

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

**Streamflow estimation at partially gaged sites using multiple
dependence conditions via vine copulas**

Kuk-Hyun Ahn¹

October 2020

¹Assistant Professor, Department of Civil and Environmental Engineering, Kongju National University, Cheon-an, South Korea; *Corresponding author*; e-mail: ahnkukhyun@gmail.com

ABSTRACT

25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

Reliable estimates of missing streamflow values are relevant for water resources planning and management. This study proposes a multiple dependence condition model via vine copulas for the purpose of estimating streamflow at partially gaged sites. The proposed model is attractive in modeling the high dimensional joint distribution by building a hierarchy of conditional bivariate copulas when provided a complex streamflow gage network. The usefulness of the proposed model is firstly highlighted using a synthetic streamflow scenario. In this analysis, the bivariate copula model and a variant of the vine copulas are also employed to show the ability of the multiple dependence structure adopted in the proposed model. Furthermore, the evaluations are extended to a case study of 54 gages located within the Yadkin-Pee Dee River Basin, the eastern U. S. Both results inform that the proposed model is better suited for infilling missing values. To be specific, the proposed multiple dependence model shows the improvement of 9.2 % on average compared to the bivariate model from the historical case study. The performance of the vine copula is further compared with six other infilling approaches to confirm its applicability. Results demonstrate that the proposed model produces more reliable streamflow estimates than the other approaches. In particular, when applied to partially gaged sites with sufficient available data, the proposed model clearly outperforms the other models. Even though the model is illustrated by a specific case, it can be extended to other regions with diverse hydro-climatological variables for the objective of infilling.

Keywords: vine copulas, multiple dependence condition model, infilling approach, and streamflow estimation at partially gaged site

49 **1. Introduction**

50 Hydrological observation records covering long-term periods are instrumental in water
51 resources planning and management including the design of flood defense systems and
52 irrigation water management (Aissia et al., 2017; Beguería et al., 2019). However, available
53 streamflow data is often limited due to several situations like equipment failures, budgetary
54 cuts, and natural hazards (Kalteh and Hjorth, 2009). Missing data is particularly observed in
55 remote catchments where equipment failures are repaired only after significant delays
56 following extreme events, which can be crucial for hydrological frequency analysis. Hence,
57 hydrologists often rely on simulated sequences to infill missing data in partially gaged
58 catchments (Booker and Snelder, 2012) by using two primary modeling approaches such as:
59 (1) process-based models (i.e., estimating streamflow based on a conceptual understanding of
60 hydrological processes), and (2) transfer-based statistical models (i.e., transferring information
61 from gaged to ungaged catchments) (Farmer and Vogel, 2016). This paper focuses on the latter,
62 which estimates historical daily streamflow at inadequately and partially gaged sites by the
63 means of a statistical relationship.

64

65 Over the past few decades, a variety of statistical models including simple drainage area scaling
66 (Croley and Hartmann, 1986), spatial interpolation technique (Pugliese et al., 2014), regression
67 model (Beauchamp et al., 1989) and flow duration curves (FDCs; Hughes and Smakhtin, 1996),
68 have been developed. In particular, the flow duration curve method has been regarded as one
69 of the most trustworthy regionalization approaches (Archfield and Vogel, 2010; Boscarello et
70 al., 2016; Castellarin et al., 2004; Li et al., 2010; Mendicino and Senatore, 2013). If the target
71 watershed is completely ungaged, FDCs can be established using regression models to

72 regionalize the parameter sets of defined distributions (e.g., Ahn and Palmer, 2016a; Blum et
73 al., 2017) or to regionalize a set of primary quantiles (Cunderlik and Ouarda, 2006; Schnier
74 and Cai, 2014; Zaman et al., 2012). On the other hand, if the target watershed is poorly or
75 partially gaged, FDC models are built using the following four steps: (1) estimating non-
76 exceedance probability for recorded streamflow from the target watershed of interest; (2)
77 selecting one or multiple donor watersheds for the target watershed; (3) transferring the time-
78 series of non-exceedance probability from the donor watersheds for missing streamflow values;
79 and (4) converting corresponding streamflow values back from the transferred non-exceedance
80 probability. When FDCs are utilized for partially gaged watersheds, how the donor watersheds
81 are selected (step 2) and how the probabilities are transferred from the donor watersheds (step
82 3) are fairly crucial in the FDC framework.

83

84 Many studies have developed diverse approaches for steps 2 and 3 in FDC modelling. While
85 the basic formulation is that non-exceedance probabilities of the target site are transferred by
86 those at the single donor site, a weighted average of non-exceedance probability from the
87 selected donor sites has been suggested by Smakhtin (1999) instead. In addition, Farmer (2015)
88 adopted a kriging model to regionalized daily standard (i.e., z-scored) probabilities based on
89 non-exceedance probabilities from many donors in a region, using the quantile function of a
90 standard normal distribution. Although these studies are promising, the joint distribution of
91 non-exceedance probability between the target and donor watersheds is modeled based on a
92 Gaussian assumption which cannot properly permit different percentile values such as extremes
93 that have different spatial dependence structures from donor sites. To circumvent this limitation,
94 Worland et al. (2019) suggested the copula theory after showing that a unifying framework of
95 copulas is equivalent to that of FDC (i.e., estimations of the conditional probabilities at the

96 target watershed given known values at the donors).

97

98 Increasing attention has been received to copulas in the field of hydrology, with applications in
99 flood frequency analysis, drought risk analysis, and multi-site streamflow generations (Ahn
100 and Palmer, 2016b; Ariff et al., 2012; Chen et al., 2015; Daneshkhah et al., 2016; Fu and Butler,
101 2014). Copulas are effective mathematical functions that are capable of combining univariate
102 marginal distribution functions of random variables into their joint cumulative distribution
103 function and allow representation of diverse dependence structures between these random
104 variables corresponding to their family members (Sklar, A., 1959). For example, Fu and Butler
105 (2014) showed that the Gumbel copula performs well in representing multiple flooding
106 characteristics as compared to the other copulas from the Archimedean family, namely the
107 Clayton and Frank copulas. To estimate streamflow (i.e., infilling missing data) at poorly and
108 partially gaged sites, Worland et al. (2019) have developed bivariate copulas with an
109 Archimedean copula, but limited their application to a single donor. Albeit the limitation, their
110 bivariate copulas may be acceptable since the higher dimension of copulas is not rich enough
111 to model all possible mutual dependencies among multisite donors (see Karmakar and
112 Simonovic, 2009 for details). Hao and Singh (2013) also describe that multivariate copulas are
113 incapable of modeling multisite data exhibiting complex patterns of dependence.

114

115 However, if the theoretical limitation of a multivariate copula is mitigated, dependency
116 information from multiple donor sites may allow more reliable predictions of regionalized
117 streamflow. Vine copulas, also known as pair copulas, offer a far efficient way to construct
118 higher dimensional dependence (Bedford et al., 2002; Joe, 2014). They have hierarchical

119 structures that sequentially apply bivariate copulas as the building local blocks for constructing
120 a higher dimensional copula. The high flexibility of vine copulas enables modeling a wide
121 range of complex data dependencies. In particular, Aas et al. (2009) have popularized two
122 classes of vine copulas, canonical vines (C-vines) and drawable vines (D-vines) by allowing
123 diverse pair-copula families such as the bivariate Student-t copula and bivariate Clayton copula.
124 After the seminal paper, those two vines have been used in many fields including economics
125 (Arreola Hernandez et al., 2017; Zimmer, 2015), finance (Dissmann et al., 2013; Lu, 2013),
126 and engineering (Bhatti and Do, 2019; Erhardt et al., 2015; Xu et al., 2017). Similarly, a few
127 studies have used vine copulas in hydrologic applications with diverse purposes (Daneshkhah
128 et al., 2016; Liu et al., 2015; Vernieuwe et al., 2015; Shafaei et al., 2017) although they have
129 not been introduced to infill missing data.

130

131 Based on the usefulness of vine copulas, Kraus and Czado (2017) have developed a promising
132 algorithm that sequentially fits such a D-vine copula model ($\mathcal{M}_{\text{Kraus}}$). The algorithm adds
133 covariates to the model with the objective of maximizing a conditional likelihood and stops
134 adding covariates to the model when none of the remaining covariates can significantly
135 increase the model's conditional likelihood. While it is promising, one challenge that can arise
136 but has not been previously discussed is overfitting when covariates are correlated with each
137 other. In this situation, the model may adopt ineffective covariates and eventually leads to poor
138 predictions. In particular, for the purpose of infilling, streamflow values at the target site are
139 often correlated by those of many donors. Although the structure of $\mathcal{M}_{\text{Kraus}}$ is potentially
140 favorable to estimate streamflow, modified model procedure is required to determine the most
141 influential covariates.

142

143 This study forwards two novel contributions to infill missing data in the field of hydrology: (1)
144 a D-vine copula-based model is introduced to estimate streamflow for poorly and partially
145 gaged watersheds and (2) the existing model ($\mathcal{M}_{\text{Kraus}}$) is further improved by incorporating a
146 new procedure to determine the optimal number of donor sites (namely $\mathcal{M}_{\text{Dvine}}$). First,
147 synthetic data are generated to compare $\mathcal{M}_{\text{Kraus}}$ and $\mathcal{M}_{\text{Dvine}}$. In this analysis, bivariate
148 copulas (namely $\mathcal{M}_{\text{Bicop}}$) is also employed to demonstrate the usefulness of a high
149 dimensional joint dependence structure. Afterwards, a real infilling example is utilized to
150 compare the proposed vine-based model with six other streamflow-transfer models adopted in
151 literatures.

152

153 **2. Methodology**

154 *2.1 D-vine copulas*

155 A copula C is k -variate cumulative distribution function on $[0, 1]^k$ with all uniform margins.
156 The C can be understood as a function that links the marginal cumulative distributions
157 (F_1, \dots, F_k) to form a joint distribution F . The C associated with joint distribution F is a
158 distribution function $C: [0, 1]^k \rightarrow [0, 1]$ such that, for all streamflow vector $\mathbf{q} =$
159 $(q_1, \dots, q_k)^T$, the C satisfies:

160

$$161 \quad F(q_1, \dots, q_k) = C(F_1(q_1), \dots, F_k(q_k)) \quad \text{Eq. (1)}$$

162

163 where C is unique if F_1, \dots, F_k are continuous.

164 Based on Sklar's theorem (Sklar, A., 1959), a multivariate distribution function is a
165 composition of a set of marginal distributions; thus, equation (1) can be expressed in terms of
166 densities,

167

$$168 \quad f(q_1, \dots, q_k) = [\prod_{i=1}^k f_i(q_i)]c(F_1(q_1), \dots, F_k(q_k)) \quad \text{Eq. (2)}$$

169

170 where c is a k -dimensional copula density acquired by partial differentiation of the copula C
171 (i.e., $c(F_1(q_1), \dots, F_k(q_k)) := \frac{\partial^k}{\partial_1 \dots \partial_k} C(F_1(q_1), \dots, F_k(q_k))$) and $f_i(\cdot)$ is the marginal density
172 corresponding to $F_i(\cdot)$.

173

174 Following Bedford and Cooke (2001), any copula density $c(F_1(q_1), \dots, F_k(q_k))$ can be
175 decomposed into a product of $k(k - 1)/2$ pair copula densities. Aas et al. (2009) adopted this
176 idea and introduced the copula class of pair copula constructions (PCCs) known as vine copulas.
177 These copulas are suitable to model various dependency structures. Vine structures established
178 by $k(k - 1)/2$ pair copulas are arranged in $k - 1$ trees (Brechmann et al., 2013) and can be
179 categorized as C-vines and D-vines (Liu et al., 2015). This study focuses on D-vines since they
180 are more widely used in practice (Daneshkhah et al., 2016).

