Copula autoregressive methodology for multi-lag, multi-site simulation of rainfall

Andres Felipe Ramirez¹ and Carlos Felipe Valencia¹

¹University of los Andes

November 26, 2022

Abstract

This work presents a methodology for the synthetic generation of rainfall time series based on the copula autoregressive methodology with multiple lags and for multiple sites. In this model, the multivariate time series is decomposed using pairwise copula functions to represent the whole cross-dependence, spatial and temporal structure of the data. We explore the advantages of using this nonlinear method over more traditional approaches that as an intermediate step transform the data to a normal distribution or usually omit the zero mass characteristics of the data. The use of copulas gives flexibility to represent the serial variability of the observed data on the simulation and allows for more control of the desired properties. We use discrete zero mass density distributions to assess the nature of rainfall, alongside a vector generalized linear model for the evaluation of time series distributions and their time dependence in multiple locations. We found that the copula autoregressive methodology models in a satisfactory manner the characteristics of the data, including its zero mass characteristics. These results will help to better understand the fluctuating nature of rainfall and also help to understand the underlying stochastic process.

Copula autoregressive methodology for multi-lag, multi-site simulation of rainfall

1

2

12

3	Andrés Felipe Ramírez ¹ , Carlos Felipe Valencia ¹
4	$^{1}\mathrm{University}$ of los Andes, Cra 1 E 19A-40, Bogota, Colombia. Center for optimization and applied
5	probability (COPA)
6	Key Points:
7	• We applied a multi-site and multi-lag methodology for simulation of rainfall time
8	series based on copula functions.
9	• We use bivariate Copula functions to model the discrete-continuous distribution,
10	allowing to model the occurrence and amount of rainfall.
11	• Marginal distributions of rainfall are estimated using a VGLM considering discrete-

continuous distributions including temporal covariables.

Corresponding author: Andrés Felipe Ramírez, af.ramirez12@uniandes.edu.co

Corresponding author: Carlos Felipe Valencia, cf.valencia@uniandes.edu.co

13 Abstract

This work presents a methodology for the synthetic generation of rainfall time series based 14 on the copula autoregressive methodology with multiple lags and for multiple sites. In 15 this model, the multivariate time series is decomposed using pairwise copula functions 16 to represent the whole cross-dependence, spatial and temporal structure of the data. We 17 explore the advantages of using this nonlinear method over more traditional approaches 18 that as an intermediate step transform the data to a normal distribution or usually omit 19 the zero mass characteristics of the data. The use of copulas gives flexibility to repre-20 sent the serial variability of the observed data on the simulation and allows for more con-21 trol of the desired properties. We use discrete zero mass density distributions to assess 22 the nature of rainfall, alongside a vector generalized linear model for the evaluation of 23 time series distributions and their time dependence in multiple locations. We found that 24 the copula autoregressive methodology models in a satisfactory manner the character-25 istics of the data, including its zero mass characteristics. These results will help to bet-26 ter understand the fluctuating nature of rainfall and also help to understand the under-27 lying stochastic process. 28

Keywords— Rainfall simulation, Copula autoregressive model, Multi-site rain fall simulation, multi-lag rainfall simulation, Vector Generalized Linear Model

31 1 Introduction

Stochastic generation of synthetic rainfall time series is a fundamental tool in areas such 32 as hydrology, agriculture, engineering, climate and energy, where it can be used for the 33 analysis of complex systems using numerical implementations (D. S. Wilks & Wilby (1999); 34 Srikanthan & McMahon (2001)). A correct representation of the stochastic nature of rainfall 35 includes its spatial and temporal properties. These characteristics allow for the study of 36 several phenomena without the necessity of real intervention; moreover, they support the 37 optimal design of larger systems that otherwise would require difficult physical or math-38 ematical models. Furthermore, synthetic generation in multiple locations is applied more 39 widely in applications than the single-location counterpart despite the requirement for the 40 representation of the spatial dependence between the locations. For example, in energy pro-41 duction and the integration of renewable sources, the correct understanding and modeling 42 of the stochastic process for rainfall is crucial for the evaluation of complementary viability 43 and reliability in power systems (Tinaikar (2013)). 44

The problem of synthetically generating daily rainfall time series has been approached 45 with different methodologies during the last 40 years. A detailed description and classifi-46 cation of the models used can be consulted in D. S. Wilks & Wilby (1999), Srikanthan & 47 McMahon (2001), Mehrotra et al. (2006) or Vaittinada Ayar et al. (2020). One taxonomy 48 is to separate the models between the single-site models that are concerned with the correct 49 specification of the statistical properties of observed time series to reproduce them with 50 simulated data and the models that extend these properties to the multi-site problem on 51 which the spatial relations have to be incorporated. 52

The difficulty of the rainfall stochastic process is that it represents two underlying 53 events, the occurrence and the amount. When seeing rainfall as a single process, the 54 marginal distribution at each time should have a mass on zero, that is, the probability 55 of the amount being equal to zero is greater than 0. Some methodologies approach this 56 problem using a two-step mechanism that first simulates the occurrence (for example using 57 a first-order discrete Markov Chain) and conditionally on this they simulate the amount of 58 rainfall (e.g. Richardson (1981), D. Wilks (1998), Mhanna & Bauwens (2012),). As a result 59 of this method, the random variables that model the amount of rainfall are conditionally 60 independent from the ones that model the occurrence. For example, if there is rain on two 61 consecutive days, the amount of rain on the second day is independent from the amount on 62

the first day. This could undermine the representation of some statistical properties such 63 as the autocorrelation function when looking at the series conditionally on the occurrence. 64 One alternative is to discretize the stochastic process and then use a first or second order 65 Markov Chain (e.g. Haan et al. (1976)); however, the main repercussion is the less accurate 66 representation of the amount of rainfall in the marginal densities and the addition in compu-67 tational effort needed when using a large number of states. Another branch of models uses 68 resampling techniques such as bootstrapping in order to approximate the distribution of 69 observed series (e.g. Buishand & Brandsma (2001)). Other strategies are known to assume 70 that the non-zero part of the dataset might have normal marginal densities (Zhou et al. 71 (2020), Ahn (2020), Ayar et al. (2020)). Once the new data is normal, the derived method-72 ologies take advantage of the well known linear ARIMA (autoregressive integrated moving 73 average) stochastic process to construct theoretically robust statistical models. With nor-74 mal time series, it is possible to incorporate temporal and spatial dependence in the first 75 two moments of the distributions; however, this strategy depends on the marginal densi-76 ties and the threshold used for censored values and positive rainfall measurements. This 77 approach might poorly estimate several characteristics of the actual process (Zhou et al. 78 (2019)). Furthermore, rainfall data is unlikely to correlate symmetrically for low and high 79 values (Bárdossy et al. (2017)) and the required transformation from the marginals to the 80 normal distribution (gaussianization) conditions the behavior of the series (Sarmiento et al. 81 (2018)).82

The use of copula functions to model the time or spatial dependence of rainfall time 83 series has been explored more during the last years, where several applications in hydrology 84 use copula functions to assess the complexity of the data (e.g. Kao & Govindaraju (2008), 85 Bárdossy & Pegram (2009), Zhang & Singh (2019), Serinaldi (2009a), Serinaldi (2009b)). 86 The benefits of using the copula function is that they can estimate the underlying stochastic 87 process and structure between several rainfall time series independently of their marginal 88 distribution functions (Vandenberghe et al. (2010), Balistrocchi & Bacchi (2011)). This 89 provides enough flexibility on the relations among variables to overcome some limitations 90 of linear Gaussian processes that rely on specific transformations (Bárdossy et al. (2017), 91 Sarmiento et al. (2018)). Several studies, such as the one conducted by Zhang & Singh 92 (2007) suggest that the copula based distribution function fits the dependence structure of 93 observed rainfall characteristics data series better than the multivariate normal probability 94 distribution. 95

