

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Discrete k-nearest neighbor resampling for simulating multisite precipitation occurrence and adaption to climate change

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

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

Taesam Lee¹ and Vijay P. Singh²

¹ Department of Civil Engineering, ERI, Gyeongsang National University,
501 Jinju-daero, Jinju, Gyeongnam, South Korea, 660-701

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

Corresponding Author :
Taesam Lee, Ph.D.
Gyeongsang National University, Dept. of Civil Engineering
Tel)+82-55-772-1797, Fax)+82-55-772-1799
Email) tae3lee@gnu.ac.kr

27

Abstract

28 Stochastic weather simulation models are commonly employed in water resources management
29 agricultural applications, forest management, transportation management, and recreational
30 activities. The data simulated by these models, such as precipitation, temperature, and wind, are
31 used as input for hydrological and agricultural models. Stochastic simulation of multisite
32 precipitation occurrence is a challenge because of its intermittent characteristics as well as spatial
33 and temporal cross-correlation. The multisite occurrence model with standard normal variate
34 (MONR) has been used for preserving key statistics and contemporaneous correlation, but it
35 cannot reproduce lagged crosscorrelation between stations and long stochastic simulation is
36 therefore required to estimate its parameters. Employing a nonparametric technique, k-nearest
37 neighbor resampling (KNNR), and coupling it with Genetic Algorithm (GA), this study proposes
38 a novel simulation method for multisite precipitation occurrence, overcoming the shortcomings of
39 the existing MONR model. The proposed discrete version of KNNR (DKNNR) model is compared
40 with an existing parametric model, called multisite occurrence model with standard normal variate
41 (MONR). The datasets simulated from both the DKNNR model and the MONR model are tested
42 using a number of statistics, such as occurrence and transition probabilities as well as temporal
43 and spatial cross-correlations. Results show that the proposed DKNNR model can be a good
44 alternative for simulating multisite precipitation occurrence, while preserving the lagged
45 crosscorrelation between sites and simulating multisite occurrence from a simple and direct
46 procedure without no parameterization. We also tested the model capability to adapt climate
47 change. It is shown that the model is capable but further improvement is required to have specific
48 variations of the occurrence probability due to climate change. Combining with the generated

49 occurrence, the multisite precipitation amount can then be simulated by any multisite amount
50 model.

51

52 **1. Introduction**

53 Stochastic simulation of weather variables has been employed for water resources
54 management, hydrological design, agricultural applications, forest management, transportation
55 management, recreation activities, filling in missing historical data, extending observed records,
56 simulating data, and simulating different weather conditions. Stochastic simulation models play a
57 key role in producing weather sequences, while preserving the statistical characteristics of
58 observed data. A number of stochastic weather simulation models have been developed using
59 parametric and nonparametric approaches (Lee, 2017; Lee et al., 2012; Wilby et al., 2003; Wilks,
60 1999; Wilks and Wilby, 1999).

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

67 Among weather variables, the precipitation variable possesses intermittency and zero values
68 between precipitation events, and to properly reproduce them is difficult and remains a challenge
69 (Beersma and Buishand, 2003; Hughes et al., 1999; Katz and Zheng, 1999). Due to this difficulty,
70 precipitation is simulated separately from other variables. The main method for reproducing
71 intermittency has been the multiplication of precipitation occurrence and an amount as $Z=X \cdot Y$,
72 where X is the occurrence (binary as either 0 or 1) and Y is the amount (Jeong et al., 2013; Lee and

73 Park, 2017; Todorovic and Woolhiser, 1975). The spatial and temporal dependence in the
74 occurrence and amount of precipitation introduces further complexity in multisite simulation.

75 Wilks (1998) presented a multisite simulation model for the occurrence process (i.e. X) using
76 the standard normal variable that is spatially dependent, representing the relation between the
77 occurrence variable and the standard normal variable with simulation data. Originally, the
78 occurrence of precipitation had been simulated with discrete Markov Chain (MC) model (Katz,
79 1977). Compared to the MC model requiring a significant number of parameters to generate
80 multisite occurrence, the multisite occurrence model proposed by Wilks (1998) transforms the
81 standard normal variate and simulates the sequence with multivariate normal distribution, and then
82 back-transforms the multivariate normal sequence to the original domain. The model is able to
83 reproduce the contemporaneous multisite dependence structure and lagged dependence only for
84 the same site while requiring a complex simulation process to estimate parameter for each site and
85 being unable to preserve lagged dependence between sites.

86 Meanwhile, Lee et al. (2010a) proposed a nonparametric-based stochastic simulation model
87 for hydrometeorological variables. They overcame the shortcoming of a previous nonparametric
88 simulation model (Lall and Sharma, 1996), called k-nearest neighbor resampling (KNNR) such
89 that the simulated data cannot produce patterns different from those of the observed data
90 (Brandsma and Buishand, 1998; Mehrotra et al., 2006; St-Hilaire et al., 2012). In addition to this
91 KNNR, Lee et al. (2010a) used a meta-heuristic algorithm Genetic Algorithm (GA) that led to the
92 reproduction of similar populations by mixing the simulated dataset. While the KNNR is employed
93 to find similar historical analogues of multisite occurrence to the current status of a simulation
94 series, GA is applied to use its skill to generate a new descendant from the historical parent chosen
95 with the KNNR. In this procedure, the multisite occurrence of the precipitation variable can be

96 simulated while preserving spatial and temporal correlations. Note that meta-heuristic techniques
97 to GA have been popularly employed in a number of hydrometeorological applications (Chau,
98 2017; Fotovatikhah et al., 2018; Taormina et al., 2015; Wang et al., 2013). A number of variants
99 of KNNR-GA have since been applied (Lee et al., 2012; Lee and Park, 2017). None of these
100 models can adopt the multisite occurrence in precipitation whose characteristics are binary and
101 temporally and spatially related.

102 Therefore, in the current study we propose a novel stochastic simulation method for multisite
103 occurrence of the precipitation variable with the KNNR-GA based nonparametric approach that
104 (1) simulates multisite occurrence with a simple and direct procedure without parameterization of
105 all the required occurrence probabilities; and (2) reproduces the complex temporal and spatial
106 correlation between stations as well as the basic occurrence probabilities. Note that the proposed
107 nonparametric model is compared with the most popularly employed model proposed by Wilks
108 (1998). Even though the multisite occurrence data from this model (Wilks, 1998) preserves various
109 statistical characteristics of the observed data well, significant underestimation of lagged cross-
110 correlation still exists. Furthermore, the relation between standard normal variable and occurrence
111 variable relies on long stochastic simulation.

112 The paper is organized as follows. The next section presents a mathematical background of
113 existing multisite occurrence modeling. The modeling procedure is discussed in section 3. The
114 study area and data are reported in section 4. The model is applied in section 5. Results of the
115 proposed model are discussed in section 6, and summary and conclusions are presented in section
116 7.

117 2. Background

118 2.1. Single site occurrence modeling

119 Let X_t^s represent the occurrence of daily precipitation for a location s ($s=1, \dots, S$) on day t
120 ($t=1, \dots, n$; n is the number observed days) and let X_t^s be either zero for dry day or one for wet day.

121 The first order Markov chain model for X_t^s is defined with the assumption that the occurrence
122 probability of a wet day is fully defined by the previous day as

$$123 \Pr\{X_t^s = 1 \mid X_{t-1}^s = 0\} = p_{01}^s \quad (1)$$

$$124 \Pr\{X_t^s = 1 \mid X_{t-1}^s = 1\} = p_{11}^s \quad (2)$$

125 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
126 with the same previous condition. This consists of a transition probability matrix (TPM) as

$$127 TPM^s = \begin{bmatrix} p_{00}^s & p_{01}^s \\ p_{10}^s & p_{11}^s \end{bmatrix} = \begin{bmatrix} 1 - p_{01}^s & p_{01}^s \\ 1 - p_{11}^s & p_{11}^s \end{bmatrix} \quad (3)$$

128 The marginal distributions of TPM (i.e. p_0 and p_1) can be expressed with TPM and its condition of
129 $p_0 + p_1 = 1$ as:

$$130 p_0^s = \frac{p_{01}^s}{1 + p_{01}^s - p_{11}^s} \quad (4)$$

$$131 p_1^s = \frac{1 - p_{11}^s}{1 + p_{01}^s - p_{11}^s} \quad (5)$$

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

$$135 \quad \rho_1(s, s) = p_{11}^s - p_{01}^s \quad (6)$$

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

$$137 \quad X_t^s = \begin{cases} 1 & \text{if } u_t^s \leq p_{i1}^s \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

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

$$141 \quad X_t^s = \begin{cases} 1 & \text{if } w_t^s \leq \Phi^{-1}(p_{i1}^s) \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

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

143 **2.2. Multisite occurrence modeling**

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

$$147 \quad \tau(q, s) = \text{corr}[w_t^q, w_t^s] \quad (9)$$

148 Also, the correlation of the original occurrence variate is

149
$$\rho(q, s) = \text{corr}[X_t^q, X_t^s] \quad (10)$$

150 Once the correlation of the standard normal variate is known, the simulation of multisite
151 precipitation occurrence is straightforward. Multivariate standard normal distribution is used with
152 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
153 elements of $\tau(q, s)$ for $q \in \{1, \dots, S\}$ and $s \in \{1, \dots, S\}$.

154 Since direct estimation of $\tau(q, s)$ is not applicable, a simulation technique is used to estimate
155 $\tau(q, s)$ from $\rho(q, s)$. A long sequence of the occurrence process is simulated with different values
156 of $\tau(q, s)$ and its corresponding correlation of the original domain $\rho(q, s)$ is estimated with the
157 simulated long sequence by the inverse standard normal cumulative function (i.e. Φ^{-1}). A curve
158 between $\tau(q, s)$ and $\rho(q, s)$ is derived from this long simulation with the MONR model and is
159 employed for the parameter estimation for real application.

160 **3. DKNNR**

161 **3.1. DKNNR modeling procedure**

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

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

168 (1) Estimate the distance between the current (i.e. time index: c) multisite occurrence
 169 X_c^s and the observed multisite occurrence x_i^s . Here, the distance is measured for
 170 $i=1, \dots, n-1$ as

$$171 \quad D_i = \sum_{s=1}^S |X_c^s - x_i^s| \quad (11)$$

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

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

$$177 \quad w_m = \frac{1/m}{\sum_{j=1}^k 1/j}, \quad m = 1, \dots, k \quad (12)$$

178 (4) Assume the selected time index from step (3) as p . Note that there are a number of
 179 values that have the same distance as the selected D_p , since D_p is a natural number
 180 between 0 and S . For example, if $S=2$ and $X_c^1=0$ and $X_c^2=1$, the two sequences have
 181 the same $D=1$ as $[x_i^1=0$ and $x_i^2=0]$ and $[x_i^1=1$ and $x_i^2=1]$. In this case, a random
 182 selection procedure is required to take into account the cases with the same quantity.
 183 One particular time index is randomly selected with the equal probabilities among
 184 the time indices of the same distances. Note that instead of the random selection, one
 185 can always use the first one. In such a case, only one historical combination of
 186 multisite occurrences will be selected.

