice of predictors based on ranking the
422 metric values. This is especially useful when the number of combinations of variables is
423 unwieldy. Optimality was determined by assigning a rank number to the RMSE and RPSS
424 values in such a way that the number 1 was assigned to the lowest RMSE value, 2 to the second
425 lowest RMSE value, and so on, and the number 1 was assigned to the largest RPSS value, 2 to
426 the second largest RPSS value, and so on. For a fixed number of cross-validated predictor
427 candidates, and for each RMSE/RPSS pair, one pair for each combination of predictors, we
428 determined an RMSE and RPSS rank and took the sum of these ranks. The smallest of all of
429 these sums corresponds to the best or optimal set of predictors among all possible sets of cross-
430 validated predictors. We then compared the rank sum along with the number of predictors to
431 choose the best set of predictors. The chosen trio of predictors mentioned above had the
432 unequivocally highest value of RPSS and second lowest RMSE value out of all possible
433 combinations of the original set of seventeen candidates, the lowest RMSE being only slightly
434 smaller at 48.92 mm. Conceptually, this procedure is similar to the “best subsets regression” or
435 “step-wise regression” (Helsel and Hirsch, 2002), but in the spirit of using kNN algorithm for
436 forecasting, we designed this selection scheme to use the kNN algorithm instead.

437 CDI forecasts were subsequently made using the selected set of predictors. The forecast
438 procedure is tested using the leave-one-out cross-validation method. Each historical observation

439 is omitted in turn, and the model is developed using the remaining years of data. A prediction of
440 the observation that was not kept in the model-building set is then made and compared with the
441 actual outcome for that year. Results from a variant of this approach are presented in the next
442 section. The CDI for the 2001 Kharif season is predicted using the model developed based on
443 data from 1901 – 2000. Similarly, the CDI for 2002 is predicted based on the model that is
444 developed using the data from 1901 – 2001. Thus, as we move from year to year, we update the
445 model observations and predict the future state.

446

447 *4.2.2: The k-Nearest Neighbors Real-Time Forecasting Model*

448 The forecasts were developed using a non-parametric k -nearest neighbors (k-NN) model.
449 This is a data-driven approach that develops a conditional probability distribution of the CDI
450 given the predictors by first identifying the k -historical climate conditions that are most similar to
451 the current values of the climate predictors and then randomly drawing the vector of CDI values
452 in the historical data that correspond to these k neighbors. The neighbors are weighted so that
453 the closer or more similar neighbors are chosen more often than those further away. The key
454 steps are as follows.

455 Let \mathbf{X} be the design matrix of size $n \times p$, where p = number of predictors selected from
456 the original pool of candidates. Let \mathbf{x}_i denote the i^{th} row of \mathbf{X} . Hence, \mathbf{x}_i is a vector containing
457 the values of each of the p predictor variables during year i . Denoting the current values of the
458 predictors by \mathbf{x}_c , the idea is to find k such predictor vectors from the historical record (i.e. find k
459 values of \mathbf{x}_i with $i < c$) that are most "similar" to the value of \mathbf{x}_c and use this information to
460 construct a sampling distribution of CDI from which we can issue probabilistic forecasts. The
461 number of neighbors in the model, or k , represents the number of degrees of freedom in the
462 model, and should be chosen with care, as the choice of k affects the skewness and level of
463 uncertainty in the sampling distributions. After trying several different values for k , we found an
464 optimal value to be $k = 25$. Rajagopalan and Lall (1999) recommend that, as a rule of thumb
465 based on asymptotic arguments, k be roughly equal to \sqrt{n} , where n = the total number of
466 observations. In our situation, it was evident that we required more neighbors than this rule
467 would allow, due to the skewness and variance apparent in the sampling distributions when using
468 only eleven or fewer neighbors. Lall and Sharma (1996) note that if their discrete kernel is used
469 for resampling the conditional bootstrap, then the weights for further neighbors decrease and
470 hence, choosing a larger k may reduce variance of estimate, while potentially increasing the bias
471 in the estimate of the conditional distribution. Cross-validation can also be used to choose an
472 optimal value for k in a given setting.

473 Let \mathbf{y} be the n -dimensional vector of seasonal CDI values, each component of which
474 represents the aggregate water deficit level over the JJAS growing season of every year in the
475 historical record. Assume that \mathbf{y} has been centered and normalized by its historical average to
476 produce mean-normalized anomalies. The first step was to consider the individual distance
477 values (under some specified metric) between \mathbf{x}_c and \mathbf{x}_i for $i = 1, \dots, c-1$. The chosen distance
478 metric for our k-NN model was the Mahalanobis distance (Mahalanobis, 1936)

479

480

$$D_M(\mathbf{x}_c, \mathbf{x}_i) = \sqrt{(\mathbf{x}_c - \mathbf{x}_i)^T \Sigma^{-1} (\mathbf{x}_c - \mathbf{x}_i)} \dots (5)$$

481

482 where Σ is the covariance matrix of the training values in \mathbf{X} . The Mahalanobis distance measure
483 judges point separations in a metric space based on statistical dissimilarity, as opposed to solely
484 physical distance. Hence, the level of similarity between predictor values across different years
485 is determined by the orientation and location of each point relative to the scatterplot of the
486 predictor data. Large distances from \mathbf{x}_c represent predictor values that are statistically anomalous
487 in the context of the predictor data.

488 After the Mahalanobis distances had been calculated, the k (with $k = 25$) smallest distance
489 values were selected and the corresponding years in which these distances occurred were noted.
490 These years, hereby referred to as the *analog years*, are the years during which the predictor
491 signals were most similar to those of the current year. The vector-valued predictors during these
492 analog years are referred to as the *neighbors* of \mathbf{x}_c .

493 The final step was to resample CDI values from the analog years. The resampling
494 technique employed is a nonparametric method known as the *bootstrap* (Efron, 1979; Efron and
495 Tibshirani, 1993). The idea behind the bootstrap component is to sample with replacement from
496 a pool of data using the underlying distribution that generated the data to guide the sampling
497 process. We chose not to assign a parametric family of distributions to the CDI data, and instead
498 estimated its underlying distribution non-parametrically using a kernel density estimator. This
499 non-parametric method of k -NN bootstrapping was first introduced in Lall and Sharma (1996).
500 Applications of the methods using different variants have since been presented (see for example,
501 Rajagopalan and Lall, 1999, Souza and Lall, 2003 and references therein). We employed the
502 same discrete resampling kernel proposed in Lall and Sharma (1996), which has the general form
503 $K(j) = 1/(j*S)$ with $S = \sum_{j=1}^k 1/j$, where j is the rank of each neighbor of \mathbf{x}_c , a rank of $j=1$
504 assigned to the closest neighbor and a rank of $j=k$ assigned to the most distant neighbor. Our
505 strategy was to build this kernel density estimator based on the ranks of the selected neighbors
506 and resample the predictand values from these analog years. We resampled from the twenty-five
507 analog CDI values 1,000 times, and each of the twenty-five values was resampled proportionally
508 to the probability of its occurrence as determined by the density estimator.

509

510 4.2.3: Analyzing the k -NN Results

511 The way in which model results are interpreted and presented is important for potential
512 stakeholders. In this case study, our interest was in forecasting the CDI for a given potato
513 growing season in Satara. The information from these forecasts can be of great use to potato
514 farmers in Satara as well as corporations with investments in these farming areas. This
515 necessitates a clear and concise communication of the forecast results.

516 The output of the k -NN model was a time series for each forecasted year consisting of
517 1,000 realizations. This is the sampling distribution for the CDI and consists of mean-
518 normalized anomaly values from the analog years converted to percentage values. As stated in
519 the previous section, the deficit value from each analog year in the sampling distribution is

520 represented proportionally to its probability of occurrence as assigned by a kernel density
521 estimator. The sampling distribution is used to issue one-season-ahead probabilistic forecasts
522 (i.e. the likelihood of a deficit for the forthcoming growing season). There are a whole slew of
523 possibilities when it comes to using these sampling distributions for probability-based forecasts.
524 Our approach includes the following for a given forecasted growing season:

- 525 1. A boxplot depicting the sampling distribution with the observed percent anomaly value
526 superimposed on the boxplot for every growing season forecasted. In using predictand
527 anomalies, the historical mean becomes the zero line in the coordinate plane of the
528 boxplot.
- 529 2. A three-category forecasting system with the categories “above normal”, “normal” and
530 “below normal”, provided that the historical mean/climatology is the threshold that is
531 desired.
- 532 3. Calculate the probabilities for the categories specified in step 2 from the sampling
533 distribution generated in step 1, and use this to evaluate the accuracy and strength of the
534 forecast based on contingency metrics such as hit rates and false alarms.
- 535 4. To get a sense of the spread/variability in the boxplot distribution, calculate the
536 Interquartile Range (IQR).
- 537 5. Compare the value of the observed percent anomaly of the predictand with the category
538 in which the majority of the probability mass of the sampling distribution lies. This is of
539 central importance in getting a basic sense of the accuracy of the forecast.

