

# 1 **A quick gap-filling of missing hydrometeorological data**

2 Christoforos Pappas<sup>1,2</sup>, Simon Michael Papalexiou<sup>2</sup>, Demetris Koutsoyiannis<sup>2</sup>

3 <sup>1</sup>Institute of Environmental Engineering, ETH Zurich, 8093 Zurich, Switzerland  
4 ([pappas@ifu.baug.ethz.ch](mailto:pappas@ifu.baug.ethz.ch))

5 <sup>2</sup>Department of Water Resources, Faculty of Civil Engineering, National Technical University of  
6 Athens, Heroon Polytechniou 5, GR-157 80 Zographou, Greece

## 7 **Abstract**

8 Data-gaps are ubiquitous in hydrometeorological time series and filling these values remains still  
9 a challenge. Since datasets without missing values may be a prerequisite in performing many  
10 analyses, a quick and efficient gap-filling methodology is required. In this study the problem of  
11 filling sporadic, single-value gaps using time-adjacent observations from the same location is  
12 investigated. The applicability of a local average (i.e., based on few neighboring in time  
13 observations) is examined and its advantages over the sample average (i.e., using the whole  
14 dataset) are illustrated. The analysis reveals that a quick and very efficient (i.e., minimum mean  
15 squared estimation error) gap-filling is achieved by combining a strictly local average (i.e., using  
16 one observation before and one after the missing value) with the sample mean.

17 **Keywords:** hydrometeorological data, missing values, gap-filling, interpolation, time series  
18 analysis

## 19 **1. Introduction**

20 Observing natural phenomena is of ultimate importance for understanding their complex  
21 characteristics. Understanding and simulating the earth-system processes requires a dense  
22 monitoring network of, not only long and reliable records, but also serially complete  
23 observations [e.g., *Butler*, 2014; *Baldocchi et al.*, 2012; *Silberstein*, 2006]. Time series of  
24 different geophysical variables (e.g., precipitation) emerge therefore from the systematic  
25 monitoring of their temporal evolution.

26

27 These instrumental time series are often plagued with a percentage of missing values (caused for  
28 example by malfunctioning of the equipment) creating sporadic and/or continuous gaps in their  
29 regular time-step. Many practical applications (e.g., extreme value analysis, continuous  
30 hydrological modeling) as well as statistical methodologies (e.g., spectral analysis, calibration  
31 (learning) algorithms, stochastic modeling and downscaling) have no tolerance to missing  
32 values. Preprocessing of raw datasets by infilling their missing values is thus a necessary  
33 procedure. Several interpolation techniques have been developed ranging from rather simple to  
34 extremely complex approaches. For example, *Henn et al.*, [2013], *Graham* [2009], *Horton and*  
35 *Kleinman* [2007], *Allison* [2003], *Roth* [1994], *Kemp et al.* [1983] provide detailed reviews of  
36 several gap-filling approaches applied to various scientific disciplines.

37

38 Several methods have been proposed for gap-filling environmental datasets e.g., linear or logistic  
39 regression, polynomial or spline interpolation, inverse distance weighting, ordinary kriging, and  
40 stochastic models that are fitted to the available records. More details on the aforementioned  
41 methodologies can be found in *Koutsoyiannis and Langousis* [2011] as well as in *Maidment*,

42 [1993; ch. 19.4]. Additional statistical techniques that have been developed in the last decade,  
43 include artificial neural networks and nearest neighbor techniques [*Elshorbagy et al.*, 2000,  
44 2002], as well as approaches based on Kalman filter [*Alavi et al.*, 2006] and nonlinear  
45 mathematical programming [*Teegavarapu*, 2012]. Hybrid methods (employing both process-  
46 based and statistical tools) have been often also applied as part of weather generators [e.g., the  
47 MicroMet meteorological model; *Liston and Elder*, 2006]. Yet, the complexity and the  
48 computational demand of such methodologies often hamper their applicability to real world  
49 applications. While data-gaps are ubiquitous in hydrometeorological time series, how these gaps  
50 were filled is not often reported, or naïve approaches have been unjustifiably selected (e.g., such  
51 as filling the gaps with a fixed value, often corresponding to the sample average).

52

53 In this study, we present a definitive argument against the use of the sample average for filling  
54 correlated hydrometeorological data. In addition, an innovative methodology, tailored for a quick  
55 filling of sporadic (i.e., single-value) gaps using information from time-adjacent values of the  
56 same location (i.e., within-station method), is presented and its advantages over other commonly  
57 used approaches are illustrated. The present study provides therefore a quick gap-filling with  
58 high efficacy and is geared towards practitioners and data analysts.

## 59 **2. Autocorrelation structure**

60 Filling missing data, irrespective of the implemented statistical technique, requires a good  
61 understanding of the underlying process and its peculiarities. Although many properties are  
62 necessary for a complete description of the observed variables, their autocorrelation structure is  
63 of great importance. Autocorrelation describes the linear dependences among different values of  
64 a time series providing therefore insights on how their persistence evolves in time. As such, it is

65 a key component in distilling information on the missing data and thus a cornerstone for the  
66 presented methodology.

67

68 Numerous studies illustrate different correlation patterns appropriate for describing several  
69 geophysical phenomena. Trying to cover the entire spectrum of autocorrelation structures widely  
70 detected and used in the hydrometeorological literature, we are focusing on: (i) processes with  
71 short-term persistence, characterized by exponential autocorrelation structure, and (ii) processes  
72 with long-term persistence, described by a power law autocorrelation function. These two  
73 structures have totally different characteristics in terms of time-dependence of the process and  
74 they are commonly present in different hydrometeorological variables (e.g., runoff,  
75 *Koutsoyiannis* [2013]; precipitation, *Marani* [2003]; sea level pressure and temperature, *Percival*  
76 *et al.* [2001], *Stephenson et al.* [2000]). Note that the suggested methodology is not limited to  
77 these two particular correlation structures. On the contrary, as it is demonstrated in the following  
78 sections, its applicability is more general, providing that the lag-1 autocorrelation can be  
79 estimated. Thus, the selection of known autocorrelation structures serves only for illustration of  
80 the theoretical framework underlying the methodology. In real-world applications, the estimation  
81 of empirical autocorrelations is enough for assuring the applicability and the efficacy of the  
82 proposed methodology.

83

84 In the following sections, hydrometeorological variables are treated as random variables modeled  
85 as stationary (more specifically, weakly stationary, i.e., with constant expected value and  
86 autocorrelation that depends only on the time lag; *Papoulis*, 1965, p.302) stochastic processes in  
87 discrete time. Regarding the notation used, the so-called Dutch notational convention is applied:

88 matrices and vectors are denoted by bold, random variables and stochastic processes are  
89 underlined, whereas their realizations (e.g. observed values) and the regular variables, are not.

## 90 **2.1 Exponential autocorrelation structure**

91 It has been very often claimed that hydrometeorological variables exhibit short-range  
92 dependence. For example, several studies have asserted that daily precipitation [*Gilman*, 1963],  
93 sea-surface temperature anomalies [*Frankignoul and Hasselmann*, 1977], Arctic sea ice  
94 [*Blanchard-Wrigglesworth et al.*, 2011; but see *Agarwal et al.*, 2012], climate variability  
95 [*Hasselmann*, 1976], as well as teleconnection patterns (such as North Atlantic Oscillation,  
96 Pacific-North American and West Pacific) [*Feldstein*, 2000; *Wunsch*, 1999; but see *Percival et*  
97 *al.*, 2001; *Stephenson et al.*, 2000] are characterized by Markovian dependence structure, i.e., the  
98 future appears to be independent of the past under the condition of known present [*Papoulis*,  
99 1965, p.535]. This dependence is theoretically justified in a few cases, but appears to be  
100 physically implausible [*Koutsoyiannis and Montanari*, 2007; *Koutsoyiannis*, 2011]. Thus, the  
101 wide use of the model can be attributed to its simplicity rather than to its sound physical basis.