Copula functions are naturally defined for continuous variables to accommodate the 96 estimation for zero inflated distributions. Authors like Serinaldi (2009a), Serinaldi (2009b) 97 and Li et al. (2013) use a partition in quadrants of the bivariate uniform distributions to 98 generate conditional distributions. These partitions are adaptation of the ideas of Shimizu 99 (1993) and Herr & Krzysztofowicz (2005) applied to the copula model for time series. Al-100 though these copula models are a great tool to model the series dependence structure up 101 to lag 1 in an autoregressive process, the spatial dependence of multi-site time series is not 102 usually modeled directly from the between series and conditional series dependence. These 103 dependence structures in daily rainfall time series make it difficult to model over the range 104 of observations using only linear correlation coefficients, such as Pearson correlation. In-105 stead, dependence methodologies such as Kendall's τ are often used to model the correlation 106 between locations (i.e Serinaldi (2009a), Serinaldi (2009b)). 107

In this paper, we propose a fully copula based methodology to model time series of 108 rainfall within an autoregressive process with multiple lags and with spatial dependence 109 based also on copula functions. That is, the whole stochastic system can be expressed as a 110 multivariate model decomposed in pairwise copula functions using the COPAR methodol-111 ogy (Brechmann et al. (2015)), which bases its calculations on generating all the pair copula 112 construction needed for the between series and conditional series dependence among the dif-113 ferent time series for each of the sites following an R-vine structure. In addition, we employ 114 a vector generalized linear model (VGLM) to fit a zero inflated continuous distribution in 115 terms of temporal variables such as the month of the year. With this method, it is possible 116 to estimate the marginal distributions in one model avoiding the possible over-fitting that 117 can be created when partitioning the sample and estimating the marginal distributions for 118 each month. 119

The structure of this paper goes as follows. The methodology used as well as the contributions made to the state of the art are described in Section 2. A case study is presented in Section 3 along with the evaluation of the results. Finally, section 4 presents a discussion about the methodology used and the future work, as well as the conclusions.

¹²⁴ 2 Methodology

The rainfall synthetic data generator is based on a model that has two main components:(i) marginal distribution estimation through a VGLM for mixtures of densities with discrete

mass distributions (Section 2.1), and (ii) a multivariate copula algorithm for modeling the temporal and spatial dependence in one statistical model using decomposition in bivariate uniform densities (Section 2.2).

The rainfall time series is not a stationary process given that, according to a specific lo-130 cation, there are seasons and inter-annual cycles that have to be considered. These seasonal 131 components of the series directly affect the distribution of the marginal random variable. 132 Contrary to Gaussian time series, on which the seasonal effect is modeled on the mean, the 133 copula model allows for a more flexible temporal effect in some (or all) the parameters of the 134 distributions. The reason for this is that the temporal and spatial dependencies are mod-135 eled through the uniform variables resulting from using the inverse cumulative probability 136 function transformation. Using a VGLM to estimate all rainfall marginal densities permits 137 the use of all available data to include the temporal effects as independent variables. Fur-138 thermore, the occurrence process (rain or no-rain) is embedded in the marginal distribution 139 by using a mixture model that combines the positive part with the probability mass on zero. 140 Once the multivariate rainfall time series is standardized to be uniform [0, 1] for all values, it 141 is modeled by the copula functions to express the whole joint probability distribution. This 142 distribution has many components (dimensions) that would make it difficult to estimate 143 without simplification rules to avoid the curse of dimensionality. We use a decomposition 144 in a R-vine structure and limit the numbers of lags that are significant, as done in a linear 145 autoregressive (AR) model. This approach extends the previous models in Serinaldi (2009a) 146 and Li et al. (2013) that only permit one lag. The same principle can be used to model 147 the multivariate distribution applied to the spatial dependence. The copula autoregressive 148 methodology models the entire dependence of the time series differing from more tradi-149 tional tools for modeling time series such as autoregressive processes (AR) that transform 150 any distribution to normal. 151

152

2.1 Rainfall marginal distribution

The marginal distributions of ground measured rainfall is modeled using a vector generalized linear model (VGLMs). Unlike the classical generalized linear models (GLMs), there is not a restriction on the number of parameters in the distributions, and they are purposely general to allow greater utility (Yee, 2015). Consider the data set expressed as (x_i, y_i) for i = 1, ..., n, where x_i is a vector of p explanatory variables and y_i is the response for the observation i. The objective is to fit a generalized regression model involving the parameter ν_j , where j = 1, ..., J. The VGLM model, estimates each one of the parameters in the generalized regression as the function of a linear combination of the explanatory variables:

$$g_j(\nu_j) = \beta_j^T x = \beta_{(j)1} X_1 + \beta_{(j)2} X_2 + \dots + \beta_{(j)p} X_p \quad , \tag{1}$$

153 154

155

where g_j is a parameter link function. In this study, we estimated the parameters of the VGLMs with the Interactively Re-weighted Least Squares (IRLS) algorithm. For a detailed explanation of the estimation of the VGLM, the reader is referred to Yee (2015)

For modeling the marginal probability distribution of rainfall, we selected the zero adjusted gamma distribution, due to the zero mass characteristics presented in the data. This is a special case of zero adjusted distributions on zero and the positive real line. The zero adjusted gamma distribution is a combination of a discrete value 0 with probability ν and a gamma GA(μ , σ) distribution with probability (1- ν), where μ is the mean and σ is the square root of the dispersion parameter according to the exponential family factorization. The probability function of the zero adjusted gamma distribution denoted by ZAGA($\mu(t), \sigma(t), \nu(t)$) is given by:

$$f_{y}(y \mid \mu(t), \sigma(t), \nu(t)) = \begin{cases} \nu(t) & y = 0\\ (1 - \nu(t))f_{w}(y \mid \mu(t), \sigma(t)) & y > 0 \end{cases}$$
(2)

For $0 \le y(t) < \infty$, where $\mu(t) > 0$, $\sigma(t) > 0$, $0 < \nu(t) < 1$ and $W \sim GA(\mu(t), \sigma(t))$ has a gamma distribution with density:

$$f_w(y \mid \mu(t), \sigma(t)) = \frac{y^{\frac{1}{\sigma^2} - 1} \exp\left(\frac{-y}{\sigma^2 \mu}\right)}{(\sigma^2 \mu) \Gamma(\frac{1}{\sigma^2})} \quad y > 0 \quad , \tag{3}$$

where the dependence of the parameters on t is left implicit for simplicity. The mean of the ZAGA distribution is $(1 - \nu(t))\mu$ and the variance is $(1 - \nu(t))\mu^2(\nu + \sigma^2)$. The independent variable t represents the yearly seasons that change the marginal distributions. In our case, we model t as the month of the year using a factor with values (1-12). The default link functions that relate the parameter (μ, σ, ν) depending on the time are:

$$\log(\mu) = \beta_{\mu}^{T} t \tag{4}$$

$$\log(\sigma) = \beta_{\sigma}^{T} t \tag{5}$$

$$\log(\nu) = \beta_{\nu}^{T} t \tag{6}$$

¹⁵⁶ 2.2 The copula autoregressive model

The COPAR model proposed by Brechmann et al. (2015) exploits the flexibility of vine 157 copulas for non-linear and asymmetric modeling of serial and between-series dependence. 158 The fundamental pieces to build the autoregressive model are the bivariate copulas, which 159 are distributions on the unit square $[0,1]^2$ such that both marginals are uniform U(0,1). 160 The theorem by Sklar (1959) explains that for any given variables X and Y with joint 161 distribution $F_{X,Y}(x,y)$ and marginals cumulative distribution functions (CDF) $F_X(x)$ and 162 $F_Y(y)$ respectively, there exists a unique copula function $C_{XY}(\cdot, \cdot)$ that connects $F_{X,Y}(\cdot, \cdot)$ 163 to $F_X(\cdot)$ and $F_Y(\cdot)$ via $F_{X,Y}(x,y) = C_{XY}(F_X(x),F_Y(y))$. Therefore, the information in the 164 joint distribution is decomposed into that of the marginal distributions and the one of the 165 copula function, where the latter captures the entire dependence structure between X and 166 Y (Chen & Fan, 2006). For a detailed explanation on copula theory the reader is referred 167 to Nelsen (1999). 168