540 In general, the construction of such a sampling distribution allows the investigator the freedom to
541 calculate probabilities on many different thresholds. The thresholds should be defined by the
542 particular application and the needs of any stakeholders involved.

543 **5. Case Study: Forecast Results and Discussion**

544 ***5.1: CDI Forecast Results and Comparison with IMD Monsoon Forecasts***

545 We hereby present the results of the CDI forecasts for the 2001 – 2013 JJAS seasons in
546 the Satara district, Maharashtra, India. Forecasts are specifically made in the interest of
547 irrigation requirements for potatoes grown in the Satara district, and we discuss the results in this
548 context. The output of the k-NN model is the forecasting distributions for CDI of the thirteen
549 years and a series of boxplots representing these forecast distributions as shown in Fig. 5. The
550 probabilities calculated from these distributions are shown in Table 1, columns 2 and 3.

551 Figure 5 shows a series of boxplot diagrams depicting the k-NN forecast distributions for
552 CDI over the years 2001 – 2013. All calculations in this Figure, including the construction of the
553 distributions themselves, were done using anomalies of the predictand rather than the raw
554 predictand values. The anomalies were calculated by subtracting the 1901 – 2013 mean from the
555 data and dividing by this mean value and converting the quotient to a percentage. The idea is to
556 gauge the level of seasonal crop water deficit in a forecasted year with respect to the level of
557 crop water deficit that has occurred on average over the entire historical record. This should
558 address the question: how “normal” or “abnormal” is a given level of deficit over a season with
559 respect to everything we have seen or experienced thus far. Given that the forecast is developed
560 one season ahead, the sign of a strong shift in the probability will alert the decision-makers to an
561 anticipated deficit or surplus event.

562 We have created two general possibilities: the observed percent anomaly values (triangles
563 in Fig. 5) can be positive or negative. As the forecasts have been carried out using anomalies
564 instead of raw values, the 1901 – 2013 historical average is re-positioned as the zero line in Fig.
565 5. We calculate the probability under the kNN forecast distribution of observing positive
566 (negative) deficit anomalies for each year in 2001 – 2013. These are retrospective forecasts in
567 the sense that these anomalies have already been observed and recorded but not used in building
568 the model. These probabilities, corresponding observed percent anomalies and IQR values are
569 presented in Table 1. The utility of these forecasts are discussed in section 5.2.

570 Given the above information, we judge the accuracy of the forecasts during any given
571 year on a few simple criteria: the directional agreement between the observed percent predictand
572 anomaly and the median of the forecast distribution (Fig. 5), joint consideration of the forecast
573 probabilities and the observed percent anomaly (Table 1, columns 2, 3 and 4) and the level of
574 uncertainty in the forecast distribution (Fig. 5 and Table 1, column 5). Uncertainty is measured
575 by the IQR of the boxplot distribution. In the present context, we say that a forecast for a given
576 year has *identical directionality* (with respect to the observation) if both the median of this
577 forecast and the observation (as a percent anomaly) are either positive (above the historical
578 average) or negative (at or below the historical average). The absence of identical directionality
579 will be called *dissimilar directionality*.

580 The box-and-whiskers plots shown in Fig. 5 for each year illustrates the range of possible
581 values of the CDI for that year. We have identical directionalities for the years 2001, 2004,
582 2005, 2006, 2007, 2010, 2011, 2012 and 2013. For the years 2001, 2011 and 2012, the model
583 correctly forecasted that the water stress conditions for the Maharastran potatoes would be above
584 the CDI climatology. We can see from Fig. 5 that both the observed percent anomalies
585 (triangles) and the medians for all of these forecasted years are positive. Additionally, Table 1,
586 column 2 shows that the majority of the probability mass of the kNN distribution is placed in the
587 “Above Mean” category for 2001, 2011 and 2012, while column 4 shows that for these years, the
588 observed CDI anomalies are positive. Similarly, for the years 2004, 2005, 2006, 2007, 2010 and
589 2013, the model correctly forecasted that water stress conditions for the potatoes would be below
590 the historical average, and this can be seen from Fig. 5, where the observed anomalies and the
591 medians for all of these forecasted years are negative. Similarly, Table 1, column 3 shows that
592 the majority of the probability mass from the kNN forecasting model was placed on the “Below
593 Mean” category for these years, and the corresponding observed CDI anomalies are also
594 negative. For the years 2002, 2003, 2008 and 2009, we have dissimilar directionalities. The
595 forecasts suggest higher probability values for below average CDI during 2002 and 2003,
596 whereas positive anomalies were observed for these years. Similarly, the forecasts for 2008 and
597 2009 placed the majority of the probability mass on higher than average CDI, suggesting that
598 these years were likely to see higher than normal potato water stress. However, the observed
599 CDI anomalies were negative, implying the opposite scenario.

600 We say that a *hit* has occurred if identical directionality is observed. A *miss* occurs if the
601 forecast implies below average water stress, but the observation shows above average water
602 stress. Finally, a *false alarm* occurs if the forecast implies above average water stress while the
603 observation shows below average water stress. Table 2 shows that the hit rate of the kNN
604 forecasts is 9/13, the miss rate is 2/13 and the false alarm rate is 2/13. Table 3 shows a
605 comparison of our CDI forecasts with seasonal total precipitation forecasts of the India

606 Meteorological Department, abbreviated IMD. The IMD forecast presented here for 2001 is
607 long-range for precipitation in the JJAS season over three climatically homogeneous regions in
608 India: Northwest India, Peninsular India, and Northeast India. Maharashtra is in Peninsular
609 India, and so we refer to this forecast. For 2001, the forecast result was categorized as either
610 normal, above normal or below normal. “Normal” is defined as being within $\pm 10\%$ of the long-
611 period average, or LPA. Beginning in 2003, IMD began offering two-stage forecasts, the first
612 released in mid-April using data up to March and an update in June using data up through May.
613 For both 2011 and 2013, we used the initial country-wide forecast, as the updated forecasts for
614 JJAS could not be found. In 2003, IMD began to divide their forecast results into five
615 categories: drought/deficient, below normal, near normal/normal, above normal and excess.
616 “Deficient” (drought) is defined as JJAS total seasonal rainfall that is less than 90% of the long
617 period average (LPA). “Below normal” is defined as JJAS rainfall that is 90% – 96% of the
618 LPA, “normal” (sometimes called “near normal”) is defined as JJAS rainfall that is 96% – 104%
619 of the LPA, “above normal” is defined as JJAS rainfall that is 104% – 110% of the LPA and
620 “excess” is defined as JJAS rainfall that is more than 110% of the LPA. The IMD forecasts are
621 reported as percentages of the LPA, as shown in column 3 of Table 3. Going by the categories
622 defined by IMD, and comparing these forecasts with actual JJAS seasonal total precipitation
623 anomalies from our gridded rainfall data set, where these anomalies have been calculated with
624 respect to the long period average defined as 1901 – 2013, we classify each forecast as a hit, miss
625 or false alarm as was done with the CDI forecasts. The hit rate for IMD is 1/9, the miss rate is
626 3/9 and the false alarm rate is 5/9. We must bear in mind that the total precipitation forecasts
627 given here are for an entire region that includes the state of Maharashtra, whereas our CDI
628 forecasts are generated based on CDI calculations from the target location of Satara,
629 Maharashtra, India. Hence, our CDI anomalies reflect the conditions of Satara on a much higher
630 resolution than the coarse IMD precipitation anomalies. Furthermore, we are comparing IMD
631 forecasts with actual precipitation totals from Satara, and computed with respect to the 1901 –
632 2013 LPA instead of the 1951 – 2000 LPA of IMD, under the reasonable assumption that the
633 LPA does not change much between those two definitions. While the IMD monsoon forecasts
634 can provide a broad regional understanding of the monsoon conditions, supplementing them with
635 targeted crop-specific forecasts such as ours will help improve agricultural planning and regional
636 water management. To conclude, we used observations for ITF and Nino 3.4 JJAS to generate
637 CDI forecasts for the years 1976 - 2000 and augmented these forecasts with the 2001 - 2013 CDI
638 forecasts depicted in Figure 5. Running the forecasts for a longer period of time, which in this
639 case is 38 years, ensures robustness of the procedure. The hit/false alarm/miss rates resulting
640 from this extended retrospective, adaptive forecast are 24/38 hits, 9/38 false alarms and 5/38
641 misses, respectively. Hence, we are observing 63% hits, which indicates a fairly good, robust
642 forecasting procedure for an informative crop water stress index.

