

Save hydrological observations! Return period estimation without data decimation

Elena Volpi^{a,*}, Aldo Fiori^a, Salvatore Grimaldi^{b,c}, Federico Lombardo^d,
Demetris Koutsoyiannis^e

^a*Department of Engineering, University of Roma Tre, Rome, Italy*

^b*Department for Innovation in Biological, Agro-food and Forest systems (DIBAF),
University of Tuscia, Viterbo, Italy*

^c*Honors Center of Italian Universities (H2CU), Sapienza University of Rome, Rome,
Italy*

^d*Dipartimento di Ingegneria Civile, Edile e Ambientale, Sapienza Università di Roma,
Rome, Italy*

^e*Department of Water Resources and Environmental Engineering, School of Civil
Engineering, National Technical University of Athens, Greece*

Abstract

The concept of return period and its estimation are pivotal in risk management for many geophysical applications. Return period is usually estimated by inferring a probability distribution from an observed series of the random process of interest and then applying the classical equation, i.e. the inverse of the exceedance probability. Traditionally, we form a statistical sample by selecting, from the "complete" time series (e.g. at the daily scale), those values that can reasonably be considered as realizations of independent extremes, e.g. annual maxima or peaks over a certain high threshold. Such a selection procedure entails that a large number of observations are discarded; this wastage of information could have important consequences in practical prob-

*Corresponding author

Email address: elena.volpi@uniroma3.it (Elena Volpi)

lems, where the reduction of the already small size of common hydrological records significantly affects the reliability of the estimates. Under such circumstances, it is crucial to exploit all the available information. To this end, we investigate the advantages of estimating the return period without any data decimation, by using the full data-set. The proposed procedure, denoted as *Complete Time-series Analysis* (CTA), exploits the property that the average interarrival time (i.e. return period) of potentially damaging events is not affected by the dependence structure of the underlying process, even for cyclo-stationary (e.g. seasonal) processes. For the sake of illustration, the CTA is compared to that based on annual maxima selection, through a simple non-parametric approach, discussing advantages and limitations of the method. Results suggest that the proposed CTA approach provides a more conservative return period estimation in an holistic implementation framework within a broader range of return period values than that pertaining to other methods, which means not only the largest extremes that are the focus of extreme value theory.

Keywords: Return period; Interarrival time; Complete Time-series Analysis; Persistence; Seasonality; Annual maxima

1. Introduction

Hydraulic risk analysis relies on finding the probability of failure of a given hydraulic structure or, more generally, system due to the occurrence of intense hydrological events, where the probability of failure is usually expressed in terms of return period. Different failure mechanisms could be considered, where each of them results from the combination of multiple characteristics

7 of the hydrological loads (Schumann, 2017). Hence, under general condi-
8 tions, the return period of structure failure should be quantified taking into
9 account the joint probability of failure mechanisms, i.e. the joint probability
10 of the random variables describing the hydrological load and the complex
11 interactions between the structure and the hydrological loads acting on it
12 (see, e.g. Volpi and Fiori, 2014).

13 In the simplest case, we have a single failure mechanism that is ruled
14 by a single random variable describing the hydrological load, for example a
15 bridge destroyed by a flood. Under such circumstances, the return period
16 of structure failure corresponds to that of the hydrological load. Once the
17 key variable representing the hydrological load is identified, the problem is
18 solved by inferring a probability distribution from a series of realizations
19 of this random variable, in order to determine the magnitude of the event,
20 corresponding to a given return period or probability of failure.

21 Given this premise, it is clear that the concept of return period and how
22 it is estimated from observations is central to risk management problems in
23 hydrological/hydraulic applications; yet this is true also in many other geo-
24 physical and engineering fields. Even if return period is a widely applied
25 and well established probabilistic tool for hydrological applications, since the
26 pioneering work of Alexander (1959) there have been few studies attempt-
27 ing to analyze the differences between estimated return periods of hydro-
28 logical extremes using different methods of estimation. Nonetheless, some
29 researchers have recently investigated the concept of return period when the
30 basic assumptions of stationarity and independence are omitted; see, among
31 others, Rootzén and Katz (2013); Obeysekera and Salas (2016); Read and

32 Vogel (2016); Fernández and Salas (1999); Douglas et al. (2002); Bunde et al.
33 (2003); Eichner et al. (2011); Volpi et al. (2015) and references therein. As
34 detailed in the following, the purpose of our work is to investigate a new sta-
35 tistical approach to infer return period from a complete record of observed
36 data; therefore, we must assume a dependence structure in time and a sta-
37 tionary framework, because the non-stationary hypothesis implies a priori
38 attributions, supported by deductive reasoning, that go beyond the scope of
39 this paper (Koutsoyiannis and Montanari, 2015; Serinaldi et al., 2018).

40 Under the assumption of stationarity, Bunde et al. (2003) and Volpi et al.
41 (2015) have shown that the independence condition is not necessary in or-
42 der to apply the classical equation of return period, i.e., the inverse of the
43 exceedance probability. Following Volpi et al. (2015), this paper highlights
44 how temporal dependence does not alter the average interarrival time for-
45 mulation, even in stochastic processes characterized by *cyclo-stationarity*,
46 a characteristic that hydrological and other geophysical processes exhibit at
47 sub-annual scales. Furthermore, we investigate here the potential application
48 of this important property of return period, derived from the full available
49 record, for frequency analysis; specifically, we show how the return period
50 can be directly estimated from raw data records of a time-dependent process,
51 regardless of its dependence structure, under stationary or cyclo-stationary
52 conditions.

53 This alternative approach for return period estimation, which is proposed
54 here for the first time and denoted as *Complete Time-series Analysis* (CTA),
55 is compared to the traditional approach based on frequency analysis of An-
56 nual Maxima (AM), which constitutes the basis of traditional extreme value

57 analysis. Indeed, we usually analyze AM to catch the tail of the distribution
58 of the parent process, where the latter is the process of interest. Hence, the
59 rationale behind CTA is to exploit all the information provided by obser-
60 vational data (Marani and Ignaccolo, 2015; Zorzetto et al., 2016), with the
61 objective of better estimating the return period in a wider range of values,
62 not only at the largest extremes that are the focus of extreme value theory.
63 Note, indeed, that small to moderate return period values are still of interest
64 in several practical problems, such as pluvial flooding. Furthermore, it is
65 important to stress that CTA provides different return period estimates with
66 respect to annual maxima by considering *all* the occurrences of the *dangerous*
67 values (e.g. exceedance of the random variable above any threshold value of
68 interest) within the observed record, as it will be discussed later on.

69 Hence, we aim at exploring the potential conveniences of CTA compared
70 to traditional approaches, and not to elaborate on the return period concept
71 (see Volpi et al., 2015). For this reason, in this article we base and limit our
72 investigation to the non-parametric approach. Two illustrative examples are
73 presented, both relying on synthetic processes: the first one makes use of a
74 very simple process whose correlation structure is known a priori; the second
75 example resembles the main characteristics of a real-world process. Future
76 research will focus on the parametric implementation of CTA for its practical
77 use in real world cases; indeed, the problem of fitting a model to the complete
78 record of observations and the evaluation of the related uncertainty deserves
79 further attention.

80 The remainder of this paper is organized as follows. Section 2 briefly re-
81 calls the definitions of return period available in the literature and illustrates

82 the properties of the interarrival time for stationary and cyclo-stationary
 83 processes. In Section 3 the new approach for return period estimation, re-
 84 lying on frequency analysis of complete times-series (CTA), is introduced as
 85 an alternative to traditional methods based on annual maxima and peaks
 86 over threshold (Section 4). In sections 5 and 6, complete time-series analy-
 87 sis is compared to the standard approach using annual maxima for a simple
 88 stationary process and for a cyclo-stationary process, that mimics the char-
 89 acteristics of a real world phenomenon, respectively. Section 7 is concerned
 90 with the potential issues related to the application of the proposed approach
 91 in real word problems, while the Conclusion Section summarizes the main
 92 findings of this paper.

93 2. Return period: Definitions and properties

94 Let $Z(\tau)$ be a stochastic process that represents a natural process evolving
 95 in continuous time τ . As observations of Z are only made in discrete time,
 96 we consider here the corresponding discrete-time process Z_j that is obtained
 97 by sampling $Z(\tau)$ at constant time intervals $\Delta\tau$, i.e. $Z_j = Z(j\Delta\tau)$ where
 98 j ($= 1, 2, \dots$) denotes discrete time. We make the only assumption that Z_j
 99 is a stationary process, fully described in terms of its marginal probability
 100 function $P_Z(z) = \Pr\{Z \leq z\}$ and, up to the second order in terms of joint dis-
 101 tribution, by its autocorrelation function $\rho_\theta = \gamma_\theta/\gamma_0$ (with $\theta = 0, \pm 1, \pm 2, \dots$),
 102 where $\gamma_\theta = \text{cov}[Z_j, Z_{j+\theta}]$ and $\rho_\theta \in [-1, 1]$; further, we denote the mean of
 103 the process as $\mu = E[Z_j]$ and its standard deviation $\sigma = \sqrt{\gamma_0}$.

104 For design and risk assessment purposes, we are interested in the occur-
 105 rence of dangerous events that might result in a system or structure failure.

106 We define here a dangerous event as the exceedance at the $\Delta\tau$ scale of a
 107 threshold level, $A = \{Z > z\}$, for instance the discharge exceeding a given
 108 high threshold level, potentially causing the flooding of an urbanized area.
 109 The probability of A is given by $\Pr A = \Pr \{Z > z\} = 1 - P_Z(z) = 1 - \Pr B$,
 110 where B denotes the complementary event of A . In hydrological applications,
 111 as well as in many other engineering fields, the rareness of the dangerous
 112 events is usually measured in terms of return period $T(z)$, thus assuming
 113 that the event A will occur on average once every T years, which is the time
 114 unit commonly used for return periods in hydrology. Mathematically, it is

$$115 \quad \frac{T}{\Delta\tau} = \mathbb{E}[X] = \sum_{x=1}^{\infty} x f_X(x) \quad (1)$$

116 where X is the number of discrete time steps to an occurrence of the event A
 117 and $f_X(x) = \Pr \{X = x\}$ is its probability mass function (pmf). Note that in
 118 eq. (1), which only refers to discrete-time processes, T is measured in units
 119 of time, i.e. $\Delta\tau$; if $\Delta\tau = 1$ year the return period is measured in years.