169

2.2.1 Pair-copula decomposition for univariate time series

Let $\{X_t\}_{t=1,\dots,T}$ be a univariate stationary time series of continuously distributed data, with marginal distribution function F_{X_t} and density function f_{X_t} . The joint distribution of $\{X_t\}$ can be decomposed by selecting models for the conditionals as:

$$f(x_1, \cdots, x_T) = f(x_1) \prod_{t=2}^T f(x_t \mid x_{t-1}, \cdots, x_1).$$
(7)

Smith et al. (2010) outline that this expression can be used to obtain a general decomposition in terms of bivariate copulas, as for every s < t exists a copula density $c_{t,s}$ such that:

$$f(x_t, x_s \mid x_{t-1}, \cdots, x_{s+1}) = c_{t,s}(F(x_t \mid x_{t-1}, \cdots, x_{s+1}), F(x_s \mid x_{t-1}, \cdots, x_{s+1})) \cdot f(x_t \mid x_{t-1}, \cdots, x_{s+1}), f(x_s \mid x_{t-1}, \cdots, x_{s+1}),$$
(8)

where $F(X_t | X_{t-1}, \dots, X_{s+1})$ and $F(X_s | X_{t-1}, \dots, X_{s+1})$ are the conditional distributions functions of X_t and X_s respectively. This expression is the density of the Sklar theorem, conditional upon $\{X_{t-1}, \dots, X_{s+1}\}$. Rearranging terms in Equation (8) gives:

$$f(x_t \mid x_{t-1}, \cdots, x_s) = c_{t,s}(F(x_t \mid x_{t-1}, \cdots, x_{s+1}), F(x_s \mid x_{t-1}, \cdots, x_{s+1})) \cdot f(x_t \mid x_{t-1}, \cdots, x_{s+1}).$$
(9)

By recursive conditioning on $s = 1; 2; \dots, t - 1$ we obtain:

$$f(x_t \mid x_{t-1}, \cdots, x_1) = \prod_{j=1}^{t-2} \{ c_{t,j}(F(x_t \mid x_{t-1}, \cdots, x_{j+1}), F(x_j \mid x_{t-1}, \cdots, x_{j+1})) \}$$

$$c_{t,t-1}(F(x_t), F(x_{t-1})) f(x_t).$$
(10)

For notation simplicity, denote $u_{t|j} = F(X_t \mid X_{t-1}, \dots, X_j)$ and $u_{j|t} = F(X_j \mid X_t, \dots, X_{j+1})$ as the projections backwards and forwards t - j steps, respectively. Also, by denoting $u_{t|t} = F(X_t)$, the joint density function in equation (7) can be written as:

$$f(x_1, \cdots, x_T) = \prod_{t=2}^{T} \left[\prod_{j=1}^{t-1} \{ c_{t,j}(u_{t|j+1}, u_{j|t-1}) \} f(X_t) \right] f(X_1).$$
(11)

This decomposition of the joint density function of the time series allows the method-170 ology to achieve very flexible models, as no restrictions are necessary in the selection of the 171 copula families. Nonetheless, copulas corresponding to the same time lag must be identical 172 to assure that the simulated time series is stationary. This particular pair-copula decompo-173 sition is known as D-vine and takes part in a more general class of decomposition known as 174 R-vines, which is a graphic theoretic model to establish which pair-copulas are included in 175 the decomposition of a time series. An R-vine on d variables is a sequence of d-1 linked 176 trees (connected acyclic graphs) that satisfy several conditions. R-vine theory is described 177 in detail by Bedford (2001) and H Joe (2011). 178

179

2.2.2 The COPAR model for multivariate time series

Although the model works in general for any number of dimensions, we explain the copula autoregressive model for two univariate time series $\{X_t\}_{t=1,\dots,T}$ and $\{Y_t\}_{t=1,\dots,T}$ jointly distributed at time point $t = 1, \dots, T$ for brevity in the exposition. A flexible multivariate distribution of $\{X_t\}$ and $\{Y_t\}$ is presented, which allows the nonlinear dependence-serial as well as between-series dependence structure to be constructed. The model is based on a particular R-vine structure that has the following components (Brechmann et al., 2015).

1. Marginal distributions F_X and F_Y of $\{X_t\}$ and $\{Y_t\}$

- 2. An R-vine for the serial and between-series dependence of $\{X_t\}$ and $\{Y_t\}$, with the selection of the following pair-copulas:
 - (a) Serial dependence of $\{X_t\}$:

$$X_s, X_t \mid X_{s+1}, \cdots, X_{t-1} \quad \forall \quad 1 \le s \le t \le T$$

$$\tag{12}$$

(b) Between series dependence:

$$X_s, Y_t \mid X_{s+1}, \cdots, X_t \quad \forall \quad 1 \le s \le t \le T$$

$$\tag{13}$$

and

$$Y_s, X_t \mid X_1, \cdots, X_{t-1}; Y_{s+1}, \cdots, Y_{t-1} \quad \forall \ 1 \le s \le t \le T$$
(14)

(c) Conditional serial dependence of $\{Y_t\}$:

$$Y_s, Y_t \mid X_1, \cdots, X_t; Y_{s+1}, \cdots, Y_{t-1} \quad \forall \quad 1 \le s \le t \le T$$

$$\tag{15}$$

As mentioned before, pair-copulas with equal lag must be identical. It is important to notice that $\{X_t\}$ has a pivotal role in this modeling approach because the serial dependence of $\{X_t\}$ is modeled unconditionally, meanwhile that of $\{Y_t\}$ is specified conditionally on $\{X_t\}$. The selection of the copulas is performed sequentially since, for example, $c_{X_tX_{t-2}|X_{t-1}}$ depends on the copula $c_{X_tX_{t-1}}$. Finally, the order of the model (k) is selected if all paircopulas corresponding to a lag length greater than k are independence copulas.

195

2.2.3 Pair-copula estimation of variables with zero mass

One of the contributions of this paper is the inclusion of variables with zero mass in 196 the copula autoregressive modeling for time series shown in section 2.2.2. In order to model 197 these particular distributions, we selected a zero adjusted gamma distribution for the daily 198 rainfall data, which is a mixture distribution that combines a Bernoulli and a Gamma prob-199 ability distribution, dependent on the parameter $\nu(t)$ which models the probability of zero 200 or non-zero values. This distribution is explained in depth in section 2.1. The zero mass 201 characteristics of the data makes the density probability function $f_Y(y)$ to have a jump 202 discontinuity on zero, (Figure 1). For modeling the copula joint distribution, however, it is 203 necessary to evaluate this function and calculate the respective distribution function $F_Y(y)$, 204 since the results of this calculation should be a continuous uniform variable. Moreover, to 205 generate synthetic data, the inverse function F_Y^{-1} has to be applied on random numbers. 206 To address this issue, we adopt the approach described in Faugeras (2012), from which 207 the bivariate copula estimation of the autoregressive model explained in section 2.2 can be 208 assessed with some fundamental considerations. In particular, for the points that have a 209 mass value at $\nu(t)$, a uniform variable is generated when evaluating F_Y to fill the gaps 210 created by the jumps. In this case, it is clear that $F_Y(0) = \nu$, and therefore, for the purpose 211 of evaluation and simulation, it is replaced by a uniform random variable between 0 and 212

²¹³ ν . Figure 2 presents the evaluation of $F_Y(y)$ for all the points in Figure 1, where $\nu = 0.6$. ²¹⁴ Also, Figure 2 presents the respective histogram. Figure 3 show the counterparts after using ²¹⁵ the randomization procedure. The proposed strategy to use copula functions with mixed ²¹⁶ random variables with probability mass in zero is explained in algorithm 1.

