

Evaluation of statistical methods for quantifying fractal scaling in water-quality time series with irregular sampling

Journal Article**Author(s):**

Zhang, Qian; Harman, Ciaran J.; Kirchner, James W.

Publication date:

2018

Permanent link:

<https://doi.org/10.3929/ethz-b-000242983>

Rights / license:

[Creative Commons Attribution 3.0 Unported](#)

Originally published in:

Hydrology and Earth System Sciences 22(2), <https://doi.org/10.5194/hess-22-1175-2018>



Evaluation of statistical methods for quantifying fractal scaling in water-quality time series with irregular sampling

Qian Zhang¹, Ciaran J. Harman², and James W. Kirchner^{3,4,5}

¹University of Maryland Center for Environmental Science, US Environmental Protection Agency Chesapeake Bay Program Office, 410 Severn Avenue, Suite 112, Annapolis, Maryland 21403, USA

²Department of Environmental Health and Engineering, Johns Hopkins University, 3400 North Charles Street, Baltimore, Maryland 21218, USA

³Department of Environmental System Sciences, ETH Zurich, Universitatstrasse 16, 8092 Zurich, Switzerland

⁴Swiss Federal Research Institute WSL, Zurcherstrasse 111, 8903 Birmensdorf, Switzerland

⁵Department of Earth and Planetary Science, University of California, Berkeley, Berkeley, California 94720, USA

Correspondence: Qian Zhang (qzhang@chesapeakebay.net)

Received: 30 May 2017 – Discussion started: 21 June 2017

Revised: 21 November 2017 – Accepted: 19 December 2017 – Published: 12 February 2018

Abstract. River water-quality time series often exhibit fractal scaling, which here refers to autocorrelation that decays as a power law over some range of scales. Fractal scaling presents challenges to the identification of deterministic trends because (1) fractal scaling has the potential to lead to false inference about the statistical significance of trends and (2) the abundance of irregularly spaced data in water-quality monitoring networks complicates efforts to quantify fractal scaling. Traditional methods for estimating fractal scaling – in the form of spectral slope (β) or other equivalent scaling parameters (e.g., Hurst exponent) – are generally inapplicable to irregularly sampled data. Here we consider two types of estimation approaches for irregularly sampled data and evaluate their performance using synthetic time series. These time series were generated such that (1) they exhibit a wide range of prescribed fractal scaling behaviors, ranging from white noise ($\beta = 0$) to Brown noise ($\beta = 2$) and (2) their sampling gap intervals mimic the sampling irregularity (as quantified by both the skewness and mean of gap-interval lengths) in real water-quality data. The results suggest that none of the existing methods fully account for the effects of sampling irregularity on β estimation. First, the results illustrate the danger of using interpolation for gap filling when examining autocorrelation, as the interpolation methods consistently underestimate or overestimate β under a wide range of prescribed β values and gap distributions. Second, the widely used Lomb–Scargle spectral method also consistently under-

estimates β . A previously published modified form, using only the lowest 5 % of the frequencies for spectral slope estimation, has very poor precision, although the overall bias is small. Third, a recent wavelet-based method, coupled with an aliasing filter, generally has the smallest bias and root-mean-squared error among all methods for a wide range of prescribed β values and gap distributions. The aliasing method, however, does not itself account for sampling irregularity, and this introduces some bias in the result. Nonetheless, the wavelet method is recommended for estimating β in irregular time series until improved methods are developed. Finally, all methods' performances depend strongly on the sampling irregularity, highlighting that the accuracy and precision of each method are data specific. Accurately quantifying the strength of fractal scaling in irregular water-quality time series remains an unresolved challenge for the hydrologic community and for other disciplines that must grapple with irregular sampling.

1 Introduction

1.1 Autocorrelations in time series

It is well known that time series from natural systems often exhibit autocorrelation; that is, observations at each time step are correlated with observations one or more time steps

in the past. This property is usually characterized by the autocorrelation function (ACF), which is defined as follows for a process X_t at lag k :

$$\gamma(k) = \text{cov}(X_t, X_{t+k}). \quad (1)$$

In practice, autocorrelation has been frequently modeled with classical techniques such as autoregressive (AR) or autoregressive moving-average (ARMA) models (Darken et al., 2002; Yue et al., 2002; Box et al., 2008). These models assume that the underlying process has short-term memory; i.e., the ACF decays exponentially with lag k (Box et al., 2008).

Although the short-term memory assumption holds sometimes, it cannot adequately describe many time series whose ACFs decay as a power law (thus much slower than exponentially) and may not reach zero even for large lags, which implies that the ACF is non-summable. This property is commonly referred to as long-term memory or fractal scaling, as opposed to short-term memory (Beran, 2010).

