

# The mode of the climacogram estimator for a Gaussian Hurst-Kolmogorov process

Panayiotis Dimitriadis and Demetris Koutsoyiannis

## ABSTRACT

Geophysical processes are often characterized by long-term persistence. An important characteristic of such behaviour is the induced large statistical bias, i.e. the deviation of a statistical characteristic from its theoretical value. Here, we examine the most probable value (i.e. mode) of the estimator of variance to adjust the model for statistical bias. Particularly, we conduct an extensive Monte Carlo analysis based on the climacogram (i.e. variance of the average process vs. scale) of the simple scaling (Gaussian Hurst-Kolmogorov) process, and we show that its classical estimator is highly skewed especially in large scales. We observe that the mode of the climacogram estimator can be well approximated by its lower quartile (25% quantile). To derive an easy-to-fit empirical expression for the mode, we assume that the climacogram estimator follows a gamma distribution, an assumption strictly valid for Gaussian white noise processes. The results suggest that when a single timeseries is available, it is advantageous to estimate the Hurst parameter using the mode estimator rather than the expected one. Finally, it is discussed that while the proposed model for mode bias works well for Gaussian processes, for higher accuracy and non-Gaussian processes, one should perform a Monte Carlo simulation following an explicit generation algorithm.

**Key words** | climacogram, long-term persistence, mode estimator, statistical bias, stochastic uncertainty

**Panayiotis Dimitriadis** (corresponding author)  
**Demetris Koutsoyiannis**  
Department of Water Resources and  
Environmental Engineering, School of Civil  
Engineering,  
National Technical University of Athens,  
Heroon Polytechniou 5, 15880 Zographou,  
Greece  
E-mail: [pandim@itia.ntua.gr](mailto:pandim@itia.ntua.gr)

## INTRODUCTION

An important attribute characterizing geophysical processes is the high spatio-temporal dependence, in the sense that a random variable of such a process at a specific time or location strongly depends on several (even infinite) past, or of different location, random variables of the same process. This type of dependence requires long samples for its identification, which is a rare case in most natural processes, and thus, for the estimation of its parameters, it is advised to use only up to the second-order statistics (Lombardo *et al.* 2014) and only in cases where very long samples are available to expand to higher orders. The above issues are further highlighted in Dimitriadis (2017), where several (overall thirteen) such processes with various lengths and

physical properties expanding from small-scale turbulence to large-scale hydrometeorological processes are analyzed in terms of their long-term behaviour using massive databases and unbiased estimators of the second-order dependence structure. Interestingly, all the examined processes exhibited long-term persistence, otherwise known as Hurst-Kolmogorov (HK) behaviour (coined by Koutsoyiannis & Cohn (2008)), i.e. power-law decay of the autocorrelation function with lag (for a literature review on long-term persistent processes in hydrometeorology, see also O'Connell *et al.* (2016)). Additionally, Koutsoyiannis (2011) provided a theoretical justification of the HK behaviour in geophysical processes, showing that it is linked to the second law of

thermodynamics (i.e. entropy extremization), and specifically, the stronger the persistence of the dependence structure of a process, the higher the entropy of the process at large scales.

The identification of the dependence structure of a process can be highly affected by the sample uncertainty and statistical bias where the true statistical properties (mean, variance etc.) of a statistic (e.g. variance) of a stochastic process may differ from the one estimated from a series with finite length. The deviations of the statistical characteristics from their true values should be taken into account not only for the marginal characteristics but also for the dependence structure of the process. Therefore, to correctly adjust the stochastic model to the observed series of the physical process, we should account for the bias effect since all series are of finite (and often short) lengths.

The second-order properties can be similarly assessed by common stochastic tools such as the autocovariance function (a function of lag), power spectrum (a function of frequency), and variation of statistics (e.g. variance) of the averaged process vs. scale, a tool known as climacogram (Koutsoyiannis 2010). It is shown that the latter estimator of the second-order dependence structure, as compared to the other two metrics, encompasses additional advantages in stochastic model building and interpretation from data; for example, it is characterized by smaller statistical uncertainty and easier to handle expressions of the statistical bias (Dimitriadis & Koutsoyiannis 2015). Therefore, it is advisable that the sample uncertainty of the second-order dependence structure be tackled with the estimator with the lower variation, such as the climacogram. When multiple sample realizations (i.e. recorded series) are known, the handling of the statistical bias arising from a selected stochastic model may be based on the unbiased estimator of the expected value of the climacogram (Dimitriadis & Koutsoyiannis 2018). However, when a single data series of observations is available for the model fitting (which is the case when geophysical processes are studied), it would be interesting to examine the mode of the climacogram, instead of the expected value; the two may differ in case of strong HK behaviour. This estimator is equivalent to a maximum-likelihood estimator (e.g. Kendzioriski *et al.* 1999) for processes with zero (i.e. white noise) or short-term (e.g. Markov) dependence structure, while here we further extend it for HK processes (see also the work of

Tyralis & Koutsoyiannis (2011) for the expectation of the climacogram). It is noted that while the climacogram is often based on the second central moment (i.e. variance), other types of moments (e.g. raw, L-moments or K-moments; Koutsoyiannis 2019) can be used to measure fluctuation in scale, and here, we focus on the central second-order climacogram (i.e. fluctuation measured by variance vs. scale).

## METHODS

In this section, we present the applied methods, namely the climacogram estimator, the statistical bias expressions for the mode and expected values of the estimator and the algorithm for the stochastic synthesis of the Gaussian HK process for the Monte Carlo analysis.

