



1

2

3

4 **Discrete k-nearest neighbor resampling for simulating multisite**

5 **precipitation occurrence and adaption to climate change**

6 : Discrete KNNR for Multisite Occurrence (DKMO version1.0) - model development

7

8 Keywords: daily precipitation, discrete, k-nearest neighbor, Markov chain, multisite, occurrence

9

10

11 Taesam Lee¹ and Vijay P. Singh²

12 ¹ Department of Civil Engineering, ERI, Gyeongsang National University,

13 501 Jinju-daero, Jinju, Gyeongnam, South Korea, 660-701

14 ² Department of Biological and Agricultural Engineering & Zachry Department of
15 Civil Engineering, Texas A&M University, 321 Scoates Hall, College Station, Texas,
16 United States, 77843

17

18

19

20

21 Corresponding Author :

22

23 Taesam Lee, Ph.D.

24 Gyeongsang National University, Dept. of Civil Engineering

25 Tel)+82-55-772-1797, Fax)+82-55-772-1799

26 Email) tae3lee@gnu.ac.kr



27

28

Abstract

29 Stochastic weather simulation models are commonly employed in water resources management
30 and agricultural applications. The data simulated by these models, such as precipitation,
31 temperature, and wind, are used as input for hydrological and agricultural models. Stochastic
32 simulation of multisite precipitation occurrence is a challenge because of its intermittent
33 characteristics as well as spatial and temporal cross-correlation. Employing a nonparametric
34 technique, k-nearest neighbor resampling (KNNR), and coupling it with Genetic Algorithm (GA),
35 this study proposes a novel simulation method for multisite precipitation occurrence. The proposed
36 discrete version of KNNR (DKNNR) model is compared with an existing parametric model, called
37 multisite occurrence model with standard normal variate (MONR). The datasets simulated from
38 both the DKNNR model and the MONR model are tested using a number of statistics, such as
39 occurrence and transition probabilities as well as temporal and spatial cross-correlations. Results
40 show that the proposed DKNNR model can be a good alternative for simulating multisite
41 precipitation occurrence. We also tested the model capability to adapt climate change. It is shown
42 that the model is capable but further improvement is required to have specific variations of the
43 occurrence probability due to climate change. Combining with the generated occurrence, the
44 multisite precipitation amount can then be simulated by any multisite amount model.

45



46 **1. Introduction**

47 Stochastic simulation of weather variables has been employed for water resources
48 management, hydrological design, agricultural applications, filling in missing historical data,
49 extending observed records, simulating data, and simulating different weather conditions.
50 Stochastic simulation models play a key role in producing weather sequences, while preserving
51 the statistical characteristics of observed data. A number of stochastic weather simulation models
52 have been developed using parametric and nonparametric approaches (Lee, 2017; Lee et al., 2012;
53 Wilby et al., 2003; Wilks, 1999; Wilks and Wilby, 1999).

54 Parametric approaches summarize the statistical characteristics of observed weather data
55 with a parameter set (Jeong et al., 2012; Lee, 2016; Zheng and Katz, 2008). The parameters fitted
56 with observed weather data are employed in simulation. In nonparametric approaches, historical
57 analogs with current conditions are searched following the weather simulation data (Buishand and
58 Brandsma, 2001; Lee et al., 2012). Furthermore, combinations of parametric and nonparametric
59 models have also been proposed (Apipattanavis et al., 2007; Frost et al., 2011).

60 Among weather variables, the precipitation variable possesses intermittency and zero values
61 between precipitation events, and to properly reproduce them is difficult and remains a challenge
62 (Beersma and Buishand, 2003; Hughes et al., 1999; Katz and Zheng, 1999). Due to this difficulty,
63 precipitation is simulated separately from other variables. The main method for reproducing
64 intermittency has been the multiplication of precipitation occurrence and an amount as $Z=X \cdot Y$,
65 where X is the occurrence (binary as either 0 or 1) and Y is the amount (Jeong et al., 2013; Lee and
66 Park, 2017; Todorovic and Woolhiser, 1975). The spatial and temporal dependence in the
67 occurrence and amount of precipitation introduces further complexity multisite simulation.



68 Wilks (1998) presented a multisite simulation model for the occurrence process (i.e. X) using
69 the standard normal variable that is spatially dependent, representing the relation between the
70 occurrence variable and the standard normal variable with simulation data. Even though the
71 multisite occurrence data simulated by this model preserves various statistical characteristics of
72 the observed data well, some drawbacks still exist, such as underestimation of lagged cross-
73 correlation. Furthermore, the relation between standard normal variable and occurrence variable
74 relies on long stochastic simulation.

75 Lall and Sharma (1996) proposed a nonparametric simulation model, called k-nearest
76 neighbor resampling (KNNR). The model has been updated to simulate multivariate hydro-
77 meteorological variables (Brandsma and Buishand, 1998; Mehrotra et al., 2006; St-Hilaire et al.,
78 2012). One of the major drawbacks of this multivariate KNNR model is that the simulated data
79 cannot produce patterns different from those of the observed data. Lee et al. (2010a) overcame this
80 shortcoming by mixing the simulated dataset with Genetic Algorithm (GA) that led to the
81 reproduction of similar populations. A number of variants of KNNR-GA have since been applied
82 (Lee et al., 2012; Lee and Park, 2017).

83 Therefore, in the current study we propose a novel simulation method for multisite
84 occurrence of the precipitation variable with a nonparametric approach. The proposed
85 nonparametric model is compared with the existing multisite model (Wilks, 1998). The paper is
86 organized as follows. The next section presents a mathematical background of existing multisite
87 occurrence modeling. The modeling procedure is discussed in section 3. The study area and data
88 are reported in section 4. The model is applied in section 5. Results of the proposed model are
89 discussed in section 6, and summary and conclusions are presented in section 7.



90 2. Background

91 2.1. Single site occurrence modeling

92 Let X_t^s represent the occurrence of daily precipitation for a location s ($s=1, \dots, S$) on day t
 93 ($t=1, \dots, n$; n is the number observed days) and let X_t^s be either zero for dry day or one for wet day.
 94 The first order Markov chain model for X_t^s is defined with the assumption that the occurrence
 95 probability of a wet day is fully defined by the previous day as

$$96 \quad \Pr\{X_t^s = 1 \mid X_{t-1}^s = 0\} = p_{01}^s \quad (1)$$

$$97 \quad \Pr\{X_t^s = 1 \mid X_{t-1}^s = 1\} = p_{11}^s \quad (2)$$

98 Also $p_{00}^s = 1 - p_{01}^s$ and $p_{10}^s = 1 - p_{11}^s$, since the summation of zero and one should be unity
 99 with the same previous condition. This consists of a transition probability matrix (TPM) as

$$100 \quad TPM^s = \begin{bmatrix} p_{00}^s & p_{01}^s \\ p_{10}^s & p_{11}^s \end{bmatrix} = \begin{bmatrix} 1 - p_{01}^s & p_{01}^s \\ 1 - p_{11}^s & p_{11}^s \end{bmatrix} \quad (3)$$

101 The marginal distributions of TPM (i.e. p_0 and p_1) can be expressed with TPM and its condition of
 102 $p_0 + p_1 = 1$ as:

$$103 \quad p_0^s = \frac{p_{01}^s}{1 + p_{01}^s - p_{11}^s} \quad (4)$$

$$104 \quad p_1^s = \frac{1 - p_{11}^s}{1 + p_{01}^s - p_{11}^s} \quad (5)$$



105 Note that p_1 represents the probability of precipitation occurrence for a day, while p_0 does non-
106 occurrence. The lag-1 autocorrelation of precipitation occurrence is the combination of transition
107 probabilities as:

$$108 \quad \rho_1(s, s) = p_{11}^s - p_{01}^s \quad (6)$$

109 The simulation can be done by comparing TPM with a uniform random number (u_t^s) as

$$110 \quad X_t^s = \begin{cases} 1 & \text{if } u_t^s \leq p_{i1}^s \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

111 where p_{i1}^s is the selected probability from TPM regarding the previous condition i (i.e. either 0 or
112 1). Wilks (1998) suggested a different method using a standard normal random number $w_t^s \sim \mathcal{N}[0,1]$
113 as

$$114 \quad X_t^s = \begin{cases} 1 & \text{if } w_t^s \leq \Phi^{-1}(p_{i1}^s) \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

115 where Φ^{-1} indicates the inverse of the standard normal cumulative function Φ .

116 **2.2. Multisite occurrence modeling**

117 Wilks (1998) suggested a multisite occurrence model using a standard normal random
118 number (here, denoted as MONR) that is spatially dependent but serially independent. The
119 correlation of the standard normal variate for a site pair of q and s can be expressed as:

$$120 \quad \tau(q, s) = \text{corr}[w_t^q, w_t^s] \quad (9)$$

121 Also, the correlation of the original occurrence variate is



122
$$\rho(q, s) = \text{corr}[X_t^q, X_t^s] \quad (10)$$

123 Once the correlation of the standard normal variate is known, the simulation of multisite
124 precipitation occurrence is straightforward. Multivariate standard normal distribution is used with
125 the parameter set of $[\mathbf{0}, \mathbf{T}]$ where $\mathbf{0}$ is the zero vector ($S \times 1$) and \mathbf{T} is the correlation matrix with the
126 elements of $\tau(q, s)$ for $q \in \{1, \dots, S\}$ and $s \in \{1, \dots, S\}$.

127 Since direct estimation of $\tau(q, s)$ is not applicable, a simulation technique is used to estimate
128 $\tau(q, s)$ from $\rho(q, s)$. A long sequence of the occurrence process is simulated with different values
129 of $\tau(q, s)$ and its corresponding correlation of the original domain $\rho(q, s)$ is estimated with the
130 simulated long sequence by the inverse standard normal cumulative function (i.e. Φ^{-1}). A curve
131 between $\tau(q, s)$ and $\rho(q, s)$ is derived from this long simulation with the MONR model and is
132 employed for the parameter estimation for real application.