Fractal scaling has been increasingly recognized in studies of hydrological time series, particularly for the common task of trend identification. Such hydrological series include river flows (Montanari et al., 2000; Khaliq et al., 2008, 2009; Ehsanzadeh and Adamowski, 2010), air and sea temperatures (Fatichi et al., 2009; Lennartz and Bunde, 2009; Franzke, 2012a, b), conservative tracers (Kirchner et al., 2000, 2001; Godsey et al., 2010), and non-conservative chemical constituents (Kirchner and Neal, 2013; Aubert et al., 2014). Because for fractal scaling processes the variance of the sample mean converges to zero much slower than the rate of n^{-1} (n : sample size), the fractal scaling property must be taken into account to avoid false positives (Type I errors) when inferring the statistical significance of trends (Cohn and Lins, 2005; Fatichi et al., 2009; Ehsanzadeh and Adamowski, 2010; Franzke, 2012a). Unfortunately, as stressed by Cohn and Lins (2005), it is “surprising that nearly every assessment of trend significance in geophysical variables published during the past few decades has failed [to do so]”, and a similar tendency is evident in the decade following that statement as well.

1.2 Overview of approaches for quantification of fractal scaling

Several equivalent metrics can be used to quantify fractal scaling. Here we provide a review of the definitions of such processes and several typical modeling approaches, including both time-domain and frequency-domain techniques, with special attention to their reconciliation. For a more comprehensive review, readers are referred to Beran et al. (2013), Boutahar et al. (2007), and Witt and Malamud (2013).

Strictly speaking, X_t is called a stationary long-memory process if the condition

$$\lim_{k \rightarrow \infty} k^\alpha \gamma(k) = C_1 > 0, \quad (2)$$

where C_1 is a constant and is satisfied by some $\alpha \in (0, 1)$ (Boutahar et al., 2007; Beran et al., 2013). Equivalently, X_t is a long-memory process if, in the spectral domain, the condition

$$\lim_{\omega \rightarrow 0} |\omega|^\beta f(\omega) = C_2 > 0 \quad (3)$$

is satisfied by some $\beta \in (0, 1)$, where C_2 is a constant and $f(\omega)$ is the spectral density function of X_t , which is related to ACF as follows (which is also known as the Wiener–Khinchin theorem):

$$f(\omega) = \frac{1}{2\pi} \sum_{k=-\infty}^{\infty} \gamma(k) e^{-ik\omega}, \quad (4)$$

where ω is angular frequency (Boutahar et al., 2007).

One popular model for describing long-memory processes is the so-called fractional autoregressive integrated moving-average model, or ARFIMA (p, q, d), which is an extension of ARMA models and is defined as follows:

$$(1 - B)^d \varphi(B) X_t = \psi(B) \varepsilon_t, \quad (5)$$

where ε_t is a series of independent, identically distributed Gaussian random numbers ($0, \sigma_\varepsilon^2$), B is the backshift operator (i.e., $BX_t = X_{t-1}$), and functions $\varphi(\bullet)$ and $\psi(\bullet)$ are polynomials of order p and q , respectively. The fractional differencing parameter d is related to the parameter α in Eq. (2) as follows:

$$d = \frac{1 - \alpha}{2} \in (-0.5, 0.5) \quad (6)$$

(Beran et al., 2013; Witt and Malamud, 2013).

In addition to a slowly decaying ACF, a long-memory process manifests itself in two other equivalent fashions. One is the so-called Hurst effect, which states that, on a log–log scale, the range of variability of a process changes linearly with the length of the time period under consideration. This power-law slope is often referred to as the Hurst exponent or Hurst coefficient H (Hurst, 1951), which is related to d as follows:

$$H = d + 0.5 \quad (7)$$

(Beran et al., 2013; Witt and Malamud, 2013).

The second equivalent description of long-memory processes, this time from a frequency-domain perspective, is fractal scaling, which describes a power-law decrease in spectral power with increasing frequency, yielding power spectra that are linear on log–log axes (Lomb, 1976; Scargle, 1982; Kirchner, 2005). Mathematically, this inverse proportionality can be expressed as

$$f(\omega) = C_3 |\omega|^{-\beta}, \quad (8)$$

where C_3 is a constant and the scaling exponent β is termed the spectral slope. In particular, for spectral slopes of zero,