### The climacogram

The analysis of a process through the variance of the averaged process vs. scale has been thoroughly applied in stochastic processes (e.g. Papoulis 1991; Vanmarcke 2010). However, its importance to the analysis of the second-order dependence structure is highlighted mainly by more recent works (see a historical review in Koutsoyiannis (2018)). Also, the simple name *climacogram* allowed its further understanding through visualization; indeed, the term originates from the Greek *climax* (meaning scale) and *gramma* (meaning written; cf. the terms autocorrelogram for the autocorrelation, scaleogram for the power spectrum and wavelets).

It has been shown that the climacogram, treated as an estimator (rather than just a tool for the identification of long-term behaviour of the second-order dependence structure), has additional advantages from the more widely applied estimators of the autocovariance and power spectrum (Dimitriadis & Koutsoyiannis 2015). Namely, the climacogram could provide a more direct, easy, and accurate means to make diagnoses from data and build stochastic models in comparison to the power spectrum and autocovariance. For example, the climacogram, compared to other tools, has the lowest standardized estimation error for processes with short- and long-term persistence, zero discretization error for averaged processes,

simple and analytical expression for the statistical bias, always positive values, well-defined and usually monotonic behaviour, smallest fluctuation of skewness on small scales while closest to zero skewness in larger scales, and mode closest to the expected (i.e. mean) value in large scales. Also, the climacogram is directly linked to the entropy production of a process (Koutsoyiannis 2011, 2016). Furthermore, the climacogram expands the notion of *variance* by making it a function of *time scale* and is *per se* further expandable for statistics different from the central estimators of fluctuation (e.g. second raw moment and second L-moment vs. scale; Koutsoyiannis 2019) for different characteristics of the estimator (e.g. mode and median) and even for moments of higher (e.g. third and fourth) orders (Dimitriadis & Koutsoyiannis 2018). Recently, Koutsoyiannis (2019) extended the notion of climacogram for orders higher than two and showed how to substitute the joint moments of a process, allowing in this manner to tackle some limitations of the latter such as the discretization effect and statistical bias.

Symbolically, the climacogram is:

$$\gamma(k) := \text{var}[\underline{x}(k)] \tag{1}$$

where  $\text{var}[\ ]$  denotes the variance and  $\underline{x}(k) := 1/k \int_0^k \underline{x}(t)dt$  is the continuous-time process at scale  $k$  (in dimensions of time), which equals the discrete-one averaged in time intervals  $\Delta$ , i.e.  $\underline{x}_\kappa := 1/\kappa \sum_{i=1}^{\kappa} \underline{x}_i$ , in the discrete-time scale  $\kappa = k/\Delta$  (dimensionless natural number whereas for real numbers see the adjustment in Koutsoyiannis (2011)).

### The Gaussian long-term persistent process and its stochastic synthesis

The most common processes employed in geophysics, and particularly in hydrology, are the white noise process, the Markov process (with an exponential decay of the autocorrelation), and long-term persistent processes, which are characterized by a power-law decay of the climacogram (or equivalently of the autocorrelation) as a function of scale (or lag). A typical representative of the latter processes is the Gaussian HK process defined as follows:

$$\underline{x}(k) - \mu =_d (k)^{(H-1)}(\underline{x}(1) - \mu) \tag{2}$$

where  $=_d$  denotes equality in distribution with  $\mu$  the mean and  $\gamma(k) = \gamma(\Delta)/\kappa^{2-2H}$  the variance of the process for each scale  $k$ ,  $H$  is the Hurst parameter ( $0 < H < 1$ ) otherwise defined as (Dimitriadis *et al.* 2016a)  $H := 1 + \frac{1}{2} \lim_{k \rightarrow \infty} d \ln(\gamma(k))/d \ln k$ ; the quantity in the limit is the derivative of  $\ln(\gamma(k))$  with respect to  $\ln(k)$ .

It is noted that this process has infinite variance at scale zero and thus, it should not be used to model the small scales of a physical process (in spite of the fact that the fractional-Gaussian-noise-fGn-process is widely used to model several processes at small scales; Koutsoyiannis *et al.* 2018). For the stochastic synthesis of the Gaussian HK model, we may use the sum of arbitrarily many independent Markov processes, thus expressing the target climacogram as follows (Dimitriadis & Koutsoyiannis 2015):

$$\gamma(\kappa\Delta) = \sum_{i=1}^l \frac{2\lambda_i}{(\kappa\Delta/q_i)^2} (\kappa\Delta/q_i + e^{-\kappa\Delta/q_i} - 1) \tag{3}$$

where  $\lambda_i$  is the variance,  $q_i$  a time scale parameter for each Markov model  $i$ , and  $l$  the total number of Markov processes. Mandelbrot (1963) has shown that for  $l \rightarrow \infty$ , the above model can adequately describe an fGn (or else Gaussian HK) process for any generated length (see also Mandelbrot & Van Ness 1968; Mandelbrot & Wallis 1968). Koutsoyiannis (2002) has analytically estimated the parameters of three AR(1) models ( $l=3$ ) to capture the HK process for  $n < 10^4$ . Dimitriadis & Koutsoyiannis (2015) have expanded this framework to the sum of arbitrarily many AR(1) models (abbreviated as SAR) for the generation of any type of process with an autoregressive dependence structure and up to any number of scales, by using a suitable function with only two parameters, namely  $p_1$  and  $p_2$ , that link the lag-1 autocorrelations of each Markov model, e.g. through the expression  $q_i = p_1 p_2^{i-1}$ , with  $i = 1, \dots, l$  and  $l$  often taken equal to the integer part of  $\log(n) + 1$ . For example, for  $n = 10^6$  and  $H = 0.8$ , we have  $l = 7$ ,  $p_1 = 0.394$  and  $p_2 = 12.356$  for a maximum standardized error between the true  $\gamma_t$  (Equation (2)) and modelled  $\gamma_m$  (Equation (3)) climacogram (i.e.  $\max|\gamma_t - \gamma_m|/\gamma_t$  for all scales) equal to 0.009 (Table 1).

