1

2 **Streamflow estimation at partially gaged sites using multiple**

3 **dependence conditions via vine copulas**

4

5 Kuk-Hyun Ahn[1]

6

7

8

9

10

11

12

13

14

15

16

17

18

19

20

21

22 [1]Assistant Professor, Department of Civil and Environmental Engineering, Kongju National
23 University, Cheon-an, South Korea; *Corresponding author;* e-mail: ahnkukhyun@gmail.com

24

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

25

26                                         ABSTRACT

27

28      Reliable estimates of missing streamflow values are relevant for water resources planning and

29      management. This study proposes a multiple dependence condition model via vine copulas for

30      the purpose of estimating streamflow at partially gaged sites. The proposed model is attractive

31      in modeling the high dimensional joint distribution by building a hierarchy of conditional

32      bivariate copulas when provided a complex streamflow gage network. The usefulness of the

33      proposed model is firstly highlighted using a synthetic streamflow scenario. In this analysis,

34      the bivariate copula model and a variant of the vine copulas are also employed to show the

35      ability of the multiple dependence structure adopted in the proposed model. Furthermore, the

36      evaluations are extended to a case study of 54 gages located within the Yadkin-Pee Dee River

37      Basin, the eastern U. S. Both results inform that the proposed model is better suited for infilling

38      missing values. After that, the performance of the vine copula is compared with six other

39      infilling approaches to confirm its applicability. Results demonstrate that the proposed model

40      produces more reliable streamflow estimates than the other approaches. In particular, when

41      applied to partially gaged sites with sufficient available data, the proposed model clearly

42      outperforms the other models. Even though the model is illustrated by a specific case, it can be

43      extended to other regions with diverse hydro-climatological variables for the objective of

44      infilling.

45

46      **Keywords: vine copulas, multiple dependence condition model, streamflow estimation**

47               **and infilling approach**

48

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 ungagged 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 watershed(s) for missing streamflow

79    values; and (4) converting corresponding streamflow values back from the transferred non-

80    exceedance probability. When FDCs are utilized for partially gaged watersheds, how the donor

81    watersheds are selected (step 2) and how the probabilities are transferred from the donor

82    watersheds (step 3) play crucial roles 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 paid to copulas in the field of hydrology, with applications in

99  flood frequency analysis, drought risk analysis, and multi-site streamflow simulations (Ahn

100  and Palmer, 2016b; Ariff et al., 2012; Chen et al., 2015; Daneshkhah et al., 2016; Fu and Butler,

101  2014). Copulas are efficient 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, 1997). 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) although they have not been introduced to

129    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 $\boldsymbol{q} =$

159   $(q_1, \dots, q_k)^T$, the $C$ satisfies:

160

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

169

170    where $c$ is a $k$-dimensional copula density acquired by partial differentiation of the copula $C$

171    (i.e., $c(F_1(q_1), \dots, F_k(q_k)) := \frac{\partial^k}{\partial_1 \cdots \partial_k} C(F_1(q_1), \dots, F_k(q_k))$) 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

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

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

193     $$\prod_{j=1}^{k-1} \prod_{\jmath=1}^{k-j} c_{j,j+\jmath|(j+1):(j+\jmath-1)}(F(q_\jmath|q_{\jmath+1}, \ldots, q_{\jmath+j-1,}), F(q_{\jmath+j}|q_{\jmath+1}, \ldots, q_{\jmath+j-1,}))$$     Eq. (3)

194

195     where $c_{j,j+\jmath|(j+1):(j+\jmath-1)}$ 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     $$f(q_1, \ldots, q_5) = [\prod_{i=1}^{5} f_i(q_i)] c_{12} \cdot c_{23} \cdot c_{34} \cdot c_{45} \cdot c_{13|2} \cdot c_{24|3} \cdot c_{24|3} \cdot c_{35|4} \cdot c_{14|23} \cdot c_{25|34} \cdot$$

200     $$c_{15|234}$$                                                                          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 $F(q_{\dot{j}}|q_{\dot{j}+1}, \ldots, q_{\dot{j}+j-1})$, 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 $q_{\dot{j}+1}, \ldots, q_{\dot{j}+j-1}$ and $\boldsymbol{v} = \{q_{\dot{j}+1}, \ldots, q_{\dot{j}+j-1}\}\backslash q_i$ the streamflow vector without $q_i$ used in the

209 following recursive relationship (Aas et al., 2009):

210

211 $$h\big(q_{\dot{j}}\big|\boldsymbol{v}\big) := F\big(q_{\dot{j}}\big|\boldsymbol{v}\big) = \frac{\partial C_{\dot{j}i|\boldsymbol{v}}(F(q_{\dot{j}}|\boldsymbol{v}), F(q_i|\boldsymbol{v}))}{\partial F(q_i|\boldsymbol{v})}$$ Eq. (5)

212

213 where the $h$-function is associated with the pair-copula $C_{\dot{j}i|\boldsymbol{v}}$.

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 ($\mathcal{M}_{\mathrm{Dvine}}$)*

217 Following Kraus and Czado (2017), a two-step estimation procedure is adopted for the

218 prediction of the streamflow value at the target watershed. The algorithm ($\mathcal{M}_{\mathrm{Dvine}}$) 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     $q_1, \ldots, q_{k-1}$ from the donor sites. In the first step, the marginal cumulative probabilities

224     $F_k(q_k)$ and $F_j(q_j)$, $j = 1, \ldots, k-1$, 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 $q_i^\zeta$, $\zeta = 1, \ldots, \xi$, at $i$th site, defined as $\widehat{F}_i(q_i) = \frac{1}{nh} \sum_{\zeta=1}^{\xi} \Omega(\frac{q_i - q_i^\zeta}{h})$.

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 $F_k(q_k) - F_{I_1}(q_{I_1}) - F_{I_2}(q_{I_2}) - \ldots - F_{I_{k-1}}(q_{I_{k-1}})$, such that $F_k(q_k)$ is

234     the first node in the first tree and the other orders of donors ($I_1$, ..., $I_{k-1}$) 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) to

238     represent diverse dependence structures. For example, a Gaussian copula is proper when the

239     non-exceedance probabilities between two watersheds are associated in the body of their

240     distribution but are not asymptotically dependent in the both tails. On the other hand, a Gumbel

241     copula may be appropriate for the situation wherein the non-exceedance probabilities exhibit

242     tail dependence, where high flows are connected by same rainfall events but low flows are not

243     correlated (e.g., due to regulation) (Salvadori and De Michele, 2004). Details of the five

244     bivariate copulas are presented in the Supporting Information. Parameters for the five bivariate

245  copulas are estimated based on Kendall rank-based correlation ($\rho^\tau$) between sites. The optimal

246  bivariate copula for each pair copula is determined based on the panelized likelihood function

247  (i.e., AIC).

248

249  The final number ($\chi_k$) of donor sites is further optimized under a cross-validation approach. In

250  this approach, 80 % of the regional data are employed for model fitting; the other 20 %, for

251  testing. Again, this procedure is conducted 5 times, each time using a different set of data for

252  testing. As a measure for the model's fit, the root mean squared error (RMSE; equation (6))

253  from observed streamflow at the target site is utilized.