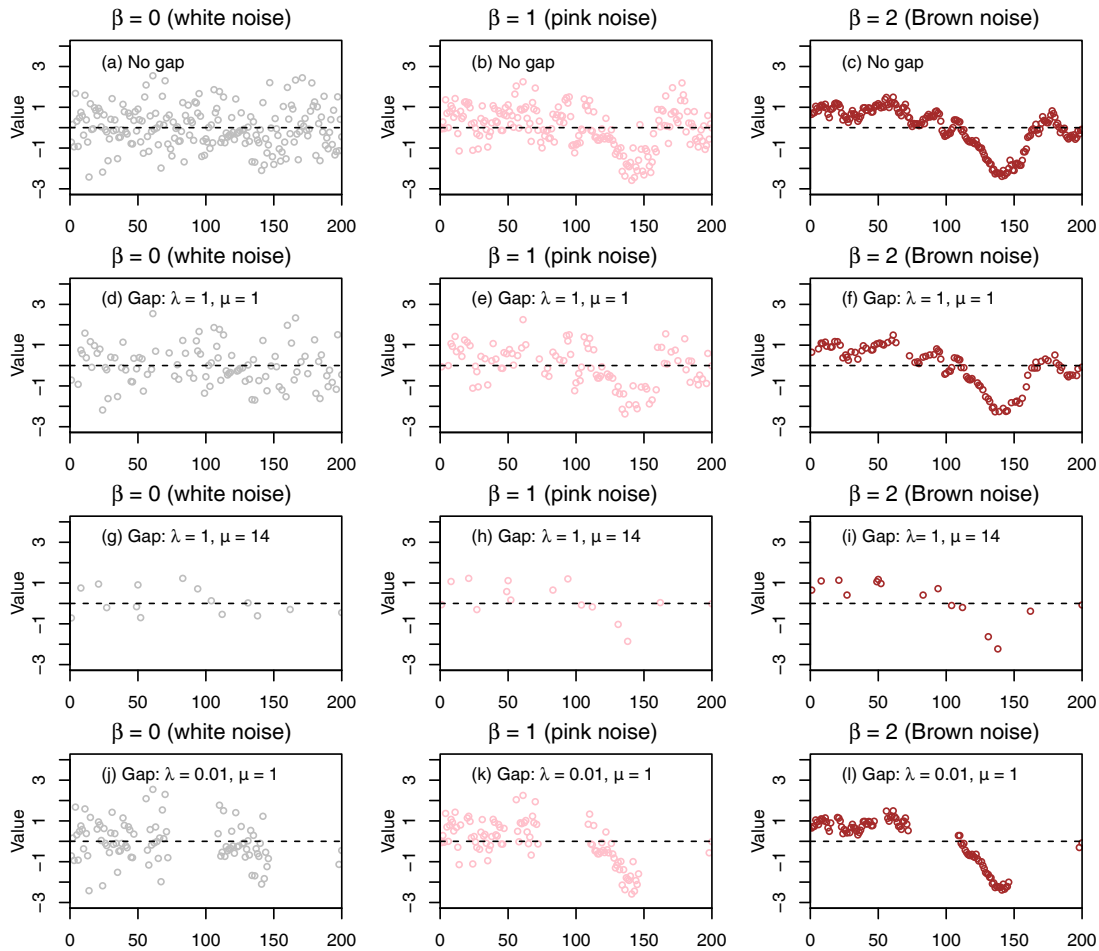


Figure 1. Synthetic time series with 200 time steps for three representative fractal scaling processes that correspond to white noise ($\beta = 0$), pink noise ($\beta = 1$), and Brown noise ($\beta = 2$). (a–c) show the simulated time series without any gap. (d–l) show the same time series as in (a–c) but with data gaps that were simulated using three different negative binomial (NB) distributions – that is, (d–f): $\text{NB}(\lambda = 1, \mu = 1)$; (g–i): $\text{NB}(\lambda = 1, \mu = 14)$; (j–l): $\text{NB}(\lambda = 0.01, \mu = 1)$.

one, and two, the underlying processes are termed as “white”, “pink” (or “flicker”), and “Brown” (or “red”) noises, respectively (Witt and Malamud, 2013). Illustrative examples of these three noises are shown in Fig. 1a–c.

In addition, it can be shown that the spectral density function for ARFIMA (p, d, q) is

$$f(\omega) = \frac{\sigma_\varepsilon^2}{2\pi} \frac{|\psi(e^{-i\omega})|^2}{|\varphi(e^{-i\omega})|^2} |1 - e^{-i\omega}|^{-2d} \quad (9)$$

for $-\pi < \omega < \pi$ (Boutahar et al., 2007; Beran et al., 2013). For $|\omega| \ll 1$, Eq. (9) can be approximated by

$$f(\omega) = C_4 |\omega|^{-2d} \quad (10)$$

with

$$C_4 = \frac{\sigma_\varepsilon^2}{2\pi} \frac{|\psi(1)|^2}{|\varphi(1)|^2}. \quad (11)$$

Equation (10) thus exhibits the asymptotic behavior required for a long-memory process given by Eq. (3). In addition, a comparison of Eqs. (10) and (8) reveals that

$$\beta = 2d. \quad (12)$$

Overall, these derivations indicate that these different types of scaling parameters (i.e., α , d , and H and β) can be used equivalently to describe the strength of fractal scaling. Specifically, their equivalency can be summarized as follows:

$$\beta = 2d = 1 - \alpha = 2H - 1. \quad (13)$$

It should be noted, however, that the parameters d , α , and H are only applicable over a fixed range of fractal scaling, which is equivalent to $(-1, 1)$ in terms of β .