217

for
$$t \leftarrow 1$$
 to T do

$$u_{t|t} = F_y = \begin{cases} \nu(t) & y = 0 \\ \nu(t) + (1 - \nu(t)) \frac{\gamma(\sigma(t)^{-2}, y\mu(t)^{-1}\sigma(t)^{-2})}{\Gamma(\sigma(t)^{-2})} & y > 0 \end{cases}$$
(16)

end

218

 $\begin{array}{ll} \mbox{if} \ u_{t|t} = \nu(t) \ \mbox{then} \\ | \ \ u_{t|t} = \textit{runif}(0,\nu(t)) \\ \mbox{end} \end{array}$

for $t \leftarrow 1$ to T do

end

Algorithm 1: Synthetic non zero uniform data generation

Where $F_y = (y \mid \mu(t), \sigma(t), \nu(t))$ is the distribution function of the zero adjusted gamma distribution for $0 \le y(t) < \infty$, with $\mu(t) > 0$, $\sigma(t) > 0$, $0 < \nu(t) < 1$. Also, runif is a random uniform variable between 0 and $\nu(t)$.



Figure 1. Scatter plot (LEFT) and histogram (RIGHT) of density function for distribution ZAGA(sigma=0.5, ν =0.6)



Figure 2. Scatter plot (LEFT) and histogram (RIGHT) of CDF function w/o Algorithm 1 for distribution ZAGA(sigma=0.5, ν =0.6)



Figure 3. Scatter plot (LEFT) and histogram (RIGHT) of CDF function w/ Algorithm 1 for distribution ZAGA(sigma=0.5, ν =0.6)

If X and Y are two jointly distributed variables with probability mass on zero, the probability distribution function (CDF) $F_{XY}(x, y)$ may be defined conditionally on four cases with respective probabilities: $P_{00} = P(X = 0, Y = 0), P_{10} = P(X > 0, Y = 0),$ $P_{01} = P(X = 0, Y > 0)$ and $P_1 = P(X > 0, Y > 0)$. Therefore, it is possible to write it as:

$$F_{XY}(x,y) = P_{00} + P_{10}H_X(x) + P_{01}H_Y(y) + P_{11}H_{XY}(x,y), \qquad (17)$$

where $H_{XY} = C_{XY}^{11}(F_X^1(x), F_Y^1(y))$ can be modeled as a bivariate continuous copula function 222 (conditional on both variables being positive), with $F_X^1(x) = P(X \le x | X > 0, Y > 0)$ and 223 $F_Y^1(y) = P(Y \le y | X > 0, Y > 0)$. According to this, the resulting bivariate copula function 224 to represent $F_{XY}(x, y)$ can be also decomposed in quadrants. That is, if X and Y are both 225 marginally distributed ZAGA with $p_x = \nu(X) = P(X = 0)$, and $p_y = P(Y = 0)$, then 226 $F_{X,Y} = C_{XY}(u,v)$, where $u = F_X(x)$, and $v = F_X(y)$. Let $c_{XY}(\cdot, \cdot)$ be the correspondent 227 bivariate density of $C_{XY}(\cdot, \cdot)$, then it would look as Figure 4. With this function, it is possible 228 to define the conditional cumulative distributions of Y given X, and their respective inverse, 229 using the splitting mechanism described in Herr & Krzysztofowicz (2005) and Serinaldi 230 (2009a). 231



Figure 4. Bivariate copula density for two variables with ZAGA distributions. In the left, the quadrants are represented. In the right, the empirical bivariate copula density for real rainfall data

Note that the same partition can be applied directly to the density $c_{XY}(u, v)$, where:

232

$$c_{XY}(u,v) = P_{00} + P_{10}c_u\left(\frac{u-p_x}{1-p_x}\right) + P_{01}c_v\left(\frac{v-p_y}{1-p_y}\right) + \quad (18)$$

$$P_{11}c_{XY}^{11}\left(\frac{1-p_x}{P_{11}}\left\lfloor\left(\frac{u-p_x}{1-p_x}\right)-C_u(u)\frac{P_{10}}{1-p_x}\right\rfloor,\frac{1-p_y}{P_{11}}\left\lfloor\left(\frac{v-p_y}{1-p_y}\right)-C_v(v)\frac{P_{01}}{1-p_y}\right\rfloor\right),$$
(19)

where $c_{XY}^{11}(u, v)$ is the corresponding density function of $C_{XY}^{11}(u, v)$, that is defined on the condition that both variables X and Y are greater than zero. C_u and C_v are the cumulative distribution functions of c_u and c_v respectively, where

$$C_u(u) = P\left(F_X(X) \le \frac{u - p_x}{1 - p_x} \Big| X > 0, Y = 0\right) = P\left(p_x Z + (1 - p_x) F_X^1\left[H_X^{-1}(Z)\right] \le \frac{u - p_x}{1 - p_x}\right),$$
(20)

where $Z \sim Uniform(0, 1)$. Analogously,

$$C_{v}(v) = P\left(F_{Y}(y) \le \frac{v - p_{y}}{1 - p_{y}} \Big| X = 0, Y > 0\right) = P\left(p_{y}Z + (1 - p_{y})F_{Y}^{1}\left[H_{Y}^{-1}(Z)\right] \le \frac{v - p_{y}}{1 - p_{y}}\right)$$
(21)

236

2.3 Discrete copula CDF and inverse CDF

The at-site rainfall daily generator is based on the simulation from mixed discrete-237 continuous conditional distribution, which is deduced from a copula-based discrete-continuous 238 conditional bivariate distribution. Moreover, this study is an extension of the research done 239 by Serinaldi (2009b), Serinaldi (2008), Serinaldi (2009a), Herr & Krzysztofowicz (2005). 240 The main contribution to these methodologies is that when we model the temporal and 241 spatial (multi-site) conditional series dependencies, we model the whole multivariate time-242 series with copula relations (COPAR methodology). The resulting methodology models the 243 dependence structure directly from the conditional dependence generated in the construc-244 tion of the R-vine pair-copula structure following the specified structures in Brechmann 245 et al. (2015). For an in depth description and mathematical background in the construc-246 tion of bivariate copula-based discrete-continuous conditional distribution and the copula 247 autoregressive methodology, the reader is referred to the previous authors mentioned. 248

The main building block for simulating time series based on a copula based autoregressive model is the conditional bivariate copula. For example, for a stationary univariate time series, let $C_{t,t-k}(v|v_k)$ be the bivariate copula of the variable v with its respective lag k, for instance, when k takes the value of 1, we define the copula of $v = F_X(X_t)$ with $v_1 = F_X(X_{t-1})$. To generate the value v conditioned on v_1 , it is necessary to define the function:

$$h_{t,t-1} = \frac{\delta C_{t,t-1}(v,v_1)}{\delta v_1} , \qquad (22)$$

and its respective inverse $h_{t,t-1}^{-1}$. Note that according to the time series decomposition presented in equation 11, when k > 1, $h_{t,t-k}$ correspond to the derivative of the conditional bivariate copula $C_{t,t-k}(v, v_k | v_{k+1}, ..., v_1)$. When the original variables have mass on zero, these derivatives have to be adapted in order to account for the partition in quadrants (Figure 4). This section illustrates the estimation of the conditional copula CDF $(h_{t,j})$ and its inverse $(h_{t,j}^{-1})$ for the mixed discrete-continuous copula, and then how to use this function in order to build the multivariate time series synthetic generation algorithm. The construction of the conditional CDF is presented in Algorithm 2, and conditional inverse CDF is explained in Algorithm 3, for $v = F_x(x_t|x_{j+1}, \dots, x_{t-1})$ and $v_j = F_x(x_j|x_{j+1}, \dots, x_{t-1})$.

else

$$\psi = \frac{P_{10}}{P_{10} + P_{11}} c_{v_j} \left(v_j \right)$$

258

$$\begin{split} \mathbf{if} \ v &\leq p_v \ \mathbf{then} \\ \left| \begin{array}{c} h_{t,j} &= \left(\frac{v_j}{p_{v_j}}\right) \psi \\ \mathbf{else} \\ \\ \tilde{v} &= \frac{1-p_v}{P_{11}} \left[\left(\frac{v-p_v}{1-p_v}\right) - C_v \left(v\right) \frac{P_{10}}{1-p_v} \right] \\ \tilde{v}_j &= \frac{1-p_{v_j}}{P_{11}} \left[\left(\frac{v_j - p_{v_j}}{1-p_{v_j}}\right) - C_{v_j} \left(v_j\right) \frac{P_{10}}{1-p_{v_j}} \right] \\ \\ h_{t,j} &= \psi + (1-\psi) \frac{\delta C_{Xt,X_j}^{11}(\tilde{v},\tilde{v}_j)}{\delta \tilde{v}_j} \\ \\ \mathbf{end} \end{split}$$

end

Algorithm 2: Continuous-Discrete bivariate copula-based CDF.