**Table 1** | Parameters  $p_1$  and  $p_2$  estimated to approximate different types of the  $N(0,1)$ -HK model (i.e.  $\mu = 0$  and  $\gamma(\Delta) = 1$ ) with  $l = 7$  and  $n \leq 10^6$

H	$p_1$	$p_2$	Maximum error (standardized)
0.51	0.022	17.122	0.001
0.60	0.091	12.607	0.006
0.70	0.124	13.317	0.009
0.80	0.394	12.356	0.009
0.90	0.395	14.708	0.005
0.99	0.548	19.465	0.001

**The mode of climacogram estimator and its statistical bias**

The climacogram can be estimated from a sample through an estimator as similarly done for the estimators of the marginal moments. Here, for the climacogram we use a classical estimator:

$$\hat{\gamma}(\kappa\Delta) = \frac{1}{[n/\kappa] - 1} \sum_{i=1}^{[n/\kappa]} (\bar{x}_i^{(\kappa)} - \bar{x})^2 \tag{4}$$

where  $[n/\kappa]$  is the integer part of  $n/\kappa$ ,  $\bar{x}_i^{(\kappa)} = \sum_{l=\kappa(i-1)+1}^{\kappa i} x_l / \kappa$  is the averaged process at scale  $\kappa = k/\Delta$  for  $i \in [1, [n/\kappa]]$ ,  $\bar{x} = \sum_{l=1}^n x_l / n = \bar{x}_1^{(n)}$  is the sample average, and  $n$  is the series length.

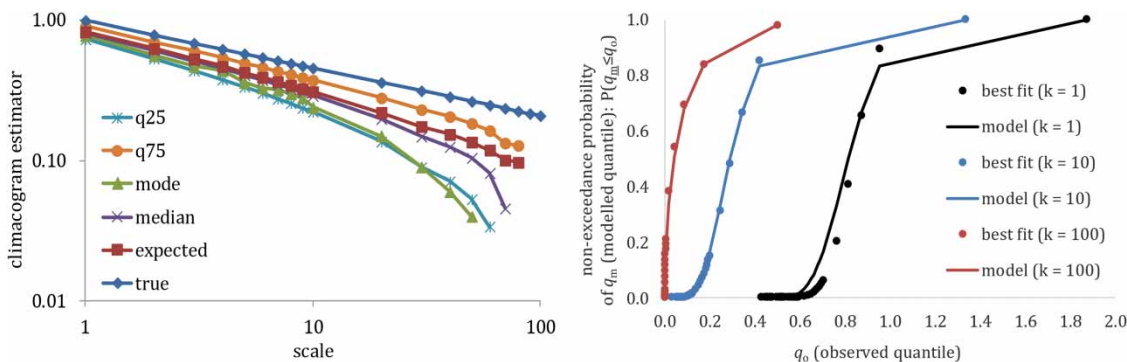
Since the above estimator is a random variable, it has a marginal distribution (see an illustration in Figure 1). The

true value of a statistical characteristic (e.g. variance) of a stochastic model may differ from the one estimated from a series with finite length. To correctly adjust the stochastic model to the observed series of the physical process, one should account for the bias effect. An important question is how the statistical bias is generally handled through the second-order dependence structure in case of long-term persistent processes. Particularly, the selected stochastic model should be adjusted for bias before it is fitted to the sample dependence structure. It is noted that neglecting the bias effect in case of a long-term persistent process leads to underestimations of the stochastic model parameters such as the Hurst parameter and to erroneous conclusions. An adjustment of the models for bias is usually done by equating the observed dependence structure to the expected value of the applied estimator. The alternative studied here is the mode, instead of the expected value, of the dependence structure, which represents the most probable value (and thus, the most expected) of the variance estimator at each scale.

The statistical bias of an estimator is the difference of the expected value of the estimator from its true value (e.g. Papoulis 1991). Thus, the bias of the climacogram is shown as follows (e.g. Koutsoyiannis 2011):

$$B_E[\hat{\gamma}(\kappa\Delta)] = E[\hat{\gamma}(\kappa\Delta)] - \gamma(\kappa\Delta) = \frac{(\kappa/n)\gamma(\kappa\Delta) - \gamma(n\Delta)}{1 - \kappa/n} \tag{5}$$

where  $B_E[\ ]$  denotes the bias of the expected value of a statistical estimator of a process. Clearly, for the mean value of a process, we have that  $B_E[\hat{\mu}] = E\left[\sum_{i=1}^n x_i / n\right] - \mu = 0$ .



**Figure 1** | An illustration for an  $N(0,1)$ -HK ( $H = 0.83$ ,  $n = 200$ ) process of (left) how several statistical characteristics of the climacogram estimator vary with scale and (right) the observed quantile ( $q_o$ ) vs. the non-exceedance probability of the modelled quantile  $P(q_m \leq q_o)$ , showing how the gamma distribution can adequately approximate the distribution of the climacogram estimator especially at large scales.

Following the same rationale, we define an expansion of the notion of bias for the mode of the above estimator of the climacogram, i.e.:

$$B_M[\hat{\gamma}(\kappa\Delta)] = M[\hat{\gamma}(\kappa\Delta)] - \gamma(\kappa\Delta) \quad (6)$$

where  $M[\underline{x}] := \arg \max [f(\underline{x})]$  denotes the mode of the variable  $\underline{x}$  with density function  $f(x)$ . We refer to  $B_M[\ ]$  as the mode bias.