181

182 A D-vine is characterized by the ordering of its variables (see Figure 1). In the first tree, the

183 dependence of the first and second variables, of the second and third, of the third and fourth,
 184 and so on, is modeled using pair-copulas. In the second tree, conditional dependence of the first
 185 and third given the second variable (i.e., $c_{1,3|2}(F(q_1|q_2), F(q_3|q_2))$), the second and fourth
 186 given the third (i.e., $c_{2,4|3}(F(q_2|q_3), F(q_4|q_3))$), and so on, is modeled. Similarly, pairwise
 187 dependencies of two variables are modeled in subsequent trees conditioned on those variables
 188 which lie between the two variables in the first tree (e.g.,
 189 $c_{1,5|2,3,4}(F(q_1|q_2, q_3, q_4), F(q_5|q_2, q_3, q_4))$). The density of the k -dimensional D-vine can be
 190 computed as follows (Aas et al., 2009):

191

$$192 \quad f(q_1, \dots, q_k) = [\prod_{i=1}^k f_i(q_i)] \times$$

$$193 \quad \prod_{j=1}^{k-1} \prod_{\ell=1}^{k-j} c_{j, j+\ell | (j+1):(j+\ell-1)}(F(q_j | q_{j+1}, \dots, q_{j+\ell-1}), F(q_{j+\ell} | q_{j+1}, \dots, q_{j+\ell-1})) \quad \text{Eq. (3)}$$

194

195 where $c_{j, j+\ell | (j+1):(j+\ell-1)}$ indicates the bivariate copula densities.

196 For the five-dimensional D-vine copula as an example in Figure 1, the corresponding vine
 197 distribution has the joint density as follows:

198

$$199 \quad f(q_1, \dots, q_5) = [\prod_{i=1}^5 f_i(q_i)] c_{12} \cdot c_{23} \cdot c_{34} \cdot c_{45} \cdot c_{13|2} \cdot c_{24|3} \cdot c_{24|3} \cdot c_{35|4} \cdot c_{14|23} \cdot c_{25|34} \cdot$$

$$200 \quad c_{15|234} \quad \text{Eq. (4)}$$

201

202 where $c_{1,2}(F_1(q_1), F_2(q_2))$ is simply denoted as $c_{1,2}$.

203

204 As presented in equation (4), the conditional distribution functions and conditional bivariate
205 copulas are required in vine copula modeling. The conditional distribution functions
206 $F(q_j|q_{j+1}, \dots, q_{j+j-1})$, also known as h -functions, in equation (4) can be addressed using the
207 pair-copulas from lower trees by using equation (5). Let q_i be a conditional value of
208 $q_{j+1}, \dots, q_{j+j-1}$ and $\mathbf{v} = \{q_{j+1}, \dots, q_{j+j-1}\} \setminus q_i$ the streamflow vector without q_i used in the
209 following recursive relationship (Aas et al., 2009):

210

$$211 \quad h(q_j|\mathbf{v}) := F(q_j|\mathbf{v}) = \frac{\partial C_{ji|\mathbf{v}}(F(q_j|\mathbf{v}), F(q_i|\mathbf{v}))}{\partial F(q_i|\mathbf{v})} \quad \text{Eq. (5)}$$

212

213 where the h -function is associated with the pair-copula $C_{ji|\mathbf{v}}$.

214 More details about D-vines can be found in Bedford et al., (2002) and Czado (2010, 2019).

215

216 *2.2 Algorithm of D-vine copula-based estimation ($\mathcal{M}_{\text{Dvine}}$)*

217 Following Kraus and Czado (2017), a two-step estimation procedure is utilized for the
218 prediction of the streamflow value at the target watershed. The algorithm ($\mathcal{M}_{\text{Dvine}}$) is
219 developed using two library packages in the R programming language (Bevacqua, 2017;
220 Schepsmeier et al., 2015).

221

222 Let q_k be the quantile of streamflow at the target watershed given the streamflow values
 223 q_1, \dots, q_{k-1} from the donor sites. In the first step, the marginal cumulative probabilities
 224 $F_k(q_k)$ and $F_j(q_j)$, $j = 1, \dots, k - 1$, are estimated using the semiparametric approach. To be
 225 specific, this study uses the continuous kernel smoothing estimator (Geenens, 2014), which is,
 226 given observed streamflow q_i^ζ , $\zeta = 1, \dots, \xi$, at i th site, defined as $\hat{F}_i(q_i) = \frac{1}{nh} \sum_{\zeta=1}^{\xi} \Omega\left(\frac{q_i - q_i^\zeta}{h}\right)$.
 227 Here, $\Omega(q_i)$ is the “kernel” function with $\omega(\cdot)$ being a symmetric probability density
 228 function and h is the parameter controlling the smoothness of the final estimate. In this study,
 229 a Gaussian kernel is used for all $\omega(\cdot)$. The estimated cumulative probabilities are then
 230 employed to model the D-vine copula in the second step.

231

232 Next, to easily estimate conditional streamflow values at the target site, the D-vine copula is
 233 fitted with fixed order $F_k(q_k) - F_{I_1}(q_{I_1}) - F_{I_2}(q_{I_2}) - \dots - F_{I_{k-1}}(q_{I_{k-1}})$, such that $F_k(q_k)$ is
 234 the first node in the first tree and the other orders of donors (I_1, \dots, I_{k-1}) are decided based
 235 on their correlations to the target site (i.e., $F_{I_1}(q_{I_1})$ showing the greatest correlations to
 236 $F_k(q_k)$). To build the D-vine copula model, five bivariate copulas (Gaussian, Student-t, Frank,
 237 Gumbel, and Clayton copulas) are considered as potential pair copulas (building blocks)
 238 although more families of Copulas such as extreme value copulas (EVC) are desirable. The
 239 five candidates may be sufficient to represent diverse dependence structures. For example, a
 240 Gaussian copula is proper when the non-exceedance probabilities between two watersheds are
 241 associated in the body of their distribution but are not asymptotically dependent in the both
 242 tails. On the other hand, a Gumbel copula may be appropriate for the situation wherein the non-
 243 exceedance probabilities exhibit tail dependence, where high flows are connected by same
 244 rainfall events but low flows are not related (e.g., due to regulation) (Salvadori and De Michele,

2004). Details of the five bivariate copulas are presented in the Supporting Information. Parameters for the five bivariate copulas are estimated based on Kendall rank-based correlation (ρ^τ) between sites. The optimal bivariate copula for each pair copula is determined based on the penalized likelihood function (i.e., AIC).

249

The final number (χ_k) of donor sites is further optimized under a cross-validation approach. In this approach, 80 % of the regional data are employed for model fitting; the other 20 %, for testing. Again, this procedure is conducted 5 times, each time using a different set of data for testing. As a measure for the model's fit, the root mean squared error (RMSE; equation (6)) from observed streamflow at the target site is utilized.

255

$$RMSE_{\chi_k} = \sqrt{\frac{1}{\xi} \sum_{\zeta=1}^{\xi} (q_k - \hat{q}_k^{\chi})^2} \quad \text{Eq. (6)}$$

257

Finally, conditional streamflow values at the target site can be estimated using the inverse form of the conditional distribution function (i.e., Eq. 5). To depict the ideas, a trivariate case (i.e., $\chi = 2$) is considered here. Based on the streamflow values at the donor sites (q_2, q_3), \hat{q}_1 can be obtained using the conditional distribution function $h(q_1|q_2, q_3)$. For some fixed probabilities ϕ (e.g., $\phi = 0.1, \dots, 0.9$), $F_1(\hat{q}_1)$ is derived from $C_{1|2,3}$ using an explicit function:

264

265 $C_{1|2,3}^{-1}(\phi|F_2(q_2), F_3(q_3)) = h_{1|2}^{-1}(h_{1|32}^{-1}(\phi|h_{2|1}(F_2(q_2)|F_1(q_1))))|F_1(q_1))$ Eq. (7)

266

267 where $C_{1|2,3}^{-1}$ is the inverse of the copula function given the ϕ quantile curve of the copula
 268 (Liu et al., 2015; Xu and Childs, 2013). Therefore, the ϕ th copula-based conditional quantile
 269 function of streamflow at the target site can be calculated as follows:

270

271 $q_1(\phi|q_2q_3) = F_1^{-1}(C_{1|2,3}^{-1}(\phi|F_2(q_2), F_3(q_3))) =$
 272 $F_1^{-1}(h_{1|2}^{-1}(h_{1|32}^{-1}(\phi|h_{2|1}(F_2(q_2)|F_1(q_1))))|F_1(q_1)))$ Eq. (8)

273

274 Similarly, for the k -dimensional case, the ϕ th copula-based conditional quantile function can
 275 be calculated along with streamflow at the $k-1$ donor sites. To acquire an estimate at the target
 276 site, 1000 samples from uniform distribution over the interval $[0, 1]$ are generated using Monte
 277 Carlo simulations. In this study, the mean value of these generations is regarded as the best
 278 estimate.

279

280 **3. Application**

281 This study first explores the performance of $\mathcal{M}_{\text{Dvine}}$ under synthetic example. In this analysis,
 282 $\mathcal{M}_{\text{Bicop}}$ and $\mathcal{M}_{\text{Kraus}}$ are also employed to show the usefulness of $\mathcal{M}_{\text{Dvine}}$. For $\mathcal{M}_{\text{Bicop}}$, the
 283 optimal bivariate copula is selected based on the AIC while the five bivariate copulas (Gaussian,
 284 Student-t, Frank, Gumbel, and Clayton copulas) are considered as its potential candidates. A

285 brief description of two additional models are presented in the supporting information. After
286 that, those three models are used for a real application to 54 stream gages located in a region
287 of the eastern United States by estimating streamflow in partially gaged locations. Finally,
288 seven infilling approaches (Table 1) are also utilized and evaluated in a cross-validated
289 framework to evaluate the performance of the proposed model.

290

291 *3.1 Synthetic simulation*

292 Synthetic streamflow data are generated using controlled Monte Carlo experiment to explore
293 how well the three copula-based models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$, $\mathcal{M}_{\text{Dvine}}$) provide streamflow
294 predictions at the target site given a complex streamflow data in a pseudo gage network. In this
295 analysis, a six-dimensional streamflow set $(q_1^\zeta, q_2^\zeta, q_3^\zeta, q_4^\zeta, q_5^\zeta, q_6^\zeta)$, $\zeta = 1, \dots, \xi =$
296 2190 (i. e. $\frac{2190}{365} = 6$ years), is modelled using four bivariate copulas (Gaussian, Student-t,
297 Flank, and Clayton copulas) and lognormal distributions for margins (see Figure 2).

298

299 The performance of each model is evaluated in a calibration-validation framework. First,
300 synthetic streamflow data are generated for six-dimensional gage network. Then, φ years of
301 data are randomly selected to be assumed known at the target gage, and the streamflow for the
302 remaining $6-\varphi$ years of data is then estimated as missing values ($\varphi = 4$ in this analysis). This
303 process is repeated 20 times to build an ensemble prediction. In particular, this study assumes
304 the fifth streamflow data (i.e., q_5) to be predicted. In this assessment, two characteristics are
305 considered to compare the three models: model prediction reliability and uncertainty
306 quantification skill. Model prediction reliability is tested using the root mean squared error

307 (RMSE; Eq. 6) and Nash-Sutcliffe efficiency (NSE), which are further described in Section 3.4.
308 Uncertainty quantification skill is judged by the ability of each model to build prediction
309 intervals (PIs) that correctly bound predictions (see Section 3.4). Here, coverage probabilities,
310 defined as the proportion of the time that true values occur into these PIs, are employed to show
311 the usefulness of the proposed model.

312

313 *3.2 Application to the Yadkin-Pee Dee River*

314 The Yadkin-Pee Dee River Basin (Figure 3), covering around 18,700 km² and one of the largest
315 river basins in North Carolina and South Carolina (Fisk, 2010), is used as real data to evaluate
316 infilling ability. The basin flows from the northwestern corner of North Carolina near Blowing
317 Rock and extends south by southeast, crossing the south-central border of North Carolina into
318 South Carolina, with slightly more than half of its watershed in North Carolina. Most of the
319 land covered within the basin is forested or used for agriculture although urban areas of the
320 basin are expanding.

321

322 Daily streamflow data at 54 gages are gathered throughout the study region from web interface
323 of the U.S. Geological Survey (USGS) National Water Information System (NWIS) (U.S.
324 Geological Survey, 2018). The 54 gages are selected based on the following criteria: (1) all
325 gages are recorded continuously for 15 years of daily streamflow over the period from January
326 2004 to December 2018, and (2) gages have non-zero daily values for the period in the first
327 criterion since gages with streamflow values equal to zero require a more flexible modeling
328 structure. Thus, it is common to model zero flows separately in regionalization studies. Based

329 on the second criterion, this study discards 10 gage stations (not shown).