$$\begin{array}{l} \mathbf{if} \ v_j \leq p_{v_j} \ \mathbf{then} \\ \\ \left| \begin{array}{c} \mathbf{if} \ v \leq P_{00}/(P_{00}+P_{01}) \ \mathbf{then} \\ \\ | \ h_{t,j}^{-1} = \frac{P_{00}+P_{01}}{P_{00}} p_v v \\ \\ \mathbf{else} \\ \\ \left| \begin{array}{c} h_{t,j}^{-1} = C_v^{-1} \left(v - \frac{P_{00}}{\frac{P_{00}+P_{01}}{1-\frac{P_{00}}{P_{00}+P_{01}}} \right) \\ \\ \mathbf{end} \\ \end{array} \right) \\ \end{array}$$

259

$$\begin{split} \xi &= \frac{P_{10}}{P_{10} + P_{11}} c_{v_j}(v_j) \\ \text{if } v &\leq \xi \text{ then} \\ \mid h_{t,j}^{-1} = p_v v / \xi \\ \text{else} \\ \mid \tilde{v} &= G_v^{-1}(v) \\ \tilde{v}_j &= \frac{1 - p_{v_j}}{P_{11}} \left[\left(\frac{v_j - p_{v_j}}{1 - p_{v_j}} \right) - C_{v_j}(v_j) \frac{P_{10}}{1 - p_{v_j}} \right] \\ \mid h_{t,j}^{-1} &= \xi + (1 - \xi) \frac{\delta C_{X_t, X_j}^{11}(\tilde{v}, \tilde{v_j})}{\delta \tilde{v_j}}^{-1} \\ \text{end} \end{split}$$

end



Where $p_v = P_{00} + P_{01}$ and $p_{v_j} = P_{00} + P_{10}$, and:

$$G_{v}(a) = \frac{1 - p_{v}}{P_{11}} \left[\left(\frac{a - p_{v}}{1 - p_{v}} \right) - C_{v}(a) \frac{P_{10}}{1 - p_{v}} \right] .$$
(23)

The discrete copula distribution function and the inverse function are conditionally esti-260 mated from the probabilities of each of the quadrants and the respective empirical distri-261 butions. That is, C_u and C_v can be estimated directly from the uniform variables. For 262 the top right quadrant (where both variables are continuous), the estimation is done with 263 parametric copulas using the correspondent conditional values. Thus, $\frac{\delta C^{11}(u,v)}{\delta u}$ is the copula 264 function where both marginals take non-zero values. We limit the parametric families used 265 for the P_{11} quadrant to be: 0 = independence copula, 1 = Gaussian copula, 2 = Student 266 t copula (t-copula), 3 = Clayton copula, 4 = Gumbel copula, 5 = Frank copula and 6 =267 Joe copula. However, the analysis could be performed with any family, in particular copula 268 families with long tails such as the Plackett copula. 269

2.3.1 Simulation algorithm

270

This section presents the simulation algorithm based on the copula autoregressive model. First, we introduce the algorithm for univariate time series from Smith et al. (2010), after that, we explain the adjustment that was assessed for two jointly observed time series. This procedure can be extended to any number of variables and eany number of lags.

Going back to the D-vine decomposition shown in section 2.2.1, the critical aspect is the evaluation of $u_{t|j+1}$ and $u_{j|t-1}$ from Equation (11). In this regard, it is worth mentioning the following property H Joe (2011): let $u_1 = F(X_1)$ and $u_2 = F(X_2)$ be conditional distribution functions, and $F(X_1, X_2) = C(u_1, u_2; \theta)$ where C is a bivariate copula function with parameters θ , then $F(X_1 \mid X_2) = h(u_1 \mid u_2; \theta)$ is defined as in equation 22. For $j \leq t$, application of this property gives the following recursive relationships:

$$u_{t|j} = F(x_t \mid x_{t-1}, \cdots, x_j) = h_{t,j}(U_{t|j+1} \mid U_{j|t-1}; \theta_{t,j})$$
(24)

$$u_{j|t} = F(x_j \mid x_t, \cdots, x_{j+1}) = h_{t,j}(U_{j|t-1} \mid U_{t|j+1}; \theta_{t,j}) .$$
(25)

From Equations (24) and (25), all values of $u_{j|t}$ and $u_{t|j}$ can be obtained, since they correspond to a forward and backward recursion respectively. A more detailed explanation of this process is presented in algorithm 4.

> for $t \leftarrow 1$ to T do | Set $u_{t|t} = F(x_t)$ end for $k \leftarrow 1$ to T - 1 do | for $i \leftarrow k + 1$ to T do | Backward recursion : $u_{i|i-k} = h_{i,i-k}(u_{i|i-k+1} \mid u_{i-k|i-1}; \theta_{i,i-k})$ Forward recursion : $u_{i-k|i} = h_{i,i-k}(u_{i-k|i+1} \mid u_{i|i-k+1}; \theta_{i,i-k})$ end

end

278

Algorithm 4: Recursion Algorithm for determining
$$u_{j|t}$$
 and $u_{t|j}$

The conditional distribution function of x_t given the previous values, can be obtained from Equation (24) as $F(x_t \mid x_{t-1}, \dots, x_1) = u_{t|1} = h_{t,1}(u_{t|2} \mid u_{1|t-1}; \theta_{t,1})$ where $u_{t|2} = F(x_t \mid x_{t-1}, \dots, x_2)$. Recursively, the conditional distribution function of X_t given the previous variables on the series can be expressed as:

$$F(x_t \mid x_{t-a}, \cdots, x_1) = h_{t,1} \quad o \quad h_{t,2} \quad o \quad \cdots \quad h_{t,t-1} \quad o \quad F(x_t) \;. \tag{26}$$

For evaluating all $h_{t,j}$ functions, the values of $u_{1|t-1}, \dots, u_{t-1|t-1}$ must be computed, hence the importance of Algorithm 4. Assuming that the series is Markovian of order p, the distribution function can be simplified as $F(x_t \mid x_{t-1}, \dots, x_1) = F(x_t \mid x_{t-1}, \dots, x_{t-q})$ where $q = \min(p, t-1)$. Therefore, uniformly distributed random numbers between 0 and 1 (w_t) are generated and a realization of the time series is computed as $x_t = F^{-1}(w_t \mid x_{t-1}, \dots, x_{t-q})$. This methodology for univariate time series simulation is presented in algorithm 5.

> for $t \leftarrow 1$ to T do Generate $w_t \sim \text{Uniform}(0,1)$ if t = 1 then $\begin{vmatrix} \text{Set } x_1 = F^{-1}(w_1) \\ \text{else} \\ \mid \text{Set } x_t = F^{-1} \ o \ h_{t,t-1}^{-1} \ o \ \cdots \ o \ h_{t,t-q}^{-1}(w_t \mid u_{i-q\mid i-1}) \\ \text{end} \end{vmatrix}$