133 **3. DKNNR**

134 **3.1. DKNNR modeling procedure**

135 In the current study, a novel multisite simulation model for discrete occurrence of precipitation
136 variable with k-nearest neighbor resampling (KNNR) technique (Lall and Sharma, 1996; Lee
137 and Ouarda, 2011; Lee et al., 2017) for discrete case (denoted as Discrete KNNR; DKNNR)
138 is proposed by combining a mixture mechanism with Genetic Algorithm (GA).

139 Provided the number of nearest neighbors, k , is known, the discrete k-nearest neighbor
140 resampling with genetic algorithm is done as follows:



141 (1) Estimate the distance between the current (i.e. time index: c) multisite occurrence
142 X_c^s and the observed multisite occurrence x_i^s . Here, the distance is measured for
143 $i=1, \dots, n-1$ as

$$144 \quad D_i = \sum_{s=1}^S |X_c^s - x_i^s| \quad (11)$$

145 (2) Arrange the estimated distances from step (1) in ascending order, select the first k
146 distances (i.e., the smallest k values), and reserve the time indices of the smallest k
147 distances.

148 (3) Randomly select one of the stored k time indices with the weighting probability
149 given by

$$150 \quad w_m = \frac{1/m}{\sum_{j=1}^k 1/j}, \quad m = 1, \dots, k \quad (12)$$

151 (4) Assume the selected time index from step (3) as p . Note that there are a number of
152 values that have the same distance as the selected D_p , since D_p is a natural number
153 between 0 and S . A random selection procedure is required to take into account the
154 cases with the same quantity. One particular time index is randomly selected with
155 the equal probabilities among the time indices of the same distances.

156 (5) Assign the binary vector of the proceeding index of the selected time as
157 $\mathbf{x}_{p+1} = [x_{p+1}^s]_{s \in \{1, S\}}$. Here, p is the finally selected time index from step (4).

158 (6) Execute the following steps for GA mixing if GA mixing is selected. Otherwise, skip
159 this step.



160 (6-1) Reproduction: Select one additional time index using steps (1) through (4) and
 161 denote this index as p^* . Obtain the corresponding precipitation occurrence
 162 values, $\mathbf{x}_{p^*+1} = [x_{p^*+1}^s]_{s \in \{1, \dots, S\}}$. The subsequent two GA operators employ the two
 163 selected vectors, \mathbf{x}_{p+1} and \mathbf{x}_{p^*+1} .

164 (6-2) Crossover: Replace each element x_{p+1}^s with $x_{p^*+1}^s$ at probability P_{cr} , i.e.,

$$165 \quad X_{c+1}^s = \begin{cases} x_{p^*+1}^s & \text{if } \varepsilon < P_{cr} \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (13)$$

166 where ε is a uniform random number between 0 and 1.

167 (6-3) Mutation: Replace each element (i.e., each station, $s=1, \dots, S$) with one selected
 168 from all the observations of this element for $i=1, \dots, n$ with probability P_m , i.e.,

$$169 \quad X_{c+1}^s = \begin{cases} x_{\xi+1}^s & \text{if } \varepsilon < P_m \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (14)$$

170 where $x_{\xi+1}^s$ is selected from $[x_i^s]_{i \in \{1, \dots, n\}}$ with equal probability for $i=1, \dots, n$ and
 171 ε is a uniform random number between 0 and 1.

172 (7) Repeat steps (1)-(6) until the required data are generated.

173 The selection of the number of nearest neighbors (k) has been investigated by Lall and
 174 Sharma (1996) and Lee and Ouarda (2011). A simple selection method was applied in the current
 175 study as $k = \sqrt{n}$. For hydrometeorological stochastic simulations, this heuristic approach of k
 176 selection has been employed (Lall and Sharma, 1996; Lee and Ouarda, 2012; Lee et al., 2010b;
 177 Prairie et al., 2006; Rajagopalan and Lall, 1999). The roles of crossover probability P_{cr} and



178 mutation probability P_m were studied by Lee et al. (2010b). Lee et al. (2010b) showed that $P_{cr}=0.1$
 179 and $P_m=0.01$ can be a reasonable parameter set which does not critically affect the performance.
 180 Therefore, this parameter set was applied in the current study. In Appendix A, an example of the
 181 DKNNR simulation procedure is explained in detail.

182 **3.2. Adaptation to climate change**

183 The capability of model to take climate change into account is critical. For example, the
 184 marginal distributions and transition probabilities in Eqs. (5) and (3) can change in future climate
 185 scenarios. It is known that nonparametric simulation models have a difficulty to adapt to climate
 186 change, since the models employ in general the current observation sequences. However, the
 187 proposed model in the current study possesses the capability to adapt to the variations of
 188 probabilities by tuning the crossover and mutation probabilities in P_{cr} (13) and P_m (14), adding
 189 the condition when applied.

190 For example, the probability of P_{11} can be increased with the cross-over probability P_{cr} by
 191 adding the condition to increase the probability of P_{11} as:

$$192 \quad X_{c+1}^s = \begin{cases} x_{p^*+1}^s & \text{if } \varepsilon < P_{cr} \text{ \& } x_{p^*+1}^s = 1 \text{ \& } X_c^s = 1 \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (15)$$

193 It is obviously possible to increase the probability of P_1 by removing the condition of $X_c^s = 1$.

194 In addition, further adjustment can be made with the mutation process in Eq. (14) as

$$195 \quad X_{c+1}^s = \begin{cases} x_{\xi+1}^s & \text{if } \varepsilon < P_m \text{ and } x_{\xi+1}^s = 1 \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (16)$$



196 This adjustment of adding the condition $x_{\xi+1}^s = 1$ can increase the marginal distribution as much as
197 $P_m \times P_1$. This has been tested in the case study.

198 **4. Study area and data description**

199 For testing the occurrence model, 12 weather stations were selected from Yeongnam province
200 which is located in the southeastern part of South Korea, as shown in Figure 1. Information on
201 longitude and latitude (fourth and fifth columns) as well as order index and the identification
202 number (first and second columns) of these stations operated by Korea Meteorological
203 Administration with the area name (third column) is shown at Table 1.

204 Figure 1 illustrates the locations of the selected weather stations. All the stations are inside
205 Yeongnam province which consists of two different regions as north and south Gyeongsang as
206 well as the self-governing cities of Busan, Deagu, and Ulsan. Most of the Yeongnam region is
207 drained to Nakdong River. It is important to analyze the impact of weather conditions for planning
208 agricultural operations and water resources management especially during the summer season,
209 because around 50-60 percent of the annual precipitation occurs during the summer season from
210 June to September. The length of daily precipitation data record ranges from 1976 to 2008 and
211 the summer season record was employed since a large number of rainy days occurs during summer
212 and it is important to preserve these characteristics. Also, the whole year dataset was tested and
213 other seasons were further applied but the correlation coefficient was relatively high and its
214 correlation matrix estimated was not a positive semi-definite matrix for the MONR model.



215 **5. Application**

216 To analyze the performance of the proposed DKNNR model, the occurrence of precipitation
217 was simulated. The DKNNR simulation was compared with that of the MONR model. For each
218 model, 100 series of daily occurrence with the same record length were simulated. The key
219 statistics of observed data and each generated series, such as transition probabilities (P_{11} , P_{01} , and
220 P_1) and cross-correlation (see Eq.(10)), were determined. The MONR model underestimated the
221 lag-1 cross-correlation, as indicated by Wilks (1998). In the current study, this statistic was
222 analyzed, since a synoptic scale weather system could result in lagged cross-correlation (Wilks,
223 1998). It was formulated as

$$224 \quad \rho_1(q, s) = \text{corr}[X_{t-1}^q, X_t^s] \quad (17)$$

225 Statistics from 100 generated series were evaluated by the root mean square error (RMSE)
226 expressed as below:

$$227 \quad RMSE = \left(\frac{1}{N} \sum_{m=1}^N (\Gamma_m^G - \Gamma^h)^2 \right)^{1/2} \quad (18)$$

228 where N is the number of series (here 100), Γ_m^G is the statistic estimated from the m^{th} generated
229 series, while Γ^h is the statistic for the observed data. Note that lower RMSE indicates better
230 performance representing the summarized error of a given statistic of generated series from the
231 statistic of the observed data.

232 The 100 simulated statistic values were illustrated with boxplots to show their variability as
233 shown in Figure 2 - Figure 4. The box of boxplot represents the interquartile range (IQR) ranging
234 25 percentile to 75 percentile. The whiskers extend to up and down $1.5 \times \text{IQR}$. Data beyond the



235 whiskers ($1.5 \times \text{IQR}$) are indicated by a plus sign (+). The horizontal line inside the box represents
236 the median of the data. The statistics of the observed data are denoted by a cross (x). The closer a
237 cross is to the horizontal line inside the box, the better the simulated data from a model reproduces
238 the statistical characteristics of the observed data.

239 6. Results

240 6.1. Occurrence and transition probabilities

241 The data simulated from the proposed DKNNR model and the existing MONR model were
242 analyzed. The estimated transition probabilities (P_{11} and P_{01} in Eq. (3)) as well as the occurrence
243 probability (P_1 in Eq. (5)) are shown in Table 2 and Figure 2 - Figure 4 for the observed data and
244 the data generated from the DKNNR and MONR models. In Table 2, the observed statistic shows
245 that P_{11} is always higher than P_{01} and P_1 is between P_{11} and P_{01} . Site 6 shows the lowest P_{11} and
246 P_1 and site 12 shows the highest P_{11} .

247 As shown in Figure 2, the probability P_{11} of the observed data shows that sites 6, 7, 8, and 9
248 located in the northern part of the region exhibited lower consistency (i.e. consecutive rainy days)
249 than did the other sites, while sites 5 and 12 had higher probability of P_{11} than did other sites. Both
250 models preserved well the observed P_{11} statistic. It seems that the MONR model had a slight better
251 performance since this statistic is parameterized in the model as shown in the section 2.2. Note
252 that the MONR model employed the transition probabilities in simulating rainfall occurrence,
253 while DKNNR model did not. The occurrence probability P_1 can be described with the
254 combination of transition probabilities as in Eq. (5). Even though the transition probabilities were
255 not employed in simulating rainfall occurrence, the DKNNR model preserved this statistic fairly
256 well.



