

1 **Stochastic periodic autoregressive to anything (SPARTA): Modelling and**
2 **simulation of cyclostationary processes with arbitrary marginal distributions**

3
4 **Ioannis Tsoukalas^{1,*}, Andreas Efstratiadis¹, and Christos Makropoulos¹**

5
6 ¹ Department of Water Resources and Environmental Engineering, School of Civil Engineering,
7 National Technical University of Athens, Heron Polytechniou 5, 15780 Zographou, Greece

8 Corresponding author: Ioannis Tsoukalas (itsoukal@mail.ntua.gr)

9
10 **Key points:**

- 11 • Simulation of periodic processes with any marginal distributions
12 • Flexibility in the selection of distribution fitting method
13 • Generation of synthetic time series in univariate or multivariate mode
14 • Accurate preservation of essential statistics and observed dependencies

15 **Keywords:**

16 stochastic simulation, hydrological processes, cyclostationarity, synthetic time series, multivariate
17 autoregressive models, arbitrary marginal distributions, Nataf joint-distribution model, normal to
18 anything, linear correlation, dependence patterns

19 **Index terms:** 1869, 3265, 3270, 1872, 1880

20
21 **Supporting information:**

22 <http://onlinelibrary.wiley.com/store/10.1002/2017WR021394/asset/supinfo/wrcr23047-sup-0001-FigS01.pdf?v=1&s=09a47c3b688bf7e9ef4e01f72aacca90c6cc5f32>

23
24
25 **Citation:**

26 Tsoukalas, I., Efstratiadis, A., & Makropoulos, C. (2017). Stochastic periodic autoregressive to
27 anything (SPARTA): Modeling and simulation of cyclostationary processes with arbitrary
28 marginal distributions. *Water Resources Research*, 53. <https://doi.org/10.1002/2017WR021394>

29
30 **Journal:** *Water Resources Research*

31 Received 26 JUN 2017

32 Accepted 17 DEC 2017

33 Accepted article online 21 DEC 2017

34 **Abstract**

35 Stochastic models in hydrology traditionally aim at reproducing the empirically-derived statistical
36 characteristics of the observed data rather than any specific distribution model that attempts to
37 describe the usually non-Gaussian statistical behavior of the associated processes. SPARTA
38 (Stochastic Periodic AutoRegressive To Anything) offers an alternative and novel approach which
39 allows the explicit representation of each process of interest with any distribution model, while
40 simultaneously establishes dependence patterns that cannot be fully captured by the typical linear
41 stochastic schemes. Cornerstone of the proposed approach is the Nataf joint-distribution model,
42 which is related with the Gaussian copula, combined with Gaussian periodic autoregressive
43 processes. In order to obtain the target stochastic structure, we have also developed a
44 computationally simple and efficient algorithm, based on a hybrid Monte-Carlo procedure that is
45 used to approximate the required equivalent correlation coefficients. Theoretical and practical
46 benefits of the proposed method, contrasted to outcomes from widely-used stochastic models, are
47 demonstrated by means of real-world as well as hypothetical monthly simulation examples
48 involving both univariate and multivariate time series.

49 **Plain Language Summary**

50 Stochastic hydrology and, particularly, the synthesis of long hydrometeorological time series (e.g.,
51 precipitation and streamflow) is of high importance in water-related studies since it enables to
52 account for the intrinsic uncertainty of the associated processes. This in turn provides the means
53 to embed uncertainty within planning and decision making. Typically, stochastic models in
54 hydrology aim in the resemblance of the empirically-derived statistical characteristics of the
55 observed time series rather than in reproducing a specific distribution model. In this work we
56 propose a novel approach termed SPARTA (Stochastic Periodic AutoRegressive To Anything)
57 that allows the simulation of multivariate cyclostationary processes with explicit reproduction of
58 the desirable marginal distributions and correlation structures. Its theoretical background is based
59 on the Nataf joint-distribution model (NDM), a procedure that emerged from operations research
60 and is also related with the Gaussian copula. The theoretical and practical benefits of the proposed
61 method are demonstrated by means of real-world and hypothetical simulation studies, involving
62 the generation of both univariate and multivariate time series.

63 **1 Introduction**

64 According to the classical classification by *Matalas [1975]*, synthetic hydrology constitutes a sub-
65 branch of stochastic hydrology, which is usually credited to the pivotal works conducted by the
66 Harvard water program [*Maass et al., 1962*] and *Thomas and Fiering [1962]*. Early attempts to
67 simulate synthetic time series were based on the theory of stochastic processes and the use of linear
68 stochastic models, accounting for the key peculiarities of hydrometeorological processes, namely
69 periodicity and skewness [*Thomas and Burden, 1963; Matalas, 1967; Fiering and Jackson, 1971;*
70 *Klemeš and Borůvka, 1974*].

71 Typically, the standard hypothesis for synthetic time series generation via such approaches does
72 not lie in the reproduction of a specific distribution, but on the resemblance of the statistical
73 characteristics of the parent historical time series. These are usually expressed in terms of low-
74 order statistics (e.g. mean, variance, skewness) and correlations in time and space [*Matalas and*
75 *Wallis, 1976; Salas, 1993*]. However, for a given set of low-order statistics multiple distribution

76 functions may be represented, thus making the simulation problem only partially defined [cf.
77 *Matalas and Wallis, 1976 p. 66*].

78 On the other hand, theoretical reasons and empirical evidence may impose the preservation of a
79 specific distribution for the modelled processes, as highlighted by *Klemeš and Borůvka [1974]*
80 (our emphasis):

81 *“Simulation of a serially correlated series **with a given marginal distribution** is one of the*
82 *important prerequisites of synthetic hydrology and of its applications to analysis of water*
83 *resource system”.*

84 The generation of synthetic data following specific, typically skewed, distributions becomes even
85 more challenging when aiming to simulate hydrometeorological processes at time scales finer than
86 annual. In that case, the stochastic model should account for all facets of cyclostationarity,
87 involving not only, the stochastic structure of the underlying processes but also their distribution,
88 which may be seasonally-varying.

89 The standard approaches to handle skewness within linear stochastic models can be classified in
90 three categories: (a) explicit methods, (b) transformation methods, and c) implicit methods that
91 produce non-Gaussian innovation terms within the generation scheme. Such approaches suffer
92 from notable, although not so apparent, limitations that in combination with the need to account
93 for non-Gaussian distributions motivated this study.

94 Explicit methods are designed (and hence constrained) to generate realizations from a specific
95 distribution family. Common approaches within stochastic hydrology are the stationary
96 multivariate lag-1 model with Log-Normal distribution, proposed by *Matalas [1967]*, and the
97 gamma-autoregressive (GAR) model of *Lawrance and Lewis [1981]*, as well as its periodic
98 extension [*Fernandez and Salas, 1986*]. We remark that so far GAR is restricted for univariate
99 cases, which is a major limitation, since in most water resources applications multiple processes
100 have to be represented simultaneously.

101 Transformation approaches initially aim to “normalize” the non-Gaussian historical data through
102 a proper transformation function; next, parameter estimation and simulation are performed on the
103 normalized data and the final product, i.e., the synthetic data, are obtained via the inverse
104 transformation [*Salas et al., 1985*]. Early attempts used relatively simple conversions, such as Box-
105 Cox, logarithmic, and alternatives, which is well-known that cannot always ensure a satisfactory
106 normalization (e.g. when the original data are too asymmetric). For this reason, for the case of
107 hydrometeorological data, exhibiting significant skewness, more complex schemes have been
108 proposed, yet involving several unknown parameters and also requiring the use of optimization
109 [e.g., *Koutsoyiannis et al., 2008; Papalexiou et al., 2011*]. In fact, the increase of complexity
110 inevitably raises several questions, namely: How many parameters should be used? How does the
111 sample size affect their estimation? In the case of multivariate and cyclostationary simulations,
112 should we use the same transformation function for all processes and seasons?

113 Nevertheless, even an accurate normalization procedure does not ensure that the inverse
114 transformation (i.e., the normalization – simulation – de-normalization scheme) will preserve both
115 the statistical characteristics and the correlation structure of the original variables [*Salas et al.,*
116 *1980 p. 73; Bras and Rodríguez-Iturbe, 1985; Lall and Sharma, 1996; Sharma et al., 1997*].
117 Actually, it is argued that a general method for normalizing all types of data does not exist
118 [*Papalexiou et al., 2011*]. We could also argue that neither an optimal transformation for each

119 specific process exists (particularly in the multivariate case). Thus, the selection and the parameters
120 of the transformation model are prone to subjectivity and indefiniteness.

121 To avoid such ill-transformations, the common practice has leaned towards incorporating
122 skewness within the generation mechanism of the stochastic model itself. In this context, several
123 *implicit* schemes have been proposed to embed non-Gaussian noise within the innovation term.
124 The first attempts are attributed to *Thomas and Burden [1963]* and *Fiering and Jackson [1971]*
125 who proposed univariate simulation schemes for skewed and periodic streamflow data. Their key
126 assumption is the preservation of the desirable statistical characteristics through the generation of
127 white noise from a given distribution, usually the three-parametric Gamma (Pearson type-III). We
128 remark that such approaches generate *explicitly* gamma-distributed variables for the white noise,
129 while the strict “explicitness” is lost when the latter are synthesized to provide the variables of
130 interest [cf. *Matalas and Wallis, 1976 p. 66*]. Hence, the desirable distribution is only
131 approximately preserved [*Koutsoyiannis and Manetas, 1996*]. After the pioneering works of
132 *Fiering [1964]* and *Matalas and Wallis [1976]*, implicit approaches have been implemented within
133 several linear stochastic models, including the multivariate periodic autoregressive model [cf.
134 *Koutsoyiannis, 1999*], the multivariate symmetric moving average model [*Koutsoyiannis, 2000*]
135 and their integration within Castalia framework [*Efstratiadis et al., 2014*].

136 A well-known alternative to all above categories of linear stochastic models is offered by the so-
137 called *non-parametric* approaches, which aim to reproduce the empirical distributions of the
138 observed processes, typically through resampling of historical data [e.g., *Lall and Sharma, 1996*;
139 *Sharma et al., 1997*; *Srinivas and Srinivasan, 2005*; *Mehrotra et al., 2006*; *Marković et al., 2015*].
140 In the literature, such approaches have gained particular attention when the marginal distributions
141 exhibit bi- or multi-modality, which is usually driven by multiple generation mechanisms [*Lall*
142 *and Sharma, 1996*; *Sharma et al., 1997*]. However, the use of the empirical distributions prohibits
143 from fitting to a theoretical model and extrapolating out of the observed data ranges. The lack of
144 theoretical basis makes also difficult to reproduce long-term persistence and cross-correlations
145 among many variables, with few exceptions [e.g., *Kirsch et al., 2013*; *Herman et al., 2016*].
146 Heuristic solutions to the above limitations, such as the recently introduced optimization-based
147 approach by *Borgomeo et al. [2015]*, are subject to extremely high computational effort, and they
148 are also prone to inherent inefficiencies of optimization algorithms.

149 Another relatively new and promising option is offered by copulas, which have recently been
150 embedded in multivariate stochastic simulation schemes in order to describe complex
151 dependencies among hydrological variables [*Hao and Singh, 2013*; *Chen et al., 2015*]. However,
152 it can be argued that copulas are not directly compatible with linear stochastic models, which rely
153 on Pearson’s correlation coefficient, since they typically employ rank correlation statistics (e.g.,
154 Spearman’s ρ_s , or Kendall’s τ) to describe the dependencies among the variables. Nevertheless,
155 they are more sensitive against sampling uncertainty than classical stochastic schemes, in their
156 attempt to describe complex (i.e., nonlinear) dependencies on the basis of usually limited
157 hydrological data. Furthermore, as many researchers argue (see discussion in the aforementioned
158 papers), they rely on quite complicated and computationally demanding generation schemes,
159 especially in high-dimensional spaces.

160 In order to tackle the aforementioned shortcomings, we have developed an *explicit* method, called
161 Stochastic Periodic AutoRegressive To Anything (SPARTA) model, which offers a *generalized*
162 procedure with solid theoretical background for the generation of cyclostationary processes from
163 *a priori* defined distribution functions that are seasonally-varying. The proposed method builds

164 upon the so-called Nataf joint-distribution model [NDM; [Nataf, 1962](#)], which is generic mapping
 165 procedure, and the AutoRegressive To Anything (ARTA) model, introduced by [Cario and Nelson](#)
 166 [1996] to represent stationary processes with arbitrary marginal distributions and autocorrelation
 167 structure. Initially, ARTA was formulated as univariate and later extended for multivariate
 168 simulations [[Biller and Nelson, 2003](#)]. Both versions involve the simulation of stationary
 169 processes, but they have not been extended to account for cyclostationarity which is *sine qua non*
 170 requirement for hydrological processes.

171 The rationale and computational procedure of SPARTA are described in the next three sections,
 172 where section 2 summarizes the overall methodology, section 3 describes the individual
 173 computational steps, while section 4 emphasizes on the Nataf joint-distribution model and the
 174 associated numerical scheme that has been developed and implemented within SPARTA. In
 175 section 5 we evaluate the proposed method by means of three case studies, involving real-world
 176 and hypothetical simulations. A broader discussion on good modelling practices, as well as the
 177 key conclusions and perspectives of this research are outlined in sections 6 and 7, respectively.

178 2 SPARTA at a glance

179 SPARTA aims at simulating periodic processes from any given marginal distribution and a given
 180 stochastic structure, typically (but not exclusively) expressed in terms of first order
 181 autocorrelations and lag zero cross-correlations. Its fundamental advantage is the explicit
 182 preservation of the theoretical marginal distributions of the processes, in contrast to existing linear
 183 stochastic approaches that preserve the marginal statistics (not the distributions themselves) up to
 184 a specific order, typically the third one (i.e., mean, standard deviation, skewness). Briefly, our
 185 approach involves the simulation of an auxiliary process from the Periodic AutoRegressive (PAR)
 186 family, in the “normal” domain (i.e., Gaussian), which allows accounting for cyclostationarity,
 187 and then its mapping to the “real” domain, via the desired inverse cumulative distribution functions
 188 (ICDFs). More specifically: Let $\underline{x}_{s,t} = [\underline{x}_{s,t}^1, \dots, \underline{x}_{s,t}^m]^T$ be a m -dimensional vector of
 189 cyclostationary stochastic processes to simulate, where $s = 1, \dots, S$ denotes the season (e.g.,
 190 month) and $t = 1, \dots, T$ denotes the aggregated time scale (e.g., year). Each element of $\underline{x}_{s,t}$ is
 191 symbolized $x_{s,t}^i$, where $i = 1, \dots, m$ denotes an individual random process, and $x_{s,t}^i$ denotes its
 192 realization. Herein, index i will be also referred to as “location” or “site”, without necessarily
 193 implying spatial reference. Let also $\rho_{s,s-\tau}^{i,j} := \text{Corr}[\underline{x}_s^i, \underline{x}_{s-\tau}^j]$ be the Pearson coefficient of
 194 correlation among processes i and j , for season s and time lag τ . For instance, when $j = i$ and $\tau \neq 0$,
 195 the quantity ρ represents the autocorrelation of the process for lag τ , while for $j \neq i$ and $\tau = 0$, ρ
 196 represents the cross-correlation between i and j , for zero time lag. Furthermore, when the
 197 superscripts or subscripts of ρ are identical (i.e., when $j = i$ or $\tau = 0$) we may omit repeating them
 198 for convenience (e.g., $\rho_{s,s-\tau}^{i,i}$ may be written as $\rho_{s,s-\tau}^i$ and $\rho_{s,s}^{i,j}$ as $\rho_s^{i,j}$).

199 For each process at each season s and each location i , we assign a specific statistical distribution,
 200 $F_{x_s^i} := P(x_s^i \leq x)$, and also assign target coefficients of correlation, $\rho_{s,s-\tau}^{i,j}$, to preserve within the
 201 proposed generation scheme. The key idea of SPARTA lies in the generation of an auxiliary
 202 process $\underline{z}_{s,t} = [\underline{z}_{s,t}^1, \dots, \underline{z}_{s,t}^m]^T$ from a standard Normal Periodic AutoRegressive process
 203 (symbolized PAR-N), with such parameters that their mapping via the corresponding inverse
 204 marginal distributions (ICDFs) results into processes with the target marginal distributions and the
 205 target correlation structure, i.e.,

$$\underline{x}_{s,t}^i = F_{\underline{x}_s^i}^{-1}[\Phi(\underline{z}_{s,t}^i)] \quad (1)$$

206 where $\Phi(\cdot)$ is the CDF of the standard Gaussian distribution and $F_{\underline{x}_s^i}^{-1}(\cdot)$ denotes the ICDFs of the
 207 target distributions of process i at season s .