102

103 The Markovian property is reproduced by an autoregressive model of order one, AR(1), and the  
104 autocorrelation for different values of time lag  $j$  is given by:

$$105 \quad \rho_j = \text{corr}[\underline{x}_i, \underline{x}_{i+j}] = \frac{\text{cov}[\underline{x}_i, \underline{x}_{i+j}]}{\text{var}[\underline{x}_i]} = \rho^{|j|} \quad (1)$$

106 where  $\rho$  is the lag-1 autocorrelation coefficient ( $|\rho| < 1$ ) which quantifies the short-range  
107 dependence. This relationship implies that the time dependence decreases exponentially as the  
108 time step (lag) increases, leading to practically negligible values of autocorrelation even for  
109 small values of lag (Figure 1).

## 110 2.2 Power law autocorrelation structure

111 There is strong empirical evidence that many natural phenomena are better characterized by  
112 highly persistent serial correlations rather than exponentially decaying autocorrelation structures.  
113 This natural behavior is often referred as Hurst phenomenon, long-term persistence, long-range  
114 dependence, or Hurst-Kolmogorov (HK) behavior [*Koutsoyiannis and Cohn, 2008*]. Here, the  
115 latter term is adopted, acknowledging the pioneering contribution of both, H. E. Hurst who first  
116 detected empirically that Nile river-level data exhibit long-term persistence [*Hurst, 1951*], and A.  
117 N. Kolmogorov who developed a basic mathematical framework describing this behavior  
118 [*Kolmogorov, 1940*].

119  
120 HK behavior is identified in many diverse geophysical quantities such as wind power [*Bakker*  
121 *and van den Hurk, 2012; Haslett and Raftery, 1989*]; precipitation [*Fatichi et al., 2012;*  
122 *Koutsoyiannis and Langousis, 2011; Montanari et al., 1996; Savina et al., 2011*, but see *Bunde et*  
123 *al., 2013*]; snow depth [*Egli and Jonas, 2009*]; temperature [*Bloomfield, 1992; Gil-Alana, 2005;*  
124 *Scafetta and West, 2005*]; river discharge [Nile, Africa, *Koutsoyiannis, 2002*; Warta, Poland,  
125 *Radziejewski and Kundzewicz, 1997*; Po, Italy, *Montanari, 2012*; Tiber, Italy *Grimaldi, 2004*;  
126 Boeotikos Kephisos, Greece, *Koutsoyiannis, 2003*]; indices of North Atlantic Oscillation  
127 [*Stephenson et al., 2000*]; solar activity [*Ogurtsov, 2004; Scafetta and West, 2005*; but see  
128 *Rypdal and Rypdal, 2012*]; extratropical atmospheric circulation anomalies [*Tsonis et al., 1999*];  
129 paleoclimate records [*Huybers and Curry, 2006; Markonis and Koutsoyiannis, 2012*].

130  
131 A power law representation of autocorrelation decay with lag appears to be more appropriate for  
132 describing the temporal dependences of these phenomena. While lag-1 autocorrelation measures

133 short-term persistence, the Hurst exponent,  $H$  ( $0.5 < H < 1$ ) is used to characterize the strength  
 134 of the HK behavior (i.e., long-term persistence). For the case of random noise  $H = 0.5$ , whereas  
 135 for real-world time series, like the examples mentioned above,  $H$  is often much higher. The  
 136 autocorrelation function for lag  $j$ , is given by:

$$137 \quad \rho_j = \frac{\text{cov}[\underline{x}_i, \underline{x}_{i+j}]}{\text{var}[\underline{x}_i]} = \frac{1}{2} \left( |j+1|^{2H} + |j-1|^{2H} \right) - |j|^{2H} \quad (2)$$

138 which for large  $j$  is proportional to  $j^{2H-2}$ . This behavior implies that the autocorrelation  
 139 decreases according to a power-type function of lag, which is much slower than the exponential  
 140 decay described by the Markovian dependence. Indeed, while the time dependence of the AR(1)  
 141 process practically (for  $\rho < 0.75$ ) vanishes for lag  $\approx 10$ , the autocorrelations of HK process are fat-  
 142 tailed, maintaining significant values even for lags orders of magnitude higher (Figure 1).

### 143 3. Filling methods

144 Gap-filling techniques can be often presented as weighted averages of existing observations  
 145 [Koutsoyiannis and Langousis, 2011] which can be summarized as follows:

$$146 \quad \underline{y} = \mathbf{w}^T \underline{\mathbf{X}} + \underline{e} \quad (3)$$

147 where  $\underline{y}$  is the missing value under examination ( $\underline{y} \equiv \underline{x}_0$ ),  $\underline{\mathbf{X}} = [\underline{x}_{-N}, \dots, \underline{x}_{-2}, \underline{x}_{-1}, \underline{x}_1, \underline{x}_2, \dots, \underline{x}_N]^T$   
 148 is a vector with the  $2N$  random variables corresponding to the available observations ( $T$  denotes  
 149 the transpose of the vector),  $\mathbf{w} = [w_{-N}, \dots, w_{-2}, w_{-1}, w_1, w_2, \dots, w_N]^T$  is a vector with the weights  
 150 assigned to each of the available observed values  $\underline{\mathbf{X}}$ , and  $\underline{e}$  is the estimation error. Different  
 151 infilling techniques provide therefore means for estimating the weighting parameters  $\mathbf{w}$ . In the

152 following sections, the Mean Squared Error (defined as  $\text{MSE} = \text{E}[\underline{e}^2] = \sigma_e^2 + \mu_e^2$ ) is used as a  
 153 performance metric for assessing different gap-filling approaches.

### 154 3.1 Optimal Local Average (OLA)

155 We assume that all the records used for the gap-filling ( $x_{-n}, \dots, x_{-2}, x_{-1}, x_1, x_2, \dots, x_n$ ) have equal  
 156 weights (i.e.  $w_{-n} = \dots = w_{-2} = w_{-1} = w_1 = w_2 = \dots = w_n = 1/2n$ ; where  $2n$  is the number of time-  
 157 adjacent values used for estimating the missing value under examination). The estimated missing  
 158 value under examination can then be expressed as  $\hat{x}_t = \left( \sum_{i=1}^n x_{t-i} + \sum_{i=1}^n x_{t+i} \right) / 2n$  (i.e., arithmetic  
 159 mean approach). Within the framework of Optimal Local Average (OLA) approach we test  
 160 which is the optimal number of neighboring records ( $n$ ) that should be used in order to have the  
 161 best estimation of the missing value (i.e., the one that minimizes MSE). The MSE is given by:

$$162 \quad \text{MSE} = \text{E}[\underline{e}^2] = \text{E}\left[(x_t - \hat{x}_t)^2\right] = \text{E}\left[\left(x_t - \frac{\sum_{i=1}^n x_{t-i} + \sum_{i=1}^n x_{t+i}}{2n}\right)^2\right] \quad (4)$$

163 Assuming that the underlying stochastic process is (weakly) stationary, we can express the MSE  
 164 as a function of the standard deviation  $\sigma$  of the process, the number of the neighboring values  $n$   
 165 used for the infilling, and the correlation coefficient  $\rho_i$  for different values of lag  $i$ :

$$166 \quad \text{MSE} = \text{E}[\underline{e}^2] = \frac{1}{2} \left( \frac{\sigma}{n} \right)^2 \left[ (2n+1) \left( n - 2 \sum_{i=1}^n \rho_i \right) + \sum_{i=1}^{2n} (2n+1-i) \rho_i \right] \quad (5)$$

167 The necessary algebraic manipulations for the derivation of Eq. (5) are detailed in the auxiliary  
 168 material (S1). The resulting MSE for the two examined autocorrelation structures and different  
 169 values of lag-1 autocorrelation (based on Eq. (5)) is illustrated in Figure 2 (see also Figure S2).