1.3 Motivation and objective of this work

To account for fractal scaling in trend analysis, one must be able to first quantify the strength of fractal scaling for a given

time series. Numerous estimation methods have been developed for this purpose, including the Hurst rescaled range analysis, Higuchi's method, Geweke and Porter-Hudak's method, Whittle's maximum likelihood estimator, detrended fluctuation analysis, and others (Taquu et al., 1995; Montanari et al., 1997, 1999; Rea et al., 2009; Stroe-Kunold et al., 2009). For brevity, these methods are not elaborated here; readers are referred to Beran (2010) and Witt and Malamud (2013) for details. While these estimation methods have been extensively adopted, they are unfortunately only applicable to regular (i.e., evenly spaced) data, e.g., daily streamflow discharge, monthly temperature. In practice, many types of hydrological data, including river water-quality data, are often sampled irregularly or have missing values, and hence their strengths of fractal scaling cannot be readily estimated with the above traditional estimation methods.

Thus, estimation of fractal scaling in irregularly sampled data is an important challenge for hydrologists and practitioners. Many data analysts may be tempted to interpolate the time series to make it regular and hence analyzable (Graham, 2009). Although technically convenient, interpolation can be problematic if it distorts the series' autocorrelation structure (Kirchner and Weil, 1998). In this regard, it is important to evaluate various types of interpolation methods using carefully designed benchmark tests and to identify the scenarios under which the interpolated data can yield reliable (or, alternatively, biased) estimates of spectral slope.

Moreover, quantification of fractal scaling in real-world water-quality data is subject to several common complexities. First, water-quality data are rarely normally distributed; instead, they are typically characterized by log-normal or other skewed distributions (Hirsch et al., 1991; Helsel and Hirsch, 2002), with potential consequences for β estimation. Moreover, water-quality data also tend to exhibit long-term trends, seasonality, and flow dependence (Hirsch et al., 1991; Helsel and Hirsch, 2002), which can also affect the accuracy of β estimates. Thus, it may be more plausible to quantify β in transformed time series after accounting for the seasonal patterns and discharge-driven variations in the original time series, which is the approach taken in this paper. For the trend aspect, however, it remains a puzzle whether the data set should be detrended before conducting β estimation. Such detrending treatment can certainly affect the estimated value of β and hence the validity of (or confidence in) any inference made regarding the statistical significance of temporal trends in the time series. This somewhat circular issue is beyond the scope of our current work – it has been previously discussed in the context of short-term memory (Zetterqvist, 1991; Darken et al., 2002; Yue et al., 2002; Noguchi et al., 2011; Clarke, 2013; Sang et al., 2014), but it is not well understood in the context of fractal scaling (or long-term memory) and hence presents an important area for future research.

In the above context, the main objective of this work was to use Monte Carlo simulation to systematically evaluate and

compare two broad types of approaches for estimating the strength of fractal scaling (i.e., spectral slope β) in irregularly sampled river water-quality time series. Specific aims of this work include the following:

1. to examine the sampling irregularity of typical river water-quality monitoring data and to simulate time series that contain such irregularity, and
2. to evaluate two broad types of approaches for estimating β in simulated irregularly sampled time series.

The first type of approach includes several forms of interpolation techniques for gap filling, thus making the data regular and analyzable by traditional estimation methods. The second type of approach includes the well-known Lomb–Scargle periodogram (Lomb, 1976; Scargle, 1982) and a recently developed wavelet method combined with a spectral aliasing filter (Kirchner and Neal, 2013). The latter two methods can be directly applied to irregularly spaced data; here we aim to compare them with the interpolation techniques. Details of these various approaches are provided in Sect. 3.1.

This work was designed to make several specific contributions. First, it uses benchmark tests to quantify the performance of a wide range of methods for estimating fractal scaling in irregularly sampled water-quality data. Second, it proposes an innovative and general approach for modeling sampling irregularity in water-quality records. Third, while this work was not intended to compare all published estimation methods for fractal scaling, it does provide and demonstrate a generalizable framework for data simulation (with gaps) and β estimation, which can be readily applied toward the evaluation of other methods that are not covered here. Last but not least, while this work was intended to help hydrologists and practitioners understand the performance of various approaches for water-quality time series, the findings and approaches may be broadly applicable to irregularly sampled data in other scientific disciplines.

The rest of the paper is organized as follows. We propose a general approach for modeling sampling irregularity in typical river water-quality data and discuss our approach for simulating irregularly sampled data (Sect. 2). We then introduce various methods for estimating fractal scaling in irregular time series and compare their estimation performance (Sect. 3). We close with a discussion of the results and implications (Sect. 4).

2 Quantification of sampling irregularity in river water-quality data

2.1 Modeling of sampling irregularity

River water-quality data are often sampled irregularly. In some cases, samples are taken more frequently during particular periods of interest, such as high flows or drought pe-

riods; here we will address the implications of the irregularity, but not the (intentional) bias, inherent in such a sampling strategy. In other cases, the sampling is planned with a fixed sampling interval (e.g., 1 day) but samples are missed (or lost, or fail quality-control checks) at some time steps during implementation. In still other cases, the sampling is intrinsically irregular because, for example, one cannot measure the chemistry of rainfall on rainless days or the chemistry of a stream that has dried up. Theoretically, any deviation from fixed-interval sampling can affect the subsequent analysis of the time series.