208 The main challenge, also encountered in the original model (i.e., ARTA), is the identification of
 209 proper parameters for the auxiliary process in the “normal” domain that reproduce the desired
 210 stochastic structure, after applying the mapping procedure. This arises from the fact that the
 211 Pearson correlation coefficient, which is used to describe all kinds of dependencies within linear
 212 stochastic models (including PAR), cannot be preserved when applying a non-linear monotonic
 213 transformation, such as the ICDF. In particular, Eq. (1) results into underestimation of target
 214 correlations, $\rho_{s,s-\tau}^{i,j}$, when they are directly applied to the auxiliary processes. The origin of this
 215 shortcoming is the fact that the Pearsons’ correlation coefficient (in contrast to rank correlation
 216 statistics) is invariant only under linear transformations [Embrechts et al., 1999 p. 7], while for
 217 any other transformation, the correlation coefficients should be properly adjusted. As we will
 218 discuss later (section 4.1), early works in stochastic hydrology were aware of this issue and
 219 attempted to provide analytical or empirical solutions to this problem, for specific distributions
 220 (e.g., Log-Normal).

221 Following the rationale of ARTA, here we ensure the representation of *any* distribution across
 222 seasons and processes by employing the so-called Nataf joint-distribution model [NDM; Nataf,
 223 1962]. NDM offers a generic solution to the mapping problem, thus assigning suitable coefficients
 224 to the auxiliary processes that will finally attain the desirable correlation after the transformation
 225 to the “real” domain. Here, we employ NDM in order to identify such “equivalent” coefficients,
 226 $\tilde{\rho}_{s,s-\tau}^{i,j}$, to be used within the PAR-N generation procedure. As will be elucidated in section 4, for
 227 their estimation we have developed a hybrid method, on the basis of target CDFs, $F_{\underline{x}_s^i}$, and target
 228 $\rho_{s,s-\tau}^{i,j}$.

229 Summarizing, the implementation of SPARTA comprises five steps:

230 **Step 1:** For each variable i and each season s , specify a suitable target marginal distribution, $F_{\underline{x}_s^i}$,
 231 and also identify the dependencies to be preserved in time and space, as well as the target values
 232 of the associated coefficients of correlation, $\rho_{s,s-\tau}^{i,j}$.

233 **Step 2:** On the basis of the desirable dependencies to preserve (in terms of auto- and cross-
 234 correlations), identify the suitable auxiliary model from the PAR-N family.

235 **Step 3:** Employ NDM to determine the equivalent coefficients of correlation, $\tilde{\rho}_{s,s-\tau}^{i,j}$, for all pairs
 236 of variables that are required by the auxiliary model.

237 **Step 4:** Estimate the parameters of the auxiliary model, on the basis of equivalent correlations, and
 238 run the model to generate the auxiliary Gaussian synthetic time series of $\underline{z}_{s,t}$.

239 **Step 5:** Map the auxiliary process $\underline{z}_{s,t}$ to the actual domain using their ICDFs, i.e., through Eq.
 240 (1), to obtain $\underline{x}_{s,t}$.

241 The above steps are described in section 3, while step 3, which is the core element of the proposed
 242 methodology, is discussed in detail in section 4.

243 3 Insights to the computational procedure

244 3.1 Selection of target marginal distributions and correlations

245 In contrast to classical stochastic approaches, which imply the use of a specific statistical model
246 for the noise, SPARTA allows to employ pre-specified distribution models, in order to describe
247 the statistical structure of the modelled processes themselves and not of the noise, which is an
248 auxiliary process. This flexibility involves the selection of the marginal distributions, $F_{\underline{x}_s^i}$, and the
249 identification of their parameters. In addition, the proposed approach allows for identifying target
250 dependencies to preserve, in time and space, expressed by means of target coefficients of
251 correlation, $\rho_{s,s-\tau}^{i,j}$. We highlight that the specification of the above inputs is not a straightforward
252 decision neither it is advised to be made automatically. As thoroughly discussed in section 6, the
253 modeler should account for multilateral information, based both on historical data and expert
254 judgment, in order to establish a realistic formulation of the stochastic simulation model.

255 3.2 The auxiliary model

256 As mentioned above, the generation procedure of SPARTA requires the synthesis of an auxiliary
257 process $\underline{z}_{s,t}$, which is then mapped to the actual one, i.e., $\underline{x}_{s,t}$. This process has to be
258 cyclostationary (since the underlying process is also cyclostationary) and normal. These premises
259 are fulfilled by standard periodic autoregressive models with normally-distributed noise (PAR-N)
260 of any order [e.g., *Salas and Pegram, 1977; Salas et al., 1985; Salas, 1993*].

261 Although any stochastic scheme from the PAR-N family may be applicable, we pay attention to
262 the PAR(1) process, in order to keep things simple and parsimonious, thus providing an easy to
263 follow narrative. In addition, it is argued that the assumption of a first-order model is well-justified
264 for most of practical applications in hydrology [*Efstratiadis et al., 2014*]. Nevertheless, higher-
265 order models may be cumbersome, because the empirical estimation of joint statistics from
266 historical samples is subject to major uncertainty, usually resulting to ill-posed conditions (e.g.,
267 due to inconsistent autocorrelation structures), which in turn leads to substantial defects within
268 parameter estimation.

269 With respect to cross-correlations, the multivariate PAR(1) model, in its full formulation, preserves
270 both the lag zero and lag one dependencies. However, as *Koutsoyiannis and Manetas [1996]* have
271 shown, for reasons of parsimony it is sufficient using the contemporaneous PAR(1) [*Salas, 1993*
272 p. 19.31], which does not explicitly accounts for lag-one cross-correlations within parameter
273 estimation. This is also advocated by an older study of *Pegram and James [1972]*. For instance, in
274 a four-variable problem with 12 seasons, the full PAR(1) model requires the specification of 264
275 parameters to describe the dependencies among the variables, while the contemporaneous one
276 entails 120.

277 3.3 Estimation of equivalent coefficients of correlation

278 In order to employ the multivariate contemporaneous PAR(1)-N within SPARTA, it is essential to
279 provide the equivalent lag-1 month-to-month correlations (i.e., autocorrelations), $\tilde{\rho}_{s,s-1}^i$, for each
280 process i and season s , as well as the equivalent zero-lag cross-correlations, $\tilde{\rho}_s^{i,j}$, for each pair of
281 processes i and j and season s . We remark that the equivalent correlations differ from the target
282 ones, and they are estimated on the basis of the NDM approach, which is described in detail in
283 section 4.

284 **3.4 Parameter estimation within PAR(1)-N process**

285 **3.4.1 Multivariate contemporaneous case**

286 Keeping the same notation for the auxiliary and actual processes, the multivariate PAR(1) reads
 287 (for convenience, time index t is omitted):

$$\underline{\mathbf{z}}_s = \tilde{\mathbf{A}}_s \underline{\mathbf{z}}_{s-1} + \tilde{\mathbf{B}}_s \underline{\mathbf{w}}_s \quad (2)$$

288 where $\underline{\mathbf{z}}_s = [\underline{z}_s^1, \dots, \underline{z}_s^m]^\top$ is a vector of m stochastic processes in season s , $\tilde{\mathbf{A}}_s, \tilde{\mathbf{B}}_s$ are $m \times m$
 289 parameter matrices that depend on season s , and $\underline{\mathbf{w}}_s = [\underline{w}_s^1, \dots, \underline{w}_s^m]^\top$ is a vector of mutually
 290 independent random variables. By definition, the random process $\underline{\mathbf{z}}_s$ is Gaussian, provided that $\underline{\mathbf{w}}_s$
 291 is generated from the standard normal distribution, i.e., $\underline{\mathbf{w}}_s \sim N(0, 1)$.

292 For each season s , the parameter matrix $\tilde{\mathbf{A}}_s$ is diagonal and contains the equivalent lag-1 month-
 293 to-month correlations, $\tilde{\rho}_{s,s-1}^i$, i.e.,

$$\tilde{\mathbf{A}}_s = \text{diag}(\tilde{\rho}_{s,s-1}^1, \dots, \tilde{\rho}_{s,s-1}^m) \quad (3)$$

294 On the other hand, parameter matrices $\tilde{\mathbf{B}}_s$ are calculated as follows:

$$\tilde{\mathbf{B}}_s \tilde{\mathbf{B}}_s^\top = \tilde{\mathbf{G}}_s \quad (4)$$

295 where $\tilde{\mathbf{G}}_s := \tilde{\mathbf{C}}_s - \tilde{\mathbf{A}}_s \tilde{\mathbf{C}}_{s-1} \tilde{\mathbf{A}}_s^\top$ and $\tilde{\mathbf{C}}_s$ is a symmetric $m \times m$ matrix that contains the equivalent
 296 lag-zero cross-correlations, $\tilde{\rho}_s^{i,j}$, i.e.,

297
$$\tilde{\mathbf{C}}_s = \begin{pmatrix} 1 & \dots & \tilde{\rho}_s^{1,m} \\ \vdots & \ddots & \vdots \\ \tilde{\rho}_s^{m,1} & \dots & 1 \end{pmatrix}$$

298 In order to estimate the parameter matrix $\tilde{\mathbf{B}}_s$, it is essential to solve a decomposition problem, also
 299 expressed as finding the square root of $\tilde{\mathbf{G}}_s$. This can be obtained with the use of typical numerical
 300 techniques, such as Cholesky or singular value decomposition [e.g., *Johnson, 1987*]. We remark
 301 that when $\tilde{\mathbf{G}}_s$ is positive definite, it has infinite number of feasible solutions, such as the solutions
 302 provided by the aforementioned numerical methods. On the other hand, if $\tilde{\mathbf{G}}_s$ is non-positive
 303 definite (this is often the case when the historical data are of different length) the problem does not
 304 have a feasible solution, thus requiring the detection of a parameter matrix $\tilde{\mathbf{B}}_s$ ensuring an
 305 approximation of the given $\tilde{\mathbf{G}}_s$, e.g., through optimization [*Koutsoyiannis, 1999; Higham, 2002*].

306 In particular, *Koutsoyiannis [1999]* has developed an optimization-based approach, paying
 307 attention on the preservation of skewness, which is a well-known trouble of multivariate stochastic
 308 models, asking for generating skewed white noise [e.g., *Todini, 1980*]. A great advantage of our
 309 approach is the assumption of normality within the auxiliary process, which substantially
 310 simplifies the optimization problem within decomposing non-positive definite matrices. More
 311 precisely, the empirical penalty term considered by *Koutsoyiannis [1999]*, in order to prohibit the
 312 generation of highly-skewed white noise, which introduces significant complexity to the
 313 optimization procedure [cf. *Efstratiadis et al., 2014*], is neglected, thus resulting to a “reduced”
 314 objective function that only contains a distance term to minimize.

315 3.4.2 Univariate case

316 The univariate model can easily be derived from the above equations. Since $m = 1$, $\tilde{\mathbf{A}}_s = \tilde{\rho}_{s,s-1}^1$
317 and $\tilde{\mathbf{C}}_s = 1$, thus $\tilde{\mathbf{B}}_s \tilde{\mathbf{B}}_s^T = 1 - \tilde{\rho}_{s,s-1}^1 \tilde{\rho}_{s,s-1}^1$, which leads to $\tilde{\mathbf{B}}_s = \sqrt{1 - \tilde{\rho}_{s,s-1}^1{}^2}$. Hence, by
318 substituting in Eq. (2) and removing the redundant indices we read:

$$\underline{z}_s = \tilde{\rho}_{s,s-1} \underline{z}_{s-1} + \sqrt{1 - \tilde{\rho}_{s,s-1}^2} \underline{w}_s \quad (5)$$

319 where \underline{w}_s are i.i.d. white noise with $N \sim (0, 1)$. We remark that since $i = 1$ the superscript of $\tilde{\rho}(\cdot)$
320 has been omitted for simplicity.

321 3.5 Mapping auxiliary processes to the actual domain

322 After generating the synthetic time series of the auxiliary processes \underline{z}_s , the last step is its mapping
323 throughout Eq. (1) to the actual domain \underline{x}_s , through the inverse CDFs. This procedure is
324 implemented for each individual process and season. Due to the use of the inverse CDF, as well
325 as the use of equivalent coefficients of correlation within the PAR(1)-N model, the resulting data
326 will preserve both the target marginal distributions, for all seasons and locations, as well as the
327 target auto- and cross-correlations. Even in case of non-positive definite correlation matrices,
328 where the desired stochastic characteristics are not explicitly preserved by the PAR(1)-N model,
329 the “reduced” optimization approach ensures a very good approximation, with minimal
330 computational burden.

331 4 Nataf joint-distribution model and computational advances

332 4.1 Historical summary and rationale

333 The problem of obtaining a joint pdf of random variables based on their individual distributions
334 and correlation has long been discussed within the statistical community. *Nataf [1962]* has
335 proposed a quite simple, yet general solution by mapping multivariate normal variables with a
336 given correlation matrix to multivariate uniform variables, which in turn are mapped to the desired
337 distributions via the corresponding inverse cumulative functions. The key challenge is to identify
338 the equivalent correlations to be applied within the generation of random variables in the normal
339 domain, in order to attain the desired correlation in the real domain. In their classical work, *Liu*
340 *and Der Kiureghian [1986]* showed that the Nataf’s Distribution Model (NDM) is suitable for
341 describing a wide range of correlation values. Later, *Cario and Nelson [1997]*, developed a
342 generalized procedure based on NDM and referred to as NORTA (NORmal To Anything), for the
343 generation of correlated random vectors with arbitrary marginal distributions, including discrete
344 and mixed ones. In fact, NDM may be considered as a specific case of copulas [*Sklar, 1973*], and
345 more specifically the Gaussian one. In fact, linear stochastics are compatible with the latter copula,
346 since both use the Pearson’s linear correlation as measure of dependence. *Lebrun and Dutfoy*
347 *[2009]*, in view of copula theory, provide an extensive and insightful discussion on the relation of
348 NDM with the Gaussian copula, as well as provide an alternative formulation of the former in
349 terms of Spearman’s ρ_s and Kendall’s τ .

350 We remark that when *Cario and Nelson [1997]* have published their work, they argued that the
351 generality of their approach came at the cost of computational efficiency (i.e., computational time),

352 since the estimation $\tilde{\rho}$ presupposed solving numerically a double integral in the infinite domain.
 353 However, this argument is far from interest now, grace to continuous advances in computing,
 354 which have significantly contributed in waiving such barriers.

355 4.2 Theoretical background

356 In the general case, let that we wish to generate a correlated random vector $\underline{\mathbf{x}} = [\underline{x}_1, \dots, \underline{x}_k, \dots, \underline{x}_m]^T$
 357 with target marginal distributions $F_{\underline{x}_k}$ and target correlation matrix:

$$358 \quad \mathbf{C}_{\underline{\mathbf{x}}} = \begin{pmatrix} 1 & \cdots & \rho_{1,m} \\ \vdots & \ddots & \vdots \\ \rho_{m,1} & \cdots & 1 \end{pmatrix}$$

359 Let also $\underline{\mathbf{z}} = [\underline{z}_1, \dots, \underline{z}_k, \dots, \underline{z}_m]^T$ be a multivariate normal vector with correlation matrix
 360 (equivalent):

$$361 \quad \tilde{\mathbf{C}}_{\underline{\mathbf{z}}} = \begin{pmatrix} 1 & \cdots & \tilde{\rho}_{1,m} \\ \vdots & \ddots & \vdots \\ \tilde{\rho}_{m,1} & \cdots & 1 \end{pmatrix}$$

362 In order to obtain $\underline{\mathbf{x}}$ through $\underline{\mathbf{z}}$ the following mapping equation is employed:

$$\underline{x}_k = F_{\underline{x}_k}^{-1}[\Phi(\underline{z}_k)] \quad (6)$$

363 where $F_{\underline{x}_k}^{-1}$ is the ICDF of variable k and $\Phi(\cdot)$ is the standard normal CDF. A direct outcome of Eq.
 364 (6) is that for two variables \underline{x}_k and \underline{x}_l their correlation is given by:

$$\text{Corr}[\underline{x}_k, \underline{x}_l] = \rho_{k,l} = \text{Corr}\left[F_{\underline{x}_k}^{-1}[\Phi(\underline{z}_k)], F_{\underline{x}_l}^{-1}[\Phi(\underline{z}_l)]\right] \quad (7)$$

365 thus the target correlations $\rho_{k,l}$ are associated with the unknowns $\tilde{\rho}_{k,l}$.