187 (5) Assign the binary vector of the proceeding index of the selected time as

188 $\mathbf{x}_{p+1} = [x_{p+1}^s]_{s \in \{1, S\}}$. Here, p is the finally selected time index from step (4).

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

191 (6-1) Reproduction: Select one additional time index using steps (1) through (4) and
192 denote this index as p^* . Obtain the corresponding precipitation occurrence
193 values, $\mathbf{x}_{p^*+1} = [x_{p^*+1}^s]_{s \in \{1, \dots, S\}}$. The subsequent two GA operators employ the two
194 selected vectors, \mathbf{x}_{p+1} and \mathbf{x}_{p^*+1} . This reproduction process is a mating process
195 by finding another individual that has similar characteristics to the current one
196 \mathbf{x}_{p+1} . With this procedure, a vector to similar the current vector will be mated
197 and will produce a new descendant.

198 (6-2) Crossover: Replace each element x_{p+1}^s with $x_{p^*+1}^s$ at probability P_{cr} , i.e.,

199
$$X_{c+1}^s = \begin{cases} x_{p^*+1}^s & \text{if } \varepsilon < P_{cr} \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (13)$$

200 where ε is a uniform random number between 0 and 1. From this crossover, a
201 new occurrence vector whose elements are similar to the historical ones is generated.

202 (6-3) Mutation: Replace each element (i.e., each station, $s=1, \dots, S$) with one selected
203 from all the observations of this element for $i=1, \dots, n$ with probability P_m , i.e.,

204
$$X_{c+1}^s = \begin{cases} x_{\xi+1}^s & \text{if } \varepsilon < P_m \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (14)$$

205 where $x_{\xi+1}^s$ is selected from $[x_i^s]_{i \in \{1, \dots, n\}}$ with equal probability for $i=1, \dots, n$
206 and ε is a uniform random number between 0 and 1. This mutation procedure
207 allows to generate a multisite occurrence combination that is totally different
208 from the historical records. Without this procedure, always similar multisite
209 occurrences to historical combinations are generated, which is not feasible for
210 a simulation purpose.

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

212 The selection of the number of nearest neighbors (k) has been investigated by Lall and
213 Sharma (1996) and Lee and Ouarda (2011). A simple selection method was applied in the current
214 study as $k = \sqrt{n}$. For hydrometeorological stochastic simulation, this heuristic approach of k
215 selection has been employed (Lall and Sharma, 1996; Lee and Ouarda, 2012; Lee et al., 2010b;
216 Prairie et al., 2006; Rajagopalan and Lall, 1999). One can use generalized cross-validation (GCV)
217 as shown in Sharma and Lall (1996) and Lee and Ouarda 2011 by treating this simulation as a
218 prediction problem. However, the current multisite occurrence simulation does not necessarily
219 require accurate value prediction and not much difference on simulation using the simple heuristic
220 approach is reported. Also, this heuristic approach of k selection has been popularly employed for
221 hydrometeorological stochastic simulations (Lall and Sharma, 1996; Lee and Ouarda, 2012; Lee
222 et al., 2010b; Prairie et al., 2006; Rajagopalan and Lall, 1999).

223 In Appendix A, an example of the DKNNR simulation procedure is explained in detail.

224 **3.2. Adaptation to climate change**

225 The capability of model to take climate change into account is critical. For example, the
226 marginal distributions and transition probabilities in Eqs. (5) and (3) can change in future climate

227 scenarios. It is known that nonparametric simulation models have a difficulty to adapt to climate
 228 change, since the models employ in general the current observation sequences. However, the
 229 proposed model in the current study possesses the capability to adapt to the variations of
 230 probabilities by tuning the crossover and mutation probabilities in P_{cr} (13) and P_m (14), adding
 231 the condition when applied.

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

$$234 \quad X_{c+1}^s = \begin{cases} x_{p^{*+1}}^s & \text{if } \varepsilon < P_{cr} \text{ \& } x_{p^{*+1}}^s = 1 \text{ \& } X_c^s = 1 \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (15)$$

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

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

$$237 \quad X_{c+1}^s = \begin{cases} x_{\xi+1}^s & \text{if } \varepsilon < P_m \text{ and } x_{\xi+1}^s = 1 \\ x_{p+1}^s & \text{otherwise} \end{cases} \quad (16)$$

238 This adjustment of adding the condition $x_{\xi+1}^s = 1$ can increase the marginal distribution as much as
 239 $P_m \times P_1$. This has been tested in the case study.

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

241 For testing the occurrence model, 12 weather stations were selected from Yeongnam province
 242 which is located in the southeastern part of South Korea, as shown in Figure 1. Information on
 243 longitude and latitude (fourth and fifth columns) as well as order index and the identification

244 number (first and second columns) of these stations operated by Korea Meteorological
245 Administration with the area name (third column) is shown at Table 1.

246 Figure 1 illustrates the locations of the selected weather stations. All the stations are inside
247 Yeongnam province which consists of two different regions as north and south Gyeongsang as
248 well as the self-governing cities of Busan, Daegu, and Ulsan. Most of the Yeongnam region is
249 drained to Nakdong River. To validate the proposed model appropriately, tested sites must be
250 highly correlated with each other as well as significant temporal relation. The employed stations
251 inside the Yeongnam area cover one of the most important watersheds, the Nakdong River basin,
252 where the Nakdong River passes through the entire basin and its hydrological assessments for
253 agriculture and climate change have particular values in water resources management such as
254 floods and droughts.

255 It is important to analyze the impact of weather conditions for planning agricultural
256 operations and water resources management especially during the summer season, because around
257 50-60 percent of the annual precipitation occurs during the summer season from June to
258 September. The length of daily precipitation data record ranges from 1976 to 2015 and the summer
259 season record was employed since a large number of rainy days occurs during summer and it is
260 important to preserve these characteristics. Also, the whole year dataset was tested and other
261 seasons were further applied but the correlation coefficient was relatively high and its correlation
262 matrix estimated was not a positive semi-definite matrix for the MONR model.

263 **5. Application**

264 To analyze the performance of the proposed DKNNR model, the occurrence of precipitation
265 was simulated. The DKNNR simulation was compared with that of the MONR model. For each

266 model, 100 series of daily occurrence with the same record length were simulated. The key
 267 statistics of observed data and each generated series, such as transition probabilities (P_{11} , P_{01} , and
 268 P_1) and cross-correlation (see Eq.(10)), were determined. The MONR model underestimated the
 269 lag-1 cross-correlation, as indicated by Wilks (1998). In the current study, this statistic was
 270 analyzed, since a synoptic scale weather system could result in lagged cross-correlation (Wilks,
 271 1998). It was formulated as

$$272 \quad \rho_1(q, s) = \text{corr}[X_{t-1}^q, X_t^s] \quad (17)$$

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

$$275 \quad RMSE = \left(\frac{1}{N} \sum_{m=1}^N (\Gamma_m^G - \Gamma^h)^2 \right)^{1/2} \quad (18)$$

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

280 The 100 simulated statistic values were illustrated with boxplots to show their variability as
 281 shown in Figure 4 - Figure 6. The box of boxplot represents the interquartile range (IQR) ranging
 282 25 percentile to 75 percentile. The whiskers extend to up and down $1.5 \times \text{IQR}$. Data beyond the
 283 whiskers ($1.5 \times \text{IQR}$) are indicated by a plus sign (+). The horizontal line inside the box represents
 284 the median of the data. The statistics of the observed data are denoted by a cross (x). The closer a

285 cross is to the horizontal line inside the box, the better the simulated data from a model reproduces
286 the statistical characteristics of the observed data.

287 The roles of crossover probability P_{cr} (Eq. (13)) and mutation probability P_m (Eq.(14)) were
288 studied by Lee et al. (2010b). In the current study, we further tested to select an appropriate
289 parameter set of these two parameters with the simulated data from the DKNNR model and the
290 record length of 100,000. RMSE (Eq. (18)) of the three transition and limiting probabilities (P_{11} ,
291 P_{01} , and P_1) between the simulated data and the observed was used, since those probabilities are
292 key statistics that the simulated data must be met with the observed and no parameterization on
293 these probabilities has been made for the current DKNNR model. The results are shown in Figure
294 2 and Figure 3 for P_{cr} and P_m , respectively. For P_{cr} in Figure 2, the probability of 0.02 shows the
295 smallest RMSE in all transition and limiting probabilities. The RMSE of P_m in Figure 3 shows a
296 slight fluctuation along with P_m . However, all three probabilities (P_{11} , P_{01} , and P_1) have relatively
297 small RMSEs in $P_m = 0.003$. Therefore, the parameter set 0.02 and 0.003 is chosen for P_{cr} and P_m ,
298 respectively, and employed in the current study.

299 **6. Results**

300 **6.1. Occurrence and transition probabilities**

301 The data simulated from the proposed DKNNR model and the existing MONR model were
302 analyzed. The estimated transition probabilities (P_{11} and P_{01} in Eq. (3)) as well as the occurrence
303 probability (P_1 in Eq. (5)) are shown in Table 2 and Figure 4 - Figure 6 for the observed data and
304 the data generated from the DKNNR and MONR models. In Table 2, the observed statistic shows
305 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
306 P_1 and site 12 shows the highest P_{11} .

307 As shown in Figure 4, the probability P_{11} of the observed data shows that sites 6, 7, 8, and 9
308 located in the northern part of the region exhibited lower consistency (i.e. consecutive rainy days)
309 than did the other sites, while sites 5 and 12 had higher probability of P_{11} than did other sites. Both
310 models preserved well the observed P_{11} statistic. It seems that the MONR model had a slightly
311 better performance since this statistic is parameterized in the model as shown in the section 2.2.
312 Note that the MONR model employed the transition probabilities in simulating rainfall occurrence,
313 while DKNNR model did not. The occurrence probability P_1 can be described with the
314 combination of transition probabilities as in Eq. (5). Even though the transition probabilities were
315 not employed in simulating rainfall occurrence, the DKNNR model preserved this statistic fairly
316 well.