254

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

256

257  Finally, conditional streamflow values at the target site can be estimated using the inverse form

258  of the conditional distribution function (i.e., Eq. 5). To depict the ideas, a trivariate case (i.e.,

259  $\chi = 2$) is considered here. Based on the streamflow values at the donor sites ($q_2, q_3$), $\widehat{q_1}$ can

260  be obtained using the conditional distribution function $h(q_1|q_2, q_3)$. For some fixed

261  probabilities $\phi$ (e.g., $\phi = 0.1, \ldots, 0.9$), $F_1(\widehat{q_1})$ is derived from $C_{1|2,3}$ using an explicit

262  function:

263

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

265

266    where $C_{1|2,3}^{-1}$ is the inverse of the copula function given the $\phi$ quantile curve of the copula

267    (Liu et al., 2015; Xu and Childs, 2013). Therefore, the $\phi$th copula-based conditional quantile

268    function of streamflow at the target site can be calculated as follows:

269

270    $$q_1(\phi|q_2q_3) = F_1^{-1}(C_{1|2,3}^{-1}(\phi|F_2(q_2), F_3(q_3))) =$$

271    $$F_1^{-1}(h_{1|2}^{-1}(h_{1|32}^{-1}(\phi|h_{2|1}(F_2(q_2)|F_1(q_1)))|F_1(q_1)))$$          Eq. (8)

272

273    Similarly, for the $k$-dimensional case, the $\phi$th copula-based conditional quantile function can

274    be calculated along with streamflow at the $k$-1 donor sites. To acquire an estimate at the target

275    site, 1000 samples from uniform distribution over the interval [0, 1] are generated using Monte

276    Carlo simulations. In this study, the mean value of these generations is regarded as the best

277    estimate.

278

279    **3. Application**

280    This study first explores the performance of $\mathcal{M}_{\text{Dvine}}$ under synthetic example. In this analysis,

281    $\mathcal{M}_{\text{Bicop}}$ and $\mathcal{M}_{\text{Kraus}}$ are also employed to show the usefulness of $\mathcal{M}_{\text{Dvine}}$. For $\mathcal{M}_{\text{Bicop}}$, the

282    optimal bivariate copula is selected based on the AIC while the five bivariate copulas (Gaussian,

283    Student-t, Frank, Gumbel, and Clayton copulas) are considered as its potential candidates. A

284    brief description of two additional models are presented in the supporting information. After

285     that, those three models are used for a real application to 54 stream gages located in a region

286     of the eastern United States by estimating streamflow in partially gaged locations. Finally,

287     seven infilling approaches (Table 1) are also utilized and evaluated in a cross-validated

288     framework to evaluate the performance of the proposed model.

289

290     *3.1 Synthetic simulation*

291     Synthetic streamflow data are generated using controlled Monte Carlo experiment to explore

292     how well the three copula-based models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$, $\mathcal{M}_{\text{Dvine}}$) provide streamflow

293     predictions at the target site given a complex streamflow data in a pseudo gage network. In this

294     analysis, a six-dimensional streamflow set ($q_1^\zeta$, $q_2^\zeta$, $q_3^\zeta$, $q_4^\zeta$, $q_5^\zeta$, $q_6^\zeta$), $\zeta = 1, \dots, \xi =$

295     2190 (i.e. $\frac{2190}{365} = 6$ years), is modelled using four bivariate copulas (Gaussian, Student-t,

296     Flank, and Clayton copulas) and lognormal distributions for margins (see Figure 2).

297

298     The performance of each model is evaluated in a calibration-validation framework. First,

299     synthetic streamflow data are generated for six-dimensional gage network. Then, $\varphi$ years of

300     data are randomly selected to be assumed known at the target gage, and the streamflow for the

301     remaining 6-$\varphi$ years of data is then estimated as missing values ($\varphi = 4$ in this analysis). This

302     process is repeated 20 times to build an ensemble prediction. In particular, this study assumes

303     the fifth streamflow data (i.e., $q_5$) to be predicted. In this assessment, two characteristics are

304     considered to compare the three models: model prediction reliability and uncertainty

305     quantification skill. Model prediction reliability is tested using the root mean squared error

306     (RMSE; Eq. 6) and Nash-Sutcliffe efficiency (NSE), which are further described in Section 3.4.

307   Uncertainty quantification skill is judged by the ability of each model to build prediction

308   intervals (PIs) that correctly bound predictions (see Section 3.4). Here, coverage probabilities,

309   defined as the proportion of the time that true values occur into these PIs, are employed to show

310   the usefulness of the proposed model.

311

312   *3.2 Application to the Yadkin-Pee Dee River*

313   The Yadkin-Pee Dee River Basin (Figure 3), covering around 18,700 km$^2$ and one of the largest

314   river basins in North Carolina and South Carolina (Fisk, 2010), is used as real data to evaluate

315   infilling ability. The basin flows from the northwestern corner of North Carolina near Blowing

316   Rock and extends south by southeast, crossing the south-central border of North Carolina into

317   South Carolina, with slightly more than half of its watershed in North Carolina. Most of the

318   land covered within the basin is forested or used for agriculture although urban areas of the

319   basin are expanding.

320

321   Daily streamflow data at 54 gages are gathered throughout the study region from web interface

322   of the U.S. Geological Survey (USGS) National Water Information System (NWIS) (U.S.

323   Geological Survey, 2018). The 54 gages are selected based on the following criteria: (1) all

324   gages are recorded continuously for 15 years of daily streamflow over the period from January

325   2004 to December 2018, and (2) gages have non-zero daily values for the period in the first

326   criterion since gages with streamflow values equal to zero require a more flexible modeling

327   structure. Thus, it is common to model zero flows separately in regionalization studies. Based

328   on the second criterion, 10 gages are discarded (not shown).

329

### *3.3 Intermodel comparison framework*

331      A set of seven infilling approaches is used in the final assessment (see Table 1): (1) $\mathcal{M}_{\text{FDC–IDW}}$,

332      (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)

333      $\mathcal{M}_{\text{Kriging–streamflow}}$, and (7) $\mathcal{M}_{\text{Dvine}}$. This set of seven models is tested in a cross-validation

334      framework under two different cases. The two cases consider situations wherein $\varphi$ have

335      values of 2 and 8 to represent relatively deficit- and sufficient-records for the target site. Similar

336      to the comparative assessment to show the usefulness of the proposed copula-based model (see

337      Section 3.1), each case is repeated 20 times by randomly selecting $\varphi$ years over the applied

338      period. The reliability of each model is evaluated using RMSE and NSE metrics over the

339      validated four-year period randomly selected in the remaining data (i.e., 4 years in 15-$\varphi$ years).

340

### *3.4 Error metrics and error decomposition*

342      As presented in Sections 3.1 and 3.3, the root mean squared error (RMSE; Eq. 6) and Nash-

343      Sutcliffe efficiency (NSE) are employed to evaluate prediction skills:

344

345          $$NSE = 1 - \frac{\sum_{\zeta=1}^{\xi}(\widehat{q^{\zeta}}-q^{\zeta})^2}{\sum_{\zeta=1}^{\xi}(q^{\zeta}-\overline{q^{\zeta}})^2}$$        Eq. (9)

346

347      The NSE (RMSE) can range from $-\infty$ to 1 (0 to $\infty$), with higher NSE (lower RMSE) implying

348      better performance. Both metrics have been commonly used in hydrology analysis (Boyle et

349    al., 2000).