To quantify sampling irregularity, we propose a simple and general approach that can be applied to any time series of monitoring data. Specifically, for a given time series with N points, the time intervals between adjacent samples are calculated; these intervals themselves make up a time series of $N-1$ points that we call Δt . For this time series, the following parameters are calculated to quantify its sampling irregularity:

- L = the length of the period of record;
- N = the number of samples in the record;
- $\Delta t_{\text{nominal}}$ = the nominal sampling interval under regular sampling (e.g., $\Delta t_{\text{nominal}} = 1$ day for daily samples);
- $\Delta t^* = \Delta t / \Delta t_{\text{nominal}}$, the sample intervals non-dimensionalized by the nominal sampling interval;
- $\Delta t_{\text{average}} = L / (N - 1)$ the average of all the entries in Δt .

The quantification is illustrated with two simple examples. The first example contains data sampled every hour from 01:00 to 11:00 UTC on 1 day. In this case, $L = 10$ h, $N = 11$ samples, $\Delta t = \{1, 1, 1, 1, 1, 1, 1, 1, 1, 1\}$ h, and $\Delta t_{\text{nominal}} = \Delta t_{\text{average}} = 1$ h. The second example contains data sampled at 01:00, 03:00, 04:00, 08:00, and 11:00. In this case, $L = 10$ h, $N = 5$ samples, $\Delta t = \{2, 1, 4, 3\}$ h, $\Delta t_{\text{nominal}} = 1$ h, and $\Delta t_{\text{average}} = 2.5$ h. It is readily evident that the first case corresponds to fixed-interval (regular) sampling that has the property of $\Delta t_{\text{average}} / \Delta t_{\text{nominal}} = 1$ (dimensionless), whereas the second case corresponds to irregular sampling for which $\Delta t_{\text{average}} / \Delta t_{\text{nominal}} > 1$.

The dimensionless set Δt^* contains essential information for determining sampling irregularity. This set is modeled as independent, identically distributed values drawn from a negative binomial (NB) distribution. This distribution has two dimensionless parameters, the shape parameter (λ) and the mean parameter (μ), which collectively represent the irregularity of the samples. The NB distribution is a flexible distribution that provides a discrete analogue of a gamma distribution. The geometric distribution, itself the discrete analogue of the exponential distribution, is a special case of the NB distribution when $\lambda = 1$.

The parameters μ and λ represent different aspects of sampling irregularity, as illustrated by the examples shown in Fig. 2. The mean parameter μ represents the fractional increase in the average interval between samples due to gaps: $\mu = \text{mean}(\Delta t^*) - 1 = (\Delta t_{\text{average}} - \Delta t_{\text{nominal}}) / \Delta t_{\text{nominal}}$. Thus, the special case of $\mu = 0$ corresponds to regular sampling (i.e., $\Delta t_{\text{average}} = \Delta t_{\text{nominal}}$), whereas any larger value of μ corresponds to irregular sampling (i.e., $\Delta t_{\text{average}} > \Delta t_{\text{nominal}}$) (Fig. 2c). The shape parameter λ characterizes the similarity of gaps to each other; that is, a small λ indicates that the samples contain gaps of widely varying lengths, whereas a large λ indicates that the samples contain many gaps of similar lengths (Fig. 2a, b).

To visually illustrate these gap distributions, representative samples of irregular time series are presented in Fig. 1 for the three special processes described above (Sect. 1.2), i.e., white noise, pink noise, and Brown noise. Specifically, three different gap distributions, namely, NB($\lambda = 1, \mu = 1$), NB($\lambda = 1, \mu = 14$), and NB($\lambda = 0.01, \mu = 1$), were simulated and each was applied to convert the three original (regular) time series (Fig. 1a–c) to irregular time series (Fig. 1d–f). These simulations clearly illustrate the effects of the two parameters λ and μ . In particular, compared with NB($\lambda = 1, \mu = 1$), NB($\lambda = 1, \mu = 14$) shows a similar level of sampling irregularity (same λ) but a much longer average gap interval (larger μ). Again compared with NB($\lambda = 1, \mu = 1$), NB($\lambda = 0.01, \mu = 1$) shows the same average interval (same μ) but a much more irregular (skewed) gap distribution that contains a few very large gaps (smaller λ).

2.2 Examination of sampling irregularity in real river water-quality data

The above modeling approach was applied to real water-quality data from two large river monitoring networks in the United States to examine sampling irregularity. One such network is the Chesapeake Bay River Input Monitoring Program, which typically samples streams roughly once or twice monthly, accompanied with additional sampling during storm flows (Langland et al., 2012; Zhang et al., 2015). These data were obtained from the US Geological Survey National Water Information System (<http://doi.org/10.5066/F7P55KJN>). The other network is the Lake Erie and Ohio Tributary Monitoring Program, which typically samples streams at a daily resolution (National Center for Water Quality Research, 2015). For each site, we determined the NB parameters to quantify sampling irregularity. The mean parameter μ can be estimated as described above, and the shape parameter λ can be calculated directly from the mean and variance of Δt^* as follows: $\lambda = \mu^2 / [\text{var}(\Delta t^*) - \mu] = (\text{mean}(\Delta t^*) - 1)^2 / [\text{var}(\Delta t^*) - \text{mean}(\Delta t^*) + 1]$. Alternatively, a maximum likelihood approach can be used, which employs the *fitdist* function in the *fitdistrplus* R package (Delignette-Muller and Dutang, 2015). In general, the two approaches produce similar results, which are summarized in