For a Gaussian white noise process of length  $n$  and variance  $\gamma(1)$ , the distribution of its sample variance follows the gamma distribution  $\Gamma((n-1)/2, 2\gamma(1)/(n-1))$  (Cochran 1934). The averaged process at scale  $\kappa$ , with a sample length of  $n/\kappa$  and variance  $\gamma(k) = \gamma(\Delta)/\kappa$ , follows  $\Gamma((n/\kappa-1)/2, 2\gamma(\Delta)/(n-\kappa))$ , with  $M[\hat{\gamma}(\kappa\Delta)] = \gamma(\Delta) \frac{n-3\kappa}{\kappa(n-\kappa)}$  for  $n/\kappa \geq 3$ , or else 0. Hence, for  $n/\kappa \gg 3$ , we have that  $(n-3\kappa)/(n-\kappa) \approx 1$ , and  $M[\hat{\gamma}(\kappa\Delta)] \approx E[\hat{\gamma}(\kappa\Delta)] = \gamma(\Delta)/\kappa = \gamma(\kappa\Delta)$ , i.e. zero bias. However, for long-term persistent processes, the mode bias is non-zero and its analytical solution is no longer easy to derive.

From the above results, it becomes evident that the statistical bias always depends on the selected model and not on the data as commonly thought. For example, consider the Gaussian HK process in the previous section with an autocorrelation function in discrete time  $\rho_v = 1/2(|v+1|^{2H} + |v-1|^{2H}) - |v|^{2H}$ , where  $v$  is the discrete-time lag. The bias of the autocorrelation is similarly defined as  $B_E[\hat{\rho}(v)] = E[\hat{\rho}(v)] - \rho(v)$ , and thus depends on the model parameter  $H$ . It is noted that the above apply even to the so-called non-parametric models, since they also involve estimation from data, and thus, these models should be similarly adjusted for statistical bias to avoid underestimation of the process variability during a Monte Carlo simulation.

For simplicity and without loss of generality, we set  $\Delta=1$  for the rest of the analysis. It is evident that  $B_M[\hat{\gamma}(\kappa)] \leq B_E[\hat{\gamma}(\kappa)] \leq 0$  or else  $|B_E[\hat{\gamma}(\kappa)]| \leq |B_M[\hat{\gamma}(\kappa)]|$ , since the sample variance is positively skewed, i.e.  $E[\hat{\gamma}(\kappa)] \geq M[\hat{\gamma}(\kappa)]$  and the equality holds when  $n \rightarrow \infty$ , where the variance of the sample variance is zero for an ergodic process. A preliminary analysis of common HK-type processes has shown that the mode climacogram is close to the low quartile (25% quantile) of the marginal

distribution of variance at each scale (Dimitriadis *et al.* 2016c; Gournary 2017). Therefore, when the mode of the variance estimator is of interest, we may use a Monte Carlo technique (as described in the next section) to accurately estimate the mode bias or, in case the marginal distribution of the climacogram is known, to calculate the 25% quantile at each scale to approximate the mode bias.

## MONTE CARLO ANALYSIS FOR THE MODE OF THE VARIANCE ESTIMATOR

We perform Monte Carlo experiments over the  $N(0,1)$ -HK model for a wide range of Hurst parameters  $H$  (i.e. 0.5 to 0.95) and for a wide range of series lengths  $n$  (i.e. 20–2,000). Specifically, we produce a number ( $N$ ) of synthetic series through the SAR model described in the section ‘The Gaussian long-term persistent process and its stochastic synthesis’, where  $N$  depends on the sample mean value to reach the expected one at scale  $\kappa=n/10$  based on the rule of thumb when using the climacogram as shown in Dimitriadis & Koutsoyiannis (2015). We found that for  $N \approx 10^6/n^{2-2H}$ , the standardized error between the theoretical expected value and the sample one (Equation (5)) is lower than 1% at scale  $\kappa=n/10$ . In this way, the mode is expected also to be well preserved with a similar error. However, caution should be given to the selection of the sample mode estimator to ensure that its variance permits a robust estimation of the true value of the mode. Since the distribution function of the estimator of variance is unknown for long-term persistent processes and given that the mode value is the most likely to occur within the sample, we calculate the sample mode from each simulated series by finding the most probable value with an accuracy of two decimal digits. Specifically, we round up each value of the time series, and for each scale, to the second decimal digit, and we estimate the most probable value of the rounded time series (for higher accuracies a larger  $N$  was required). Also, other estimators for the sample mode (e.g. Bickel & Fruwirth 2006) could be used and compared to the proposed one in future research to optimize the performance of the analysis.

Here, to derive an easy-to-fit empirical expression to approximate the mode bias, we adopt the assumption that the above distribution is nearly gamma for smaller scales

(see also a similar analysis in Gournary (2017) and Dimitriadis et al. (2018)). Using the results from the Monte Carlo analysis, we then evaluate the parameter of the gamma distribution for each  $H, n$ , and  $\kappa$ , and we build a model for the mode, then later test its performance. Although the true autocorrelation function of the averaged process for a long-term persistent process does not vary with scale, the sample autocorrelation will be also prone to bias (e.g. Dimitriadis & Koutsoyiannis 2015) affecting the distribution function of the sample variance at each scale. To minimize the sample error for the fitting of the two-parameter gamma distribution, we use the theoretical expression for the expected value of the sample climacogram, i.e.  $E[\hat{\gamma}(\kappa)]$ , and the variance of the sample climacogram, i.e.  $\text{Var}[\hat{\gamma}(\kappa)]$ , as evaluated from the Monte Carlo analysis, which exhibits the lowest variability in estimation among the four central moments (Dimitriadis & Koutsoyiannis 2018; Figure 2). Based on these two measures, we estimate the two parameters of the gamma distribution.