350

351    Following derivations suggested in Gupta et al. (2009), the RMSE can be further decomposed

352    into three components:

353

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

355

356    where $\mu$ $(\hat{\mu})$ and $\sigma$ $(\hat{\sigma})$ represent the average and standard deviation for the observed

357    (estimated) streamflow, respectively, and $r$ indicates the estimated correlation coefficient.

358    The first component $(\hat{\mu} - \mu)^2$ is a measure of how well the average of the observed

359    streamflow represents the average of the estimated streamflow; the second component

360    $(\hat{\sigma} - \sigma)^2$ is a measure of how well the variance of the prediction represents the variance of the

361    observed streamflow; and the third component $2\sigma\hat{\sigma}(1 - r)$ is dominated by the correlation

362    and is defined as the "timing" component (Worland et al., 2019). Using these three defined

363    components, their absolute contributions are explored in this study.

364

365    In addition, the accuracy of the uncertainty quantification skill is also evaluated for the copula-

366    based models ($\mathcal{M}_{Bicop}$, $\mathcal{M}_{Kraus}$, $\mathcal{M}_{Dvine}$). To be specific, this study utilizes the PI coverage

367    probability (PICP), which a common metric for this purpose (He et al., 2017; Niemierko et al.,

368    2019). It provides the relative number of data points that fall between the defined bounds as

369    expressed follows:

370

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

372

373  where $\Theta_\zeta$ is the indicator variable if $q^\zeta$ is covered by the $\zeta$th PI defined by the lower bound

374  $L^\zeta$ and upper bound $U^\zeta$. This study examines the prediction accuracy of single quantiles.

375  Therefore, the lower bound is defined as $L^\zeta = -\infty$ and $U^\zeta = \widehat{q^{\zeta,\varpi}}$ where $\varpi$ is the

376  estimated quantile at time $\zeta$. Accordingly, the upper bound is not a constant, but is re-assigned.

377  By subtracting the nominal confidence $\varpi$ from PICP, the average coverage error (ACE) is

378  obtained as follows:

379

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

381

382  The metric clearly indicates if the predicted quantile is underestimated (ACE < 0) or

383  overestimated (ACE > 0) while taking small values around 0 for ideal case.

384

385  **4. Results**

386  *4.1 Results for synthetic experiment*

387  Prediction results from out-of-samples for the RMSE and NSE metrics are presented for the

388  three copula-based models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$, $\mathcal{M}_{\text{Dvine}}$) in Table 2. The ACE scores are also

389    described for $\varpi \in \{0.05, 0.10, 0.50, 0.90, 0.95\}$ in Table 3. When compared to the other

390    models, $\mathcal{M}_{\text{Bicop}}$ achieves lower RMSE values in the right tail of the RMSE distribution over

391    the validation periods, but severely underperforms the majority of the designed experiment,

392    suggesting this model formulation relying on a single donor leads to poor predictions. $\mathcal{M}_{\text{Kraus}}$

393    provides higher RMSE values for all the RMSE distribution, particularly for the right tail of

394    the RMSE distribution. The model utilizes streamflow data from all donors (i.e., five donor

395    sites) although the first two gages (Gages 1 and 2) show insignificant associations to the target

396    site ($r_{1,5} = 0.11$ and $r_{2,5} = 0.14$). $\mathcal{M}_{\text{Dvine}}$ unequivocally produces the best predictions.

397    $\mathcal{M}_{\text{Dvine}}$ adopts streamflow data from two or three donors (Gages 3, 4 and 6) without utilizing

398    streamflow data from the first two donors when a multiple dependence structure is established

399    to build an ensemble prediction. It outperforms $\mathcal{M}_{\text{Bicop}}$ and $\mathcal{M}_{\text{Kraus}}$ across all validation

400    periods, besides a few with the worst performance. Even in this case, the maximum RMSE of

401    $\mathcal{M}_{\text{Dvine}}$ is fairly less than the maximum RMSE of $\mathcal{M}_{\text{Kraus}}$.

402

403    In addition, the ACE results present how the three models characterize prediction uncertainty.

404    $\mathcal{M}_{\text{Dvine}}$ is capable of properly covering the predications across the entire distribution while

405    slight overestimation occurs for the smallest two quantiles. The remaining upper quantiles also

406    tend to slightly overestimate the true values but the overestimations are less than the other

407    models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$). Taken together, the results of the synthetic experiment suggest that

408    $\mathcal{M}_{\text{Dvine}}$ yields the best predictions among the copula-based models tested.

409

410    *4.2 Performance of the copula-based models in the Yadkin-Pee Dee River*

411     Using the insights developed from the synthetic experiment above, the three copula-based

412     models are applied to the streamflow data for the Yadkin-Pee Dee River. At first, upper and

413     lower tail dependences ($\lambda_{upper}$ and $\lambda_{lower}$) are examined for all two pairs of sites (see Figure

414     4) using the approach of Schmid and Schmidt (2007). Theoretical background is described in

415     the Supporting Information (Text S3). Note that in this analysis, the dependences become more

416     obvious as the values approach unity. Two major insights emerge from this figure. First, many

417     site-pairs exhibit strong upper tail dependence, suggesting that streamflow variability has a

418     tendency to be more correlated under high-flow conditions compared to under low-flow

419     conditions (i.e., asymmetric dependence). The lack of lower-tail dependence may be due to

420     contributions governing low streamflow such as river regulation. Next, even under high- or

421     low-flow conditions, there is a wide range of tail dependence across the study basin (i.e.,

422     heterogeneous dependence). To sum up, a wide range of complex dependencies is observed in

423     the streamflow data over the study basin. The complex dependences suggest, when streamflow

424     is estimated from multiple donors, the potential usefulness of considering a multiple

425     dependence structure, which is one of the main features of vine copulas.

426

427     Figure 5 shows the RMSE and NSE results for the three copula-based models under a "leave-

428     one-out" cross validation framework. This process is repeated 20 times to build an ensemble

429     prediction by using test periods randomly defined. For this analysis, five years of data are

430     selected to be assumed as the observed period at the target gage, and another four years are

431     randomly selected in the remaining data for the test period. Similar to the results from the

432     synthetic experiment, $\mathcal{M}_{Kraus}$ performs poorly in both the RMSE and NSE metrics (median

433     RMSE = 1.549 and NSE = 0.652). The bivariate copula performs well (median RMSE =

434    1.496), indicating that this approach efficiently leverages available information even though

435    the information is limited to single donor. Particularly, $\mathcal{M}_{\text{Bicop}}$ achieves lowest RMSE values

436    in the upper side of the RMSE box (e.g., third quartile), providing a strong uncertainty

437    quantification skill for the upper bound. However, $\mathcal{M}_{\text{Dvine}}$ yields the best median RMSE and

438    NSE values (= 1.359 and 0.719). Given the heterogeneous dependence conditions (see Figure

439    4), the high dimensional structures are effective in modeling a complex streamflow gage

440    network. This feature can substantially improve prediction of target site flows.

441

442    Figure 6a presents the ACE scores described for principal quantiles, $\varpi \in$

443    $\{0.05, 0.10, 0.20, \dots, 0.90, 0.95\}$, across all target sites under the cross validation framework.