120 As highlighted by previous literature studies (see, e.g. Fernández and
 121 Salas, 1999; Douglas et al., 2002), the return period can be defined as the
 122 average of (i) the *waiting time*, that is the time interval ranging from the
 123 present to the next threshold exceedance, or (ii) the *interarrival time*, that
 124 is the time elapsing between any two successive realizations of the dangerous
 125 event. As explained later on, we adopt here the second definition, which im-
 126 plicitly assumes that a dangerous event has just occurred. We remark that
 127 such a definition is customary in hydrological applications (see, e.g. Chow
 128 et al., 1988; Kottegoda and Rosso, 1997; Salvadori et al., 2007). For conve-
 129 nience, herein we express discrete time as $t = j - j_0$, where j_0 is the current
 130 instant of time (when the dangerous event has just occurred); therefore, the

131 discrete-time process is indicated as Z_t , and $t = 0$ denotes the present. As a
 132 consequence, the pmf of the interarrival time can be written as (Fernández
 133 and Salas, 1999)

$$\begin{aligned}
 f_X(x) &= \Pr(B_1, B_2, \dots, B_{x-1}, A_x | A_0) \\
 &= \frac{\Pr(A_0, B_1, B_2, \dots, B_{x-1}, A_x)}{\Pr A_0}
 \end{aligned}
 \tag{2}$$

135 Then, the definition of return period based on the concept of interarrival time
 136 relies on a conditional probability. The average interarrival time is obtained
 137 by substituting Equation (2) into (1).

138 If Z_t is a purely random process, then the return period T is given by
 139 (e.g. Stedinger et al., 1993)

$$\frac{T(z)}{\Delta\tau} = \frac{1}{1 - P_Z(z)}
 \tag{3}$$

141 regardless of the definition used for X in Eq.(1), i.e. waiting time or inter-
 142 arrival time. The above relationship holds true even if the stationary and
 143 independent process is not sampled at constant time intervals; in this case
 144 $\Delta\tau$ is the average time interval between consecutive samples (Koutsoyiannis,
 145 2008).

146 Although the independence condition is typically assumed as a necessary
 147 condition for Equation (3), it has been recently demonstrated by Volpi et al.
 148 (2015) - independently from the conceptual arguments presented by Bunde
 149 et al. (2003) - that the return period $T(z)$, defined as the average interarrival
 150 time, is expressed by Eq. (3) even for processes correlated in time, with any
 151 type of dependence structure of Z , an important and not very well-known
 152 fact that is exploited in the following development.

153 Conversely, when based on the concept of waiting time, the formulation of
154 the return period strictly depends on the correlation structure of the process;
155 specifically, for any dangerous events A (or threshold levels z), the average
156 waiting time is an increasing function of the correlation of the process, thus
157 resulting in values larger than the average interarrival time (which ignores
158 correlation). Hence, for processes which are positively correlated in time
159 (such as most hydrological processes) Equation (3) returns a *lower bound* in
160 terms of $T(z)$.

161 Although the average interarrival time $T(z)$ remains the same for cor-
162 related and independent processes, the probability that the threshold z is
163 exceeded in a given period can be very different in the two cases. In fact, if
164 a dangerous event occurs at present time, then the conditional probability of
165 occurrence of another dangerous event at successive instants of time will be
166 greater than the independent case; this yields that the probability mass func-
167 tion of the interarrival time (that corresponds to the probability of failure)
168 will have a larger mass for small temporal values and a lower mass elsewhere,
169 hence a larger variance with respect to the independent case. Since the aver-
170 age waiting time is an increasing function of the variance of the interarrival
171 time, as shown in Volpi et al. (2015), the latter characteristic of the interar-
172 rival time distribution explains why the average waiting time is larger than
173 the average interarrival time. As a further consequence, the definition of
174 return period based on the interarrival time might result in higher values of
175 the probability of failure with respect to the independent case; the reader is
176 referred to Volpi et al. (2015) for further details on the theoretical properties
177 of both definitions of T .

178 *2.1. Cyclo-stationary processes*

179 The property of interarrival time mentioned above has been derived based
 180 on the assumption that the underlying process Z is stationary (see the Ap-
 181 pendix B in Volpi et al. (2015) for details). However, many natural processes
 182 exhibit statistical properties that are invariant to a shift of the time origin by
 183 integral multiples of a certain period Π , due to e.g. the seasonal variability
 184 of environmental phenomena at sub-annual scales (Koutsoyiannis, 2016). In
 185 stochastic hydrology, such processes are usually modelled by cyclo-stationary
 186 processes with period Π .

187 Let us consider a cyclo-stationary process that is characterized by a joint
 188 distribution function that varies within the time period Π (typically equal to
 189 one year), such that $\Pr B_t = \Pr B_{\Pi+t}$, $\Pr(B_t, B_{t+1}) = \Pr(B_{\Pi+t}, B_{\Pi+t+1})$ and
 190 so on. For such a process, the pmf given in Equation (2) can be regarded as
 191 the pmf of the interarrival time conditional to the occurrence of the dangerous
 192 event at time $t = 0$, i.e. $f_X(x|t = 0)$. For any value of t , this conditional pmf
 193 can be written as

$$194 \quad f_X(x|t) = \frac{\Pr(A_t, B_{t+1}, B_{t+2}, \dots, B_{t+x-1}, A_{t+x})}{\Pr A_t} \quad (4)$$

195 To account for the possible occurrence of A at every instant of time t within
 196 the period Π , we marginalize the above conditional probability by summing
 197 the pmf in Eq.(4) with respect to all the possible values of time $t' \in [t, t+\Pi-1]$
 198 according to their probability of occurrence. The latter quantity is nothing

199 else than the conditional probability $\Pr(A_{t'}|t' \in [t, t + \Pi - 1])$. Hence, it is

$$\begin{aligned}
 f_X(x|t) &= \sum_{t'=t}^{t+\Pi-1} \Pr(A_{t'}|t' \in [t, t + \Pi - 1]) \frac{\Pr(A_{t'}, B_{t'+1}, B_{t'+2}, \dots, B_{t'+x-1}, A_{t'+x})}{\Pr A_{t'}} \\
 &= \sum_{t'=t}^{t+\Pi-1} \frac{\Pr A_{t'}}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \frac{\Pr(A_{t'}, B_{t'+1}, B_{t'+2}, \dots, B_{t'+x-1}, A_{t'+x})}{\Pr A_{t'}}
 \end{aligned}
 \tag{5}$$

201 Finally, the average interarrival time of the cyclo-stationary process is
 202 obtained by substituting in Equation (1) the pmf given in Eq. (5), thus
 203 obtaining

$$\frac{T(z)}{\Delta\tau} = \frac{1}{1 - \overline{P_Z(z)}}
 \tag{6}$$

205 where $\overline{P_Z(z)} = \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} = \frac{1}{1 - \sum_{r=t}^{t+\Pi-1} \Pr B_r}$ is the marginal probability of
 206 non-exceeding the threshold value z within any period $[t, t + \Pi - 1]$; since
 207 we are dealing with a cyclo-stationary processes, $\overline{P_Z(z)}$ remains the same
 208 for any t . The derivation of Equation (6) is given in Appendix A. While
 209 for cyclo-stationary processes the exceeding probability $\Pr A_{t'} = 1 - \Pr B_{t'}$
 210 varies with time $t' \in [t, t + \Pi - 1]$, the return period of the dangerous event
 211 A is a constant value, independent of time t , and it is expressed again by the
 212 classical equation of return period.

213 **3. Novel return period estimation: Complete Time-series Analysis** 214 **(CTA)**

215 The property that the average interarrival time (i.e. return period) is
 216 not affected by the dependence structure of the underlying process, even for
 217 cyclo-stationary processes, sheds a new light on the problem of return period
 218 estimation in practical problems. Under the hypothesis of stationarity, or

219 cyclo-stationarity in the case of processes exhibiting seasonality and sam-
 220 pled at the sub-annual scale, the return period can be estimated by using
 221 Equation (3) starting from any kind of observational data, independent or
 222 correlated in time, thereby potentially exploiting all the available information
 223 on the underlying process. Hence, the only necessary assumption is that of
 224 stationarity; but stationarity is also related to ergodicity, which in turn is a
 225 prerequisite to make inference from data. As previously mentioned, the sta-
 226 tionarity issue goes beyond the scope of this work; the interested readers are
 227 referred to the work of Koutsoyiannis (2014), Montanari and Koutsoyiannis
 228 (2014), Serinaldi and Kilsby (2015), Koutsoyiannis (2016), Serinaldi et al.
 229 (2018) and Luke et al. (2017) for a comprehensive discussion.

230 Let $\mathbf{z} = z_1, \dots, z_n$ be an observed realization of the stochastic process Z_t ,
 231 where n is the length of the data series. Following its formal definition based
 232 on the interarrival time, the empirical return period could be estimated by
 233 first deriving from \mathbf{z} the sample of the interarrival time for each value of the
 234 threshold z , i.e. $\mathbf{x}(z) = x_1, x_2, \dots, x_{\eta(z)}$, and then averaging in time, as in
 235 Equation (1). This approach can be applied only in the case of very long
 236 time-series (very large n), for which the size of the interarrival series $\eta(z)$ is
 237 large enough to return reliable estimates for high threshold values z .