643 We define a *strong forecast* as a forecast in which the probability assigned to one of the
644 two categories is at least 60%. In our situation, ten out of the thirteen years witnessed strong
645 forecasts. A weak forecast runs the risk of being less informative to decision-makers, whereas a
646 strong forecast is much more assertive and definitive, and hence decisions can be made more
647 easily with a strong forecast. The forecasts were also correct for seven of these ten years, as seen
648 in Table 2. The forecasts were correct, but barely weak, for two years (2001 and 2011). If one
649 considers acting only if the probability associated with a CDI forecast is at least 60%, then the
650 forecast is correct seven out of ten times. Raising this to 66% leads to four out of six years
651 classified correctly.

652 It is important to point out that one should also consider the uncertainty (column five in
653 Table 1) when evaluating the power of the forecasts. Knowing the uncertainty is useful since
654 years in which the uncertainty in the forecast is low and there is a strong indication for CDI may
655 lead to different risk management actions than years in which the forecast has strong directional
656 change but is also marked by high uncertainty.

657

658 *5.2: Discussion of Results: The Utility of Targeted Forecasts*

659 It is natural to ask how one might go about using CDI forecasts. Here is a short example
660 of how these forecasts can facilitate decision-making. In 2001, irrigating, or ensuring water
661 storage equal to 294,745 gallons per acre for the potatoes would have been the ideal situation, as
662 this is equivalent to being 14.4% above the average CDI value of 241 mm of water storage
663 equivalent. However, this exact amount cannot be known in the absence of the observed CDI
664 anomaly, which is found in column four of Table 1. Using the median as a plausible estimate for
665 the true anomaly value, roughly 268,980 gallons per acre would have been irrigated or stored
666 instead. A more risk-averse decision-maker may choose to use the upper quartile or even
667 maximum of the kNN-generated sampling distribution as a proxy for the true anomaly value.
668 Such decisions are often made on the basis of prior experience.

669 Although total seasonal rainfall is sometimes used for agricultural water planning, CDI
670 boasts a significant advantage over total seasonal rainfall in this capacity. CDI reliably accounts
671 for water stress incurred by haphazard and erratic patterns of rainfall during the season. A total
672 seasonal rainfall forecast that indicates a growing season with sufficient rainfall will not be
673 reliable when rain throughout the season is erratically distributed in clusters of rainy days,
674 whereby all of the rainfall in a given season occurs within a portion of the season, and the
675 remainder of the season is virtually dry. This is a common occurrence in monsoonal climates,
676 and may have deleterious effects on crops that are vulnerable to prolonged dry periods and/or
677 chunks of time during which rainfall is excessive. Long dry spells throughout the season that
678 can be detrimental to drought-sensitive crops are not accounted for in a measure of total seasonal
679 rainfall, making it possible for the seasonal rainfall to appear sufficient due to sporadic
680 occurrences of large precipitation events. Consequently, it can also serve as a better indicator
681 than regional rainfall to devise index insurance products for agriculture, where crop specific
682 indices can be developed (Skees, 2016). These characteristics of crop water stress must be
683 accounted for in the proper planning and management of agricultural water resources.

684 To illustrate the above point further, we appeal to Figure 6. In this figure, the varying
685 rainfall distribution is indicated by the vertical bars, the crop demand is given by the horizontal
686 line (primary y-axis), and the time series shows the cumulative deficit. The second panel shows
687 two distinct years during which the total seasonal rainfall was 590 mm (vertical line). During
688 one of these two years, the CDI value was 111 mm of water deficit for the potato crop, while the
689 CDI value for the other year was 228 mm. This indicates that the water stress for a particular
690 crop relies on both the magnitude and frequency of seasonal rainfall. When daily seasonal
691 rainfall is more uniform, the daily deficit values do not have the chance to accumulate as much
692 as when rainfall is less uniform and, as a result, when there are persistent dry spells or long
693 precipitation-inactive periods. Panel three shows the resulting cumulative deficit when daily
694 rainfall occurs with greater frequency during the JJAS season and hence the total seasonal

695 rainfall is distributed among the days of the growing season fairly uniformly. The fourth panel,
696 immediately to the right of the third panel, shows the resulting cumulative deficit when rainfall is
697 dominant during the first and last months of the JJAS season. While rainfall events do occur in
698 between, the magnitude of the rainfall is quite low, allowing the seasonal daily CDI time series
699 to spike to a considerably higher maximum value (228 mm) than the CDI time series in panel
700 three (111 mm maximum). The CDI time series recedes and recovers at the end of the season
701 when the rainfall increases in magnitude. Hence, CDI can discriminate between two monsoon
702 seasons which have the same total rainfall, but differ in that one may have rainfall distributed
703 uniformly over the season through modest rainfall events, while the other may have a few intense
704 rain events separated by long dry periods. As we can see, the latter gives rise to a much higher
705 CDI.

706 An interesting and excellent discussion concerning the usability of such science is found
707 in Dilling and Lemos (2011) and several papers cited therein. In the context of that discussion,
708 we find that our forecasting procedure combines the “science push” and “demand pull”
709 approaches to creating scientific usability. The impetus for crafting the CDI and, prior to that,
710 independently developing the k-NN algorithm, was scientific. However, the decision to combine
711 them and apply them as we have to seasonal forecasting was made with agricultural stakeholder
712 interests in mind. As discussed in Dilling and Lemos (2011), the problem of overcoming
713 informal institutional barriers to availing of such seasonal forecasts, namely the idea that current
714 methods of forecasting through weather and climate prediction centers are the only reliable
715 methods, is one potentially faced by our methodology. If this is the case, this is unfortunate, as
716 we feel that our targeted forecasting system is potentially very useful to stakeholders and
717 decision-makers in relevant sectors.

718

719 **6. Summary and Conclusion**

720 A novel crop water stress index, the CDI, was developed here as a way of estimating
721 water storage and irrigation requirements in the interest of agricultural water resources. As
722 management of water resources requires advance knowledge of water risk, the main task
723 accomplished here was the forecasting of CDI as an effective method for understanding and
724 hedging risk. This concept of forecasting CDI for evaluating irrigation requirements was applied
725 to a case study in the Satara district of Maharashtra, India in which the CDI pertaining to
726 potatoes grown in Satara during the Southwest monsoon season was forecasted using large-scale
727 climate indices as predictors in a semi-parametric k-nearest neighbors stochastic model that
728 issues probabilistic forecasts. The climate indices used were defined either concurrent to the
729 monsoon season or three to six months prior. Based on the hit and false alarm rates, the results
730 achieved using our methodology were more favorable than precipitation forecasts conducted by
731 the India Meteorological Department. We also observed in our method a greater tendency
732 towards strong and informative forecasts.

733 This study developed a framework for quantifying and analyzing climate-induced
734 agricultural risks. It is based on (a) developing CDI for assessing crop-specific water risk,
735 irrigation requirements and water storage needs for the agricultural sector; (b) investigating the
736 sources of predictability for this indicator, and (c) developing statistically verifiable models for
737 issuing season-ahead probabilistic forecasts for evaluating water risk and irrigation needs. We

738 can conclude that this is a useful approach to investigating irrigation requirements and that
739 bootstrap-based uncertainty estimation is useful for developing probability-based management
740 models for optimizing agricultural decisions.

741

742 **Acknowledgements**

743 This research was supported by:

744 (a) NSF grant 1360446 (Water Sustainability and Climate, Category 3)

745 (b) PSC-CUNY award 69729-00 47

746 Partial support for the third and fourth authors is provided from PepsiCo Inc. through the
747 WATER RISKS AND SUSTAINABILITY grant. The statements contained within the
748 manuscript/research article are not the opinions of the funding agency or the U.S. government
749 but reflect the authors' opinions.