330

331 *3.3 Intermodel comparison framework*

332 A set of seven infilling approaches is used in the final assessment (see Table 1): (1) $\mathcal{M}_{\text{FDC-IDW}}$,
333 (2) $\mathcal{M}_{\text{IDW-streamflow}}$, (3) $\mathcal{M}_{\text{Rho-streamflow}}$, (4) $\mathcal{M}_{\text{FDC-highestrho}}$, (5) $\mathcal{M}_{\text{DAR-streamflow}}$, (6)
334 $\mathcal{M}_{\text{Kriging-streamflow}}$, and (7) $\mathcal{M}_{\text{Dvine}}$. This set of seven models is tested in a cross-validation
335 framework under two different cases. The two cases consider situations wherein φ have
336 values of 2 and 8 to represent relatively deficit- and sufficient-records for the target site. Similar
337 to the comparative assessment to show the usefulness of the proposed copula-based model (see
338 Section 3.1), each case is repeated 20 times by randomly selecting φ years over the applied
339 period. The reliability of each model is evaluated using RMSE and NSE metrics over the
340 validated four-year period randomly selected in the remaining data (i.e., 4 years in $15-\varphi$ years).

341

342 *3.4 Error metrics and error decomposition*

343 As presented in Sections 3.1 and 3.3, the root mean squared error (RMSE; Eq. 6) and Nash-
344 Sutcliffe efficiency (NSE) are employed to evaluate prediction skills:

345

$$346 \quad NSE = 1 - \frac{\sum_{\zeta=1}^{\xi} (\widehat{q^{\zeta}} - q^{\zeta})^2}{\sum_{\zeta=1}^{\xi} (q^{\zeta} - \bar{q}^{\zeta})^2} \quad \text{Eq. (9)}$$

347

348 The NSE (RMSE) can range from $-\infty$ to 1 (0 to ∞), with higher NSE (lower RMSE) implying
349 better performance. Both metrics have been commonly used in hydrology analysis (Boyle et
350 al., 2000).

351

352 Following derivations suggested in Gupta et al. (2009), the RMSE can be further decomposed
353 into three components:

354

$$355 \quad RMSE^2 = MSE = (\hat{\mu} - \mu)^2 + (\hat{\sigma} - \sigma)^2 + 2\sigma\hat{\sigma}(1 - r) \quad \text{Eq. (10)}$$

356

357 where μ ($\hat{\mu}$) and σ ($\hat{\sigma}$) represent the average and standard deviation for the observed
358 (estimated) streamflow, respectively, and r indicates the estimated correlation coefficient.
359 The first component $(\hat{\mu} - \mu)^2$ is a measure of how well the average of the observed
360 streamflow represents the average of the estimated streamflow; the second component
361 $(\hat{\sigma} - \sigma)^2$ is a measure of how well the variance of the prediction represents the variance of the
362 observed streamflow; and the third component $2\sigma\hat{\sigma}(1 - r)$ is dominated by the correlation
363 and is defined as the “timing” component (Worland et al., 2019). Using these three defined
364 components, their absolute contributions are explored in this study.

365

366 In addition, the accuracy of the uncertainty quantification skill is also evaluated for the copula-
367 based models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$, $\mathcal{M}_{\text{Dvine}}$). To be specific, this study utilizes the PI coverage
368 probability (PICP), which a common metric for this purpose (He et al., 2017; Niemierko et al.,

2019). It provides the relative number of data points that fall between the defined bounds as expressed follows:

371

$$\text{PICP} = \frac{1}{\xi} \sum_{\zeta=1}^{\xi} \Theta_{\zeta} \quad \text{with} \quad \Theta_{\zeta} = \begin{cases} 1, & \text{if } q^{\zeta} \in [L^{\zeta}, U^{\zeta}] \\ 0, & \text{else} \end{cases} \quad \text{Eq. (11)}$$

373

where Θ_{ζ} is the indicator variable if q^{ζ} is covered by the ζ th PI defined by the lower bound L^{ζ} and upper bound U^{ζ} . This study examines the prediction accuracy of single quantiles. Therefore, the lower bound is defined as $L^{\zeta} = -\infty$ and $U^{\zeta} = \widehat{q^{\zeta, \varpi}}$ where ϖ is the estimated quantile at time ζ . Accordingly, the upper bound is not a constant, but is re-assigned. By subtracting the nominal confidence ϖ from PICP, the average coverage error (ACE) is obtained as follows:

380

$$\text{ACE} = \text{PICP} - \varpi \quad \text{Eq. (12)}$$

382

The metric clearly indicates if the predicted quantile is underestimated ($\text{ACE} < 0$) or overestimated ($\text{ACE} > 0$) while taking small values around 0 for ideal case.

385

386 **4. Results**

387 *4.1 Results for synthetic experiment*

388 Prediction results from out-of-samples for the RMSE and NSE metrics are presented for the
389 three copula-based models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$, $\mathcal{M}_{\text{Dvine}}$) in Table 2. The ACE scores are also
390 described for $\varpi \in \{0.05, 0.10, 0.50, 0.90, 0.95\}$ in Table 3. When compared to the other
391 models, $\mathcal{M}_{\text{Bicop}}$ achieves lower RMSE values in the right tail of the RMSE distribution over
392 the validation periods, but severely underperforms the majority of the designed experiment,
393 suggesting this model formulation relying on a single donor leads to poor predictions. $\mathcal{M}_{\text{Kraus}}$
394 provides higher RMSE values for all the RMSE distribution, particularly for the right tail of
395 the RMSE distribution. The model utilizes streamflow data from all donors (i.e., five donor
396 sites) although the first two gages (Gages 1 and 2) show insignificant associations to the target
397 site ($r_{1,5} = 0.11$ and $r_{2,5} = 0.14$). $\mathcal{M}_{\text{Dvine}}$ unequivocally produces the best predictions.
398 $\mathcal{M}_{\text{Dvine}}$ adopts streamflow data from two or three donors (Gages 3, 4 and 6) without utilizing
399 streamflow data from the first two donors when a multiple dependence structure is established
400 to build an ensemble prediction. It outperforms $\mathcal{M}_{\text{Bicop}}$ and $\mathcal{M}_{\text{Kraus}}$ across all validation
401 periods, besides a few with the worst performance. Even in this case, the maximum RMSE of
402 $\mathcal{M}_{\text{Dvine}}$ is fairly less than the maximum RMSE of $\mathcal{M}_{\text{Kraus}}$.

403

404 In addition, the ACE results present how the three models characterize prediction uncertainty.
405 $\mathcal{M}_{\text{Dvine}}$ is capable of properly covering the predications across the entire distribution while
406 slight overestimation occurs for the smallest two quantiles. The remaining upper quantiles also
407 tend to slightly overestimate the true values but the overestimations are less than the other
408 models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$). Taken together, the results of the synthetic experiment suggest that
409 $\mathcal{M}_{\text{Dvine}}$ yields the best predictions among the copula-based models tested.

410

411 *4.2 Performance of the copula-based models in the Yadkin-Pee Dee River*

412 Using the insights developed from the synthetic experiment above, the three copula-based
413 models are applied to the streamflow data for the Yadkin-Pee Dee River. At first, upper and
414 lower tail dependences (λ_{upper} and λ_{lower}) are examined for all two pairs of sites (see Figure
415 4) using the approach of Schmid and Schmidt (2007). Theoretical background is described in
416 the Supporting Information (Text S3). Note that in this analysis, the dependences become more
417 obvious as the values approach unity. Two major insights emerge from this figure. First, many
418 site-pairs exhibit strong upper tail dependence, suggesting that streamflow variability has a
419 tendency to be more correlated under high-flow conditions compared to under low-flow
420 conditions (i.e., asymmetric dependence). The lack of lower-tail dependence may be due to
421 contributions governing low streamflow such as river regulation. Next, even under high- or
422 low-flow conditions, there is a wide range of tail dependence across the study basin (i.e.,
423 heterogeneous dependence). To sum up, a wide range of complex dependencies is observed in
424 the streamflow data over the study basin. The complex dependences suggest, when streamflow
425 is estimated from multiple donors, the potential usefulness of considering a multiple
426 dependence structure, which is one of the main features of vine copulas.

427

428 Figure 5 shows the RMSE and NSE results for the three copula-based models under a “leave-
429 one-out” cross validation framework. This process is repeated 20 times to build an ensemble
430 prediction by using test periods randomly defined. For this analysis, five years of data are
431 selected to be assumed as the observed period at the target gage, and another four years are
432 randomly selected in the remaining data for the test period. Similar to the results from the
433 synthetic experiment, \mathcal{M}_{Kraus} performs poorly in both the RMSE and NSE metrics (median

434 RMSE = 1.549 and NSE = 0.652). The bivariate copula performs well (median RMSE =
435 1.496), indicating that this approach efficiently leverages available information even though
436 the information is limited to single donor. Particularly, $\mathcal{M}_{\text{Bicop}}$ achieves lowest RMSE values
437 in the upper side of the RMSE box (e.g., third quartile), providing a strong uncertainty
438 quantification skill for the upper bound. However, $\mathcal{M}_{\text{Dvine}}$ yields the best median RMSE and
439 NSE values (= 1.359 and 0.719). Given the heterogeneous dependence conditions (see Figure
440 4), the high dimensional structures are effective in modeling a complex streamflow gage
441 network. This feature can substantially improve prediction of target site flows.

442

443 Figure 6a presents the ACE scores described for principal quantiles, $\varpi \in$
444 $\{0.05, 0.10, 0.20, \dots, 0.90, 0.95\}$, across all target sites under the cross validation framework.
445 Figure 6b presents 95% PIs for each model for an example time period (1 May 2018 to 31 July
446 2018) for one target site (USGS site ID: 02143500). Note that the ACE would ideally take zero
447 value, regardless of the quantiles. The ACE scores for the three models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$,
448 $\mathcal{M}_{\text{Dvine}}$) range from 0.004 to 0.0007 when considering all the quantiles together. However, the
449 scores vary depending on the quantiles. For instance, the ACE score for $\mathcal{M}_{\text{Kraus}}$ is noticeably
450 positive but is almost zero around the median streamflow, indicating that the model properly
451 represent uncertainty of the median streamflow. $\mathcal{M}_{\text{Bicop}}$ and $\mathcal{M}_{\text{Dvine}}$ result in very similar
452 ACE scores although $\mathcal{M}_{\text{Dvine}}$ performs slightly better than $\mathcal{M}_{\text{Bicop}}$. The differences in
453 characterization of prediction uncertainty can be confirmed from a particular target site (Figure
454 6b).

455

456 Based on the results in Figures 5 and 6, $\mathcal{M}_{\text{Dvine}}$ outperforms the other copula models (as
457 judged by model prediction reliability and uncertainty quantification skill) and is thus selected
458 as an appropriate copula model to infill missing data in partially gaged. Figure 7 shows an
459 example application of $\mathcal{M}_{\text{Dvine}}$ including the optimal donor sites, proper bivariate copulas
460 and their parameters for one target site (USGS site number #214645022) when the model is
461 calibrated using the full 15-year record.

462

463 *4.3 Intermodel comparison for streamflow estimation*

464 To assess the predictive skill of the proposed vine copula model, it is compared with six other
465 statistical models (see Table 1). Figure 8 shows RMSE and NSE for the seven models where
466 the streamflow values are estimated based on the available data defined by the two different
467 cases, labeled “deficit record” and “sufficient record” (see Section 3.3). Under all cases, the
468 vine copula approach outperforms the other infilling approaches. For example, for the
469 “sufficient record” case, median NSE for $\mathcal{M}_{\text{Dvine}}$ is 0.673 whereas those for
470 $\mathcal{M}_{\text{IDW-streamflow}}$ and $\mathcal{M}_{\text{rho-streamflow}}$ are 0.462 and 0.649, respectively. In this analysis, the
471 approaches, which are based on streamflow values of the donor sites without utilizing non-
472 exceedance probability including DAR-streamflow and Kriging-streamflow, yield relatively
473 increased bias in their predictions. On the other hand, an application of FDC models offers
474 reliable predictions. For instance, for the “sufficient record” case, median RMSE for
475 $\mathcal{M}_{\text{FDC-highestrho}}$ is 1.603 compared to that of a direct of using streamflow (e.g., median RMSE
476 of $\mathcal{M}_{\text{FDC-streamflow}} = 3.422$ for the sufficient record). Similar interpretation can be found in
477 the comparison between $\mathcal{M}_{\text{FDC-IDW}}$ and $\mathcal{M}_{\text{IDW-streamflow}}$. The results from these approaches
478 suggest that utilizing FDC process leads to a reliable estimation, which is a primary structure

479 in the vine copula. The other noticeable feature is that available data length provides a
480 significant influence on performance of some infilling methods. In particular, this is quite
481 evident for the vine copula model (median RMSEs: 1.598 and 1.379 for deficit and sufficient
482 records, respectively).