Table 1. Quantification of sampling irregularity for selected water-quality constituents at nine sites of the Chesapeake Bay River Input Monitoring Program and six sites of the Lake Erie and Ohio Tributary Monitoring Program. (λ : shape parameter estimated using maximum likelihood; λ' : shape parameter estimated using the direct approach (see Sect. 2.2); μ : mean parameter; Δ average: average gap interval; N : total number of samples.)

I. Chesapeake Bay River Input Monitoring program												
Site ID	River and station name	Drainage area (km ²)	Total nitrogen (TN)					Total phosphorus (TP)				
			λ	λ'	μ	Δ average (days)	N	λ	λ'	μ	Δ average (days)	N
01578310	Susquehanna River at Conowingo, MD	70 189	0.8	1.1	13.5	14.5	876	0.8	1.0	13.4	14.4	881
01646580	Potomac River at Chain Bridge, Washington D.C.	30 044	0.9	0.6	9.5	10.5	1385	1.1	1.0	24.4	25.4	579
02035000	James River at Cartersville, VA	16 213	0.8	1.0	13.9	14.9	960	0.8	1.1	13.7	14.7	974
01668000	Rappahannock River near Fredericksburg, VA	4144	0.8	0.6	15.6	16.6	776	0.8	0.6	15.2	16.2	796
02041650	Appomattox River at Matoaca, VA	3471	0.8	0.8	15.1	16.1	798	0.8	0.8	14.9	15.9	810
01673000	Pamunkey River near Hanover, VA	2774	0.8	0.9	15.1	16.1	873	0.8	1.0	14.7	15.7	894
01674500	Mattaponi River near Beulahville, VA	1557	0.7	0.9	14.3	15.3	810	0.8	0.9	14.2	15.2	820
01594440	Patuxent River at Bowie, MD	901	0.9	1.1	15.3	16.3	787	0.8	0.8	14.0	15.0	861
01491000	Choptank River near Greensboro, MD	293	1.2	1.5	19.6	20.6	680	1.1	1.0	20.5	21.5	690
II. Lake Erie and Ohio tributary monitoring program												
Site ID	River and station name	Drainage area (km ²)	Nitrate-plus-nitrite (NO _x)					Total phosphorus (TP)				
			λ	λ'	μ	Δ average (days)	N	λ	λ'	μ	Δ average (days)	N
04193500	Maumee River at Waterville, OH	16 395	0.005	0.0003	0.19	1.19	9101	0.005	0.0003	0.19	1.19	9101
04198000	Sandusky River near Fremont, OH	3245	0.01	0.003	0.22	1.22	9641	0.01	0.003	0.22	1.22	9655
04208000	Cuyahoga River at Independence, OH	1834	0.007	0.006	0.13	1.13	7421	0.007	0.006	0.13	1.13	7426
04212100	Grand River near Painesville, OH	1777	0.01	0.005	0.21	1.21	5023	0.01	0.005	0.22	1.22	4994
04197100	Honey Creek at Melmore, OH	386	0.007	0.005	0.06	1.06	9914	0.007	0.005	0.06	1.06	9914
04197170	Rock Creek at Tiffin, OH	90	0.007	0.008	0.06	1.06	8422	0.007	0.008	0.06	1.06	8440