750

751 **Data Availability**

752 The CDI data used in this paper is available upon request of the contact author.

753

754 **References**

- 755 1. Allen, R.G., Pereira, L.S., Raes, D., Smith, M. Crop Evapotranspiration --- Guidelines
756 for Computing Crop Water Requirements. *FAO Irrigation and Drainage Paper 56*, FAO
757 of the UN, Rome, 15 pp., 1998.
- 758 2. Asoka, A. and Mishra, V. Prediction of Vegetation Anomalies to Improve Food Security
759 and Water Management in India, *GEOPHYS RES LETT*, V. 42, pp. 5290 – 5298, 2015.
- 760 3. Belayneh, A. and Adamowski, J. Standard Precipitation Index Drought Forecasting
761 Using Neural Networks, Wavelet Neural Networks and Support Vector Regression,
762 *Applied Computational Intelligence and Soft Computing*, 13 pp., 2012.
- 763 4. Candille, G. and Talagrand, O. Evaluation of Probabilistic Prediction Systems for a
764 Scalar Variable, *Q J ROY METEOR SOC*, V. 131, pp. 2131 – 2150, 2005.
- 765 5. Chen,
- 766 6. Devineni, N., Perveen, S., and Lall, U. Assessing Chronic and Climate Induced Water
767 Risk Through Spatially Distributed Cumulative Deficit Measures: A New Picture of
768 Water Sustainability in India, *WATER RESOUR RES*, V. 49, pp. 2135-2145, 2013.
- 769 7. Devineni, N., Lall, U., Etienne, E., Shi, D., and Xi, C. America's water risk: Current
770 demand and climate variability, *GEOPHYS RES LETT*, V. 42, 2015.
- 771 8. Dilling, L. and Lemos, M.C. Creating usable science: Opportunities and constraints for
772 climate knowledge use and their implications for science policy, *Global Environmental*
773 *Change*, V. 21, Issue 2, pp. 680 – 689, <https://doi.org/10.1016/j.gloenvcha.2010.11.006>,
774 2011.
- 775 9. Doorenbos, J., Pruitt, W.O. Guidelines for Predicting Crop Water Requirements:
776 *Irrigation and Drainage Paper 24*, FAO of the UN, Rome, 154 pp., 1977.

- 777 10. The Economic Times, India Times, [http://articles.economictimes.indiatimes.com/2013-](http://articles.economictimes.indiatimes.com/2013-09-25/news/42394669_1_drip-irrigation-farming-market)
778 [09-25/news/42394669_1_drip-irrigation-farming-market](http://articles.economictimes.indiatimes.com/2013-09-25/news/42394669_1_drip-irrigation-farming-market) (Accessed: 3/1/2018), 2013.
- 779 11. Efron B. Bootstrap Methods: Another Look at the Jackknife. *Annals of Statistics*, V. 7, pp.
780 1-26, 1979.
- 781 12. Efron, B. and Tibishirani, R. *An Introduction to the Bootstrap*. Chapman and Hall, New
782 York, 456 pages, 1993.
- 783 13. Epstein, E.S. A Scoring System for Probability Forecasts of Ranked Categories, *J APPL*
784 *METEROL*, V. 8, pp. 985 – 987, [https://journals.ametsoc.org/doi/pdf/10.1175/1520-](https://journals.ametsoc.org/doi/pdf/10.1175/1520-0450(1969)008%3C0985:ASSFPF%3E2.0.CO%3B2)
785 [0450\(1969\)008%3C0985:ASSFPF%3E2.0.CO%3B2](https://journals.ametsoc.org/doi/pdf/10.1175/1520-0450(1969)008%3C0985:ASSFPF%3E2.0.CO%3B2), 1969.
- 786 14. Etienne, E., Devineni, N., Khanbilvardi, R., and Lall, U. Development of a Demand
787 Sensitive Drought Index and its Application for Agriculture over the Conterminous
788 United States, *J HYDROL*, V. 534, 219–229,
789 <http://dx.doi.org/10.1016/j.jhydrol.2015.12.060>, 2016.
- 790 15. Gadgil, S., Rajeevan, M., and Francis, P.A. Monsoon Variability: Links to Major
791 Oscillations Over the Equatorial Pacific and Indian Oceans, *CURR SCI INDIA*, V. 93,
792 pp. 182 – 194, 2007.
- 793 16. Gordon, A.L., Ma, S., Olson, D.B., Hacker, P., Ffield, A., Talley, L.D., Wilson, D., and
794 Baringer, M. Advection and diffusion of Indonesian throughflow water within the Indian
795 Ocean South Equatorial Current. *GEOPHYS RES LETT*, V. 24, pp. 2573-2576,
796 <http://dx.doi.org/10.1029/97GL01061>, 1997.
- 797 17. Hargreaves, G.H. & Samani, Z.A. Estimating Potential Evapotranspiration. *Journal of the*
798 *Irrigation and Drainage Division*, V. 108, pp. 225-230, 1982.
- 799 18. Heim Jr., R.R. A Review of Twentieth-Century Drought Indices Used in the United
800 States. *Bulletin of the American Meteorological Society*, V. , pp. 1149 – 1165, 2002.
- 801 19. Helsel, D.R. & Hirsch, R.M. *Statistical Methods in Water Resources*, US Geological
802 Survey, 467 pages, 2002.
- 803 20. Ho, M., Parthasarathy, V., Etienne, E., Russo, T., Devineni, N., & Lall, U. America's
804 water: Agricultural water demands and the response of groundwater. *GEOPHYS RES*
805 *LETT*, V. 43, pp. 7546–7555. <http://dx.doi.org/10.1002/2016GL069797>, 2016.
- 806 21. Ihara, C., Kushnir, Y., Cane, M.A., & de la Peña, V.H. Indian summer monsoon rainfall
807 and its link with ENSO and Indian Ocean climate indices. *INT J CLIMATOL*, V. 27, pp.
808 179-187, <http://dx.doi.org/10.1002/joc.1394>, 2007.
- 809 22. Kar, S., Acharya, N., Mohanty, U.C. & Kulkarni, M.A. Skill of Monthly Rainfall
810 Forecasts Over India Using Multi-Model Ensemble Schemes. *INT J CLIMATOL*, V. 32,
811 pp. 1271 – 1286, <http://dx.doi.org/10.1002/joc.2334>, 2012.
- 812 23. KNMI Climate Explorer, <https://climexp.knmi.nl>, 1/1/2014
- 813 24. Kumar, K.K., Sonam, M.K. & Kumar, R.K. Seasonal Forecasting of Indian Summer
814 Monsoon Rainfall: A Review. *WEATHER*, V. 50, pp. 449 – 467,
815 <http://dx.doi.org/10.1002/j.1477-8696.1995.tb06071.x>, 1995.
- 816 25. Lall, U. & Sharma, A. A Nearest Neighbor Bootstrap for Resampling Hydrologic Time
817 Series. *WATER RESOUR RES*, V. 32, pp. 679 – 693, [http://dx.doi.org/](http://dx.doi.org/10.1029/95WR02966)
818 [10.1029/95WR02966](http://dx.doi.org/10.1029/95WR02966), 1996.
- 819 26. Li, S. & Goddard, L. Retrospective Forecasts with the ECHAM4.5 AGCM IRI Tech.
820 Report 05 – 02 December 2005, 2005.