238 In hydrological applications the length of the observed series is usually
 239 less than one hundred years (e.g. at the daily or at the hourly scale); hence
 240 it is important to exploit all the available information for sample estimates.
 241 The empirical return period can be obtained by directly applying Equation
 242 (3) to the entire time-series \mathbf{z} , where the cumulative probability function P_Z
 243 is substituted by a probability model inferred from the observed data. For

244 simplicity, we adopt here a non-parametric approach, by substituting P_Z with
 245 its empirical counterpart. The empirical distribution function (edf), denoted
 246 in the following as \hat{F}_n , provides an estimate of the distribution function of
 247 the underlying stochastic process (e.g. Kolmogorov, 1933). Let $\{Z_i\}_{i=1}^n$ be a
 248 sequence of dependent and equally distributed random variables, then F_n is
 249 defined as the fraction of random variables that are less than or equal to the
 250 specified value z , i.e.

$$251 \quad F_n(z) = \frac{1}{n} \sum_{i=1}^n I_{\{Z_i \leq z\}} \quad (7)$$

252 where $I_{Z_i \leq z}$ is the indicator of the event $\{Z_i \leq z\}$; the estimate, $\hat{F}_n(z)$, is
 253 obtained by considering in Eq. (7) the outcome \mathbf{z} instead of the random
 254 variables.

255 The edf in Eq. (7) is an unbiased estimator of the marginal probability
 256 function, i.e. $E[F_n(z)] = P_Z(z)$; this means that the dependence structure
 257 of the underlying process does not affect the expectation of the edf, while it
 258 affects its covariance, as reported by Azriel and Schwartzman (2015). More-
 259 over, these Authors remark that the estimator in Eq. (7) is consistent for
 260 all Gaussian stationary ergodic processes, that are characterized by an au-
 261 tocorrelation function decreasing to zero as the lag-time goes to infinity (i.e.
 262 the necessary condition for ergodicity). The latter case includes both short
 263 and long-range dependent processes, such as the Hurst-Kolmogorov process
 264 (Koutsoyiannis, 2016). The consistency property of the edf continues to hold
 265 for non-Gaussian distributions under various forms of dependence (see, e.g.,
 266 Dedecker and Merlevéde, 2007; Wu, 2006). Note that in the case of cyclo-
 267 stationary processes, where the marginal probability distribution changes
 268 with time within the period Π , \hat{F}_n directly provides an estimate of $\overline{P_Z(z)}$,

269 which takes into account the variability during the period of the marginal
270 probability function (i.e. the alternation of events characterized by a different
271 probability of exceedance).

272 The computational approach based on the complete edf (complete time-
273 series analysis, CTA), returns on average the same results to those based
274 on the observed interarrival time-series, provided that n (hence $\eta(z)$) is very
275 large. An illustration example showing the latter property will be presented
276 in the following sections.

277 The estimation approach proposed in this Section makes use of the whole
278 available information on the process Z , i.e. the complete time-series \mathbf{z} , with
279 the aim of returning a reliable and robust estimate of the return period of
280 the dangerous event $\{Z > z\}$ for any threshold value z , i.e. in a broader
281 range of values with respect to the traditional approaches used in extreme
282 value analysis. It is important to stress that CTA requires the availability of
283 an uninterrupted record of observations (i.e. \mathbf{z}) where the process is sampled
284 at constant time intervals ($\Delta\tau$); indeed, this is a necessary condition for the
285 property of the average interarrival time to hold true, as explained in Volpi
286 et al. (2015).

287 In the following sections, we compare CTA with the traditional sampling
288 methods used in extreme value analysis, thus highlighting advantages and
289 limitations of both the strategies.

290 **4. Traditional approaches for return period estimation**

291 Since independence of Z_t is usually invoked for the derivation of Equation
292 (3) (e.g. Benjamin and Cornell (1970), p. 233, Kottegoda and Rosso (1997),

293 p. 190 and Chow et al. (1988), p. 383), it is common practice in hydrological
294 applications to implement some techniques for data selection aimed to allow
295 the assumption of the statistical independence of the observations. These
296 techniques constitute the basis for the extreme value analysis, whose objec-
297 tive is to quantify the stochastic behavior of a process at unusually large
298 or small levels that potentially lead to the failure of a system (Salvadori
299 et al., 2007; Coles, 2001). Classical methods for extreme (and independent)
300 value selection are the block maxima approach, where the block generally
301 coincides with the year (Annual Maxima, AM), and the more complex Peak-
302 Over-Threshold approach (POT).

303 The wide popularity of AM relies on its simplicity, but also on the limited
304 access in the past to regularly-sampled, long observed series of the random
305 variable of interest. Although it is necessary to have available the complete
306 series to establish if an extreme event is an annual maximum, it was common
307 practice to take note essentially of the values reached during extreme events,
308 especially in old times when the observation of the hydrological variables was
309 not systematic (apart from a few noteworthy cases such as that described in
310 Calenda et al. (2005)).

311 When a complete time-series of observations is available, POT is adopted
312 in practical applications (especially when the length of the observed period is
313 short) to select the largest possible amount of data, yet respecting the inde-
314 pendence assumption. POT originated in hydrology with the rationale that
315 if additional information about the extreme upper tail were used besides the
316 annual maxima, then more accurate and reliable estimates of the parame-
317 ters and quantiles of extreme value distributions would be obtained (see, e.g.

318 Katz et al., 2002). This is the same rationale behind CTA, as described in
 319 the previous section; however, it is important to stress that while CTA does
 320 not require any data selection, POT asks for additional computational efforts
 321 to select only the peaks over the threshold (i.e. the maximum of a cluster
 322 of values all exceeding the threshold) and necessitates the introduction of
 323 further information or parameters in order to select among those peaks only
 324 independent ones (see, e.g. Coles, 2001).

325 For simplicity, in this work we compare CTA to AM. To avoid confusion,
 326 we use Y to indicate the annual maximum of the random variable Z . The
 327 annual maximum time-series is derived from \mathbf{z} as $\mathbf{y} = \{y_1, \dots, y_{n/n_Y}\}$, where
 328 n is the number of observations, $y_i = \max\{z_{n_Y(i-1)+1}, \dots, z_{n_Y i}\}$ and n_Y is
 329 the number of time-intervals $\Delta\tau$ per year (e.g. $n_Y = 365$ in the case Z is
 330 observed at the daily scale). The empirical return period of the values in
 331 \mathbf{y} can be evaluated by using the same procedures described above for the
 332 time-series \mathbf{z} and based on the edf (since $(n/n_Y) \leq n$).

333 4.1. Purely random and stationary processes

334 The probability distribution function of Y is by definition different from
 335 that of Z . Since we aim at exploring the difference between the two under
 336 general conditions (cyclo-stationarity and persistence), for the sake of clarity
 337 we start from the well known stationary and independent case (then $\rho_\theta = 0$
 338 for $\theta \neq 0$) by introducing a general framework that is instrumental to the
 339 discussion reported in the following sections. Given a threshold value z , the
 340 probability of annual maxima $P_Y(z) = \Pr\{Y \leq z\}$ can be easily derived from
 341 that of the parent process as (e.g. Coles, 2001)

$$342 \quad P_Y(z) = P_Z(z)^{n_Y} \quad (8)$$

343 By using Equation (3) the corresponding return period is derived

$$344 \quad T_Y(z) = \frac{\Delta\tau n_Y}{1 - P_Z(z)^{n_Y}} \quad (9)$$

345 where $\Delta\tau n_Y$ is one year. Note that here n_Y is not a random variable, since
346 \mathbf{z} is an observed series of the stochastic process Z sampled at a constant time
347 intervals $\Delta\tau$; hence, the number of observations in each year is constant, be-
348 ing uniquely determined by the time interval $\Delta\tau$. Conversely, in traditional
349 extreme value theory the exponent in Eq. (8) is not a constant, being the
350 number of peaks of clusters of values, but rather can be regarded as a realiza-
351 tion of a Poisson distributed random variable; this yields a different form for
352 the probability distribution of annual maxima, which gives numerical values
353 not significantly different from those provided by Eq. (8) for large $P_Z(z)$
354 (Koutsoyiannis, 2004).

355 If $n_Y = 1$ CTA obviously gives the same results of AM. If $n_Y > 1$ (which
356 means that $\Delta\tau < 1$ year), Equation (9) results in larger values with respect
357 to $1/(1 - P_Z(z))$, as shown in Figure 1. Note that the figure depicts the
358 theoretical return period of annual maxima as function of that of the con-
359 tinuous process, i.e. when assuming $n_Y \rightarrow \infty$, for any kind of process Z ;
360 in other words, it is the case of infinite sample length ($n \rightarrow \infty$). For con-
361 venience, the return period of annual maxima is denoted by T_Y while that
362 of the parent process by T_Z ; both are measured in years from now on. The
363 figure shows to what extent AM results in larger values of the return period
364 with respect to that of the underlying process, or in other words, the prob-
365 ability of $\{Y > z\}$ is smaller than that of the dangerous events $\{Z > z\}$.
366 This is due to a *wastage* of information. During any year, additional events
367 may have occurred that are excluded by the analysis because such data are

368 not the annual maximum in the year they arose, as in the example depicted
 369 in Figure 2a. For the sake of illustration, Figure 2 shows one year of a log-
 370 normal AR(1) daily time-series with mean $\mu = 1$, variance $\sigma^2 = 1$ and lag-1
 371 correlation coefficient $\rho_1 = 0$ (panel a) and $\rho_1 = 0.85$ (panel b); note that
 372 seasonal variability is not considered in this independent and stationary ex-
 373 ample. The figure depicts with red dots the information discarded by the
 374 annual maxima approach (red circles), in the estimation of the return period
 375 of the event $Z > 4$. Furthermore, these events might be possibly larger than
 376 the maximum in other years. It also follows that the minimum value of the
 377 return period of annual maxima is equal to 1 year, which means that based
 378 on annual maxima analysis we cannot measure the rareness of events that
 379 occur more frequently than once every year.

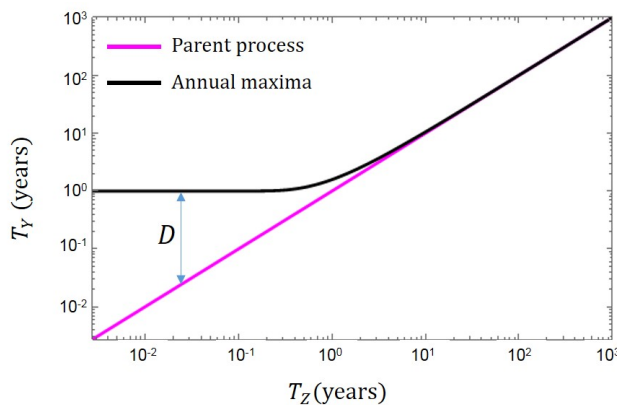


Figure 1: Return period of the annual maxima (T_Y) as function of that of the continuous parent process (T_Z , $n_Y \rightarrow \infty$). The difference between the two (D) significantly reduces ($D \leq 0.05$) only when when T_Z becomes larger than about 10 years.