We first set the scale parameter of the gamma distribution such as to simulate the sample ratio of the aforementioned parameters, i.e.  $b(H, n, \kappa) = \text{Var}[\hat{\gamma}(\kappa)]/E[\hat{\gamma}(\kappa)]$  and so, the shape parameter can be also estimated as  $a(H, n, \kappa) = E[\hat{\gamma}(\kappa)]/b(H, n, \kappa)$ .

We observe (Figure 2) that for  $a(H, n, \kappa) > 1$ , the shape parameter  $a(H, n, \kappa)$  is approximately proportional, by a function  $c(H)$ , to the corresponding shape parameter for the white noise process  $a(0.5, n, \kappa) = (n/\kappa - 1)/2$  raised to a function  $p(H)$ , i.e.:

$$a(H, n, \kappa) = c(H)((n/\kappa - 1)/2)^{p(H)} \tag{7}$$

where  $a(H, n, \kappa) > 1$  is a function corresponding to the shape parameter of the gamma distribution function, while for  $a(H, n, \kappa) \leq 1$  or  $\kappa \geq n/3$ , the mode is considered close to zero.

The two functions of the above expression are fitted as follows (Figure 2):

$$c(H) = 1.68(H - 0.5)^2 - 0.3025(H - 0.5) + 1 \tag{8}$$

and

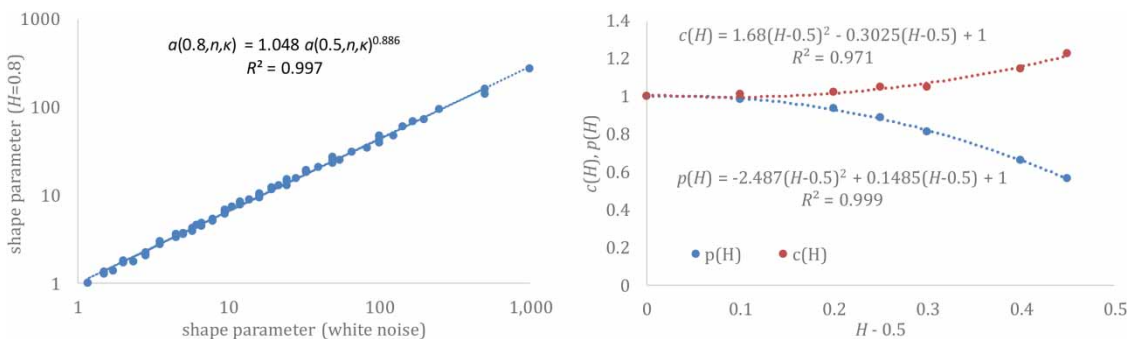
$$p(H) = -2.4865(H - 0.5)^2 + 0.1485(H - 0.5) + 1 \tag{9}$$

The above two adjustments allow us to empirically express the mode of the climacogram estimator as a function of  $H, n$ , and  $\kappa$ :

$$\begin{aligned} M[\hat{\gamma}(\kappa)] &= (a(H, n, \kappa) - 1)b(H, n, \kappa) \\ &= (1 - 1/a(H, n, \kappa))E[\hat{\gamma}(\kappa)] \end{aligned} \tag{10}$$

It is noted that based on the above assumptions, the standard deviation, and the skewness and excess kurtosis coefficients of the climacogram estimator can be estimated as  $b(H, n, \kappa)\sqrt{a(H, n, \kappa)}$ ,  $2/\sqrt{a(H, n, \kappa)}$ , and  $6/a(H, n, \kappa)$ , respectively. Since  $a(H, n, \kappa) \leq a(0.5, n, \kappa)$  all the above measures will be larger than those in case of a white noise process.

The above expression can approximate the mode by an absolute difference of 0.005 from the Monte Carlo estimates, while for better approximations it is advised to implement a new Monte Carlo analysis (see also discussion and application in the section ‘Applications to annual



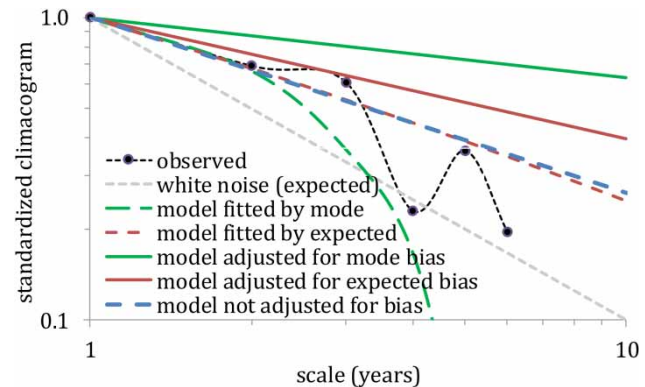
**Figure 2** | (left) The shape parameter assuming a gamma distribution for the mode estimator of the climacogram of an  $N(0,1)$ -HK process (for  $H = 0.8$  and for all  $n$  and  $\kappa$  simulated in the Monte Carlo analysis) vs. the theoretical shape parameter of the white noise process. (right) Proposed model for the  $c(H)$  and  $p(H)$  functions for all examined  $H$  from the Monte Carlo analysis.

streamflow). Interestingly, the standardized error between the mode and expected values of the estimator, i.e.  $\varepsilon = \left| \frac{E[\hat{\gamma}(\kappa)] - M[\hat{\gamma}(\kappa)]}{E[\hat{\gamma}(\kappa)]} \right|$ , is calculated from the Monte Carlo analysis to reach a maximum value of 67% corresponding to cases with  $H \geq 0.6$  and  $n/\kappa \leq 10$ , while for the white noise process it can be theoretically estimated as  $\varepsilon = 2/(n/\kappa - 1)$ , which for  $\kappa = n/10$  is approximately 20%.