257 As shown in Figure 3, the P_{01} probability showed a slightly different behavior such that sites
258 1, 2, and 3 located in the middle part of the Yeongnam province showed a higher probability than
259 did other sites. A slight underestimation was seen for sites 2 and 11 but it was not critical, since its
260 observed value with a cross mark was close to the upper IQR representing 75 percentile.

261 The behavior of P_1 was found to be same as that of the P_{11} probability. It can be seen in
262 Figure 4 that no significant underestimation is seen for the DKNNR model (top panel). The P_1
263 statistic was fairly preserved by both DKNNR and MONR models. Note that the MONR model
264 parameterized the P_1 statistic through the transition probabilities as in Eq. (5), while DKNNR
265 model did not. Although the DKNNR model did not use any parameter for simulation, the P_1
266 statistic was preserved fairly well.

267 **6.2. Cross-correlation**

268 Cross-correlation is a measure of relationship between sites. Preservation of cross-
269 correlation is important for the simulation of precipitation occurrence and is required in the
270 regional analysis for water resources management or agricultural applications. Furthermore,
271 lagged cross-correlation is also essential as much as is cross-correlation (i.e. contemporaneous
272 correlation). For example, the amount of streamflow for a watershed from a certain precipitation
273 event is highly related with lagged cross-correlation. It is accepted that precipitation event is not
274 significantly correlated with more than one day. Therefore, only lag-1 cross-correlation was
275 analyzed in the current study.

276 The cross-correlation of observed data is shown in Table 3. High cross-correlation among
277 grouped sites, such as sites 6, 7, and 8 (northern part) and sites 3, 4, and 5 as well as 12 (southeast
278 coastal area, 0.68-0.87), was found. As expected, sites 5 and 12 had the highest cross-correlation



279 (0.87) due to the proximity. The northern sites and coastal sites showed low cross-correlation. This
280 observed cross-correlation was well preserved in the data generated from both DKNNR and
281 MONR models, as shown in Figure 5 as well as Table 4 and Table 5. However, consistently slight
282 but significant underestimation of cross-correlation was seen for the data generated by the MONR
283 model (see the bottom panel of Figure 5). Note that the errorbars are extended to upper and lower
284 lines of the circles to $1.95 \times$ standard deviation. The difference of RMSE in Table 6 showed this
285 characteristic, as most of the values were positive, to be indicating that the proposed DKNNR
286 model performed better for cross-correlation.

287 The lag-1 cross-correlation of observed data, as shown in Table 7, ranged from 0.22-0.35.
288 The lag-1 cross-correlation for the same site (i.e. $\rho_1(q, s)$, $q=s$) was autocorrelation and was highly
289 related with P_{01} and P_{11} as in Eq. (6). All the lag-1 cross-correlations exhibited similar magnitudes
290 even for autocorrelation. This implies that the lag-1 cross-correlation among the selected sites was
291 as strong as the autocorrelation and as much as the transition probabilities P_{01} and P_{11} , thereof.
292 Relatively low lag-1 cross-correlation was observed between northern sites (6, 7, and 8) and
293 coastal sites (3, 4, and 5), as shown in Table 7.

294 The observed cross-correlation was well preserved in the data generated by the DKNNR
295 model, as shown in the top panel of Figure 6, while the MONR model showed significant
296 underestimation, as seen in the bottom panel of Figure 6. The difference of RMSE shown in Table
297 8 reflects this behavior. In the bottom panel of Figure 6, some of the lag-1 cross-correlations were
298 well preserved, that was aligned with the base line. From Table 8, the MONR model reproduced
299 the autocorrelations well with the shaded values. It is because the lag-1 autocorrelation was
300 indirectly parameterized with the transition probabilities of P_{11} and P_{01} as in Eq. (6). Other than



301 this autocorrelation, the lag-1 cross-correlation was not reproduced with the MONR model. This
302 shortcoming was mentioned by Wilks (1998). Meanwhile, the proposed DKNNR model preserved
303 this statistic without any parameterization.

304 Also, the whole year data instead of the summer season data was tested for model fitting.
305 Note that all the results presented above were with the summer season data (June-September) as
306 mentioned in section 4 on the data description. The lag-1 cross-correlation is shown in Figure 7
307 which indicates that the same characteristic was observed as for the summer season, such that the
308 proposed DKNNR model preserved better the lagged cross-correlation than did the existing
309 MONR model. Other statistics, such as correlation matrix and transition probabilities, exhibited
310 the same results (not shown). Also, other seasons were tried but the estimated correlation matrix
311 was not a positive semi-definite matrix and its inverse cannot be made for multivariate normal
312 distribution in the MONR model. It was because the selected stations were close to each other
313 (around 50-100 km) and produced high cross-correlation, especially in the occurrence during dry
314 seasons. Special remedy should be applied, such as decreasing cross-correlation by force, but
315 further remedy was not applied in the current study since it was not within the current scope and
316 focus.

317 **6.3. Adaptation to climate change**

318 Model adaptability to climate change in hydro-meteorological simulation models is a critical
319 factor, since one of the major applications of the models is to assess the impact of climate change.
320 Therefore, we tested the capability of the proposed model in the current study by adjusting the
321 probabilities of cross-over and mutation as in Eqs.(15) and (16). A number of variations can be
322 made with different conditions.



323 In Figure 8, the changes of transition and marginal probabilities are shown along with
324 increasing the crossover probability P_{cr} from 0.01 to 0.2 with the condition that the candidate
325 value is one and the previous value is also one as in Eq. (15) for the selected 5 stations among the
326 12 stations (from station 1 to station 5, see Table 1 for the detail). The stations were limited in this
327 analysis due to computational time. At each case 100 series were simulated. The average value of
328 the simulated statistics is presented in the figure. It is obvious that the transition probability P_{11}
329 increased as intended along with the increase of P_{cr} . As expected from Eq. (5), P_1 presents that
330 the change of P_1 is highly related to P_{11} . However, the probability P_{01} fluctuated along with the
331 increase of P_{cr} . Elaborate work to adjust all the probabilities is however required.

332 The changes in transition and marginal probabilities are presented in Figure 9 with increasing
333 mutation probability P_m from 0.01 to 0.2 under the condition that the candidate value is one so that
334 the marginal probability P_1 increased. P_{01} also increased along with increasing P_1 . The change of
335 P_{11} was not related with other probabilities. The combination of the adjustment of P_{cr} and P_m with
336 a certain condition to the previous state will allow the specific adaptation for simulating future
337 climatic scenarios.

338 7. Conclusions

339 In the current study, a nonparametric simulation model, based on discrete KNNR and
340 DKNNR, is proposed. The proposed DKNNR model is compared with the existing MONR model.
341 Occurrence and transition probabilities and cross-correlation as well as lag-1 cross-correlation are
342 estimated for both models. Better preservation of cross-correlation and lag-1 cross-correlation with
343 the DKNNR model than the MONR model is observed. For some cases (i.e., the whole year data
344 and other seasons than the summer season), the estimated cross-correlation matrix is not a positive



345 semi-definite matrix so the multivariate normal simulation is not applicable for the MONR model
346 because the tested sites are close to each other with high cross-correlation.

347 Results of this study indicate that the proposed DKNNR model reproduces the occurrence
348 and transition probabilities fairly well and preserves the cross-correlations better than the existing
349 MONR model. Thus, the proposed DKNNR model can be a good alternative for simulating
350 multisite precipitation occurrence.

351 We tested further enhancement of the proposed model for adapting climate change through
352 modifying the mutation and crossover probability P_m and P_{cr} with the current and previous states.
353 The results show that the current model has the capability to adapt to the climate change scenarios
354 but elaborate work is required however. Further study on improving the model adaptability to
355 climate change will be followed in near future.

356 Also, the simulated multisite occurrence can be coupled with a multisite amount model to
357 produce precipitation events, including zero values. Further development can be made for multisite
358 amount models with a nonparametric technique, such as KNNR and bootstrapping.

359 **Code and Data Availability**

360 DKNNR code is written in Matlab and is available at the supplement.

361 The precipitation data employed in the current study is downloadable through
362 <http://www.weather.go.kr/weather/main.jsp>

363 **Acknowledgment**

364 This work was supported by the National Research Foundation of Korea (NRF) grant (NRF-
365 2018R1A2B6001799) funded by the Korean Government (MEST).



366 **Appendix A: Example of DKNNR**

367 In this appendix, one example of DKNNR simulation is presented with observed dataset in
368 Table A 1 (i.e. $\mathbf{x}_i = [x_i^s]_{s \in \{1, S\}}$ for $i=1, \dots, n$; here $S=12$ and $n=16$). The upper part of the table
369 presents the observed precipitation (unit: mm). Its occurrence data is presented in the bottom part
370 of this table. The current precipitation occurrence $\mathbf{X}_c = [X_c^s]_{s \in \{1, \dots, 12\}}$ is shown in the second row of
371 Table A 2. The number of nearest neighbors $k = \sqrt{n} = \sqrt{16} = 4$ and the parameters for GA (i.e. P_c
372 and P_m) are 0.1 and 0.01, respectively. The simulation can be made as follows:

373 (1) Estimate the distance D_i between \mathbf{x}_i and \mathbf{X}_c for $i=1, \dots, n-1$ as in Eq.(11). For example,
374 for $i=1$,

$$375 \quad D_1 = \sum_{s=1}^S |X_c^s - x_1^s| = |0-1| + |1-1| + \dots + |0-1| = 6$$

376 All the estimated distances are shown in the last column of Table A 2.

377 (2) The daily index values are sorted according to the smallest distances shown in the first
378 two columns of Table A 3. The sorted day indices and their corresponding distances are
379 shown in the third and fourth columns of Table A 3. Among k number of sorted indices,
380 one is selected with the weight probability (see Eq.(12)), which is shown in the last
381 column of Table A 3.

382 (3) Simulate a uniform random number (u) between 0 and 1. Say $u=0.321$. This value must
383 be compared with the cumulative weighted probabilities in the last column of Table A 3
384 as [0 0.48 0.72 0.88 1.0]. The corresponding day index is assigned as to where the
385 simulated uniform number falls in the cumulative weighted probabilities, here [0 0.48].