380 As the threshold z increases (we look to more and more intense events),
 381 the return period of annual maxima, T_Y , tends to that of the parent process,

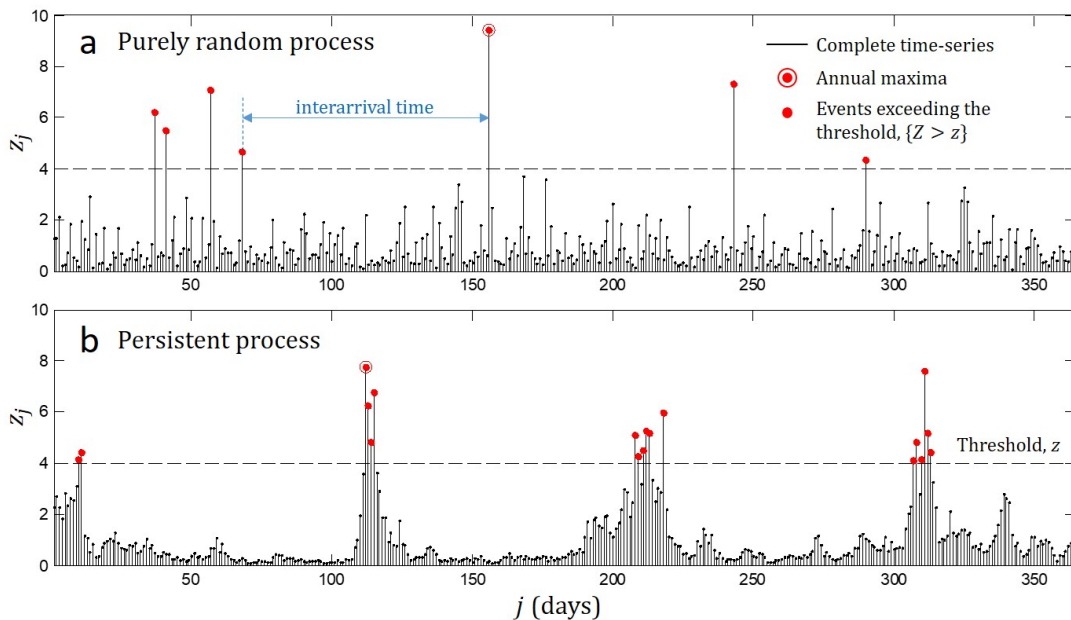


Figure 2: One year of a lognormal AR(1) daily time-series with mean $\mu = 1$, variance $\sigma^2 = 1$ and lag-1 correlation coefficient a) $\rho_1 = 0$ and b) $\rho_1 = 0.85$ (no seasonal variability). The events exceeding the threshold (red dots) that are not annual maxima (red circle) are discarded in AM, resulting in an overestimation of the average interarrival time of the parent process, i.e. $T_Y(z) > T_z(z)$.

382 T_Z ; indeed, very large events are expected to be always selected as annual
 383 maxima, hence the wastage of information is reduced toward zero as the
 384 threshold increases. The difference $T_Y - T_Z$ (denoted in Figure 1 as D)
 385 reduces to less than 5% of T_Z only when T_Z becomes larger than about 10
 386 years. This means that relatively frequent events, that might be of interest
 387 when the expected damage is modest, are generally underestimated when
 388 only the annual maxima are available. The asymptotic convergence of the
 389 annual maxima distribution to that of the parent process for high threshold
 390 values, about $T_Z \geq 10$ years according to the above figure, is at the basis of

391 extreme value theory.

392 **5. Persistent stationary processes**

393 The difference between annual maxima and the parent process for those
394 events that are characterized by small to medium values of the return period
395 is expected to worsen in the case of time-dependent, positively correlated
396 processes, where the dangerous events tend to occur in clusters; see the ex-
397 ample shown in Figure 2b, where the number of events neglected by annual
398 maxima approach (the red dots) increases with respect to the independent
399 case depicted figure 2a. The latter condition is that usually matched in hy-
400 drological applications; for instance, the rainfall amount observed at given
401 time-scale exhibits a complex, persistent behavior in time, affected by sea-
402 sonality, which depends on the time-scale itself. Hence, we compare here the
403 empirical return period of the complete time-series to that of annual maxima
404 for time persistent processes by making use of a simple synthetic example.

405 In this section we consider an autoregressive process of order one, AR(1),
406 supposed to represent a stationary and persistent natural process observed
407 at the daily scale (i.e. $\Delta\tau = 1$ day). The process analyzed here is charac-
408 terized by a marginal lognormal probability distribution function with mean
409 μ and variance σ^2 , while its time-dependence structure is ruled by the lag-1
410 correlation coefficient, ρ_1 . We assume here $\mu = 1$ and two different values
411 for the variance, $\sigma^2 = 1$ or $\sigma^2 = 5$; further, we let the correlation coefficient
412 vary between 0 (independent process) and 0.99 (persistent process), noting
413 that we are ignoring seasonality here.

414 *5.1. Theoretical difference between CTA and AM for persistent processes*

415 We first compare the empirical return period of the annual maxima to
416 that of the complete series by analyzing a very long series, specifically $n =$
417 365×10^5 days (i.e. 10^5 years). This analysis is intended to investigate the
418 theoretical difference between CTA and AM for persistent processes, when
419 the accuracy of the return period estimate is essentially not affected by the
420 length of the observed period (as in the independent case discussed in Section
421 4.1). Results are represented in Figure 3 for several values of ρ_1 ranging
422 between 0.5 and 0.99 and for $\sigma^2 = 1$ (Figure 3a) and $\sigma^2 = 5$ (Figure 3b); the
423 independent case ($\rho_1 = 0$) is reported as a reference. Values of $\rho_1 \in [0, 0.5]$
424 are not considered since the difference with the independent case is negligible;
425 finally, we look at results for return period values included between 1 day,
426 that is the the minimum value that can be explored based on the temporal
427 resolution of the available series, and 1000 years, such that estimates are not
428 influenced by the finite length of the simulated series.

429 We recall here that the theoretical return period (according to Equation
430 (3)) of the parent process is fully determined by the lognormal probability
431 distribution, which is represented in Figure 3 by the magenta curve. The
432 return period estimated from the complete series returns for any ρ_1 the the-
433 oretical distribution (black dashed curves that overlap for all ρ_1 the magenta
434 curve); thus, it is not affected at all by the correlation structure of the process
435 (as demonstrated by Volpi et al. (2015)).

436 Since in this numerical experiment n is very large, also the empirical
437 return period computed as the average of the interarrival time between suc-
438 cessive events $\{Z > z\}$, i.e. by strictly following Equation (1), returns the

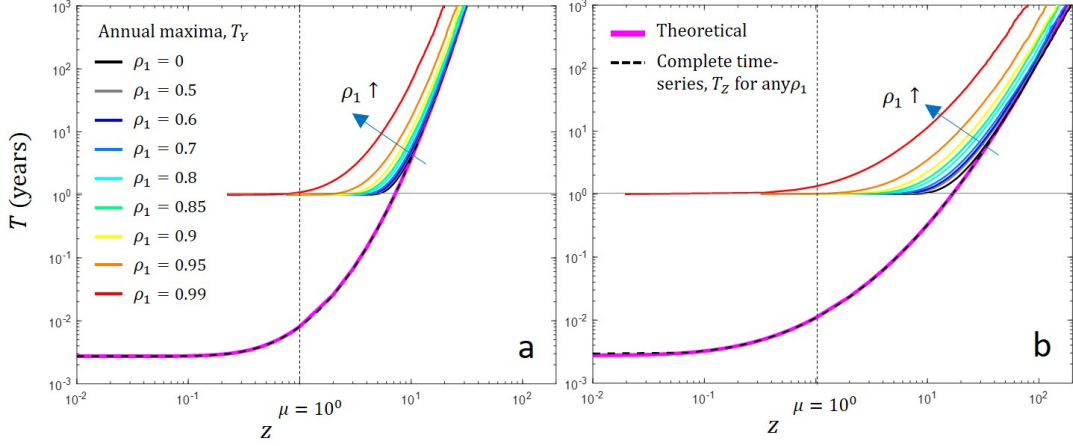


Figure 3: Lognormal AR(1) daily process with $\mu = 1$, variance $\sigma^2 = 1$ (a) or $\sigma^2 = 5$ (b), and lag-1 correlation coefficient ρ_1 : empirical return periods of the complete series (T_Z , black dashed curves indistinguishable from the theoretical, thick magenta curves) and of the annual maxima (T_Y , colored continuous curves) for several values of ρ_1 ranging between 0.5 and 0.99 compared to the theoretical one (magenta curve). The independent case ($\rho_1 = 0$) is depicted as a reference.

439 theoretical value for any ρ_1 and for threshold values z up to that represented
 440 in the figure. The latter result, which numerically demonstrates the theoret-
 441 ical finding by Volpi et al. (2015), is illustrated in Figure 4 for the specific
 442 process characterized by the parameter combination $\sigma^2 = 5$ and $\rho_1 = 0.85$.