444    Figure 6b presents 95% PIs for each model for an example time period (1 May 2018 to 31 July

445    2018) for one target site (USGS site ID: 02143500). Note that the ACE would ideally take zero

446    value, regardless of the quantiles. The ACE scores for the three models ($\mathcal{M}_{\text{Bicop}}$, $\mathcal{M}_{\text{Kraus}}$,

447    $\mathcal{M}_{\text{Dvine}}$) range from 0.004 to 0.0007 when considering all the quantiles together. However, the

448    scores vary depending on the quantiles. For instance, the ACE score for $\mathcal{M}_{\text{Kraus}}$ is noticeably

449    positive but is almost zero around the median streamflow, indicating that the model properly

450    represent uncertainty of the median streamflow. $\mathcal{M}_{\text{Bicop}}$ and $\mathcal{M}_{\text{Dvine}}$ result in very similar

451    ACE scores although $\mathcal{M}_{\text{Dvine}}$ performs slightly better than $\mathcal{M}_{\text{Bicop}}$. The differences in

452    characterization of prediction uncertainty can be confirmed from a particular target site (Figure

453    6b).

454

455    Based on the results in Figures 5 and 6, $\mathcal{M}_{\text{Dvine}}$ has the most reliable overall performance (as

Hydrology and
Earth System
Sciences
Discussions

Open Access

EGU

456 judged by model prediction reliability and uncertainty quantification skill) and is selected as

457 an appropriate copula model to infill missing data in partially gaged. Figure 7 shows an

458 example application of $\mathcal{M}_{\text{Dvine}}$ including the optimal donor sites, proper bivariate copulas

459 and their parameters for one target site (USGS site number #214645022) when the model is

460 calibrated using the full 15-year record.

461

462 *4.3 Intermodel comparison for streamflow estimation*

463 To assess the predictive skill of the proposed vine copula model, it is compared with six other

464 statistical models (see Table 1). Figure 8 shows RMSE and NSE for the seven models where

465 the streamflow values are estimated based on the available data defined by the two different

466 cases, labeled "deficit record" and "sufficient record" (see Section 3.3). Under all cases, the

467 vine copula approach outperforms the other infilling approaches. For example, for the

468 "sufficient record" case, median NSE for $\mathcal{M}_{\text{Dvine}}$ is 0.673 whereas those for

469 $\mathcal{M}_{\text{IDW−streamflow}}$ and $\mathcal{M}_{\text{rho−streamflow}}$ are 0.462 and 0.649, respectively. In this analysis, the

470 approaches, which are based on streamflow values of the donor sites without utilizing non-

471 exceedance probability including DAR-streamflow and Kriging-streamflow, yield relatively

472 increased bias in their predictions. On the other hand, an application of FDC models offers

473 reliable predictions. For instance, for the "sufficient record" case, median RMSE for

474 $\mathcal{M}_{\text{FDC−highestrho}}$ is 1.603 compared to that of a direct of using streamflow (e.g., median RMSE

475 of $\mathcal{M}_{\text{FDC−streamflow}}$ = 3.422 for the sufficient record). Similar interpretation can be found in

476 the comparison between $\mathcal{M}_{\text{FDC−IDW}}$ and $\mathcal{M}_{\text{IDW-streamflow}}$. The results from these approaches

477 suggest that utilizing FDC process leads to a reliable estimation, which is a primary structure

478 in the vine copula. The other noticeable feature is that available data length provides a

479   significant influence on performance of some infilling methods. In particular, this is quite

480   evident for the vine copula model (median RMSEs: 1.598 and 1.379 for deficit and sufficient

481   records, respectively).

482

483   *4.4 Prediction error decomposition*

484   The RMSE is decomposed into their components (bias, variance, and timing components) for

485   both the "deficit record" and "sufficient record" predictions (Figure 9). For the both cases,

486   prediction errors for all seven models are caused largely by timing components. In particular,

487   models estimating directly streamflow values (IDW-streamflow, DAR-streamflow, Kriging-

488   streamflow) produce a somewhat biased component, which increases when a shorter record is

489   employed in the model. For instance, the timing component for $\mathcal{M}_{\text{IDW-streamflow}}$ is 4.11 and

490   3.75 for the "deficit record" and "sufficient record", respectively. Moreover, timing

491   components dominate the error metric for all cases. However, the importance of variance

492   component is increased, especially in three models (FDC-IDW, DAR-streamflow, Kriging-

493   streamflow). Lastly, the results inform that if the proposed vine copulas approach is adapted,

494   variance and timing components are better captured, leading to better streamflow estimations,

495   which is beneficial in the practical applications of water resources management.

496

497   Finally, two predictions are further produced using two additional experiments: (1) the

498   observed marginal cumulative probabilities (i.e., using all 15 years) and conditional streamflow

499   values constructed from the partial record (i.e., based on $\varphi$ years), and (2) the estimated

500   marginal cumulative probabilities (i.e., based on $\varphi$ years) and conditional streamflow values

501    constructed from the full record (i.e., all 15 years). Their prediction abilities are evaluated over

502    the validated four-year period randomly selected in the remaining data. Similar to the previous

503    analysis, each analysis is tested 20 times. The results from these experiments provide an

504    inference to better isolate how error components from the two-step procedure (see section 2.2)

505    influence prediction skill.

506

507    Figure 10 shows the ACE scores from the out-of-sampled predictions using the proposed Dvine

508    model under the two scenarios. When considering all the quantiles together, the ACE scores

509    for the two scenarios are 0.003 (scenario #1) and 0.006 (scenario #2) on average under the

510    "deficit record" prediction. Also, the scores under the "sufficient record" prediction are all

511    nearly 0.003. Those results of the scores are sufficiently closed to zero, implying that both

512    predictions are reliable. Yet, compared to the predictions estimated by the cumulative

513    probabilities estimated by the partial record, and conditional models constructed by full records

514    (i.e., scenario #2), the ACE scores are achieved better, if the cumulative probabilities are

515    determined by the full record, except for some of the low and high quantiles. Similar

516    interpretation can be found in the NSE performance of two scenarios (see insets of Figure 10).

517    It may suggest that the first procedure (i.e., how to determine the cumulative probabilities for

518    the target site and its donors) is needed to pay careful attention when $\mathcal{M}_{\text{Dvine}}$ is utilized.

519    Nevertheless, the procedure to construct the conditional model in a streamflow gage network

520    is obviosuly crucial since the over or under-estimations are observed in many quantiles when

521    the insufficient sampling is employed in this process.

522

523    **5. Conclusion**

524   This study introduces a multiple dependence conditional model (i.e., vine copulas) to produce

525   streamflow estimates at partially gaged sites. The model includes a flexible high dimensional

526   joint dependence structure and conditional bivariate copula simulations. In order to confirm the

527   usefulness of a multiple dependence structure and the procedure for an appropriate number of

528   donor sites in the final vine copula model, the bivariate copula model and two types of vine

529   copulas with their unique procedure to determine the optimal number of donor sites are first

530   investigated using the generated data. These analyses were further extended in a case study of

531   the Yadkin-Pee Dee River Basin, the eastern United States by estimating streamflow in partially

532   gaged locations. In this analysis, six statistical infilling approaches were also employed to

533   represent applicability of the proposed model.