366 An apparent approach could be setting $\tilde{\mathbf{C}}_{\underline{\mathbf{z}}} \equiv \mathbf{C}_{\underline{\mathbf{x}}}$. However, both theoretical and empirical evidence
 367 have indicated that this assumption will result to systematically underestimated correlations within
 368 the synthetic data. The theoretical justification of this behavior stems from the Pearson correlation
 369 coefficient itself, since it is not invariant under non-linear monotonic transformations, such as
 370 those imposed by the ICDFs [Embrechts et al., 1999 p. 8]. More specifically, the largest the
 371 departure of the actual distribution, $F_{\underline{x}_k}$, from the normal one, the largest will be the
 372 underestimation. Therefore, and except the trivial normal case, in order to eliminate biases, we
 373 should assign *a priori* larger values to $\tilde{\rho}_{k,l}$.

374 Hopefully, NDM and its theoretical background can provide a theoretical solution to the above
 375 problem by means of an appropriate correlation matrix $\tilde{\mathbf{C}}_{\underline{\mathbf{z}}}$ that leads to the target correlation matrix
 376 $\mathbf{C}_{\underline{\mathbf{x}}}$. As highlighted by Liu and Der Kiureghian [1986], in order to employ NDM it is essential to
 377 ensure 1) one to one mapping of Eq. (6), and 2) positive definite correlation matrix $\tilde{\mathbf{C}}_{\underline{\mathbf{z}}}$. The former
 378 requirement is by definition valid in typical case of continuous distributions used in hydrology,
 379 while the latter is also usually satisfied, since the distances $\varepsilon_{k,l} := |\rho_{k,l} - \tilde{\rho}_{k,l}|$ are expected to be
 380 generally small (provided, of course, that the target matrix $\mathbf{C}_{\underline{\mathbf{x}}}$ is positive definite).

381 The following procedure is applied to each specific pair of variables \underline{x}_k and \underline{x}_l (i.e., $m(m -$
 382 $1)/2$ times). Given that

$$\text{Corr}[\underline{x}_k, \underline{x}_l] = \rho_{k,l} = \frac{E[\underline{x}_k, \underline{x}_l] - E[\underline{x}_k]E[\underline{x}_l]}{\sqrt{\text{Var}[\underline{x}_k]\text{Var}[\underline{x}_l]}} \quad (8)$$

383 where $E[\underline{x}_k], E[\underline{x}_l]$ and $\text{Var}[\underline{x}_k], \text{Var}[\underline{x}_l]$ are the mean and variance of \underline{x}_k and \underline{x}_l respectively,
 384 which are obviously known since the associated marginal distributions are already specified (and
 385 have finite moments, otherwise the Pearson correlation coefficient cannot be defined) the
 386 computational procedure is limited to identifying $E[\underline{x}_k, \underline{x}_l]$. Since the corresponding variables to
 387 be mapped, \underline{z}_k and \underline{z}_l , respectively, are by definition normally distributed, with
 388 correlation $\text{Corr}[\underline{z}_k, \underline{z}_l] = \tilde{\rho}_{k,l}$, then, using (6) and the first cross-product moment of \underline{x}_k and \underline{x}_l we
 389 get:

$$\begin{aligned} E[\underline{x}_k, \underline{x}_l] &= E \left[F_{\underline{x}_k}^{-1}[\Phi(\underline{z}_k)] F_{\underline{x}_l}^{-1}[\Phi(\underline{z}_l)] \right] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} F_{\underline{x}_k}^{-1}[\Phi(z_k)] F_{\underline{x}_l}^{-1}[\Phi(z_l)] \varphi_2(z_k, z_l, \tilde{\rho}_{k,l}) dz_k dz_l \end{aligned} \quad (9)$$

390 where $\varphi_2(z_k, z_l, \tilde{\rho}_{k,l})$ is the bivariate standard normal probability density function. Therefore, each
 391 target $\rho_{k,l}$ is a function of $\tilde{\rho}_{k,l}$, which is embedded in $\varphi_2(z_k, z_l, \tilde{\rho}_{k,l})$, and the given marginal
 392 distributions $F_{\underline{x}_k}$ and $F_{\underline{x}_l}$, i.e.,

$$\rho_{k,l} = \mathcal{F}(\tilde{\rho}_{k,l} | F_{\underline{x}_k}, F_{\underline{x}_l}). \quad (10)$$

393 Unfortunately, Eq. (10) cannot be analytically derived from Eq. (9), with the exception of few
 394 special cases [Li and Hammond, 1975; Cario and Nelson, 1997]. Among them the Log-Normal
 395 case [Mostafa and Mahmoud, 1964] which is of particular interest in hydrology. The
 396 aforementioned researchers, as well as Liu and Der Kiureghian [1986], provided several Lemmas
 397 that can be useful in order to approximate Eq. (10). Among them,

398 **Lemma 1:** $\rho_{k,l}$ is a strictly increasing function of $\tilde{\rho}_{k,l}$.

399 **Lemma 2:** $\tilde{\rho}_{k,l} = 0$ for $\rho_{k,l} = 0$ as well as, $\tilde{\rho}_{k,l} \geq (\leq) 0$ if $\rho_{k,l} \geq (\leq) 0$.

400 **Lemma 3:** $|\rho_{k,l}| \leq |\tilde{\rho}_{k,l}|$.

401 Note that in Lemma 3, the equality sign is valid when $\rho_{k,l} = 0$ or when both marginal distributions
 402 are normal. Furthermore, the minimum and maximum attainable value of $\rho_{k,l}$ is given for $\tilde{\rho}_{k,l} =$
 403 -1 and $\tilde{\rho}_{k,l} = 1$, respectively. The literature offers a variety of approaches to establish $\mathcal{F}(\cdot)$,
 404 including crude search procedures [Cario and Nelson, 1996], methods based on the Gauss-Kronrod
 405 quadrature rule [Cario, 1996], root finding methods [Li and Hammond, 1975; Chen, 2001] as well
 406 as Gauss-Hermite quadrature and Monte-Carlo methods [Xiao, 2014]. Herein, we propose a
 407 simple and easy to implement method based on hybrid combination of Monte-Carlo simulation
 408 and polynomial interpolation.

409 **4.3 Hybrid procedure for solving $\mathcal{F}(\cdot)$**

410 As already mentioned, in order to preserve the target correlations $\rho_{k,l}$ in the actual domain, after
 411 mapping the generated Gaussian values with their prescribed distributions, using Eq. (6), it is
 412 essential to establish a suitable relationship between $\tilde{\rho}_{k,l}$ and $\rho_{k,l}$. In this context, we have
 413 developed the following procedure (indices k and l are omitted for simplicity):

414 **Step 1:** Create a q -dimensional vector $\tilde{\mathbf{r}} = [\tilde{r}^1, \dots, \tilde{r}^q]$ of equally spaced values in the interval
 415 $[r_{\min}, r_{\max}]$. Here, lemma 2 can be accounted for in order to determine the boundaries r_{\min} and r_{\max} ,
 416 since it provides insights regarding the sign of $\tilde{\rho}$. For example, if the target correlation ρ is
 417 positive, then we set $r_{\min} = 0$ and $r_{\max} = 1$.

418 **Step 2:** For each element of $\tilde{\mathbf{r}}$, generate N samples from the bivariate standard normal distribution,
 419 with correlation \tilde{r}^i .

420 **Step 3:** Map the synthetic data to the actual domain through Eq. (6), using the associated target
 421 marginal distribution,

422 **Step 4:** Calculate the empirical correlations r^i and store them in the vector $\mathbf{r} = [r^1, \dots, r^q]$.

423 **Step 5:** Approximate the relationship between target (ρ) and equivalent ($\tilde{\rho}$) correlation by
 424 establishing a polynomial function of order p , among the values of $\tilde{\mathbf{r}}$ and \mathbf{r} i.e.:

$$\rho = \mathcal{F}(\tilde{\rho} | F_{\tilde{x}_k}, F_{\tilde{x}_l}) \cong r = a_p \tilde{r}^p + a_{p-1} \tilde{r}^{p-1} + \dots + a_1 \tilde{r}^1 + a_0 \quad (11)$$

425 **Step 6:** Evaluate the equivalent correlation $\tilde{\rho}_{k,l}$ by inverting the relationship between the fitted
 426 polynomial and the target correlation $\rho_{k,l}$.

427 We highlight that, according to Weierstrass approximation theorem, the formulation of the
 428 polynomial expression (11) is theoretically feasible, since $\mathcal{F}(\cdot)$ is continuous and $\tilde{\mathbf{r}}$ is bounded on
 429 the interval $[-1, 1]$. Moreover, we remark that the constant term a_0 could be omitted, as indicated
 430 by Lemma 2.

431 The above procedure, which is hybrid combination of Monte Carlo simulation and numerical
 432 interpolation through polynomial regression, uses three input arguments, i.e., the vector dimension
 433 q , the sample size N , and the polynomial order p . The first two influence the accuracy and
 434 computational effort of the Monte Carlo procedure, while the third influences the accuracy of the
 435 interpolation approach. Preliminary analysis detected that a good balance between accuracy and
 436 computational efficiency is ensured for q around 10 — 20, and N around 50 000 — 100 000 trials.
 437 Regarding the polynomial order, [Xiao \[2014\]](#) conducted an extensive analysis, with distributions
 438 exhibiting a wide range of skewness and kurtosis coefficients, and concluded that $\mathcal{F}(\cdot)$ can be
 439 accurately approximated by a polynomial of less than ninth degree ($p \leq 9$). Apparently, for $p = q$
 440 $- 1$, the polynomial passes exactly through all simulated points, yet, in order to ensure parsimony,
 441 it may be preferable employing a less complicated expression. In this vein, in order to avoid over-
 442 fitting, we propose adjusting the order of the polynomial with the use of cross-validation
 443 techniques or the Akaike information criterion [[Akaike, 1974](#)]. We note that in the basis of a
 444 systematic study one could identify alternative functions instead of polynomials in order to
 445 describe the relationship $\mathcal{F}(\cdot)$.

446 The key advantage of the proposed methodology, which is applicable for continuous, discrete or
 447 mixed-type distributions, is its simplicity and the fact that it doesn't depend on specialized
 448 algorithms to solve the double integral of Eq. (9), in order to obtain a valid expression $\mathcal{F}(\cdot)$. It is

449 noteworthy that despite the iterative nature of the algorithm, its implementation in high-level
 450 programming languages, such as R or MATLAB, requires less than 1 second (assuming $N =$
 451 $150\,000$ and $m = 20$) on a modest 3.0 GHz Intel Dual-Core i5 processor with 4 GB RAM.

452 4.4 Numerical example

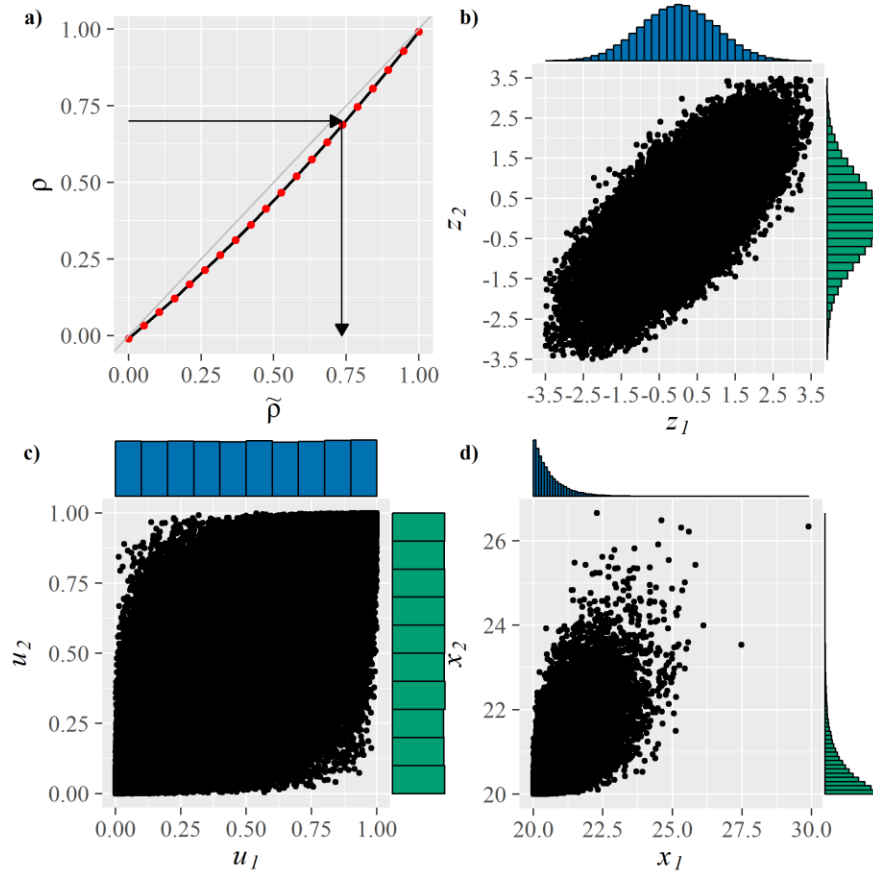
453 Consider a bivariate example with variables \underline{x}_1 and \underline{x}_2 , representing either the same process at two
 454 adjacent seasons or two simultaneous processes at the same season. We assume that the two
 455 variables follow the same target marginal distribution ($F_{\underline{x}_1} \equiv F_{\underline{x}_2}$), i.e., the Pearson type-III, with
 456 probability density function:

$$f(x|a, b, c) = \frac{1}{|b|\Gamma(a)} \left(\frac{x-c}{b}\right)^{a-1} \exp\left(-\frac{x-c}{b}\right) \quad (12)$$

457 where $\Gamma(\cdot)$ is the gamma function, a , b and c are the shape, scale and location parameters
 458 respectively. For both variables, we assume the common parameter values $a = 1$, $b = 20$ and $c =$
 459 0.60 , as well as a target correlation $\rho_{1,2} = 0.70$. Employing the NDM approach and the numerical
 460 method of section 4.3, with $q = 20$, $N = 50\,000$ and $p = 2$, we approximate $\mathcal{F}(\cdot)$ through the
 461 following polynomial (quadratic) function (indices are omitted for simplicity):

$$\rho = \mathcal{F}(\tilde{\rho}|F_{\underline{x}_1}, F_{\underline{x}_2}) \cong 0.2049\tilde{\rho}^2 + 0.7963\tilde{\rho} - 0.0009 \quad (13)$$

462 Given the relationship (13), it is easy to solve for the equivalent correlation $\tilde{\rho}_{1,2}$ which can be used
 463 for the generation of standard normal variables, \underline{z}_1 and \underline{z}_2 , that results to the target value $\rho_{1,2}$. In
 464 particular, for $\rho_{1,2} = 0.7$ and inverting (13), we get $\tilde{\rho}_{1,2} = 0.739$.



465

466 **Figure 1:** Hypothetical example of mapping two correlated variables, where the target and equivalent
 467 correlations represented through Eq. (13) are shown in panel (a). Panels (b), (c) and (d) illustrate the data
 468 in the normal, uniform and actual domain, respectively.

469 The mapping procedure of the numerical example, is shown in Figure 1 for the generation of
 470 100 000 correlated values. In panel (a) we depict the relationship between target and equivalent
 471 correlations as established via Eq. (13). In panel (b) we illustrate the simulated auxiliary Gaussian
 472 variables, \underline{z}_1 and \underline{z}_2 , which are generated by assigning the equivalent correlation $\tilde{\rho}_{1,2}$. Initially,
 473 these variables are mapped to the uniform domain through function $\Phi(\cdot)$ (panel c), and then they
 474 are mapped to the actual domain (panel d), via the corresponding inverse functions, $F_{\underline{x}_1}^{-1}$ and $F_{\underline{x}_2}^{-1}$.
 475 Within the two mapping procedures, the equivalent correlation $\tilde{\rho}_{1,2}$ is progressively decreased,
 476 down to the target value $\rho_{1,2}$.