386 Therefore, the selected day, p , is 14. The occurrences of the following day $p+1=15$ for 12
387 stations are selected as in the second row of Table A 4.

388 (4) For GA mixture, another set must be chosen as in step (3). Say $u=0.561$, which falls in
389 $[0.48 \ 0.72]$. The second one should be selected. However, there are a number of days with
390 the same distances. Specifically, six days have the same distances with $D_i=4$. In this case,
391 one among all six days is selected with equal probability. Assume that $p=4$ is selected and
392 the following occurrences are selected as shown in the third row of Table A 4.

393 (5) With two sets, crossover and mutation process is performed as follows:

394 (5-1) Crossover: For each station, a uniform random number (ε) is generated and
395 compared with $P_c=0.1$ here. Say $\varepsilon =0.345$, then skip since $\varepsilon =0.345 > P_c=0.1$. For
396 $s=6$, assume the generated random number, $\varepsilon (=0.051) < P_c(=0.1)$ and then switch
397 the 6th station value of Set 1 into the value of Set 2 (see Table A 4). The occurrence
398 state of X_{c+1}^s turns into 1 from 0 as shown in the fourth row of Table A 4 as well as
399 station 8.

400 (5-2) Mutation: For each station, a uniform random number (ε) is generated and compared
401 with $P_m=0.01$. For $s=12$, assume $\varepsilon =0.009 < P_m=0.01$ and switch the 12th station
402 value of Set 1 with the one selected among all the observed 12th station values with
403 equal probability (here the last column, $s=12$, of the bottom part of Table A 1, [1 1
404 0 0 ... 1]). The occurrence state of X_{c+1}^{12} turns into 0 from 1 as shown in the fourth
405 column of Table A 4.

406 (6) Repeat steps (1)-(5) until the target simulation length is reached.



408 **References**

- 409
410 Apipattanavis, S., Podesta, G., Rajagopalan, B., and Katz, R. W.: A semiparametric
411 multivariate and multisite weather generator, *Water Resources Research*, 43, Artn W11401, 2007.
- 412 Beersma, J. J. and Buishand, A. T.: Multi-site simulation of daily precipitation and
413 temperature conditional on the atmospheric circulation, *Climate Research*, 25, 121-133, 2003.
- 414 Brandsma, T. and Buishand, T. A.: Simulation of extreme precipitation in the Rhine basin
415 by nearest-neighbour resampling, *Hydrology and Earth System Sciences*, 2, 195-209, 1998.
- 416 Buishand, T. A. and Brandsma, T.: Multisite simulation of daily precipitation and
417 temperature in the Rhine basin by nearest-neighbor resampling, *Water Resources Research*, 37,
418 2761-2776, 2001.
- 419 Frost, A. J., Charles, S. P., Timbal, B., Chiew, F. H. S., Mehrotra, R., Nguyen, K. C.,
420 Chandler, R. E., McGregor, J. L., Fu, G., Kirono, D. G. C., Fernandez, E., and Kent, D. M.: A
421 comparison of multi-site daily rainfall downscaling techniques under Australian conditions,
422 *Journal of Hydrology*, 408, 1-18, 2011.
- 423 Hughes, J. P., Guttorp, P., and Charles, S. P.: A non-homogeneous hidden Markov model
424 for precipitation occurrence, *Journal of the Royal Statistical Society. Series C: Applied Statistics*,
425 48, 15-30, 1999.
- 426 Jeong, D. I., St-Hilaire, A., Ouarda, T. B. M. J., and Gachon, P.: A multi-site statistical
427 downscaling model for daily precipitation using global scale GCM precipitation outputs,
428 *International Journal of Climatology*, 33, 2431-2447, 2013.
- 429 Jeong, D. I., St-Hilaire, A., Ouarda, T. B. M. J., and Gachon, P.: Multisite statistical
430 downscaling model for daily precipitation combined by multivariate multiple linear regression and
431 stochastic weather generator, *Climatic Change*, 114, 567-591, 2012.



- 432 Katz, R. W. and Zheng, X.: Mixture model for overdispersion of precipitation, Journal of
433 Climate, 12, 2528-2537, 1999.
- 434 Lall, U. and Sharma, A.: A nearest neighbor bootstrap for resampling hydrologic time
435 series, Water Resources Research, 32, 679-693, 1996.
- 436 Lee, T.: Multisite stochastic simulation of daily precipitation from copula modeling with
437 a gamma marginal distribution, Theoretical and Applied Climatology, doi: 10.1007/s00704-017-
438 2147-0, 2017. 1-10, 2017.
- 439 Lee, T.: Stochastic simulation of precipitation data for preserving key statistics in their
440 original domain and application to climate change analysis, Theoretical and Applied Climatology,
441 124, 91-102, 2016.
- 442 Lee, T. and Ouarda, T. B. M. J.: Identification of model order and number of neighbors
443 for k-nearest neighbor resampling, Journal of Hydrology, 404, 136-145, 2011.
- 444 Lee, T. and Ouarda, T. B. M. J.: Stochastic simulation of nonstationary oscillation hydro-
445 climatic processes using empirical mode decomposition, Water Resources Research, 48, 1-15,
446 2012.
- 447 Lee, T., Ouarda, T. B. M. J., and Jeong, C.: Nonparametric multivariate weather generator
448 and an extreme value theory for bandwidth selection, Journal of Hydrology, 452-453, 161-171,
449 2012.
- 450 Lee, T., Ouarda, T. B. M. J., and Yoon, S.: KNN-based local linear regression for the
451 analysis and simulation of low flow extremes under climatic influence, Climate Dynamics, doi:
452 10.1007/s00382-017-3525-0, 2017. 1-19, 2017.



453 Lee, T. and Park, T.: Nonparametric temporal downscaling with event-based population
454 generating algorithm for RCM daily precipitation to hourly: Model development and performance
455 evaluation, *Journal of Hydrology*, 547, 498-516, 2017.

456 Lee, T., Salas, J. D., and Prairie, J.: An enhanced nonparametric streamflow
457 disaggregation model with genetic algorithm, *Water Resources Research*, 46, 2010a.

458 Lee, T., Salas, J. D., and Prairie, J.: An Enhanced Nonparametric Streamflow
459 Disaggregation Model with Genetic Algorithm, *Water Resources Research*, 46, W08545, 2010b.

460 Mehrotra, R., Srikanthan, R., and Sharma, A.: A comparison of three stochastic multi-site
461 precipitation occurrence generators, *Journal of Hydrology*, 331, 280-292, 2006.

462 Prairie, J. R., Rajagopalan, B., Fulp, T. J., and Zagona, E. A.: Modified K-NN model for
463 stochastic streamflow simulation, *Journal of Hydrologic Engineering*, 11, 371-378, 2006.

464 Rajagopalan, B. and Lall, U.: A k-nearest-neighbor simulator for daily precipitation and
465 other weather variables, *Water Resources Research*, 35, 3089-3101, 1999.

466 St-Hilaire, A., Ouarda, T. B. M. J., Bargaoui, Z., Daigle, A., and Bilodeau, L.: Daily river
467 water temperature forecast model with a k-nearest neighbour approach, *Hydrological Processes*,
468 26, 1302-1310, 2012.

469 Todorovic, P. and Woolhiser, D. A.: Stochastic model of n-day precipitation *Journal of*
470 *Applied Meteorology*, 14, 17-24, 1975.

471 Wilby, R. L., Tomlinson, O. J., and Dawson, C. W.: Multi-site simulation of precipitation
472 by conditional resampling, *Climate Research*, 23, 183-194, 2003.

473 Wilks, D. S.: Multisite downscaling of daily precipitation with a stochastic weather
474 generator, *Climate Research*, 11, 125-136, 1999.



475 Wilks, D. S.: Multisite generalization of a daily stochastic precipitation generation model,

476 Journal of Hydrology, 210, 178-191, 1998.

477 Wilks, D. S. and Wilby, R. L.: The weather generation game: a review of stochastic

478 weather models, Progress in Physical Geography, 23, 329-357, 1999.

479 Zheng, X. and Katz, R. W.: Simulation of spatial dependence in daily rainfall using

480 multisite generators, Water Resources Research, 44, 2008.

481

482



483

484 Table 1. Information on 12 selected stations from Yeongnam province, South Korea.

Order	Station Number [†]	Name	Longitude	Latitude
1	138	Pohang	129.3797	36.0327
2	143	Daegu	128.6189	35.8850
3	152	Ulsan	129.3200	35.5600
4	159	Busan	129.0319	35.1044
5	162	Tongyeong	128.4356	34.8453
6	277	Youngdeok	129.4092	36.5331
7	278	Uisung	128.6883	36.3558
8	279	Gumi	128.3206	36.1306
9	281	Youngcheon	128.9514	35.9772
10	285	Hapcheon	128.1697	35.5650
11	288	Milyang	128.7439	35.4914
12	294	Geojae	128.6044	34.8881

485 [†]The station number indicates the identification number operated by Korea Meteorological
486 Administration (KMA).

487

488



489 Table 2. Occurrence and transition probabilities of observed data and data simulated by DKNNR
 490 and MONR for 12 stations from Yeongnam province, South Korea, during the summer season.
 491 Note that 100 sets with the same record length as the observed data were simulated and the
 492 statistics of 100 sets were averaged.