## APPLICATIONS TO ANNUAL STREAMFLOW

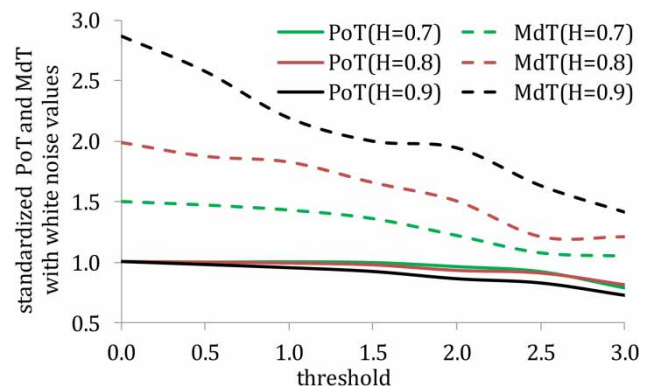
For illustrations of possible implications of the above results, we apply a stochastic analysis based on the expected and the mode values of the climacogram to a streamflow process at the Peneios river (Thessaly, Greece), where a historical streamflow annual time series is available at the upstream station of Ali Efenti with only a 13-year length (for more information on the study area, see Dimitriadis *et al.* (2016b)). For the identification of the stochastic model, we adjust for statistical bias and, in particular, we fit the mode of the estimator rather than its expectation. It is noted that the proposed empirical model for the mode bias (Equation (10)) is derived from a Monte Carlo analysis for sample lengths of  $n \geq 20$ , and so for this application, we perform a new Monte Carlo analysis to fit the observed climacogram for scales  $1 \leq \kappa \leq n/10$  (rule of thumb; Dimitriadis & Koutsoyiannis 2015) and so here, for the first two scales (Figure 3). We find that an HK model can adequately simulate the observed standardized climacogram, i.e.  $\hat{\gamma}(\kappa)/\hat{\gamma}(1)$ , with  $H = 0.9$ . We also estimate the Hurst parameter with the expectation of the estimator, and we find  $H' \approx 0.8$  and  $H'' \approx 0.7$ , with or without adjusting for bias, respectively. Evidently, both latter values underestimate the long-term persistence behaviour (Figure 3).

It is noted that the dependence structure of a process (e.g. streamflow) will have a small effect at the risk imposed by the expected number of peaks over threshold (e.g. for the design of a dam or for flood risk mapping) as compared to the effect of the marginal distribution of the process (Volpi *et al.* 2015; Serinaldi & Kilsby 2018). However, the dependence structure will have a great effect (especially for processes with long-term behaviour) at the duration of successive peaks over threshold (e.g. maximum duration of wet/dry periods or of flood inundation), which may highly affect urban as well as agricultural areas and insurance



**Figure 3** | Standardized climacogram estimations of the observed standardized time series (black line), the white noise model (grey line), and the three fitted  $N(0,1)$ -HK stochastic processes: (a) adjusting for bias of the mode of the estimator (green line), i.e.  $M[\hat{\gamma}(\kappa)]/M[\hat{\gamma}(1)]$ , and of its expectation (red line), i.e.  $E[\hat{\gamma}(\kappa)]/E[\hat{\gamma}(1)]$ , and (b) not adjusting for bias (blue line), i.e.  $\hat{\gamma}(\kappa) = \gamma(\kappa)$ , also corresponding to the non-parametric model configuration. Please refer to the online version of this paper to see this figure in colour: <http://dx.doi.org/10.2166/hydro.2019.038>.

policies (e.g. Serinaldi & Kilsby 2016; Goulianou *et al.* 2019). To illustrate this, we generate an adequate number  $N$  (see the section ‘Monte Carlo analysis for the mode of the variance estimator’) of HK synthetic timeseries with  $H = 0.5$  ( $N = 5 \times 10^3$ ),  $H = 0.7$  ( $N = 4 \times 10^4$ ),  $H = 0.8$  ( $N = 10^5$ ), and  $H = 0.9$  ( $N = 3 \times 10^5$ ). For convenience, we assume an  $N(0,1)$  distribution for all processes. We then estimate the expected frequency of the number of peaks over various thresholds (PoT) as well as the expected frequency of the maximum duration of successive peaks over various thresholds (MdT), and we standardize them with the PoT and MdT values of the white noise process (Figure 4).



**Figure 4** | Expected frequency of peak over threshold (PoT) and expected maximum duration of successive peaks over threshold (MdT) standardized with the PoT and MdT values of the  $N(0,1)$  white noise process for various HK- $N(0,1)$  processes.

We find that the MdT varies with threshold and long-term persistence, while the PoT stays almost unaffected by both. Additional analyses and quantifications on the reflection of long-term term persistence in terms of clustering in time can be found in Iliopoulou & Koutsoyiannis (2019).

The results from this study suggest that the sample estimator of the variance can be skewed even for long samples in the presence of long-term persistence behaviour as opposed to the white noise process. Therefore, the mode is different from the expectation and more suitable to use in estimation. We propose that when a single recorded series is available and a Gaussian HK process is fitted with small sample size and relatively high Hurst parameter, it is advantageous to employ the mode of the estimator as calculated from the empirical model of Equation (10), rather than its expectation (Equation (5)), so as to avoid underestimation of the Hurst parameter (and thus, the uncertainty of the process). In case of a non-Gaussian distribution, larger accuracy, or a different estimator of the second-order dependence structure (e.g. other climacogram estimator, autocovariance, power spectrum, variogram etc.), we should employ the Monte Carlo technique and test whether the mode of the estimator used is close enough to its expected value. If this is true, then the expected value can be used to adjust the model for bias, whereas if the two values vary then the model should be adjusted for bias based on the mode estimator. For Monte Carlo analysis of a non-Gaussian-correlated process, an explicit algorithm should be preferred (Dimitriadis & Koutsoyiannis 2018) since the mode value is expected to highly depend on higher-order moments in case of long-term persistent processes.

