

1 **Stochastic evaluation of simple pairing approaches**
2 **to reconstruct incomplete rainfall time series**

3 **Daniele Pedretti · Roger D. Beckie**

4
5 Received: date / Accepted: date

6 **Abstract** Two-station pairing approaches are routinely used to infill missing
7 information in incomplete rainfall databases. We evaluated the performance
8 of three simple methodologies to reconstruct incomplete time series in pres-
9 ence of variable nonlinear correlation between data pairs. Nonlinearity stems
10 from the statistics describing the marginal peak-over-threshold (POT) values
11 of rainfall events. A Monte Carlo analysis was developed to quantitatively as-
12 sess expected errors from the use of chronological pairing (CP) with linear
13 and nonlinear regression and frequency pairing (FP). CP is based on *a priori*
14 selection of regression functions, while FP is based on matching the probabili-
15 ty of non-exceedance of an event from one time series with the probability of
16 non-exceedance of a similar event from another time series. We adopted a Gen-
17 eralized Pareto (GP) model to describe POT events, and a t-copula algorithm
18 to generate reference nonlinearly correlated pairs of random temporal distri-
19 butions distributed according with the GP model. The results suggest that the
20 optimal methodology strongly depends on GP statistics. In general, CP seems
21 to provide the lowest errors when GP statistics were similar and correlation

D. Pedretti, R. Beckie
Earth, Ocean and Atmospheric Sciences University of British Columbia (UBC) 2207 Main
Mall, Vancouver BC V6T1Z4 E-mail: dpedretti@eos.ubc.ca

22 became linear; we found that a power-2 function performs well for the selected
23 statistics when the number of missing points is limited. FP outperforms the
24 other methods when POT statistics are different and variables are markedly
25 nonlinearly correlated. Ensemble-based results seem to be supported by the
26 analysis of observed precipitation at two real-world gauge stations.

27 **Keywords** Missing data · Rainfall · Chronological pairing · Frequency
28 pairing · Nonlinearity · Copulas

29 1 Introduction

30 The analysis of missing rainfall data has been long addressed in the litera-
31 ture (e.g. [12,27,30]). An increasing number of sophisticated approaches have
32 blossomed in the recent years in the attempt to obtain optimal methodolo-
33 gies to reconstruct missing data. Examples are kriging algorithms, artificial
34 intelligence or neural networks (e.g. [14,15,20,21]), which allow for multiple
35 stations to be simultaneously used to infill missing points at a specific location.
36 Yet, it seems that simpler approaches such as two-station chronological and
37 frequency pairing are still widely adopted among practitioners (e.g. [24]). The
38 popularity of these methods may stem from the limited number of available
39 and complete gauge stations, and the practicality of these algorithms, which
40 are often quickly programmable in a spreadsheet.

41 The most common and widely known approach is possibly chronologically
42 pairing (CP) with linear regression. It consists of regression-based estimation
43 of missing measurements at a primary location (showing missing data) using
44 observations at a nearby secondary, more complete location. Least-squares or
45 other best-fitting techniques are applied to obtain regression coefficients. A
46 schematic example of linear regression is depicted in Figure 1-1.

47 There are several potential weaknesses associated with this approach. We
48 highlight for instance that reconstructed data series from regression-based ap-
49 proaches often underestimate the variance of the original data series [13]. A
50 suite of methods, including the maintenance of variance extension (MOVE)
51 techniques, the Generalized MOVE (GMOVE), the modified Kendall-Theil
52 Robust Line (KTRL2) and the Robust Line of Organic Correlation (RLOC),
53 do not suffer from this problem (e.g. [11–13,32,17,18]). Another difficulty is
54 that variables may be nonlinearly correlated. This intrinsically renders a lin-
55 ear regression estimation error-prone. The paired data from rainfall time series
56 are also typically much more scattered (i.e. more noisy and poorly correlated)
57 than those shown in Figure 1-1 (e.g. [24]). In small mountainous catchments,
58 this can be also linked to topographic control of local convective cells, which
59 strongly affects precipitation even at short spatial scales [10].

60 We show in this work that the nonlinear correlation between variables
61 stems from the difference in the probabilistic distribution of rainfall events
62 recorded at each individual station. In general, rainfall distributions are pos-
63 itively skewed and display heavy-tailed distributions which are typically de-
64 scribed in terms of the probability of rainfall peaks exceeding a fixed threshold
65 (in short, peak-over-threshold or POT). Examples of POT distributions pro-
66 posed in the past are Generalized Extreme Values (GEV), Gamma or Gener-
67 alized Pareto (GP) [19,22,24,26,31]. Below the thresholds, no specific distri-
68 bution is usually observed.

69 To deal with nonlinearity, one may opt to embed nonlinear regression mod-
70 els into CP estimations. Identifying optimal nonlinear functions is often not
71 trivial, as a link between nonlinear regression functions and rainfall statistics
72 may be non trivial. Moreover, scattering may hide the presence of nonlinear
73 structures. These aspects are explored in detail in the first part of this work.
74 More complex approaches like variable transformations or data regularization

75 have been proposed in the past (e.g. [3]). However, it seems that the actual
76 implementation of these concepts is still not common in routine applications.

77 An appealing alternative methodology to effectively deal with nonlinearity
78 and that avoids the selection of a regression model is frequency pairing (FP).
79 The approach is derived from engineering applications such as prediction of
80 flooding return periods or rainfall depths at ungauged locations. Most of these
81 techniques were developed at the end of the 1980s (e.g. [1,23,29]) but only
82 recently applied to the estimation of missing rainfall (e.g. [24]). FP is based
83 on the probability of non-exceedance (P_E) of an event at the primary location
84 compared to P_E of another event at a secondary location. The methodology is
85 graphically explained in Figure 1-2. In each dataset, an event of given rainfall
86 intensity is uniquely associated with a P_E value; a missing datum at a specific
87 time at a primary location is assigned a value corresponding to P_E of the same
88 event measured at a secondary location. In the provided example (Figure 1-2),
89 station 2 records an event $r_2 = 2.5$ at a generic time; if station 1 is missing
90 an event at the same time, the latter is assigned a value $r_1 = 0.25$ since r_1
91 and r_2 match the same probability ($P_{E1} = P_{E2} = 0.6$). It is worth mentioning
92 that reconstructed datasets using FP preserve all the distributional moments
93 of the original distribution, including the variance.

94 There are also potential weaknesses associated with the FP methodology.
95 This approach relies on a direct relationship between probabilities at differ-
96 ent stations; as such, a first requirement may be to obtain reliable empirical
97 distributions from finite datasets, a task that can be complicated in case of
98 limited data. This is especially true when the data are power-law distributed
99 and the observations only cover a limited range and do not represent the full
100 variability of the distribution, specifically the tails (e.g. [7]).

101 In this work we addressed two relevant questions regarding two-station
102 pairing approaches. (a) What errors should we expect from simple approaches

103 to infill missing databases? (b) How reliable are such predictions? We at-
104 tempted to provide an answer by first developing a Monte-Carlo (MC) frame-
105 work, where the ability of CP and FP is tested against synthetic reference
106 time series. Benchmarking data are generated using a t-copula approach to
107 follow statistical distributions commonly associated with rainfall databases.
108 We assessed each method by evaluating the ensemble average of the difference
109 between reference and estimated missing values as a function of the number of
110 missing data to be estimated and for different statistical rainfall distributions.
111 In addition we evaluated the methods through two real-world case studies with
112 a similar geographical configuration and comparable daily rainfall statistics.

113 Our purpose was to provide a direct indication for practitioners aiming
114 to quantitatively evaluate the efficiency of unsophisticated two-station pairing
115 approaches when data pairs can be affected by nonlinearity and the number of
116 missing points is variable. We emphasize that the tested methodologies may
117 not be the best available algorithms to optimally reconstruct incomplete time
118 series. Indeed, we intentionally did not test more sophisticated approaches
119 (including multiple stations or modified CP methods using MOVE or similar
120 corrections). While an exhaustive study of all methods goes beyond the scope
121 of this analysis, we also note that the comparison of different infilling methods
122 have been already evaluated (e.g. [3,4,8,14,16,20,25]).

123 The paper is structured as follows. Section 2 introduces initial relevant
124 aspects of this analysis such as the presence of nonlinearity in the corre-
125 lation between data pairs and theoretical aspects regarding the behavior of
126 the GP model. Section 3 presents the methodologies to be evaluated; it first
127 addresses CP methods and then the FP methods. Section 4 describes the MC-
128 based analysis, including the generation of synthetic realizations based on the
129 t-copula approach. Section 5 presents the results for the ensemble-based anal-
130 ysis, emphasizing the difference between estimates made in presence of linear

131 correlation and those based on nonlinear correlations. Section 6 describes the
 132 application of the tested methods to two real-world case studies and exam-
 133 ines the results. The paper ends with the general conclusions drawn from this
 134 analysis.

135 **2 Nonlinearity between data pairs associated with POT statistics**

136 The role of different POT statistics to generate nonlinear correlations between
 137 data pairs is illustrated. Compiled reviews of POT statistics revealed that the
 138 Generalized Pareto (GP) model is able to fit daily precipitation observed in
 139 several experimental sites in the world [31]. While other studies indicated that
 140 alternative models can be applied to fit POT statistics [26], we adopted the
 141 GP as a working model to develop our analysis. This is done without loss
 142 of generality, as any power-law model can be applied to generate nonlinear
 143 correlations.

144 Understanding the behavior of the GP function and a few theoretical as-
 145 pects about convergence criteria related with subsampling errors is required
 146 before addressing the case of partially sampled distributions and the use of
 147 infilling methodologies.

148 **2.1 Behavior of the GP model**

149 The three-parameter GP model is part of the family of power-law parametric
 150 models (e.g. [2, 6, 28]). The probability density function for a rainfall event (r)
 151 can be written as

$$g(r|k, \sigma, \theta) = \left(\frac{1}{\sigma}\right) \left(1 + k \frac{(r - \theta)}{\sigma}\right)^{\left(-1 - \frac{1}{k}\right)}, \quad (1)$$

152 for $\theta < r$, when $k > 0$ or $\theta < r < \theta - \sigma/k$ when $k < 0$. The special condition
 153 $k = 0$ has not been explored in this work. According to common formalisms,
 154 k is termed the "shape" parameter, σ is the "scale" parameter and θ is the
 155 "location" parameter.