170

171 When processes with exponential autocorrelation structure are analyzed (Figure 2a and S2a), the  
172 strictly local average (i.e.,  $n = 1$ ) provides the minimum MSE for a wide range of lag-1  
173 autocorrelations. There is a critical value of lag-1 autocorrelation ( $\rho_{cr}^{AR} = 0.29$ ) above which the  
174 strictly local average provides the best estimate (Table 1 and Figure 2a and S2a). This manifests  
175 the fundamental Markovian property underlying the AR(1) process, i.e., when  $\rho$  becomes non-  
176 negligible ( $\rho \geq 0.29$ ) the information content of neighboring values should only be used,  
177 otherwise the MSE is larger (Figure 2a and S2a). For the case of power law autocorrelation  
178 structure, the time-adjacent values required for a minimum MSE decrease gradually as the lag-1  
179 autocorrelation increases (Figure 2b and S2b), as opposed to the sharp response of the AR(1)  
180 processes (Figure 2a and S2a). A critical value of  $\rho$  above which the strictly local average is  
181 preferable still exist but it is higher than the one of AR(1) processes ( $\rho_{cr}^{HK} = 0.52$ ; Table 1 and  
182 Figure 2b and S2b).

183

184 In summary, no matter which is the underlying autocorrelation structure (exponential or power-  
185 type) when  $\rho \geq 0.52$  the strictly local average (i.e., using one observation before and one after  
186 the missing value) provides the best estimate. Moreover, as Figure 2 and S2 illustrate, for a wide  
187 range of lag-1 autocorrelations, the sample average inflates the MSE, and therefore should be  
188 avoided when correlated data have to be infilled.

### 189 **3.2 Weighted Sum of local and total Average (WSA)**

190 Building upon the aforementioned findings (i.e., the time-adjacent values used for an efficient  
191 gap-filling can be determined by the lag-1 autocorrelation of the examined data), we provide a  
192 generalized framework distilling information from both, local and sample (global) average. The

193 sum of the strictly local (i.e., one value before and one after the missing record) and sample  
 194 average, weighted according to the lag-1 autocorrelation of the examined data, is used as an  
 195 estimation of the missing value. The estimated missing value is then given by:

$$196 \quad \hat{x}_t = \lambda \frac{\sum_{i=1}^N (x_{t-i} + x_{t+i})}{2N} + (1-\lambda) \frac{x_{t-1} + x_{t+1}}{2} \quad (6)$$

197 where  $\lambda$  is the weighting factor assigned to the sample average. In essence, parameter  $\lambda$  reflects  
 198 the strength of the temporal correlation i.e., low (high) values of  $\lambda$  imply low (high) contribution  
 199 of the sample average and thus high (low) temporal autocorrelation. Under the assumption of  
 200 (weak) stationarity, the MSE is given by:

$$201 \quad \text{MSE} = \text{E}[\underline{e}^2] = \text{E}[(x_t - \hat{x}_t)^2] = \text{E} \left[ \left( x_t - \left( \lambda \frac{\sum_{i=1}^N (x_{t-i} + x_{t+i})}{2N} + (1-\lambda) \frac{x_{t-1} + x_{t+1}}{2} \right) \right)^2 \right] \quad (7)$$

202 where  $2N$  is the length of the available observations. After some algebraic manipulations,  
 203 detailed in the auxiliary material (S3), the following expression is obtained:

$$204 \quad \text{MSE} = \frac{1}{2} \sigma^2 (3 - 4\rho_1 + \rho_2) - 2\lambda \sigma^2 \left[ \frac{1}{N} \sum_{i=1}^N \rho_i - \frac{1}{2N} \left( \sum_{i=1}^{N-1} \rho_i - \sum_{i=2}^{N+1} \rho_i + 1 \right) - \rho_1 + \frac{\rho_2}{2} + 0.5 \right] \\
 + \lambda^2 \sigma^2 \left[ \frac{1}{2N^2} \left( 2 \sum_{i=1}^{N-1} (N-i) \rho_i + \sum_{i=2}^{N+1} (i-1) \rho_i + \sum_{i=N+2}^{2N} (2N+1-i) \rho_i + N \right) \right. \\
 \left. + \frac{\rho_2}{2} + \frac{1}{2} - \frac{1}{N} \left( \sum_{i=1}^{N-1} \rho_i + \sum_{i=2}^{N+1} \rho_i + 1 \right) \right] \quad (8)$$

205 The MSE is therefore expressed as a function of  $\sigma$ ,  $\rho_i$ ,  $N$ , and the weighting factor  $\lambda$ . For a given  
 206 autocorrelation structure (i.e., known  $\rho_i$ ), we seek the value of  $\lambda$  that yields the minimum MSE  
 207 ( $\lambda_{\text{opt}}$ ; black dots in Figure 3; see also Figure S5 and S6a). For both Markovian and HK behavior,

208 as lag-1 autocorrelation increases the contribution of the local average increases (i.e.,  $\lambda_{\text{opt}}$   
209 decreases, Figure 3; see also Figure S5 and S6a).

210

211 An interesting property of the HK behavior is reflected in Figure 3b (see also Figure S6a): for  
212 high values of lag-1 autocorrelation there is a discontinuity in the values of  $\lambda_{\text{opt}}$ . More  
213 specifically, while it is expected that for high lag-1 autocorrelation the sample average does not  
214 contribute at all to the estimation of the missing value (i.e.,  $\lambda_{\text{opt}} = 0$ ),  $\lambda$  does not reach zero  
215 gradually and has non-zero values even for high values of lag-1 autocorrelation (Figure 3b and  
216 S6a). The rationale behind this behavior is that it takes time for a process with long-range  
217 dependence to reveal its characteristics. More specifically, when the available time series length  
218 is relatively small, the estimated sample average is in essence a local rather than a global average  
219 (see also detailed discussion in the auxiliary material (S4)).

220

221 In order to assess the influence of sample size in the MSE estimation and thus in  $\lambda_{\text{opt}}$ , a  
222 sensitivity analysis was conducted checking sample sizes from  $2 \times 5$  to  $2 \times 10^7$  (auxiliary material  
223 S4 and S5). For processes with exponential autocorrelation structure the relationship of  $\lambda_{\text{opt}}$  with  
224  $\rho$  does not vary much with the time series length (Figure S4 and S5) and it is thus approximated  
225 by:

$$226 \quad \lambda_{\text{opt}}^{\text{AR}(1)} = (1 - \rho)^{2.26} \quad (9)$$

227 For processes with HK behavior  $\lambda_{\text{opt}}$  depends highly on the time series length (Figure S4 and S6).

228 To mimic this type of dependence the  $\lambda_{\text{opt}}$  vs  $\rho$  relationship is approximated using two additional  
229 parameters ( $\lambda_1, \gamma$ ):

230 
$$\lambda_{\text{opt}}^{\text{HK}} = \left(1 - (1 - \lambda_1^\gamma) \rho\right)^{1/\gamma} \quad (10)$$

231 where  $\lambda_1 = 0.70 / (1 + \ln^2(N))^{0.69}$  and  $\gamma = 0.44 - 0.33 / (1 + \ln^2(1 + 0.03 \ln^2(N)))$ .