- 821 27. Liu, X. & Pan, Y. Agricultural Drought Monitoring: Progress, Challenges, and
822 Prospects. *J GEOGR SCI*, V. 26, pp. 750 – 767, [http://dx.doi.org/10.1007/s11442-016-](http://dx.doi.org/10.1007/s11442-016-1297-9)
823 [1297-9](http://dx.doi.org/10.1007/s11442-016-1297-9), 2016.
- 824 28. Mahalanobis, P.C. On the Generalized Distance in Statistics. *Proceedings of the*
825 *National Institute of Sciences of India*, V. 2, pp. 49 – 55, 1936.
- 826 29. McKee, T.B., Doesken, N.J. & Kleist, J. The Relationship of Drought Frequency and
827 Duration to Time Scales. Eighth Conference on Applied Climatology, Anaheim,
828 California, 17 – 22 January 1993, 1993.
- 829 30. Mishra, A.K. & Desai, V.R. Drought Forecasting Using Stochastic Models. *STOCH*
830 *ENV RES RISK A*, V. 19, pp. 326 – 339, <http://dx.doi.org/10.1007/s00477-005-0238-4>,
831 2005.
- 832 31. Mishra, A.K. & Desai, V.R. Drought Forecasting Using Feed-Forward Recursive Neural
833 Network. *ECOL MODEL*, V. 198, pp. 127 – 138,
834 <http://dx.doi.org/10.1016/j.ecolmodel.2006.04.017>, 2006.
- 835 32. Mishra, A.K. & Singh, V.P. A Review of Drought Concepts. *J HYDROL*, V. 391, pp.
836 202 – 216, <http://dx.doi.org/10.1016/j.hydrol.2010.07.012>, 2010.
- 837 33. Murphy, A.H. On the “ranked probability score”. *J APPL METEOROL*, V. 8, pp. 988 –
838 989, [https://doi.org/10.1175/1520-0450\(1969\)008<0988:OTPS>2.0.CO%3B2](https://doi.org/10.1175/1520-0450(1969)008<0988:OTPS>2.0.CO%3B2), 1969.
- 839 34. Murphy, A.H. A Note on the Ranked Probability Score. *J APPL METEOROL*, V. 10,
840 pp. 155 – 156, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0450(1971)010<0155:ANOTRP>2.0.CO%3B2)
841 [0450\(1971\)010<0155:ANOTRP>2.0.CO%3B2](https://doi.org/10.1175/1520-0450(1971)010<0155:ANOTRP>2.0.CO%3B2), 1971.
- 842 35. Ngo-Duc, T., Polcher, J. & Laval, K. A 53-year Forcing Data Set for Land Surface
843 Models. *J GEOPHYS RES*, V. 110, 13 pp., <http://dx.doi.org/10.1029/2004JD005434>,
844 2005.
- 845 36. Nikam, A.V., Shendage, P.N., Jadhav, K.L. & Deokate, T.B. Economics of Production
846 of *Kharif* Potato in Satara, India. *International Journal of Agricultural Science*, V. 4, pp.
847 274 – 279, 2008.
- 848 37. Palmer, W.C. Meteorological Drought. Research Paper No. 45, U.S. Department of
849 Commerce, Washington, D.C., 65 pp., 1965.
- 850 38. Parthasarathy, B. & Pant, G.B. Seasonal Relationships Between Indian Summer
851 Monsoon Rainfall and the Southern Oscillation. *J CLIMATOL*, V. 5, pp. 369 – 378,
852 [http://dx.doi.org/551.513.7:551.553.11:551.577.32\(540\)](http://dx.doi.org/551.513.7:551.553.11:551.577.32(540)), 1985.
- 853 39. Parthasarathy, B., Diaz, H.F. & Escheid, J.K. Prediction of All-India Summer Monsoon
854 Rainfall with Regional and Large-Scale Parameters. *J GEOPHYS RES*, V. 93, pp. 5341 –
855 5350, <http://dx.doi.org/10.1029/JD093iD05p05341>, 1988.
- 856 40. R Core Team (2018). R: A language and environment for statistical computing. R
857 Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org/>.
- 858 41. Rajagopalan, B. & Lall, U. A k-nearest neighbor simulator for daily precipitation and
859 other weather variables. *WATER RESOUR RES*, V. 35, pp. 3089 – 3101,
860 [http://dx.doi.org/1999WR9000280043-1397/99/1999WR900028\\$09.00](http://dx.doi.org/1999WR9000280043-1397/99/1999WR900028$09.00), 1999.
- 861 42. Rajeevan, M., Bhate, J., Kale, J.D. & Lal, B. High Resolution Daily Gridded Rainfall
862 Data for the Indian Region: Analysis of Break and Active Monsoon Spells. *CURR SCI*
863 *INDIA*, V. 91, pp. 296 – 306, 2006.
- 864 43. Rasmusson, E.M. & Carpenter, T.H. The Relationship Between Eastern Equatorial
865 Pacific Sea Surface Temperature and Rainfall Over India and Sri Lanka. *MON*

- 866 WEATHER REV, V. 111, pp. 517 – 528, <http://dx.doi.org/10.1175/1520->
867 [0493\(1983\)111%3C0517:TRBEEP%3E2.0.CO;2](http://dx.doi.org/10.1175/1520-0493(1983)111%3C0517:TRBEEP%3E2.0.CO;2), 1983.
- 868 44. Rockstrom, J., Karlberg, L., Wani, S.P., Barron, J., Hatibu, N., Oweis, T., Bruggeman,
869 A., Farahani, J. & Qiang, Z. Managing Water in Rainfed Agriculture – The Need for a
870 Paradigm Shift, AGR WATER MANAGE, V. 97, pp. 543 – 550,
871 <http://dx.doi.org/10.1016/j.agwat.2009.09.009>, 2009.
- 872 45. Roeckner, E. and Coauthors. The atmospheric general circulation model ECHAM5:
873 Model description and simulation of present-day climate. Max-Planck-Institut für
874 Meteorologie Rep. 218, Hamburg, Germany, 90, 1996.
- 875 46. Serrano-Vicente, S.M., Beguería, S. & López-Moreno, J.I. A Multiscalar Drought Index
876 Sensitive to Global Warming: The Standardized Precipitation Evapotranspiration Index. J
877 CLIMATE, V. 23, pp. 1696 – 1718, <http://dx.doi.org/10.1175/2009JCLI2909.1>, 2010.
- 878 47. Shah, R.D. & Mishra, V. Utility of Global Ensemble Forecast System (GEFS)
879 Reforecast for Medium-Range Drought Prediction in India. J HYDROMETEOROL, V.
880 17, pp. 1781 – 1800, <http://dx.doi.org/10.1175/JHM-D-15-0050.1>, 2016.
- 881 48. Shah, R.D., Sahai, A.K. & Mishra, V. Short to Sub-Seasonal Hydrologic Forecast to
882 Manage Water and Agricultural Resources in India. Hydrol. Earth Syst. Sci., V. 21, pp.
883 707 – 720, <http://dx.doi.org/10.5194/hess-21-707-2017>, 2017.
- 884 49. Shukla, J. & Paolino, D.A. The Southern Oscillation and Long-Range Forecasting of the
885 Summer Monsoon Rainfall over India. MON WEATHER REV, V. 111, pp. 1830 – 1837,
886 [http://dx.doi.org/10.1175/1520-0493\(1983\)111%3C1830:TSOALR%3E1.0.CO;2](http://dx.doi.org/10.1175/1520-0493(1983)111%3C1830:TSOALR%3E1.0.CO;2), 1983.
- 887 50. Skees, J.R. Innovations in Index Insurance for the Poor in Lower Income Countries.
888 Agriculture and Resource Economics Review, V. 37, pp. 1 – 15, [http://doi.org/](http://doi.org/10.1017/S1068280500002094)
889 [10.1017/S1068280500002094](http://doi.org/10.1017/S1068280500002094), 2016.
- 890 51. Souza, F.A. & Lall, U. Seasonal to Interannual Ensemble Streamflow Forecasts for
891 Ceara, Brazil: Applications of Multivariate, Semiparametric Algorithm. WATER
892 RESOUR RES, V. 39, 13 pp., <http://dx.doi.org/10.1029/2002WR001373>, 2003.
- 893 52. Thapliyal, V. Prediction of Indian Monsoon Variability Evaluation and Prospects
894 Including Development of a New Model. China Ocean Press, pp. 397 – 416, 1987.
- 895 53. Irrigation Water Use, <https://water.usgs.gov/edu/wuir.html>, accessed 3/14/2018, 2017.
- 896 54. van den Dool, H.M. Empirical Methods in Short-Term Climate Prediction, Oxford
897 University Press, 215 pp., 2007.
- 898 55. Walker, G.T. Correlations in seasonal variations of weather, IX: A further study of world
899 weather (World Weather II), Memoirs of India Meteorological Department, V. 24, pp.
900 275 – 332, 1924.