477 We remark that due to the very large sample size, the empirical correlation between the auxiliary
 478 synthetic variables \underline{z}_1 and \underline{z}_2 coincides the theoretical one, i.e., $\tilde{\rho}_{1,2} = 0.739$, while the empirical
 479 correlation between the actual variables \underline{x}_1 and \underline{x}_2 is 0.707, thus practically identical to the target
 480 value $\rho_{1,2} = 0.70$. Moreover, the empirically estimated parameters of the derived distributions are
 481 $a = 0.947$, $b = 20.001$ and $c = 0.622$, for the synthetic variable \underline{x}_1 and $a = 0.921$, $b = 20.000$ and c
 482 $= 0.671$ for \underline{x}_2 . The aforementioned values, which were computed through the maximum
 483 likelihood estimation method (MLE), are in agreement with the theoretical ones.

484 4.5 Coupling SPARTA and NDM

485 It is apparent that in order to align NDM with SPARTA, we just have to set $\underline{x}_k := \underline{x}_s^i$ and $\underline{x}_l :=$
486 $\underline{x}_{s-\tau}^j$ throughout equations (7) to (10), and approximate the required (by the auxiliary model)
487 equivalent correlation coefficients $\tilde{\rho}_{s,s-\tau}^{i,j}$ of the target correlations $\rho_{s,s-\tau}^{i,j}$. For the estimation of the
488 equivalent correlations across all processes and seasons, we also offer the aforementioned hybrid
489 computational procedure to approximate the relationship of Eq. (10), i.e., $\mathcal{F}(\cdot)$.

490 4.6 Previous applications of NDM in hydrology

491 NDM-based approaches have been widely applied in industrial, financial and operations research
492 applications, as indicated from the popularity of the original article by *Nataf [1962]* and the
493 relevant publications [e.g., *Liu and Der Kiureghian, 1986; Cario and Nelson, 1996, 1997; Biller*
494 *and Nelson, 2003*].

495 While hydrological community does not make direct reference to NDM and the associated models,
496 such as NORTA, ARTA, VARTA, etc., it actually shares the same rationale, even from the geneses
497 of hydrological stochastics. Loosely speaking, the core idea of NDM comprises the initiation from
498 the Gaussian domain, with properly adjusted correlation coefficients, and then a mapping to the
499 desirable domain.

500 In particular, *Matalas [1967]* has studied the effects of logarithmic transformations in the context
501 of synthesizing log-normally distributed processes, concluding that the so far prevailing
502 transformation approach failed to resemble the historical statistics. To reestablish consistency, he
503 developed a framework based on the generation of normal processes, and provided a set of
504 theoretical equations to estimate the statistical parameters (including adjusted correlation
505 coefficients) in the Log-Normal domain. Later, *Klemeš and Borůvka [1974]* developed a
506 generation scheme for gamma-distributed univariate first-order Markov chains, through a mapping
507 procedure of Gaussian processes with the use of adjusted correlation coefficients. More recently,
508 *Kelly and Krzysztofowicz [1997]* proposed and illustrated through several hydrology-related
509 applications, a flexible bivariate distribution model, termed meta-Gaussian, which builds upon the
510 bivariate standard normal distribution and the normal quantile transformation. Furthermore, *Wilks*
511 *[1998]*, in the context of his widely known weather generation model, has also employed a
512 transformation procedure initiating from the standard Gaussian distribution, coupled with an
513 empirical method to estimate the adjusted correlations for the simulation of multivariate daily
514 precipitation with mixed exponential distributions. This seminal work has triggered the
515 development of improved schemes, supporting more distributions and correlation structures.
516 Detailed reviews are provided by *Wilks and Wilby [1999]* and *Ailliot et al. [2015]*. Additionally,
517 running advances in stochastic hydrology are also in alignment with NDM. In particular, in a
518 similar vein, *Serinaldi and Lombardo [2017]* proposed a fast procedure for autocorrelated
519 univariate binary processes, while *Lee [2017]* introduced a simulation-based method for Gamma-
520 distributed precipitation. Finally, *Papalexiou [2017]* proposes an elegant and unified overview for
521 synthetic data generation using autoregressive models.

522 5 Case studies

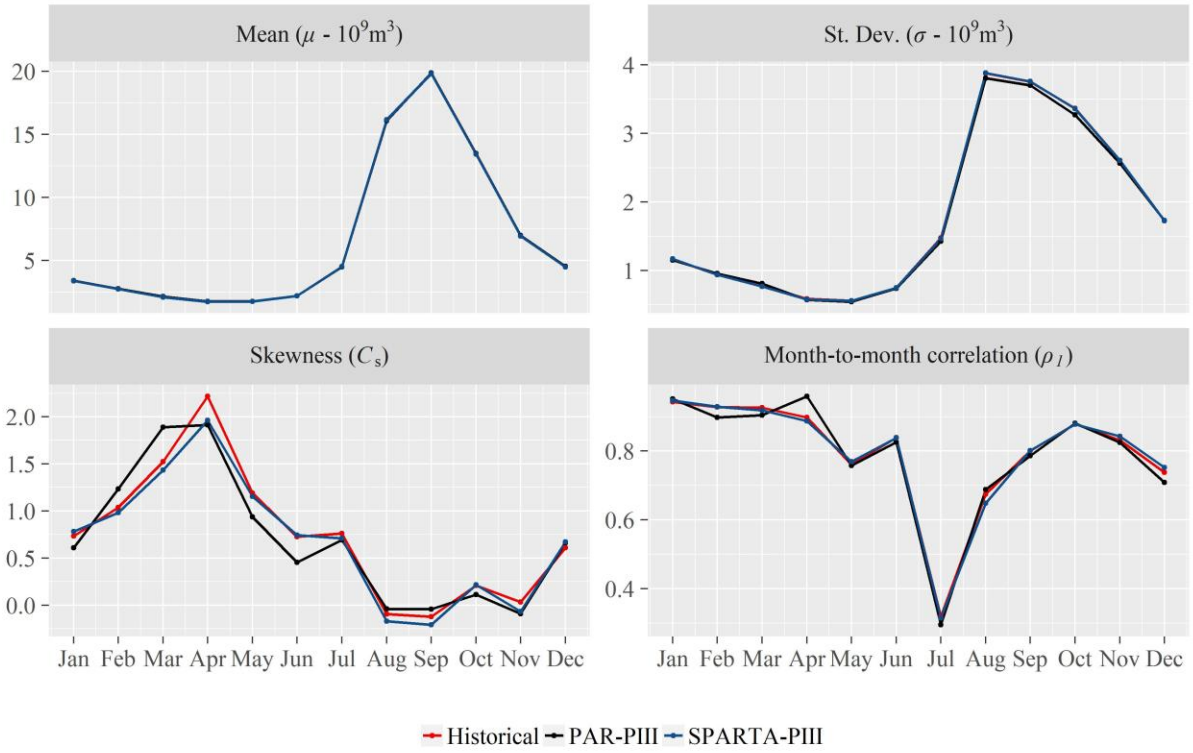
523 5.1 Univariate simulation with common distribution models

524 The first case study involves the simulation of monthly flow of Nile River at Aswan dam, based
525 on a historical dataset from March 1870 to December 1945 [Hipel and McLeod, 1994]. The flows
526 are characterized by strong seasonality and high correlations across all subsequent months (Figure
527 2). In order to demonstrate the performance of SPARTA against PAR, we compare the outcomes
528 of a stochastic simulation scenario of 2 000 years length, which has been used several times in the
529 past for providing synthetic flows [e.g., Koutsoyiannis et al., 2008]. Since PAR(1) is typically
530 coupled with Pearson type-III distribution for white noise generation (referred to as PAR-PIII
531 model), in order to conduct a fair and meaningful evaluation, within SPARTA we also set this
532 distribution as target one for all months (referred to as SPARTA-PIII model). We remind that
533 SPARTA explicitly accounts for the marginal distribution of each season, while PAR-PIII,
534 similarly to most linear stochastic models, attempts to resemble the statistical characteristics via
535 implicitly representing the marginal distributions into the innovation term. We note that the
536 multivariate formulation of PAR-PIII of order 1 is given in Appendix A.

537 It is remarked that due to the use of Pearson type-III distribution, which allows for negative
538 location parameters, the two models can produce negative values that would not be acceptable in
539 a real-world hydrological study. A typical way to address this inconsistency within both models is
540 the artificial truncation of all synthetic values to zero, which would yet introduce bias to the
541 stochastic structure of the synthetic processes. However, among the two models, SPARTA also
542 offers a much more rigorous alternative, since the data are generated via the corresponding ICDFs.
543 The latter property enables fitting another positively bounded distribution model (e.g., Gamma,
544 Log-Normal, etc.) to the observed data that explicitly prohibits the generation of negative values.

545 The two models are evaluated through visual inspection of simulated against observed values of
546 their monthly statistical characteristics, in terms of calculated values of mean, μ , standard
547 deviation, σ , skewness coefficient, C_s , and lag-1 month-to-month correlation, ρ_1 (Figure 2), as well
548 as in terms of their monthly marginal distributions (Figure 3). It is noted that the latter statistics
549 were calculated after truncation of negative values. Except for the trivial case of means and
550 standard deviations, which are perfectly reproduced by both models, for the skewness and month-
551 to-month correlations, only SPARTA-PIII ensures full consistency with the target values across
552 all seasons. In addition, SPARTA-PIII fits perfectly the target theoretical distribution models,
553 which is a direct outcome of employing the inverse mapping, while PAR-PIII occasionally
554 deviates from the target distributions, and particularly their tails (e.g., in February, March, April
555 and May).

556 To further highlight the advantages of SPARTA over PAR-PIII, we also investigate the derived
557 dependence forms, by focusing on the scatter plots of the 12 pairs of adjacent monthly data sets
558 (Figure 4). Interestingly, PAR-PIII, although it preserves quite satisfactory the key statistical
559 characteristics, including the observed coefficients of correlation, it fails to capture the full extent
560 of the observed patterns, in contrast to SPARTA-PIII, which generates well-spread data pairs
561 which are in compliance with the observations. In particular, in the scatter plots of pairs December
562 – January, January – February, February – March and March – April, it is evident that PAR-PIII
563 not only fails to capture the dependence patterns of the historical data, but also seems fails to
564 produce synthetic pairs out of a lower boundary. Therefore, the synthetic dependencies are not in
565 good agreement with the observed ones, although the correlation coefficients themselves are
566 reproduced with high accuracy.

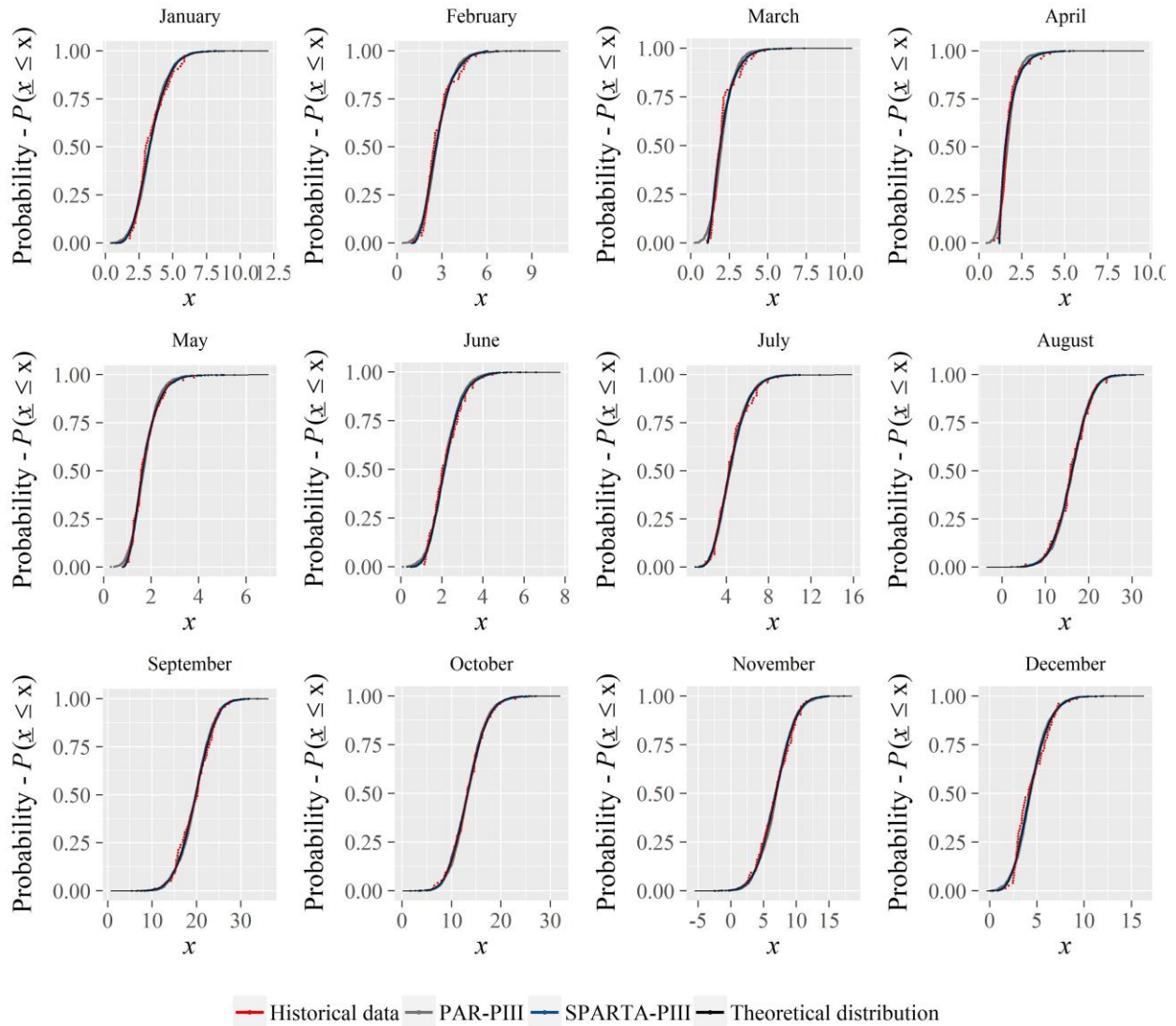


567

568

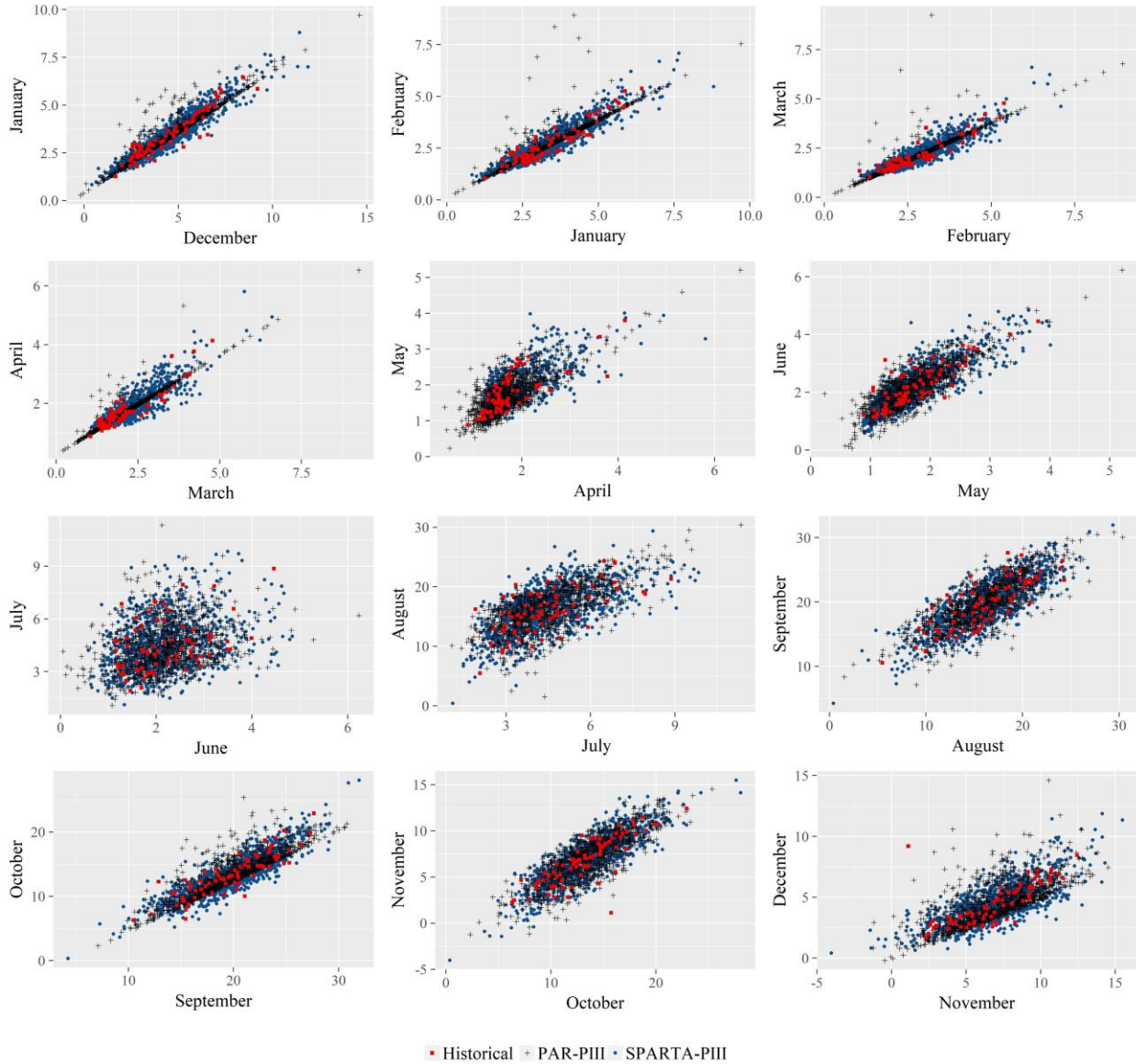
569

Figure 2: Comparison of key statistics (μ , σ , C_s and ρ_1) between historical and simulated flow data of Nile River (PAR and SPARTA).