443 In Figure 3, the return period estimated by selecting the annual maxima,
 444 T_Y (colored curves) assumes values larger than the theoretical ones pertaining
 445 to the parent process or the independent case ($\rho_1 = 0$, black curve); this
 446 implies that the corresponding z -values are smaller. We also notice that for
 447 $\rho_1 \geq 0.9$ the annual maxima span over a wide range, covering values that
 448 are even smaller than the mean of the process ($\mu = 1$, vertical dashed line

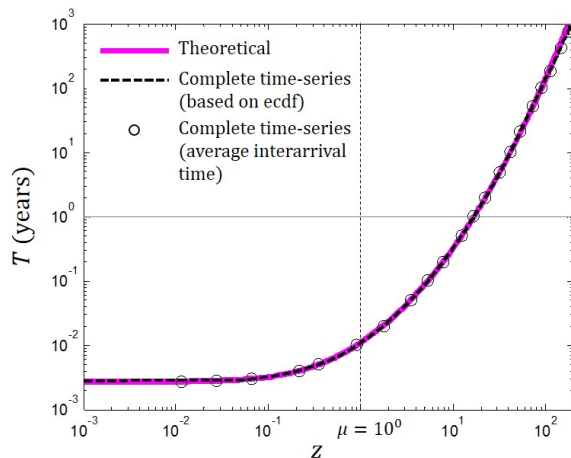


Figure 4: Lognormal AR(1) daily process with $\mu = 1$, variance $\sigma^2 = 5$ and lag-1 correlation coefficient $\rho_1 = 0.85$: empirical return period of the complete series T_Z , derived based on the ecdf (black dashed curves) and as the average of the empirical interarrival times (black circles) compared to the theoretical one (magenta curve).

449 in the figure). On average, the number of daily data exceeding the threshold
 450 $\{Z > \min(\mathbf{y})\}$ that are discarded by AM per year ranges in between 10 and
 451 350 values when ρ_1 increases from 0 to 0.99 respectively, for both $\sigma^2 = 1$ and
 452 $\sigma^2 = 5$.

453 Further, the return period estimate converges only for large values of
 454 z to the theoretical value, with a rate of convergence that depends on the
 455 persistence of the process. Hence, the larger is ρ_1 the slower is the rate of
 456 convergence of T_Y to T_Z . An important role is also played by the variance of
 457 the process; in the case $\sigma^2 = 5$ (Figure 3b) the convergence of the complete
 458 distribution to the theoretical one is even slower than in the case depicted
 459 in Figure 3a ($\sigma^2 = 1$). To give a quantitative measure of the deviation of
 460 the annual maxima estimate from the theoretical return period of the parent

461 process, the difference $D = T_Y - T_Z$ for $T_Z = 10$ years moves from 0.5 years
462 for $\rho_1 = 0$ (i.e. $\sim 5\%$, as in the theoretical independent case depicted in
463 Figure 1), to 150 (200) years, for $\rho_1 = 0.99$ and $\sigma^2 = 1$ (5).

464 Note that return period estimate based on CTA is compared here only
465 to that pertaining to annual maxima, but a similar comparison could be
466 made by considering the POT approach. It is expected that return period
467 estimates based on POT result in intermediate values between AM and CTA,
468 as a function of the threshold used to select peaks, but closer to AM estimates.
469 Note indeed, that CTA considers all the values exceeding the threshold z (see,
470 e.g., figure 2b), while POT uses only the independent maxima of clusters of
471 values exceeding z . The difference between the two approaches might be
472 relevant for practical purposes, as discussed later in Section 7.

473 *5.2. Effects of finite time-series length*

474 The situation depicted in Figure 3 is not met in practical applications,
475 when the limited length of the observed series significantly affects the re-
476 turn period estimation in terms of both accuracy and uncertainty. Generally
477 speaking, accuracy is expected to improve while uncertainty decreases when
478 increasing the length of the dataset of observations of a given process; note
479 that here the number of observations is not a direct measure of the amount of
480 information provided by data because of the correlation among the observed
481 values in complete time-series. Hence, we aim at comparing the overall ro-
482 bustness of return period estimates for small sample lengths, obtained by
483 CTA instead of the selected annual maxima (which are commonly assumed
484 to be independent). To investigate the latter issue, we repeated the above
485 analysis for different values of n within the range from 10 to 200 years for a

486 large number a synthetic time-series ($M = 10,000$).

487 Results obtained when assuming $\sigma^2 = 5$ and $\rho_1 = 0.6$ are depicted in
488 Figure 5; the figure summarizes the empirical probability estimates together
489 with their 95% uncertainty bounds derived by using both methods for some
490 values of n ranging in between 10 and 200 years. Note that results are pre-
491 sented here in terms of edf to avoid infinite values of return period that might
492 occur due to the use of Equation (1) when z is larger than the maximum ob-
493 served value in the dataset. If the edf in Equation (7) is normalized with
494 respect to $n + 1$ instead of n , we obtain the classic Weibull plotting position
495 formula (Makkonen, 2006); indeed, the latter is usually adopted to avoid in-
496 finite values of the estimated return period for the sample maximum. This
497 issue goes beyond the scope of this analysis, which is intended to discuss the
498 variability of return period estimate due to finite sample lengths; it will be
499 considered in future works together with the problem of model fitting.

500 Figure 5 shows how the AM estimate converges on the average only with
501 increasing T (moving from Figure 5a to c) to that of the complete time-series,
502 which overlaps the theoretical one for any sample length n . As expected, the
503 uncertainty bounds reduce with the sample size n . Uncertainty also reduces
504 as T increases; however, this unexpected behavior is a consequence of the
505 fact that probability is upper bounded to unity and due to the adoption
506 of the edf given in Eq. (7). Notwithstanding this, it is worth noting that
507 AM uncertainty bounds are narrower than those pertaining to CTA for any
508 value of T . This means that the selection of annual maxima results in an
509 undersampling effect, that manifests itself in terms of underestimation of the
510 exceeding probability of the parent process (i.e. overestimation of the non-

511 exceeding probability or of the return period as discussed in Figure 3), and
 512 of its sampling variability.

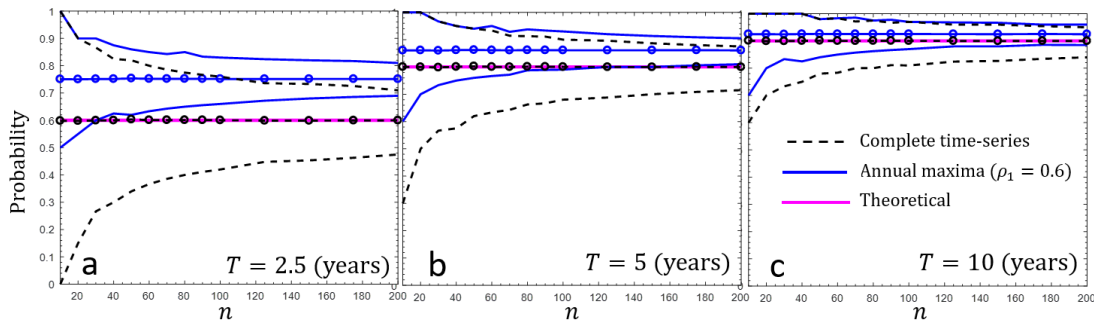


Figure 5: Lognormal AR(1) daily process with $\mu = 1$, variance $\sigma^2 = 5$ and lag-1 correlation coefficient $\rho_1 = 0.6$: CTA (black curves) and AM (blue curves) empirical probability for some values of the return period T as function of the time-series length, n ; the estimated probabilities are represented in terms of average values (dotted curve) and 95% uncertainty bounds. In each panel the theoretical probability is reported as a reference (magenta curve).

513 Correlation in time, which in this case is fully represented by the lag-1
 514 correlation coefficient ρ_1 , significantly affects the accuracy and the variability
 515 of the return period estimates obtained by both methods. Results for all
 516 values of ρ_1 considered in this illustrative example, are summarized in Figure
 517 6 for $T = 5$ years. Figure 6a and b depict the behavior of the average
 518 probability estimates based on AM (panel a) and CTA (panel b). It can
 519 be noticed that the *bias* resulting from AM is strongly enhanced by high
 520 values of ρ_1 . Conversely, the average probability estimates based on CTA
 521 are unbiased for any ρ_1 and n .

522 Further, the underestimation of the sampling variability which is observed
 523 in Figure 5 for $\rho_1 = 0.6$ when using AM, magnifies in the cases of processes

524 strongly correlated in time. To illustrate the latter issue we depict in Figure
525 6 also the coefficient of variation, C_V of the probability estimates for $T = 5$
526 years, computed by analyzing annual maxima (Figure 6c) and the complete
527 time-series (Figure 6d). The figure clearly shows that while C_V is comparable
528 for small values of ρ_1 , large differences arise when ρ_1 approaches to one.
529 While C_V of CTA estimate increases with the persistence of the process,
530 that of AM reduces; the latter behaviour is a consequence of the fact that
531 probability, that is upper bounded by one, is overestimated when analyzing
532 annual maxima.

533 A similar analysis could be performed in terms of quantiles, by investi-
534 gating how the order statistics of annual maxima and complete time-series
535 are influenced by the correlation structure of the process. However, slightly
536 different results (not shown) are obtained in terms of empirical return period
537 quantiles with respect to those obtained in terms of edf (as in Figure 6). In
538 fact, the probability distribution of the order statistics does not correspond on
539 average to the theoretical probability distribution of the underlying process
540 (David and Nagaraja, 2003); moreover, it is affected by the autocorrelation
541 structure of the process. Conversely, the edf expressed by Equation (7) is
542 an unbiased estimator regardless of the type and strength of the correlation
543 structure.

544 **6. Persistent and cyclo-stationary processes**

545 In order to provide some insights into the use of CTA in applications,
546 we analyze here a synthetic process which resembles the main statistical
547 properties of an hydrological observed series. Specifically, we analyze a daily

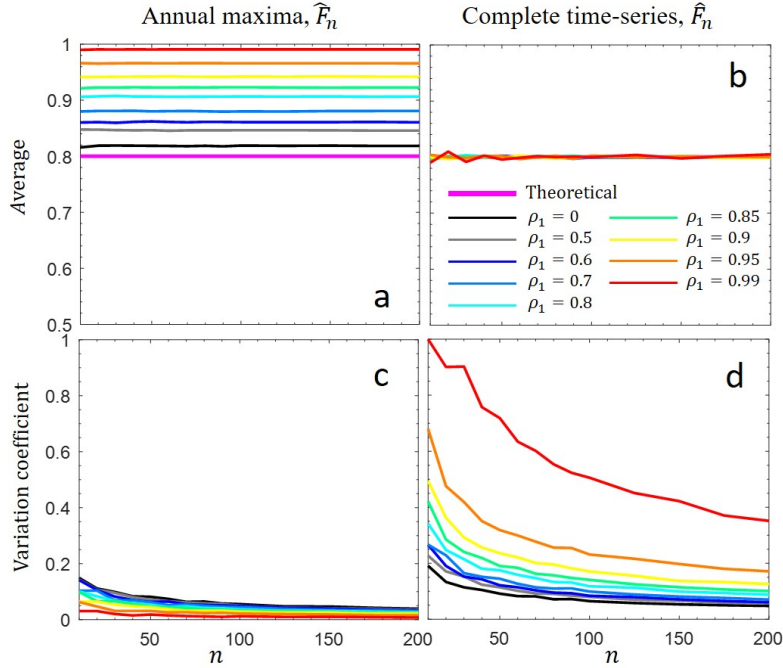


Figure 6: Lognormal AR(1) daily time-series, mean $\mu = 1$, variance $\sigma^2 = 5$, and lag-1 correlation coefficient $\rho_1 \in [0.50, 0.99]$: AM (left panels) and CTA (right panels) ecdf \hat{F}_n for $T = 5$ years; average values (upper panels) and variation coefficient (lower panels) as function of the time-series length, n . In panels (a) and (b) the theoretical probability (magenta line) and the independent case ($\rho_1 = 0$, black curve) are reported as reference.