232

233 In practice, for Markovian processes, once the lag-1 autocorrelation is estimated from the data, it  
234 can be plugged-in to Eq. (9) and estimate the value of  $\lambda_{\text{opt}}$ , then, Eq. (6) can be applied for gap-  
235 filling the examined missing value. For data with HK behavior, one additional (but very simple)  
236 step is needed, i.e., the calculation of  $\lambda_1$  and  $\gamma$  given the length of the available records.

## 237 **4. Methods intercomparison and discussion**

238 The presented methods (OLA and WSA) are in essence generalizations of the widely applied  
239 concept of arithmetic mean, enhanced with information from the lag-1 autocorrelation of the  
240 examined time series. Their performance is tested against the sample average and the strictly  
241 local average (i.e., a linear interpolation between the values adjacent to the missing value) by  
242 comparing the resulting estimation error (Figure 4 and Figure 5).

243

### 244 **4.1 Monte-Carlo simulations**

245 Synthetic time series with 100000 values were generated from AR(1) and HK processes with  
246 zero mean and standard deviation equal to one and with lag-1 autocorrelation coefficients  
247 covering the entire range of possible values. The time series with HK dynamics were simulated  
248 using the function `SimulateFGN` from the R package `FGN` [Veenstra and McLeod, 2012]. Each  
249 value of the time series was then sequentially removed and the artificial data gap was then filled  
250 with the OLA and WSA approach, as well as we the sample and the strictly local average.

251

252 For the entire range of lag-1 autocorrelation for both exponential and power-type autocorrelation  
253 structures, the WSA methodology provides the minimum MSE (Figure 4). When uncorrelated  
254 data are examined (i.e., lag-1 autocorrelation tends to zero) the WSA method converge to the  
255 sample average approach, while for highly correlated time series, it converges to the strictly local  
256 average approach. For relatively strongly correlated data (e.g., with lag-1 autocorrelation higher  
257 than 0.5) WSA and OLA methods have similar performance but the flexible character of WSA  
258 approach makes it appropriate for gap-filling sporadic gaps across the entire spectrum of  
259 temporal dependences (minimum MSE; Figure 4). Similarly to the inverse distance weighting  
260 (where information from neighboring in space station is used), in WSA approach the notion of  
261 similarity (in time) between data points is crucial. The WSA methodology is therefore based on  
262 Tobler's first law in geography i.e., "everything is related to everything else, but near things are  
263 more related than distant things" [Tobler, 1970].

264

265 The sample average provides the worst results (highest MSE; Figure 4). Replacing a missing  
266 value with the sample average is a simple and easy approach for dealing with missing data, but  
267 as our analysis reveals much better results can be obtained by applying tools of similar  
268 complexity (i.e., WSA approach). It is also worth mentioning that a great advantage of the WSA  
269 methodology is that the probability distribution of the observations and the temporal  
270 relationships (i.e., autocorrelation) remain relatively undisturbed, avoiding induction of biases in  
271 the mean, variance or autocorrelation of the final time series. As already underlined elsewhere  
272 [e.g., *Little and Rubin*, 2002], replacing all missing values in a dataset with a single value (e.g.,  
273 using the sample mean) apart from reducing the variance, can often artificially inflate the  
274 significance of any statistical test that is based on these statistics. The proposed method is not

275 free of these problems, but if the autocorrelation is strong and the percentage of missing values  
276 low, the reduction of variance is not substantial. In addition, with the WSA methodology no  
277 assumption was made regarding the distribution underlying the examined dataset. Moreover, as  
278 demonstrated in Eq. (8), there is no dependence of the MSE to the mean properties of the  
279 analyzed dataset, therefore there is no need for data preprocessing (e.g., normalization). Given  
280 that the underlying autocorrelation structure is identified, this simple method does not impose  
281 any requirement for calculation of other statistical quantities apart from the sample mean and the  
282 lag-1 autocorrelation for its application.

## 283 **4.2 Real-world applications**

284 An additional illustration of the aforementioned finding is summarized in Figure 5. Real-world  
285 time series from the Global Historical Climatology Network, (GHCN version 2.60;  
286 [www.ncdc.noaa.gov/oa/climate/ghcn-daily](http://www.ncdc.noaa.gov/oa/climate/ghcn-daily)) and from the Roda Nilometer, near Cairo (minimum  
287 water levels of Nile; [Toussoun, 1925]) are presented (Figure 5).

288  
289 Time series of annual precipitation (GHCN, station ID: CA003031093) behaving as Markovian  
290 process with lag-1 autocorrelation  $\rho = 0.29$ , and temperature (GHCN, station ID:  
291 GM000003342) presenting HK dynamics with Hurst exponent  $H = 0.72$  (estimated using the  
292 slope of the climacogram, which is a double-logarithmic plot of the standard deviation of the  
293 sample at an aggregate timescale vs the timescale; Koutsoyiannis, 2003, 2010) are examined.  
294 Observed records spanning from 1893 to 2011, were infilled removing sequentially every single  
295 value and gap-filling with the presented methodologies (Figure 5). The original and the infilled  
296 time series, as well as the efficiency of each infilling approach (defined as  
297  $1 - \text{nRMSE} = 1 - \sqrt{\text{MSE}}/s$  where nRMSE is the normalized Root Mean Squared Error and  $s$  is the

298 standard deviation of the observed time series) are presented in Figure 5a,c and Figure 5b,d  
299 respectively.

300

301 The longest instrumental record of the water levels of Nile is also examined (focusing on the  
302 period 622 AD to 1470 AD that the record is almost uninterrupted), illustrating the advantages of  
303 within-station gap-filling approaches. In this case, the use of within-station information is  
304 apparently the only solution since neighboring stations, covering the same time period, are not  
305 available. The annual minimum water levels of Nile are characterized by HK dynamics with  
306 Hurst exponent  $H = 0.87$  [Koutsoyiannis, 2013]. For clarity in the illustration, only 200 years are  
307 presented, covering the period 800 AD to 1000 AD (Figure 5e,f), but the processing was made,  
308 and the efficiency was calculated, for the entire series. The WSA approach yields the highest  
309 efficiency (Figure 5f). Since the data present strong autocorrelation, gap-filling with OLA and  
310 local average approach converges to the same results (Figure 5e,f; see also Table 1). As in the  
311 previous two real-world examples (i.e., annual precipitation and temperature time series), the  
312 sample average has no skill (i.e., efficiency tends to 0; Figure 5b,d,f) manifesting that its use in  
313 gap-filling hydrometeorological variables is not only unjustified, but also seriously flawed.

314

315 In accordance with the presented theory (Section 3) and Monte-Carlo simulations (Figure 4),  
316 WSA approach provides the highest efficiency ( $1 - nRMSE$ ; Figure 5b,d,f). It is worth underline  
317 that, when the examined time series do not present high autocorrelation, the other three  
318 approaches (i.e., sample and local average as well as OLA) lead to comparable results in terms of  
319 overall efficiency (Figure 5b,d), but the distribution of the infilled time series varies significantly  
320 (Figure 5a,c). When data with strong autocorrelation are examined (Figure 5f), the use of

321 neighboring in time values improves significantly the performance of the gap-filling approach  
322 (Figure 5e,f). As expected, the use of sample average vanishes the variability presented in the  
323 original record, while the local average preserves many interesting features of the original  
324 records (Figure 5a,c,e).