317 As shown in Figure 5, the P_{01} probability showed a slightly different behavior such that sites
318 1, 2, and 3 located in the middle part of the Yeongnam province showed a higher probability than
319 did other sites. A slight underestimation was seen for sites 2 and 11 but it was not critical, since its
320 observed value with a cross mark was close to the upper IQR representing 75 percentile.

321 The behavior of P_1 was found to be the same as that of the P_{11} probability. It can be seen in
322 Figure 6 that no significant underestimation is seen for the DKNNR model (top panel). The P_1
323 statistic was fairly preserved by both DKNNR and MONR models. Note that the MONR model
324 parameterized the P_1 statistic through the transition probabilities as in Eq. (5), while DKNNR
325 model did not. Although the DKNNR model used almost no parameters for simulation, the P_1
326 statistic was preserved fairly well.

327 **6.2. Cross-correlation**

328 Cross-correlation is a measure of relationship between sites. The preservation of cross-
329 correlation is important for the simulation of precipitation occurrence and is required in the
330 regional analysis for water resources management or agricultural applications. Furthermore,
331 lagged cross-correlation is also essential as much as is cross-correlation (i.e. contemporaneous
332 correlation). For example, the amount of streamflow for a watershed from a certain precipitation
333 event is highly related with lagged cross-correlation. It is accepted that precipitation is not
334 significantly correlated with that for more than one day. Therefore, only lag-1 cross-correlation
335 was analyzed in the current study.

336 The cross-correlation of observed data is shown in

337 Table 3. High cross-correlation among grouped sites, such as sites 6, 7, and 8 (northern part)
338 and sites 3, 4, and 5 as well as 12 (southeast coastal area, 0.68-0.87), was found. As expected, sites
339 5 and 12 had the highest cross-correlation (0.87) due to the proximity. The northern sites and
340 coastal sites showed low cross-correlation. This observed cross-correlation was well preserved in
341 the data generated from both DKNNR and MONR models, as shown in Figure 7 as well as Table
342 4 and Table 5. However, consistently slight but significant underestimation of cross-correlation
343 was seen for the data generated by the MONR model (see the bottom panel of Figure 7). Note that
344 the errorbars are extended to upper and lower lines of the circles to $1.95 \times$ standard deviation. The
345 difference of RMSE in Table 6 showed this characteristic, as most of the values were positive, to
346 be indicating that the proposed DKNNR model performed better for cross-correlation.

347 The lag-1 cross-correlation of observed data, as shown in Table 7, ranged from 0.22-0.35.
348 The lag-1 cross-correlation for the same site (i.e. $\rho_1(q, s)$, $q=s$) was autocorrelation and was highly
349 related with P_{01} and P_{11} as in Eq. (6). All the lag-1 cross-correlations exhibited similar magnitudes
350 even for autocorrelation. This implies that the lag-1 cross-correlation among the selected sites was
351 as strong as the autocorrelation and as much as the transition probabilities P_{01} and P_{11} , thereof.

352 The observed lag-1 cross-correlations were well preserved in the data generated by the
353 DKNNR model, as shown in the top panel of Figure 8, while the MONR model showed significant
354 underestimation, as seen in the bottom panel of Figure 8. The difference of RMSE shown in Table
355 8 reflects this behavior. In the bottom panel of Figure 8, some of the lag-1 cross-correlations were
356 well preserved, that was aligned with the base line. From Table 8, the MONR model reproduced
357 the autocorrelations well with the shaded values. It is because the lag-1 autocorrelation was
358 indirectly parameterized with the transition probabilities of P_{11} and P_{01} as in Eq. (6). Other than

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

362 We further tested the performance measurements of MAE and Bias. The estimates showed
363 that MAE had no difference from RMSE. In addition, Bias of the lag-1 correlation presents
364 significant negative values implying its underestimation for the simulated data of the MONR
365 model as shown in Table 9, while Table 10 of the DKNNR model shows a much smaller bias.

366 Also, the whole year data instead of the summer season data was tested for model fitting.
367 Note that all the results presented above were with the summer season data (June-September) as
368 mentioned in section 4 on the data description. The lag-1 cross-correlation is shown in Figure 9
369 which indicates that the same characteristic was observed as for the summer season, such that the
370 proposed DKNNR model preserved better the lagged cross-correlation than did the existing
371 MONR model. Other statistics, such as correlation matrix and transition probabilities, exhibited
372 the same results (not shown). Also, other seasons were tried but the estimated correlation matrix
373 was not a positive semi-definite matrix and its inverse cannot be made for multivariate normal
374 distribution in the MONR model. It was because the selected stations were close to each other
375 (around 50-100 km) and produced high cross-correlation, especially in the occurrence during dry
376 seasons. Special remedy for the existing MONR model should be applied, such as decreasing
377 cross-correlation by force, but further remedy was not applied in the current study since it was not
378 within the current scope and focus.

379 **6.3. Adaptation to climate change**

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

385 In Figure 10, the changes of transition and marginal probabilities are shown along with
386 increasing the crossover probability P_{cr} from 0.01 to 0.2 with the condition that that the candidate
387 value is one and the previous value is also one as in Eq. (15) for the selected 5 stations among the
388 12 stations (from station 1 to station 5, see Table 1 for details). The stations were limited in this
389 analysis due to computational time. At each case 100 series were simulated. The average value of
390 the simulated statistics is presented in the figure. It is obvious that the transition probability P_{11}
391 increased as intended along with the increase of P_{cr} . As expected from Eq. (5), P_1 presents that
392 the change of P_1 is highly related to P_{11} . However, the probability P_{01} fluctuated along with the
393 increase of P_{cr} . Elaborate work to adjust all the probabilities is however required.

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

400 **7. Conclusions**

401 In the current study, a nonparametric simulation model, based on discrete KNNR and
402 DKNNR, is proposed to overcome the shortcomings of the existing MONR model such as long
403 stochastic simulation for the parameter estimation and underestimation of the lagged
404 crosscorrelation between sites. Occurrence and transition probabilities and cross-correlation as
405 well as lag-1 cross-correlation are estimated for both models. Better preservation of cross-
406 correlation and lag-1 cross-correlation with the DKNNR model than the MONR model is observed.
407 For some cases (i.e., the whole year data and other seasons than the summer season), the estimated
408 cross-correlation matrix is not a positive semi-definite matrix so the multivariate normal
409 simulation is not applicable for the MONR model, because the tested sites are close to each other
410 with high cross-correlation.

411 Results of this study indicate that the proposed DKNNR model reproduces the occurrence
412 and transition probabilities fairly well and preserves the cross-correlations better than the existing
413 MONR model. Furthermore, not much effort is required to estimate the parameters in the DKNNR
414 model, while the MONR model requires a long stochastic simulation just to estimate each
415 parameter. Thus, the proposed DKNNR model can be a good alternative for simulating multisite
416 precipitation occurrence.

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

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

425 **Code and Data Availability**

426 DKNNR code is written in Matlab and is available as a supplement.

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

429 **Acknowledgment**

430 This work was supported by the National Research Foundation of Korea (NRF) grant (NRF-
431 2018R1A2B6001799) funded by the Korean Government (MEST).

432 **Appendix A: Example of DKNNR**

433 In this appendix, one example of DKNNR simulation is presented with observed dataset in
434 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
435 presents the observed precipitation (unit: mm). Its occurrence data is presented in the bottom part
436 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
437 Table A 2. The number of nearest neighbors $k = \sqrt{n} = \sqrt{16} = 4$ and the parameters for GA (i.e. P_c
438 and P_m) are 0.1 and 0.01, respectively. Simulation can be made as follows:

439 (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,
440 for $i=1$,

441
$$D_1 = \sum_{s=1}^S |X_c^s - x_1^s| = |0-1| + |1-1| + \dots + |0-1| = 6$$

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

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

448 (3) Simulate a uniform random number (u) between 0 and 1. Say $u=0.321$. This value must
 449 be compared with the cumulative weighted probabilities in the last column of Table A 3
 450 as [0 0.48 0.72 0.88 1.0]. The corresponding day index is assigned as to where the
 451 simulated uniform number falls in the cumulative weighted probabilities, here [0 0.48].
 452 Therefore, the selected day, p , is 14. The occurrences of the following day $p+1=15$ for 12
 453 stations are selected as in the second row of Table A 4.

454 (4) For GA mixture, another set must be chosen as in step (3). Say $u=0.561$, which falls in
 455 [0.48 0.72]. The second one should be selected. However, there are a number of days with
 456 the same distances. Specifically, six days have the same distances with $D_i=4$. In this case,
 457 one among all six days is selected with equal probability. Assume that $p=4$ is selected and
 458 the following occurrences are selected, as shown in the third row of Table A 4.

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

460 (5-1) Crossover: For each station, a uniform random number (ε) is generated and
 461 compared with $P_c=0.1$ here. Say $\varepsilon =0.345$, then skip since $\varepsilon =0.345 > P_c=0.1$. For

462 $s=6$, assume the generated random number, $\varepsilon (=0.051) < P_c (=0.1)$ and then switch
463 the 6th station value of Set 1 into the value of Set 2 (see Table A 4). The occurrence
464 state of X_{c+1}^s turns into 1 from 0 as shown in the fourth row of Table A 4 as well as
465 station 8.

466 (5-2) Mutation: For each station, a uniform random number (ε) is generated and compared
467 with $P_m=0.01$. For $s=12$, assume $\varepsilon =0.009 < P_m=0.01$ and switch the 12th station
468 value of Set 1 with the one selected among all the observed 12th station values with
469 equal probability (here the last column, $s=12$, of the bottom part of Table A 1, [1 1
470 0 0 ... 1]). The occurrence state of X_{c+1}^{12} turns into 0 from 1 as shown in the fourth
471 column of Table A 4.

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

473

474 **References**

475
476 Apipattanavis, S., Podesta, G., Rajagopalan, B., and Katz, R. W.: A semiparametric
477 multivariate and multisite weather generator, *Water Resources Research*, 43, Artn W11401, 2007.

478 Beersma, J. J. and Buishand, A. T.: Multi-site simulation of daily precipitation and
479 temperature conditional on the atmospheric circulation, *Climate Research*, 25, 121-133, 2003.

480 Brandsma, T. and Buishand, T. A.: Simulation of extreme precipitation in the Rhine basin
481 by nearest-neighbour resampling, *Hydrology and Earth System Sciences*, 2, 195-209, 1998.