534

535   Results of the synthetic experiment and application to the Yadkin-Pee Dee River Basin

536   demonstrate that the propose model has benefits in some aspects. First, a multiple dependence

537   structure adopted in the proposed model is beneficial. From the massive evaluation experiments,

538   this study shows that multiple dependence structure clearly outperforms a single dependence

539   structure although there is the risk of overfitting when too many dependence structures are

540   employed. Moreover, this study confirms that the proposed multiple dependence structure

541   model with their optimum number of donor sites produces more reliable streamflow estimation

542   than other common infilling models. Next, the proposed model allows the development of

543   confidence intervals to consider prediction uncertainty, which is fairly attractive compared to

544   other models. For example, Bárdossy and Pegram (2013) argue that confidence intervals

545   obtained using an ordinary kriging model do not reflect the prediction uncertainty well

546   particularly on a daily scale. Overall, this study exhibits that a vine copula is potentially an

547   effective tool to support water resource management planners for objectives like gap-filling or

548    extending missing streamflow records.

549

550    While the results of the proposed model are favorable, there are possible limitations worthy of

551    further discussion. First, the assessment illustrated in this study focuses on model performance

552    under cross-validation at partially gaged basins, but additional work is needed to extend the

553    proposed model to ungagged basins. One possible way is to build a regression based model

554    with spatial proximity and physical basin characteristics to define associations between the

555    target and donor sites (e.g., Ahn and Steinschneider, 2019). Second, this study does not

556    consider potential nonstationarity in FDCs and correlations caused by the influence of

557    anthropogenic activity and change in land use. Nonstationarity may not be problematic in this

558    analysis since the assessment is limited to 15 years across the gaging network. However, if

559    longer records were used, it would be beneficial to consider the potential nonstationarity. This

560    exploration is left for future work.

561

562    There are several opportunities to improve the model structure. For instance, a vine copula is

563    able to incorporate more additional conditioning variables. One feasible approach is to add a

564    time series of climate data (e.g., precipitation) or to decompose a time series of streamflow

565    from the donor sites into a number of periodic components at different frequency levels through

566    the wavelet decomposition approach (Kisi and Cimen, 2011).

567

568    Lastly, the results presented here are specific to a study basin used in a case study. The proposed

569    model has not restricted to other watersheds around the world and its application is further

570   required towards drawing more generalized conclusions. In addition, the model could be used

571   for the purpose of infilling missing values of other hydrometeorological variables besides

572   streamflow (e.g., precipitation and soil moisture). For this application, the implementation of a

573   vine copula with combined discrete and continuous margins (i.e., to account for no rainfall

574   days) should be explored (e.g., Stoeber et al., 2013).

575

581

582                  **REFERENCES**

583

584   Aas, K., Czado, C., Frigessi, A. and Bakken, H.: Pair-copula constructions of multiple
585   dependence, Insur. Math. Econ., 44(2), 182–198, 2009.

586   Ahn, K.-H. and Palmer, R.: Regional flood frequency analysis using spatial proximity and basin
587   characteristics: Quantile regression vs. parameter regression technique, J. Hydrol., 540, 515–
588   526, doi:http://dx.doi.org/10.1016/j.jhydrol.2016.06.047, 2016a.

589   Ahn, K.-H. and Palmer, R. N.: Use of a nonstationary copula to predict future bivariate low
590   flow frequency in the Connecticut river basin, Hydrol. Process., 30(19), 3518–3532,
591   doi:10.1002/hyp.10876, 2016b.

592   Ahn, K.-H. and Steinschneider, S.: Hierarchical Bayesian Model for Streamflow Estimation at
593   Ungauged Sites via Spatial Scaling in the Great Lakes Basin, J. Water Resour. Plan. Manag.,
594   145(8), 04019030, 2019.

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

597   Archfield, S. A. and Vogel, R. M.: Map correlation method: Selection of a reference streamgage

598   to estimate daily streamflow at ungaged catchments, Water Resour. Res., 46(10), 2010.

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

601   Arreola Hernandez, J., Hammoudeh, S., Nguyen, D. K., Al Janabi, M. A. and Reboredo, J. C.:
602   Global financial crisis and dependence risk analysis of sector portfolios: a vine copula approach,
603   Appl. Econ., 49(25), 2409–2427, 2017.

604   Bárdossy, A. and Pegram, G.: Interpolation of precipitation under topographic influence at
605   different time scales, Water Resour. Res., 49(8), 4545–4565, 2013.

606   Bárdossy, A. and Pegram, G.: Infilling missing precipitation records–A comparison of a new
607   copula-based method with other techniques, J. Hydrol., 519, 1162–1170, 2014.

608   Beauchamp, J., Downing, D. and Railsback, S.: Comparison of regression and time-series
609   methods for synthesizing missing streamflow records, JAWRA J. Am. Water Resour. Assoc.,
610   25(5), 961–975, 1989.

611   Bedford, T. and Cooke, R. M.: Probability density decomposition for conditionally dependent
612   random variables modeled by vines, Ann. Math. Artif. Intell., 32(1–4), 245–268, 2001.

613   Bedford, T., Cooke, R. M. and others: Vines–a new graphical model for dependent random
614   variables, Ann. Stat., 30(4), 1031–1068, 2002.

615   Beguería, S., Tomas-Burguera, M., Serrano-Notivoli, R., Peña-Angulo, D., Vicente-Serrano, S.
616   M. and González-Hidalgo, J.-C.: Gap filling of monthly temperature data and its effect on
617   climatic variability and trends, J. Clim., 32(22), 7797–7821, 2019.

618   Bevacqua, E.: CDVineCopulaConditional: Sampling from Conditional C-and D-Vine Copulas,
619   R package version 0.1. 0., 2017.

620   Bhatti, M. I. and Do, H. Q.: Recent development in copula and its applications to the energy,
621   forestry and environmental sciences, Int. J. Hydrog. Energy, 44(36), 19453–19473, 2019.

622   Blum, A. G., Archfield, S. A. and Vogel, R. M.: On the probability distribution of daily
623   streamflow in the United States, Hydrol. Earth Syst. Sci., 21(6), 3093–3103, 2017.

624   Booker, D. and Snelder, T.: Comparing methods for estimating flow duration curves at
625   ungauged sites, J. Hydrol., 434, 78–94, 2012.

626   Boscarello, L., Ravazzani, G., Cislaghi, A. and Mancini, M.: Regionalization of flow-duration
627   curves through catchment classification with streamflow signatures and physiographic–climate
628   indices, J. Hydrol. Eng., 21(3), 05015027, 2016.

629   Boyle, D. P., Gupta, H. V. and Sorooshian, S.: Toward improved calibration of hydrologic
630   models: Combining the strengths of manual and automatic methods, Water Resour. Res., 36(12),
631   3663–3674, 2000.

632   Brechmann, E. C., Hendrich, K. and Czado, C.: Conditional copula simulation for systemic

633  risk stress testing, Insur. Math. Econ., 53(3), 722–732, 2013.

634  Castellarin, A., Galeati, G., Brandimarte, L., Montanari, A. and Brath, A.: Regional flow-
635  duration curves: reliability for ungauged basins, Adv. Water Resour., 27(10), 953–965, 2004.