	Obs			DKNNR			MONR		
	P11	P01	P1	P11	P01	P1	P11	P01	P1
S1	0.57	0.26	0.38	0.55	0.26	0.37	0.58	0.27	0.39
S2	0.56	0.27	0.38	0.57	0.26	0.37	0.58	0.27	0.39
S3	0.57	0.26	0.38	0.56	0.25	0.37	0.58	0.27	0.39
S4	0.58	0.25	0.37	0.56	0.24	0.36	0.60	0.25	0.39
S5	0.58	0.25	0.37	0.57	0.24	0.36	0.60	0.25	0.38
S6	0.52	0.24	0.33	0.50	0.24	0.32	0.53	0.25	0.35
S7	0.54	0.25	0.35	0.53	0.24	0.34	0.56	0.26	0.37
S8	0.55	0.25	0.36	0.52	0.24	0.34	0.57	0.26	0.38
S9	0.54	0.24	0.35	0.54	0.23	0.34	0.55	0.26	0.36
S10	0.58	0.24	0.37	0.56	0.23	0.35	0.57	0.26	0.38
S11	0.56	0.24	0.36	0.55	0.23	0.34	0.57	0.25	0.37
S12	0.59	0.24	0.37	0.58	0.24	0.36	0.61	0.25	0.39

493
 494
 495



496 Table 3. Cross-correlation of observed data for 12 stations from Yeongnam province, South
497 Korea.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	1.00	0.71	0.71	0.65	0.58	0.71	0.65	0.64	0.76	0.65	0.67	0.60
S2	0.71	1.00	0.67	0.64	0.61	0.65	0.69	0.72	0.80	0.72	0.73	0.62
S3	0.71	0.67	1.00	0.75	0.68	0.62	0.57	0.58	0.68	0.67	0.75	0.70
S4	0.65	0.64	0.75	1.00	0.79	0.57	0.56	0.56	0.66	0.67	0.74	0.82
S5	0.58	0.61	0.68	0.79	1.00	0.52	0.54	0.56	0.61	0.65	0.70	0.87
S6	0.71	0.65	0.62	0.57	0.52	1.00	0.70	0.66	0.68	0.59	0.60	0.55
S7	0.65	0.69	0.57	0.56	0.54	0.70	1.00	0.79	0.71	0.64	0.63	0.56
S8	0.64	0.72	0.58	0.56	0.56	0.66	0.79	1.00	0.71	0.68	0.65	0.57
S9	0.76	0.80	0.68	0.66	0.61	0.68	0.71	0.71	1.00	0.69	0.72	0.62
S10	0.65	0.72	0.67	0.67	0.65	0.59	0.64	0.68	0.69	1.00	0.77	0.66
S11	0.67	0.73	0.75	0.74	0.70	0.60	0.63	0.65	0.72	0.77	1.00	0.71
S12	0.60	0.62	0.70	0.82	0.87	0.55	0.56	0.57	0.62	0.66	0.71	1.00

498

499

500



501 Table 4. Averaged cross-correlation of the 100 simulated series from the DKNNR model for 12
502 stations from Yeongnam province, South Korea.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	1.00	0.69	0.69	0.63	0.57	0.69	0.64	0.63	0.74	0.63	0.66	0.59
S2	0.69	1.00	0.65	0.63	0.61	0.63	0.68	0.70	0.77	0.71	0.72	0.61
S3	0.69	0.65	1.00	0.73	0.66	0.60	0.56	0.57	0.67	0.66	0.73	0.68
S4	0.63	0.63	0.73	1.00	0.77	0.56	0.55	0.56	0.64	0.65	0.72	0.80
S5	0.57	0.61	0.66	0.77	1.00	0.51	0.53	0.55	0.60	0.64	0.69	0.84
S6	0.69	0.63	0.60	0.56	0.51	1.00	0.68	0.65	0.66	0.58	0.59	0.54
S7	0.64	0.68	0.56	0.55	0.53	0.68	1.00	0.76	0.70	0.63	0.61	0.55
S8	0.63	0.70	0.57	0.56	0.55	0.65	0.76	1.00	0.70	0.67	0.64	0.56
S9	0.74	0.77	0.67	0.64	0.60	0.66	0.70	0.70	1.00	0.68	0.71	0.61
S10	0.63	0.71	0.66	0.65	0.64	0.58	0.63	0.67	0.68	1.00	0.75	0.65
S11	0.66	0.72	0.73	0.72	0.69	0.59	0.61	0.64	0.71	0.75	1.00	0.70
S12	0.59	0.61	0.68	0.80	0.84	0.54	0.55	0.56	0.61	0.65	0.70	1.00

503

504



505

506

507

Table 5. Averaged cross-correlation of 100 simulated series from the MONR model for 12 stations from Yeongnam province.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	1.00	0.69	0.69	0.59	0.57	0.67	0.63	0.62	0.74	0.62	0.63	0.57
S2	0.69	1.00	0.63	0.62	0.60	0.64	0.66	0.69	0.76	0.70	0.70	0.60
S3	0.69	0.63	1.00	0.71	0.65	0.59	0.55	0.56	0.64	0.64	0.72	0.68
S4	0.59	0.62	0.71	1.00	0.78	0.54	0.54	0.54	0.63	0.62	0.70	0.78
S5	0.57	0.60	0.65	0.78	1.00	0.51	0.52	0.55	0.59	0.62	0.66	0.84
S6	0.67	0.64	0.59	0.54	0.51	1.00	0.67	0.63	0.67	0.56	0.59	0.52
S7	0.63	0.66	0.55	0.54	0.52	0.67	1.00	0.76	0.67	0.61	0.59	0.53
S8	0.62	0.69	0.56	0.54	0.55	0.63	0.76	1.00	0.69	0.65	0.62	0.54
S9	0.74	0.76	0.64	0.63	0.59	0.67	0.67	0.69	1.00	0.65	0.70	0.59
S10	0.62	0.70	0.64	0.62	0.62	0.56	0.61	0.65	0.65	1.00	0.73	0.62
S11	0.63	0.70	0.72	0.70	0.66	0.59	0.59	0.62	0.70	0.73	1.00	0.68
S12	0.57	0.60	0.68	0.78	0.84	0.52	0.53	0.54	0.59	0.62	0.68	1.00

508

509

510

511



512 Table 6. The difference of RMSE of cross-correlation between MONR and DKNNR. Note that
 513 the positive value indicates that the DKNNR model better performs in preserving the cross-
 514 correlation, while a negative value (underlined) shows that the MONR model better performs.

MONR- DKNNR	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	0.000	0.000	0.007	0.040	0.005	0.016	0.004	0.008	0.006	0.016	0.026	0.018
S2	0.000	0.000	0.016	0.016	0.005	<u>0.005</u> [†]	0.016	0.011	0.018	0.009	0.016	0.010
S3	0.007	0.016	0.000	0.016	0.011	0.009	0.004	0.005	0.025	0.020	0.014	0.001
S4	0.040	0.016	0.016	0.000	<u>0.002</u>	0.018	0.013	0.015	0.016	0.027	0.023	0.020
S5	0.005	0.005	0.011	<u>0.002</u>	0.000	0.007	0.012	0.007	0.006	0.016	0.021	0.007
S6	0.016	<u>0.005</u>	0.009	0.018	0.007	0.000	0.009	0.014	<u>0.006</u>	0.019	0.001	0.016
S7	0.004	0.016	0.004	0.013	0.012	0.009	0.000	0.008	0.023	0.014	0.018	0.010
S8	0.008	0.011	0.005	0.015	0.007	0.014	0.008	0.000	0.010	0.017	0.024	0.015
S9	0.006	0.018	0.025	0.016	0.006	<u>0.006</u>	0.023	0.010	0.000	0.023	0.007	0.017
S10	0.016	0.009	0.020	0.027	0.016	0.019	0.014	0.017	0.023	0.000	0.018	0.026
S11	0.026	0.016	0.014	0.023	0.021	0.001	0.018	0.024	0.007	0.018	0.000	0.020
S12	0.018	0.010	0.001	0.020	0.007	0.016	0.010	0.015	0.017	0.026	0.020	0.000

515 [†]Underline represents a negative value implying that the MONR model better performs.

516

517

518

519



520 Table 7. Lag-1 cross-correlation of observed data for 12 stations from Yeongnam province,
 521 South Korea.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	0.30 [‡]	0.27	0.31	0.27	0.25	0.30	0.26	0.25	0.28	0.27	0.29	0.26
S2	0.28	0.29	0.29	0.27	0.25	0.29	0.28	0.27	0.30	0.29	0.32	0.26
S3	0.29	0.26	0.31	0.30	0.27	0.27	0.25	0.24	0.27	0.27	0.30	0.27
S4	0.29	0.28	0.32	0.34	0.31	0.29	0.27	0.26	0.28	0.29	0.31	0.32
S5	0.30	0.29	0.32	0.34	0.34	0.29	0.27	0.27	0.30	0.30	0.34	0.35
S6	0.25	0.22	0.26	0.24	0.23	0.28	0.24	0.22	0.25	0.23	0.25	0.23
S7	0.26	0.26	0.27	0.26	0.25	0.29	0.30	0.27	0.27	0.27	0.28	0.26
S8	0.29	0.29	0.29	0.27	0.26	0.30	0.31	0.30	0.30	0.30	0.31	0.27
S9	0.29	0.29	0.30	0.28	0.26	0.29	0.27	0.27	0.30	0.29	0.32	0.27
S10	0.29	0.31	0.32	0.31	0.29	0.29	0.30	0.30	0.32	0.33	0.34	0.29
S11	0.27	0.29	0.31	0.30	0.27	0.28	0.27	0.26	0.29	0.29	0.32	0.28
S12	0.30	0.30	0.33	0.35	0.33	0.30	0.28	0.27	0.30	0.31	0.35	0.35

522 [‡]Shaded values represents lag-1 autocorrelation (i.e. the one lagged correlation for the same site).

523

524



525 Table 8. The difference of RMSE of lag-1 cross-correlation between MONR and DKNNR. Note
526 that a positive value indicates that the DKNNR model better performs in preserving lag-1 cross-
527 correlation, while a negative value (underlined) shows that the MONR model better performs.