570

571 **Figure 3:** Comparison between simulated flow data (10^9 m^3), through PAR-PIII and SPARTA-PIII,
 572 empirical and theoretical cumulative distribution functions (Weibull plotting position). Simulated negative
 573 values are also included to avoid the distortion of the established CDFs.



574

575 **Figure 4:** Month-to-month scatter plots of historical and simulated flow data (10^9 m^3), through PAR-PIII
 576 and SPARTA-PIII. Simulated negative values are also included to avoid the distortion of the established
 577 dependence patterns.

578 **5.2 Toy simulation with seasonally-varying distribution models**

579 The second case study involves the simulation of a hypothetical seasonal process, $\underline{x}_{s,t}$, with
 580 different marginal distribution per season (for convenience, 12 seasons are considered). The target
 581 distribution models and the associated parameters across seasons are given in Table 1. In addition,
 582 we assume the target lag-1 (i.e., season-to-season) correlation coefficients equal to $\boldsymbol{\rho} =$
 583 $[\rho_{12,1}, \rho_{1,2}, \dots, \rho_{s,s-1}, \dots, \rho_{10,11}, \rho_{11,12}] = [0.93, 0.90, 0.76, 0.84, 0.32, 0.67, 0.80, 0.88, 0.83, 0.74,$
 584 $0.94, 0.93]$. Using SPARTA we generated $1\,000 \times 12 = 12\,000$ synthetic values of $\underline{x}_{s,t}$ and
 585 compared their statistical characteristics against the target ones. We remark that in contrast to the
 586 previous case study, we do not compare against another linear stochastic model (e.g., PAR-PIII),

587 given that we have specified different statistical distributions across seasons, which cannot be
 588 represented by such models.

589 The theoretical and simulated values of the key statistical characteristics of the modelled process
 590 are illustrated in Table 2. The former were calculated through the corresponding theoretical
 591 equations of each distribution. As shown, SPARTA is very efficient, since it reproduces all key
 592 statistics, including the kurtosis coefficient, C_k . Furthermore, SPARTA preserves the parameters
 593 of the target marginal distributions (Table 1, upper part), which are estimated through the MLE
 594 method. Actually, as shown in Table 1 (lower part), there is close agreement between the target
 595 and simulated parameter values for all seasons. This is also visually confirmed by plotting the
 596 associated CDFs (Figure 5), as the discrepancies between the theoretical and empirical distributions
 597 are almost indistinguishable. It is noted that the distributions employed for season 5 and 10 allowed
 598 the generation of negative values since we assigned to the former a Gaussian one (which is
 599 unbounded) and in the latter a Pearson Type-III with location parameter $c = -50$ which coincides
 600 with its theoretical lower bound (given that $b > 0$). All other distributions are defined in the positive
 601 real axis, hence they don't allow the generation of negative values.

602 Furthermore, the stochastic structure of the hypothetical process, by means of season-to-season
 603 correlations, ρ_1 , is reproduced, despite the fact that it exhibits significant variability, also
 604 comprising some very high ρ_1 values. In order to shed further light on the seasonal dependence
 605 patterns, we provide scatter plots combined with histograms for four adjacent seasons, from which
 606 it becomes evident that SPARTA can reproduce a plethora of marginal distributions and
 607 simultaneously account for dependence patterns of different complexity (Figure 6).

608 **Table 1:** Theoretical distributions and associated parameters of hypothetical process across seasons, as well
 609 as MLE estimation of simulated data.

Season	1	2	3	4	5	6	7	8	9	10	11	12
Distribution/ Parameters	PIII	Exp	Gam	Norm	LoNo	Wei	Wei	LoNo	Exp	PIII	Wei	Gam
	Theoretical Values											
a	1.7	0.015	10	85	5	4.5	6	6	0.003	11	3	9
b	10	-	0.15	30	0.3	680	820	0.25	-	19	155	0.2
c	40	-	-	-	-	-	-	-	-	-50	-	-
	Simulated Values											
a	1.72	0.015	10.01	85	5	4.47	5.99	6	0.003	9.12	2.97	9.09
b	9.88	-	0.15	29.98	0.29	680.03	819.91	0.25	-	20.98	154.90	0.20
c	39.94	-	-	-	-	-	-	-	-	-51.39	-	-

*Distribution abbreviations: PIII: Pearson type-III ($a =$ shape, $b =$ scale, $c =$ location), Exp: Exponential ($a =$ rate), Gam: Gamma ($a =$ shape, $b =$ rate), Norm: Normal ($a =$ mean, $b =$ st. dev.), LoNo: Log-Normal ($a =$ log mean, $b =$ log st. dev.), Wei: Weibull ($a =$ shape, $b =$ scale).

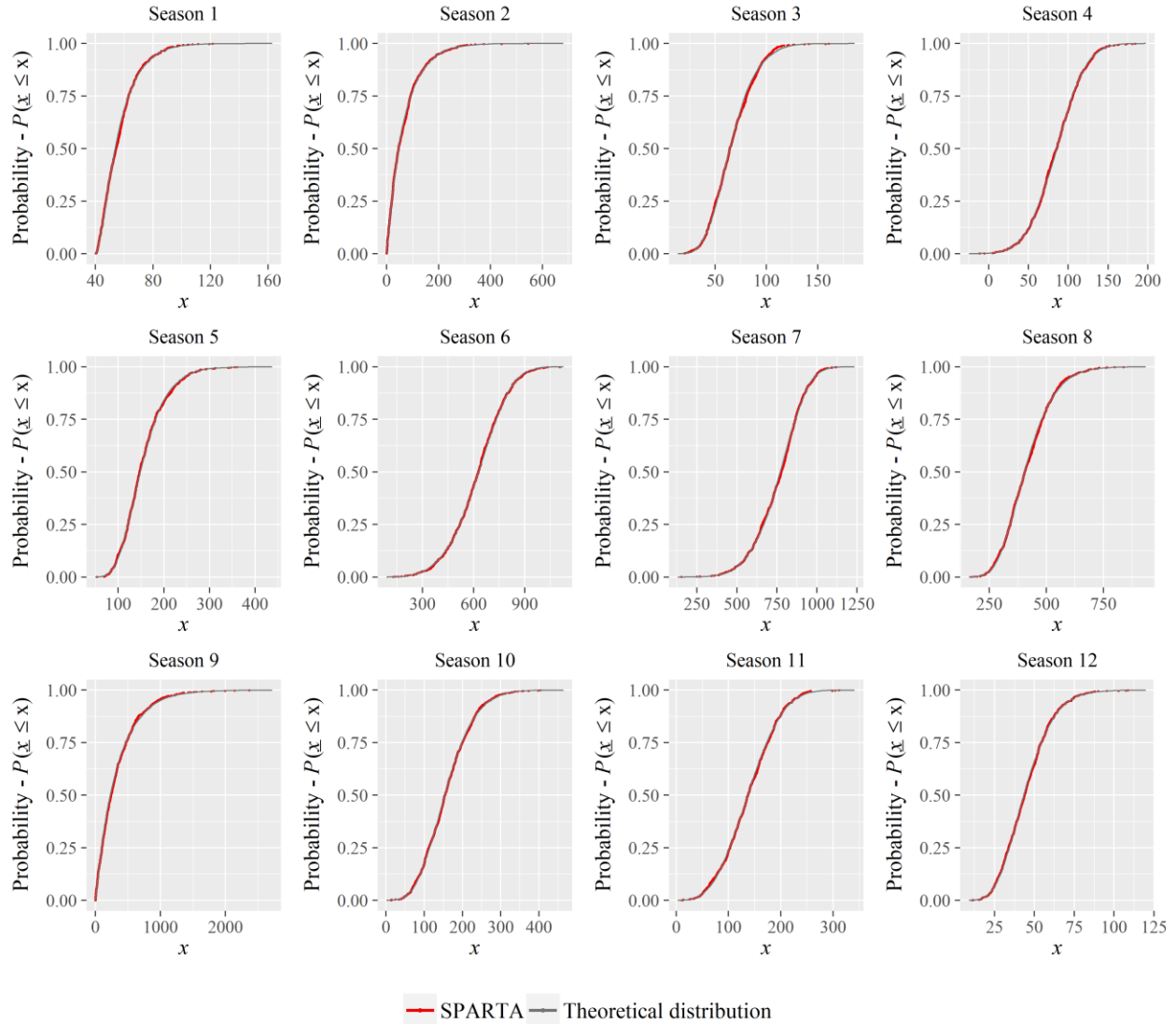
610

611

Table 2: Simulated and theoretical values of key statistical characteristics of hypothetical process.

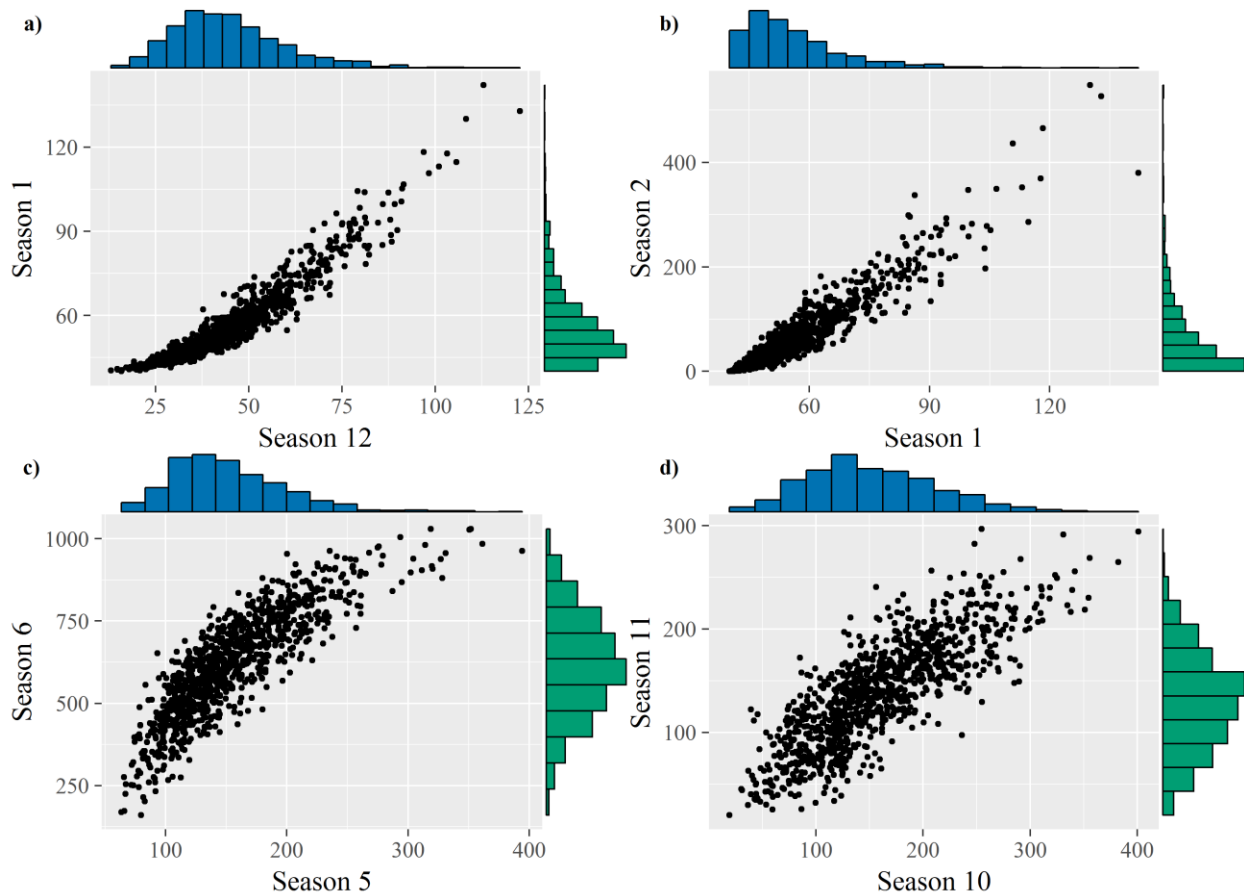
Season/ Statistic	1	2	3	4	5	6	7	8	9	10	11	12
μ (Theor.)	57.00	66.67	66.67	85.00	155.24	620.55	760.72	416.23	333.33	159.00	138.41	45.00
μ (Sim.)	56.99	66.56	66.67	85.00	155.27	620.53	760.81	416.34	333.23	159.01	138.37	45.00
σ (Theor.)	13.03	66.67	21.08	30.00	47.64	156.45	147.40	105.70	333.33	63.02	50.30	15.00
σ (Sim.)	13.26	66.96	21.20	30.00	48.18	156.02	147.18	107.38	335.69	63.80	50.24	15.14
C_s (Theor.)	1.53	2.00	0.63	0.00	0.97	-0.17	-0.37	0.88	2.00	0.60	0.16	0.66
C_s (Sim.)	1.75	1.98	0.72	-0.04	1.09	-0.13	-0.39	0.94	1.89	0.75	0.27	0.82
C_k (Theor.)	6.53	9.00	3.60	3.00	4.99	2.80	3.03	4.06	9.00	3.54	2.72	3.66
C_k (Sim.)	7.62	8.01	3.84	2.98	5.20	2.88	3.20	4.46	7.32	3.85	3.05	4.20
ρ_1 (Theor.)	0.93	0.90	0.76	0.84	0.32	0.67	0.80	0.88	0.83	0.74	0.94	0.93
ρ_1 (Sim.)	0.94	0.90	0.76	0.82	0.31	0.66	0.80	0.87	0.85	0.77	0.95	0.93
$\tilde{\rho}_1$ (Equiv.)	0.95	0.91	0.80	0.85	0.32	0.70	0.80	0.90	0.88	0.78	0.96	0.94

*Table abbreviations: Theor: Theoretical value, Sim: Simulated value, Equiv: Equivalent value.



613

614 **Figure 5:** Comparison between simulated (SPARTA) and theoretical cumulative distribution functions
 615 (Weibull plotting position) of hypothetical process. Simulated negative values (season 5 and 10) are also
 616 included to avoid the distortion of the established CDFs.



617
 618 **Figure 6:** Scatter plots with histograms for a) season 12 vs. 1 b) season 1 vs. 2, c) season 5 vs. 6, and d)
 619 season 10 vs. 11.

620 5.3 Multivariate simulation

621 The third case study involves the simultaneous generation of monthly runoff and rainfall data at
 622 two major reservoirs of the water supply system of Athens, i.e., Evinos and Mornos (details about
 623 the system are provided by [Koutsoyiannis et al, \[2003\]](#)). The historical data cover a 29-year period
 624 (Oct/1979 – Sep/2008), which is marginally adequate for estimating up to third moment statistics
 625 with acceptable accuracy. For convenience, herein we will refer to Evinos runoff and rainfall as
 626 “sites” A and B, respectively, and to Mornos runoff and rainfall as “sites” C and D, respectively
 627 (here term “site” denotes a specific hydrological process at a specific location).