901
902
903
904
905
906
907
908
909
910
911

912

Tables913 **Table 1**

Year	Probability of Above Mean	Probability of Below Mean	Observed CDI Anomaly (%)	Boxplot IQR (vertical axis units of %-anomalies)
2001	0.59	0.41	+14.4	10.9
2002	0.42	0.58	+15.5	21.0
2003	0.20	0.80	+37.8	23.1
2004	0.35	0.65	-20.1	7.70
2005	0.25	0.75	-51.3	12.1
2006	0.37	0.63	-47.9	10.0
2007	0.37	0.63	-20.5	2.60
2008	0.75	0.25	-6.33	19.1
2009	0.64	0.36	-30.0	5.10
2010	0.18	0.82	-56.4	31.1
2011	0.58	0.42	+2.72	0.19
2012	0.68	0.32	+25.4	9.90
2013	0.18	0.82	-9.36	24.6

914

915

916

917

918

919

920

921

922 **Table 2**

Year	Forecast	Actual Observation	Result
2001	AM (59%)	AM	Hit
2002	BM (58%)	AM	Miss
2003	BM (80%)	AM	Miss
2004	BM (65%)	BM	Hit
2005	BM (75%)	BM	Hit
2006	BM (63%)	BM	Hit
2007	BM (63%)	BM	Hit
2008	AM (75%)	BM	False Alarm
2009	AM (64%)	BM	False Alarm
2010	BM (82%)	BM	Hit
2011	AM (58%)	AM	Hit
2012	AM (68%)	AM	Hit
2013	BM (82%)	BM	Hit

923

924

925

926

927

928

929

930

931

932

933

934 **Table 3**

Year	CDI Forecast Results	IMD Precipitation Forecast	Actual Precipitation	IMD Forecast Results
2001	Hit	96% of LPA	93% of LPA	Hit
2002	Miss	Not Available	68% of LPA	NA
2003	Miss	99% of LPA	40% of LPA	Miss
2004	Hit	103% of LPA	160% of LPA	False Alarm
2005	Hit	Not Available	160% of LPA	NA
2006	Hit	90% of LPA	141% of LPA	False Alarm
2007	Hit	96% of LPA	163% of LPA	False Alarm
2008	False Alarm	Not Available	95% of LPA	NA
2009	False Alarm	Not Available	212% of LPA	NA
2010	Hit	99% of LPA	199% of LPA	False Alarm
2011	Hit	98% of LPA	85% of LPA	Miss
2012	Hit	96% of LPA	46% of LPA	Miss
2013	Hit	98% of LPA	150% of LPA	False Alarm