636  Chen, L., Singh, V. P., Guo, S., Zhou, J. and Zhang, J.: Copula-based method for multisite
637  monthly and daily streamflow simulation, J. Hydrol., 528, 369–384, 2015.

638  Croley, T. and Hartmann, H.: NOAA Technical Memorandum ERL GLERL-61: Near-Real-
639  Time Forecasting of Large-Lake Water Supplies: A User's Manual, Ann Arbor MI, 1986.

640  Cunderlik, J. M. and Ouarda, T. B.: Regional flood-duration–frequency modeling in the
641  changing environment, J. Hydrol., 318(1), 276–291, 2006.

642  Czado, C.: Pair-copula constructions of multivariate copulas, in Copula theory and its
643  applications, pp. 93–109, Springer., 2010.

644  Czado, C.: Analyzing Dependent Data with Vine Copulas, Lect. Notes Stat. Springer, 2019.

645  Daneshkhah, A., Remesan, R., Chatrabgoun, O. and Holman, I. P.: Probabilistic modeling of
646  flood characterizations with parametric and minimum information pair-copula model, J.
647  Hydrol., 540, 469–487, 2016.

648  Dissmann, J., Brechmann, E. C., Czado, C. and Kurowicka, D.: Selecting and estimating
649  regular vine copulae and application to financial returns, Comput. Stat. Data Anal., 59, 52–69,
650  2013.

651  Erhardt, T. M., Czado, C. and Schepsmeier, U.: R-vine models for spatial time series with an
652  application to daily mean temperature, Biometrics, 71(2), 323–332, 2015.

653  Farmer, W.: Estimating records of daily streamflow at ungaged locations in the southeast
654  United States, PhD Disertation Tufts Univ. MA USA, 2015.

655  Farmer, W. H. and Vogel, R. M.: On the deterministic and stochastic use of hydrologic models,
656  Water Resour. Res., 52(7), 5619–5633, 2016.

657  Fisk, J.: Reproductive Ecology and Habitat Use of the Robust Redhorse in the Pee Dee River,
658  North Carolina and South Carolina., 2010.

659  Fu, G. and Butler, D.: Copula-based frequency analysis of overflow and flooding in urban
660  drainage systems, J. Hydrol., 510, 49–58, 2014.

661  Geenens, G.: Probit transformation for kernel density estimation on the unit interval, J. Am.
662  Stat. Assoc., 109(505), 346–358, 2014.

663  Gupta, H. V., Kling, H., Yilmaz, K. K. and Martinez, G. F.: Decomposition of the mean squared
664  error and NSE performance criteria: Implications for improving hydrological modelling, J.
665  Hydrol., 377(1), 80–91, 2009.

666  Hao, Z. and Singh, V. P.: Modeling multisite streamflow dependence with maximum entropy

667    copula, Water Resour. Res., 49(10), 7139–7143, 2013.

668    He, Y., Liu, R., Li, H., Wang, S. and Lu, X.: Short-term power load probability density
669    forecasting method using kernel-based support vector quantile regression and Copula theory,
670    Appl. Energy, 185, 254–266, 2017.

671    Hughes, D. and Smakhtin, V.: Daily flow time series patching or extension: a spatial
672    interpolation approach based on flow duration curves, Hydrol. Sci. J., 41(6), 851–871, 1996.

673    Joe, H.: Multivariate Models and Multivariate Dependence Concepts, Taylor & Francis. [online]
674    Available from: http://books.google.com/books?id=iJbRZL2QzMAC, 1997.

675    Kalteh, A. M. and Hjorth, P.: Imputation of missing values in a precipitation–runoff process
676    database, Hydrol. Res., 40(4), 420–432, 2009.

677    Karmakar, S. and Simonovic, S.: Bivariate flood frequency analysis. Part 2: a copula-based
678    approach with mixed marginal distributions, J. Flood Risk Manag., 2(1), 32–44, 2009.

679    Kisi, O. and Cimen, M.: A wavelet-support vector machine conjunction model for monthly
680    streamflow forecasting, J. Hydrol., 399(1–2), 132–140, 2011.

681    Kraus, D. and Czado, C.: D-vine copula based quantile regression, Comput. Stat. Data Anal.,
682    110, 1–18, 2017.

683    Li, M., Shao, Q., Zhang, L. and Chiew, F. H.: A new regionalization approach and its
684    application to predict flow duration curve in ungauged basins, J. Hydrol., 389(1), 137–145,
685    2010.

686    Liu, Z., Zhou, P., Chen, X. and Guan, Y.: A multivariate conditional model for streamflow
687    prediction and spatial precipitation refinement, J. Geophys. Res. Atmospheres, 120(19), 10–
688    116, 2015.

689    Lu, W.: A high-dimensional vine copula approach to comovement of China's financial markets,
690    in 2013 International Conference on Management Science and Engineering 20th Annual
691    Conference Proceedings, pp. 1538–1543, IEEE., 2013.

692    Mendicino, G. and Senatore, A.: Evaluation of parametric and statistical approaches for the
693    regionalization of flow duration curves in intermittent regimes, J. Hydrol., 480, 19–32, 2013.

694    Niemierko, R., Töppel, J. and Tränkler, T.: A D-vine copula quantile regression approach for
695    the prediction of residential heating energy consumption based on historical data, Appl. Energy,
696    233, 691–708, 2019.

697    Pugliese, A., Castellarin, A. and Brath, A.: Geostatistical prediction of flow–duration curves in
698    an index-flow framework, Hydrol. Earth Syst. Sci., 18(9), 3801–3816, 2014.

699    Salvadori, G. and De Michele, C.: Frequency analysis via copulas: Theoretical aspects and
700    applications to hydrological events, Water Resour. Res., 40(12), 2004.

701    Schepsmeier, U., Stoeber, J., Brechmann, E. C., Graeler, B., Nagler, T., Erhardt, T., Almeida,

702  C., Min, A., Czado, C., Hofmann, M. and others: Package 'VineCopula,' R Package Version,
703  2(5), 2015.

704  Schmid, F. and Schmidt, R.: Multivariate conditional versions of Spearman's rho and related
705  measures of tail dependence, J. Multivar. Anal., 98(6), 1123–1140, 2007.

706  Schnier, S. and Cai, X.: Prediction of regional streamflow frequency using model tree
707  ensembles, J. Hydrol., 517, 298–309, 2014.

708  Sklar, A.: Fonctions de Répartition À N Dimensions Et Leurs Marges, Université Paris 8.
709  [online] Available from: http://books.google.com/books?id=nreSmAEACAAJ, 1959.

710  Smakhtin, V. Y.: Generation of natural daily flow time-series in regulated rivers using a non-
711  linear spatial interpolation technique, Regul. Rivers Res. Manag. Int. J. Devoted River Res.
712  Manag., 15(4), 311–323, 1999.

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

715  U.S. Geological Survey: National Water Information System (NWISWeb): U.S. Geological
716  Survey database, [online] Available from: http://waterservices.usgs.gov/ (Accessed 1 January
717  2018), 2018.

718  Vernieuwe, H., Vandenberghe, S., De Baets, B. and Verhoest, N.: A continuous rainfall model
719  based on vine copulas, Hydrol. Earth Syst. Sci., 19(6), 2685–2699, 2015.