628 In this problem we employed the multivariate version of SPARTA and compared against the
 629 contemporaneous PAR(1) model with Pearson type-III white noise, again, referred as PAR-PIII
 630 model (Appendix A). Similarly to case study 1, in the context of specifying the underlying
 631 marginal distributions of SPARTA, and in order to ensure fair comparisons, we decided fitting the
 632 Pearson type-III model at all sites and for all months, and estimating its parameters via the method
 633 of moments. Under this premise, the generating scheme will be next referred to as SPARTA-PIII.
 634 Although we remark, that in an operational, “real-world study” one could take advantage of
 635 SPARTA model flexibility and select appropriate distributions models that are positively bounded,
 636 thus directly surpass the problem of negative values generation (see also the previous sections).

637 The performance of both models was assessed in a monthly basis, by contrasting the statistical
638 characteristics of historical data that should be theoretically preserved by the corresponding
639 generating schemes (i.e., monthly means, standard deviations, and skewness coefficients, lag-1
640 correlations across months, and zero-lag cross-correlations between all sites) against the simulated
641 ones.

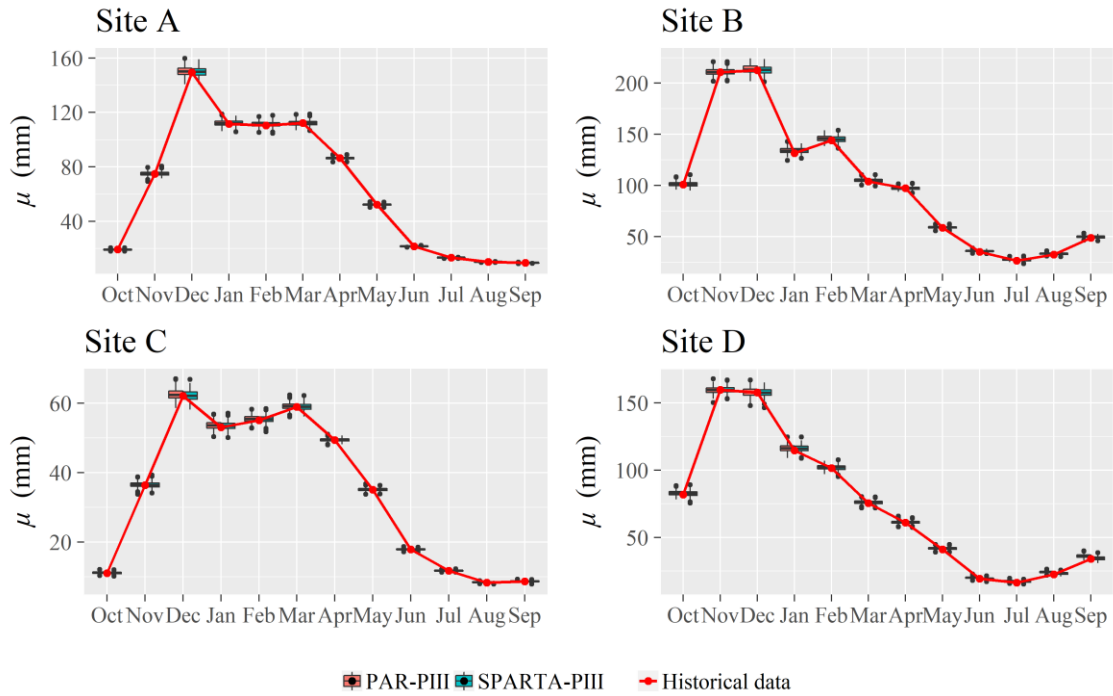
642 It is well-known that while the theoretical equations of any stochastic model are built in order to
643 explicitly reproduce a specific set of statistical characteristics, this preservation is only ensured for
644 very long (theoretically infinite) simulation horizons [Efstratiadis *et al.*, 2014]. If we consider
645 relatively small horizons and repeat the simulation many times, the smaller the length of the
646 synthetic sample, the larger is expected to be the variability of the simulated against the theoretical
647 values of these characteristics. In this context, the stochastic model that ensures the minimum
648 variability will be recognized as the most robust, since its performance will be the less sensitive
649 against the simulation length. In this context, we employed two experiments, the first one by
650 employing a single simulation of 500 000 years length, and the second one by running each model
651 500 times, to obtain independent synthetic samples of 1 000 years length. This Monte Carlo
652 approach allowed for evaluating the uncertainty of the simulated statistical characteristics (after
653 truncation of negative values to zero), which is depicted by means of box-plots (Figure 7 to Figure
654 11).

655 As shown in supplementary material (SM; Figure S1-S5), the estimated statistical characteristics
656 from the large (i.e., 500 000 years) synthetic sample perfectly agree with the historical ones, thus
657 confirming the solid theoretical background of SPARTA-PIII. As expected, PAR-PIII also ensures
658 perfect fitting of the simulated to the observed statistics, except for skewness, which are slightly
659 underestimated. Probably, this systematic deviation is due to the simplified method employed for
660 covariance matrix decompositions (namely, the Cholesky technique), as already mentioned in
661 section 3.4.1.

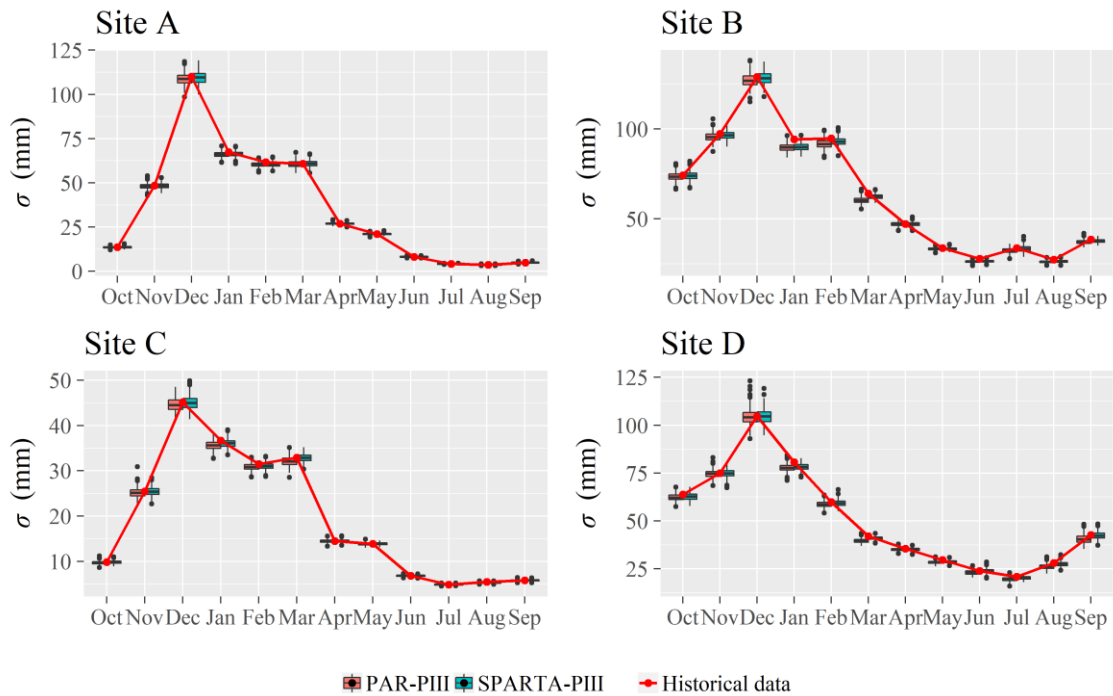
662 The superiority of SPARTA-PIII against PAR-PIII is further revealed when evaluating the fitting
663 of synthetic data to the theoretical distribution that has been adopted in this simulation experiment,
664 i.e., Pearson type III. The latter is mathematically defined through Eq. (12) comprising three
665 parameters, i.e., shape, a , scale, b , and location, c , which have been estimated for each site and
666 each month with the method of moments (SM, Table S1). It is clearly shown that the estimated
667 parameter values originated by SPARTA-PIII are very close to the theoretical ones, thus the
668 desirable distributions are accurately reproduced. On the other hand, there are several cases where
669 the PAR-derived parameters, and consequently the derived distributions, oscillate significantly
670 from the theoretical model. This becomes even more evident when expressing these deviations in
671 terms of root mean square error, per site and parameter. As shown in SM, Table S2, this error is
672 up to three times larger than the error induced by SPARTA-PIII.

673 With respect to the second (i.e., Monte Carlo) experiment, from Figure 7 and Figure 8 it is shown
674 that both SPARTA-PIII and PAR-PIII are able to reproduce the observed monthly means and
675 standard deviations, respectively, since their variability is generally low across all sites and
676 seasons. Regarding the reproduction of monthly coefficients of skewness (Figure 9), it seems that
677 SPARTA-PIII slightly outperforms PAR-PIII in terms of statistical uncertainty, as indicated by
678 the narrower box-plots that are provided in several cases (e.g., October, March, August and
679 September for site A, October, November and March for site B, November, December and March
680 for site C, and March, August and September for site D). Finally, in terms of lag-1 month-to-month

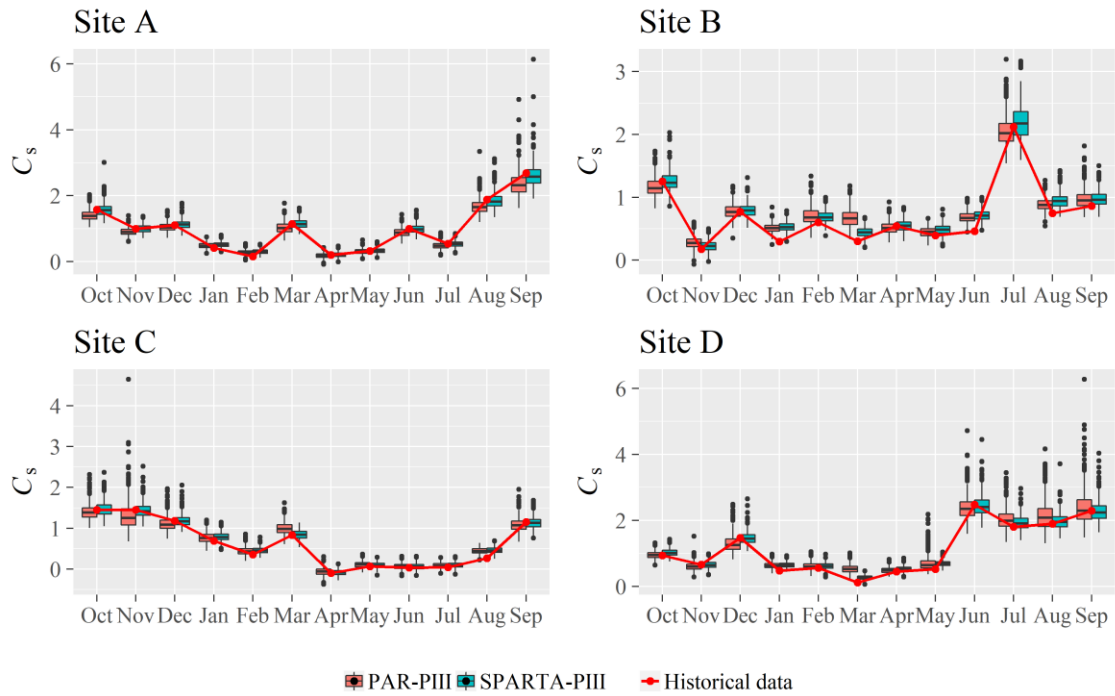
681 and lag-0 cross-correlations, both schemes ensure robustness, as illustrated in Figure 10 and Figure
 682 11, respectively.



683
 684 **Figure 7:** Comparison of monthly mean values, μ , of historical and synthetic data.

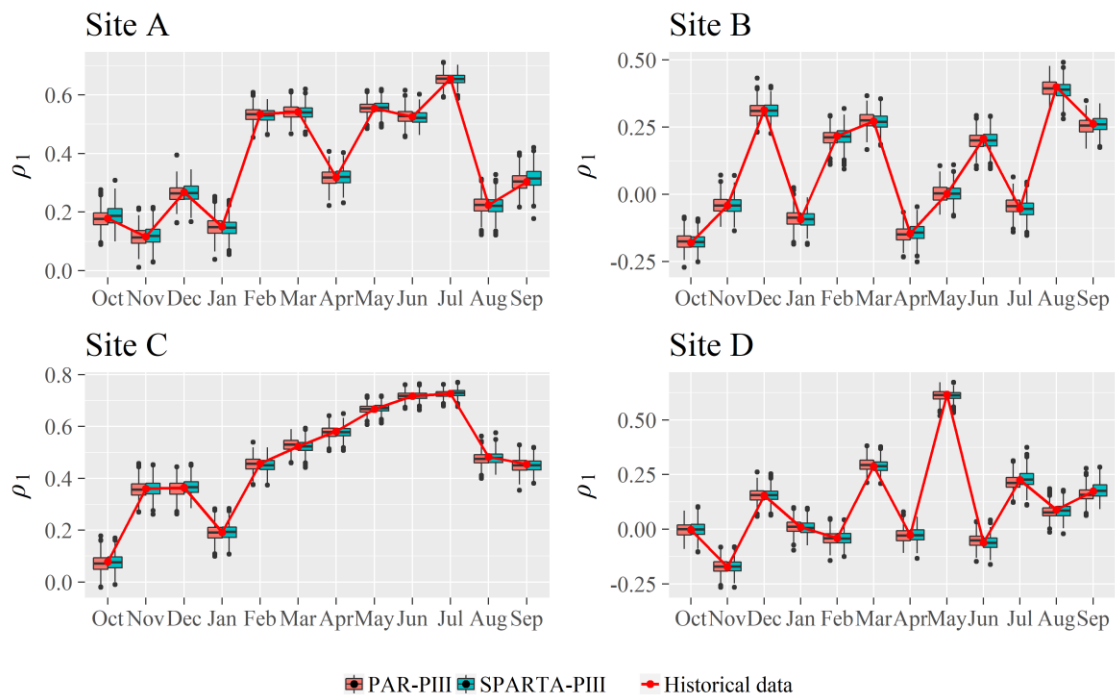


685
 686 **Figure 8:** Comparison of monthly standard deviation values, σ , of historical and synthetic data.



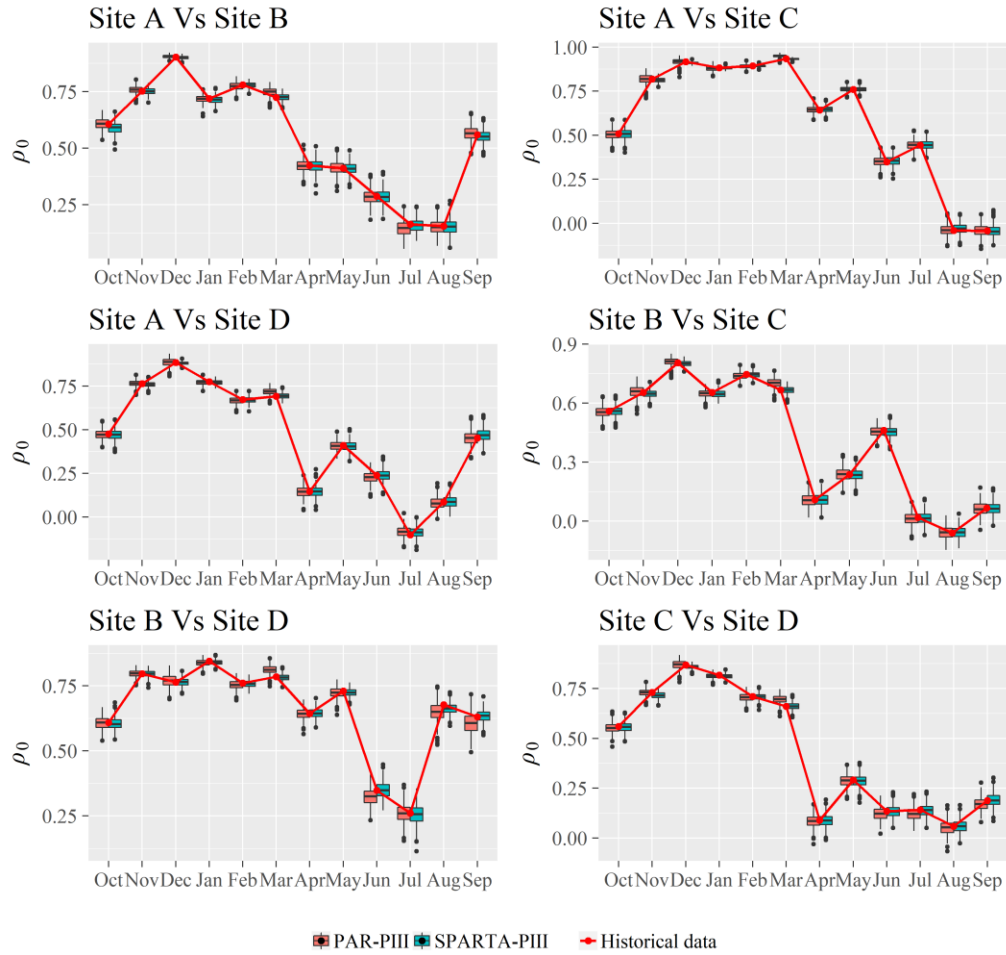
687

688 **Figure 9:** Comparison of monthly skewness coefficients, C_s , of historical and synthetic data.