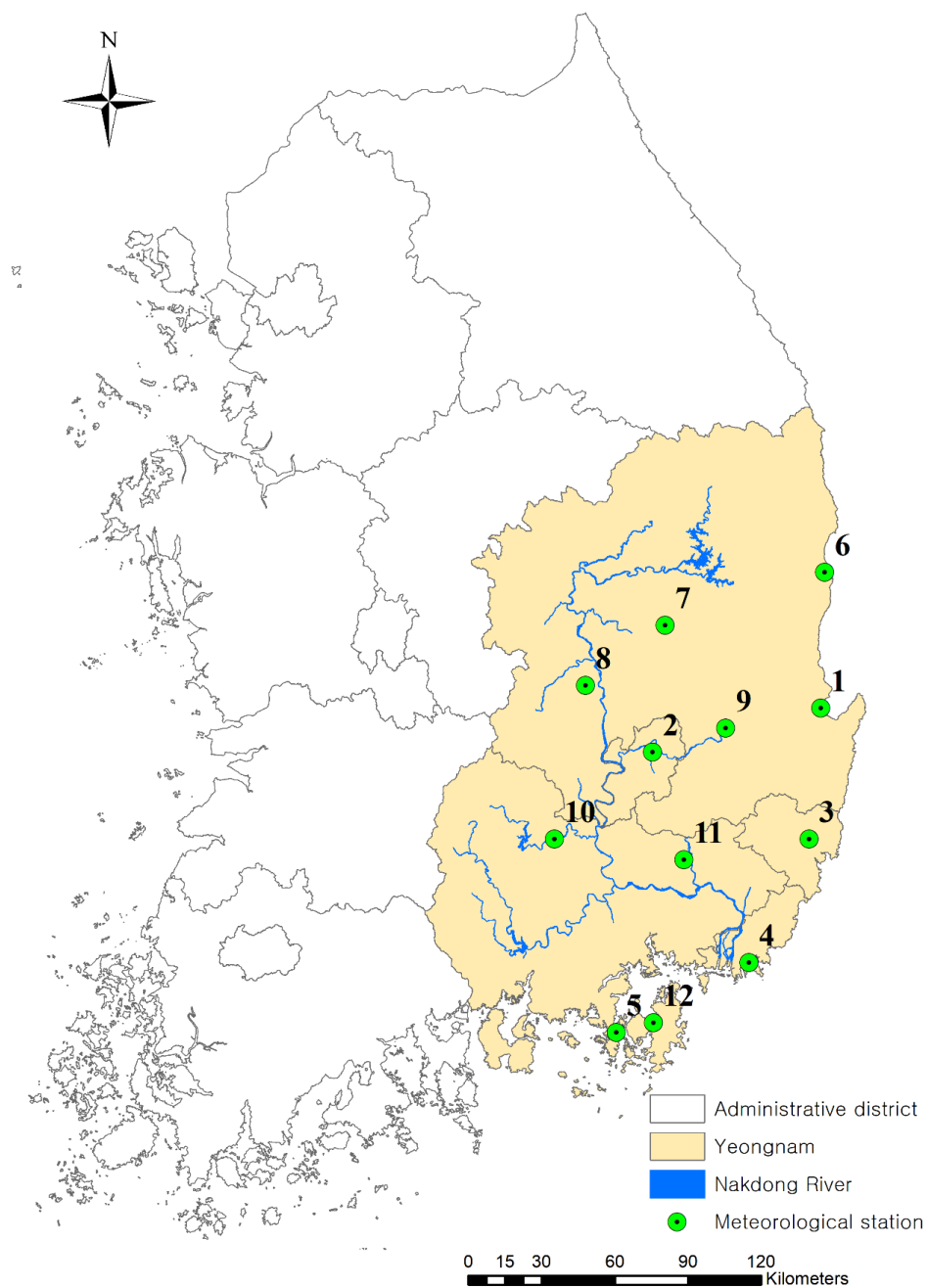
MONR- DKNNR	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	<u>0.003</u>	0.050	0.081	0.062	0.044	0.098	0.060	0.046	0.048	0.050	0.076	0.046
S2	0.056	<u>0.004</u> [†]	0.078	0.053	0.036	0.092	0.065	0.053	0.064	0.043	0.078	0.037
S3	0.065	0.053	<u>0.002</u>	0.048	0.041	0.096	0.070	0.054	0.062	0.045	0.060	0.019
S4	0.093	0.084	0.087	<u>0.001</u> [‡]	0.040	0.123	0.089	0.083	0.081	0.065	0.078	0.034
S5	0.109	0.096	0.111	0.074	<u>0.002</u>	0.129	0.106	0.088	0.110	0.076	0.120	0.045
S6	0.031	0.016	0.062	0.043	0.044	<u>0.001</u>	0.020	0.017	0.031	0.029	0.046	0.029
S7	0.053	0.048	0.081	0.063	0.057	0.085	<u>0.003</u>	0.025	0.060	0.048	0.078	0.056
S8	0.089	0.077	0.096	0.080	0.063	0.111	0.070	<u>0.001</u>	0.084	0.070	0.101	0.063
S9	0.049	0.047	0.091	0.064	0.052	0.088	0.055	0.050	<u>0.004</u>	0.064	0.084	0.055
S10	0.085	0.094	0.107	0.090	0.065	0.123	0.107	0.093	0.106	<u>0.000</u>	0.095	0.061
S11	0.065	0.064	0.076	0.054	0.036	0.096	0.081	0.062	0.064	0.032	<u>0.001</u>	0.034
S12	0.118	0.102	0.105	0.080	0.043	0.138	0.108	0.096	0.115	0.093	0.120	<u>0.000</u>

528 [†]Underline represents a negative value implying that the MONR model better performs.

529 [‡]Shaded values represent lag-1 autocorrelation (i.e. the lagged-1 correlation for the same site).

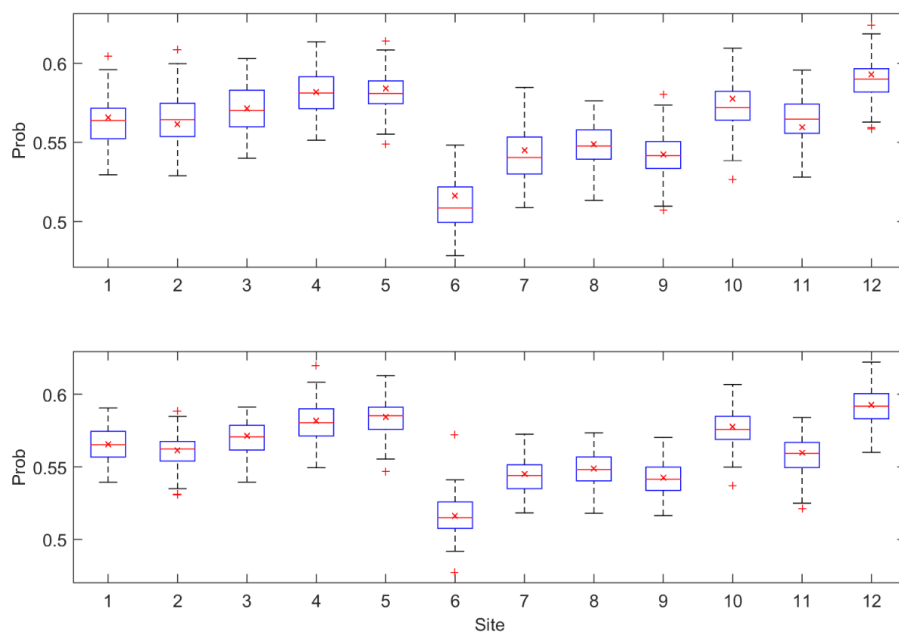
530

531



532

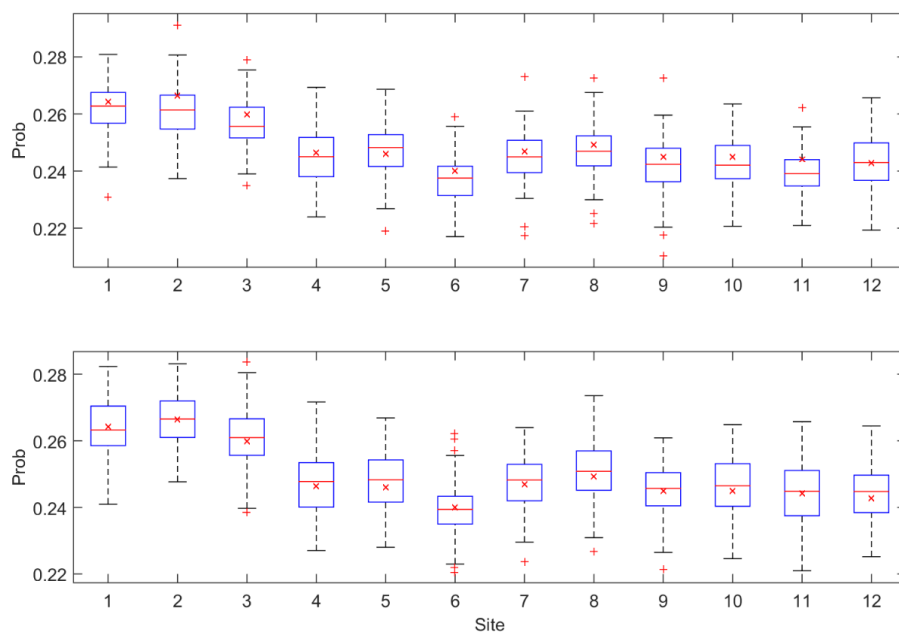
533 Figure 1. Locations of 12 selected weather stations at the Yeongnam province. See Table 1 for
534 further information about the stations.



535

536 Figure 2. Boxplots of the P11 probability for the simulated data from the DKNNR model (top
537 panel) and the MONR model (bottom panel) as well as the observed (x marker) for the 12
538 selected weather stations from the Yeongnam province.

539



540

541 Figure 3. Boxplots of the P01 probability for the data simulated from the DKNNR model (top
542 panel) and the MONR model (bottom panel) as well as the observed (x marker) for the 12
543 selected weather stations from the Yeongnam province.