548 discharge process that is characterized by non-normality, a strong seasonal
 549 pattern and by long-range persistence.

550 The type of analysis envisaged here requires a very long series of data.
 551 Hence, for the sake of illustration we consider a fractional autoregressive
 552 moving average process, FARMA(p, d, q), which models the Tiber River daily
 553 discharge time-series observed at Roma-Ripetta station. Observations cover
 554 a period of 54 years, from 1930 to 1983, but only the first 15 years were used
 555 to calibrate the linear parametric model; as an example, Figure 7a shows the

556 observed series (black line) for a time window of three years, from 1933 to
 557 1936. The seasonal pattern clearly emerges from the structure of the auto-
 558 correlation function, as depicted in Figure 7b (black line). The model was
 559 calibrated after normalizing the data (based on a log-normal transforma-
 560 tion) and removing seasonality; Figure 7 shows a subsample of the simulated
 561 series compared to the observed one (panel a) and the corresponding auto-
 562 correlation functions (panel b), thus highlighting the capability of the model
 563 in reproducing the complex behavior of the real world process. The reader
 564 is referred to Grimaldi (2004) for a detailed description of model structure,
 565 calibration and performance. We remark again that the model employed
 566 here is for the sole sake of illustration, and other general and more parsimo-
 567 nious methods could have been used to generate synthetic series from the
 568 observed process with any arbitrary autocorrelation structure, as discussed
 569 by Koutsoyiannis (2016), yet this goes beyond the scope of this work.

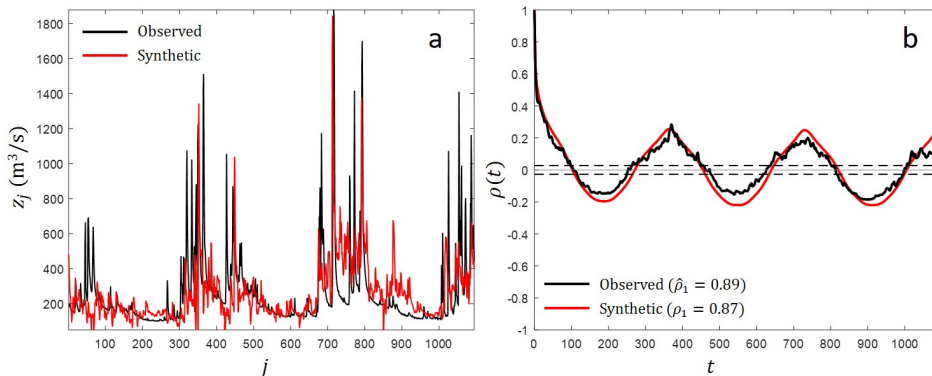


Figure 7: Synthetic series compared to that of the Tiber River (1930-1985) for a time window of three years: a) discrete-time daily discharge, and b) autocorrelation function, where dashed lines show 95% Gaussian confidence band.

570 We estimated the empirical return periods for both the annual maxima
 571 and the complete time-series by substituting in Eq. (3) the non-exceedance
 572 probability with the edf calculated using Equation (7), which gives the av-
 573 erage non-exceeding probability within the year (II). The empirical return
 574 periods are depicted in Figure 8a; since the synthetic series is very long, the
 575 figure depicts the theoretical difference between the two (unaffected by sam-
 576 ple length). AM significantly overestimates the return period of the complete
 577 time-series, if the latter is considered as a benchmark, for return periods up
 578 to 100 years. Here the bias of AM with respect to CTA ($D = T_Y - T_Z$)
 579 for $T_Z = 10$ years is equal to ~ 10 years, which means a 100% relative
 580 difference. The latter value is very close to that pertaining to the AR(1)
 581 lognormal process with similar value of ρ_1 and $\sigma^2 = 5$, discussed in Section
 582 5.2, although the variation coefficient here is smaller (about 0.7) with respect
 583 to that pertaining to the AR(1) process (about 2.3).

584 If a finite length sample is used to estimate return periods, the difference
 585 between AM and CTA might be enhanced. The effects of finite sample length
 586 for this cyclo-stationary, long-range persistent process is illustrated in panels
 587 b)-d) of Figure 8; in panels b) and c) the empirical return periods of two sam-
 588 ples of 54 years (equal to the length of the observed time-series) drawn from
 589 the FARMA calibrated model are compared to the corresponding *theoretical*
 590 ones (the full length samples). The empirical return period estimates for the
 591 event $\{Z > 1500 \text{ m}^3/s\}$ are depicted in panel (d); the average return period
 592 of annual maxima overestimates that pertaining to the complete time-series,
 593 also showing a smaller dispersion around its average value.

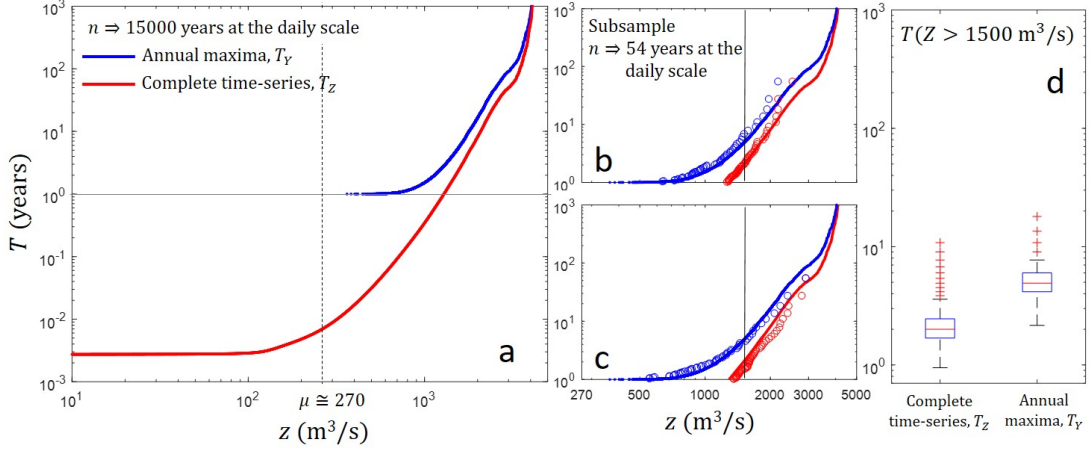


Figure 8: Synthetic daily process resembling the main statistics of the Tiber River daily discharge series (1930-1985) with mean $\mu = 270$ (m^3/s), standard deviation $\sigma = 181$ (m^3/s), lag-1 correlation coefficient $\rho_1 = 0.87$ and Hurst coefficient $H = 0.9$: empirical return period of the complete series (T_Z , red curve) and the annual maxima (T_Y , blue curve) considering the whole time-series (a) or analyzing a sub-sample of length equal to that of the observed series (b, c). Note that panels b) and c) focus on $z \geq \mu$ and $T \geq 1$ year. Panel d) shows the boxplot of the estimated return periods for the events $\{Z > 1500 \text{ m}^3/\text{s}\}$, based on 54 years sample length and both the methods.

594 7. Discussion on CTA application in real-world cases

595 It is important to note that the CTA approach to return period estimate
 596 considered here generally differs from the common methods used in hydrology
 597 (e.g. in flood frequency analysis), as explained in the following. CTA gives
 598 the return period of the event A defined as the exceedance of a threshold value
 599 (i.e. $Z_t > z$) at the temporal resolution $\Delta\tau$ at which the continuous process Z
 600 is sampled. In other words, it accounts for all the interarrival times between
 601 any successive values exceeding z at the $\Delta\tau$ scale, including those elapsing

602 between successive values of Z_t that remain above the threshold z (i.e. simply
603 equal to $\Delta\tau$, see figure 2); instead, in the conventional flood analysis the
604 above interarrival times are usually not considered. The significance of such
605 a return period estimate depends on the particular goal at hand and on the
606 temporal resolution $\Delta\tau$ that should be comparable (not much smaller) with
607 respect to the temporal scale that characterizes the natural phenomenon.
608 The temporal scale should not be confused with the characteristic scale of
609 the correlation structure (i.e. the integral scale), which itself depends on the
610 temporal resolution $\Delta\tau$.

611 For instance, in the example provided in Section 6, CTA accounts for
612 the consecutive exceedance of any threshold value of flow discharge at the
613 daily scale ($\Delta\tau = 1$ day), which is a small temporal scale with respect to the
614 average duration of a flood event. In this case, if the purpose of the analysis
615 is the assessment of the levee system, it might be not of interest to know if the
616 levee height is exceeded the day after once exceedance is already occurred
617 the day before. On the other hand, it might be important to account for
618 successive exceedances at the temporal resolution of the process $\Delta\tau$, e.g.
619 when we are evaluating the return period of daily rainfall to design a urban
620 drainage system against pluvial flooding and the critical duration of the
621 system (that maximizes the peak discharge) is approximately one day.