483

484 *4.4 Prediction error decomposition*

485 The RMSE is decomposed into their components (bias, variance, and timing components) for
486 both the “deficit record” and “sufficient record” predictions (Figure 9). For the both cases,
487 timing components primarily bring about the majority of prediction errors for all seven models.
488 In particular, models estimating directly streamflow values (IDW-streamflow, DAR-
489 streamflow, Kriging-streamflow) produce a somewhat biased component, which increases
490 when a shorter record is employed in the model. For instance, the timing component for
491 $\mathcal{M}_{\text{IDW-streamflow}}$ is 4.11 and 3.75 for the “deficit record” and “sufficient record”, respectively.
492 Moreover, timing components dominate the error metric for all cases. However, the importance
493 of variance component is increased, especially in three models (FDC-IDW, DAR-streamflow,
494 Kriging-streamflow). Lastly, the results inform that if the proposed vine copulas approach is
495 adapted, variance and timing components are better captured, leading to better streamflow
496 estimations, which is beneficial in the practical applications of water resources management.

497

498 Finally, two predictions are further produced using two additional experiments: (1) the
499 observed marginal cumulative probabilities (i.e., using all 15 years) and conditional streamflow
500 values constructed from the partial record (i.e., based on φ years), and (2) the estimated

501 marginal cumulative probabilities (i.e., based on φ years) and conditional streamflow values
502 constructed from the full record (i.e., all 15 years). Their prediction abilities are evaluated over
503 the validated four-year period randomly selected in the remaining data. Similar to the previous
504 analysis, each analysis is tested 20 times. The results from these experiments provide an
505 inference to better isolate how error components from the two-step procedure (see section 2.2)
506 influence prediction skill.

507

508 Figure 10 shows the ACE scores from the out-of-sampled predictions using the proposed Dvine
509 model under the two scenarios. When considering all the quantiles together, the ACE scores
510 for the two scenarios are 0.003 (scenario #1) and 0.006 (scenario #2) on average under the
511 “deficit record” prediction. Also, the scores under the “sufficient record” prediction are all
512 nearly 0.003. Those results of the scores are sufficiently closed to zero, implying that both
513 predictions are reliable. Yet, compared to the predictions estimated by the cumulative
514 probabilities estimated by the partial record, and conditional models constructed by full records
515 (i.e., scenario #2), the ACE scores are achieved better, if the cumulative probabilities are
516 determined by the full record, except for some of the low and high quantiles. Similar
517 interpretation can be found in the NSE performance of two scenarios (see insets of Figure 10).
518 It may suggest that the first procedure (i.e., how to determine the cumulative probabilities for
519 the target site and its donors) is needed to pay careful attention when $\mathcal{M}_{\text{Dvine}}$ is utilized.
520 Nevertheless, the procedure to construct the conditional model in a streamflow gage network
521 is obviously crucial since the over or under-estimations are observed in many quantiles when
522 the insufficient sampling is employed in this process.

523

524 **5. Conclusion**

525 This study introduces a multiple dependence conditional model (i.e., vine copulas) to produce
526 streamflow estimates at partially gaged sites. The model includes a flexible high dimensional
527 joint dependence structure and conditional bivariate copula simulations. In order to confirm the
528 usefulness of a multiple dependence structure and the procedure for an appropriate number of
529 donor sites in the final vine copula model, the bivariate copula model and two types of vine
530 copulas with their unique procedure to determine the optimal number of donor sites are first
531 investigated using the generated data. These analyses were further extended in a case study of
532 the Yadkin-Pee Dee River Basin, the eastern United States by estimating streamflow in partially
533 gaged locations. In this analysis, six statistical infilling approaches were also employed to
534 represent applicability of the proposed model.

535

536 Results of the synthetic experiment and application to the Yadkin-Pee Dee River Basin
537 demonstrate that the propose model has benefits in some aspects. First, a multiple dependence
538 structure adopted in the proposed model is beneficial. From the massive evaluation experiments,
539 this study shows that multiple dependence structure clearly outperforms a single dependence
540 structure although there is the risk of overfitting when too many dependence structures are
541 employed. For example, the proposed model shows the improvement of 9.2 % on average
542 compared to the bivariate model from the evaluation experiment over the historical case study.
543 Moreover, this study confirms that the proposed multiple dependence structure model with
544 their optimum number of donor sites produces more reliable streamflow estimation than other
545 common infilling models. To be specific, for the “sufficient record” case, the proposed model
546 shows the improvement of 13.9 % on average compared to the FDC-highestrho model. Next,

547 the proposed model allows the development of confidence intervals to consider prediction
548 uncertainty, which is fairly attractive compared to other models. For example, Bárdossy and
549 Pegram (2013) argue that confidence intervals obtained using an ordinary kriging model do not
550 reflect the prediction uncertainty well particularly on a daily scale. Overall, this study exhibits
551 that a vine copula is potentially an effective tool to support water resource management
552 planners for objectives like gap-filling or extending missing streamflow records.

553

554 While the results of the proposed model are favorable, there are possible limitations worthy of
555 further discussion. First, the proposed method is computationally expensive, even after
556 adopting the multicore processing to reduce the computational burden. This becomes more
557 problematic when the method is applied to a larger, more complex streamflow gaging network.
558 Nevertheless, because local water managers do not need to build the model repeatedly
559 whenever they face missing values once the model is calibrated for a specific site, this
560 computational burden may be a minor issue. Second, the assessment illustrated in this study
561 focuses on model performance under cross-validation at partially gaged basins, but additional
562 work is needed to extend the proposed model to ungaged basins. One possible way is to build
563 a regression based model with spatial proximity and physical basin characteristics to define
564 associations between the target and donor sites (e.g., Ahn and Steinschneider, 2019). Lastly,
565 this study does not consider potential nonstationarity in FDCs and correlations caused by the
566 influence of anthropogenic activity and change in land use. Nonstationarity may not be
567 problematic in this analysis since the assessment is limited to 15 years across the gaging
568 network. However, if longer records were used, it would be beneficial to consider the potential
569 nonstationarity. This exploration is left for future work.

570

571 There are several opportunities to improve the model structure. For instance, a vine copula is
572 able to incorporate more additional conditioning variables. One feasible approach is to add a
573 time series of climate data (e.g., precipitation) or to decompose a time series of streamflow
574 from the donor sites into a number of periodic components at different frequency levels through
575 the wavelet decomposition approach (Kisi and Cimen, 2011). Moreover, although the proposed
576 model providing a more flexible way to model multivariate dependences, it can be further
577 improved by not adopting the standard assumption (i.e., simplifying assumption) that the
578 conditional pair-copulas depend on the conditioning variables through the conditional margins
579 (Acar et al., 2012). One possible alternative is the use of the semi-parametric estimation of a
580 conditional copula (Acar et al., 2012; Vatter and Chavez-Demoulin, 2015). This semi-
581 parametric approach enables an estimate of the dependence parameters which do not rely on
582 the simplifying assumption, eventually leading to more reliable infilling estimations. I believe
583 this provides an interesting avenue for future research.

584

585 Lastly, the results presented here are specific to a study basin used in a case study. The proposed
586 model has not restricted to other watersheds around the world and its application is further
587 required towards drawing more generalized conclusions. In addition, the model could be used
588 for the purpose of infilling missing values of other hydrometeorological variables besides
589 streamflow (e.g., precipitation and soil moisture). For this application, the implementation of a
590 vine copula with combined discrete and continuous margins (i.e., to account for no rainfall
591 days) should be explored (e.g., Stoeber et al., 2013).

592

593

Acknowledgements

594 This work was supported by the National Research Foundation of Korea (NRF) grant funded
595 by the Korea government (MSIT) (No. 2019R1C1C1002438). Also, the author would like to
596 acknowledge Scott Steinschneider for his helpful comments during the development of this
597 paper.

598

599

REFERENCES

600

601 Aas, K., Czado, C., Frigessi, A., and Bakken, H.: Pair-copula constructions of multiple
602 dependence, *Insur. Math. Econ.*, 44, 182–198, 2009.

603 Acar, E. F., Genest, C., and Nešlehová, J.: Beyond simplified pair-copula constructions, *J.*
604 *Multivar. Anal.*, 110, 74–90, 2012.

605 Ahn, K.-H. and Palmer, R.: Regional flood frequency analysis using spatial proximity and basin
606 characteristics: Quantile regression vs. parameter regression technique, *J. Hydrol.*, 540, 515–
607 526, <http://dx.doi.org/10.1016/j.jhydrol.2016.06.047>, 2016a.

608 Ahn, K.-H. and Palmer, R. N.: Use of a nonstationary copula to predict future bivariate low
609 flow frequency in the Connecticut river basin, *Hydrol. Process.*, 30, 3518–3532,
610 <https://doi.org/10.1002/hyp.10876>, 2016b.

611 Ahn, K.-H. and Steinschneider, S.: Hierarchical Bayesian Model for Streamflow Estimation at
612 Ungauged Sites via Spatial Scaling in the Great Lakes Basin, *J. Water Resour. Plan. Manag.*,
613 145, 04019030, 2019.

614 Aissia, M.-A. B., Chebana, F., and Ouarda, T. B.: Multivariate missing data in hydrology–
615 Review and applications, *Adv. Water Resour.*, 110, 299–309, 2017.

616 Archfield, S. A. and Vogel, R. M.: Map correlation method: Selection of a reference streamgage
617 to estimate daily streamflow at ungauged catchments, *Water Resour. Res.*, 46, 2010.

618 Ariff, N., Jemain, A., Ibrahim, K., and Wan Zin, W.: IDF relationships using bivariate copula
619 for storm events in Peninsular Malaysia, *J. Hydrol.*, 470, 158–171, 2012.

620 Arreola Hernandez, J., Hammoudeh, S., Nguyen, D. K., Al Janabi, M. A., and Reboredo, J. C.:
621 Global financial crisis and dependence risk analysis of sector portfolios: a vine copula approach,
622 *Appl. Econ.*, 49, 2409–2427, 2017.

623 Bárdossy, A. and Pegram, G.: Interpolation of precipitation under topographic influence at
624 different time scales, *Water Resour. Res.*, 49, 4545–4565, 2013.

- 625 Bárdossy, A. and Pegram, G.: Infilling missing precipitation records—A comparison of a new
626 copula-based method with other techniques, *J. Hydrol.*, 519, 1162–1170, 2014.
- 627 Beauchamp, J., Downing, D., and Railsback, S.: Comparison of regression and time-series
628 methods for synthesizing missing streamflow records, *JAWRA J. Am. Water Resour. Assoc.*,
629 25, 961–975, 1989.
- 630 Bedford, T. and Cooke, R. M.: Probability density decomposition for conditionally dependent
631 random variables modeled by vines, *Ann. Math. Artif. Intell.*, 32, 245–268, 2001.
- 632 Bedford, T., Cooke, R. M., and others: Vines—a new graphical model for dependent random
633 variables, *Ann. Stat.*, 30, 1031–1068, 2002.
- 634 Beguería, S., Tomas-Burguera, M., Serrano-Notivoli, R., Peña-Angulo, D., Vicente-Serrano, S.
635 M., and González-Hidalgo, J.-C.: Gap filling of monthly temperature data and its effect on
636 climatic variability and trends, *J. Clim.*, 32, 7797–7821, 2019.
- 637 Bevacqua, E.: *CDVineCopulaConditional: Sampling from Conditional C-and D-Vine Copulas*,
638 R package version 0.1. 0, 2017.
- 639 Bhatti, M. I. and Do, H. Q.: Recent development in copula and its applications to the energy,
640 forestry and environmental sciences, *Int. J. Hydrog. Energy*, 44, 19453–19473, 2019.
- 641 Blum, A. G., Archfield, S. A., and Vogel, R. M.: On the probability distribution of daily
642 streamflow in the United States, *Hydrol. Earth Syst. Sci.*, 21, 3093–3103, 2017.
- 643 Booker, D. and Snelder, T.: Comparing methods for estimating flow duration curves at
644 ungauged sites, *J. Hydrol.*, 434, 78–94, 2012.
- 645 Boscarello, L., Ravazzani, G., Cislighi, A., and Mancini, M.: Regionalization of flow-duration
646 curves through catchment classification with streamflow signatures and physiographic–climate
647 indices, *J. Hydrol. Eng.*, 21, 05015027, 2016.
- 648 Boyle, D. P., Gupta, H. V., and Sorooshian, S.: Toward improved calibration of hydrologic
649 models: Combining the strengths of manual and automatic methods, *Water Resour. Res.*, 36,
650 3663–3674, 2000.
- 651 Brechmann, E. C., Hendrich, K., and Czado, C.: Conditional copula simulation for systemic
652 risk stress testing, *Insur. Math. Econ.*, 53, 722–732, 2013.
- 653 Castellarin, A., Galeati, G., Brandimarte, L., Montanari, A., and Brath, A.: Regional flow-
654 duration curves: reliability for ungauged basins, *Adv. Water Resour.*, 27, 953–965, 2004.
- 655 Chen, L., Singh, V. P., Guo, S., Zhou, J., and Zhang, J.: Copula-based method for multisite
656 monthly and daily streamflow simulation, *J. Hydrol.*, 528, 369–384, 2015.
- 657 Croley, T. and Hartmann, H.: NOAA Technical Memorandum ERL GLERL-61: Near-Real-
658 Time Forecasting of Large-Lake Water Supplies: A User’s Manual, Ann Arbor MI, 1986.
- 659 Cunderlik, J. M. and Ouarda, T. B.: Regional flood-duration–frequency modeling in the