544

545

546

547

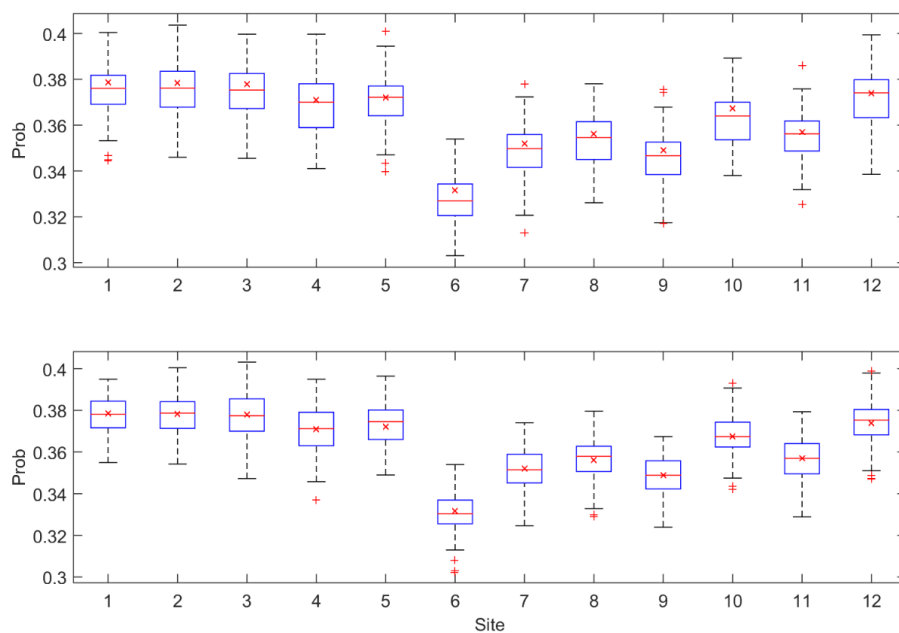
548

549

550

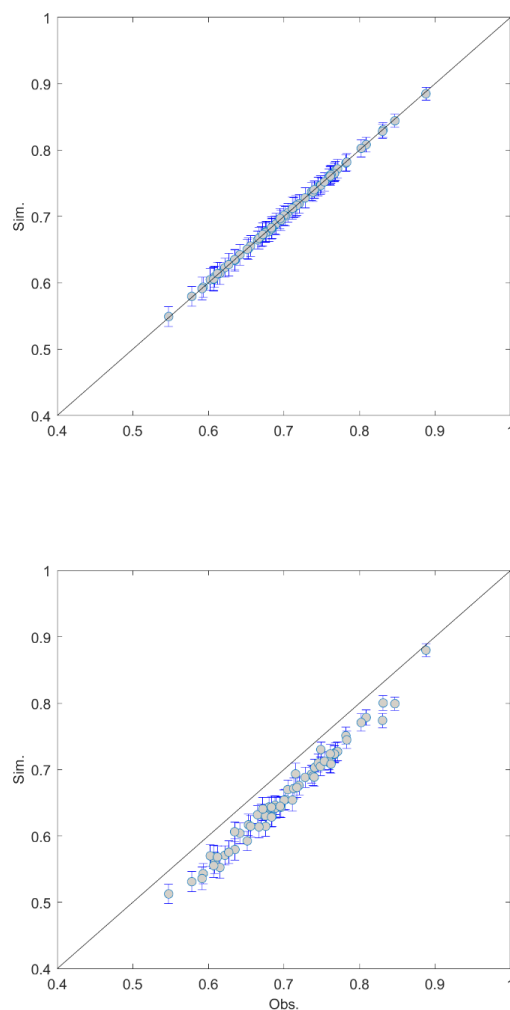
551

552



553

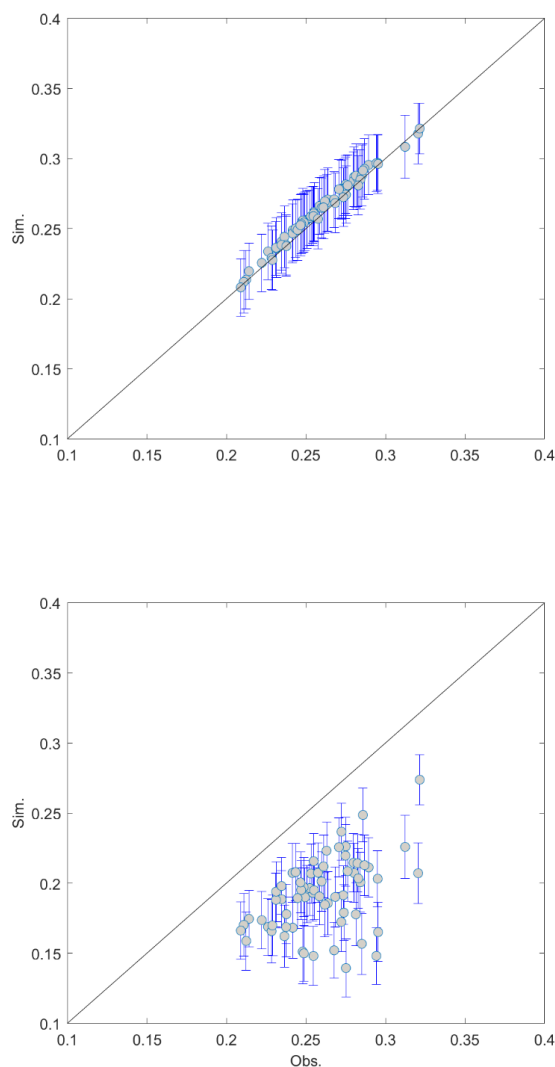
554 Figure 4. Boxplots of the P1 probability for the data simulated from the DKNNR model (top
555 panel) and the MONR model (bottom panel) as well as the observed (x marker) for the 12
556 selected weather stations from the Yeongnam province.



557

558 Figure 5. Scatterplot of cross-correlations between 12 weather stations for the observed data (X
559 coordinate) and the generated data (Y coordinate) generated from the DKNNR model (top panel)
560 and the MONR model (bottom panel). The cross-correlations from 100 generated series are
561 averaged for the filled circle and the errorbars upper and lower extended lines indicate the range
562 of $1.95 \times$ standard deviation.

563

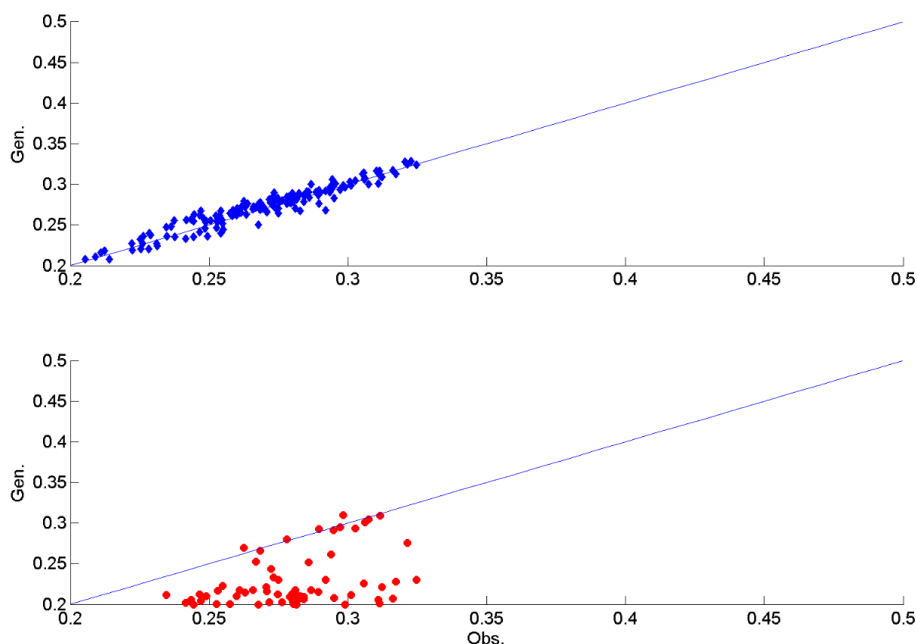


564

565 Figure 6. Scatterplot of lag-1 cross-correlations between 12 weather stations for the observed
566 data (X coordinate) and the generated data (Y coordinate) generated from the DKNNR model
567 (top panel) and the MONR model (bottom panel). The cross-correlations from 100 generated
568 series are averaged for the filled circle and the errorbars upper and lower extended lines indicate
569 the range of $1.95 \times$ standard deviation.



570



571

572 Figure 7. Scatterplot of lag-1 cross-correlations between 12 weather stations for the observed
573 data (X coordinate) and the generated data (Y coordinate) generated from the DKNNR model
574 (top panel) and the MONR model (bottom panel) with the whole year data not with the summer
575 season. The cross-correlations from 100 generated series are averaged.