935

936

937

938

939

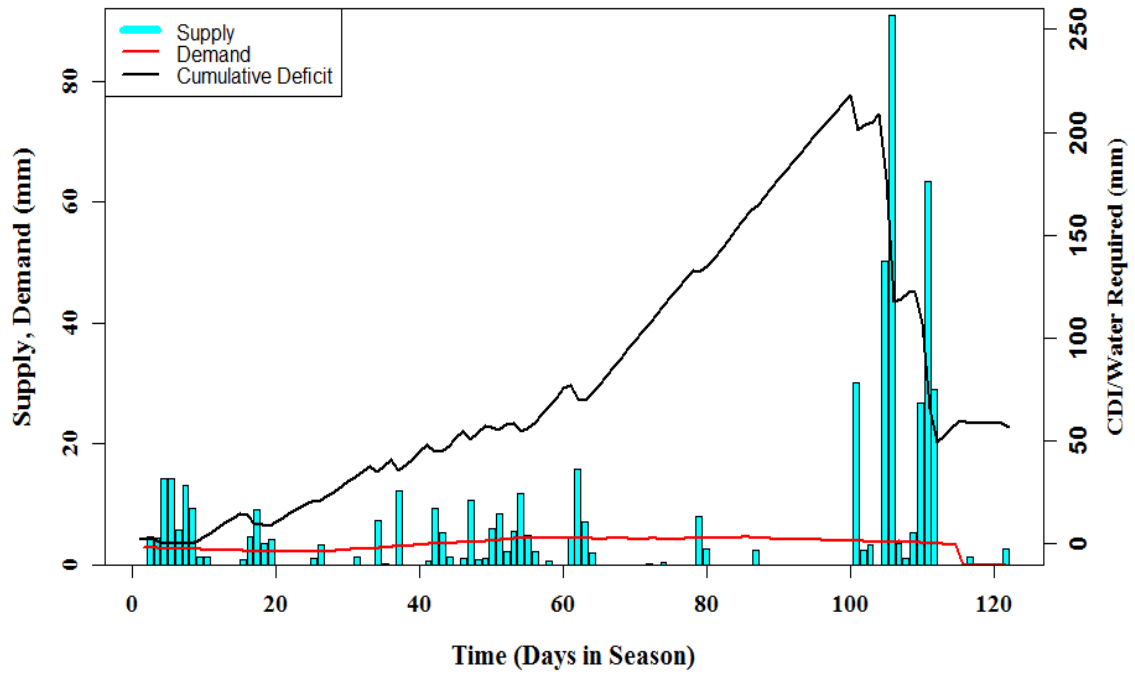
940

941

942

Figures

Seasonal Cumulative Deficit Index



944

945 **Figure 1**

946

947

948

949

950

951

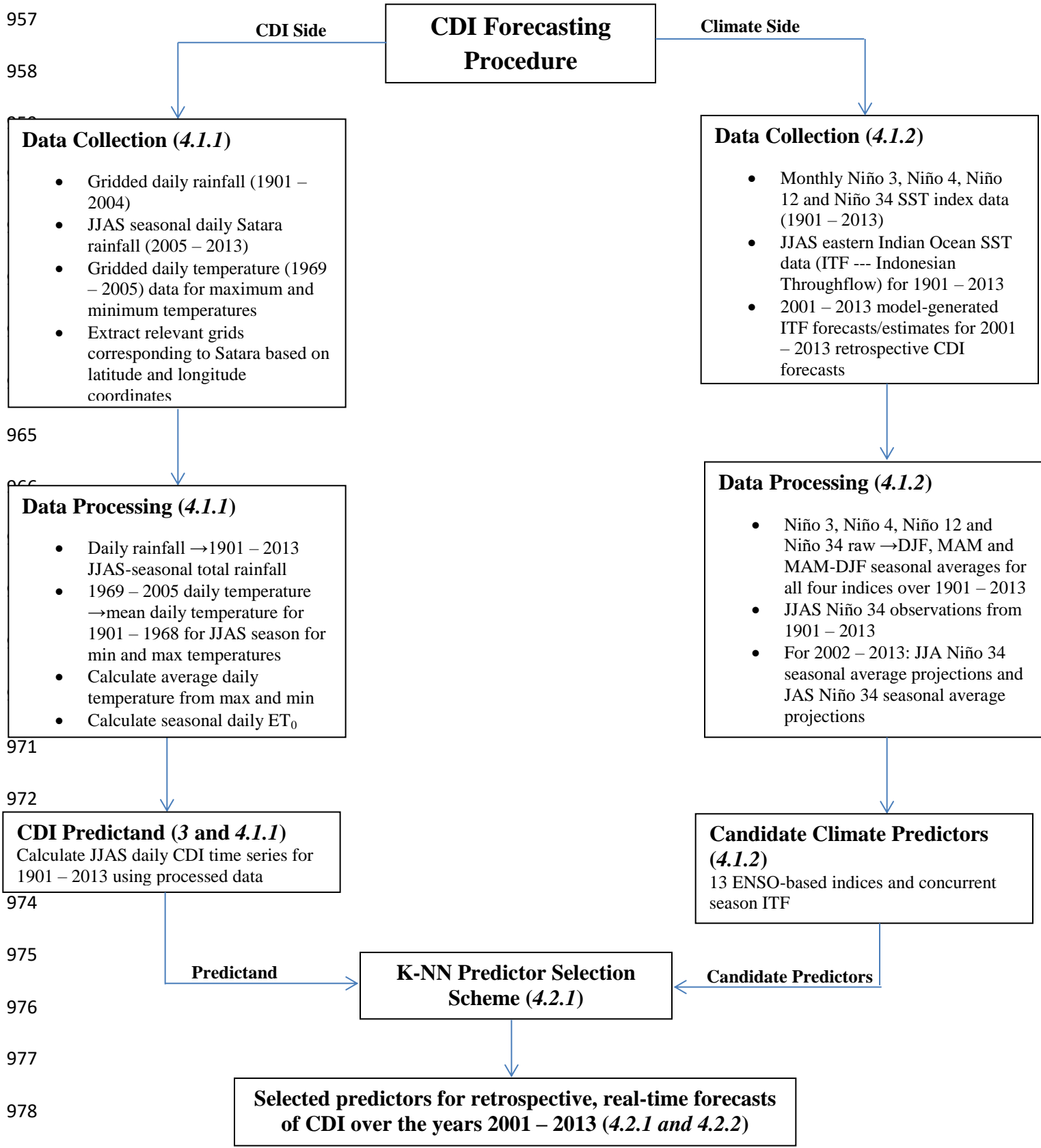
952

953

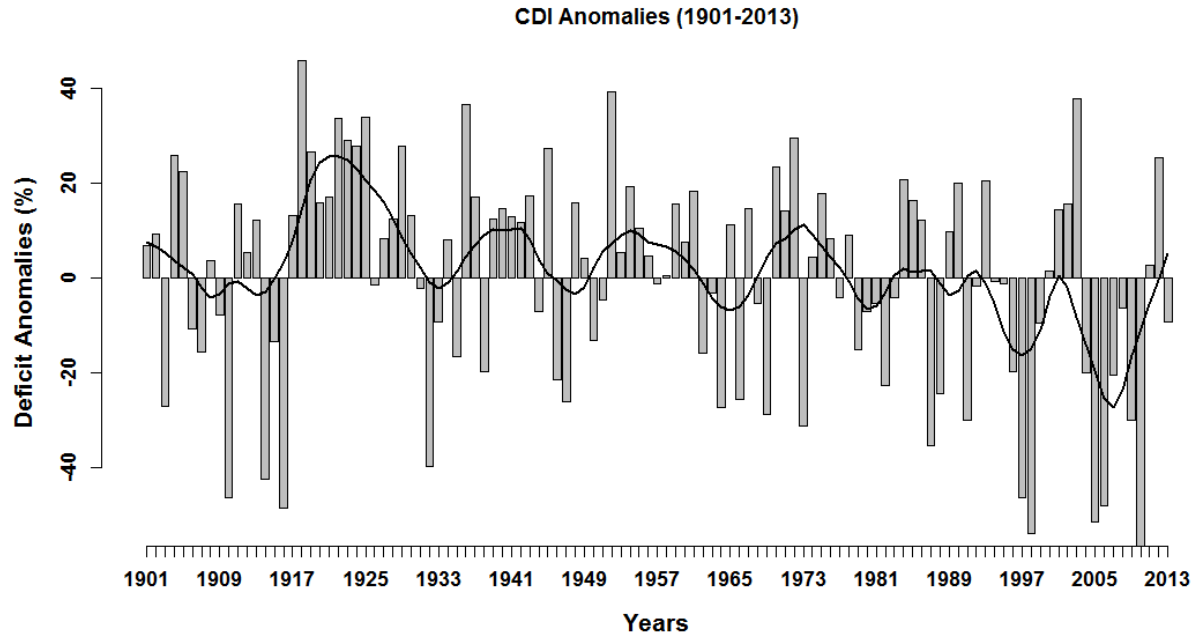
954

955

956



979 **Figure 2**



980

981 **Figure 3**

982

983

984

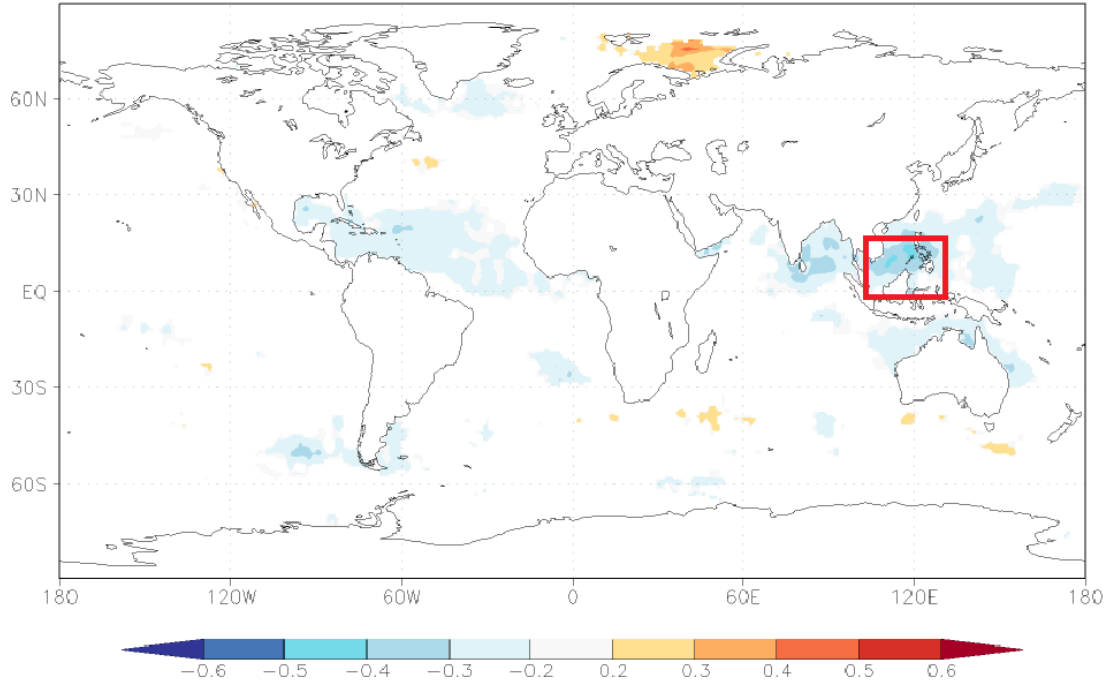
985

986

987

988

rank corr Jun-Sep averaged fcd ma1 index
with Jun-Sep averaged HadSST1 SST (detrend) 1901:2000 $p < 10\%$