## 325 **5. Limitations and further improvements**

326 While our analysis is focusing on sporadic, single-value data-gaps, generalizations of the  
327 presented approach for a wider gap-window are possible. Continuous missing values can be  
328 infilled by applying sequentially the presented framework (cascade process). More specifically,  
329 the available observations, at the gap-window boundaries, are used for the estimation of the  
330 missing value in the middle of the gap-window. This value is then used as a proxy, applying  
331 again the WSA approach for the new, restricted gap-window. A cascade-based procedure can be  
332 thus applied gap-filling sequentially continuous missing values. However, it is worth mentioning  
333 that the wider the data-gaps the more uncertain the estimated first- and higher-order statistics of  
334 the examined time series and thus the estimated missing values. A more elegant (but  
335 computationally more demanding) approach for dealing with multiple sequential data-gaps can  
336 emerge from Eq. (3). More specifically, instead of assigning equal weights to the available  
337 observations (as is the case for the presented methods; see Section 3), for each of the examined  
338 missing values, specific weighting factors ( $\mathbf{w}$ ) can be assigned to the available observations by  
339 solving explicitly Eq. (3); results of ongoing research on this issue are planned to be presented  
340 soon.



341 **6. Conclusions**

342 Conventional methods for handling missing data (such as sample average or linear interpolation  
343 of values adjacent to the missing record) are seriously flawed in the hydrometeorological time  
344 series, where the time series autocorrelation is non-negligible. Taking advantage of the  
345 information content of the lag-1 autocorrelation, a new flexible and equally simple framework  
346 for a quick gap filling of sporadic, single-value, gaps is proposed. The conclusions of our study  
347 are twofold: (i) a definitive argument against the use of the sample average for infilling  
348 correlated data is provided and demonstrated theoretically; and (ii) a new gap-filling  
349 methodology, equivalently simple but significantly more efficient, using a weighted sum of  
350 sample and strictly local average, is developed and its advantages are illustrated. The estimation  
351 of the sample mean and the lag-1 autocorrelation is the only necessity for assuring the  
352 applicability of WSA approach. The presented methodology is therefore a valuable tool for a  
353 quick filling of a small number of missing measurements tailored for hydrometeorological data  
354 as well as for a efficient gap-filling of missing paleoclimatic records, where neighboring station  
355 are not available.

356

357 **Acknowledgments**

358 We would like to thank the eponymous reviewer Prof. Salvatore Grimaldi and two anonymous  
359 reviewers for their constructive comments that helped us improve the presentation of our  
360 analysis. Fruitful discussions with Y. Dialynas, P. Kossieris, K. Kyriakidis, A. Lykou, and Y.  
361 Markonis, at an early stage of this analysis, are also gratefully acknowledged. The computer  
362 code used for the presented methodologies is available upon request from the corresponding  
363 author. All statistical analyses and graphics were performed using the R environment [R

364 *CoreTeam*, 2012] as well as additional packages as *ggplot2* [Wickham, 2009] and *SimulateFGN*  
365 [Veenstra and McLeod, 2012]. The authors would like to thank the developers for making them  
366 freely available. Global Historical Climatology Network is also acknowledged for providing the  
367 presented time series.

368

369 **References**

- 370 Agarwal, S., W. Moon, and J. S. Wettlaufer (2012), Trends, noise and re-entrant long-  
371 term persistence in Arctic sea ice, *Proc. R. Soc. A Math. Phys. Eng. Sci.*, 468(2144),  
372 2416–2432, doi:10.1098/rspa.2011.0728.
- 373 Alavi, N., J. S. Warland, and A. a. Berg (2006), Filling gaps in evapotranspiration  
374 measurements for water budget studies: Evaluation of a Kalman filtering approach,  
375 *Agric. For. Meteorol.*, 141(1), 57–66, doi:10.1016/j.agrformet.2006.09.011.
- 376 Allison, P. (2003), Missing data techniques for structural equation modeling, *J. Abnorm.*  
377 *Psychol.*, 112(4), 545–557, doi:10.1037/0021-843X.112.4.545.
- 378 Bakker, A. M. R., and B. J. J. M. van den Hurk (2012), Estimation of persistence and  
379 trends in geostrophic wind speed for the assessment of wind energy yields in  
380 Northwest Europe, *Clim. Dyn.*, 39(3-4), 767–782, doi:10.1007/s00382-011-1248-1.
- 381 Baldocchi, D., M. Reichstein, D. Papale, L. Koteen, R. Vargas, D. Agarwal, and R. Cook  
382 (2012), The Role of Trace Gas Flux Networks in the Biogeosciences, *EOS Trans.*,  
383 93(23), 217–224.
- 384 Blanchard-Wrigglesworth, E., K. C. Armour, C. M. Bitz, and E. DeWeaver (2011),  
385 Persistence and Inherent Predictability of Arctic Sea Ice in a GCM Ensemble and  
386 Observations, *J. Clim.*, 24(1), 231–250, doi:10.1175/2010JCLI3775.1.
- 387 Bloomfield, P. (1992), Trends in global temperature, *Clim. Change*, 21(1), 1–16,  
388 doi:10.1007/BF00143250.
- 389 Bunde, A., U. Büntgen, J. Ludescher, J. Luterbacher, and H. von Storch (2013), Is there  
390 memory in precipitation?, *Nat. Clim. Chang.*, 3(3), 174–175,  
391 doi:10.1038/nclimate1830.
- 392 Butler, D. G. (2014), Earth observation enters next phase, *Nature*, 508, 160–161.
- 393 Egli, L., and T. Jonas (2009), Hysteretic dynamics of seasonal snow depth distribution in  
394 the Swiss Alps, *Geophys. Res. Lett.*, 36(2), L02501, doi:10.1029/2008GL035545.
- 395 Elshorbagy, A., U. Panu, and S. Simonovic (2000), Group-based estimation of missing  
396 hydrological data: I. Approach and general methodology, *Hydrol. Sci. J.*, 45(6), 849–  
397 866.
- 398 Elshorbagy, A., S. Simonovic, and U. Panu (2002), Estimation of missing streamflow  
399 data using principles of chaos theory, *J. Hydrol.*, 255, 123–133.