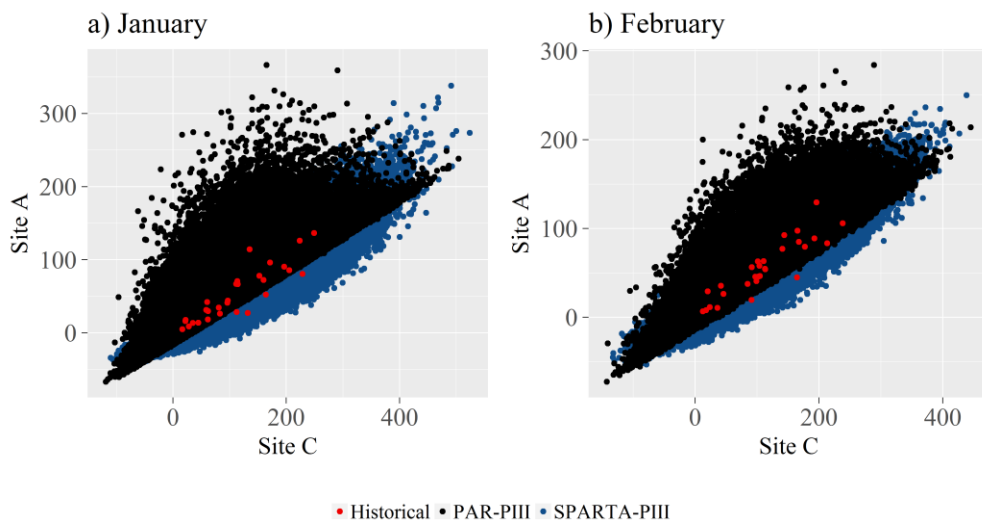


689

690 **Figure 10:** Comparison of month-to-month lag-1 correlations, ρ_1 , of historical and synthetic data.



691
 692 **Figure 11:** Comparison of monthly lag-0 cross-correlations, ρ_0 , between sites of historical and synthetic
 693 data.



694
 695 **Figure 12:** Scatter plots of 500 000 synthetic data for sites A and C, representing monthly runoff (mm)
 696 processes at Evinos and Mornos reservoirs, respectively, for (a) January and (b) February. Simulated
 697 negative values are also included to avoid the distortion of the established dependence patterns.

698 As already highlighted, a great advantage of SPARTA over linear stochastic schemes, such as
699 PAR-PIII, is its ability to reproduce realistic dependence patterns, in compliance to the observed
700 ones. This is also empirically confirmed in the current case study, which aims to reproduce both
701 temporal and spatial dependencies (i.e., dependencies between different processes). A
702 characteristic example is given in Figure 12, illustrating the scatter plots of historical and simulated
703 runoff values of at Evinos (site A) and Mornos (site C), for months January and February, from
704 the long-term experiment (i.e., 500 000 years). It becomes now even more evident that the
705 SPARTA-PIII generation scheme provides reasonably-distributed data, while the synthetic data
706 by PAR-PIII are again bounded within a specific range, which is far from truthful and does not
707 capture the full extent of the observed scatter (notice the incompatibility between the synthetic
708 series of PAR-PIII and the historical data in Figure 12).

709 **6 Discussion**

710 As briefly discussed in the introduction, and demonstrated through three case studies, the need for
711 generic simulation schemes that allow producing synthetic data from multiple distributions
712 primarily originates from the fact that the statistical behavior of many of hydroclimatic processes
713 is not satisfactory captured by classical stochastic models. Such models cannot reproduce
714 significant statistical aspects of the simulated processes (e.g., maxima and minima, associated with
715 the tails of the distribution), although the “essential”, low-order statistical characteristics of the
716 parent data may be well-preserved.

717 However, to our opinion, the overall question is not just a technical issue, i.e., providing better
718 stochastic models, but, in a more general context, revisiting the “essentials” of synthetic data. In
719 particular, we suggest moving from the preservation of a specific set of statistical characteristics,
720 which are exclusively inferred from the observed data, to the preservation of *a priori* specified
721 theoretical distributions that are hypothesized to be consistent with the anticipated stochastic
722 behavior of the underlying processes.

723 We recognize that the assignment of a specific distribution model for each modelled process is not
724 a straightforward task, since the true distribution will always be unknown. Obviously, for a given
725 data sample one can fit a plethora of distributions, combined with different parameter estimation
726 procedures (e.g., classical moments, L-moments, maximum likelihood), and use typical statistical
727 tests to assess the “optimal” scheme. Even for a given set of statistical characteristics, multiple
728 distributions may be used. However, theoretical reasons, such as the central limit theorem and the
729 principle of maximum entropy, may induce the selection of a different distribution, even when the
730 latter is not so favored by the data [e.g., Koutsoyiannis, 2005; Papalexiou and Koutsoyiannis,
731 2012]. In any case, particularly when the historical samples are short or not so much reliable, the
732 selection of the most suitable distribution may be supported by hydrological evidence. For
733 instance, one may take advantage of the statistical behavior of the underlying processes in the
734 broader area, as validated by large-scale regional studies [e.g., Blum et al., 2017].

735 A final remark involves the treatment of historical data themselves. Actually, the observed
736 statistics are subject to biases and uncertainties induced by their estimation from relatively short
737 records (e.g., unreasonably high skewness values, due to outliers). Several times, the use of data
738 as the sole means for extracting the statistical characteristics of the process of interest may also
739 result to severe inconsistencies, such as negative autocorrelations that do not have physical
740 meaning in hydrology [Koutsoyiannis, 2000]. Particularly, in the latter case it may be wise to
741 follow the paradigm of the aforementioned author and fit a theoretical model on the empirically

742 derived autocorrelation coefficients. Nevertheless, it may be preferable to assign, even manually,
743 realistic values to the “suspicious” parameters rather than leave the model employing erroneous
744 values. Moreover, due to changing environmental and hydroclimatic conditions, the statistical
745 information contained in historical data may not be fully representative of the “projected” future
746 conditions. In this context, aiming to explore the effects of change, several researchers suggest
747 perturbing the values of the statistical characteristics to be reproduced within synthetic data [e.g.,
748 *Nazemi et al., 2013; Borgomeo et al., 2015*], which obviously imply employing parameters
749 different than the data-driven ones. Nevertheless, wherever it is necessary to manually assign target
750 input values, these have to be checked against both physical consistency and hydrological
751 evidence. In this vein, we remark that NDM-based models (e.g., ARTA, VARTA and SPARTA)
752 are able to synthesize data from any distribution hence allowing their straightforward use in such
753 studies. This can be easily accomplished by changing the parameters of the distribution functions
754 (even the distribution functions themselves) or the correlation structure of the process and
755 subsequently investigate the effects of such changes to the system under study.

756 **7 Conclusions**

757 This work presents a novel approach, termed SPARTA, for the explicit stochastic simulation of
758 univariate and multivariate cyclostationary (i.e., periodic) processes with arbitrary marginal
759 distributions. SPARTA uses an auxiliary Gaussian PAR process with properly identified
760 parameters, such as after its mapping to the actual domain through the ICDFs, it results to a process
761 with the target correlation structure and *a priori* specified marginal distributions. Since the
762 temporal and spatial dependencies are typically expressed by means of Pearson correlation
763 coefficients, we focus on the identification of equivalent correlation coefficients of the auxiliary
764 processes to be used in the Gaussian domain, in order to attain the target correlations in the actual
765 domain. In this context, we use the Nataf joint distribution model, originated from statistical
766 sciences for the generation of correlated random variables with prescribed distributions. Based on
767 the theoretical background of NDM, we have developed a simple, yet efficient Monte-Carlo based
768 approach that allows for identifying the equivalent correlation coefficients, $\tilde{\rho}$, with low
769 computational effort.

770 Despite the obvious benefit of simulating processes with any marginal distributions, the proposed
771 approach is also flexible in implementing any distribution fitting method, offered by recent
772 advances in statistical sciences. This flexibility also offers the capability of explicitly ensuring the
773 generation of non-negative values within simulations, through selecting appropriate distributions
774 that are positively bounded. This very important potential, which is not offered by most of known
775 stochastic schemes used in hydrology, is attributed to the use of the ICDF; if the latter is positively
776 bounded, the generated values will be by definition non-negative.

777 The advantages of SPARTA in practice, i.e., in the context of generating monthly synthetic data,
778 have been illustrated through three stochastic simulation studies, emphasizing different aspects of
779 the proposed methodology. Furthermore, in two out of three studies, SPARTA has been contrasted
780 to the well-established linear stochastic model PAR-PIII, i.e., PAR(1) with Pearson type-III white
781 noise. The major outcomes of our analyses are:

- 782 • Both models reproduced almost perfectly the essential statistical characteristics of the
783 simulated processes up to second order (means, standard deviations, lag-1 month-to-month
784 correlations (i.e., autocorrelations), zero-lag cross-correlations);

- 785
- 786
- 787
- 788
- 789
- 790
- 791
- 792
- 793
- 794
- SPARTA was also able to preserve with high accuracy the third order statistics, expressed in terms of skewness coefficients, while in several cases PAR-PIII provided quite underestimated skewness, which varied significantly across independently generated synthetic samples;
 - SPARTA was able not only to preserve the theoretical statistical characteristics of the observed data but also the parameters of the prescribed marginal distributions, which is in fact the primary goal of simulation (see discussion);
 - SPARTA produced dependence structures in time and space that are in agreement with the observed patterns, while, in some cases, PAR-PIII provided rather irregular scatter patterns that were fragmented out of the observed ranges.

795 To this end, it is argued, that SPARTA is a convenient way to simulate cyclostationary
796 processes, either univariate or multivariate, yet it should not be regarded as a *panacea* for
797 all kind of simulation problems, since it inherits the characteristics of the auxiliary process
798 from the periodic autoregressive family. In this context, it cannot preserve the statistical
799 characteristics at aggregated time scales, e.g., annual, including long-range dependence
800 (Hurst phenomenon). For this reason, future research involves the integration of SPARTA within
801 a multi-scale stochastic framework, allowing us to reproduce the desirable distribution and
802 desirable correlation structures at multiple time scales, and also reproduce the peculiarities of
803 different scales. As shown in the literature, an effective and efficient way to address this is through
804 disaggregation techniques. For instance, the hybrid Monte Carlo procedure by *Koutsoyiannis and*
805 *Manetas* [1996], which has been successfully implemented within advanced simulation schemes
806 [e.g., *Efstratiadis et al.*, 2014; *Kossieris et al.*, 2016], can be easily aligned with SPARTA to ensure
807 statistical consistency across scales.

808 As a concluding remark, and following the discussion of section 6, the authors would like to
809 highlight the fact that the blind use of stochastic models, with overconfidence on historical data,
810 may create a distorted “reality”, thus feeding operational hydrological and water management
811 studies with inconsistent synthetic inputs. In this vein, we recommend to turn our efforts into the
812 selection of the suitable distribution model, as well as the careful assessment of the sample
813 statistics, with emphasis to high order moments and correlations that are prone to uncertainties.
814 Therefore, the flexibility of the proposed approach contributes towards the establishment of a new
815 paradigm in hydrological stochastics.

816 Acknowledgments

817 The authors would like to thank the Associate Editor and the three anonymous reviewers, for their
818 constructive comments, suggestions and critique, which helped providing a much improved
819 manuscript. **Data availability:** Nile streamflow data at Aswan dam was retrieved from an external
820 source (<http://www.stats.uwo.ca/faculty/mcleod/epubs/mhsets/>). The historical dataset (runoff and
821 rainfall) of Evinos and Mornos reservoirs is available at: <http://www.itia.ntua.gr/1746/>. **Code**
822 **availability:** The developed R scripts and functions that implement the SPARTA model are
823 available upon request to the authors.

824 **References**

- 825 Ailliot, P., D. Allard, V. Monbet, and P. Naveau (2015), Stochastic weather generators: an
826 overview of weather type models, *J. la Société Française Stat.*, 156(1), 101–113.
- 827 Akaike, H. (1974), A new look at the statistical model identification, *IEEE Trans. Automat. Contr.*,
828 19(6), 716–723.
- 829 Biller, B., and B. L. Nelson (2003), Modeling and generating multivariate time-series input
830 processes using a vector autoregressive technique, *ACM Trans. Model. Comput. Simul.*,
831 13(3), 211–237, doi:10.1145/937332.937333.
- 832 Blum, A. G., S. A. Archfield, and R. M. Vogel (2017), On the probability distribution of daily
833 streamflow in the United States, *Hydrol. Earth Syst. Sci.*, 21(6), 3093–3103,
834 doi:10.5194/hess-21-3093-2017.
- 835 Borgomeo, E., C. L. Farmer, and J. W. Hall (2015), Numerical rivers: A synthetic streamflow
836 generator for water resources vulnerability assessments, *Water Resour. Res.*, 51(7), 5382–
837 5405, doi:10.1002/2014WR016827.
- 838 Bras, R. L., and I. Rodríguez-Iturbe (1985), *Random functions and hydrology*, Addison-Wesley,
839 Reading, Mass.
- 840 Cario, M. C. (1996), Modeling and simulating time series input processes with ARTAFACETS and
841 ARTAGEN, in *Proceedings of the 28th conference on Winter simulation - WSC '96*, edited
842 by J. M. Charnes, D. J. Morrice, D. T. Brunner, and J. J. Swain, pp. 207–213, ACM Press,
843 New York, New York, USA.
- 844 Cario, M. C., and B. L. Nelson (1996), Autoregressive to anything: Time-series input processes
845 for simulation, *Oper. Res. Lett.*, 19(2), 51–58, doi:10.1016/0167-6377(96)00017-X.
- 846 Cario, M. C., and B. L. Nelson (1997), Modeling and generating random vectors with arbitrary
847 marginal distributions and correlation matrix, *Ind. Eng.*, 1–19.
- 848 Chen, H. (2001), Initialization for NORTA: Generation of Random Vectors with Specified
849 Marginals and Correlations, *INFORMS J. Comput.*, 13(4), 312–331,
850 doi:10.1287/ijoc.13.4.312.9736.
- 851 Chen, L., V. P. Singh, S. Guo, J. Zhou, and J. Zhang (2015), Copula-based method for multisite
852 monthly and daily streamflow simulation, *J. Hydrol.*, 528, 369–384,
853 doi:10.1016/j.jhydrol.2015.05.018.
- 854 Efstratiadis, A., Y. G. Dialynas, S. Kozanis, and D. Koutsoyiannis (2014), A multivariate
855 stochastic model for the generation of synthetic time series at multiple time scales
856 reproducing long-term persistence, *Environ. Model. Softw.*, 62(July), 139–152,
857 doi:10.1016/j.envsoft.2014.08.017.
- 858 Embrechts, P., A. J. McNeil, and D. Straumann (1999), Correlation and Dependence in Risk
859 Management: Properties and Pitfalls, in *Risk Management*, edited by M. A. H. Dempster, pp.
860 176–223, Cambridge University Press, Cambridge.
- 861 Fernandez, B., and J. D. Salas (1986), Periodic Gamma Autoregressive Processes for Operational
862 Hydrology, *Water Resour. Res.*, 22(10), 1385–1396.
- 863 Fiering, B., and B. Jackson (1971), *Synthetic Streamflows*, Water Resources Monograph,
864 American Geophysical Union, Washington, D. C.