- 660 changing environment, *J. Hydrol.*, 318, 276–291, 2006.
- 661 Czado, C.: Pair-copula constructions of multivariate copulas, in: *Copula theory and its*
662 *applications*, Springer, 93–109, 2010.
- 663 Czado, C.: *Analyzing Dependent Data with Vine Copulas*, Lect. Notes Stat. Springer, 2019.
- 664 Daneshkhan, A., Remesan, R., Chatrabgoun, O., and Holman, I. P.: Probabilistic modeling of
665 flood characterizations with parametric and minimum information pair-copula model, *J.*
666 *Hydrol.*, 540, 469–487, 2016.
- 667 Dissmann, J., Brechmann, E. C., Czado, C., and Kurowicka, D.: Selecting and estimating
668 regular vine copulae and application to financial returns, *Comput. Stat. Data Anal.*, 59, 52–69,
669 2013.
- 670 Erhardt, T. M., Czado, C., and Schepsmeier, U.: R-vine models for spatial time series with an
671 application to daily mean temperature, *Biometrics*, 71, 323–332, 2015.
- 672 Farmer, W.: Estimating records of daily streamflow at ungauged locations in the southeast
673 United States, PhD Dissertation Tufts Univ. MA USA, 2015.
- 674 Farmer, W. H. and Vogel, R. M.: On the deterministic and stochastic use of hydrologic models,
675 *Water Resour. Res.*, 52, 5619–5633, 2016.
- 676 Fisk, J.: *Reproductive Ecology and Habitat Use of the Robust Redhorse in the Pee Dee River,*
677 *North Carolina and South Carolina.*, 2010.
- 678 Fu, G. and Butler, D.: Copula-based frequency analysis of overflow and flooding in urban
679 drainage systems, *J. Hydrol.*, 510, 49–58, 2014.
- 680 Geenens, G.: Probit transformation for kernel density estimation on the unit interval, *J. Am.*
681 *Stat. Assoc.*, 109, 346–358, 2014.
- 682 Gupta, H. V., Kling, H., Yilmaz, K. K., and Martinez, G. F.: Decomposition of the mean squared
683 error and NSE performance criteria: Implications for improving hydrological modelling, *J.*
684 *Hydrol.*, 377, 80–91, 2009.
- 685 Hao, Z. and Singh, V. P.: Modeling multisite streamflow dependence with maximum entropy
686 copula, *Water Resour. Res.*, 49, 7139–7143, 2013.
- 687 He, Y., Liu, R., Li, H., Wang, S., and Lu, X.: Short-term power load probability density
688 forecasting method using kernel-based support vector quantile regression and Copula theory,
689 *Appl. Energy*, 185, 254–266, 2017.
- 690 Hughes, D. and Smakhtin, V.: Daily flow time series patching or extension: a spatial
691 interpolation approach based on flow duration curves, *Hydrol. Sci. J.*, 41, 851–871, 1996.
- 692 Joe, H.: *Dependence modeling with copulas*, CRC press, 2014.
- 693 Kalteh, A. M. and Hjorth, P.: Imputation of missing values in a precipitation–runoff process

- 694 database, *Hydrol. Res.*, 40, 420–432, 2009.
- 695 Karmakar, S. and Simonovic, S.: Bivariate flood frequency analysis. Part 2: a copula-based
696 approach with mixed marginal distributions, *J. Flood Risk Manag.*, 2, 32–44, 2009.
- 697 Kisi, O. and Cimen, M.: A wavelet-support vector machine conjunction model for monthly
698 streamflow forecasting, *J. Hydrol.*, 399, 132–140, 2011.
- 699 Kraus, D. and Czado, C.: D-vine copula based quantile regression, *Comput. Stat. Data Anal.*,
700 110, 1–18, 2017.
- 701 Li, M., Shao, Q., Zhang, L., and Chiew, F. H.: A new regionalization approach and its
702 application to predict flow duration curve in ungauged basins, *J. Hydrol.*, 389, 137–145, 2010.
- 703 Liu, Z., Zhou, P., Chen, X., and Guan, Y.: A multivariate conditional model for streamflow
704 prediction and spatial precipitation refinement, *J. Geophys. Res. Atmospheres*, 120, 10–116,
705 2015.
- 706 Lu, W.: A high-dimensional vine copula approach to comovement of China’s financial markets,
707 in: 2013 International Conference on Management Science and Engineering 20th Annual
708 Conference Proceedings, 1538–1543, 2013.
- 709 Mendicino, G. and Senatore, A.: Evaluation of parametric and statistical approaches for the
710 regionalization of flow duration curves in intermittent regimes, *J. Hydrol.*, 480, 19–32, 2013.
- 711 Niemierko, R., Töppel, J., and Tränkler, T.: A D-vine copula quantile regression approach for
712 the prediction of residential heating energy consumption based on historical data, *Appl. Energy*,
713 233, 691–708, 2019.
- 714 Pugliese, A., Castellarin, A., and Brath, A.: Geostatistical prediction of flow–duration curves
715 in an index-flow framework, *Hydrol. Earth Syst. Sci.*, 18, 3801–3816, 2014.
- 716 Salvadori, G. and De Michele, C.: Frequency analysis via copulas: Theoretical aspects and
717 applications to hydrological events, *Water Resour. Res.*, 40, 2004.
- 718 Schepsmeier, U., Stoeber, J., Brechmann, E. C., Graeler, B., Nagler, T., Erhardt, T., Almeida,
719 C., Min, A., Czado, C., Hofmann, M., and others: Package ‘VineCopula,’ R Package Version,
720 2, 2015.
- 721 Schmid, F. and Schmidt, R.: Multivariate conditional versions of Spearman’s rho and related
722 measures of tail dependence, *J. Multivar. Anal.*, 98, 1123–1140, 2007.
- 723 Schnier, S. and Cai, X.: Prediction of regional streamflow frequency using model tree
724 ensembles, *J. Hydrol.*, 517, 298–309, 2014.
- 725 Shafaei, M., Fakheri-Fard, A., Dinpashoh, Y., Mirabbasi, R., and De Michele, C.: Modeling
726 flood event characteristics using D-vine structures, *Theor. Appl. Climatol.*, 130, 713–724, 2017.
- 727 Sklar, A.: *Fonctions de Répartition À N Dimensions Et Leurs Marges*, Université Paris 8, 1959.

728 Smakhtin, V. Y.: Generation of natural daily flow time-series in regulated rivers using a non-
729 linear spatial interpolation technique, *Regul. Rivers Res. Manag. Int. J. Devoted River Res.*
730 *Manag.*, 15, 311–323, 1999.

731 Stoeber, J., Joe, H., and Czado, C.: Simplified pair copula constructions—limitations and
732 extensions, *J. Multivar. Anal.*, 119, 101–118, 2013.

733 U.S. Geological Survey: National Water Information System (NWISWeb): U.S. Geological
734 Survey database, 2018.

735 Vatter, T. and Chavez-Demoulin, V.: Generalized additive models for conditional dependence
736 structures, *J. Multivar. Anal.*, 141, 147–167, 2015.

737 Vernieuwe, H., Vandenberghe, S., De Baets, B., and Verhoest, N.: A continuous rainfall model
738 based on vine copulas, *Hydrol. Earth Syst. Sci.*, 19, 2685–2699, 2015.

739 Worland, S. C., Steinschneider, S., Farmer, W., Asquith, W., and Knight, R.: Copula theory as
740 a generalized framework for flow-duration curve based streamflow estimates in ungaged and
741 partially gaged catchments, *Water Resour. Res.*, 55, 9378–9397, 2019.

742 Xu, D., Wei, Q., Elsayed, E. A., Chen, Y., and Kang, R.: Multivariate degradation modeling of
743 smart electricity meter with multiple performance characteristics via vine copulas, *Qual. Reliab.*
744 *Eng. Int.*, 33, 803–821, 2017.

745 Xu, Q. and Childs, T.: Evaluating forecast performances of the quantile autoregression models
746 in the present global crisis in international equity markets, *Appl. Financ. Econ.*, 23, 105–117,
747 2013.

748 Zaman, M. A., Rahman, A., and Haddad, K.: Regional flood frequency analysis in arid regions:
749 A case study for Australia, *J. Hydrol.*, 475, 74–83, 2012.

750 Zimmer, D. M.: Analyzing comovements in housing prices using vine copulas, *Econ. Inq.*, 53,
751 1156–1169, 2015.

752

753

754

755

756

757

758

759

760

761

762

763 **List of Figures**

- 764 Figure 1 Example of D-vine structures with 5 variables, 4 trees and 10 edges
765
- 766 Figure 2 Structure of the 6-Dimensional vine model and marginal for the synthetic
767 simulation. $LN(\pi, \sigma^2)$ denotes the log normal distribution with its mean (π) and
768 variance (σ^2). The target gage is highlighted.
769
- 770 Figure 3 Map of the Yadkin-Pee Dee Basin with 54 stream gage stations
771
- 772 Figure 3 Density of the Akaike Information Criterion (AIC) for the six univariate
773 distributions applied across sites in the study basin. Mean values for each AIC
774 density are also presented.
775
- 776 Figure 4 Pairwise upper and lower tail dependence for watersheds in the Yadkin-Pee Dee
777 River Basin. The upper triangular matrix shows values for the upper-tail
778 dependence and the lower triangular matrix presents values for the lower-tail
779 dependence. The metrics can range from 0 to 1, with higher values suggesting
780 greater interdependence of two streamflows for each upper- and lower-tail.
781
- 782 Figure 5 Model performance for the Yadkin-Pee Dee river under a cross-validation
783 framework based on RMSE (dark squares) and NSE (light squares). Here, the
784 RMSE (NSE) can range from 0 to ∞ ($-\infty$ to 1), with lower RMSE (higher NSE)
785 implying better performance.
786
- 787 Figure 6 (a) Average coverage error from three copula-based models for the Yadkin-Pee
788 Dee River Basin across exemplarily quantiles, and (b) 95% PIs for three models
789 for an example period (1 May 2018 to 31 July 2018) for a specific target gauge
790 (USGS site ID: 02143500). Observed streamflow (black solid line) is also
791 presented in each figure.
792
- 793 Figure 7 Structure of the Dvine copula applied for a particular target site (USGS site ID:
794 214645022) with the defined bivariate copulas and their parameters.
795
- 796 Figure 8 Inter-model comparison using cross-validation experiments based on RMSE
797 (upper) and NSE (lower). Here, lower RMSE suggests more accurate estimations
798 for infilling missing values.
799
- 800 Figure 9 Three contributions from the decomposed mean squared error (MSE) for the
801 cross-validation experiment with (a) the deficit record and (b) sufficient record
802 scenarios.
803
- 804 Figure 10 Average coverage error of the Dvine model for two scenarios under (a) the
805 “deficit” and (b) “sufficient” cases.
806
807
808
809