156 Figure 2a depicts the behavior of a few selected analytical CDFs of the GP
 157 model characterized by different statistics. We set $\theta = 0$, as this parameter
 158 does not affect the shape of resulting distributions. It can be observed that,
 159 for a constant σ , the total number of large or extreme values (i.e. the tail
 160 of the distribution) is uniquely controlled by k , while for a constant k the
 161 tailing depends on σ . When both parameters change, tailing is controlled by
 162 the combination of both k and σ . For instance, tailing is more pronounced
 163 (i.e. there is more probability in the tail of the distributions) for the pair
 164 $k = 0.01, \sigma = 5$ than for the pair $k = 0.5, \sigma = 1$.

165 Figure 2b illustrates the expected errors related with the statistical sub-
 166 sampling. This is calculated as the ensemble mean of root mean square errors
 167 between the analytical solution of the GP models (G) using predefined pa-
 168 rameters k and σ and the estimated distributions (G') calculated from best-fit
 169 parameters from a sample of n randomly generated values following the same
 170 GP distributions. For one individual realization of size n , the error (ε') is
 171 calculated as

$$\varepsilon'(m) = \sqrt{\frac{1}{n} \left(\sum G(x|\Omega) - G[x|\Omega'(m)] \right)^2} \quad (2)$$

172 where Ω' is the vector containing the estimated parameters (k' and σ') and
 173 m is the random realization. The ensemble average of individual errors is
 174 calculated as

$$\varepsilon = \frac{1}{N_{MC}} \sum \varepsilon'(m) \quad (3)$$

175 where $m = 1, \dots, N_{MC}$. We set $N_{MC} = 10^4$. The GP random numbers were
 176 generated using the "gprnd.m" native function in the MATLAB environment.

177 The results indicate that ε decreases by a factor \sqrt{n} . More interestingly, the
 178 required number of sampling points for a specific amount of error is strongly
 179 controlled by the tailing. For low k , the number of points to obtain a similar ε
 180 is independent from σ ; however, for large k , the errors are larger for larger σ .
 181 For instance, to obtain an average error of $\varepsilon = 10^{-2}$, the sample size should
 182 be close to $n \approx 8 \times 10^1$ for $k = 0.5, \sigma = 1$ but increasing to $n \approx 2 \times 10^2$ for
 183 $k = 0.5, \sigma = 5$. Rearranging these results, Figure 2c shows the distribution of
 184 ε depending on k for different sample sizes and different values of σ . It can be
 185 observed that for negative k values the difference between $\sigma = 1$ and $\sigma = 5$
 186 are very small and $\varepsilon \rightarrow 0$ as $k \rightarrow -\infty$.

187 2.2 Correlation of synthetic rainfall pairs embedding GP distributions

188 We now illustrate how the specific combination of k and σ between stations
 189 results in a nonlinear correlation between data pairs. We assumed uncorre-
 190 lated uniformly (U) distributed rainfall events below threshold and correlated,
 191 GP distributed POT events among stations. Correlated multivariate variables
 192 are generated using a copula-based approach (e.g. [5,9]). Copula approaches
 193 provide a suitable method to account for any type of marginal distribution,
 194 including non-parametric ones.

195 The specific methodology used to generate correlated random variables uses
 196 t-copulas. Similar to Gaussian copulas, t-copulas also allow for any marginal
 197 distributions to be imposed, and involve the definition of a correlation matrix,
 198 which expresses the correlation coefficients and cross-correlation coefficients
 199 between pairs of variables (ρ_{ij}), and a degree of freedom (ν). For $i = j$ (auto-
 200 correlation), $\rho_{ij} = 1$. The larger ρ_{ij} for $i \neq j$, the better correlated becomes

201 a specific synthetic pair of variables. On the contrary, the lower ρ_{ij} , the more
202 'noisy' or scattered becomes the correlation. We set $\nu = 1$ for all simulations;
203 the specific value does not affect our conclusions (ν has to be positive to gen-
204 erate non-Gaussian multivariate distributions). In this work, we used the algo-
205 rithm suggested by Mathworks and easily programmable in a MATLAB envi-
206 ronment; details about this implementation can be found in software house web
207 page ([http://www.mathworks.com/help/stats/examples/simulating-dependent-](http://www.mathworks.com/help/stats/examples/simulating-dependent-random-variables-using-copulas.html)
208 [random-variables-using-copulas.html](http://www.mathworks.com/help/stats/examples/simulating-dependent-random-variables-using-copulas.html)).

209 In Figure 3, the resulting correlation between three hypothetical stations,
210 S1 ($k = 0.5$, $\sigma = 2.5$), S2 ($k = -0.1$, $\sigma = 5$) and S3 ($k = 0.1$, $\sigma = 1$) is shown.
211 These values are similar to those experimentally found in a reference mining
212 site in Peru (described below) and in the range of values observed in typical
213 gauge stations [31]. We set $\rho_{12} = 0.75$ and $\rho_{13} = 0.5$, respectively for pairs
214 S1-S2 and S1-S3, to illustrate the effects of data scattering on the resulting
215 correlation between data pairs.

216 Three important aspects deserve special attention. (a) Extreme rainfall
217 events are much larger for S1 than for S2 and S3, due to different shape factors,
218 while the range of S2 is larger than S1 due to the difference in scale factor
219 (consistent with the theoretical behavior shown in Figure 2). (b) Scattering
220 is quite small for the pair S1-S2 when compared with the pair S1-S3. This
221 is a consequence of the larger ρ_{12} than ρ_{13} , which determines that S1 has a
222 stronger relationship with S2 than with S3. (c) Using different GP distributions
223 to generate the correlated variables, data pairs show nonlinear correlation -
224 that is correlation that depends upon the value of the variable.

225 To examine the latter point in detail, we evaluated the coefficients of de-
226 termination (R^2) for each pair of variable, which helps to identify the shape of
227 this correlation. R^2 is estimated from the application of a linear or nonlinear
228 regression model to these data and defined from the sum of the squared dif-

229 fference between each fitted and observed value and their mean. For the pairs
 230 S1-S2, a coefficient of determination based on linear regression ($R_L^2 = 0.61$)
 231 provides much lower values than a coefficient of regression based on a power-2
 232 model ($R_P^2 = 0.90$). This occurs since the correlation between pairs closely
 233 resembles a power-law function with exponent 2. For the pairs S1-S3, scatter-
 234 ing is more intense and both linear and nonlinear regression coefficients drop
 235 ($R_L^2 \approx R_P^2 \approx 0.3$). This suggests that the higher noise partially masks the non-
 236 linearity between variables, despite the use of different POT statistics between
 237 S1 and S3. These empirical observations may be valid only for the adopted
 238 POT statistics. The combination of other statistics may generate a different
 239 shape in the correlation, which can be fit for instance by another power-law
 240 function.

241 Although not graphically shown, it is worth mentioning that the use of the
 242 same POT coefficients generates a linear correlation between variables. Also
 243 for this case, a lower ρ may generate random scattering and thus hide the
 244 presence of linearly correlated structures. Note that t-copulas implemented
 245 in this work do not generate autocorrelation among values in the individual
 246 time series. In this sense, we are implicitly not considering clustering of rain-
 247 fall events, but independent identically distributed events within each rainfall
 248 series. This assumption is made without any loss of generality, since the tested
 249 methodology does not require any autocorrelation or clustering in individual
 250 time series.

251 **3 Tested pairing methodologies**

252 **3.1 Chronological pairing**

253 The CP method is based on a regression model between two time series (Figure
 254 1). We considered two random vectors (r.v.) collecting random rainfall events

255 (r) at two different stations in the catchment. One station (primary) is missing
 256 data, which can be estimated using existing measurements from the secondary
 257 station. We assumed for simplicity that in both time series, rainfall events are
 258 simultaneously collected (except missing data) and that the time lag between
 259 each measurement is constant and equal to 1 day. Each r.v. (r_{s1} and r_{s2}) is
 260 composed by N_S samples. A general regression equation can be defined as

$$r_{s1}(\tau) = f(r_{s2}(\tau)|\Omega_c) + \epsilon \quad (4)$$

261 where $\tau = 1, \dots, N_S$, f is a generic linear or nonlinear regression function,
 262 ϵ is the vector containing the residuals, and Ω_c is the vector containing the
 263 regression model parameters. Least-square methods can be applied to obtain
 264 the best set of Ω_c minimizing ϵ .

265 Two specific CP models were used in this analysis. The first is the linear-
 266 regression-based approach (CP-L) with two regression coefficients (m and q)
 267 and in which the rainfall at the secondary station is estimated from the rainfall
 268 measurements at first station as

$$r_{s2}(\tau) = mr_{s1}(\tau) + q + \epsilon_L \quad (5)$$

269 where ϵ_L is the error associated with the linear model. The second is a power-
 270 law-based approach (CP-P), with two regression coefficients (a and b) and
 271 in which the rainfall at the secondary station is estimated from the rainfall
 272 measurements at first station as

$$r_{s2}(\tau) = ar_{s1}^b(\tau) + \epsilon_P \quad (6)$$

273 where ϵ_P is the error associated with the power-law model. The quality of
 274 these regression models can be also checked by calculating R_L^2 and R_P^2 .

275 3.2 FP approach

276 FP establishes a direct relationship between events displaying the same P_E
 277 (Figure 1). Similarly to CP, we assumed that the primary station is missing
 278 an event at a time τ , and that a second (for simplicity, simultaneous) event is
 279 observed at a secondary station. For each r.v., P_E can be defined as

$$P_E = Pr(r(\tau) < \max(r)), \quad (7)$$

280 where $\max(r)$ is the maximum recorded event of each r.v. After sorting the
 281 two r.v. in increasing order, P_E can be empirically found as

$$P_E(r) = \frac{1}{N_s} \sum I\{r(\tau) \leq r\}, \quad (8)$$

282 where

$$I\{r(\tau) \leq r\} = \begin{cases} 1 & \text{if } r(\tau) \leq r \\ 0 & \text{otherwise} \end{cases}.$$

283 Using (8), each event is associated with a unique P_E , specific for each
 284 r.v.; hence, the same P_E could define two events $(r_1(\tau), r_2(\tau))$ which could
 285 display the same or different rainfall intensity. The FP approach postulates
 286 that a missing event in the primary database can be estimated knowing the
 287 intensity of the rainfall event at the secondary stations, and knowing P_E at
 288 both locations. Thus, imposing $P_E(r_{s1}) = P_E(r_{s2})$, one can easily find r_{s1}
 289 from the observation of r_{s2} .