622 Further, as anticipated in the Introduction, for simplicity reasons we
623 based our analyses on a non-parametric estimation approach, i.e. using the
624 edf in Eq. (7) instead of a model fitted to the data. However, direct es-
625 timation of the statistics of a process is generally not possible merely from
626 the data and data alone do not enable extrapolation of estimates, as often

627 required for planning and design purposes (Koutsoyiannis, 2016). Thus, the
628 issue of fitting appropriate models in the context of the proposed approach
629 deserves further investigation. Although uncertainty is inherent in any sta-
630 tistical model, such uncertainty could be reduced by the utilization of all the
631 available information, as well as by judicious choices of model.

632 The use of CTA naturally implements seasonality handling in frequency
633 analysis. As recently discussed by Allamanno et al. (2011) (see also Ras-
634 mussen and Rosbjerg (1991) and Strupczewski et al. (2012)), disregarding
635 seasonality in hydro-climatic extreme value analysis, based on annual max-
636 ima or POT, leads to an overestimation of return period, which is less safe.
637 The problem is solved by taking into account the events that occurred in
638 all the seasons by fitting different distributions to the maxima in separate
639 seasons or months and mixing the seasonal distributions according to their
640 probability of occurrence (see, e.g. Mascaro, 2018) or by directly including
641 the seasonal rate of occurrence of the exceedance events in the POT ap-
642 proach as in Rasmussen and Rosbjerg (1991) and Allamanno et al. (2011).
643 CTA implements the former method, by considering for frequency analysis
644 all the observed values in each of the seasons; this could eventually require
645 the adoption of complex probability models (e.g. mixed models).

646 The adoption of complex probability models could also help handle the
647 heterogeneity due to the possible superposition of different physical processes
648 ruling the statistical behaviour of the random variable of interest (rainfall,
649 floods, etc.). In fact, the general understanding appears to be that low and
650 ordinary hydrological events could be dominated by a different process (see,
651 e.g. Merz and Blöschl, 2008), thus having little or no contribution to the

652 larger events. This might emerge from the edf of the annual maxima, by
653 manifesting a different statistical behavior for ordinary and extreme events;
654 this heterogeneity is expected to emerge more strongly when analyzing the
655 complete time-series that brings a larger number of values (also on the upper
656 tail of the probability distribution function) with respect to annual maxima.
657 In this regard, we believe that a priori physical knowledge about the un-
658 derlying processes, if available, could be included in the analysis to support
659 the assumption of complex mixed models for modeling (and extrapolation)
660 purposes (see, e.g. Calenda et al., 2009). In absence of additional knowl-
661 edge on the physical process, the heterogeneity assumption cannot be truly
662 tested; however, all events generally occur under the combination of numer-
663 ous factors, so that the probabilistic treatment of processes is by definition
664 a macroscopic approach that does not care about each of the specific factors
665 and reduces to fitting the most appropriate model to the empirical distribu-
666 tion of the complete data set (or to its part that is of interest for the specific
667 problem at hand).

668 Finally, we remark that CTA is based on the availability of a discrete-time
669 uninterrupted record of observations with an adequate temporal resolution,
670 which according to our analysis affects the results. In the case of few missing
671 data within long observational records, some (temporal or spatial) interpola-
672 tion techniques could be adopted to fill the gaps; in general, large gaps could
673 affect CTA more than AM or POT. The improvement of large datasets of
674 environmental observations is expected to favor CTA approach in the future.

675 8. Conclusions

676 A new approach for return period estimation is proposed, denoted as
677 Complete Time-series Analysis (CTA). This approach relies on the property
678 that the average interarrival time between successive events (e.g. $\{Z > z\}$)
679 is not affected by the correlation structure of the underlying process, regard-
680 less of the persistence of the process. This means that independence is not
681 a necessary condition when return period is defined as the average of the in-
682 terarrival time; this also implies that no data selection techniques are needed
683 to assure independence of the data for frequency analysis. Hence, once sta-
684 tionarity can be assumed, the return period can be computed by using the
685 classical equation of return period (the inverse of the exceedance probabili-
686 ty) starting from any kind of observational data, independent or correlated
687 in time, thereby potentially exploiting all the available information on the
688 parent process. This property is extended herein to include cyclo-stationary
689 processes, since hydrological and other geophysical processes typically man-
690 ifest a weaker form of stationarity at sub-annual scales connected to the
691 seasonal variability of the environmental phenomena.

692 We compare the proposed approach to the simple method of Annual Max-
693 ima (AM), typically adopted in extreme value analysis. Complete time-series
694 (observed in discrete time, at constant time intervals) and annual maxima are
695 inherently different processes, that give subtly different information on the
696 same underlying continuous process; specifically, CTA describes the marginal
697 behaviour of the whole parent process (sampled in discrete time at a given
698 temporal resolution), while AM describes the statistical behavior of its ex-
699 tremes.

700 The difference between CTA and AM are discussed herein by making
701 use of two illustrative examples. In both cases, we adopt for the sake of
702 illustration a non-parametric approach by using the empirical probability
703 distribution function of the annual maxima and the complete time-series as an
704 estimate of the marginal probability distribution function of the underlying
705 process. Some general conclusions can be drawn from the analyses, as listed
706 in the following.

- 707 • CTA results in an accurate estimate of return period of the parent
708 process for any intensity of the event (i.e. threshold value z) and,
709 on the average, for any sample length, regardless of the correlation
710 structure and the seasonality of the parent process, thus allowing to
711 investigate a wider range of return period values, not only the largest
712 extremes that are the focus of extreme value theory.
- 713 • AM leads to an overestimation of the return period (and an underes-
714 timation of its sampling variability) of the parent process for small to
715 moderate return period values, converging to CTA estimates for large
716 events. This behaviour, which is a consequence of data selection and is
717 well known in the literature for independent processes, is enhanced by
718 time-persistence of the underlying process; further, it is independent
719 on average from the sample length.
- 720 • Return period estimates provided by CTA are generally different with
721 respect to that pertaining to annual maxima because CTA considers
722 all the occurrences of the dangerous events within the observed record;
723 their significance depends on the particular goal at hand and on the

724 temporal resolution of the process, that should be comparable with re-
725 spect to the temporal scale that characterizes the natural phenomenon.

- 726 • CTA could be easily applied in the case of complex hydrological time-
727 series (such as that discussed here that reproduces the main features of
728 the Tiber River daily discharge as observed at Rome-Ripetta station),
729 which are typically characterized by non-normality, seasonality, long-
730 range persistence and, possibly, heterogeneity.

731 We found that the estimation of the return period using CTA could be
732 a convenient alternative to existing methods as function of the problem at
733 hand, for a few reasons. First, the method is easy to implement; it can
734 be employed for any sample length, without any data selection (e.g. events
735 selection in flood analysis). Moreover, CTA always results in more conser-
736 vative return period estimates, i.e. smaller estimated values with higher
737 uncertainty, by exploiting all the information content of the observed data,
738 i.e. low, ordinary and extreme discharge values that make-up the complete
739 time-series and fully describe the seasonal pattern.

740 The difference between CTA and AM tends toward zero as we look at
741 events that are more and more extreme simply because very large events are
742 expected to be always selected as annual maxima. Thus, CTA and extreme
743 value analysis are expected to give the same results in terms of the upper tail
744 of the distribution, thus supporting the adoption of extreme value analysis
745 when the interest is in large return period values. Note, however, that very
746 intense events typically pertain to the extrapolation range, where differences
747 among the methods could emerge when a probability distribution model is
748 fitted to the sample.

749 Hence, additional work is needed to fully understand advantages and
750 limitations of CTA in engineering practice. Since the main scope of this work
751 was to explore the potential advantages of the complete time-series approach
752 compared to traditional ones, we have not addressed the important issue
753 of the inference of the statistical distribution of the hydrological variable
754 of interest. Future work will investigate the problem of fitting appropriate
755 candidate models able to reproduce the complex, potentially heterogeneous
756 statistical behavior of complete time-series.

757 **Acknowledgments**

758 We acknowledge the Associate Editor, Geoff Pegram and two anonymous
759 Reviewers for thoughtful comments that helped improve the original version
760 of the manuscript. E.V. and A.F. acknowledge funding from the Italian Min-
761 istry of Education, University and Research (MIUR), in the frame of the
762 Departments of Excellence Initiative 2018-2022, attributed to the Depart-
763 ment of Engineering of Roma Tre University.

764 **Appendix A. Mean interarrival time of cyclo-stationary and per-** 765 **sistent processes**

766 The average interarrival time is obtained by substituting in the general
767 expression (1) the pmf given in Eq. (5). Note that the latter pmf depends on
768 time t ; hence, in the following we derive the expression of the return period

769 T conditional on t , i.e. $\frac{T}{\Delta\tau} = \mathbb{E}[X|t]$

$$\begin{aligned}
\frac{T}{\Delta\tau} &= \sum_{x=1}^{\infty} x \sum_{t'=t}^{t+\Pi-1} \frac{\Pr A_{t'}}{\sum_{\eta=0}^{\Pi-1} \Pr A_{\eta}} \frac{\Pr(A_{t'}, B_{t'+1}, B_{t'+2}, \dots, B_{t'+x-1}, A_{t'+x})}{\Pr A_{t'}} \\
&= \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \sum_{t'=t}^{t+\Pi-1} \sum_{x=1}^{\infty} x \Pr(A_{t'}, B_{t'+1}, B_{t'+2}, \dots, B_{t'+x-1}, A_{t'+x}) \\
770 &= \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \sum_{t'=t}^{t+\Pi-1} [1 \Pr(A_{t'}, A_{t'+1}) + 2 \Pr(A_{t'}, B_{t'+1}, A_{t'+2}) + \\
&\quad + 3 \Pr(A_{t'}, B_{t'+1}, B_{t'+2}, A_{t'+3}) + \dots]
\end{aligned} \tag{A.1}$$

771 By making use of the identity $\Pr(CA) = \Pr(C) - \Pr(CB)$, where B always
772 denotes the opposite event of A , we obtain (as in Volpi et al. (2015))

$$\begin{aligned}
\frac{T}{\Delta\tau} &= \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \sum_{t'=t}^{t+\Pi-1} [(\Pr A_{t'} - \Pr(A_{t'}, B_{t'+1})) + 2(\Pr(A_{t'}, B_{t'+1}) - \Pr(A_{t'}, B_{t'+1}, B_{t'+2})) + \dots] \\
773 &= \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \sum_{t'=t}^{t+\Pi-1} [\Pr A_{t'} + \Pr(A_{t'}, B_{t'+1}) + \Pr(A_{t'}, B_{t'+1}, B_{t'+2}) + \dots]
\end{aligned} \tag{A.2}$$