286

end

Algorithm 5: Simulation Algorithm for univariate time series

As the objective of this study is to simulate multi-site rainfall time series using the 287 estimation of the COPAR model, algorithm 5 must be modified for simulating a multivariate 288 time series. Although we present the methodology for a bivariate case, this procedure can be 289 extended to any number of variables as stated before. First consider $\{X_t\}$ as the time series 290 of a particular rainfall station and $F_X(\cdot)$ as its marginal distribution function. Also consider 291 Y_t as the series of a different rainfall station for $t = 1, \dots, T$ and $F_V(\cdot)$ as its marginal 292 distribution function. We re-order the elements of both time series into a univariate series 293 $W = (W_1, \dots, W_N)$ with N = 2T in which we interpolate the values of both series i.e. W =294 $(X_1, Y_1, \dots, X_T, Y_T)$. Recalling the model from section 2.2.2, we fit a COPAR (k) model to 295 the data in which the serial dependence of $\{X_t\}$ is modeled unconditionally and that of Y_t is 296 modeled conditionally on $\{X_t\}$. Then, we store the copulas in two categories. The first one 297 $(C_{t,i}^{(1)})$ corresponds to the copulas from Equations (12) and (14), i.e. the ones for modeling 298 $\{X_t\}$ given previous values of $\{X_t\}$ and $\{Y_t\}$, whereas the second category $(C_{t,j}^{(2)})$ corresponds 299 to the copulas from Equations (13) and (15), i.e. the ones for modeling $\{Y_t\}$ given previous 300

values of $\{X_t\}$ and $\{Y_t\}$. The simulation algorithm is modified in the following manner:

for $t \leftarrow 1$ to 2T do Generate $\omega_t \sim \text{Uniform}(0,1)$ Set $q = \min(((k+1)*2) - 1, t - 1)$ Set $m = t \mod 2$ if t = 1 then | Set $x_1 = F_V^{-1}(\omega_1)$ else if m = 1 then | Set $x_t = F_V^{-1} \circ h_{t,t-1}^{-1(1)} \circ \cdots \circ h_{t,t-q}^{-1(1)}(\omega_t \mid u_{i-q|i-1})$ else | Set $x_t = F_{\theta}^{-1} \circ h_{t,t-1}^{-1(2)} \circ \cdots \circ h_{t,t-q}^{-1(2)}(\omega_t \mid u_{i-q|i-1})$ end end

end

302

Algorithm 6: Simulation Algorithm for bivariate time series

³⁰³ **3** Implementation and Analysis of Results

In this section, we present the case study as well as the results and the comparisons 304 between the observed and simulated rainfall time series. We analyze rainfall time series from 305 ground measured data in three stations in the province of Trentino, Italy, with locations 306 shown in 5. The first site is located at Pergine Valsugana (Convento) (Variable Z) in 307 Latitude (46.01°N) and Longitude (11.23°E) with an elevation of 475m. The second site is 308 located at Cles (Convento) (Variable Y) in Latitude (46.35°N) and Longitude (11.02°E) 309 with an elevation of 665m. Finally, the last site is located at Trento (Laste) (Variable X) 310 in Latitude (46.07°N) and Longitude (11.01°E) with an elevation of 312m. The province 311 of Trentino is a region with multiple variable meteorological characteristics. These ground 312 stations measured data that consist of daily time series from 1961 to 1990. The data 313 set trentino is available in the R software package RMAWGEN (Cordano & Eccel, 2017). 314 These stations were chosen in particular because of the heterogeneity of the statistics of 315 the records and also due to their high spatial correlation. The resolution of the data is 316 0.1mm, so that any observation less than this threshold is taken to be a dry day. The data 317 was preprocessed and cleaned for outliers and missing observations. The methodologies or 318

- algorithms proposed in this paper do not depend on the selected data; this data was selected
- with the purpose of showcasing the characteristics of the observed data compared to the
- 321 simulated data.



Figure 5. Case study stations in Trento, Italy

The first part of the results consists of the analysis of the marginal distribution estima-322 tions. Figure 6 shows the results for the parameter estimation for each of the parameters 323 needed in the ZAGA Distribution using the VGLM methodology. Parameters μ and σ are 324 related to the gamma distribution (shape and scale respectively), and parameter ν is re-325 lated to the probability of the Bernoulli distribution that models the occurrence portion of 326 the process. These results show the seasonality of the process throughout the months of 327 the year; moreover, we can see that this site in particular presents a lower probability of 328 rainfall at midyear, and a higher amount of rainfall in the rest of the months of the year. 329 This phenomena is also present in the observed and simulated data. Figure 7 present the 330 density plots for the non-zero values of the rainfall data for the observed and simulated time 331 series. We can see that the simulation data correctly models the density function of the 332 monthly non-zero rainfall time series. Although the max values per month are sometimes 333 underestimated due to the restrictions imposed by the gamma distribution used to model 334 the non-zero characteristics, we can see that the methodology correctly models the reduc-335

tion of rainfall amount at midyear and the increase of rainfall amount for the other months.
Further improvement in this matter can be assessed and a further discussion can be found
in Section 4.

339

A further analysis can be conducted when analyzing the mean, standard deviation and 340 density of the data alongside its characteristics. Table 2 presents the sample mean and 341 standard deviation of the rainfall time series for the Pergine Valsugana site (variable Z) 342 for both observed (Obs) and simulated (Sim) values discriminated for every month for all 343 the years of the study. We also calculated the column m(mean) and m(sd) which are the 344 relative error of the simulation compared to the observed values for the mean and standard 345 deviation on each of the months (Sarmiento et al., 2018). The discrepancies of the mean 346 and standard deviation (SD) could be caused by the difference of the maximum values mod-347 eled by the zero adjusted gamma distribution as explained before; however, for most of the 348 months the mean and standard deviation of rainfall is captured adequately. The respective 349 results for the Cles and Trento sites (Variables Y and X) can be found in Section 5. An in 350 depth analysis of the monthly max values for the observed and simulated data is presented 351 in Figure 8. This figure shows the empirical distributions of the monthly maximum values 352 for the observed and simulated rainfall for the variable Z. When looking at the monthly 353 maximum values, the simulation method has a tendency to produce some sparse values. 354 Although more research needs to be done on this issue, we hypothesize that this is an effect 355 of the type of continuous gamma distribution used to model the non-zero values of the dis-356 tribution, since the actual process of the data has longer tails. The respective results for the 357 Cles and Trento sites (Variables Y and X) can be found in Section 5; moreover, as stated 358 before, an in depth discussion of this matter is conducted in Section 4. Figure 9 include 359 a time-series plot, a box-plot and an histogram of the rainfall for the daily, monthly and 360 yearly observed data as well as the simulated data for the Pergine Valsugana site using the 361 methodology proposed in this paper. These results show that the distribution of the vari-362 ables could be reproduced using this methodology and also the methodology could model 363 the non-stationarity characteristics of the series. The respective results of Variables Y and 364 X for this analysis can be found in Section 5. 365

366

Following these results, another key factor for comparing the simulated and observed data is observing the autocorrelation and partial autocorrelation functions as well as the correlation between sites. Figure 10 show the autocorrelation and partial autocorrelation
function for the sites of the study. This figure shows that the two lags modeled are simulated correctly; moreover, the correlation between sites is presented in Table 1. Finally,
Figure 11 show the 2D density plots of the bivariate distribution for variables Z and Y. The
correspondent 2D density plots for the other bivariate distributions (Z - X and Y - X) can
be found in Section 5.

375

	Pergine Valsugana (Variable Z)	Cles (Variable Y)	Trento (Variable X)
(Variable Z)	1.0000	0.8270	0.7605
(Variable Y)	0.8270	1.0000	0.8535
(Variable X)	0.7605	0.8535	1.0000
(Variable Z)	1.0000	0.7732	0.7143
(Variable Y)	0.7732	1.0000	0.8522
(Variable X)	0.7143	0.8522	1.0000

Table 1. Observed (TOP) and simulated (BOTTOM) data linear correlation for the 3 sites



Figure 6. VGAM parameter estimation of ZAGA distribution for Pergine Valsugana



Figure 7. Observed (TOP) and Simulated (BOTTOM) non-zero rainfall mean density plot Pergine Valsugana

	Month	Obs mean	Obs sd	Sim Mean	Sim sd	m(mean)	m(sd)
1	Jan	1.82	6.68	1.46	4.55	-0.20	-0.32
2	Feb	2.09	8.70	1.99	7.43	-0.05	-0.15
3	Mar	2.40	7.60	2.43	6.92	0.01	-0.09
4	Apr	3.21	8.95	2.05	5.62	-0.36	-0.37
5	May	3.77	9.16	2.95	7.73	-0.22	-0.16
6	Jun	3.23	6.99	2.53	6.68	-0.22	-0.04
7	Jul	2.90	7.11	2.66	7.37	-0.08	0.04
8	Aug	3.16	7.97	3.67	8.74	0.16	0.10
9	Sep	2.79	8.96	2.42	7.90	-0.13	-0.12
10	Oct	3.03	9.85	2.94	10.71	-0.03	0.09
11	Nov	3.36	10.03	3.18	9.60	-0.05	-0.04
12	Dec	1.75	6.75	0.98	5.02	-0.44	-0.26

Table 2. Monthy mean and sd comparisson Pergine Valsugana (Convento)



Figure 8. Monthly maximum rain fall plot Pergine Valsugana



Figure 9. Observed (TOP) and simulated (BOTTOM) monthly plots Pergine Valsugana (Convento)



Figure 10. Simple (ACF) and Partial autocorrelation function (PACF) plots for observed and simulated data.