290 P_E can be modelled by a parametric function $G(r|\Omega_{Fi})$, where G is a
 291 cumulative distribution function (CDF) of any specific parametric model. If
 292 Ω_{Fi} are the parameters for that specific model for the i th station, then the
 293 missing value r_{s1} can be estimated in four steps using the FP approach:

- 294 1. fit $P_E(r_{s2})$ with a model $G_2 = G(r_{s2}|\Omega_{F2})$ to obtain the parameters de-
- 295 scribing the distribution of events at r_{s2} ;
- 296 2. find the analytical value of G_2 corresponding the event $r_{s2}(\tau)$;
- 297 3. fit $P_E(r_{s1})$ with a model $G_1 = G(r_{s1}|\Omega_{F1})$ to obtain the parameters de-
- 298 scribing the distribution of events at r_{s1} ;
- 299 4. invert G_1 to obtain $r_{s1}(\tau)$, imposing $P_E(r_{s1}(\tau)) = P_E(r_{s2}(\tau))$.

300 The last step can be done either analytically (knowing G_1^{-1}) or numerically,
 301 for instance using a minimization approach. Although not directly addressed
 302 in this work, the methodology could be also extended to non-parametric dis-
 303 tributions.

304 **4 Monte-Carlo analysis**

305 A general stochastic approach based on a classical MC analysis was developed
 306 to quantitatively test each approach. The testing framework is general, can be
 307 applied to any specific parametric model and is independent of the presence
 308 of clustered rainfall events. We chose four sets of representative parameters,
 309 generating four different scenarios. In scenario 1 and 2, the same statistics
 310 (k and σ) are applied to both stations, and the correlation between pairs
 311 becomes linear. In scenarios 3 and 4, statistics are different between stations
 312 and correlation becomes nonlinear. The specific values for each scenario are
 313 listed alongside with the results. The four analyzed scenarios cover a broad
 314 range of linear and nonlinear cases which give rise to complex correlations and
 315 require detailed analysis to be fully understood and explained.

316 We used the t-copula approach described above to create a total of $N_{MC} =$
 317 10^4 realizations for each combination of these parameters. For each set of
 318 parameters, we fixed the statistical parameters that define the distributions

of the primary and secondary rainfall data, the length of the time series (N_S)
and the number of missing data to be estimated (N_R). For each scenario,

1. we generated two correlated random time series of rainfalls, one to be used
as a primary time series, and the other as a secondary time series.
2. we selected a random reference time window (RTW) on the primary time
series, of size N_R ; these values were removed from the time series and used
as reference (but unknown) values to compare with estimated values;
3. we estimated the best-fitting regression parameters between the two distri-
butions using two linear and nonlinear regression models, with the RTW
values removed from both time series (for the same period of time); subse-
quently, we applied the CP to reconstruct the missing RTW of the primary
station using the regression model and the removed RTW data from the
secondary information;
4. we found best-fitting CDF of both distributions without RTW; subse-
quently, we applied the FP approach between paired stations to estimate
the RTW missing data.

Each realization was assumed to contain five cycles of two distinct hydro-
logical seasons (dry and wet). We assumed for simplicity that no rainfall events
occurs in the dry seasons, and focused on the reconstruction of missing rainfall
events concentrated in the wet seasons and that are characterized by the POT
distribution. We rounded the size of the database length to $N_S = 10^3$, and
assumed that $N_S/2$ events fell in the range $0 < r \leq \theta$, while the remaining
 $N_S/2$ events fell within the range $r > \theta$ (POT events) and followed GP models
with different k and σ . We set $\rho_{ij} = 0.5$, which led to a scattering between
pairs of variables which is similar to the typical observations (Figure 3).

After an initial sensitivity analysis using $\rho_{ij} = 1$, we observed that scenario
1 and 2 led to linear correlation between variables, while scenario 3 and 4 gave

346 rise to nonlinear correlation resembling a power-2 function (similar to the
 347 example in Figure 3a). This was especially true when $\sigma = 5$, for which we
 348 observed a larger R_P^2 than for $\sigma = 1$. Based on this empirical observation
 349 and recalling (6), we tested a CP-P approach imposing $b=2$. We did not test
 350 other b factors; the results from this specific selection are discussed in the next
 351 section.

352 From the MC analysis, an error is defined as the difference between the
 353 estimated value and the true but unknown reference values. For each m -th
 354 realization ($m = 1, \dots, N_{MC}$), the root mean square errors (RMSE) and sum
 355 of absolute errors (SE) were calculated respectively as

$$\text{RMSE}(m) = \sqrt{\frac{1}{N_R} \sum_{i=1}^{N_R} \left(\frac{r'_i(m) - r_i}{r_i} \right)^2} \quad (9)$$

$$\text{SE}(m) = \sum_{i=1}^{N_R} |r'_i(m) - r_i| \quad (10)$$

356 where r'_i and r_i refer to the estimated rainfall the reference values, respectively.
 357 Note that, because of the different normalization in RMSE and SE, the two
 358 errors measure two different quantities. In (9), the fluctuations $r'_i(m) - r_i$ are
 359 normalized by the local values r_i ; thus, RMSE provides the average fraction
 360 of error of the set of estimates, which includes both POT events and rainfall
 361 below threshold. On the other hand, SE is a useful measure for water balance
 362 calculations as it can be interpreted as an average cumulative error of rainfall
 363 depths due to the presence of missing values.

364 5 Results

365 We initially analyzed scenarios where the selection of POT statistics results
 366 in linear correlation between pairs. In scenario 1 ($k = 0.5$ and $\sigma = 1$; Figure

367 4a and 4e) we observed that the performance of each method depends on the
 368 number of missing points: for $N_R < 10$ all methods perform similarly.; for
 369 $10 < N_R < 10^2$ CP-L and CP-P tend to outperform FP; for $N_R > 10^2$, CP-P
 370 generates very large errors (RMSE $\gg 2$), while the errors for FP and CP-L
 371 remain below RMSE=2. The better performance of CP-L may be explained
 372 considering that the selection of the same statistics generating copula-based
 373 random time series leads to linear correlation between variables. In this case,
 374 a linear regression becomes more suitable than nonlinear models.

375 For small N_R , CP-P also works well. This likely occurs since power-law
 376 models can be approximated by linear regression (e.g. by a first-order Taylor
 377 expansion) for small errors. However, as N_R exceeds a critical value (which
 378 depends on the selected POT parameters) CP-P provides much larger errors
 379 compared to CP-L. In this case, the threshold was estimated through MC
 380 analysis as $N_R \approx 10^2$.

381 When N_R approaches the data record size, fewer data points are available
 382 for the estimation of the distribution statistics or regression coefficients and
 383 their estimates are very uncertain (RMSE $\rightarrow 10^3$). Interestingly, for CP-L and
 384 FP this effect is reached only when only a few points remain ($n \approx 20$). It is
 385 unlikely that in practice one would try to estimate long time series based on
 386 such limited data. Thus, subsampling should not be a problem for this type
 387 of approach, at least in the range of values tested here.

388 The departure of FP from CP-L is not dramatic for this first analyzed case.
 389 FP does not suffer from the problem of CP-P on large N_R , explaining why FP
 390 errors remain comparable with those of CP-L (ranging between $0.5 < \text{RMSE} <$
 391 1). This is also confirmed by the observed behaviour of SE (Figure 4e), which
 392 is very similar for both cases, being only slightly larger for FP than for CP-L.
 393 Consistent with the observation about RMSE, the cumulative errors for CP-P
 394 also tend to high values when $N_R > 10^2$. It can be observed that SE increases

395 linearly with the number of missing points, which is expected considering that
396 this error expresses a cumulative sum of errors that linearly depends on the
397 number of missing points.

398 In the scenario 2 ($k = 0.5$ and $\sigma = 5$; Figure 4b and 4f), we obtained
399 comparable errors with the previous scenario from both CP-L and CP-P for
400 $N_R < 10^2$. However, both methods tend to more markedly outperform FP,
401 which provides larger errors than in scenario 1. This is especially true observing
402 SE (Figure 4f) and may suggest that a higher σ provides larger errors for FP
403 than in the case of lower σ . Also for this case CP-L performs quite similarly
404 to CP-P in the limit of a reduced number of missing points; this could again
405 refer to the accuracy of a first-order approximation of the power-law model for
406 low number of missing points. Note that the departure of FP from the other
407 methods also occurs for a limited number of missing values ($N_R < 10$).

408 We then analyzed scenarios where the selection of POT statistics results in
409 a nonlinear correlation between pairs. In scenario 3, we set $k_1 = 0.5$, $k_2 = 0.1$
410 (suffixes indicate the number of the station) and $\sigma = 1$ for both stations. The
411 results are illustrated by Figure 4c and 4g. We observed that FP clearly outper-
412 forms CP-L and CP-P, especially when the number of missing data increases.
413 The poor ability of CP-L is explained considering the nonlinearity between
414 variables; the poor ability of CP-P should be linked to both the selection of ρ
415 and to the fact that a poor correlation hides the effects of nonlinearity associ-
416 ated with the selection of different POT coefficients (as in Figure 3). Resulting
417 correlations have a different power-2 structure than for higher ρ values. Being
418 independent from the selection of the regression model, FP seems therefore
419 recommended for this case.

420 In scenario 4 ($k_1 = 0.5$, $k_2 = 0.1$ and $\sigma = 5$; Figures 4d and 4h) a different
421 behavior was observed. For a wide range of missing data, both FP and CP-
422 P perform reasonably better (especially when compared to CP-L). Errors in

423 general are larger than in the third scenario, which is associated with a larger
424 σ (which requires more simulations to converge to similar errors than $\sigma = 1$,
425 as shown in Figure 2c). It should be noted, however, that the relative differ-
426 ence in cumulative errors (SE) for each individual method are slightly reduced
427 compared to the previously analyzed case with lower σ . Interestingly, CP-P
428 performs (somewhat) better than FP, although this is true only for an inter-
429 mediate range of total missing rainfall data. This is due to the better identified
430 power-2-like correlated variability of synthetic distributions for larger σ when
431 setting $\rho = 1$, and thus statistically more likely to occur for lower ρ values.