482 Buishand, T. A. and Brandsma, T.: Multisite simulation of daily precipitation and
483 temperature in the Rhine basin by nearest-neighbor resampling, *Water Resources Research*, 37,
484 2761-2776, 2001.

485 Chau, K. W.: Use of meta-heuristic techniques in rainfall-runoffmodelling, *Water*
486 (Switzerland), 9, 2017.

487 Fotovatikhah, F., Herrera, M., Shamshirband, S., Chau, K. W., Ardabili, S. F., and Piran,
488 M. J.: Survey of computational intelligence as basis to big flood management: Challenges, research
489 directions and future work, *Engineering Applications of Computational Fluid Mechanics*, 12, 411-
490 437, 2018.

491 Frost, A. J., Charles, S. P., Timbal, B., Chiew, F. H. S., Mehrotra, R., Nguyen, K. C.,
492 Chandler, R. E., McGregor, J. L., Fu, G., Kirono, D. G. C., Fernandez, E., and Kent, D. M.: A
493 comparison of multi-site daily rainfall downscaling techniques under Australian conditions,
494 *Journal of Hydrology*, 408, 1-18, 2011.

495 Hughes, J. P., Guttorp, P., and Charles, S. P.: A non-homogeneous hidden Markov model
496 for precipitation occurrence, *Journal of the Royal Statistical Society. Series C: Applied Statistics*,
497 48, 15-30, 1999.

498 Jeong, D. I., St-Hilaire, A., Ouarda, T. B. M. J., and Gachon, P.: A multi-site statistical
499 downscaling model for daily precipitation using global scale GCM precipitation outputs,
500 *International Journal of Climatology*, 33, 2431-2447, 2013.

501 Jeong, D. I., St-Hilaire, A., Ouarda, T. B. M. J., and Gachon, P.: Multisite statistical
502 downscaling model for daily precipitation combined by multivariate multiple linear regression and
503 stochastic weather generator, *Climatic Change*, 114, 567-591, 2012.

504 Katz, R. W.: Precipitation as a Chain-Dependent Process, *Journal of Applied Meteorology*,
505 16, 671-676, 1977.

506 Katz, R. W. and Zheng, X.: Mixture model for overdispersion of precipitation, *Journal of*
507 *Climate*, 12, 2528-2537, 1999.

508 Lall, U. and Sharma, A.: A nearest neighbor bootstrap for resampling hydrologic time
509 series, *Water Resources Research*, 32, 679-693, 1996.

510 Lee, T.: Multisite stochastic simulation of daily precipitation from copula modeling with
511 a gamma marginal distribution, *Theoretical and Applied Climatology*, doi: 10.1007/s00704-017-
512 2147-0, 2017. 1-10, 2017.

513 Lee, T.: Stochastic simulation of precipitation data for preserving key statistics in their
514 original domain and application to climate change analysis, *Theoretical and Applied Climatology*,
515 124, 91-102, 2016.

516 Lee, T. and Ouarda, T. B. M. J.: Identification of model order and number of neighbors
517 for k-nearest neighbor resampling, *Journal of Hydrology*, 404, 136-145, 2011.

518 Lee, T. and Ouarda, T. B. M. J.: Stochastic simulation of nonstationary oscillation hydro-
519 climatic processes using empirical mode decomposition, *Water Resources Research*, 48, 1-15,
520 2012.

521 Lee, T., Ouarda, T. B. M. J., and Jeong, C.: Nonparametric multivariate weather generator
522 and an extreme value theory for bandwidth selection, *Journal of Hydrology*, 452-453, 161-171,
523 2012.

524 Lee, T., Ouarda, T. B. M. J., and Yoon, S.: KNN-based local linear regression for the
525 analysis and simulation of low flow extremes under climatic influence, *Climate Dynamics*, doi:
526 10.1007/s00382-017-3525-0, 2017. 1-19, 2017.

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

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

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

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

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

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

540 St-Hilaire, A., Ouarda, T. B. M. J., Bargaoui, Z., Daigle, A., and Bilodeau, L.: Daily river
541 water temperature forecast model with a k-nearest neighbour approach, *Hydrological Processes*,
542 26, 1302-1310, 2012.

543 Taormina, R., Chau, K. W., and Sivakumar, B.: Neural network river forecasting through
544 baseflow separation and binary-coded swarm optimization, *Journal of Hydrology*, 529, 1788-1797,
545 2015.

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

548 Wang, W. C., Xu, D. M., Chau, K. W., and Chen, S.: Improved annual rainfall-runoff
549 forecasting using PSO-SVM model based on EEMD, *Journal of Hydroinformatics*, 15, 1377-1390,
550 2013.

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

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

555 Wilks, D. S.: Multisite generalization of a daily stochastic precipitation generation model,
556 *Journal of Hydrology*, 210, 178-191, 1998.

557 Wilks, D. S. and Wilby, R. L.: The weather generation game: a review of stochastic
558 weather models, *Progress in Physical Geography*, 23, 329-357, 1999.

559 Zheng, X. and Katz, R. W.: Simulation of spatial dependence in daily rainfall using
560 multisite generators, *Water Resources Research*, 44, 2008.

561

562

563

564 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

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

567

568

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

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

573
 574
 575

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

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	1.00	0.70	0.70	0.64	0.58	0.70	0.65	0.63	0.75	0.64	0.66	0.59
S2	0.70	1.00	0.67	0.64	0.61	0.64	0.70	0.72	0.79	0.72	0.74	0.62
S3	0.70	0.67	1.00	0.75	0.68	0.61	0.57	0.57	0.68	0.67	0.74	0.70
S4	0.64	0.64	0.75	1.00	0.79	0.56	0.56	0.55	0.65	0.66	0.73	0.82
S5	0.58	0.61	0.68	0.79	1.00	0.51	0.54	0.55	0.61	0.65	0.70	0.87
S6	0.70	0.64	0.61	0.56	0.51	1.00	0.69	0.65	0.68	0.59	0.59	0.54
S7	0.65	0.70	0.57	0.56	0.54	0.69	1.00	0.78	0.71	0.65	0.63	0.55
S8	0.63	0.72	0.57	0.55	0.55	0.65	0.78	1.00	0.71	0.68	0.65	0.56
S9	0.75	0.79	0.68	0.65	0.61	0.68	0.71	0.71	1.00	0.68	0.71	0.62
S10	0.64	0.72	0.67	0.66	0.65	0.59	0.65	0.68	0.68	1.00	0.77	0.66
S11	0.66	0.74	0.74	0.73	0.70	0.59	0.63	0.65	0.71	0.77	1.00	0.70
S12	0.59	0.62	0.70	0.82	0.87	0.54	0.55	0.56	0.62	0.66	0.70	1.00

578
 579
 580

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

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	1.00	0.68	0.69	0.64	0.60	0.69	0.64	0.62	0.73	0.63	0.65	0.61
S2	0.68	1.00	0.67	0.63	0.62	0.63	0.68	0.72	0.77	0.74	0.73	0.63
S3	0.69	0.67	1.00	0.74	0.69	0.60	0.58	0.59	0.66	0.68	0.74	0.70
S4	0.64	0.63	0.74	1.00	0.79	0.55	0.55	0.56	0.62	0.65	0.71	0.81
S5	0.60	0.62	0.69	0.79	1.00	0.51	0.56	0.58	0.60	0.66	0.70	0.86
S6	0.69	0.63	0.60	0.55	0.51	1.00	0.68	0.64	0.65	0.59	0.58	0.53
S7	0.64	0.68	0.58	0.55	0.56	0.68	1.00	0.78	0.69	0.65	0.63	0.56
S8	0.62	0.72	0.59	0.56	0.58	0.64	0.78	1.00	0.70	0.69	0.67	0.58
S9	0.73	0.77	0.66	0.62	0.60	0.65	0.69	0.70	1.00	0.67	0.69	0.60
S10	0.63	0.74	0.68	0.65	0.66	0.59	0.65	0.69	0.67	1.00	0.77	0.67
S11	0.65	0.73	0.74	0.71	0.70	0.58	0.63	0.67	0.69	0.77	1.00	0.71
S12	0.61	0.63	0.70	0.81	0.86	0.53	0.56	0.58	0.60	0.67	0.71	1.00

583
 584
 585

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

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	1.00	0.63	0.67	0.58	0.54	0.66	0.62	0.60	0.68	0.55	0.62	0.53
S2	0.63	1.00	0.61	0.60	0.57	0.59	0.68	0.68	0.75	0.66	0.72	0.58
S3	0.67	0.61	1.00	0.71	0.67	0.57	0.56	0.53	0.65	0.61	0.71	0.69
S4	0.58	0.60	0.71	1.00	0.78	0.50	0.52	0.52	0.61	0.62	0.69	0.78
S5	0.54	0.57	0.67	0.78	1.00	0.48	0.51	0.53	0.57	0.62	0.67	0.81
S6	0.66	0.59	0.57	0.50	0.48	1.00	0.67	0.62	0.63	0.54	0.54	0.49
S7	0.62	0.68	0.56	0.52	0.51	0.67	1.00	0.75	0.70	0.61	0.62	0.52
S8	0.60	0.68	0.53	0.52	0.53	0.62	0.75	1.00	0.66	0.64	0.61	0.52
S9	0.68	0.75	0.65	0.61	0.57	0.63	0.70	0.66	1.00	0.63	0.69	0.57
S10	0.55	0.66	0.61	0.62	0.62	0.54	0.61	0.64	0.63	1.00	0.72	0.61
S11	0.62	0.72	0.71	0.69	0.67	0.54	0.62	0.61	0.69	0.72	1.00	0.66
S12	0.53	0.58	0.69	0.78	0.81	0.49	0.52	0.52	0.57	0.61	0.66	1.00

588
 589
 590
 591

592 Table 6. The difference of RMSE of cross-correlation between MONR and DKNNR. Note that
 593 the positive value indicates that the DKNNR model better performs in preserving the cross-
 594 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	0.014	0.004	0.013	0.012	0.012	0.008	0.005	0.024	0.031	0.011	0.035
S2	0.014	0.000	0.023	0.013	0.021	0.009	0.010	0.013	0.018	0.027	0.011	0.020
S3	0.004	0.023	0.000	0.015	0.004	0.014	0.003	0.022	0.009	0.028	0.011	0.004
S4	0.013	0.013	0.015	0.000	0.002	0.017	0.018	0.014	0.018	0.018	0.027	0.024
S5	0.012	0.021	0.004	0.002	0.000	0.014	0.021	0.014	0.015	0.013	0.015	0.012
S6	0.012	0.009	0.014	0.017	0.014	0.000	0.006	0.010	0.030	0.018	0.029	0.021
S7	0.008	0.010	0.003	0.018	0.021	0.006	0.000	0.005	0.008	0.024	0.012	0.023
S8	0.005	0.013	0.022	0.014	0.014	0.010	0.005	0.000	0.032	0.019	0.022	0.023
S9	0.024	0.018	0.009	0.018	0.015	0.030	0.008	0.032	0.000	0.019	0.005	0.027
S10	0.031	0.027	0.028	0.018	0.013	0.018	0.024	0.019	0.019	0.000	0.020	0.021
S11	0.011	0.011	0.011	0.027	0.015	0.029	0.012	0.022	0.005	0.020	0.000	0.022
S12	0.035	0.020	0.004	0.024	0.012	0.021	0.023	0.023	0.027	0.021	0.022	0.000

595 Note that no negative value can be found implying that the DKNNR model preserves the
 596 crosscorrelation better than the MONR model.