989

990 **Figure 4**

991

992

993

994

995

996

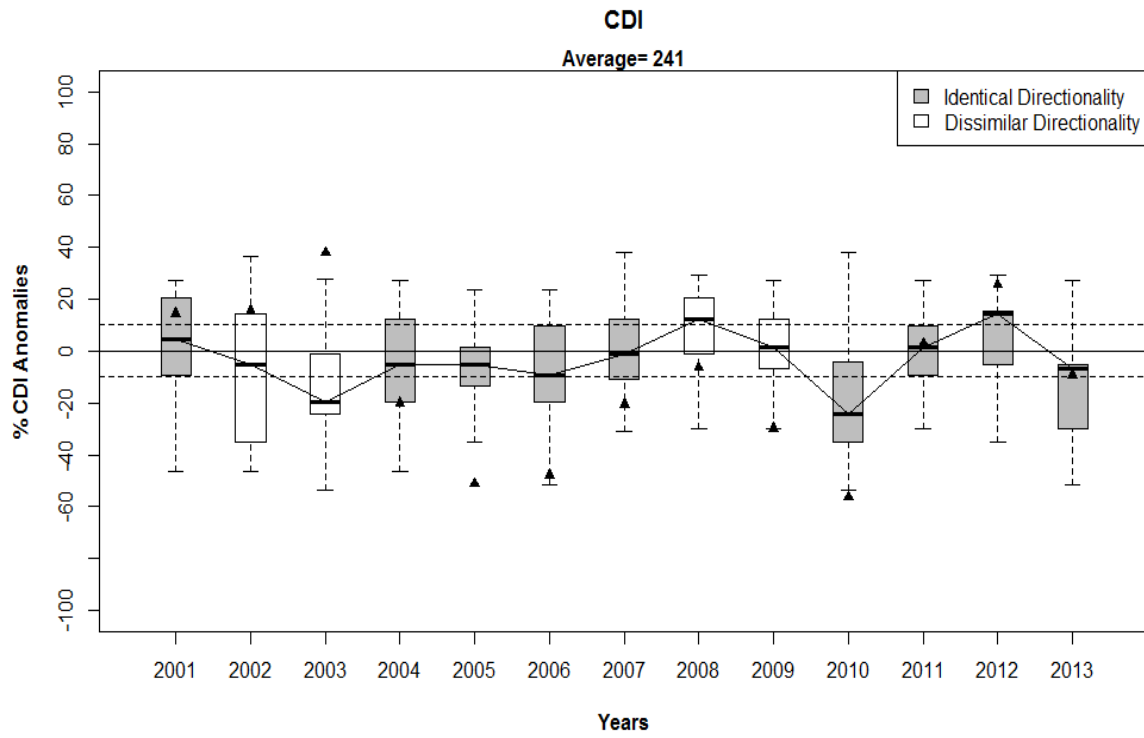
997

998

999

1000

1001



1002

1003 **Figure 5**

1004

1005

1006

1007

1008

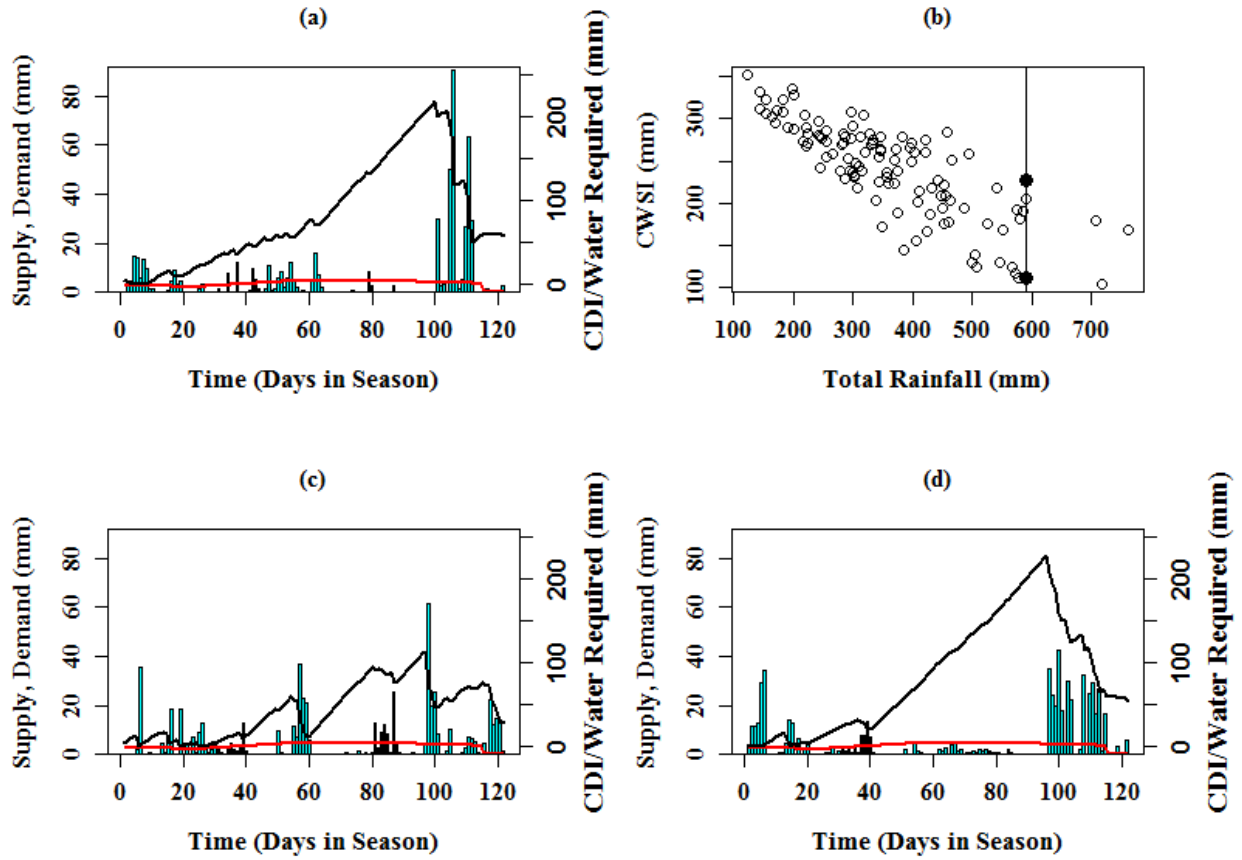
1009

1010

1011

1012

1013



1014

1015 **Figure 6**

1016

1017

1018

1019

1020

1021

1022

1023

1024

1025

1026

1027

Figure and Table Captions

1028

1029 **Table 1:** The table below shows important statistics calculated from kNN forecasts of CDI. In
1030 particular, column 2 displays the probabilities of the CDI for a particular season being above the
1031 CDI climatology. These probabilities are calculated from the kNN sampling distribution, which
1032 in turn is simulated from historical values of the CDI based on the nearest neighbors determined
1033 in the predictor variable space. Column 3 shows the complementary probabilities of being below
1034 this historical average. The forecasts for years 2001-2013 are retrospective and may serve as a
1035 cross-validation for the kNN model. Column 4 shows the values of the actual (observed) CDI
1036 anomalies with respect to the 1901-2013 climatology as percentages. A negative value implies
1037 that the actual CDI value was below the historical average by the given percentage. The rounded
1038 IQR values are shown in the final column of the table.

1039

1040 **Table 2:** The results of the kNN-generated CDI forecasts, including the most likely category
1041 (AM = Above Mean, BM = Below Mean) along with the corresponding kNN-assigned
1042 probability value expressed as a percentage in parentheses next to it (column 2), the category in
1043 which the observed anomaly value resides (column 3), and the hit/miss/false alarm designations
1044 corresponding to these results (column 4).

1045

1046 **Table 3:** A comparison of the CDI forecasts and the JJAS total seasonal precipitation forecasts
1047 generated by the India Meteorological Department (IMD). Column 2 is a repeat of column 4 in
1048 Table 2; a record of the accuracy of CDI forecasts expressed in terms of hits and misses.
1049 Column 3 contains the forecasts issued by IMD, and column 4 is the actual observations of JJAS
1050 seasonal total rainfall using rainfall data from the Satara district itself. The fifth and final
1051 column of Table 3 shows the accuracy of the IMD forecasts in terms of hits and misses using
1052 their own 5-category system.

1053 **Figure 1:** A plot of the cumulative deficit index (CDI) for the JJAS season in a randomly
1054 selected year in our data set. The plot depicts the change in CDI as rainfall distribution and crop
1055 water requirement varies over the given monsoon season. The vertical cyan bars are the daily
1056 rainfall magnitudes, the slowly-changing red line is the crop water requirement (demand) and the
1057 black time series is the CDI itself. Notice how CDI increases as rainfall is either low in
1058 magnitude or sparsely distributed in certain blocks of time in the season.

1059 **Figure 2:** Flowchart depicting the entire forecasting procedure for potato-based CDI in Satara,
1060 Maharashtra, India. The steps are categorized as data collection, data processing,
1061 predictand/predictor calculation, all of which converge to predictor selection and forecast
1062 modeling. The section number of the paper in which these steps are covered is written in italics
1063 next to the category. A brief summary of each step is given, one for the steps used in CDI
1064 calculation and one for the steps used in processing the candidate predictors from climate.

1065 **Figure 3:** Bar plot showing the CDI percent deficit anomalies for each of the years/growing
1066 seasons under consideration (1901 – 2013). The black, smooth time series is produced by an 11-

1067 year LOWESS smoothing of the CDI percent deficit anomalies and is meant to show the critical
1068 trends in the CDI over the entire 1901 – 2013 period.

1069

1070 **Figure 4:** Spearman rank correlation between CDI in Satara and SST field during the same JJAS
1071 season. SST region in the Indian Ocean (red box) that influences the CDI has a statistically
1072 significant correlation at the 95% significance level.

1073

1074 **Figure 5:** Boxplot diagrams depicting the kNN forecast distributions for CDI over the years
1075 2001 – 2013 for potatoes grown in the Satara district, Maharashtra, India. Longer, more
1076 stretched out boxes indicate a greater degree of variability, or uncertainty, in the forecast
1077 distribution. Boxes in which the median is grossly off-center indicates that the forecast
1078 distribution is heavily skewed. Anomalies with respect to the climatology of the predictand were
1079 used in the boxplot calculations. As the results are presented in terms of the percent anomalies,
1080 the historical average is located at zero. The triangles represent the observations as percent
1081 anomalies about the mean. Boxes that have been shaded in gray indicate years during which
1082 identical directionality was observed, whereas boxes that are white indicate years during which
1083 dissimilar directionality was observed.

1084

1085 **Figure 6:** The four panels pictured here depict the CDI in various ways. In panels (a), (c) and
1086 (d), the blue bars represent daily seasonal rainfall levels (in mm), the red curve represents crop
1087 evaporative water demand (ET_0) and the black time series is the CDI calculated based on this
1088 data. Panel (a) illustrates the basic nature of CDI using the daily seasonal CDI time series from
1089 the JJAS growing season of 2013. Note that this time series is specifically calculated for
1090 potatoes grown in the Satara district of Maharashtra, India during the 2013 JJAS growing season.
1091 Panel (b) shows a scatterplot of total rainfall across all growing seasons (1901 – 2013) and CDI
1092 across all growing seasons. A significant negative correlation between them is apparent from
1093 this scatterplot (Pearson correlation is -0.8, Spearman rank correlation is -0.812, Kendall rank
1094 correlation is -0.623). This panel demonstrates two different growing seasons, with two different
1095 CDI values, during which the total seasonal rainfall was the same. Panel (c) is a seasonal CDI
1096 time series plot corresponding to the growing season with the lower CDI value on the vertical
1097 line in panel (b). Panel (d) is a seasonal CDI time series plot corresponding to the growing
1098 season with the higher CDI value on the vertical line in panel (b).