432 In light of these results, the performance of each method seem to be strictly
433 associated with the estimated POT statistics. CP methods may be preferred to
434 FP methods in the case of similar POT distributions between the two stations,
435 which creates linear correlation among data pairs. However, FP would not pro-
436 vide excessive additional errors compared with the other tested methodologies
437 for linear correlation. On the other hand, when the POT statistics between
438 stations are different and the variables are nonlinearly correlated, FP is pre-
439 ferred over the other methods. From our results, this is especially true when
440 the average error per measurement has to be estimated (measured through the
441 RMSE). While the selection of CP-L is implicitly not adequate, a selection of
442 a specific nonlinear regression mode seems hazardous. The power-2-based CP
443 performed reasonably well only when the specific combination of variables de-
444 termine nonlinear correlations with similar shapes; however, it may be not easy
445 to determine *a priori* which type of regression model performed better depend-
446 ing on the combination of parameters underlying the two POT distributions.
447 This is especially true when ρ is low and data are scattering.

448 For all tested cases, we found low coefficients of variation of the MC simu-
449 lations ($CV < 1$). This suggested that the variability among fluctuations is not
450 affecting the ensemble means; as such we did not directly address this aspect

451 in the present analysis but rather chose to focus on the difference in terms of
452 expected values. It should be noted that our conclusions aimed to emphasize
453 the effect of nonlinearity on the selection of simple pairing approaches. The re-
454 sults may not be universally applicable, as we did not analyze other scenarios
455 (for which the shape of correlated pairs can be different than a power-2-like
456 structure, or the number of missing points is larger than in our analysis).

457 The impact of the distribution of values below the threshold has not been
458 directly addressed in the MC analysis. Rainfall events below θ are typically
459 characterized by an uncorrelated noise $U = (0, \theta]$, and they should not affect
460 the impact of POT statistics on nonlinear correlation. Although not tested,
461 we argue that the presented MC framework could be also easily extended to
462 incorporate values below this threshold. We finally remark that the tested
463 CP methodologies are intentionally simple and not designed, for instance, to
464 deal with maintenance of variance or other statistical effects on the resulting
465 estimated outputs. Testing of more sophisticated techniques was beyond the
466 scope of this work.

467 **6 Analysis of Two Real-World Scenarios**

468 The MC analysis provided an estimate of the impact of POT statistics on the
469 selection of the pairing approach with the lowest ensemble-based error. We
470 evaluated if these results also apply to the analysis of real-world observations,
471 which can be considered as individual realizations from an ensemble of equally
472 plausible, alternative realizations characterized by similar POT statistics. To
473 this end, we performed a cross-validation to evaluate the behavior of CP and
474 FP approaches to reconstruct missing rainfall measurements in two different
475 sites. The first site is a small catchment at a mine site located in the Peruvian
476 Andes, while the second site is located nearby Vancouver, BC.

477 6.1 Real-World Scenario 1

478 The Peruvian site is located at about 4300 m and characterized by two dis-
479 tinct hydrological seasons (Figure 5a). The wet season generally takes place
480 from October to April and the dry season from May to September. In the
481 dry season, precipitation becomes sporadic and the total cumulative rain is
482 negligible compared to the wet season. Daily rainfall events (r) were recorded
483 from 2007 to 2013 at three gauge stations (S-B, S-Y and S-Q). The gauge
484 station S-B suffered from a partial mechanical failure for a period of about
485 60-70 days during the wet season October 2011 - April 2012. Several gaps also
486 exist since April 2012 and available data are not considered reliable since then
487 (Figure 5a). In the period 2007-2011, which is considered highly reliable, the
488 existing S-B database consists of 1710 daily rainfall values (out of 1825 possi-
489 ble values, thus leading to 115 missing data). 592 days experience no detected
490 rainfall (excluding missing data).

491 Station S-Y, located approximately 4 km from S-B, was active and con-
492 tinuously collecting rainfall data since 2007, including for the period where
493 gaps occurred at S-B. An additional rainfall station, S-Q, is located approxi-
494 mately 5 km away from S-B, although only collected reliable data from 2007
495 to 2008. Application of a linear fitting model between daily measurements at
496 S-Y and S-B yields $R_L^2 = 0.35$ (using information from 2007 to 2011) and of
497 $R_L^2 = 0.3$ for S-Q/S-B pair (using information from 2007 to 2008). Although
498 the catchment is quite small (about $5 - 10\text{km}^2$) the poor R_L^2 can be explained
499 by the site's irregular topography, which influences local rainfall convective
500 cells (S-Y is located at a 200m lower elevation than S-B, while S-Q is located
501 at a 300m higher elevation than S-B). Similar values were obtained using a
502 nonlinear regression model embedding a power-2 model.

503 Three statistical models (GEV, Gamma and GP) were fit to the experi-
504 mental data to obtain empirical P_E distributions using a maximum likelihood
505 algorithm and the embedded fitting functions in the MATLAB environment.
506 The results are reported in Figure 5b. Fitting performed to the full distribu-
507 tion of observations (including zeros) are shown by black lines with squares
508 (GP) and dots (GEV). As expected, they performed relatively poorly when
509 compared to the better fits for the POT distribution (rainfall events above a
510 selected threshold). A sensitivity analysis indicated that all models performed
511 reasonably well when setting a minimum threshold of $\theta = 2\text{mm/d}$. Note that
512 the Gamma model could not be applied to the full distribution, being defined
513 for non-zero positive values only.

514 The 679 rainfall values above the 2-mm threshold summed to a total rain-
515 fall of about 7560 mm. Of the remaining values, 439 observations were below
516 the threshold, summing to only about 478 mm of cumulative rain, illustrating
517 the dominant role of POT events (hereafter, a short name for rainfall values
518 above the threshold) in the total water balance at this site. No specific sta-
519 tistical distribution and autocorrelation between pairs were found for values
520 below threshold, which was assumed to be uncorrelated noise with uniform
521 distribution $U = [0, \theta]$.

522 For the POT fitting, GEV model (grey line) tends to underestimate the
523 probability of POT events above $P_E \approx 0.9$, while Gamma (blue line) and GP
524 (red line) models slightly underestimate values over $P_E \approx 0.95$. GP slightly
525 overestimates values above $P_E = 0.999$. The GP model provides a satisfactory
526 fit to the rainfall distributions up to a threshold of $P_E \approx 0.99$ for the three rain-
527 fall stations (see Supplementary Electronic Material). However, the reliability
528 of values above $P_E = 0.99$ in the experimental database is considered poor,
529 as the number of samples was too few to provide robust estimates for very
530 rare events. We were not interested here to describe the difference stemming

531 from the selection of different parametric models, and thus assumed that the
 532 GP model can satisfactorily describe the general distribution of POT events
 533 at the three sample locations. This was also done to compare the results from
 534 the MC analysis with those from the analysis of the experimental database.

535 For the three experimental stations, best-fitting analysis of the POT dis-
 536 tributions (imposing $\theta = 2$) resulted in the following best-fit GP statistics:

- 537 – S-B: $k=0.449$, $\sigma=2.356$;
- 538 – S-Y: $k=-0.082$; $\sigma=5.867$;
- 539 – S-Q: $k=0.100$; $\sigma=1.160$.

540 It is noted that these statistics are quite similar to those used in the example of
 541 Figure 3. However, the maximum rainfall events for the theoretical S1 in Figure
 542 3 are much larger ($r \approx 140$ mm/d) than those found that the experimental site
 543 at S-B ($r \approx 35$ mm/d). This occurs despite S1 and S-B showing quite similar
 544 k and σ values. This difference is associated with the effects of subsampling,
 545 since the population sizes for both experimental and synthetic distributions
 546 are not large enough to cover the full variability of possible events associated
 547 with these statistics. This may suggest that when CP and FP methods are
 548 applied to individual realizations, the efficiency of the tested methods could be
 549 different than what is inferred from the analysis of ensemble-averaged results
 550 from MC simulations.

551 We focused on reconstructing the values exceeding $\theta = 2$ mm for the S-B
 552 station observed data for the period 2007-2011 using S-Y as secondary informa-
 553 tion. The results from the direct application of FP and CP-L to the observed
 554 database are first analyzed from the visual inspection of the scatter plot shown
 555 in Figure 5c. It can be observed that the distribution of CP-L estimated values
 556 tends to follow the best-fit line (with equation $r_{S-B} = 0.61r_{S-Y} + 5.9$), rather
 557 than the 1 : 1 distribution line. This issue gives rise to marked scattering

558 between estimated and observed values. FP is independent of the regression
559 model, and seems to be more able to reconstruct values between $r = \theta$ and
560 $r \approx 5\text{-}6$ mm/d, being closer to the 1:1 line, although with a wide scattering in
561 the estimation.

562 A more quantitative estimate of the errors between observed and recon-
563 structed data was calculated in the form of SE. Using a cross-correlation ap-
564 proach, we removed a RTW of different size from existing observations and
565 estimated the departure of estimated values from observation. The procedure
566 was then repeated by changing the initial location of the RTW, for a number
567 of times equals to the number of available time slots. The errors from each
568 RTW were averaged over the number of repetitions, to provide an average
569 cumulative error similar to (10). Figure 5d displays the calculated SE against
570 N_R . The results clearly show that FP outperforms CP-L, especially for larger
571 gaps. This is consistent with our previous MC-based conclusions: since the
572 POT statistics change between experimental stations S-B and S-Y, correla-
573 tion between pairs becomes nonlinear and FP approaches are more suitable
574 for the reconstruction of missing points. This occurs even if the GP model
575 poorly estimated events above $P_E > 0.95$. This observation may suggest that
576 FP still performs better than CP methods even if distribution tails are not
577 correctly fitted above very high thresholds, as usually occurs in practice for
578 power-law models [7].

579 We did not report the results from CP-P, as this method performed remark-
580 ably worse than FP and CP-L. Indeed, we found $\text{RMSE} > 10^4$ when $N_R > 50$,
581 and very large SE values from a limited number of missing values. The poor
582 performance of CP-P can be ascribed to the wrong selection of the underlying
583 nonlinear regression models, reflecting the difficulty in the adequate selection
584 of an *a priori* nonlinear regression model for CP approaches. It is possible that
585 for this specific combination of POT statistics the shape of nonlinear corre-

586 lation is different than a power-2 model. We did not test if CP-P provides a
587 better estimate than FP using an alternative station (based on another com-
588 bination POT statistics and thus another type of nonlinear correlation); the
589 information from the other available secondary station (S-Q) was too limited
590 and incomplete to exhaustively corroborate this hypothesis.