720  Worland, S. C., Steinschneider, S., Farmer, W., Asquith, W. and Knight, R.: Copula theory as a
721  generalized framework for flow-duration curve based streamflow estimates in ungaged and
722  partially gaged catchments, Water Resour. Res., 55(11), 9378–9397, 2019.

723  Xu, D., Wei, Q., Elsayed, E. A., Chen, Y. and Kang, R.: Multivariate degradation modeling of
724  smart electricity meter with multiple performance characteristics via vine copulas, Qual. Reliab.
725  Eng. Int., 33(4), 803–821, 2017.

726  Xu, Q. and Childs, T.: Evaluating forecast performances of the quantile autoregression models
727  in the present global crisis in international equity markets, Appl. Financ. Econ., 23(2), 105–
728  117, 2013.

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

731  Zimmer, D. M.: Analyzing comovements in housing prices using vine copulas, Econ. Inq.,
732  53(2), 1156–1169, 2015.

733

734

735

736

737

738     **List of Figures**

777
778
779
780
781
782

783
784 **List of Tables**

796
797
798
799
800
801
802
803
804
805
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
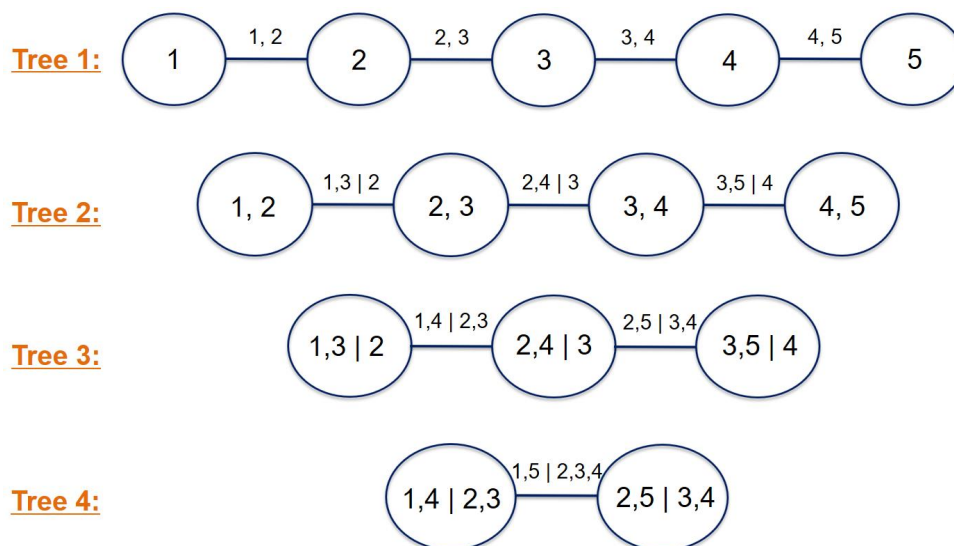
**Tree 1:** 1 — 1, 2 — 2 — 2, 3 — 3 — 3, 4 — 4 — 4, 5 — 5

**Tree 2:** 1, 2 — 1,3 | 2 — 2, 3 — 2,4 | 3 — 3, 4 — 3,5 | 4 — 4, 5

**Tree 3:** 1,3 | 2 — 1,4 | 2,3 — 2,4 | 3 — 2,5 | 3,4 — 3,5 | 4

**Tree 4:** 1,4 | 2,3 — 1,5 | 2,3,4 — 2,5 | 3,4

831
832    Figure 1 Example of D-vine structures with 5 variables, 4 trees and 10 edges
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860

861



862
863 Figure 2 Structure of the 6-Dimensional vine model and marginal for the synthetic simulation.
864 $LN(\pi, \sigma^2)$ denotes the log normal distribution with its mean ($\pi$) and variance ($\sigma^2$). The target
865 gage is highlighted.
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883

884



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

885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902

903



904
905 Figure 4 Pairwise upper and lower tail dependence for watersheds in the Yadkin-Pee Dee River
906 Basin. The upper triangular matrix shows values for the upper-tail dependence and the lower
907 triangular matrix presents values for the lower-tail dependence.
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926

927



928
929 Figure 5 Model performance for the Yadkin-Pee Dee river under a cross-validation framework
930 based on RMSE (dark squares) and NSE (light squares).
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
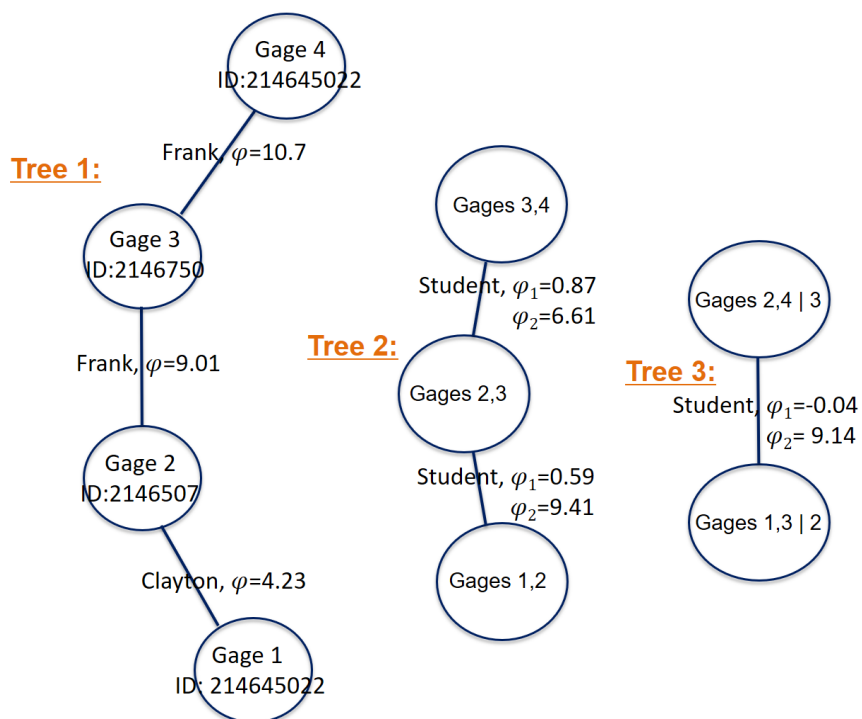946
947
948
949
950
951
952
953
954
955
956
957
958
959