- 865 Fiering, M. B. (1964), Multivariate technique for synthetic hydrology, *J. Hydraul. Div.*, 90(5), 43–
866 60.
- 867 Hao, Z., and V. P. Singh (2013), Modeling multisite streamflow dependence with maximum
868 entropy copula, *Water Resour. Res.*, 49(10), 7139–7143, doi:10.1002/wrcr.20523.
- 869 Herman, J. D., H. B. Zeff, J. R. Lamontagne, P. M. Reed, and G. W. Characklis (2016), Synthetic
870 Drought Scenario Generation to Support Bottom-Up Water Supply Vulnerability
871 Assessments, *J. Water Resour. Plan. Manag.*, 4016050, doi:10.1061/(ASCE)WR.1943-
872 5452.0000701.
- 873 Higham, N. J. (2002), Computing the nearest correlation matrix--a problem from finance, *IMA J.*
874 *Numer. Anal.*, 22(3), 329–343, doi:10.1093/imanum/22.3.329.
- 875 Hipel, K. W., and A. I. McLeod (1994), *Time series modelling of water resources and*
876 *environmental systems*, Elsevier.
- 877 Johnson, M. E. (1987), *Multivariate Statistical Simulation*, John Wiley, New York, NY, USA.
- 878 Kelly, K. S., and R. Krzysztofowicz (1997), A bivariate meta-Gaussian density for use in
879 hydrology, *Stoch. Hydrol. Hydraul.*, 11(1), 17–31, doi:10.1007/BF02428423.
- 880 Kirsch, B. R., G. W. Characklis, and H. B. Zeff (2013), Evaluating the Impact of Alternative
881 Hydro-Climate Scenarios on Transfer Agreements: Practical Improvement for Generating
882 Synthetic Streamflows, *J. Water Resour. Plan. Manag.*, 139(4), 396–406,
883 doi:10.1061/(ASCE)WR.1943-5452.0000287.
- 884 Klemeš, V., and L. Borůvka (1974), Simulation of Gamma-Distributed First-Order Markov Chain,
885 *Water Resour. Res.*, 10(1), 87–91, doi:10.1029/WR010i001p00087.
- 886 Kossieris, P., C. Makropoulos, C. Onof, and D. Koutsoyiannis (2016), A rainfall disaggregation
887 scheme for sub-hourly time scales: Coupling a Bartlett-Lewis based model with adjusting
888 procedures, *J. Hydrol.*, doi:10.1016/j.jhydrol.2016.07.015.
- 889 Koutsoyiannis, D. (1999), Optimal decomposition of covariance matrices for multivariate
890 stochastic models in hydrology, *Water Resour. Res.*, 35(4), 1219–1229,
891 doi:10.1029/1998WR900093.
- 892 Koutsoyiannis, D. (2000), A generalized mathematical framework for stochastic simulation and
893 forecast of hydrologic time series, *Water Resour. Res.*, 36(6), 1519–1533,
894 doi:10.1029/2000WR900044.
- 895 Koutsoyiannis, D. (2005), Uncertainty, entropy, scaling and hydrological stochastics. 1. Marginal
896 distributional properties of hydrological processes and state scaling / Incertitude, entropie,
897 effet d'échelle et propriétés stochastiques hydrologiques. 1. Propriétés distributionnel,
898 *Hydrol. Sci. J.*, 50(3), 381–404, doi:10.1623/hysj.50.3.381.65031.
- 899 Koutsoyiannis, D., and A. Manetas (1996), Simple disaggregation by accurate adjusting
900 procedures, *Water Resour. Res.*, 32(7), 2105–2117, doi:10.1029/96WR00488.
- 901 Koutsoyiannis, D., G. Karavokiros, A. Efstratiadis, N. Mamassis, A. Koukouvinos, and A.
902 Christofides (2003), A decision support system for the management of the water resource
903 system of Athens, *Phys. Chem. Earth, Parts A/B/C*, 28(14–15), 599–609, doi:10.1016/S1474-
904 7065(03)00106-2.
- 905 Koutsoyiannis, D., H. Yao, and A. Georgakakos (2008), Medium-range flow prediction for the

- 906 Nile: a comparison of stochastic and deterministic methods, *Hydrol. Sci. Journal-Journal Des*
907 *Sci. Hydrol.*, 53(1), 142–164, doi:10.1623/hysj.53.1.142.
- 908 Lall, U., and A. Sharma (1996), A nearest neighbor bootstrap for resampling hydrologic time
909 series, *Water Resour. Res.*, 32(3), 679–693, doi:10.1029/95WR02966.
- 910 Lawrance, A. J., and P. A. W. Lewis (1981), A new autoregressive time series model in exponential
911 variables (NEAR (1)), *Adv. Appl. Probab.*, 13(4), 826–845.
- 912 Lebrun, R., and A. Dutfoy (2009), An innovating analysis of the Nataf transformation from the
913 copula viewpoint, *Probabilistic Eng. Mech.*, 24(3), 312–320,
914 doi:10.1016/j.probengmech.2008.08.001.
- 915 Lee, T. (2017), Multisite stochastic simulation of daily precipitation from copula modeling with a
916 gamma marginal distribution, *Theor. Appl. Climatol.*, doi:10.1007/s00704-017-2147-0.
- 917 Li, S. T., and J. L. Hammond (1975), Generation of Pseudorandom Numbers with Specified
918 Univariate Distributions and Correlation Coefficients, *IEEE Trans. Syst. Man. Cybern.*, SMC-
919 5(5), 557–561, doi:10.1109/TSMC.1975.5408380.
- 920 Liu, P. L., and A. Der Kiureghian (1986), Multivariate distribution models with prescribed
921 marginals and covariances, *Probabilistic Eng. Mech.*, 1(2), 105–112, doi:10.1016/0266-
922 8920(86)90033-0.
- 923 Maass, A., M. M. Hufschmidt, R. Dorfman, H. A. Thomas, S. A. Marglin, G. M. Fair, B. T. Bower,
924 W. W. Reedy, D. F. Manzer, and M. P. Barnett (1962), *Design of water-resource systems*,
925 Cambridge: Harvard University Press.
- 926 Marković, Đ., J. Plavšić, N. Ilich, and S. Ilić (2015), Non-parametric Stochastic Generation of
927 Streamflow Series at Multiple Locations, *Water Resour. Manag.*, 29(13), 4787–4801,
928 doi:10.1007/s11269-015-1090-z.
- 929 Matalas, N. C. (1967), Mathematical assessment of synthetic hydrology, *Water Resour. Res.*, 3(4),
930 937–945, doi:10.1029/WR003i004p00937.
- 931 Matalas, N. C. (1975), Developments in stochastic hydrology, *Rev. Geophys.*, 13(3), 67,
932 doi:10.1029/RG013i003p00067.
- 933 Matalas, N. C., and J. R. Wallis (1976), *Generation of synthetic flow sequences, Systems Approach*
934 *to Water Management*, edited by A. K. Biswas, McGraw-Hill, New York, New York.
- 935 Mehrotra, R., R. Srikanthan, and A. Sharma (2006), A comparison of three stochastic multi-site
936 precipitation occurrence generators, *J. Hydrol.*, 331(1–2), 280–292,
937 doi:10.1016/j.jhydrol.2006.05.016.
- 938 Mostafa, M. D., and M. W. Mahmoud (1964), On the Problem of Estimation for the Bivariate
939 Lognormal Distribution, *Source Biometrika Biometrika Trust*, 51(34), 522–527.
- 940 Nataf, A. (1962), Statistique mathématique-determination des distributions de probabilités dont
941 les marges sont données, *C. R. Acad. Sci. Paris*, 255(1), 42–43.
- 942 Nazemi, A., H. S. Wheater, K. P. Chun, and A. Elshorbagy (2013), A stochastic reconstruction
943 framework for analysis of water resource system vulnerability to climate-induced changes in
944 river flow regime, *Water Resour. Res.*, 49(1), 291–305, doi:10.1029/2012WR012755.
- 945 Papalexiou, S. M. (2017), A unified theory for exact stochastic modelling of univariate and

946 multivariate processes with continuous, mixed type, or discrete marginal distributions and
947 any correlation structure,

948 Papalexiou, S. M., and D. Koutsoyiannis (2012), Entropy based derivation of probability
949 distributions: A case study to daily rainfall, *Adv. Water Resour.*, 45, 51–57,
950 doi:10.1016/j.advwatres.2011.11.007.

951 Papalexiou, S. M., D. Koutsoyiannis, and A. Montanari (2011), Can a simple stochastic model
952 generate rich patterns of rainfall events?, *J. Hydrol.*, 411(3–4), 279–289,
953 doi:10.1016/j.jhydrol.2011.10.008.

954 Pegram, G. G. S., and W. James (1972), Multilag multivariate autoregressive model for the
955 generation of operational hydrology, *Water Resour. Res.*, 8(4), 1074–1076,
956 doi:10.1029/WR008i004p01074.

957 Salas, J. D. (1993), Analysis and modeling of hydrologic time series, in *Handbook of hydrology*,
958 edited by D. R. Maidment, p. Ch. 19.1-19.72, Mc-Graw-Hill, Inc.

959 Salas, J. D., and G. G. S. Pegram (1977), A seasonal multivariate multilag autoregressive model
960 in hydrology, in *Proc. Third Int. Symp. on Theoretical and Applied Hydrology, Colorado*
961 *State Univ., Fort Collins, CO, USA*.

962 Salas, J. D., J. W. Delleur, V. Yevjevich, and W. L. Lane (1980), *Applied modeling of hydrologic*
963 *time series*, 2nd Print., Water Resources Publication, Littleton, Colorado.

964 Salas, J. D., G. Q. Tabios, and P. Bartolini (1985), Approaches to multivariate modeling of water
965 resources time series, *J. Am. Water Resour. Assoc.*, 21(4), 683–708, doi:10.1111/j.1752-
966 1688.1985.tb05383.x.

967 Serinaldi, F., and F. Lombardo (2017), BetaBit: A fast generator of autocorrelated binary processes
968 for geophysical research, *EPL (Europhysics Lett.)*, 118(3), 30007, doi:10.1209/0295-
969 5075/118/30007.

970 Sharma, A., D. G. Tarboton, and U. Lall (1997), Streamflow simulation: A nonparametric
971 approach, *Water Resour. Res.*, 33(2), 291–308, doi:10.1029/96WR02839.

972 Sklar, A. (1973), Random variables, joint distribution functions, and copulas, *Kybernetika*, 9(6),
973 449–460.

974 Srinivas, V. V., and K. Srinivasan (2005), Hybrid moving block bootstrap for stochastic simulation
975 of multi-site multi-season streamflows, *J. Hydrol.*, 302(1–4), 307–330,
976 doi:10.1016/j.jhydrol.2004.07.011.

977 Thomas, H. A., and R. P. Burden (1963), *Operations research in water quality management*,
978 HARVARD UNIV CAMBRIDGE MASS DIV OF ENGINEERING AND APPLIED
979 PHYSICS.

980 Thomas, H. A., and M. B. Fiering (1962), Mathematical synthesis of streamflow sequences for the
981 analysis of river basins by simulation, *Des. water Resour. Syst.*, 459–493.

982 Todini, E. (1980), The preservation of skewness in linear disaggregation schemes, *J. Hydrol.*,
983 47(3–4), 199–214, doi:10.1016/0022-1694(80)90093-1.

984 Wilks, D. S. (1998), Multisite generalization of a daily stochastic precipitation generation model,
985 *J. Hydrol.*, 210(1–4), 178–191, doi:10.1016/S0022-1694(98)00186-3.

- 986 Wilks, D. S., and R. L. Wilby (1999), The weather generation game: a review of stochastic weather
987 models, *Prog. Phys. Geogr.*, 23(3), 329–357, doi:10.1191/030913399666525256.
- 988 Xiao, Q. (2014), Evaluating correlation coefficient for Nataf transformation, *Probabilistic Eng.*
989 *Mech.*, 37, 1–6, doi:10.1016/j.probengmech.2014.03.010.
- 990

991 **Appendix A**

992 We briefly present the contemporaneous PAR(1) model with Pearson type-III (3-parameter
 993 Gamma) white noise (referred as PAR-PIII), for multivariate simulation of monthly time series.
 994 The model is able to preserve the essential statistics (i.e., mean, variance and skewness coefficient)
 995 as well as the lag-1 month-to-month correlations (i.e., autocorrelations) and the lag-0 cross-
 996 correlations between locations. Following the notation of *Koutsoyiannis [1999]*, let $\underline{\mathbf{x}}_s =$
 997 $[\underline{\mathbf{x}}_s^1, \dots, \underline{\mathbf{x}}_s^m]^T$ be a vector which of m stochastically dependent processes at season s . The generating
 998 scheme is:

$$\underline{\mathbf{x}}_s = \mathbf{A}_s \underline{\mathbf{x}}_{s-1} + \mathbf{B}_s \underline{\mathbf{w}}_s \quad (\text{A.1})$$

999 where $\mathbf{A}_s, \mathbf{B}_s$ are seasonally-varying $m \times m$ parameter matrices and $\underline{\mathbf{w}}_s = [\underline{w}_s^1, \dots, \underline{w}_s^m]^T$ is a
 1000 vector of independent random variables generated from Pearson type-III distribution. The matrices
 1001 \mathbf{A}_s are calculated as follows:

$$\mathbf{A}_s = \text{diag} \left(\frac{\text{Cov}[\underline{\mathbf{x}}_s^1, \underline{\mathbf{x}}_{s-1}^1]}{\text{Var}[\underline{\mathbf{x}}_{s-1}^1]}, \dots, \frac{\text{Cov}[\underline{\mathbf{x}}_s^m, \underline{\mathbf{x}}_{s-1}^m]}{\text{Var}[\underline{\mathbf{x}}_{s-1}^m]} \right) \quad (\text{A.2})$$

1002 while matrices \mathbf{B}_s are given by:

$$\mathbf{B}_s \mathbf{B}_s^T = \mathbf{G}_s \quad (\text{A.3})$$

1003 where

$$\mathbf{G}_s = \text{Cov}[\underline{\mathbf{x}}_s, \underline{\mathbf{x}}_s] - \mathbf{A}_s \text{Cov}[\underline{\mathbf{x}}_{s-1}, \underline{\mathbf{x}}_{s-1}] \mathbf{A}_s^T \quad (\text{A.4})$$

1004 where $\text{Cov}[\underline{\xi}, \underline{\psi}]$ denotes the covariance of vectors $\underline{\xi}$ and $\underline{\psi}$, i.e., $\text{Cov}[\underline{\xi}, \underline{\psi}] = \text{E} \left\{ \left(\underline{\xi} - \right. \right.$
 1005 $\left. \text{E}[\underline{\xi}] \right) \left(\underline{\psi}^T - \text{E}[\underline{\psi}]^T \right) \}$. At each season s , the parameter matrix \mathbf{B}_s can be estimated either through
 1006 typical decomposition techniques (e.g., Cholesky or singular value decomposition) or numerically
 1007 approximated, e.g., through optimization approaches [*Koutsoyiannis, 1999; Higham, 2002*].

1008 Regarding the white noise vector $\underline{\mathbf{w}}_s$, its statistical structure is associated with the seasonal
 1009 statistical characteristics of the parent process, through the following equations:

$$\text{E}[\underline{\mathbf{w}}_s] = \mathbf{B}_s^{-1} \{ \text{E}[\underline{\mathbf{x}}_s] - \mathbf{A}_s \text{E}[\underline{\mathbf{x}}_{s-1}] \} \quad (\text{A.5})$$

$$\text{Var}[\underline{\mathbf{w}}_s] = [1, \dots, 1]^T \quad (\text{A.6})$$

$$\mu_3[\underline{\mathbf{w}}_s] = (\mathbf{B}_s^{(3)})^{-1} \{ \mu_3[\underline{\mathbf{x}}_s] - \mathbf{A}_s^{(3)} \mu_3[\underline{\mathbf{x}}_{s-1}] \} \quad (\text{A.7})$$

1010 where $\mathbf{B}_s^{(k)}$ is a matrix whose elements are raised to power k while $\mu_3[\underline{\mathbf{w}}_s]$ and $\mu_3[\underline{\mathbf{x}}_s]$ are vectors
 1011 that denote the third central moments of $\underline{\mathbf{w}}_s$ and $\underline{\mathbf{x}}_s$ respectively. The white noise is produced by a
 1012 suitable random number generator, in particular the Pearson type-III distribution, which can
 1013 explicitly preserve $\text{E}[\underline{\mathbf{w}}_s]$, $\text{Var}[\underline{\mathbf{w}}_s]$ and $\mu_3[\underline{\mathbf{w}}_s]$.