591 6.2 Real-World Scenario 2

592 Two other gauge stations nearby Vancouver (BC) were used to further val-
593 idate the MC-based analysis and assess the actual efficiency of the CP and
594 FP methods against a much longer real-world data set. The Vancouver data
595 included complete daily rainfall values covering a period of 20 years (1995-
596 2015) and can be freely accessed through the Environment Canada webpage
597 (<http://climate.weather.gc.ca/>). The first analyzed station (Vancouver Har-
598 bour, VH) is located nearby the sea side, while the second station is located
599 a few km away from VH and in the proximity of Vancouver Grouse mountain
600 (VG), at an elevation of approximately 1500m.

601 Likely due to the distinct orographical configuration, the statistics char-
602 acterizing the data pairs are significantly different. The correlation between
603 the two datasets is $\rho = 0.55$ (which is similar to the Peruvian data set, where
604 the elevation difference is smaller). At VG, the number of records exceeding
605 a threshold of $\theta=1\text{mm}$ is 2136, which comprised about 30% of the total data
606 series length. The total rainfall volume of associated with POT events is about
607 71% of total precipitation during the analyzed time. For VH, the POT events
608 were 21% of the total rainfall events, summing up to 55% of the cumulative
609 rainfall collected over the analyzed period.

610 The results from POT fitting are graphically reported in the Supplementary
611 Online Material. We observed that, for VH, the GP was satisfactorily fitting

612 the POT values up to $P_E = 0.999$, after which the model slightly overestimated
613 the observed values. For VG, the model was accurate up to $P_E = 0.99$, after
614 which it overestimated the observations. Consistent with the Peruvian-based
615 scenario, the Gamma and GEV models provided less accurate fits than the GP
616 model. Based on these observations and to directly compare this scenario with
617 the previous one, we used the GP as a working model to fit POT values and
618 reconstruct the missing values using the FP approach. The resulting GP best-
619 fit parameters for VH were $k=0.07$ and $\sigma=10.17$, while for VG we obtained
620 $k=0.01$ and $\sigma = 20.25$.

621 We evaluated the efficiency of each infilling methodology through a stochastic
622 cross-validation approach similar to that used in the Peruvian scenario. In
623 particular, we used the VG dataset to reconstruct the POT rainfall events at
624 VH. The resulting SE for a number of missing values up to $N_R=1000$ is shown
625 in Figure 6. The plot shows that FP clearly outperforms the CP-L, consistent
626 with the analysis of the Peruvian site and in line with the stochastic infer-
627 ence from the MC analysis. Indeed, the strong difference in POT statistics
628 determines a clear nonlinearity between data pairs, confirming thus that FP
629 is more suitable than the CP-L for infilling missing data. Once again, the re-
630 sults from the CP-P methods are not reported, strongly overestimating the
631 errors compared to the other tested methodologies, possibly because the non-
632 linear structure between data pairs significantly departs from a power-2-like
633 behavior.

634 Conclusions

635 Missing rainfall values are routinely estimated using simple approaches such
636 as two-station chronological pairing (CP) and frequency pairing (FP). We
637 developed a statistical framework to quantitatively evaluate expected errors

638 and reliability of predictions from the application of simple versions of these
639 methods in the presence of nonlinear correlation between rainfall events and
640 with different number of missing points.

641 For CP approaches, we analyzed one formulation with linear regression
642 model (CP-L) and one with a nonlinear power-2 regression model (CP-P),
643 which seemed to resemble the shape of nonlinear correlation between data
644 pairs when scattering is low and POT statistics different. Synthetic realiza-
645 tions were generated with the aid of a copula-based algorithm, in order to
646 specify a correlation between variables characterized by non-Gaussian (Gen-
647 eralized Pareto, GP) marginal distributions. In total, we analyzed four syn-
648 thetic scenarios and two real world scenarios, each characterized by different
649 GP statistics and different number of missing points (N_R).

650 The Monte Carlo results showed that when the POT statistics were iden-
651 tical, CP-L and CP-P outperform the FP. This behavior is consistent with
652 the fact that similar statistics give rise to linear correlations between vari-
653 ables. This is especially true for the estimation of a reduced number of values
654 ($N_R < 10^2$, for the specific selection of values used in this work). Despite be-
655 ing based on a nonlinear regression function, CP-P did not provide excessive
656 additional errors compared with CP-L; this may occur because a power-law
657 model can be approximated by a linear function in presence of low expected
658 variability in data fluctuation.

659 When the POT statistics were different by pairs and variables were non-
660 linearly correlated, FP more markedly outperformed the other methods. This
661 behavior is justified since the use of CP-L approaches is intrinsically not appro-
662 priated when variables are nonlinearly correlated, while CP-P performs well
663 only when data resemble the power-law structure. Due to the high scatter-
664 ing, however, this occurs only for specific selections of POT parameters (e.g.
665 increasing σ in our tested scenarios). Our exercise suggests however that iden-

666 tifying the optimal parameters *a priori* may be cumbersome and uncertain.
667 Indeed, the application of the three approaches to real-case scenarios seems
668 to confirm that the FP approach may provide more robust results than CP
669 methods, FP being independent of the selection of a regression model.

670 While we only addressed a limited number of scenarios, we generally con-
671 clude that the optimal selection of the simple pairing approach must follow an
672 initial statistical analysis of existing data, which may reveal nonlinearity be-
673 tween variables which can be hidden by poor ρ and high scattering. Only when
674 nonlinearity is properly analyzed, one may obtain reliable quantitative infor-
675 mation about expected errors in missing databases conditioned to site-specific
676 empirical POT distributions.

677 **Acknowledgements** We acknowledge the useful suggestions provided by the Associate
678 Editor and three anonymous reviewers, who significantly helped to improve the manuscript.
679 All the data and additional information used and cited in this paper can be provided by the
680 corresponding author at specific requests.

Fig. 1 Conceptual examples of chronologically-pairing (CP) and frequency-based pairing (FP) approaches. In CP (1), a linear regression model is applied to fit matching values in time series 1 (r_1) and time series 2 (r_2); m and q are two fitting parameters. In FP (2), each time series is considered as an independent random space; the same probability matches two different data values in the two time series. A missing value in r_2 is calculated using a measurement in r_1 and on the basis of the probability of that specific event in r_2 .

Fig. 2 Theoretical behaviour of the GP model. (a) Analytical CDF for six different pairs of parameters. (b) Rate of convergence of expected errors versus sampling population. (c) Rate of convergence of expected errors for different shape parameters and scale parameters. We set $\theta = 0$ in all cases.

Fig. 3 Illustrative individual realizations of synthetic correlated rainfall distributions, characterized by different statistics and correlation coefficients (ρ) and $\theta = 2$. (a) Scatter diagrams for the pairs S1-S2. (b) Scatter diagrams for the pairs S1-S3. (c) Empirical POT rainfall distributions.

Fig. 4 Results of the Monte Carlo analysis, expressing the distribution over time of the mean root mean square errors (RMSE, figure a,b,c,d), and sum of cumulative errors (SE, figures e,f,g,h,) for different pairs of statistics.

Fig. 5 (a) Rainfall timeseries at two stations in the experimental site in Peru (S-B and S-Y) between 2007 and 2013; (b) empirical rainfall distribution of all recorded values, including zeros and values below the threshold (FD) and empirical POT distribution at S-B, along with three different fitting models: Gamma, Generalized Pareto (GP) and Generalized Extreme Values (GEV); (c) scatter plots for the CP-L and FP infilling approaches; (d) cumulative errors (SE) for the two infilling approaches.

Fig. 6 Cumulative errors (SE) for the two infilling approaches for the Vancouver-based scenario.