597

598

599

600

601

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

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

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

605

606

607 Table 8. The difference of RMSE of lag-1 cross-correlation between MONR and DKNNR. Note
 608 that a positive value indicates that the DKNNR model better performs in preserving lag-1 cross-
 609 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	0.048	0.075	0.049	0.041	0.095	0.059	0.036	0.047	0.055	0.063	0.052
S2	0.070	0.000	0.079	0.057	0.046	0.104	0.068	0.047	0.066	0.058	0.073	0.047
S3	0.067	0.054	0.000	0.046	0.031	0.096	0.072	0.056	0.055	0.052	0.056	0.025
S4	0.086	0.075	0.083	0.002	0.037	0.117	0.089	0.077	0.078	0.062	0.077	0.040
S5	0.111	0.096	0.098	0.074	0.002	0.124	0.103	0.085	0.105	0.070	0.108	0.049
S6	0.039	0.024	0.060	0.038	0.043	-0.002	0.028	0.017	0.045	0.034	0.055	0.037
S7	0.055	0.045	0.077	0.061	0.062	0.084	0.000	0.023	0.051	0.052	0.071	0.064
S8	0.092	0.078	0.104	0.079	0.068	0.115	0.079	0.000	0.094	0.078	0.101	0.074
S9	0.060	0.052	0.084	0.066	0.056	0.106	0.057	0.056	0.001	0.069	0.076	0.064
S10	0.091	0.094	0.105	0.081	0.062	0.123	0.107	0.085	0.100	0.001	0.092	0.063
S11	0.064	0.061	0.071	0.057	0.033	0.109	0.084	0.063	0.062	0.043	-0.002	0.043
S12	0.121	0.099	0.096	0.077	0.036	0.130	0.101	0.086	0.107	0.082	0.109	0.003

610

611

612
613 Table 9. Bias of lag-1 cross-correlation of the generated data from the DKNNR model. Note that
614 a positive value indicates the overestimation of lag-1 cross-correlation, while a negative value
615 shows underestimation. Note that $Bias = 1/N \sum_{m=1}^N \Gamma_m^G - \Gamma^h$ and see Eq. (18) for the details of each
616 term.
617

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	0.000	0.009	0.001	0.003	0.006	-0.002	0.010	0.011	0.006	0.010	0.010	0.006
S2	0.005	0.009	0.010	0.006	0.008	0.006	0.011	0.011	0.004	0.009	0.009	0.010
S3	0.002	0.010	0.001	-0.002	0.003	0.002	0.007	0.008	0.006	0.009	0.006	0.007
S4	0.006	0.009	0.004	0.001	0.007	0.003	0.008	0.008	0.009	0.010	0.010	0.005
S5	0.004	0.005	0.000	-0.001	-0.001	0.007	0.005	0.006	0.002	0.008	0.000	-0.001
S6	-0.002	0.006	0.000	0.002	-0.001	-0.002	0.004	0.003	0.002	0.005	0.004	0.001
S7	0.004	0.008	0.003	0.003	0.001	0.004	0.002	0.006	0.007	0.007	0.007	0.002
S8	0.000	0.005	0.004	0.001	0.004	-0.003	-0.003	0.000	0.001	0.004	0.006	0.003
S9	0.005	0.007	0.006	0.003	0.006	0.004	0.010	0.007	0.004	0.007	0.006	0.007
S10	0.003	0.005	0.001	-0.001	-0.001	0.001	0.001	0.001	0.003	0.000	0.002	0.001
S11	0.010	0.010	0.008	0.004	0.008	0.009	0.009	0.009	0.010	0.010	0.011	0.008
S12	0.003	0.006	0.001	-0.001	0.004	0.003	0.008	0.008	0.005	0.005	0.002	0.001

618
619

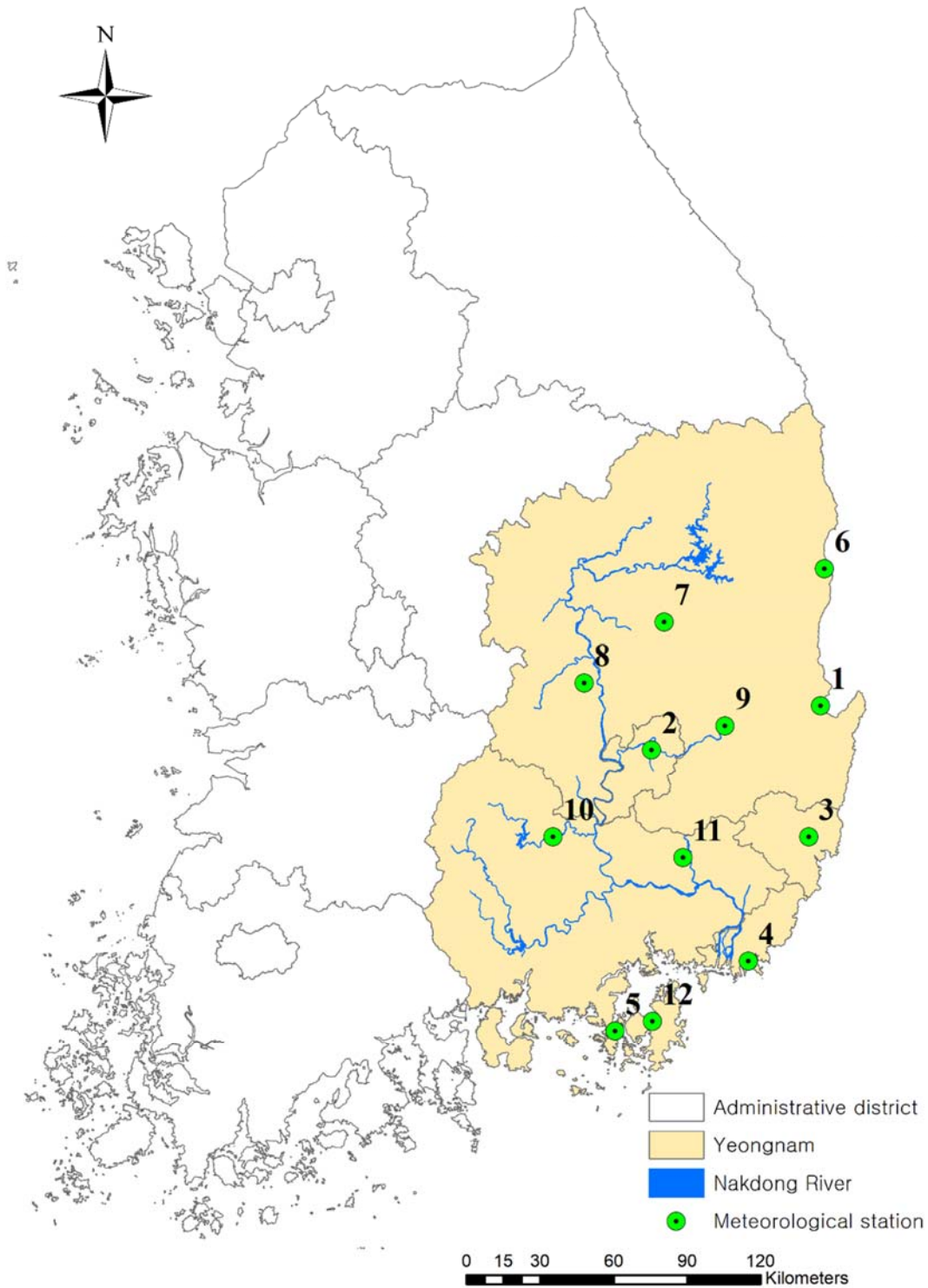
620 Table 10. Bias of lag-1 cross-correlation of the generated data from the Wilks model. Note that a
 621 positive value indicates the overestimation of lag-1 cross-correlation, while a negative value
 622 shows underestimation.

	S1	S2	S3	S4	S5	S6	S7	S8	S9	S10	S11	S12
S1	-0.001	-0.062	-0.089	-0.063	-0.055	-0.106	-0.074	-0.052	-0.060	-0.070	-0.080	-0.067
S2	-0.084	0.000	-0.096	-0.072	-0.061	-0.117	-0.083	-0.063	-0.079	-0.072	-0.089	-0.063
S3	-0.080	-0.070	0.001	-0.059	-0.043	-0.110	-0.086	-0.072	-0.069	-0.066	-0.071	-0.037
S4	-0.100	-0.090	-0.097	-0.001	-0.048	-0.129	-0.103	-0.093	-0.093	-0.077	-0.092	-0.051
S5	-0.125	-0.110	-0.111	-0.087	-0.001	-0.138	-0.117	-0.100	-0.118	-0.084	-0.121	-0.060
S6	-0.053	-0.037	-0.074	-0.051	-0.057	-0.001	-0.039	-0.030	-0.060	-0.047	-0.070	-0.049
S7	-0.068	-0.058	-0.091	-0.077	-0.077	-0.098	-0.002	-0.038	-0.065	-0.065	-0.086	-0.079
S8	-0.106	-0.091	-0.119	-0.094	-0.084	-0.128	-0.093	0.001	-0.108	-0.091	-0.116	-0.088
S9	-0.074	-0.064	-0.098	-0.080	-0.070	-0.119	-0.072	-0.070	-0.001	-0.082	-0.091	-0.078
S10	-0.105	-0.107	-0.120	-0.096	-0.075	-0.136	-0.119	-0.097	-0.113	-0.001	-0.106	-0.076
S11	-0.078	-0.074	-0.085	-0.070	-0.047	-0.123	-0.097	-0.077	-0.076	-0.056	-0.001	-0.057
S12	-0.134	-0.112	-0.108	-0.088	-0.046	-0.142	-0.116	-0.101	-0.121	-0.095	-0.122	0.000