810
811
812
813

814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856

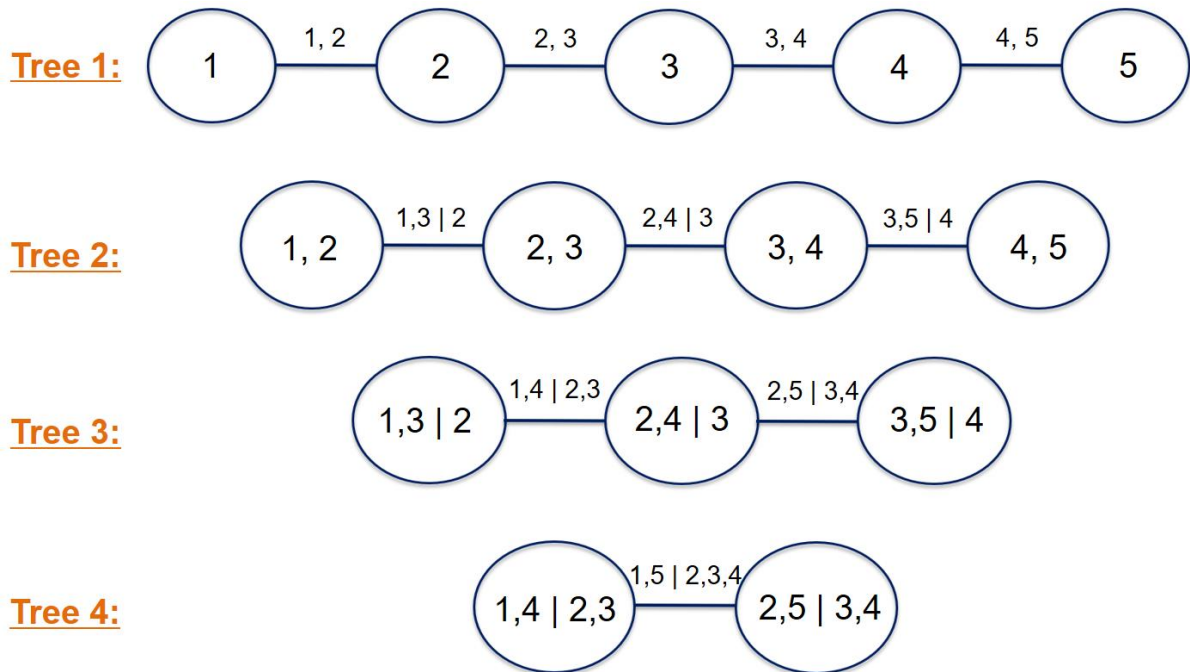
List of Tables

Table 1 Seven infilling approaches discussed in the study.

Table 2 RMSE and NSE results over the validation periods under synthetic experiment for comparing copula-based model formulations. Best metric values for each quantile are italicized and bolded.

Table 3 Results of average coverage error (ACE) over the validation periods under synthetic experiment for comparing copula-based model formulations. Best metric values for each quantile are italicized and bolded.

857



858

859 Figure 1 Example of D-vine structures with 5 variables, 4 trees and 10 edges

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

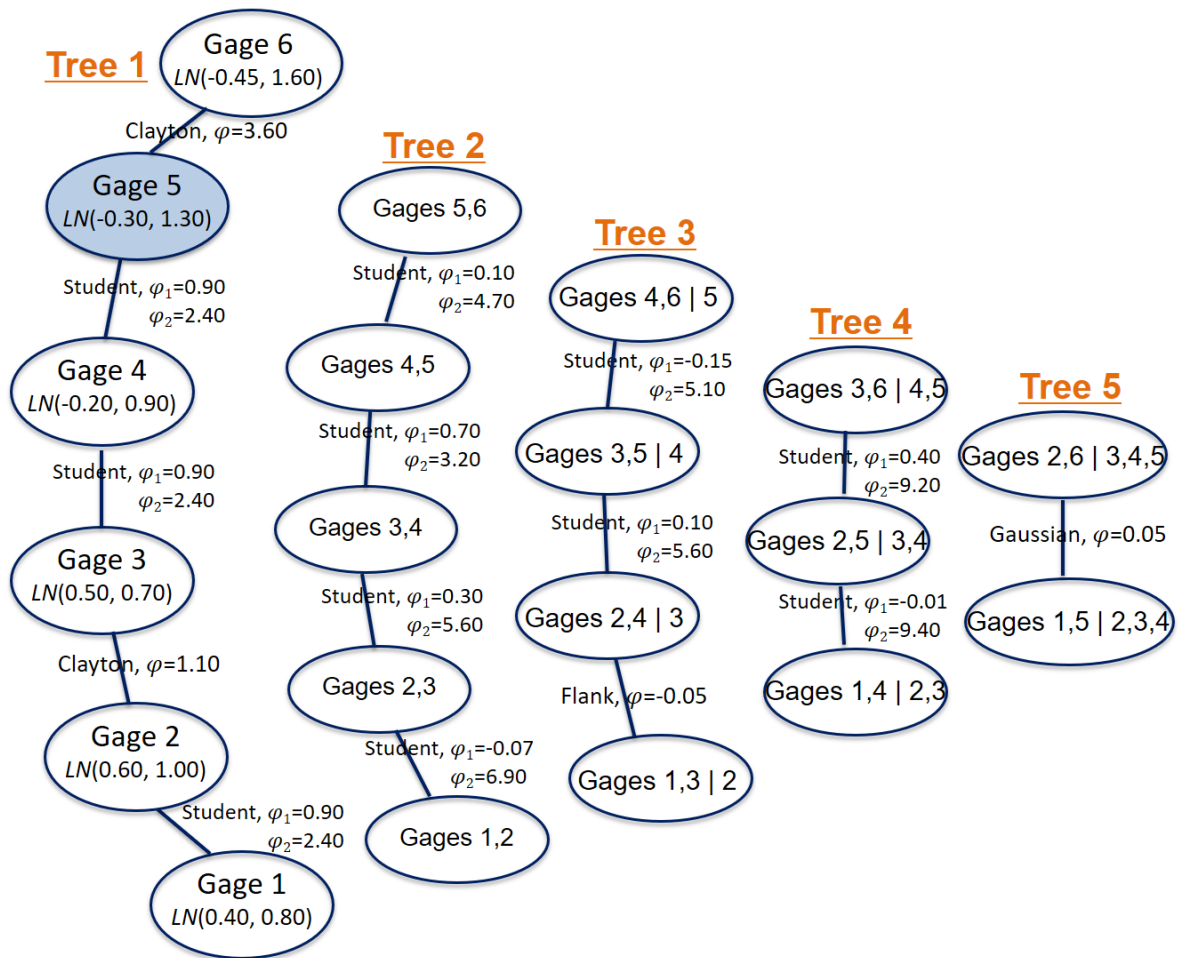
883

884

885

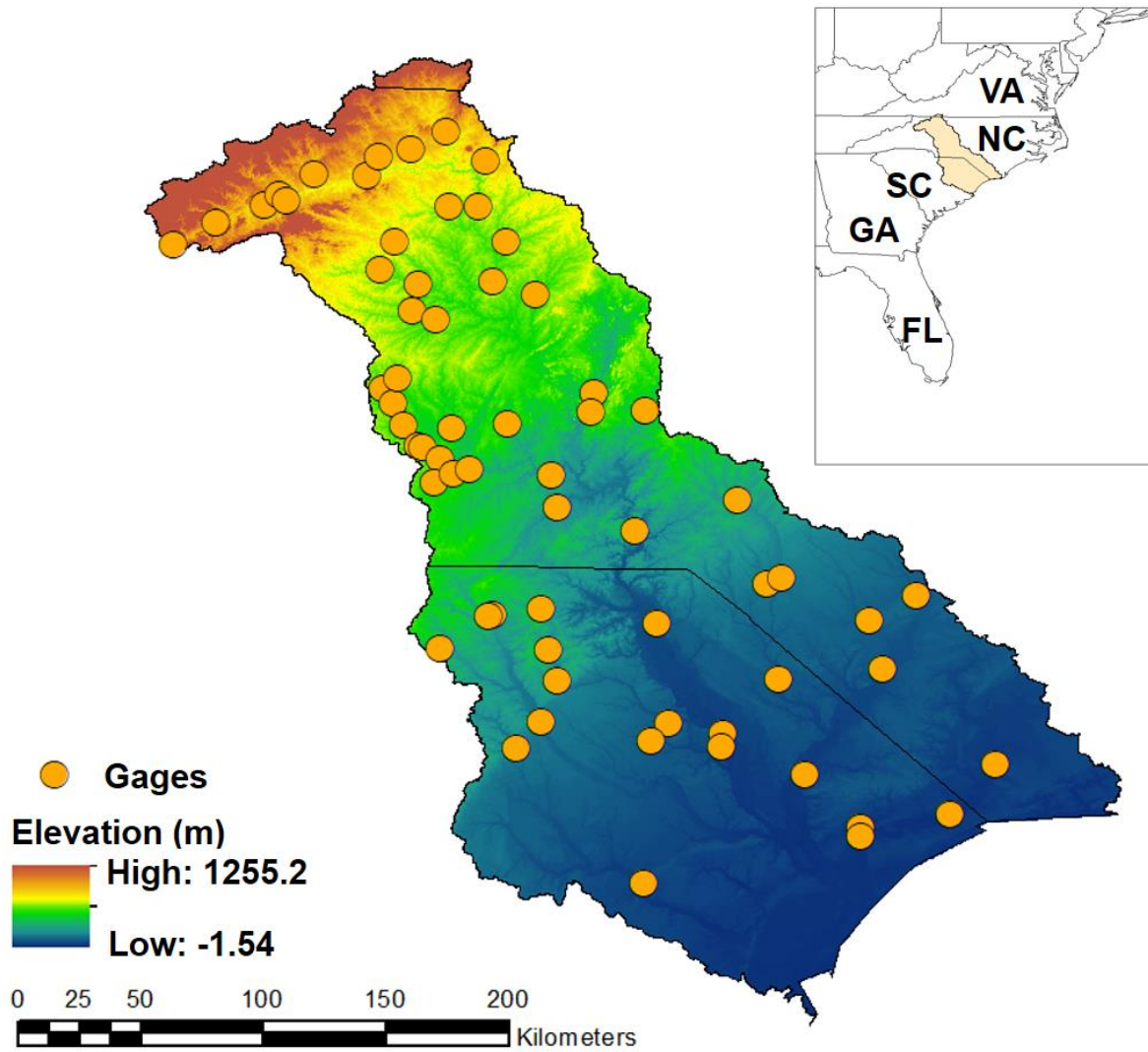
886

887



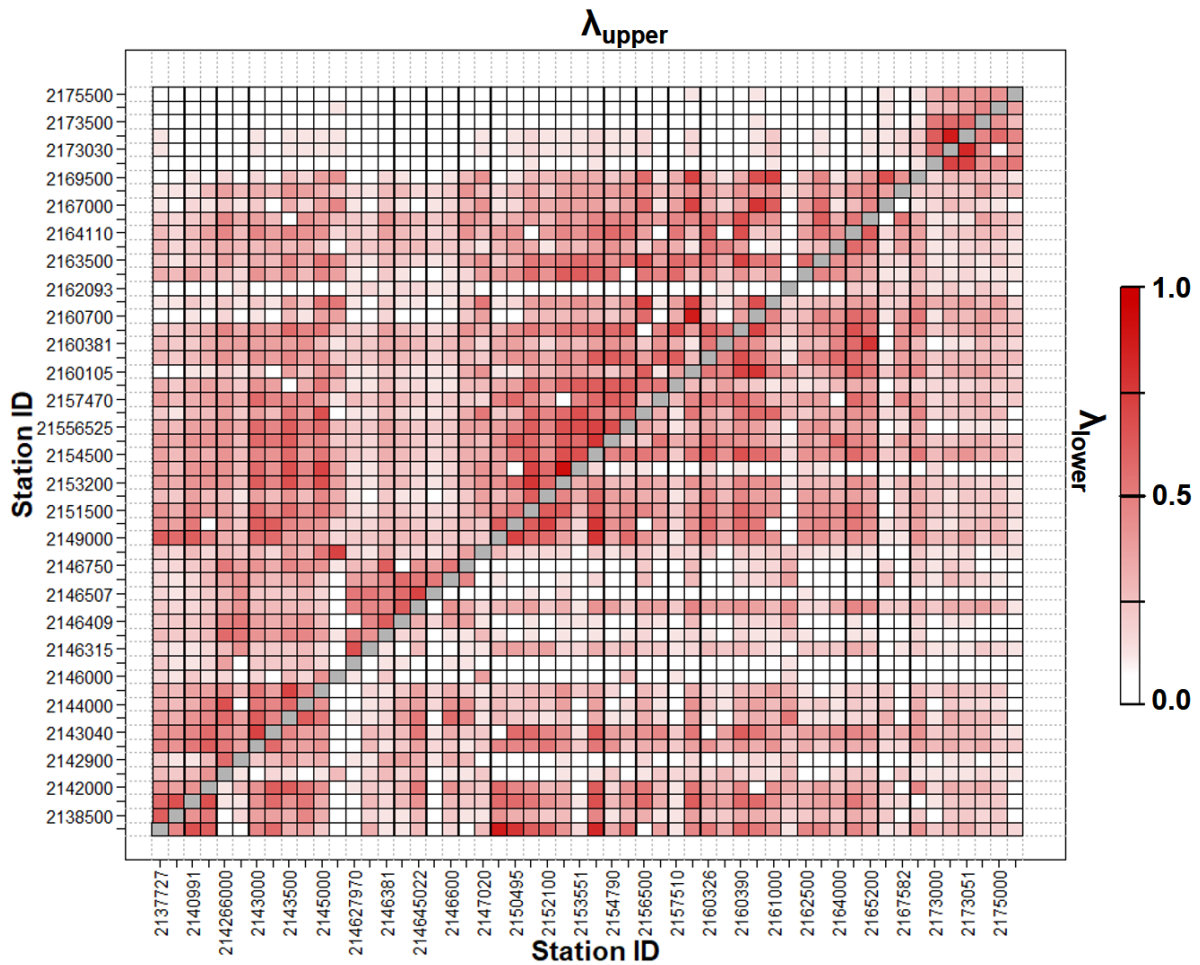
889
 890
 891
 892
 893
 894
 895
 896
 897
 898
 899
 900
 901
 902
 903
 904
 905
 906
 907
 908
 909
 910

Figure 2 Structure of the 6-Dimensional vine model and marginal for the synthetic simulation. $LN(\pi, \sigma^2)$ denotes the log normal distribution with its mean (π) and variance (σ^2). The target gage is highlighted.