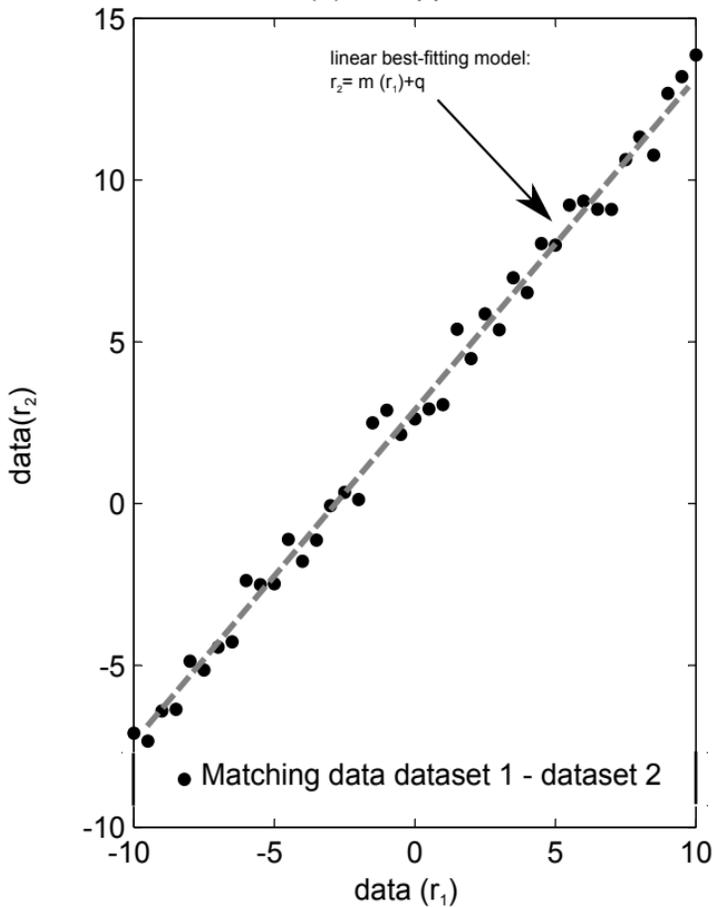
681 **References**

- 682 1. Arnell, N.W.: Unbiased estimation of flood risk with the GEV distribution. *Stochastic*
683 *Hydrology and Hydraulics* **2**(3), 201–212 (1988). DOI 10.1007/BF01550842
- 684 2. Arnold, B.C.: Pareto Distribution. In: *Encyclopedia of Statistical Sciences*. John Wiley
685 & Sons, Inc. (2004)
- 686 3. Beauchamp, J., Downing, D., Railsback, S.: Comparison of Regression and Time-Series
687 Methods for Synthesizing Missing Streamflow Records. *Journal of the American Water*
688 *Resources Association* **25**(5), 961–975 (1989). DOI 10.1111/j.1752-1688.1989.tb05410.x
- 689 4. Brdossy, A., Pegram, G.: Infilling missing precipitation records A comparison of a
690 new copula-based method with other techniques. *Journal of Hydrology* **519, Part A**,
691 1162–1170 (2014). DOI 10.1016/j.jhydrol.2014.08.025
- 692 5. Cantet, P., Arnaud, P.: Extreme rainfall analysis by a stochastic model: impact of the
693 copula choice on the sub-daily rainfall generation. *Stochastic Environmental Research*
694 *and Risk Assessment* **28**(6), 1479–1492 (2014). DOI 10.1007/s00477-014-0852-0
- 695 6. Castillo, E., Hadi, A.S.: Fitting the Generalized Pareto Distribution to Data. *Jour-*
696 *nal of the American Statistical Association* **92**(440), 1609–1620 (1997). DOI
697 10.1080/01621459.1997.10473683
- 698 7. Clauset, A., Shalizi, C.R., Newman, M.E.: Power-law distribution in empirical data.
699 *SIAM Review* **51**, 661–703 (2007). DOI 10.1137/070710111
- 700 8. Coulibaly, P., Evora, N.D.: Comparison of neural network methods for infilling miss-
701 ing daily weather records. *Journal of Hydrology* **341**(12), 27–41 (2007). DOI
702 10.1016/j.jhydrol.2007.04.020
- 703 9. Demarta, S., McNeil, A.J.: The t Copula and Related Copulas. *International Statistical*
704 *Review* **73**(1), 111–129 (2005). DOI 10.1111/j.1751-5823.2005.tb00254.x
- 705 10. Goovaerts, P.: Geostatistical approaches for incorporating elevation into the spatial
706 interpolation of rainfall. *Journal of Hydrology* **228**(12), 113–129 (2000). DOI
707 10.1016/S0022-1694(00)00144-X
- 708 11. Grygier, J.C., Stedinger, J.R., Yin, H.B.: A generalized maintenance of variance ex-
709 tension procedure for extending correlated series. *Water Resources Research* **25**(3),
710 345–349 (1989). DOI 10.1029/WR025i003p00345
- 711 12. Hirsch, R.M.: An evaluation of some record reconstruction techniques. *Water Resources*
712 *Research* **15**(6), 1781–1790 (1979). DOI 10.1029/WR015i006p01781
- 713 13. Hirsch, R.M.: A comparison of four streamflow record extension techniques. *Water*
714 *Resources Research* **18**(4), 1081–1088 (1982). DOI 10.1029/WR018i004p01081

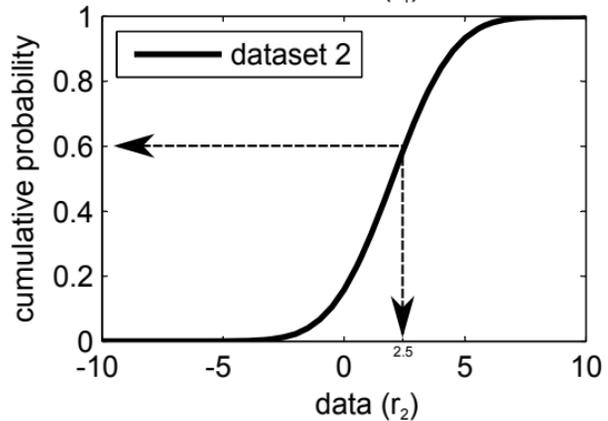
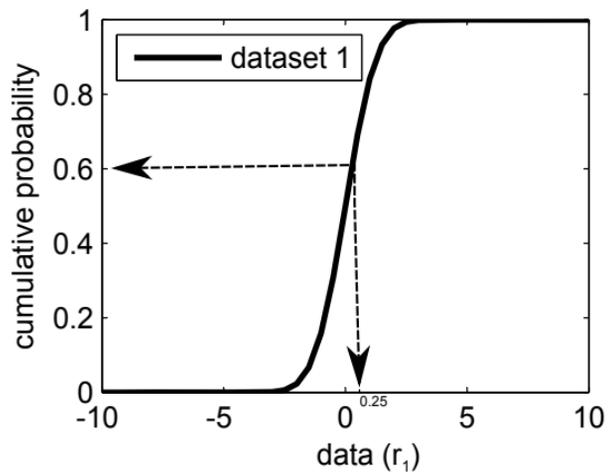
- 715 14. Kajornrit, J., Wong, K.W., Fung, C.C.: Estimation of Missing Precipitation Records
716 Using Modular Artificial Neural Networks. In: T. Huang, Z. Zeng, C. Li, C.S. Leung
717 (eds.) *Neural Information Processing*, no. 7666 in *Lecture Notes in Computer Science*,
718 pp. 52–59. Springer Berlin Heidelberg (2012)
- 719 15. Kashani, M.H., Dinpashoh, Y.: Evaluation of efficiency of different estimation methods
720 for missing climatological data. *Stochastic Environmental Research and Risk Assess-*
721 *ment* **26**(1), 59–71 (2011). DOI 10.1007/s00477-011-0536-y
- 722 16. Kebaili Bargaoui, Z., Chebbi, A.: Comparison of two kriging interpolation methods
723 applied to spatiotemporal rainfall. *Journal of Hydrology* **365**(12), 56–73 (2009). DOI
724 10.1016/j.jhydrol.2008.11.025
- 725 17. Khalil, B., Adamowski, J.: Record extension for short-gauged water quality parameters
726 using a newly proposed robust version of the line of organic correlation technique.
727 *Hydrol. Earth Syst. Sci.* **16**, 2253–2266 (2012). DOI 10.5194/hessd-9-4667-2012
- 728 18. Khalil, B., Ouarda T.B.M.J., and St-Hilaire, A. (2012). Comparison of record-extension
729 techniques for water quality variables, *Water Resources Management*. *Water Resources*
730 *Management* **26**(14), 4259–4280 (2012). DOI 10.5194/hessd-9-4667-2012
- 731 19. Kim, D., Olivera, F., Cho, H.: Effect of the inter-annual variability of rainfall statistics
732 on stochastically generated rainfall time series: part 1. Impact on peak and extreme
733 rainfall values. *Stochastic Environmental Research and Risk Assessment* **27**(7), 1601–
734 1610 (2013). DOI 10.1007/s00477-013-0696-z
- 735 20. Kim, J.W., Pachepsky, Y.A.: Reconstructing missing daily precipitation data using re-
736 gression trees and artificial neural networks for SWAT streamflow simulation. *Journal*
737 *of Hydrology* **394**(34), 305–314 (2010). DOI 10.1016/j.jhydrol.2010.09.005
- 738 21. Kim, T.W., Ahn, H.: Spatial rainfall model using a pattern classifier for estimating
739 missing daily rainfall data. *Stochastic Environmental Research and Risk Assessment*
740 **23**(3), 367–376 (2008). DOI 10.1007/s00477-008-0223-9
- 741 22. Li, Z., Li, C., Xu, Z., Zhou, X.: Frequency analysis of precipitation extremes in Heihe
742 River basin based on generalized Pareto distribution. *Stochastic Environmental Re-*
743 *search and Risk Assessment* **28**(7), 1709–1721 (2013). DOI 10.1007/s00477-013-0828-5
- 744 23. Lye, L.M.: Bayes estimate of the probability of exceedance of annual floods. *Stochastic*
745 *Hydrology and Hydraulics* **4**(1), 55–64 (1990). DOI 10.1007/BF01547732
- 746 24. Millar, R.: A statistical approach for deriving project design rainfall. pp. 273–276. *The*
747 *Australasian Institute of Mining and Metallurgy*, Melbourne, Australia (2013)
- 748 25. Musial, J.P., Verstraete, M.M., Gobron, N.: Comparing the effectiveness of recent al-
749 gorithms to fill and smooth incomplete and noisy time series. *Atmos. Chem. Phys.*