## CONCLUSIONS AND DISCUSSION

Awareness of uncertainty in assessing the dependence structure of a process is of paramount importance as it may critically affect the interpretation of results. Estimation uncertainty may introduce large statistical bias, which can be additionally magnified in the presence of long-term persistence (Dimitriadis & Koutsoyiannis 2015). In addition, if the uncertainty is underestimated, then a regular cluster of events could be erroneously regarded as an extreme cluster. Although the mode of the examined classical estimator for

variance is close to its expectation for small Hurst parameters and large lengths, we show that for larger values of the Hurst parameter and small sample lengths, equating the expected climacogram to the observed one may lead to underestimation of the long-term persistence and thus the uncertainty of the process.

We propose that when the available series have short lengths or when the empirical Hurst parameter is estimated larger than 0.5, we should always account for statistical bias. Particularly for the bias adaptation, when information is available on only a single series/realization of the process, it is advantageous to equate the mode instead of the expectation of the climacogram estimator to the sample values. Interestingly, in case of an  $N(0,1)$ -HK process, the absolute difference between the mode and expected values of the estimator is calculated (from a Monte Carlo analysis performed in this study) to reach a maximum value of 67% of the expected value, corresponding to cases with  $H \geq 0.6$  and  $n/\kappa \leq 10$ , while for the white noise process the value is approximately 20% for  $\kappa = n/10$ . In cases of different stochastic processes or estimators or when a larger accuracy of the mode bias is of interest, one should employ a Monte Carlo technique through an explicit generation algorithm (Dimitriadis & Koutsoyiannis 2018) to estimate the mode climacogram estimator or use the lower quartile (25% quantile) of the estimator (in case its distribution is known) as an approximation.

From the Monte Carlo analysis performed in this study, it is also observed that for an  $N(0,1)$ -HK process with variance  $\gamma(\kappa) = \gamma(1)/\kappa^{2-2H}$  and for large  $n$  and small  $n/\kappa$ , the distribution of the climacogram estimator tends to that of  $\Gamma((n/\kappa - 1)/2, 2\gamma(\kappa)/(n/\kappa - 1))$ , with a mean value equal to  $\gamma(\kappa)$ , i.e. zero bias. However, given the estimation uncertainty present in records exhibiting persistence, the autocorrelation of the averaged process is independent of the scale, and thus, the above distribution will never be truly reached. The underestimation of the persistence of the parent process also has critical implications for the estimation of the properties of its extremes, since it was shown that the maximum duration of successive peaks over threshold is greatly affected by the degree of dependence. Additional analyses and quantifications on the reflection of long-term term persistence in terms of clustering in time can be found in Iliopoulou & Koutsoyiannis (2019).



A final remark for discussion, considering the etymology of the terms, is that the expected value of a random process is less expected to occur than its mode (i.e. most probable value; a term coined by Pearson (1895, p. 345)), where the two coincide only in symmetrical distributions. Therefore, when only one value is known (here, only one realization of the climacogram estimator), it is more accurate to fit the model and evaluate the Hurst parameter based on the proposed mode estimator rather than the expected one.

## ACKNOWLEDGEMENT

The authors would like to thank the editor Luigi Berardi for handling the paper, one anonymous reviewer for useful comments, and Federico Lombardo for his fruitful discussion, comments, and suggestions that helped us improve the paper.

## CODE AVAILABILITY

The MATLAB script for the SAR generation algorithm is available as well as the script for a fast estimation algorithm of the sample climacogram in very long timeseries and in many scales.