- 400 Fatichi, S., V. Y. Ivanov, and E. Caporali (2012), Investigating Interannual Variability of  
401 Precipitation at the Global Scale: Is There a Connection with Seasonality?, *J. Clim.*,  
402 25(16), 5512–5523, doi:10.1175/JCLI-D-11-00356.1.
- 403 Feldstein, S. B. (2000), The Timescale, Power Spectra, and Climate Noise Properties of  
404 Teleconnection Patterns, *J. Clim.*, 13(24), 4430–4440, doi:10.1175/1520-  
405 0442(2000)013<4430:TTPSAC>2.0.CO;2.
- 406 Frankignoul, C., and K. Hasselmann (1977), Stochastic climate models, Part II  
407 Application to sea- • surface temperature anomalies and thermocline variability,  
408 *Tellus*, 29, 289–305.
- 409 Gil-Alana, L. (2005), Statistical modeling of the temperatures in the Northern  
410 Hemisphere using fractional integration techniques, *J. Clim.*, 18(24), 5357–5369.
- 411 Gilman, D. (1963), On the power spectrum of “red noise,” *J. Atmos. Sci.*, 20, 182–184.
- 412 Graham, J. (2009), Missing data analysis: Making it work in the real world, *Annu. Rev.*  
413 *Psychol.*, 60, 549–576, doi:10.1146/annurev.psych.58.110405.085530.
- 414 Grimaldi, S. (2004), Linear parametric models applied to daily hydrological series, *J.*  
415 *Hydrol. Eng.*, 9, 383–391.
- 416 Haslett, J., and A. Raftery (1989), Space-time modelling with long-memory dependence:  
417 Assessing Ireland’s wind power resource, *J. R. Stat. Soc. Ser. C (Applied Stat.)*,  
418 38(1), 1–50.
- 419 Hasselmann, K. (1976), Stochastic climate models Part I. Theory, *Tellus*, 28(6), 473–485,  
420 doi:10.1111/j.2153-3490.1976.tb00696.x.
- 421 Henn, B., M. S. Raleigh, A. Fisher, and J. D. Lundquist (2013), A Comparison of  
422 Methods for Filling Gaps in Hourly Near-Surface Air Temperature Data, *J.*  
423 *Hydrometeorol.*, 14(3), 929–945, doi:10.1175/JHM-D-12-027.1.
- 424 Horton, N. J., and K. P. Kleinman (2007), Much ado about nothing: A comparison of  
425 missing data methods and software to fit incomplete data regression models, *Am.*  
426 *Stat.*, 61(1), 79–90, doi:10.1198/000313007X172556.
- 427 Hurst, H. (1951), Long-term storage capacity of reservoirs, *Trans. Am. Soc. Civ. Eng.*,  
428 (2447).
- 429 Huybers, P., and W. Curry (2006), Links between annual, Milankovitch and continuum  
430 temperature variability, *Nature*, 441(7091), 329–32, doi:10.1038/nature04745.

- 431 Kemp, W., D. Burnell, D. Everson, and A. Thomson (1983), Estimating missing daily  
432 maximum and minimum temperatures, *J. Clim. Appl. Meteorol.*, 22, 1587–1593.
- 433 Kolmogorov, A. (1940), Wiener'sche Spiralen und einige andere interessante Kurven im  
434 Hilbert'schen Raum, *Dokl. Acad. Sci. URSS*, 26, 115–118.
- 435 Koutsoyiannis, D. (2002), The Hurst phenomenon and fractional Gaussian noise made  
436 easy, *Hydrol. Sci. J.*, 47(4), 37–41.
- 437 Koutsoyiannis, D. (2003), Climate change, the Hurst phenomenon, and hydrological  
438 statistics, *Hydrol. Sci. J.*, 48(1), 3–24, doi:10.1623/hysj.48.1.3.43481.
- 439 Koutsoyiannis, D. (2010), HESS Opinions “A random walk on water,” *Hydrol. Earth  
440 Syst. Sci.*, 14(3), 585–601, doi:10.5194/hess-14-585-2010.
- 441 Koutsoyiannis, D. (2011), Hurst–Kolmogorov dynamics as a result of extremal entropy  
442 production, *Phys. A Stat. Mech. its Appl.*, 390(8), 1424–1432,  
443 doi:10.1016/j.physa.2010.12.035.
- 444 Koutsoyiannis, D. (2013), Hydrology and Change, *Hydrol. Sci. J.*,  
445 doi:10.1080/02626667.2013.804626.
- 446 Koutsoyiannis, D., and T. Cohn (2008), The Hurst phenomenon and climate, *EGU Gen.  
447 Assem., Session IS*(Climatic and hydrological perspectives on long-term changes).
- 448 Koutsoyiannis, D., and A. Langousis (2011), Precipitation, in *Treatise on Water Science*,  
449 edited by P. Wilderer and S. Uhlenbrook, pp. 27–78, Academic Press, Oxford.
- 450 Koutsoyiannis, D., and A. Montanari (2007), Statistical analysis of hydroclimatic time  
451 series: Uncertainty and insights, *Water Resour. Res.*, 43(5), W05429,  
452 doi:10.1029/2006WR005592.
- 453 Liston, G., and K. Elder (2006), A meteorological distribution system for high-resolution  
454 terrestrial modeling (MicroMet), *J. Hydrometeorol.*, 7, 217–234.
- 455 Little, R., and D. Rubin (2002), *Statistical analysis with missing data*, Second Edi., John  
456 Wiley and Sons, New York.
- 457 Maidment, D. (1993), *Handbook of hydrology*, McGraw-Hill, New York.
- 458 Marani, M. (2003), On the correlation structure of continuous and discrete point rainfall,  
459 *Water Resour. Res.*, 39(5), 1128, doi:10.1029/2002WR001456.

- 460 Markonis, Y., and D. Koutsoyiannis (2012), Climatic Variability Over Time Scales  
461 Spanning Nine Orders of Magnitude: Connecting Milankovitch Cycles with Hurst-  
462 Kolmogorov Dynamics, *Surv. Geophys.*, doi:10.1007/s10712-012-9208-9.
- 463 Montanari, A. (2012), Hydrology of the Po River: looking for changing patterns in river  
464 discharge, *Hydrol. Earth Syst. Sci.*, 16(10), 3739–3747, doi:10.5194/hess-16-3739-  
465 2012.
- 466 Montanari, A., R. Rosso, and M. Taquu (1996), Some long-run properties of rainfall  
467 records in Italy, *J. Geophys. Res.*, 101(D23), 29431–29438.
- 468 Ogurtsov, M. (2004), New Evidence for Long-Term Persistence in the Sun’s Activity,  
469 *Sol. Phys.*, 220, 93–105.
- 470 Papoulis, A. (1965), *Probability, Random Variables, and Stochastic Processes*, McGraw-  
471 Hill.
- 472 Percival, D., J. Overland, and H. Mofjeld (2001), Interpretation of North Pacific  
473 Variability as a Short-and Long-Memory Process, *J. Clim.*, 14, 4545–4559.
- 474 Radziejewski, M., and Z. Kundzewicz (1997), Fractal analysis of flow of the river Warta,  
475 *J. Hydrol.*, 200, 280–294.
- 476 RCoreTeam (2012), *R: A Language and Environment for Statistical Computing*, Vienna,  
477 Austria.
- 478 Roth, P. L. (1994), Missing Data: a Conceptual Review for Applied Psychologists, *Pers.*  
479 *Psychol.*, 47(3), 537–560, doi:10.1111/j.1744-6570.1994.tb01736.x.
- 480 Rypdal, M., and K. Rypdal (2012), Is there long-range memory in solar activity on  
481 timescales shorter than the sunspot period?, *J. Geophys. Res.*, 117(A4), A04103,  
482 doi:10.1029/2011JA017283.
- 483 Savina, M., P. Molnar, and P. Burlando (2011), Seasonal long-term persistence in radar  
484 precipitation in complex terrain, *Water Resour. Res.*, 47(10), W10506,  
485 doi:10.1029/2010WR010170.
- 486 Scafetta, N., and B. J. West (2005), Multiscaling comparative analysis of time series and  
487 geophysical phenomena, *Complexity*, 10(4), 51–56, doi:10.1002/cplx.20076.
- 488 Silberstein, R. P. (2006), Hydrological models are so good, do we still need data?,  
489 *Environ. Model. Softw.*, 21(9), 1340–1352, doi:10.1016/j.envsoft.2005.04.019.