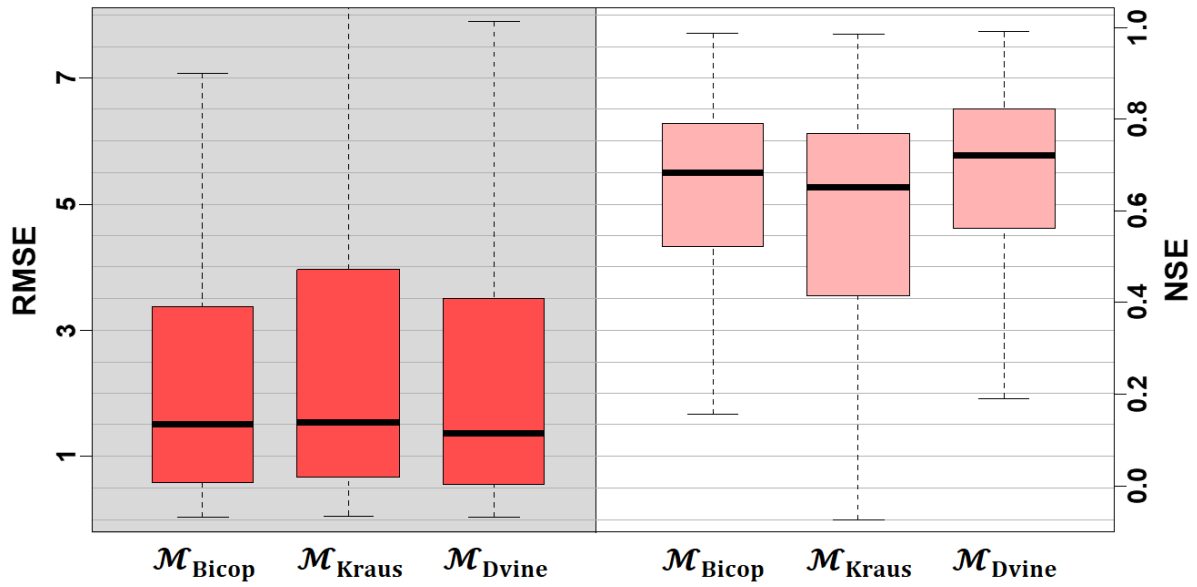
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929

Figure 3 Map of the Yadkin-Pee Dee Basin with 54 stream gage stations



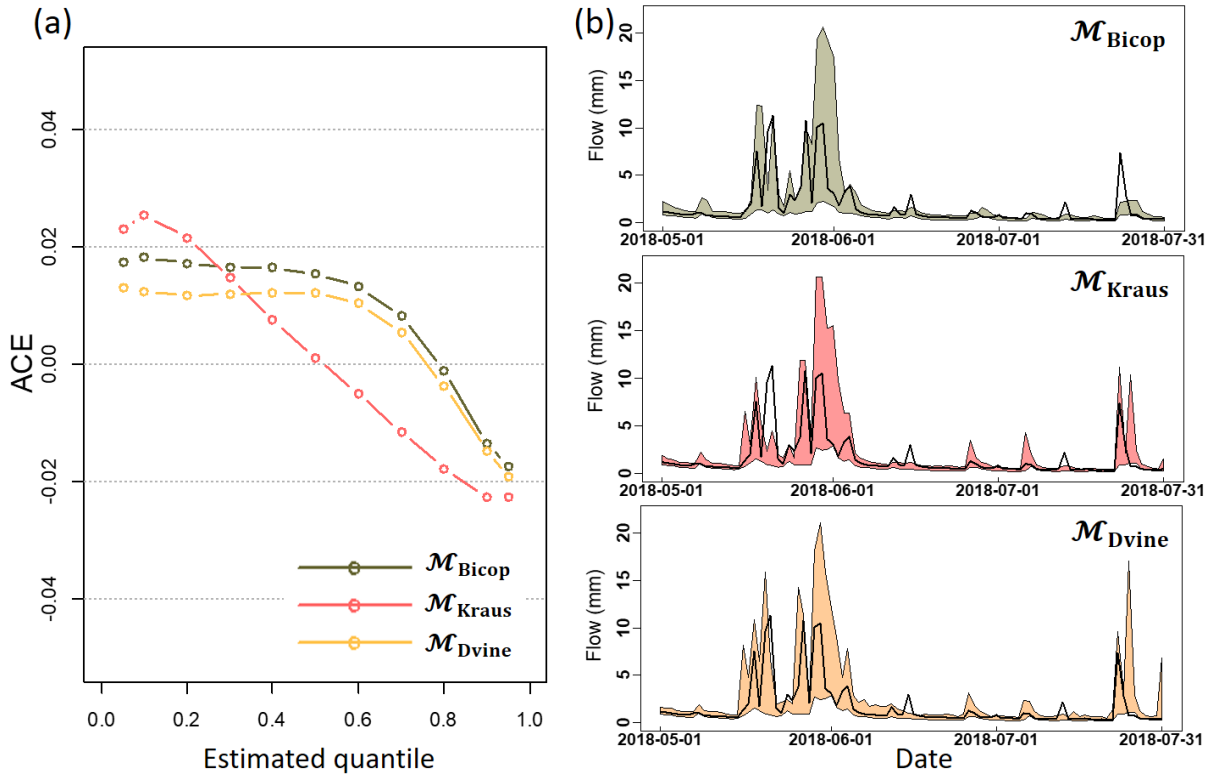
931
 932
 933
 934
 935
 936
 937
 938
 939
 940
 941
 942
 943
 944
 945
 946
 947
 948
 949
 950
 951
 952

Figure 4 Pairwise upper and lower tail dependence for watersheds in the Yadkin-Pee Dee River Basin. The upper triangular matrix shows values for the upper-tail dependence and the lower triangular matrix presents values for the lower-tail dependence. The metrics can range from 0 to 1, with higher values suggesting greater interdependence of two streamflows for each upper- and lower-tail.



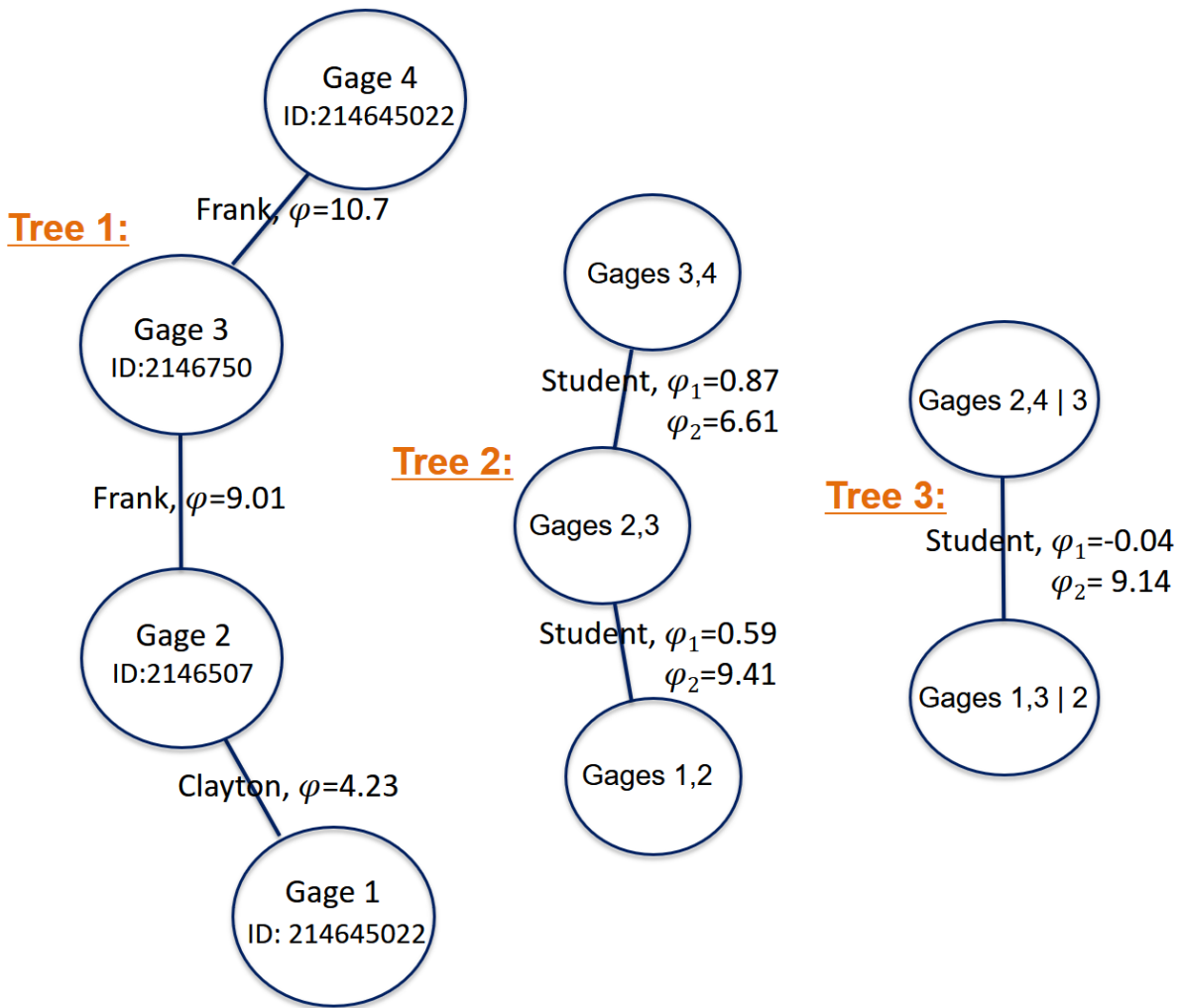
954 Figure 5 Model performance for the Yadkin-Pee Dee river under a cross-validation framework
 955 based on RMSE (dark squares) and NSE (light squares). Here, the RMSE (NSE) can range
 956 from 0 to ∞ ($-\infty$ to 1), with lower RMSE (higher NSE) implying better performance.
 957

958
 959
 960
 961
 962
 963
 964
 965
 966
 967
 968
 969
 970
 971
 972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985



987
988 Figure 6 (a) Average coverage error from three copula-based models for the Yadkin-Pee Dee
989 River Basin across exemplarily quantiles, and (b) 95% PIs for three models for an example
990 period (1 May 2018 to 31 July 2018) for a specific target gauge (USGS site ID: 02143500).
991 Observed streamflow (black solid line) is also presented in each figure.