623

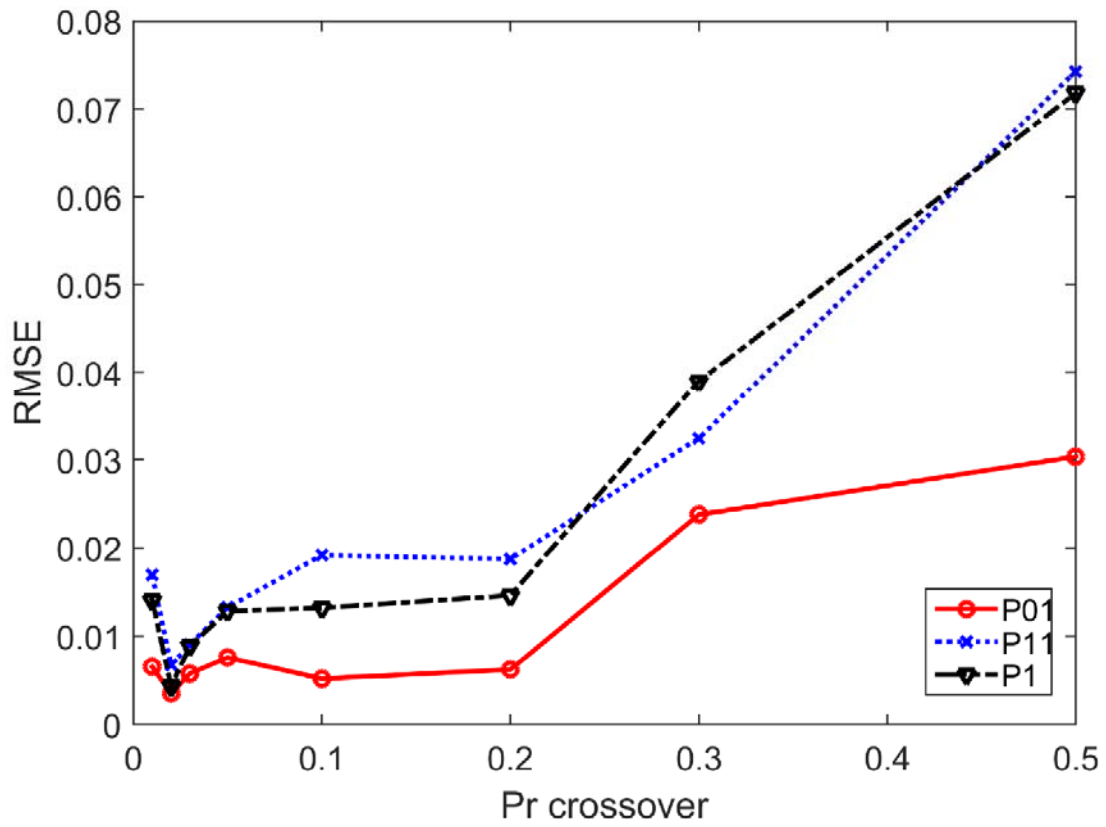
624

625



626

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

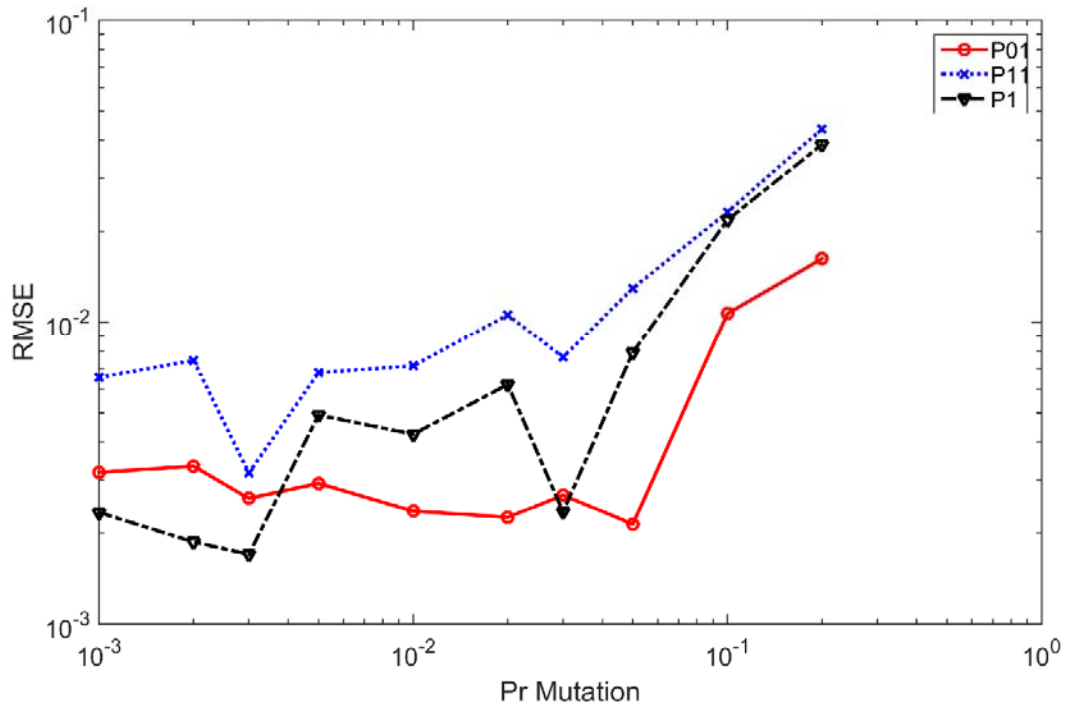


630
 631 Figure 2. Testing for different probabilities of crossover P_{cr} . RMSE is estimated for all the tested
 632 12 stations for each transition and limiting probability of the simulated data with the record
 633 length of 100,000.

634

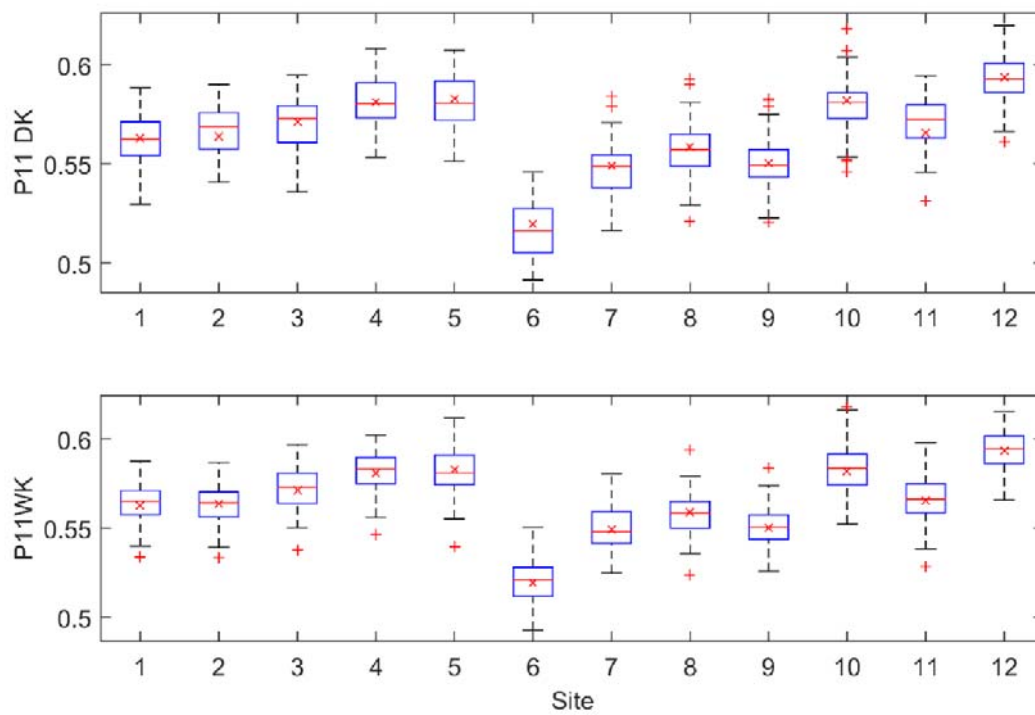
635

636
637



638
639
640
641
642
643

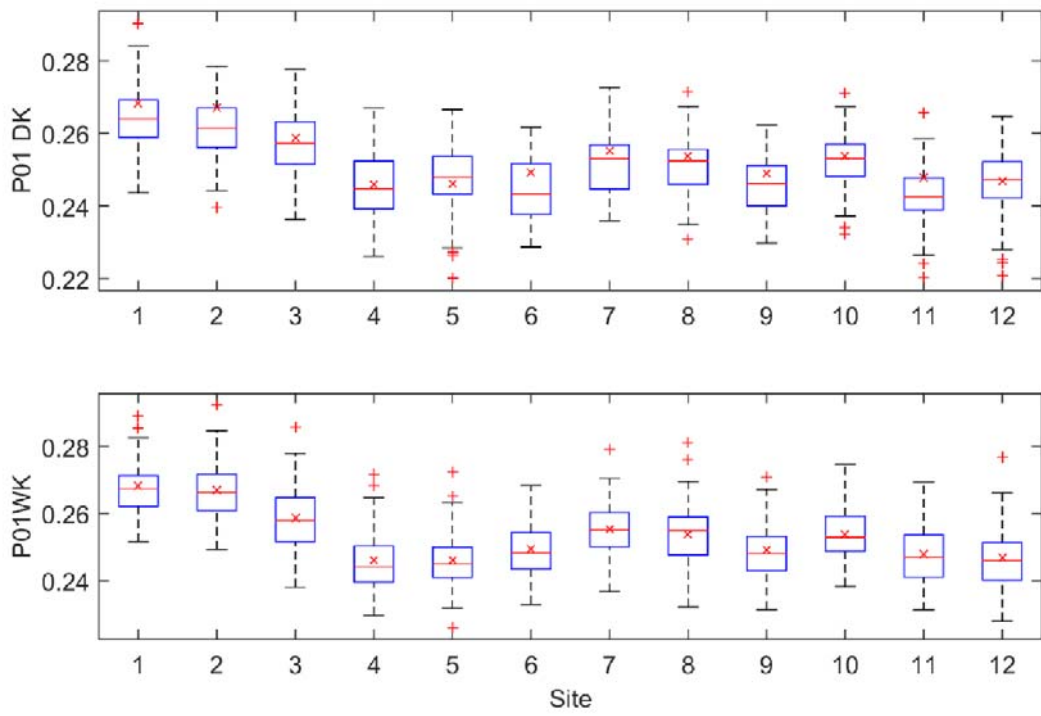
Figure 3. Testing for different probabilities of mutation P_m . RMSE is estimated for all the tested 12 stations for each transition and limiting probability of the simulated data with the record length of 100,000.