960



961
962 Figure 6 (a) Average coverage error from three copula-based models for the Yadkin-Pee Dee
963 River Basin across exemplarily quantiles, and (b) 95% PIs for three models for an example
964 period (1 May 2018 to 31 July 2018) for a specific target gauge (USGS ID: 02143500).
965 Observed streamflow (black solid line) is also presented in each figure.
966
967
968
969
970
971
972
973
974
975
976
977
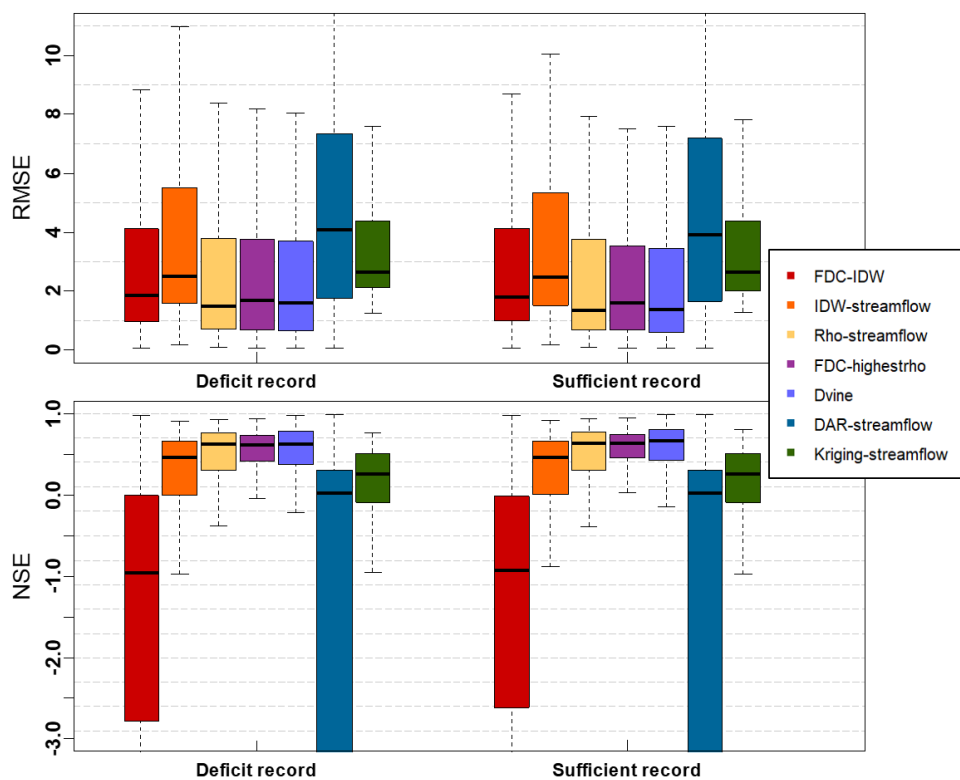978
979
980
981
982
983
984
985
986
987

988



989
990 Figure 7 Structure of the Dvine copula applied for a particular target site (USGS site number
991 214645022) with the defined bivariate copulas and their parameters.
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012

1013



1014
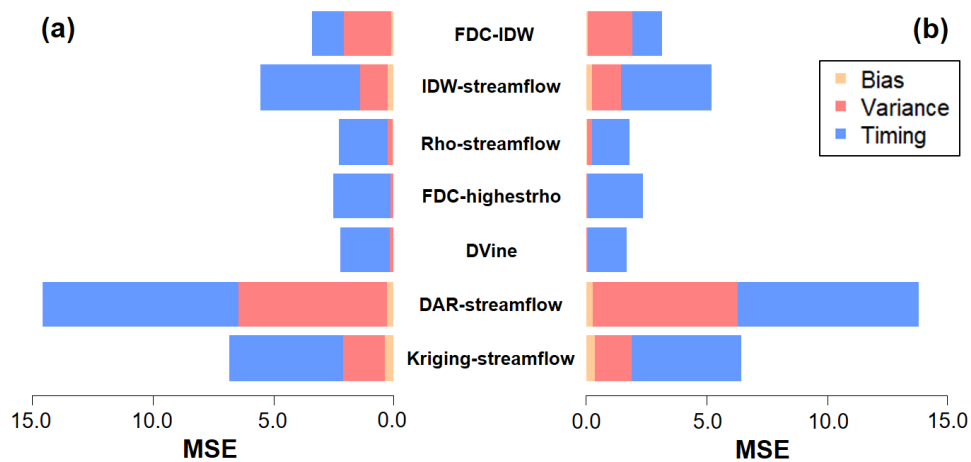1015 Figure 8 Inter-model comparison using cross-validation experiments based on RMSE (upper)
1016 and NSE (lower).
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035

1036



Figure 9 Three contributions from the decomposed mean squared error (MSE) for the cross-validation experiment with (a) the deficit record and (b) sufficient record scenarios.

1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
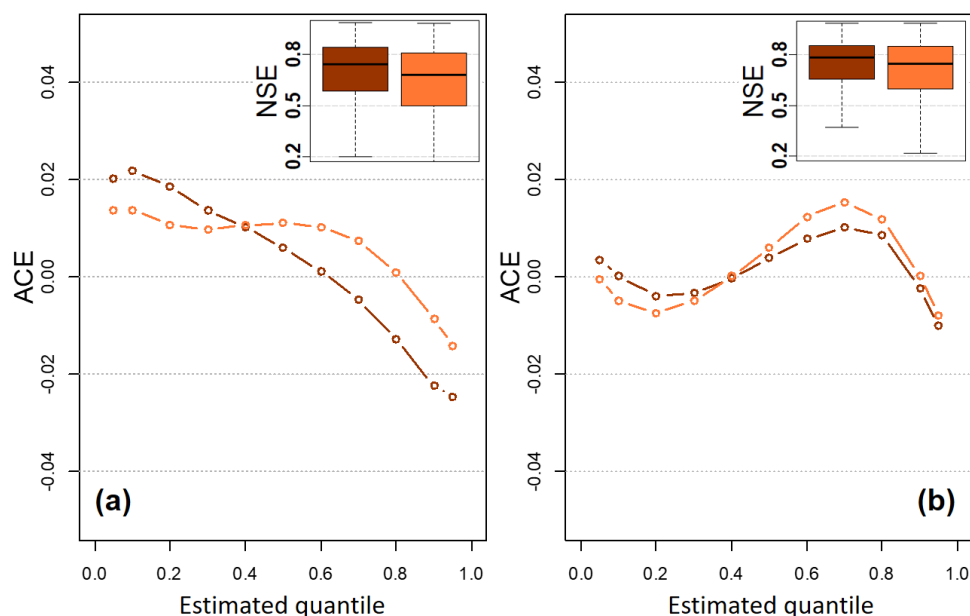1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068

1069



1070
1071 Figure 10 Average coverage error of the Dvine model for two scenarios under (a) the "deficit"
1072 and (b) "sufficient" cases. In each case, the dark line represents the scenario by the marginal
1073 cumulative probabilities using all years and conditional streamflow values constructed from
1074 the partial record. On the other hand, the light line illustrates the scenario by the marginal
1075 cumulative probabilities estimated by the partial record and conditional streamflow values
1076 constructed from the full record. Inset: NSE performance of the Dvine model for the two
1077 scenarios in each case.
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096

1097
1098    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 |

1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121

1122
1123 Table 2 RMSE and NSE results over the validation periods under synthetic experiment for
1124     comparing copula-based model formulations. Best metric values for each quantile are
1125     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 |
| | $\mathcal{M}_{\text{Dvine}}$ | ***0.895*** | ***1.046*** | ***1.112*** | ***1.391*** | 4.119 |
| Nash-Sutcliffe efficiency (NSE) | $\mathcal{M}_{\text{Bicop}}$ | ***0.464*** | 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 | ***0.805*** | ***0.838*** | ***0.869*** | ***0.905*** |

1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150

1151

1152 Table 3 Results of average coverage error (ACE) over the validation periods under synthetic
1153         experiment for comparing copula-based model formulations. Best metric values for
1154         each quantile are italicized and bolded.

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

1155
1156
1157
1158
1159
1160
1161
1162

1163