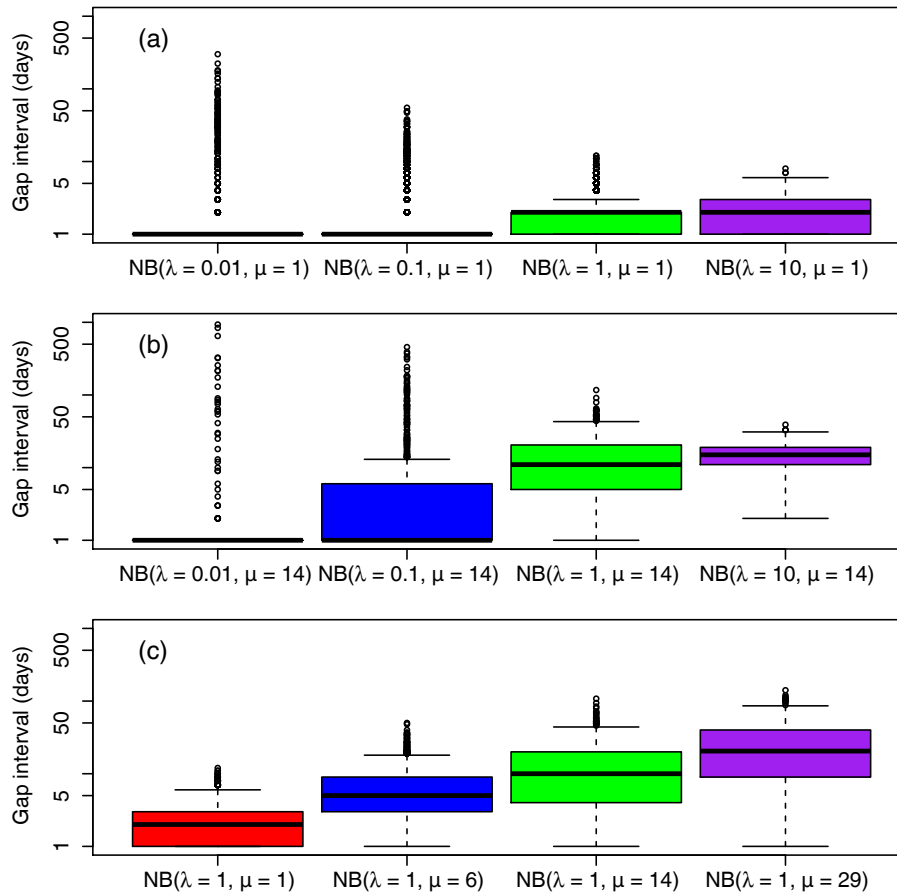


Figure 2. Examples of gap-interval simulation using negative binomial distributions, NB (shape λ , mean μ). Simulation parameters: $L=9125$ days, $\Delta t_{\text{nominal}}=1$ day. The three panels show simulation with fixed (a) $\mu=1$, (b) $\mu=14$, and (c) $\lambda=1$. Note that $\Delta t_{\text{average}}/\Delta t_{\text{nominal}}=\mu+1$.

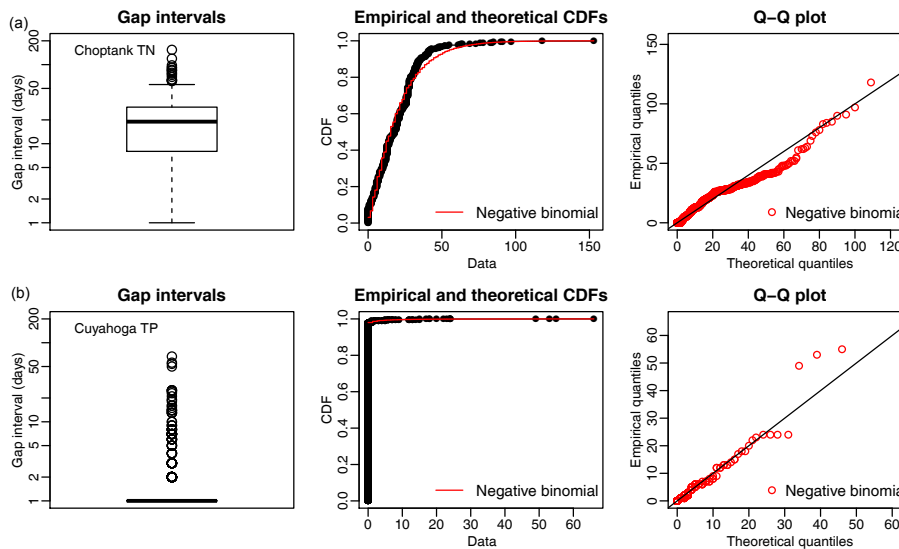


Figure 3. Examples of quantified sampling irregularity with negative binomial (NB) distributions: total nitrogen in Choptank River (a) and total phosphorus in Cuyahoga River (b). Theoretical CDF (cumulative distribution function) and quantiles are based on the fitted NB distributions. See Table 1 for estimated mean and shape parameters.

Table 1, with two examples of fitted NB distributions shown in Fig. 3.

For the Chesapeake Bay River Input Monitoring Program (nine sites), total nitrogen (TN) and total phosphorus (TP) are taken as representatives of water-quality constituents. According to the maximum likelihood approach, the shape parameter λ varies between 0.7 and 1.2 for TN and between 0.8 and 1.1 for TP (Table 1). These λ values are around 1.0, reflecting the fact that these sites have relatively even gap distributions (i.e., relatively balanced counts of large and small gaps). The mean parameter μ varies between 9.5 and 19.6 for TN and between 13.4 and 24.4 for TP in the Chesapeake monitoring network, corresponding to $\Delta t_{\text{average}}$ of 10.5–20.6 days for TN and 14.4–25.4 days for TP, respectively. This is consistent with the fact that these sites have typically been sampled roughly once or twice monthly, along with additional sampling during storm flows (Langland et al., 2012; Zhang et al., 2015).

For the Lake Erie and Ohio Tributary Monitoring Program (six sites), records of nitrate plus nitrite (NO_x) and TP were examined. According to the maximum likelihood approach, the shape parameter λ is approximately 0.01 for both constituents (Table 1). These very low λ values occur because these time series contain a few very large gaps, ranging from 35 days to 1109 days (~ 3 years). The mean parameter μ varies between 0.06 and 0.22, corresponding to $\Delta t_{\text{average}}$ of 1.06 and 1.22 days, respectively. This is consistent with the fact that these sites have been sampled at a daily resolution with occasional missing values on some days (Zhang and Ball, 2017).