## REFERENCES

- Bickel, D. R. & Fruwirth, R. 2006 [On a fast, robust estimator of the mode: comparisons to other robust estimators with applications](#). *Computational Statistics & Data Analysis* **50**, 3500–3530.
- Cochran, W. G. 1934 [The distribution of quadratic forms in a normal system, with applications to the analysis of covariance](#). *Mathematical Proceedings of the Cambridge Philosophical Society* **30** (2), 178–191. doi:10.1017/S0305004100016595.
- Dimitriadis, P. 2017 [Hurst-Kolmogorov Dynamics in Hydrometeorological Processes and in the Microscale of Turbulence](#). PhD Thesis, National Technical University of Athens, p. 167.
- Dimitriadis, P. & Koutsoyiannis, D. 2015 [Climacogram versus autocovariance and power spectrum in stochastic modelling for Markovian and Hurst-Kolmogorov processes](#). *Stochastic Environmental Research & Risk Assessment* **29** (6), 1649–1669.
- Dimitriadis, P. & Koutsoyiannis, D. 2018 [Stochastic synthesis approximating any process dependence and distribution](#). *Stochastic Environmental Research & Risk Assessment* **32** (6), 1493–1515. doi:10.1007/s00477-018-1540-2.
- Dimitriadis, P., Koutsoyiannis, D. & Papanicolaou, P. 2016a [Stochastic similarities between the microscale of turbulence and hydrometeorological processes](#). *Hydrological Sciences Journal* **61** (9), 1623–1640. doi:10.1080/02626667.2015.1085988.
- Dimitriadis, P., Tegos, A., Oikonomou, A., Pagana, V., Koukouvinos, A., Mamassis, N., Koutsoyiannis, D. & Efstratiadis, A. 2016b [Comparative evaluation of 1D and quasi-2D hydraulic models based on benchmark and real-world applications for uncertainty assessment in flood mapping](#). *Journal of Hydrology* **534**, 478–492.
- Dimitriadis, P., Gournari, N. & Koutsoyiannis, D. 2016c [Markov vs. Hurst-Kolmogorov behaviour identification in hydroclimatic processes](#). *European Geosciences Union General Assembly*. Geophysical Research Abstracts, Vol. 18, European Geosciences Union, Vienna, EGU2016-14577-4. doi:10.13140/RG.2.2.21019.05927.
- Dimitriadis, P., Gournary, N., Petsiou, A. & Koutsoyiannis, D. 2018 [How to adjust the fGn stochastic model for statistical bias when handling a single time series; application to annual flood inundation](#). In: *13th Hydroinformatics Conference*, 1–6 July 2018, Palermo, Italy.
- Goulianou, T., Papoulakos, K., Iliopoulou, T., Dimitriadis, P. & Koutsoyiannis, D. 2019 [Stochastic characteristics of flood impacts for agricultural insurance practices](#). *European Geosciences Union General Assembly 2019*, Geophysical Research Abstracts, Vol. 21, European Geosciences Union, Vienna, EGU2019-5891.
- Gournary, N. 2017 [Probability Distribution of the Climacogram Using Monte Carlo Techniques](#). Diploma Thesis, Department of Water Resources and Environmental Engineering, National Technical University of Athens, Athens (in Greek), p. 108.
- Iliopoulou, T. & Koutsoyiannis, D. 2019 [Revealing hidden persistence in maximum rainfall records](#). *Hydrological Sciences Journal*. doi.org/10.1080/02626667.2019.1657578.
- Kendzioriski, C. M., Bassingthwaite, J. B. & Tonellato, P. J. 1999 [Evaluating maximum likelihood estimation methods to determine the Hurst coefficient](#). *Physica A* **273** (3–4), 439–451.
- Koutsoyiannis, D. 2002 [The Hurst phenomenon and fractional Gaussian noise made easy](#). *Hydrological Sciences Journal* **47** (4), 573–595.
- Koutsoyiannis, D. 2010 [HESS opinions ‘A random walk on water’](#). *Hydrology and Earth System Sciences* **14**, 585–601.
- Koutsoyiannis, D. 2011 [Hurst-Kolmogorov dynamics as a result of extremal entropy production](#). *Physica A: Statistical Mechanics and its Applications* **390** (8), 1424–1432.
- Koutsoyiannis, D. 2016 [Generic and parsimonious stochastic modelling for hydrology and beyond](#). *Hydrological Sciences Journal* **61** (2), 225–244.
- Koutsoyiannis, D. 2018 [Climate Change Impacts on Hydrological Science: A Comment on the Relationship of the Climacogram with Allan Variance and Variogram](#). ResearchGate.

- Koutsoyiannis, D. 2019 Knowable moments for high-order stochastic characterization and modelling of hydrological processes. *Hydrological Sciences Journal* **64** (1), 19–33.
- Koutsoyiannis, D. & Cohn, T. A. 2008 The Hurst phenomenon and climate (solicited), *European Geosciences Union General Assembly 2008, Geophysical Research Abstracts*, Vol. 10, Vienna, 11804, European Geosciences Union. doi:10.13140/RG.2.2.13303.01447.
- Koutsoyiannis, D., Dimitriadis, P., Lombardo, F. & Stevens, S. 2018 From fractals to stochastics: seeking theoretical consistency in analysis of geophysical data. In: *Advances in Nonlinear Geosciences* (A. A. Tsonis, ed.). Springer, Cham, Switzerland, pp. 237–278.
- Lombardo, F., Volpi, E., Koutsoyiannis, D. & Papalexiou, S. M. 2014 Just two moments! A cautionary note against use of high-order moments in multifractal models in hydrology. *Hydrology and Earth System Sciences* **18**, 243–255. doi:10.5194/hess-18-243-2014.
- Mandelbrot, B. B. 1963 The variation of certain speculative prices. *The Journal of Business* **36**, 394–419.
- Mandelbrot, B. B. & Van Ness, J. W. 1968 Fractional Brownian motions, fractional noises and applications. *SIAM Review* **10**, 422–437.
- Mandelbrot, B. B. & Wallis, J. R. 1968 Noah, Joseph and operational hydrology. *Water Resource Research* **4**, 909–918.
- O'Connell, P. E., Koutsoyiannis, D., Lins, H. F., Markonis, Y., Montanari, A. & Cohn, T. A. 2016 The scientific legacy of Harold Edwin Hurst (1880–1978). *Hydrological Sciences Journal* **61** (9), 1571–1590. doi:10.1080/02626667.2015.1125998.
- Papoulis, A. 1991 *Probability, Random Variables and Stochastic Processes*, 3rd edn. McGraw-Hill, New York.
- Pearson, K. 1895 Contributions to the mathematical theory of evolution – II, Skew variation in homogeneous material. *Philosophical Transactions of the Royal Society of London* **186**, 343–414. Available from: <https://royalsocietypublishing.org/doi/pdf/10.1098/rsta.1895.0010>.
- Serinaldi, F. & Kilsby, C. G. 2016 Understanding persistence to avoid underestimation of collective flood risk. *Water* **8**, 152.
- Serinaldi, F. & Kilsby, C. G. 2018 Unsurprising surprises: the frequency of record-breaking and over-threshold hydrological extremes under spatial and temporal dependence. *Water Resources Research* **54** (9), 6460–6487.
- Tyralis, H. & Koutsoyiannis, D. 2011 Simultaneous estimation of the parameters of the Hurst-Kolmogorov stochastic process. *Stochastic Environmental Research & Risk Assessment* **25** (1), 21–33.
- Vanmarcke, E. 2010 *Random Fields: Analysis and Synthesis*. World Scientific, New Jersey, USA.
- Volpi, E., Fiori, A., Grimaldi, S., Lombardo, F. & Koutsoyiannis, D. 2015 One hundred years of return period: strengths and limitations. *Water Resources Research* **51** (10), 8570–8585.

First received 11 February 2019; accepted in revised form 6 August 2019. Available online 16 September 2019