774 Using once more the same identity, we find

$$\begin{aligned}
\frac{T}{\Delta\tau} &= \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \sum_{t'=t}^{t+\Pi-1} [(1 - \Pr B_{t'}) + (\Pr B_{t'+1} - \Pr(B_{t'}, B_{t'+1})) + \\
&\quad + (\Pr(B_{t'+1}, B_{t'+2}) - \Pr(B_{t'}, B_{t'+1}, B_{t'+2})) + \dots] \\
&= \frac{1}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} \left[\sum_{t'=t}^{t+\Pi-1} 1 - \sum_{t'=t}^{t+\Pi-1} \Pr B_{t'} + \sum_{t'=t}^{t+\Pi-1} \Pr B_{t'+1} - \sum_{t'=t}^{t+\Pi-1} \Pr(B_{t'}, B_{t'+1}) + \right. \\
775 &\quad \left. + \sum_{t'=t}^{t+\Pi-1} \Pr(B_{t'+1}, B_{t'+2}) - \dots \right] \\
&= \frac{\Pi}{\sum_{r=t}^{t+\Pi-1} \Pr A_r} = \frac{1}{1 - \frac{1}{\Pi} \sum_{r=t}^{t+\Pi-1} \Pr B_r}
\end{aligned} \tag{A.3}$$

776 which simplifies in Eq.(6) thanks to the periodic property of the cyclo-
777 stationary process, such that $\sum_{t'=t}^{t+\Pi-1} \Pr B_{t'} = \sum_{t'=t}^{t+\Pi-1} \Pr B_{t'+1}$, $\sum_{t'=t}^{t+\Pi-1} \Pr(B_{t'}, B_{t'+1}) =$
778 $\sum_{t'=t}^{t+\Pi-1} \Pr(B_{t'+1}, B_{t'+2})$ and so on, and that marginal probability of non-
779 exceeding the threshold value z within any period $[t, t + \Pi - 1]$, i.e. $\overline{P_Z(z)} =$
780 $\frac{1}{1 - \frac{1}{\Pi} \sum_{r=t}^{t+\Pi-1} \Pr B_r}$, is independent on t .

References

- A. Schumann, Flood safety versus remaining risks - Options and limitations of probabilistic concepts in flood management, *Water Resources Management*, 31, 2017, 3131–3145.
- E. Volpi, A. Fiori, Hydraulic structures subject to bivariate hydrological loads: Return period, design, and risk assessment, *Water Resources Research*, 2014, 885–897.

- G.N. Alexander, Return period relationships, *Journal of Geophysical Research*, 64, 1959, 675–982.
- H. Rootzén, R. Katz, Design life level: Quantifying risk in a changing climate, *Water Resources Research*, 49, 2013, 5964–5972.
- J. Obeysekera, J. D. Salas, Frequency of Recurrent Extremes under Nonstationarity, *Journal of Hydrologic Engineering*, 21, 2016.
- L. K. Read, R. M. Vogel, Hazard function analysis for flood planning under nonstationarity, *Water Resources Research*, 2016.
- D. Koutsoyiannis, A. Montanari, Negligent killing of scientific concepts: the stationarity case, *Hydrological Sciences Journal*, 60, 2015, 2–22.
- F. Serinaldi, C. Kilsby, F. Lombardo, Untenable nonstationarity: An assessment of the fitness for purpose of trend tests in hydrology, *Advances in Water Resources*, 111, 2018, 132–155.
- F. Serinaldi, C. Kilsby, Understanding Persistence to Avoid Underestimation of Collective Flood Risk, *Water*, 8, 2016, 152.
- H. Tyralis, P. Dimitriadis, D. Koutsoyiannis, P. O’Connell, K. Tzouka, T. Iliopoulou, On the long-range dependence properties of annual precipitation using a global network of instrumental measurements, *Advances in Water Resources*, 2017.
- B. Fernández, J. D. Salas, Return period and risk of hydrologic events. II: Applications, *Journal of Hydrologic Engineering*, 4, 1999, 308–316.

- E. M. Douglas, R. M. Vogel, C. N. Kroll, Impact of streamflow persistence on hydrologic design, *Journal of Hydrologic Engineering*, 7, 2002, 220–227.
- A. Bunde, J. F. Eichner, S. Havlin, J. W. Kantelhardt, The effect of long-term correlations on the return periods of rare events, *Physica A: Statistical Mechanics and its Applications* 330, 2003, 1–7.
- J. F. Eichner, J. W. Kantelhardt, A. Bunde, S. Havlin, The Statistics of return intervals, maxima, and centennial events under the influence of long term correlation, in: *In Extremis*, Springer-Verlag, Berlin Heidelberg, 2011, pp. 3–43.
- E. Volpi, A. Fiori, S. Grimaldi, F. Lombardo, D. Koutsoyiannis, One hundred years of return period: Strengths and limitations, *Water Resources Research*, 2015, 1–16.
- M. Marani, M. Ignaccolo, A metastatistical approach to rainfall extremes, *Advances in Water Resources*, 79, 2015, 121–126.
- E. Zorzetto, G. Botter, M. Marani, On the emergence of rainfall extremes from ordinary events, *Geophysical Research Letters*, 2016, 8076–8082.
- V. T. Chow, D. R. Maidment, L. W. Mays, *Applied hydrology*, McGraw-Hill, New York, 1988.
- N. T. Kottegoda, R. Rosso, *Probability, statistics, and reliability for civil and environmental engineers*, McGraw-Hill, Milan, 1997.
- G. Salvadori, C. De Michele, N. Kottegoda, R. Rosso, *Extremes in Nature – An Approach Using Copulas*, Springer, New York, 2007.

- J. R. Stedinger, R. M. Vogel, E. Foufoula-Georgiou, Frequency analysis of extreme events, in: D. Maidment (Ed.), Handbook of Hydrology, McGraw-Hill, New York, 1993.
- D. Koutsoyiannis, Probability and statistics for geophysical processes, National Technical University of Athens, Athens, 2008.
- D. Koutsoyiannis, Generic and parsimonious stochastic modelling for hydrology and beyond, Hydrological Sciences Journal, 61, 2016, 225–244.
- D. Koutsoyiannis, Reconciling hydrology with engineering, Hydrology Research, 45, 2014, 1174–1183.
- A. Montanari, D. Koutsoyiannis, Modeling and mitigating natural hazards: Stationarity is immortal!, Water Resources Research, 50, 2014, 9748–9756.
- F. Serinaldi, C. Kilsby, Stationarity is undead: Uncertainty dominates the distribution of extremes, Advances in Water Resources, 77, 2015, 17–36.
- A. Luke, J. A. Vrugt, A. AghaKouchak, R. Matthew, B. F. Sanders, Predicting nonstationary flood frequencies: Evidence supports an updated stationarity thesis in the United States, Water Resources Research, 53, 2017, 5469–5494.
- A. N. Kolmogorov, Sulla determinazione empirica di una legge di distribuzione, Giornale dell’Istituto Italiano degli Attuari, 4, 1933, 83–91.
- D. Azriel, A. Schwartzman, The empirical distribution of a large number of correlated normal variables, Journal of the American Statistical Association, 110, 2015, 1217–1228.

- J. Dedecker, F. Merlevéde, The empirical distribution function for dependent variables: asymptotic and nonasymptotic results in L^p , *ESAIM: Probability and Statistics*, 11, 2007, 102–114.
- W.B. Wu, Oscillations of empirical distribution functions under dependence, *High Dimensional Probability*, Institute of Mathematical Statistics, 2006, 53-61.
- J. R. Benjamin, A. A. Cornell, *Probability, Statistics, and Decision for Civil Engineers*, McGraw-Hill, New York, 1970.
- S. Coles, *An introduction to statistical modeling of extreme values*, Springer, London, 2001.
- G. Calenda, C. Mancini, E. Volpi, Distribution of the extreme peak floods of the Tiber River from the XV century, *Advances in Water Resources*, 28, 2005, 615–625.
- R. W. Katz, M. B. Parlange, P. Naveau, Statistics of extremes in hydrology, *Advances in Water Resources*, 25, 2002, 1287–1304.
- D. Koutsoyiannis, Statistics of extremes and estimation of extreme rainfall: I. Theoretical investigation, *Hydrological Sciences Journal*, 49, 2004, 574–590.
- L. Makkonen, Plotting Positions in Extreme Value Analysis, *Journal of Applied Meteorology and Climatology*, 45, 2006, 335–340.
- H. A. David, N. H. Nagaraja, *Order statistics*, John Wiley & Sons, New York, 2003.

- S. Grimaldi, Linear parametric models applied to daily hydrological series, *Journal of Hydrologic Engineering*, 9, 2004, 383–391.
- P. Allamanno, F. Laio, P. Claps, Effects of disregarding seasonality on the distribution of hydrological extremes, *Hydrology and Earth System Sciences*, 15, 2011, 3207–3215.
- P. Rasmussen, D. Rosbjerg, Prediction uncertainty in seasonal partial duration series, *Water Resources Research*, 27, 1991, 2875–2883.
- W. Strupczewski, K. Kochanek, E. Bogdanowicz, I. Markiewicz, On seasonal approach to flood frequency modelling. Part I: Two-component distribution revisited, *Hydrological Processes*, 2012.
- G. Mascaro, On the distributions of annual and seasonal daily rainfall extremes in central Arizona and their spatial variability, *Journal of Hydrology*, 559, 2018.
- B. Merz, G. Blöschl, Flood frequency hydrology: 1. Temporal, spatial, and causal expansion of information, *Water Resources Research*, 44(8), 2008.
- G. Calenda, C. Mancini, E. Volpi, Selection of the probabilistic model of extreme floods: The case of the River Tiber in Rome, *Journal of Hydrology*, 371, 2009, 1–11.