Figure 11. Uniform bivariate density observed (LEFT) and simulated (RIGHT) data variables Z - Y

³⁷⁶ 4 Discussion and Conclusions

We proposed a rainfall time series simulation strategy using the copula autoregressive 377 methodology (COPAR) proposed by Brechmann et al. (2015) implemented into the discrete-378 continuous conditional bivariate estimation problem. The results presented in this paper are 379 promising and might help to better understand the rainfall simulation phenomena. Overall, 380 the stochastic properties of the series are replicated in terms of the spatial and temporal 381 dependence, as well as the marginal distributions. Using copula functions to model the 382 pairwise dependence of all variables provides a flexible framework to generate more realistic 383 series. This is true in particular because: (i) the model can be estimated in one stage, 384 without separating the occurrence and the amount processes; (ii) copula functions may 385 exploit non-linear relations that are flexible and solve one of the fundamental problems 386 when using transformations of linear (gaussian) time series; and (iii) the whole multivariate 387 time series can be integrally modeled using the same copula based model, no matter the 388 number of dimensions (locations) and the number of lags in the autoregressive process, 389 meaning that no post-estimation introduction of spatial dependency is necessary. 390

It is important to notice that in our research, we found that the gamma distribution 391 underestimates the maximum values of the simulated rainfall series. This could be improved 392 by using a different zero adjusted distribution that better fits the non-zero long tailed 393 characteristics of the data. Future research in this matter would be to explore the properties 394 of the proposed model with hybrid exponential and generalized Pareto distributions, as 395 proposed in Li et al. (2013). In addition, the exploration of copula functions could be more 396 flexible to account for spatial and temporal tail dependence, such as meta-elliptical copulas 397 (Genest et al., 2007) or v-transformed copulas (Bárdossy & Pegram, 2009). 398

A main difference from previous approaches used to model rainfall time series is that the methodology we propose allows us to estimate the inner and cross dependencies between the different sites directly from the conditional dependencies of the rainfall observations. This estimation is assessed by one model that accounts for all the dependence structures without the need of partitioning the estimation procedure into several models. Although this is an advantage in terms of the reduction of overfitting, the estimation of the pair copula construction can be computationally expensive.

The construction of all the pair copula functions needed for this study as well as the simulation procedure had a total computational time of around 30 minutes using a 2.2 Ghz

Intel core i7-4702MQ 7 Gen processor with 16Gb of RAM. An automated algorithm for 408 generating the complete pair-copula R-vine structure for each month and the simulation 409 structure required was also generated. When modeling the discrete-continuous conditional 410 bivariate distribution, we use the empirical estimation of the distribution function for the 411 joint marginal P_{01} and P_{10} distributions. A further improvement in the method could be 412 gained by upgrading this approach, for instance, by using right skewed distributions with 413 longer tails. However, it is important to state that the results found using the empirical 414 approach were satisfactory. 415

Finally, we performed different experiments using non parametric Bernstein copula 416 functions, but the computational time was severely compromised; moreover, the grid used 417 for the discretization of the unitary square may have resulted in some sparse large values 418 that could imply over estimations. Nevertheless, after using parametric copula functions or 419 combinations of these, we mitigated the issue and improved the computational time signifi-420 cantly. For future work, this methodology could be applied to a random fields model where 421 rainfall series could be interpolated in a certain space correctly maintaining the stochastic 422 properties and nature of the data. 423

424 Acknowledgments

This research did not receive any specific grant from funding agencies in the public, commercial, or non-for-profit sectors. However, the authors express their gratitude to the colleagues in the Industrial Engineering Department at University of los Andes for their valuable comments. The dataset used to validated the methodology proposed in the paper is available in the R software package RMAWGEN available in the CRAN repository (Cordano & Eccel, 2017). The results provided in this study can be replicated with this data as explained in the text.

432 Declaration of interests: None.

433 References

- Ahn, K.-H. (2020). Coupled annual and daily multivariate and multisite stochastic weather
 generator to preserve low- and high-frequency variability to assess climate vulnerability.
 Journal of Hydrology, 581, 124443. Retrieved from http://www.sciencedirect.com/
 science/article/pii/S0022169419311783 doi: https://doi.org/10.1016/j.jhydrol.2019
 .124443
- Ayar, P. V., Blanchet, J., Paquet, E., & Penot, D. (2020). Space-time simulation of
 precipitation based on weather pattern sub-sampling and meta-gaussian model. *Journal* of Hydrology, 581, 124451. Retrieved from http://www.sciencedirect.com/science/
- 442 article/pii/S0022169419311862 doi: https://doi.org/10.1016/j.jhydrol.2019.124451
- Balistrocchi, M., & Bacchi, B. (2011). Modelling the statistical dependence of rainfall event
 variables through copula functions. *Hydrology & Earth System Sciences*, 15(6).
- Bárdossy, A., et al. (2017). Asymmetric dependence based spatial copula models: empir-*ical investigations and consequences on precipitation fields*. Stuttgart: Eigenverlag des
 Instituts für Wasser-und Umweltsystemmodellierung
- Bárdossy, A., & Pegram, G. (2009). Copula based multisite model for daily precipitation
 simulation. Hydrology & Earth System Sciences, 13(12).
- Bedford, C. R. M., Tim. (2001). Probability density decomposition for conditionally dependent random variables modeled by vines. Journal of Econometrics, 32(1), 245-268.
 Retrieved from http://www.rogermcooke.net/rogermcooke_files/aivines.pdf
- 453 Brechmann, Christian, E., & Czado, C. (2015). Copar—multivariate time series model-
- ing using the copula autoregressive model. Applied Stochastic Models in Business and
 Industry, 31(4). Retrieved from https://doi.org/10.1002/asmb.2043
- Buishand, T. A., & Brandsma, T. (2001). Multisite simulation of daily precipitation and
 temperature in the rhine basin by nearest-neighbor resampling. Water Resources Research, 37(11), 2761–2776.
- Chen, X., & Fan, Y. (2006). Estimation of copula-based semiparametric time series models.
 Journal of Econometrics, 130(2), 307-335. Retrieved from https://doi.org/10.1016/
 j.jeconom.2005.03.004
- 462 Cordano, E., & Eccel, E. (2017). Rmawgen: Multi-site auto-regressive weather generator
 463 [Computer software manual]. Retrieved from https://CRAN.R-project.org/package=
 464 RMAWGEN (R package version 1.3.7)
- ⁴⁶⁵ Faugeras, O. P. (2012, November). *Probabilistic constructions of discrete copulas*. Retrieved