490 Stephenson, D. B., V. Pavan, and R. Bojariu (2000), Is the North Atlantic Oscillation a  
491 random walk?, *Int. J. Climatol.*, 20(1), 1–18, doi:10.1002/(SICI)1097-  
492 0088(200001)20:1<1::AID-JOC456>3.0.CO;2-P.

493 Teegavarapu, R. (2012), Spatial interpolation using nonlinear mathematical programming  
494 models for estimation of missing precipitation records, *Hydrol. Sci. J.*, 57(3), 37–41.

495 Tobler, W. (1970), A computer movie simulating urban growth in the Detroit region,  
496 *Econ. Geogr.*, 46, 234–240.

497 Toussoun, O. (1925), *Mémoire sur l'histoire du Nil*.

498 Tsonis, A. A., P. J. Roebber, and J. B. Elsner (1999), Long-Range Correlations in the  
499 Extratropical Atmospheric Circulation: Origins and Implications, *J. Clim.*, 12(5),  
500 1534–1541, doi:10.1175/1520-0442(1999)012<1534:LRCITE>2.0.CO;2.

501 Veenstra, J., and A. I. McLeod (2012), Hyperbolic Decay Time Series Models, *Press*.

502 Wickham, H. (2009), *ggplot2: elegant graphics for data analysis*, Springer New York.

503 Wunsch, C. (1999), The Interpretation of Short Climate Records, with Comments on the  
504 North Atlantic and Southern Oscillations, *Bull. Am. Meteorol. Soc.*, 80(2), 245–255,  
505 doi:10.1175/1520-0477(1999)080<0245:TIOSCR>2.0.CO;2.

506

507

508 **Tables**

509 **Table 1.** Time-adjacent values needed for an optimal infilling (minimum Mean Squared Error)  
 510 of missing observations according to the Optimal Local Average methodology, for different  
 511 autocorrelation structures (short or long-term persistence) and the lag-1 autocorrelations.

Optimal Local Average			
Short-term persistence		Long-term persistence	
$\rho \leq 0.25$	$n = n_{\max}$	$\rho < 0.3$	$n = n_{\max}$
$0.26 \leq \rho \leq 0.28$	$n = 2$	$0.30 \leq \rho \leq 0.32$	$n = 4$
		$0.33 \leq \rho \leq 0.38$	$n = 3$
$\rho \geq 0.29$	$n = 1$	$0.39 \leq \rho \leq 0.51$	$n = 2$
		$\rho \geq 0.52$	$n = 1$

$\rho$ : lag-one autocorrelation coefficient

$n$ : time-adjacent values used for the infilling

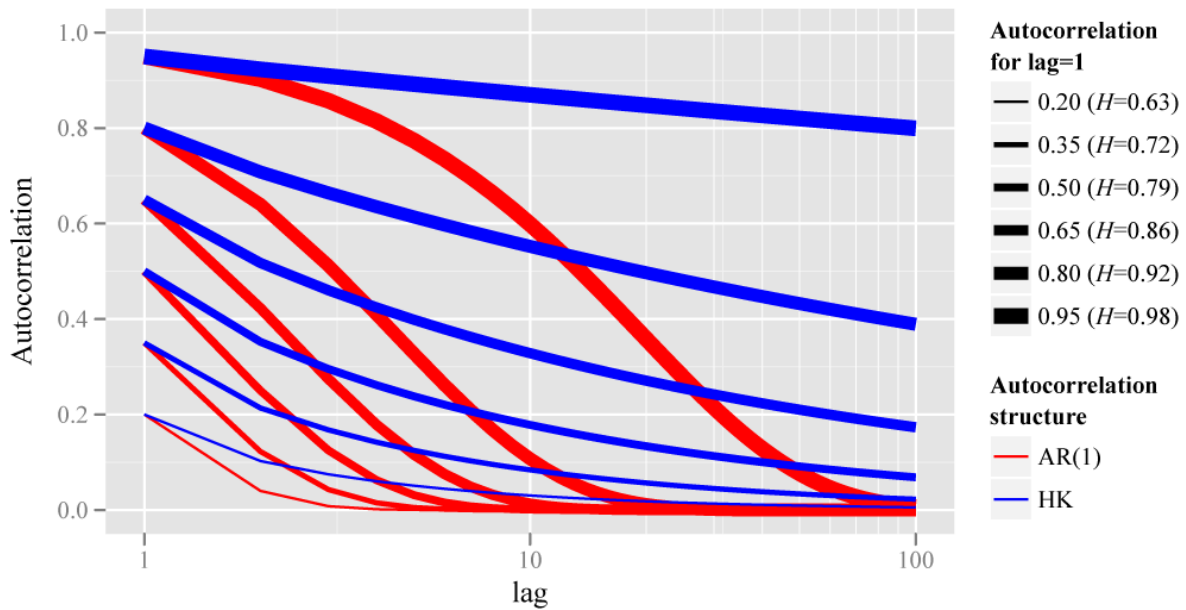
$n_{\max}$ : all the available observed values, i.e., sample average