644

645 Figure 4. Boxplots of the P11 probability for the simulated data from the DKNNR model (top
 646 panel) and the MONR model (bottom panel) as well as the observed (x marker) for the 12
 647 selected weather stations from the Yeongnam province.

648



649

650 Figure 5. Boxplots of the P01 probability for the data simulated from the DKNNR model (top
 651 panel) and the MONR model (bottom panel) as well as the observed (x marker) for the 12
 652 selected weather stations from the Yeongnam province.

653

654

655

656

657

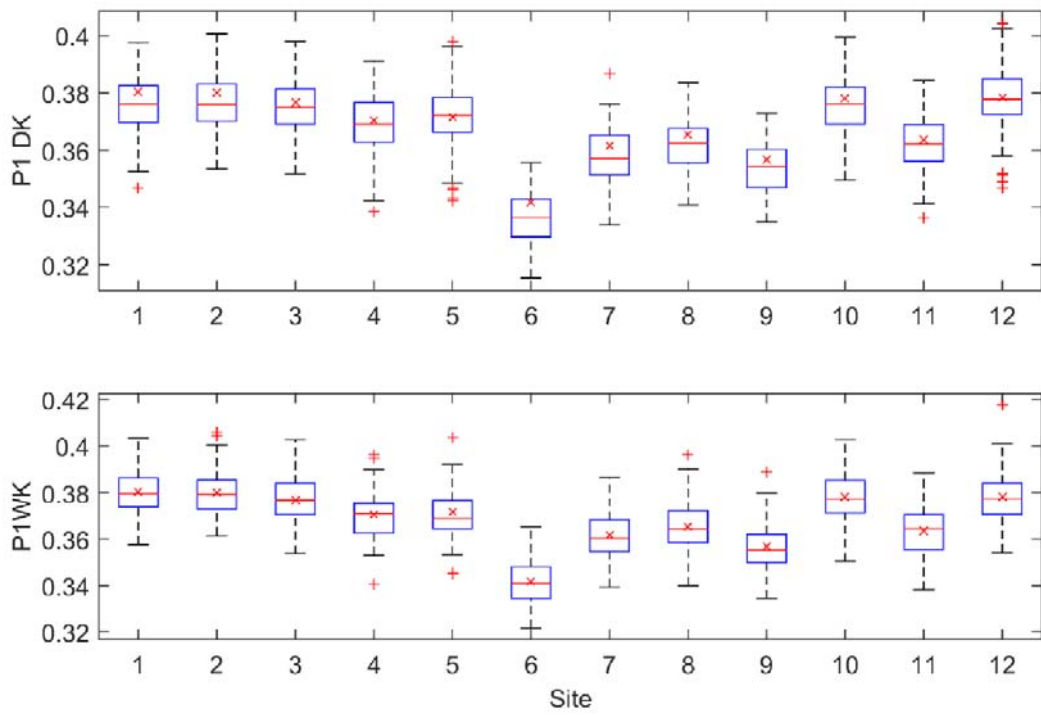
658

659

660

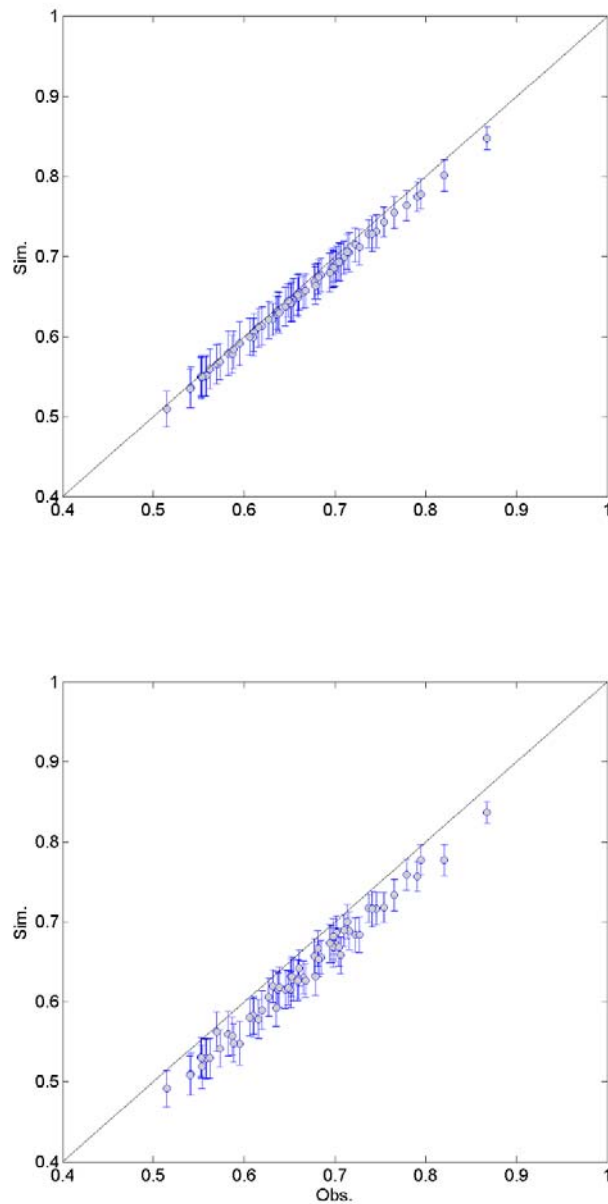
661

662



663

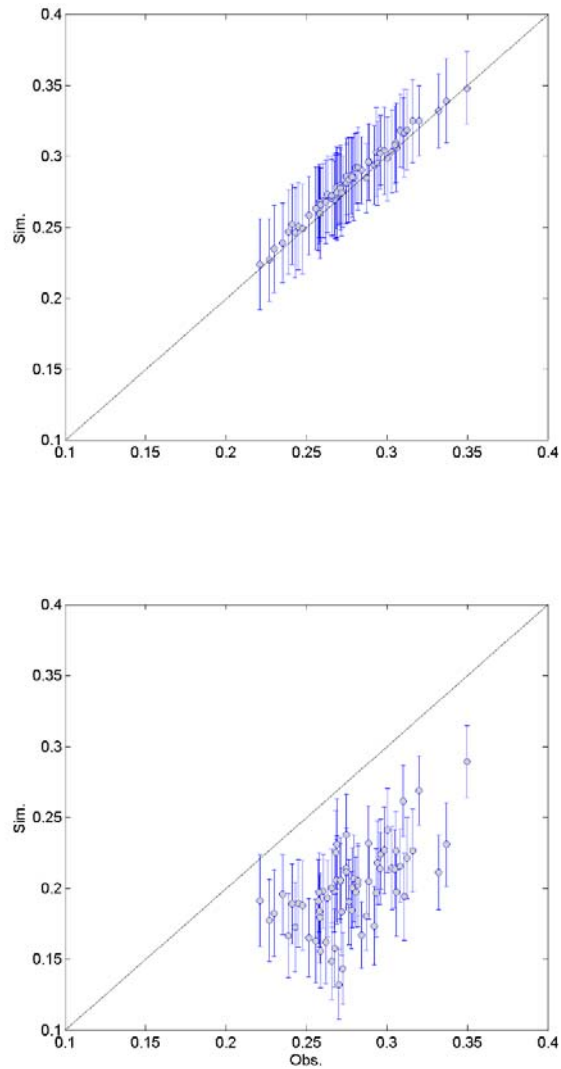
664 Figure 6. Boxplots of the P1 probability for the data simulated from the DKNNR model (top
 665 panel) and the MONR model (bottom panel) as well as the observed (x marker) for the 12
 666 selected weather stations from the Yeongnam province.



667

668 Figure 7. Scatterplot of cross-correlations between 12 weather stations for the observed data (X
 669 coordinate) and the generated data (Y coordinate) generated from the DKNNR model (top panel)
 670 and the MONR model (bottom panel). The cross-correlations from 100 generated series are
 671 averaged for the filled circle and the errorbars upper and lower extended lines indicate the range
 672 of $1.95 \times$ standard deviation.

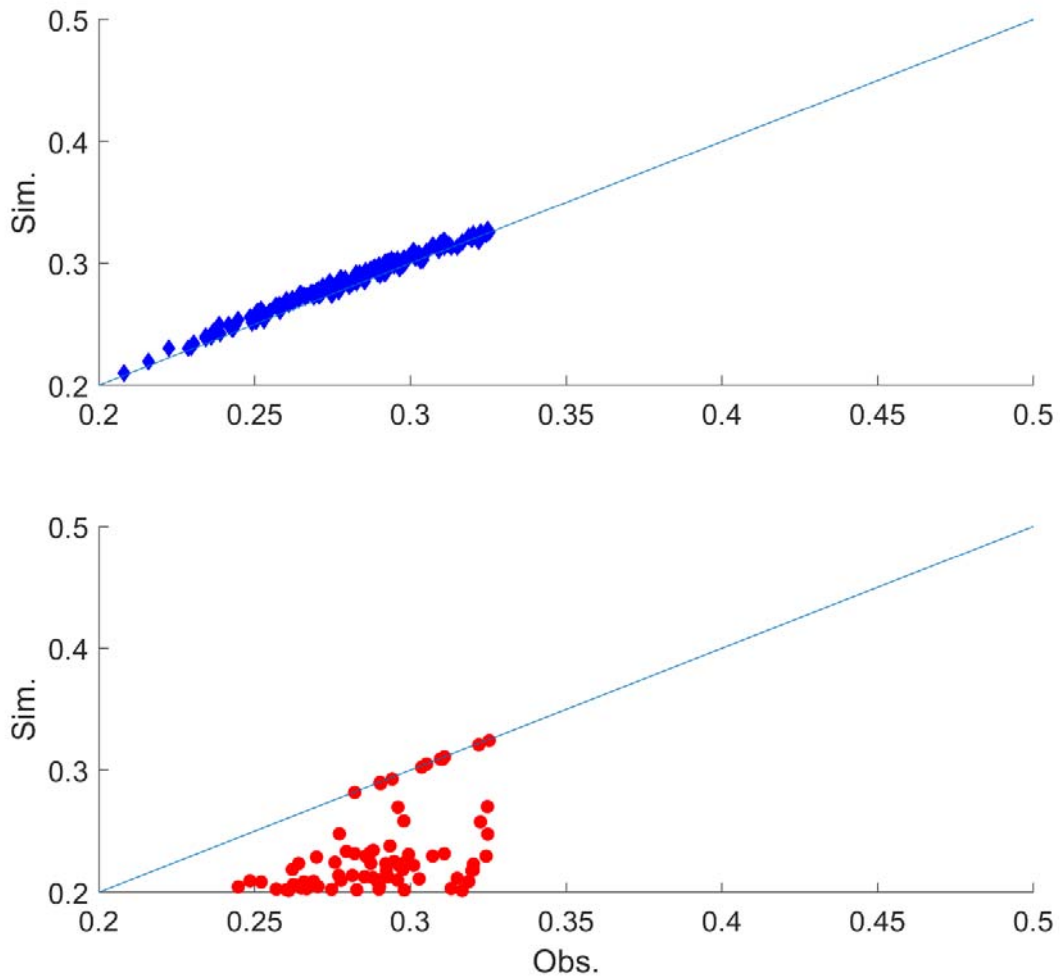
673



675

676 Figure 8. Scatterplot of lag-1 cross-correlations between 12 weather stations for the observed
 677 data (X coordinate) and the generated data (Y coordinate) generated from the DKNNR model
 678 (top panel) and the MONR model (bottom panel). The cross-correlations from 100 generated
 679 series are averaged for the filled circle and the errorbars upper and lower extended lines indicate
 680 the range of $1.95 \times$ standard deviation.