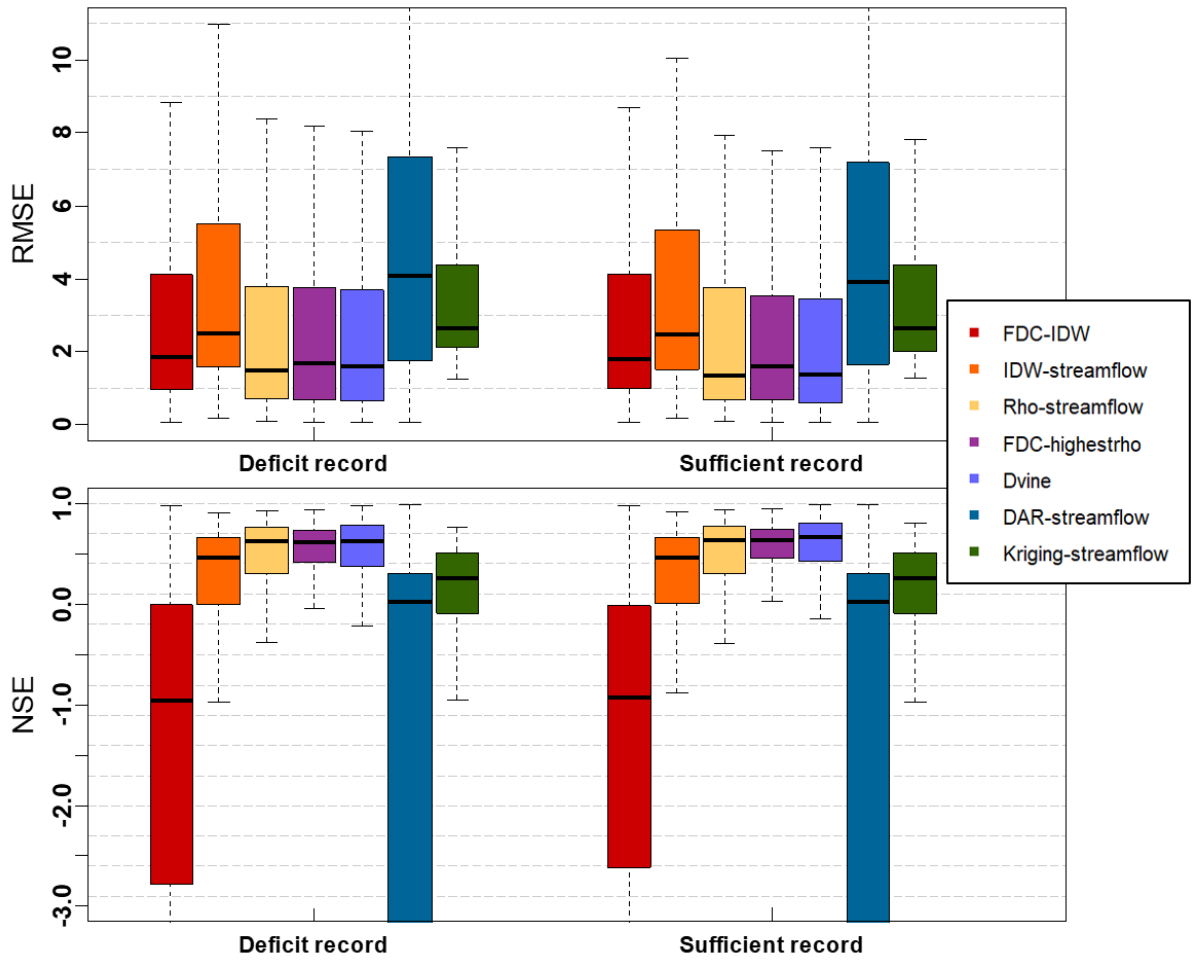
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013



1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025
 1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034

Figure 7 Structure of the Dvine copula applied for a particular target site (USGS site ID: 214645022) with the defined bivariate copulas and their parameters.

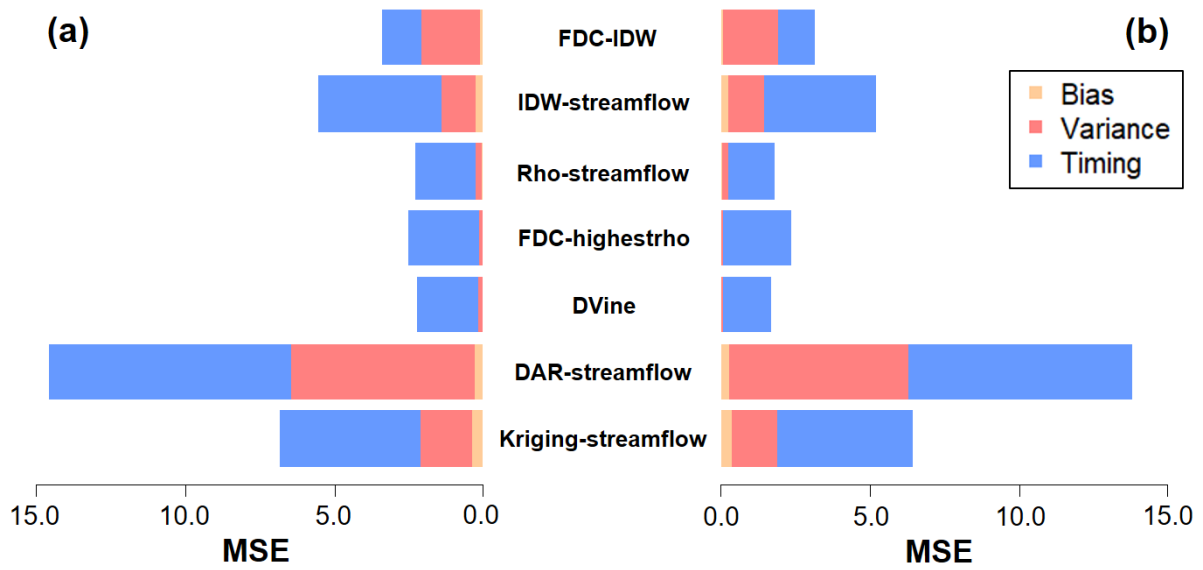
1035



1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057

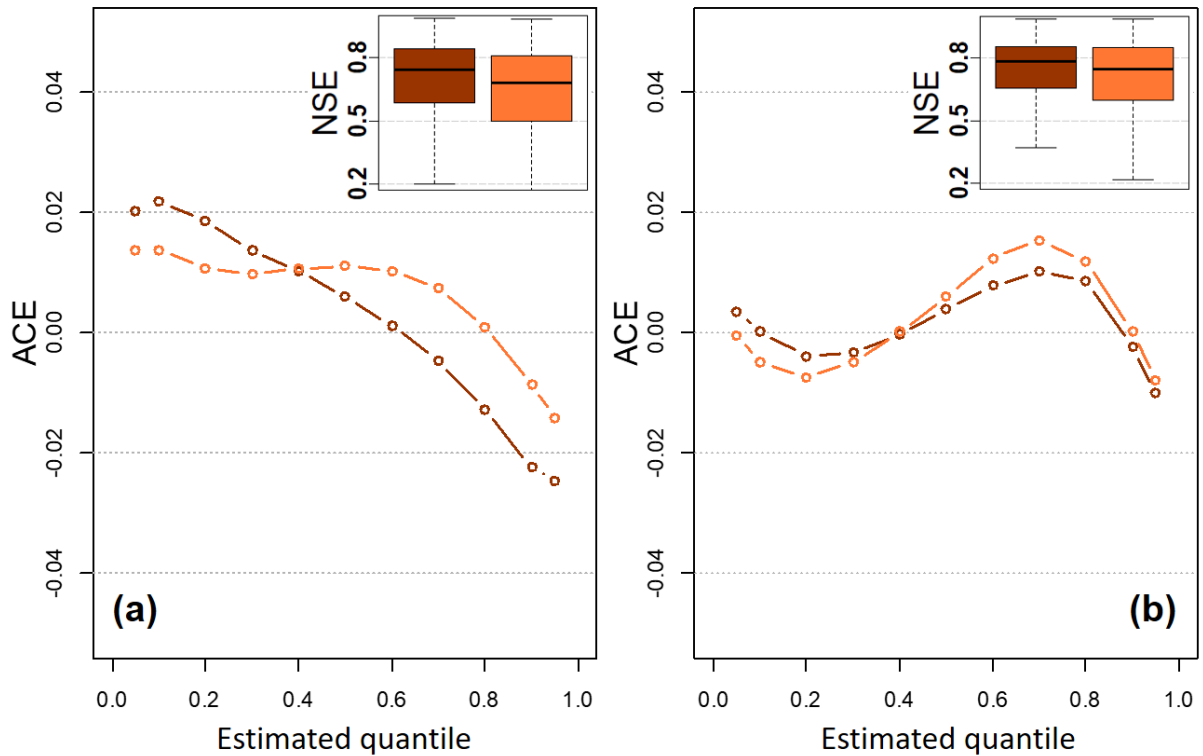
Figure 8 Inter-model comparison using cross-validation experiments based on RMSE (upper) and NSE (lower). Here, lower RMSE suggests more accurate estimations for infilling missing values.

1058
1059



1060
1061 Figure 9 Three contributions from the decomposed mean squared error (MSE) for the cross-
1062 validation experiment with (a) the deficit record and (b) sufficient record scenarios.
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090

1091
1092



1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118

Figure 10 Average coverage error of the Dvine model for two scenarios under (a) the “deficit” and (b) “sufficient” cases. In each case, the dark line represents the scenario by the marginal cumulative probabilities using all years and conditional streamflow values constructed from the partial record. On the other hand, the light line illustrates the scenario by the marginal cumulative probabilities estimated by the partial record and conditional streamflow values constructed from the full record. Inset: NSE performance of the Dvine model for the two scenarios in each case.

1119
1120
1121

Table 1 Seven infilling approaches discussed in the study

No.	Method	Description
1	FDC-IDW	Inverse distance-weighted estimate of non-exceedance probability from those of all donors.
2	IDW-streamflow	Inverse distance-weighted estimate using streamflow from all donors.
3	Rho-streamflow	Correlation-weighted streamflow estimate from the selected donors for each time step. The optimal number of donors is determined in a cross-validation framework.
4	FDC-highestrho	Estimate non-exceedance probability from the gage with the highest correlation.
5	DAR-streamflow	Drainage-area (DA) ratio for streamflow using the DA from the nearest neighbor gage.
6	Kriging-streamflow	Geostatistical interpolation method to estimate streamflow from all donors for each time step.
7	DVine	Vine copula-based estimate from the selected donors

1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143

1144
 1145
 1146
 1147
 1148

Table 2 RMSE and NSE results over the validation periods under synthetic experiment for comparing copula-based model formulations. Best metric values for each quantile are italicized and bolded.

Metric	Model formulation	Min	First quantile	Median	Third quantile	Max
Root mean squared error (RMSE)	$\mathcal{M}_{\text{Bicop}}$	0.912	1.119	1.258	1.363	3.353
	$\mathcal{M}_{\text{Kraus}}$	0.990	1.140	1.386	1.660	4.273
	$\mathcal{M}_{\text{Dvine}}$	<i>0.895</i>	<i>1.046</i>	<i>1.112</i>	<i>1.391</i>	4.119
Nash-Sutcliffe efficiency (NSE)	$\mathcal{M}_{\text{Bicop}}$	<i>0.464</i>	0.779	0.826	0.856	0.902
	$\mathcal{M}_{\text{Kraus}}$	0.198	0.724	0.782	0.825	0.885
	$\mathcal{M}_{\text{Dvine}}$	0.248	<i>0.805</i>	<i>0.838</i>	<i>0.869</i>	<i>0.905</i>

1149
 1150
 1151
 1152
 1153
 1154
 1155
 1156
 1157
 1158
 1159
 1160
 1161
 1162
 1163
 1164
 1165
 1166
 1167
 1168
 1169
 1170
 1171
 1172

1173
1174
1175
1176
1177

Table 3 Results of average coverage error (ACE) over the validation periods under synthetic experiment for comparing copula-based model formulations. Best metric values for each quantile are italicized and bolded.

Model formulation	Estimated quantile (ϖ)				
	0.05	0.10	0.50	0.90	0.95
$\mathcal{M}_{\text{Bicop}}$	0.027	0.063	0.079	0.014	0.002
$\mathcal{M}_{\text{Kraus}}$	<i>0.003</i>	<i>0.011</i>	0.055	0.024	0.001
$\mathcal{M}_{\text{Dvine}}$	0.029	0.048	<i>0.042</i>	<i>0.001</i>	<i>0.000</i>

1178
1179
1180
1181
1182
1183
1184
1185
1186