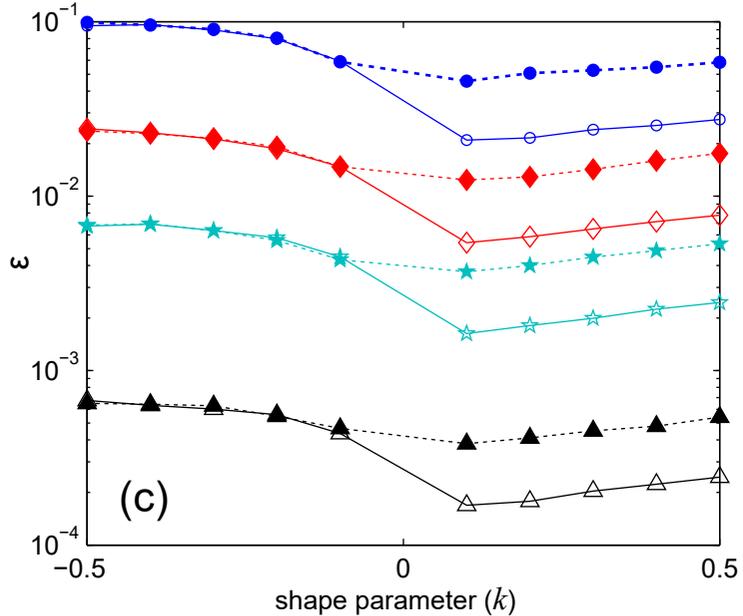
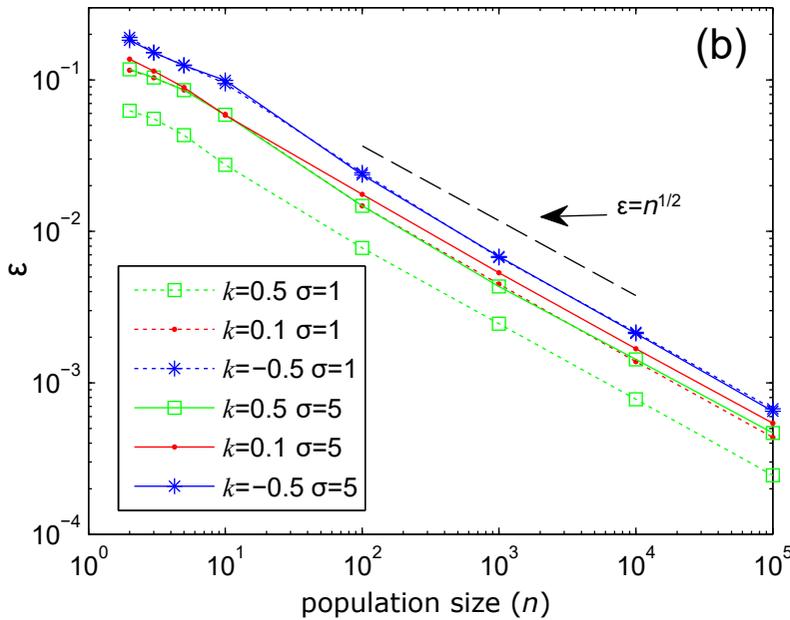
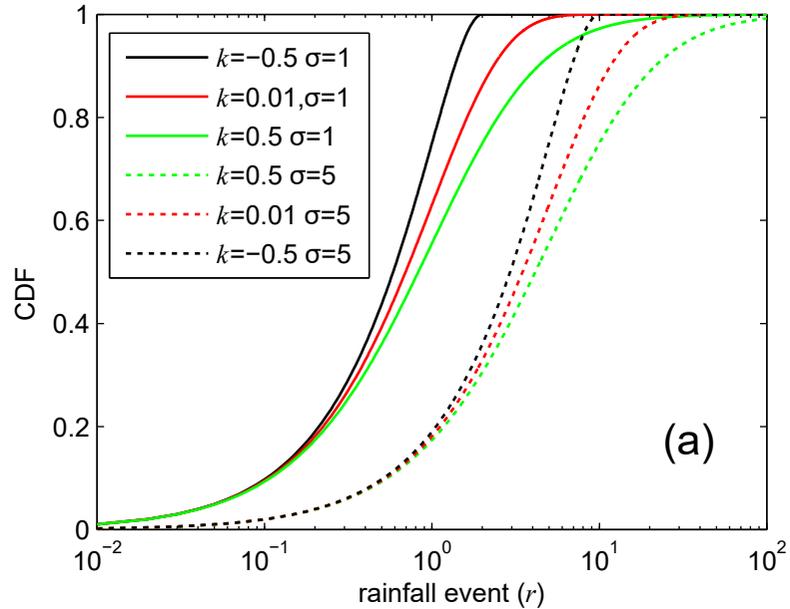
- 750 Discuss. **11**(5), 14,259–14,308 (2011). DOI 10.5194/acpd-11-14259-2011.
- 751 26. Papalexiou, S.M., Koutsoyiannis, D.: Battle of extreme value distributions: A global
752 survey on extreme daily rainfall. *Water Resources Research* **49**(1), 187–201 (2013).
753 DOI 10.1029/2012WR012557.
- 754 27. Paulhus, J.L.H., Kohler, M.A.: {Interpolation of missing precipitation records}. *Monthly*
755 *Weather Review* **80**(8), 129–133 (1952)
- 756 28. Pickands III, J.: Statistical Inference Using Extreme Order Statistics. *The Annals of*
757 *Statistics* **3**(1), 119–131 (1975)
- 758 29. Porporato, A., Ridolfi, L.: Influence of weak trends on exceedance probability. *Stochastic*
759 *Hydrology and Hydraulics* **12**(1), 1–14 (1998). DOI 10.1007/s004770050006
- 760 30. Rubin, D.B.: Inference and Missing Data. *Biometrika* **63**(3), 581–592 (1976). DOI
761 10.2307/2335739
- 762 31. Serinaldi, F., Kilsby, C.G.: Rainfall extremes: Toward reconciliation after the bat-
763 tle of distributions. *Water Resources Research* pp. n/a–n/a (2014). DOI
764 10.1002/2013WR014211
- 765 32. Vogel, R.M., Stedinger, J.R.: Minimum variance streamflow record augmenta-
766 tion procedures. *Water Resources Research* **21**(5), 715–723 (1985). DOI
767 10.1029/WR021i005p00715

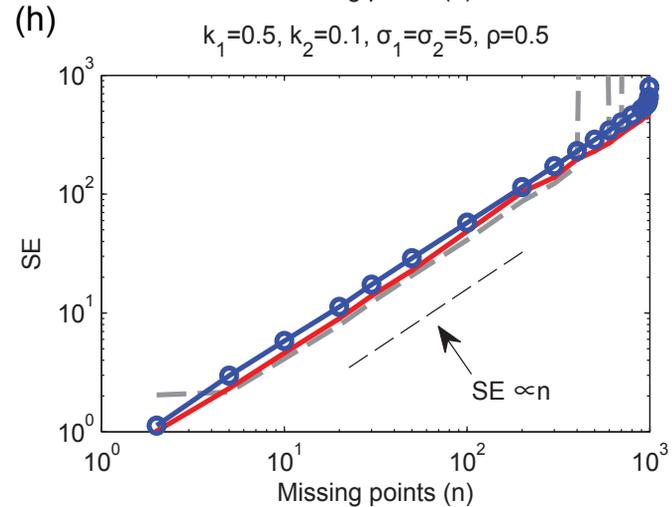
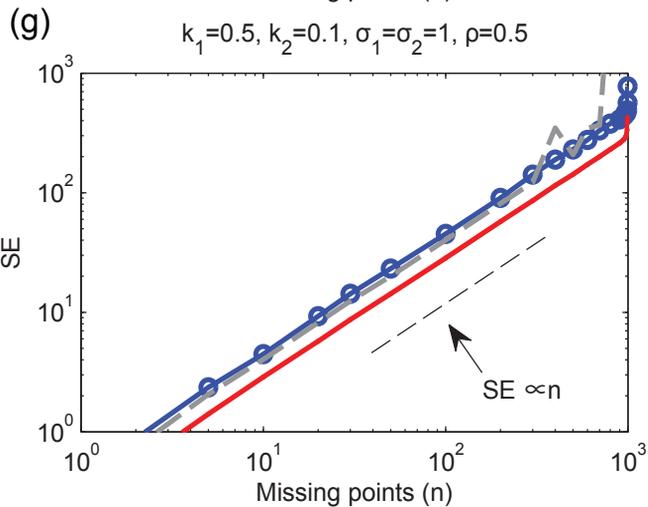
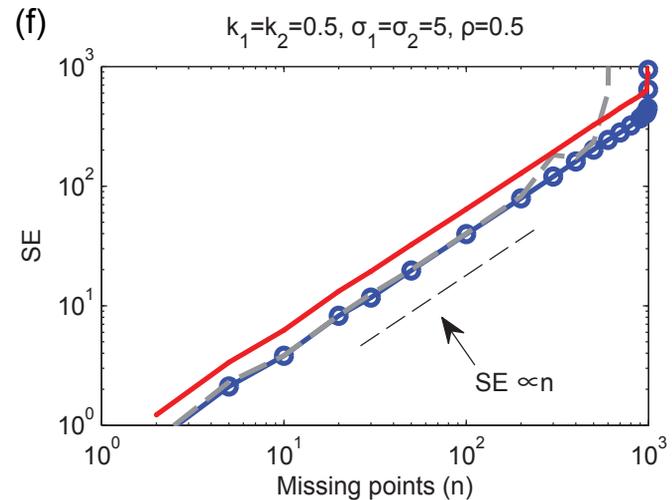
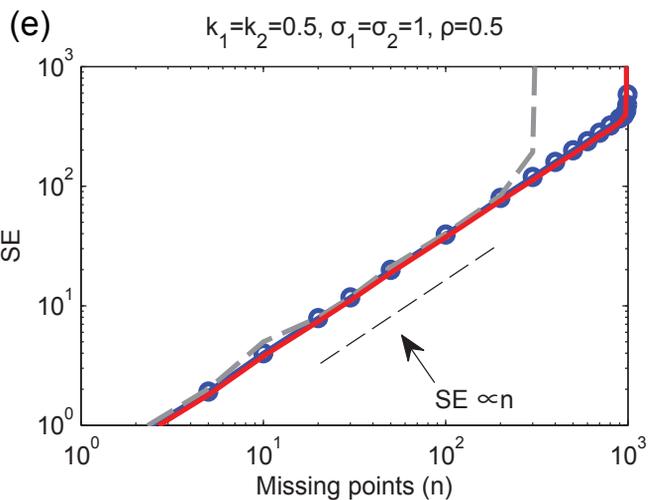
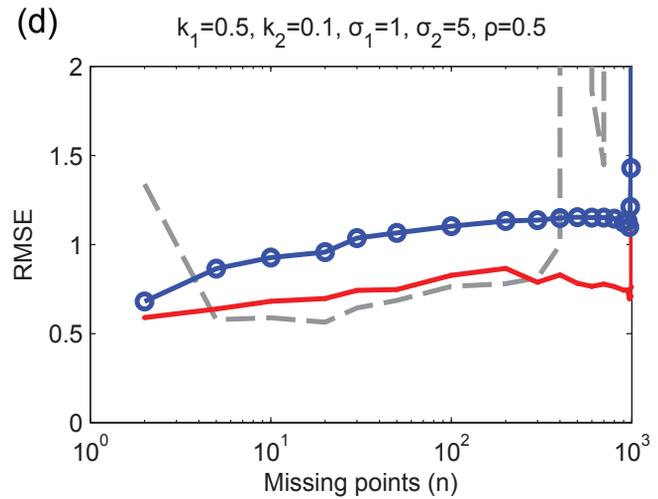
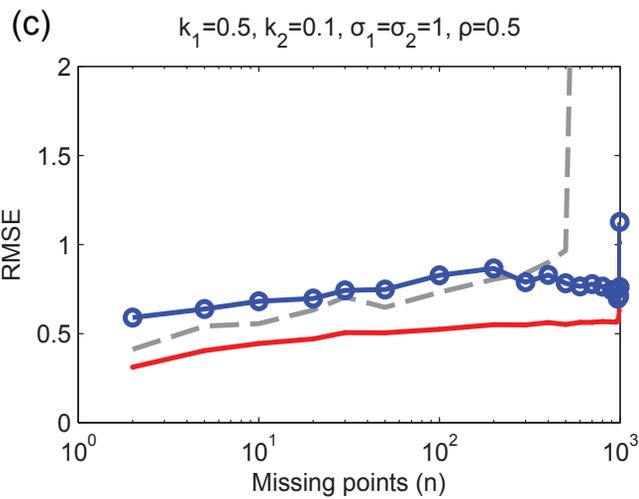
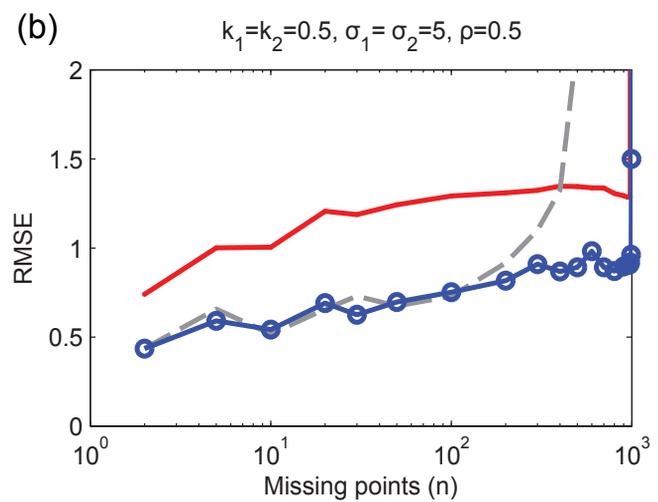
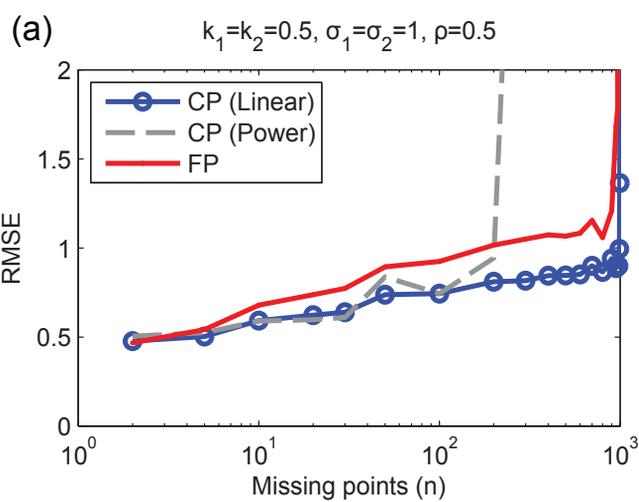
(1) CP approach