from https://hal.archives-ouvertes.fr/hal-00751393 (working paper or preprint) 466 Genest, C., Favre, A.-C., Béliveau, J., & Jacques, C. (2007). Metaelliptical copulas and their 467 use in frequency analysis of multivariate hydrological data. Water Resources Research, 468 43(9).469 Haan, C., Allen, D., & Street, J. (1976). A markov chain model of daily rainfall. Water 470 Resources Research, 12(3), 443-449. 471 Herr, H. D., & Krzysztofowicz, R. (2005). Generic probability distribution of rainfall in 472 space: The bivariate model. Journal of Hydrology, 306(1-4), 234–263. 473 H Joe, D. K. (2011). Dependence modeling: Vine copula handbook. World Scientific. 474 Kao, S.-C., & Govindaraju, R. S. (2008). Trivariate statistical analysis of extreme rainfall 475 events via the plackett family of copulas. Water Resources Research, 44(2). 476 Li, C., Singh, V. P., & Mishra, A. K. (2013). A bivariate mixed distribution with a 477 heavy-tailed component and its application to single-site daily rainfall simulation. Water 478 Resources Research, 49(2), 767–789. 479 Mehrotra, R., Srikanthan, R., & Sharma, A. (2006). A comparison of three stochastic 480 multi-site precipitation occurrence generators. Journal of Hydrology, 331(1-2), 280–292. 481 Mhanna, M., & Bauwens, W. (2012). Stochastic single-site generation of daily and monthly 482 rainfall in the middle east. Meteorological Applications, 19(1), 111–117. 483 Nelsen, R. (1999). Introduction. in: An introduction to copulas. lecture notes in statistics. 484 New York, NY: Springer. Retrieved from https://doi.org/10.1007/978-1-4757-3076 485 -0_{-1} 486 Richardson, C. W. (1981). Stochastic simulation of daily precipitation, temperature, and 487 solar radiation. Water resources research, 17(1), 182–190. 488 Sarmiento, C., Valencia, C., & Akhavan-Tabatabaei, R. (2018). Copula autoregressive 489 methodology for the simulation of wind speed and direction time series. Journal of Wind 490 Engineering and Industrial Aerodynamics, 174, 188–199. 491 Serinaldi, F. (2008). Analysis of inter-gauge dependence by kendall's τ k, upper tail de-492 pendence coefficient, and 2-copulas with application to rainfall fields. Stochastic Environ-493 mental Research and Risk Assessment, 22(6), 671–688. 494 Serinaldi, F. (2009a). Copula-based mixed models for bivariate rainfall data: an empirical 495 study in regression perspective. Stochastic environmental research and risk assessment, 496 23(5), 677-693.497 Serinaldi, F. (2009b). A multisite daily rainfall generator driven by bivariate copula-based 498

- ⁴⁹⁹ mixed distributions. Journal of Geophysical Research: Atmospheres, 114(D10).
- Shimizu, K. (1993). A bivariate mixed lognormal distribution with an analysis of rainfall
 data. Journal of Applied Meteorology, 32(2), 161–171.
- Sklar, A. (1959). Fonctions de répartition à n dimensions et leurs marges. Publications de l'Institut Statistique de l'Université de Paris, 8.
- Smith, M., Min, A., Almeida, C., & Czado, C. (2010). Modeling longitudinal data using a
 pair-copula decomposition of serial dependence. Journal of the American Statistical As-
- *sociation*, 1467-1479. Retrieved from https://doi.org/10.1198/jasa.2010.tm09572
- Srikanthan, R., & McMahon, T. (2001). Stochastic generation of annual, monthly and daily
 climate data: A review. *Hydrol. Earth Syst. Sci*, 5, 653–670.
- Tinaikar, A. (2013, 01). Harvesting energy from rainfall. International Journal of Renewable
 and Sustainable Energy, 2, 130. doi: 10.11648/j.ijrse.20130203.18
- Vaittinada Ayar, P., Blanchet, J., Paquet, E., & Penot, D. (2020). Space-time simulation
 of precipitation based on weather pattern sub-sampling and meta-gaussian model. JHyd,
 581, 124451.
- Vandenberghe, S., Verhoest, N., & De Baets, B. (2010). Fitting bivariate copulas to the
 dependence structure between storm characteristics: A detailed analysis based on 105
 year 10 min rainfall. Water resources research, 46(1).
- Wilks, D. (1998). Multisite generalization of a daily stochastic precipitation generation
 model. *journal of Hydrology*, 210(1-4), 178–191.
- ⁵¹⁹ Wilks, D. S., & Wilby, R. L. (1999). The weather generation game: a review of stochastic ⁵²⁰ weather models. *Progress in physical geography*, 23(3), 329–357.
- Yee, T. W. (2015). Vector generalized linear and additive models. Springer. Retrieved from https://doi.org/10.1007/978-1-4939-2818-7
- Zhang, L., & Singh, V. P. (2007). Bivariate rainfall frequency distributions using
 archimedean copulas. Journal of Hydrology, 332(1), 93 109. Retrieved from http://
 www.sciencedirect.com/science/article/pii/S0022169406003386 doi: https://
 doi.org/10.1016/j.jhydrol.2006.06.033
- Zhang, L., & Singh, V. P. (2019). Rainfall frequency analysis. In Copulas and their
 applications in water resources engineering (p. 367–395). Cambridge University Press.
 doi: 10.1017/9781108565103.011
- Zhou, L., Meng, Y., & Abbaspour, K. C. (2019). A new framework for multi-site stochas tic rainfall generator based on empirical orthogonal function analysis and hilbert-huang

- transform. Journal of Hydrology, 575, 730–742.
- ⁵³³ Zhou, L., Meng, Y., Lu, C., Yin, S., & Ren, D. (2020). A frequency-domain nonstationary
- ⁵³⁴ multi-site rainfall generator for use in hydrological impact assessment. Journal of Hydrol-
- ogy, 585, 124770. Retrieved from http://www.sciencedirect.com/science/article/
- ⁵³⁶ pii/S0022169420302304 doi: https://doi.org/10.1016/j.jhydrol.2020.124770

537 5 Supporting information



Figure 12. Observed and Simulated monthly plots Y variable



Figure 13. Observed and Simulated monthly plots X variable



Figure 14. Monthly maximum rain fall plot Cles and Trento sites



Figure 15. Observed and simulated non-zero rainfall mean density plot Cles



Figure 16. Observed and simulated non-zero rainfall mean density plot Trento

	Month	Obs mean	Obs sd	Sim Mean	Sim sd	m(mean)	m(sd)
1	Jan	1.66	5.97	1.42	4.72	-0.15	-0.21
2	Feb	1.62	6.35	1.66	6.38	0.03	0.01
3	Mar	1.85	5.65	1.79	4.98	-0.03	-0.12
4	Apr	2.54	6.50	1.69	4.69	-0.34	-0.28
5	May	2.97	6.91	2.42	6.81	-0.19	-0.01
6	Jun	3.04	6.91	2.33	6.92	-0.23	0.00
7	Jul	2.43	6.31	2.37	6.17	-0.02	-0.02
8	Aug	2.80	6.95	3.14	7.09	0.12	0.02
9	Sep	2.55	9.10	2.67	10.24	0.05	0.13
10	Oct	3.02	9.33	3.16	11.46	0.05	0.23
11	Nov	3.14	9.14	3.22	9.92	0.02	0.09
12	Dec	1.70	6.09	1.09	5.07	-0.36	-0.17

 Table 3.
 Monthly mean and sd comparisson Y variable

	Month	Obs mean	Obs sd	Sim Mean	Sim sd	m(mean)	m(sd)
1	Jan	1.71	5.61	1.68	4.97	-0.02	-0.11
2	Feb	1.81	5.94	1.75	5.93	-0.03	-0.00
3	Mar	1.96	5.99	2.24	6.47	0.14	0.08
4	Apr	2.66	6.53	2.05	5.95	-0.23	-0.09
5	May	3.65	8.03	2.76	7.23	-0.24	-0.10
6	Jun	3.55	7.02	2.85	7.58	-0.20	0.08
7	Jul	2.67	7.03	3.15	8.48	0.18	0.21
8	Aug	3.35	8.47	4.38	10.03	0.31	0.18
9	Sep	2.82	9.39	2.90	10.10	0.03	0.08
10	Oct	3.12	9.09	3.35	11.33	0.07	0.25
11	Nov	3.35	9.48	3.53	10.42	0.05	0.10
12	Dec	1.74	5.72	1.07	4.45	-0.39	-0.22

 Table 4.
 Montly mean and sd comparisson X variable



Figure 17. Uniform bivariate density observed (LEFT) simulated (RIGHT) data variables Y - X



Figure 18. Uniform bivariate density observed (LEFT) and simulated (RIGHT) data variables Y - X