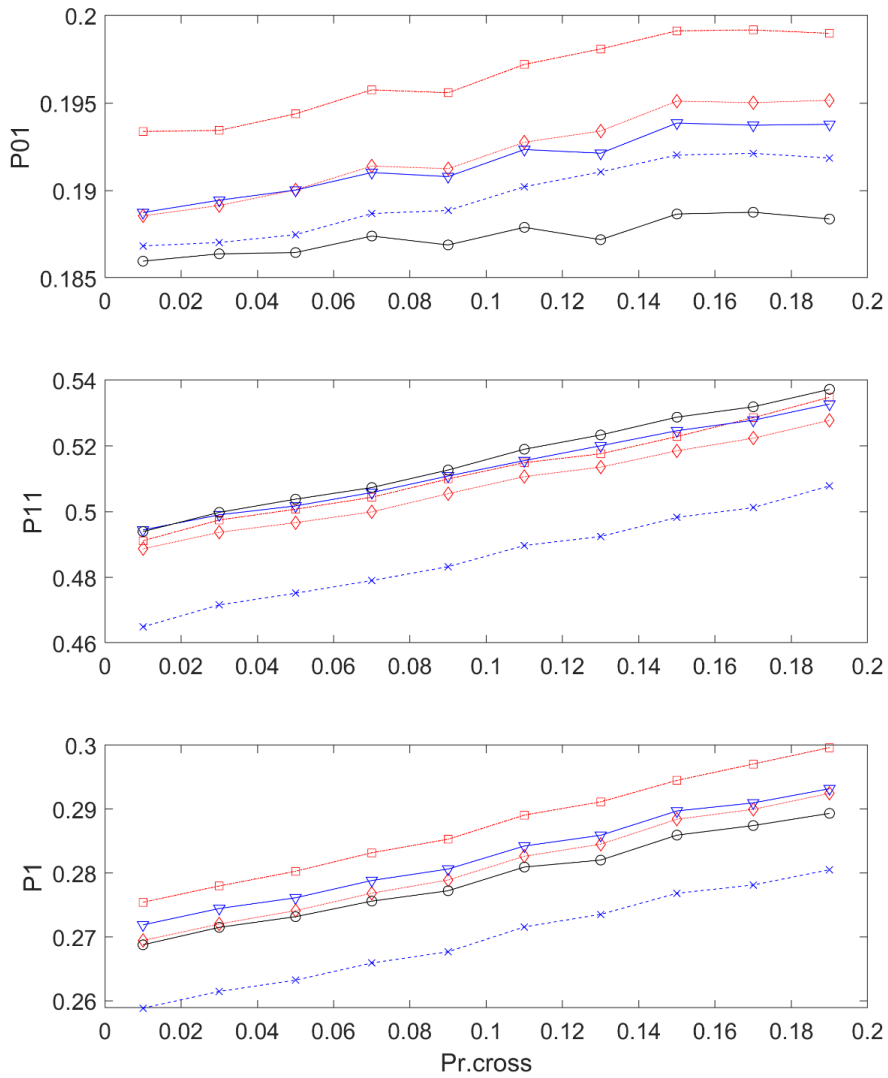
576

577

578

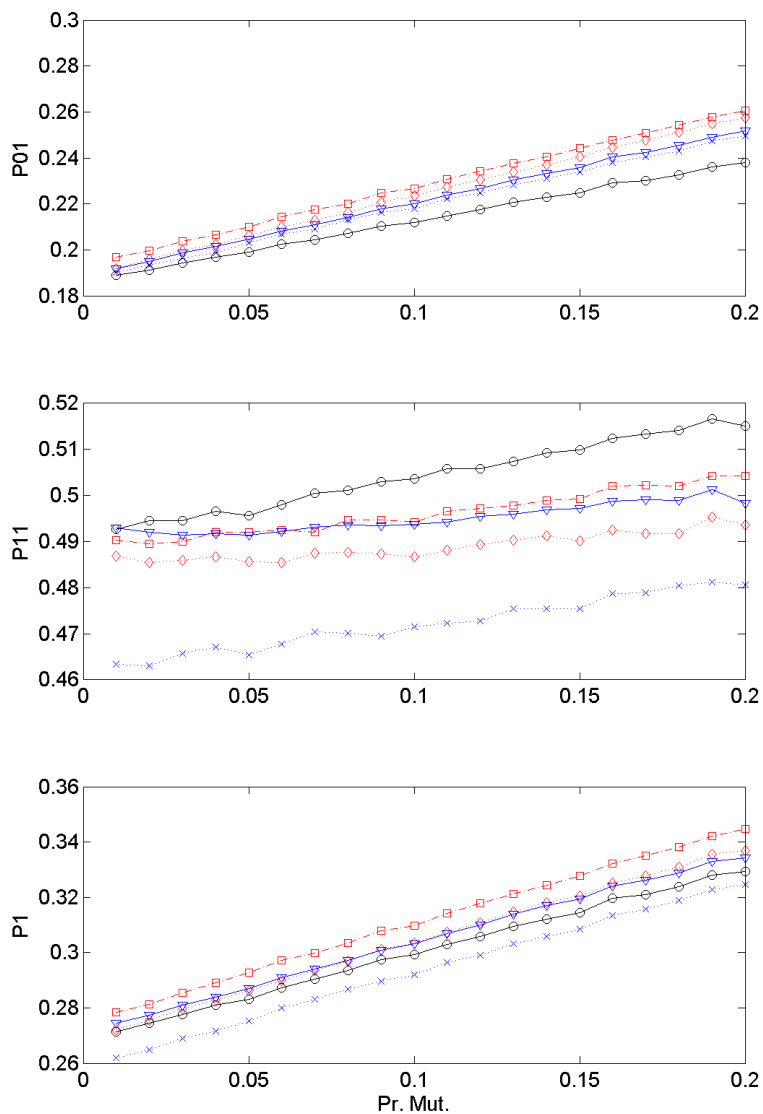
579

580



581

582 Figure 8. Transition probabilities and marginal distribution for the selected five stations along
 583 with changing the cross-over probability P_{cr} with the condition that the candidate value is one
 584 and the previous value is also one. See Eq.(15) for the detail.
 585



586

587 Figure 9. Transition probabilities and marginal distribution along with changing the cross-over
 588 probability with the condition that the mutation is processed only if the candidate value is one.
 589 See Eq.(16) for the detail.

590

591



592

593 Table A 1. Example dataset of daily rainfall with 12 weather stations and 16 days for measured
594 rainfall (mm) in the upper part of this table and its corresponding occurrences in the bottom part
595 of this table.

Day	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
1	2.0	2.9	1.2	0.0	0.0	1.8	4.0	8.9	2.0	4.6	1.3	0.6
2	52.6	39.8	47.2	17.4	11.8	31.0	30.0	33.7	52.0	57.8	37.0	17.5
3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
4	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
5	0.2	1.0	1.4	1.9	12.3	0.0	0.0	0.0	0.7	3.1	3.5	8.1
6	14.8	0.2	0.8	0.2	5.0	0.0	0.0	18.0	0.0	0.0	0.6	3.1
7	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
8	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
9	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
10	0.0	1.0	0.0	0.4	0.0	3.8	0.0	0.1	0.0	0.0	0.0	0.0
11	7.1	6.4	12.8	12.8	13.6	2.3	2.0	5.4	6.0	7.3	16.4	20.3
12	0.0	0.0	0.0	0.0	5.5	0.0	0.0	0.0	0.0	0.0	0.0	4.3
13	10.0	1.6	11.6	14.3	1.5	5.4	0.0	0.0	2.5	0.0	2.7	16.1
14	2.3	0.0	0.7	0.0	0.0	1.4	0.0	0.0	0.0	0.0	0.0	0.0
15	31.5	4.3	30.6	12.7	14.4	25.8	3.5	0.8	5.0	2.7	6.5	20.3
16	37.0	7.8	30.1	11.2	9.6	36.8	2.5	4.7	13.5	1.7	10.1	14.1
Day	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
1	1	1	1	0	0	1	1	1	1	1	1	1
2	1	1	1	1	1	1	1	1	1	1	1	1
3	0	0	0	0	0	0	0	0	0	0	0	0
4	0	0	0	0	0	0	0	0	0	0	0	0
5	1	1	1	1	1	0	0	0	1	1	1	1
6	1	1	1	1	1	0	0	1	0	0	1	1
7	0	0	0	0	0	0	0	0	0	0	0	0
8	0	0	0	0	0	0	0	0	0	0	0	0
9	0	0	0	0	0	0	0	0	0	0	0	0
10	0	1	0	1	0	1	0	1	0	0	0	0
11	1	1	1	1	1	1	1	1	1	1	1	1
12	0	0	0	0	1	0	0	0	0	0	0	1
13	1	1	1	1	1	1	0	0	1	0	1	1
14	1	0	1	0	0	1	0	0	0	0	0	0
15	1	1	1	1	1	1	1	1	1	1	1	1
16	1	1	1	1	1	1	1	1	1	1	1	1



596

597 Table A 2. Example dataset for estimating distances. The second row presents the current daily
 598 precipitation occurrences for 12 stations and the rows below show the absolute difference
 599 between the current occurrences (X_c) and the observed data in Table A 1. The last column
 600 presents the distances in Eq. (11).

day	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12	Dist
X_c	0	1	1	0	0	1	1	0	0	0	0	0	
1	1	0	0	0	0	0	0	1	1	1	1	1	6
2	1	0	0	1	1	0	0	1	1	1	1	1	8
3	0	1	1	0	0	1	1	0	0	0	0	0	4
4	0	1	1	0	0	1	1	0	0	0	0	0	4
5	1	0	0	1	1	1	1	0	1	1	1	1	9
6	1	0	0	1	1	1	1	1	0	0	1	1	8
7	0	1	1	0	0	1	1	0	0	0	0	0	4
8	0	1	1	0	0	1	1	0	0	0	0	0	4
9	0	1	1	0	0	1	1	0	0	0	0	0	4
10	0	0	1	1	0	0	1	1	0	0	0	0	4
11	1	0	0	1	1	0	0	1	1	1	1	1	8
12	0	1	1	0	1	1	1	0	0	0	0	1	6
13	1	0	0	1	1	0	1	0	1	0	1	1	7
14	1	1	0	0	0	0	1	0	0	0	0	0	3
15	1	0	0	1	1	0	0	1	1	1	1	1	8
16	1	0	0	1	1	0	0	1	1	1	1	1	8

601

602



603

604 Table A 3. Example for selecting one sequence for \mathbf{X}_{c+1} . The second row presents the distances
 605 in Table A 2. The third and fourth columns show the sorted days and distances for the smallest
 606 distances to the largest in the second column. The fourth row presents the probabilities estimated
 607 with Eq. (12). Note that there are six days whose distances are the same with each other. In this
 608 case all the days are included and among six days, one is selected with equal probabilities.

Day	Dist.	Sorted Day	Sorted Dist	Prob
1	6	14	3	0.48
2	8	3	4	0.24
3	4	4	4	0.16
4	4	7	4	0.12
5	9	8	4	
6	8	9	4	
7	4	10	4	
8	4	1	6	
9	4	12	6	
10	4	13	7	
11	8	2	8	
12	6	6	8	
13	7	11	8	
14	3	15	8	
15	8	16	8	
16	8	5	9	

609

610



611 Table A 4. Example for GA mixture for \mathbf{X}_{c+1} . The second and third rows present two selected
 612 sets, while the third row shows the final set for \mathbf{X}_{c+1} with the crossover at S6 and S8 and the
 613 mutation for S12.

	Assigned day, p	Selected day, $p+1$	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
Set1	14	15	1	0	0	1	1	0	0	1	1	1	1	1
Set2	4	5	1	0	0	1	1	1	1	0	1	1	1	1
Final			1	0	0	1	1	<u>1</u>	0	<u>0</u>	1	1	1	0

614

615

616