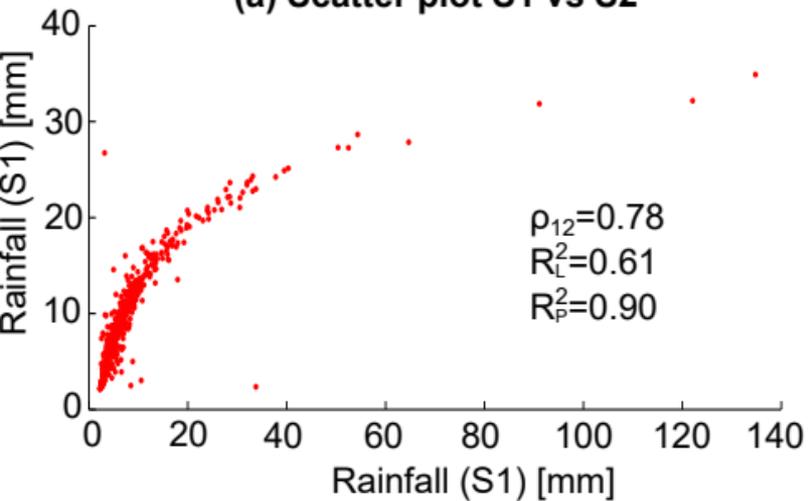
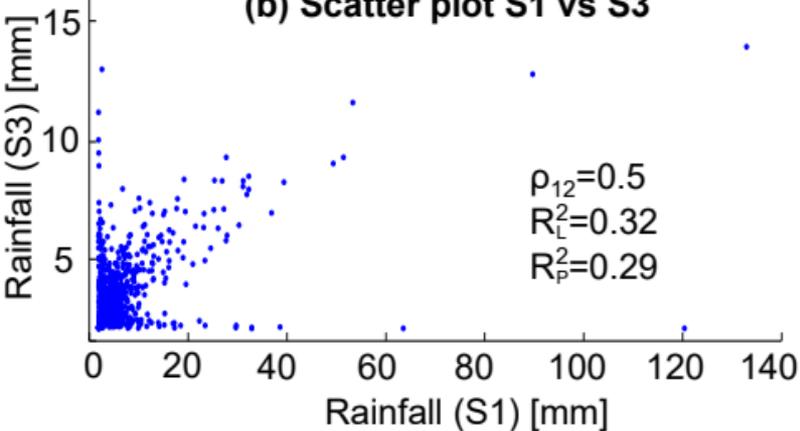
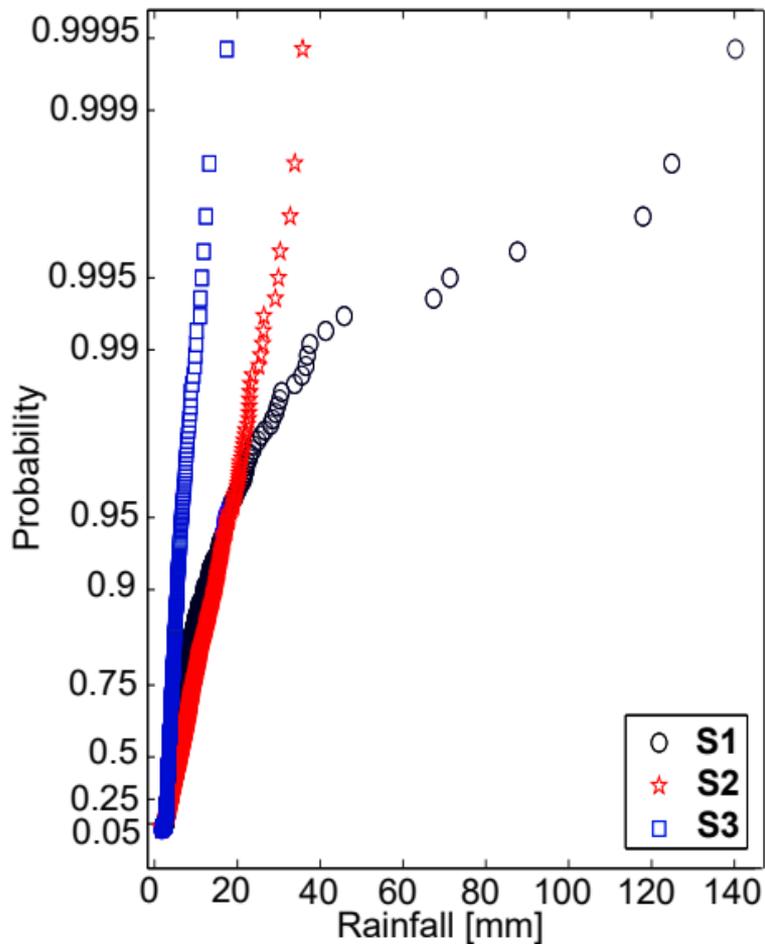


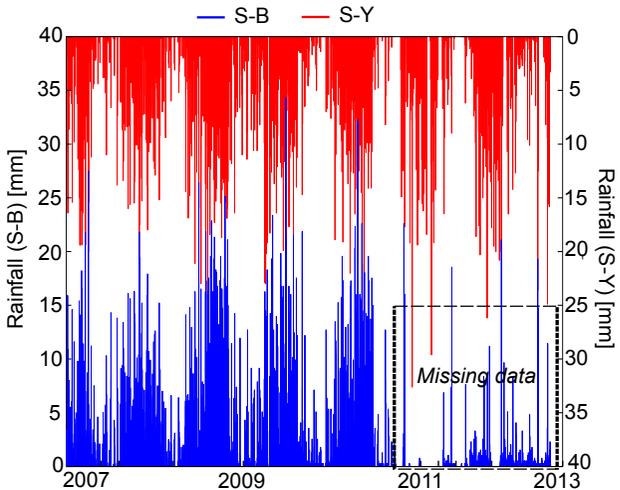
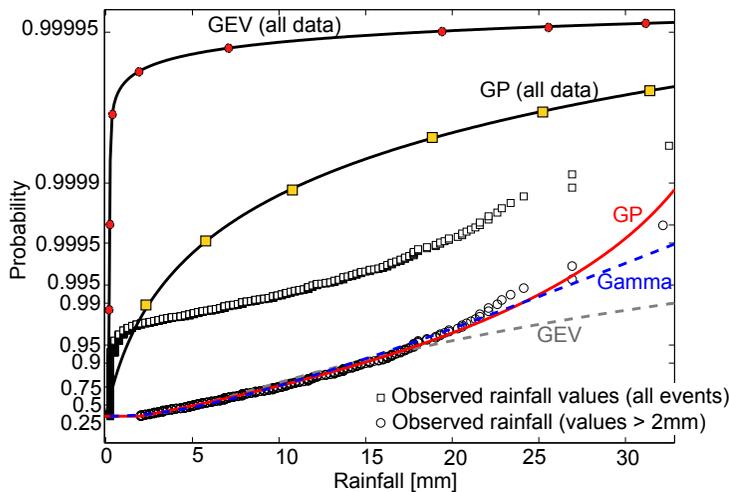
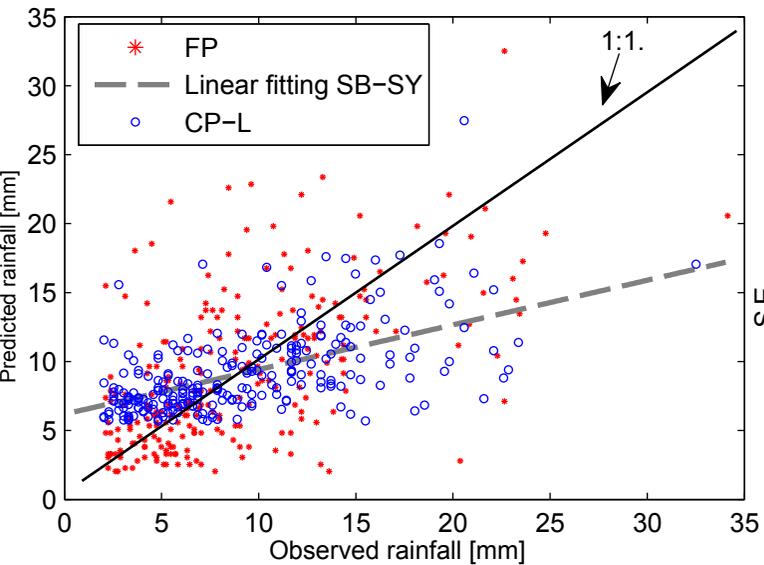
(2) FP approach







(a) Scatter plot S1 vs S2**(b) Scatter plot S1 vs S3****(c) Empirical POT distributions**

(a) Observed rainfall**(b) Statistical distributions****(c) Scatter plot observed vs reconstructed****(d) Cumulative errors (SE)**