512

513

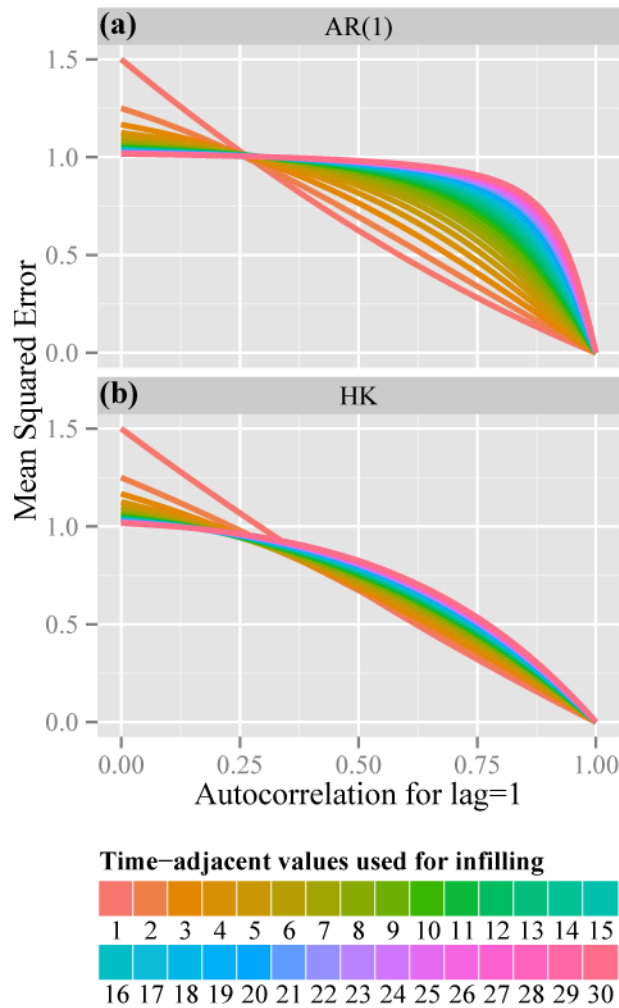


514 **Figures**



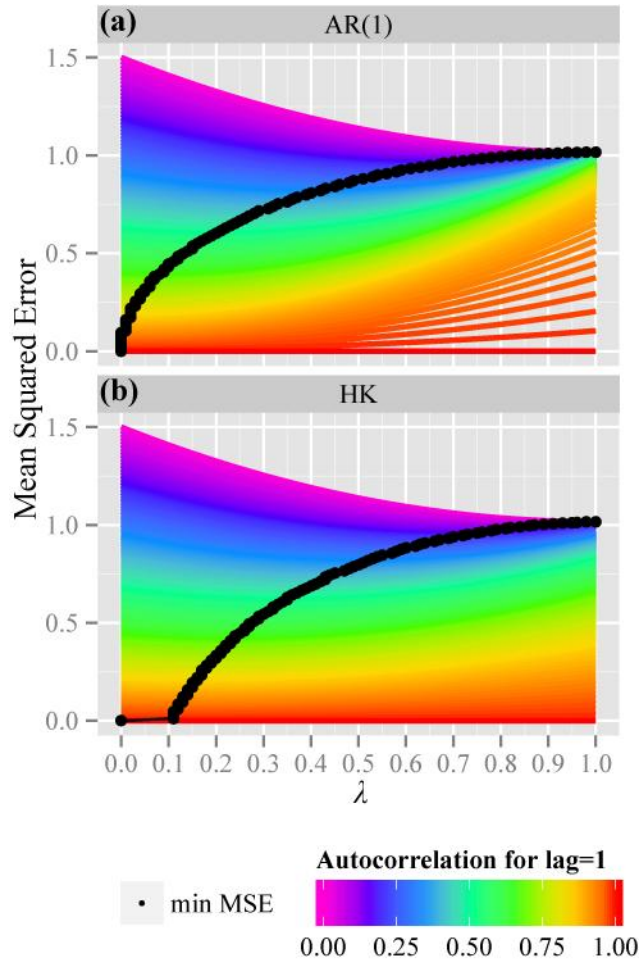
515

516 **Figure 1.** Theoretical autocorrelation functions for: (i) Markovian processes, AR(1), with  
 517 exponential decay of autocorrelation with lag (Eq. (1)) and (ii) processes with HK behavior,  
 518 described by the Hurst exponent  $H$ , with a power law relationship of autocorrelation with lag  
 519 (Eq. (2)). The lag-1 autocorrelation,  $\rho$ , characterizes the strength of short-term persistence while  
 520 the Hurst exponent,  $H$ , quantifies long-term dependences. Note that Eq. (2) implies that  $H$  and  $\rho$   
 521 are related as  $H = 0.5[\log_2(\rho+1)+1]$ .



522

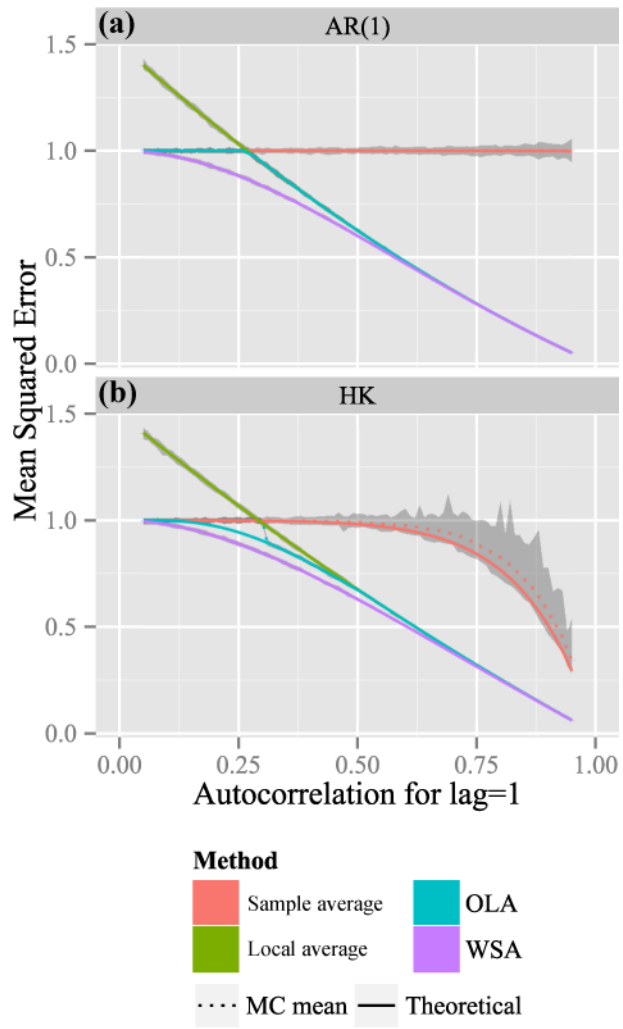
523 **Figure 2.** Illustration of the rationale underlying the Optimal Local Average (OLA)  
 524 methodology based on Eq. (5) for processes with (a) exponential, and with (b) power-law  
 525 autocorrelation structure. The Mean Squared Error of an estimated missing value, based on local  
 526 averages with different range (i.e., different number of neighboring values), for hypothetical time  
 527 series with different lag-1 autocorrelation and standard deviation equal to 1, is depicted (Eq. (5)).  
 528 When the number of time-adjacent values used for the local average estimation equals 1, one  
 529 value before and one after the missing observation are used for estimating the missing value,  
 530 while when this number equals 30, the sample average is used (i.e., the average of all available  
 531 observations, here for illustration assumed to be 30 before and 30 after the missing value).



532

533 **Figure 3.** Surface plots of the Mean Squared Error (MSE) estimated according to the Weighted  
 534 Sum of local and total Average (WSA) methodology (based on Eq. (8), and an hypothetical time  
 535 series length of  $2 \times 30$  and standard deviation equal to 1), for different values of parameter  $\lambda$ , for  
 536 processes with (a) exponential, and (b) power-law autocorrelation structure. The optimal values  
 537 of parameter  $\lambda$ , i.e., the ones that minimize the MSE are also highlighted (black dots). As the lag-  
 538 1 autocorrelation increases, the optimal values of parameter  $\lambda$ , which indicates the overall  
 539 contribution of the global average, decreases.

540



542 **Figure 4.** Estimated Mean Squared Error (MSE) based on different infilling methodologies  
543 (sample average i.e., using all the available values (here for illustration purposes  $2 \times 30$  values are  
544 used); strictly local average using one observation before and one after the missing record;  
545 Optimal Local Average methodology, OLA; Weighted Sum of local and total Average approach  
546 (WSA). Results correspond to processes with (a) exponential, and (b) power-law (b)  
547 autocorrelation structure for different values of lag-1 autocorrelation. The solid lines depict the  
548 theoretical values of MSE (see Eq. (5) and Eq. (8)) while the dashed lines and uncertainty  
549 bounds correspond to the ensemble of the Monte-Carlo simulations, filling artificial data gaps.  
550 For the entire range of lag-1 autocorrelations, the WSA approach significantly outperforms other  
551 infilling methods, providing the smallest MSE.

552



554 **Figure 5.** Real-world examples of time series with Markovian behavior (AR(1); annual  
555 precipitation, panel a) and with HK dynamics (annual temperature, panel c, annual minimum  
556 water depth, panel e). Original data are depicted in white circles, while the infilled time series are  
557 depicted in continuous colored lines. Each record was removed and infilled with the four  
558 examined approaches, calculating each time the new sample statistics. Bar-plots (panels b, d, f)  
559 illustrate the efficiency (defined as  $1 - \text{nRMSE}$  where nRMSE is the normalized Root Mean  
560 Squared Error) of each gap-filling approach (i.e., sample average, local average, Optimal Local  
561 Average, OLA, and Weighted Sum of local and total Average, WSA). Since the sample average  
562 is re-calculated each time a value is removed, the efficiency of the sample average approach is  
563 not always equal to 0 (as expected theoretically, i.e., the nRMSE once the sample average should  
564 be 100 %).