681



682

683 Figure 9. Scatterplot of lag-1 cross-correlations between 12 weather stations for the observed
684 data (X coordinate) and the generated data (Y coordinate) generated from the DKNNR model
685 (top panel) and the MONR model (bottom panel) with the whole year data not with the summer
686 season. The cross-correlations from 100 generated series are averaged.

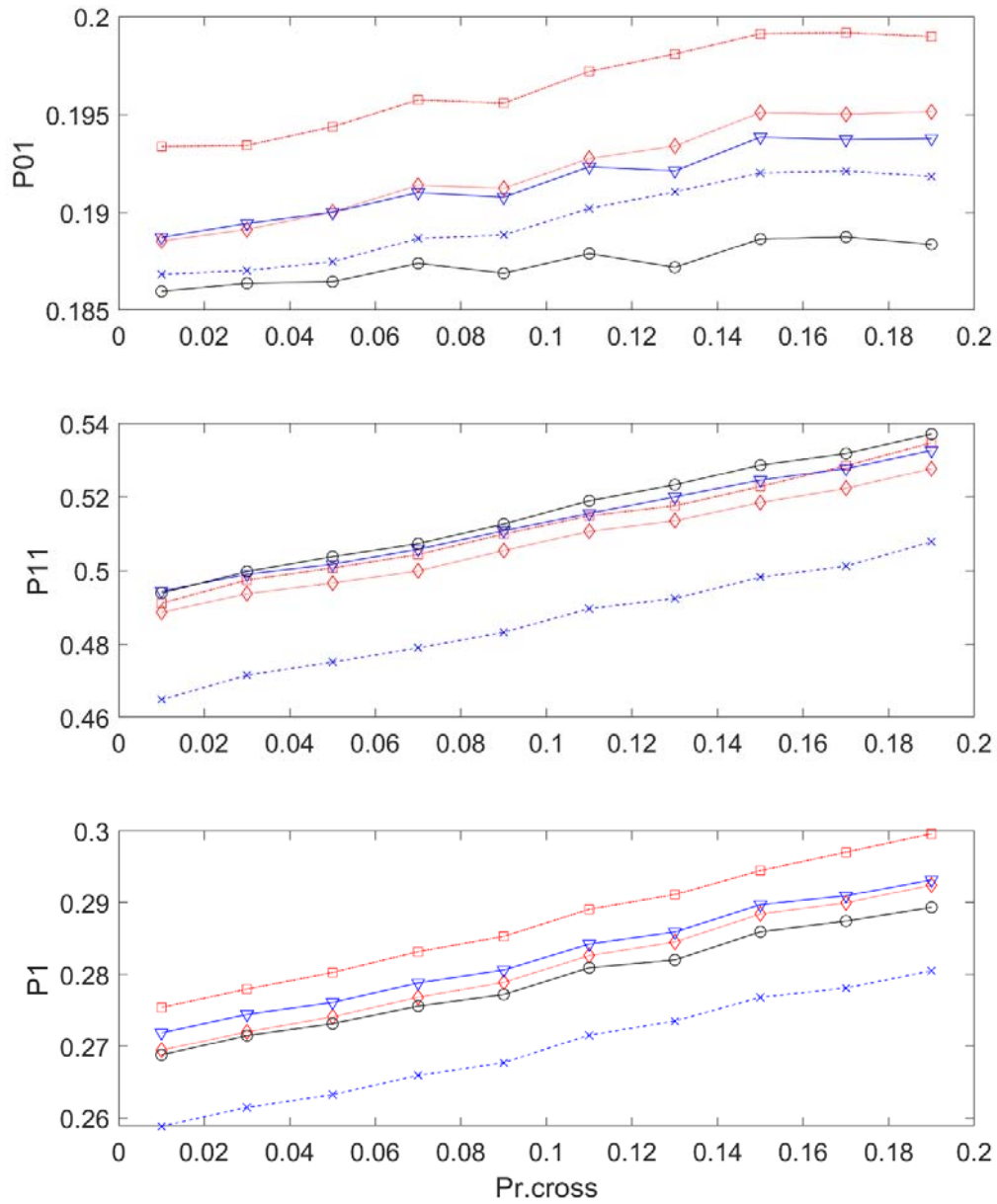
687

688

689

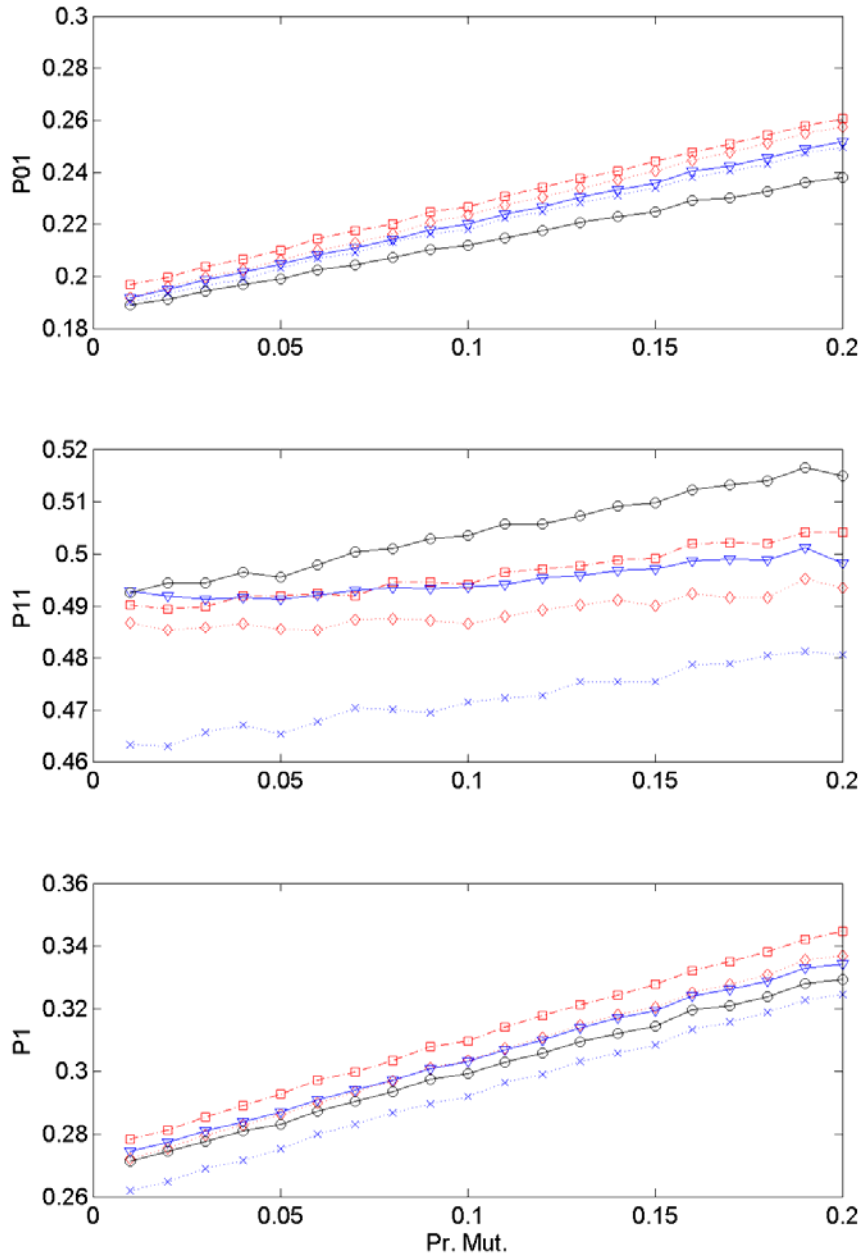
690

691



692

693 Figure 10. Transition probabilities and marginal distribution for the selected five stations along
 694 with changing the cross-over probability P_{cr} with the condition that the candidate value is one
 695 and the previous value is also one. See Eq.(15) for the detail.
 696



697

698 Figure 11. Transition probabilities and marginal distribution along with changing the cross-over
 699 probability with the condition that the mutation is processed only if the candidate value is one.
 700 See Eq.(16) for the detail.

701

702

703

704 Table A 1. Example dataset of daily rainfall with 12 weather stations and 16 days for measured
 705 rainfall (mm) in the upper part of this table and its corresponding occurrences in the bottom part
 706 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

707

708 Table A 2. Example dataset for estimating distances. The second row presents the current daily
709 precipitation occurrences for 12 stations and the rows below show the absolute difference
710 between the current occurrences (X_c) and the observed data in Table A 1. The last column
711 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

712

713

714

715 Table A 3. Example for selecting one sequence for \mathbf{X}_{c+1} . The second row presents the distances
716 in Table A 2. The third and fourth columns show the sorted days and distances for the smallest
717 distances to the largest in the second column. The fourth row presents the probabilities estimated
718 with Eq. (12). Note that there are six days whose distances are the same with each other. In this
719 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	

720

721

722 Table A 4. Example for GA mixture for \mathbf{X}_{c+1} . The second and third rows present two selected
 723 sets, while the third row shows the final set for \mathbf{X}_{c+1} with the crossover at S6 and S8 and the
 724 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

725
726

727