2.3 Simulation of time series with irregular sampling

To evaluate the various β estimation methods, our first step was to use the Monte Carlo simulation to produce time series that mimic the sampling irregularity observed in real water-quality monitoring data. We began by simulating regular (gap-free) time series using the fractional noise simulation method of Witt and Malamud (2013), which is based on inverse Fourier filtering of white noises. Our analysis showed this method performed reasonably well compared to other simulation methods for β values between 0 and 1 (see the Supplement). In addition, this method can also simulate β values beyond this range. The noises simulated by the Witt and Malamud method, however, are band limited to the Nyquist frequency (half of the sampling frequency) of the underlying white noise time series, whereas true fractional noises would contain spectral power at all frequencies, extending well above the Nyquist frequency for any sampling. Thus, these band-limited noises will be less susceptible to spectral aliasing than true fractional noises would be (see Kirchner, 2005, for detailed discussions of the aliasing issue).

A total of 100 replicates of regular (gap-free) time series were produced for nine prescribed spectral slopes, which

vary from $\beta = 0$ (white noise) to $\beta = 2$ (Brownian motion or “random walk”) with an increment of 0.25 (i.e., 0, 0.25, 0.5, 0.75, 1.0, 1.25, 1.5, 1.75, and 2). These regular time series each have a length (N) of 9125, which can be interpreted as 25 years of regular daily samples (that is, $\Delta t_{\text{nominal}} = 1$ day).

The simulated regular time series were converted to irregular time series using gap intervals that were simulated with NB distributions. To make these gap intervals mimic those in typical river water-quality time series, representative NB parameters were chosen based on results from Sect. 2.2. Specifically, μ was set at 1 and 14, corresponding to $\Delta t_{\text{average}}$ of 2 and 15 days, respectively. For λ , we chose four values that span 3 orders of magnitude, i.e., 0.01, 0.1, 1, and 10. Note that when $\lambda = 1$ the generated time series corresponds to a Bernoulli process. With the chosen values of μ and λ , a total of eight scenarios were generated, which were implemented using the *rbinom* function in the *stats* R package (R Development Core Team, 2014):

1. $\mu = 1$ (i.e., $\Delta t_{\text{average}} / \Delta t_{\text{nominal}} = 2$), $\lambda = 0.01$,
2. $\mu = 1$, $\lambda = 0.1$,
3. $\mu = 1$, $\lambda = 1$,
4. $\mu = 1$, $\lambda = 10$,
5. $\mu = 14$ (i.e., $\Delta t_{\text{average}} / \Delta t_{\text{nominal}} = 15$), $\lambda = 0.01$,
6. $\mu = 14$, $\lambda = 0.1$,
7. $\mu = 14$, $\lambda = 1$,
8. $\mu = 14$, $\lambda = 10$.

Examples of these simulations are shown with box plots in Fig. 2.

3 Evaluation of proposed estimation methods for irregular time series

3.1 Summary of estimation methods

For the simulated irregular time series, β was estimated using the aforementioned two types of approaches. The first type includes 11 different interpolation methods (designated as B1–B11 below) to fill the data gaps, thus making the data regular and analyzable by traditional methods.

- B1 Global mean: all missing values replaced with the mean of all observations.
- B2 Global median: all missing values replaced with the median of all observations.
- B3 Random replacement: all missing values replaced with observations randomly drawn (with replacement) from the time series.

- B4 Next observation carried backward (NOCB): each missing value replaced with the next available observation.
- B5 Last observation carried forward (LOCF): each missing value replaced with the preceding available observation.
- B6 Average of the two nearest samples: each missing value replaced with the mean of its next and preceding available observations.
- B7 LOWESS (locally weighted scatterplot smoothing) with a smoothing span of 1: missing values replaced using fitted values from a LOWESS model determined using all available observations (Cleveland, 1981).
- B8 LOWESS with a smoothing span of 0.75: same as B7 except that the smoothing span is 75 % of the available data (similar distinction follows for B9–B11).
- B9 LOWESS with a smoothing span of 50 %.
- B10 LOWESS with a smoothing span of 30 %.
- B11 LOWESS with a smoothing span of 10 %.

B4 and B5 were implemented using the *na.locf* function in the *zoo* R package (Zeileis and Grothendieck, 2005). B7–B11 were implemented using the *loess* function in the *stats* R package (R Development Core Team, 2014). An illustration of these interpolation methods is provided in Fig. 4. The interpolated data, along with the original regular data (designated as A1) were analyzed using Whittle’s maximum likelihood method for β estimation, which was implemented using the *FDWhittle* function in the *fractal* R package (Constantine and Percival, 2014).

The second type of approaches estimates β directly from the irregularly sampled data, using several variants of the Lomb–Scargle periodogram (designated as C1a–C1c below), and a recently developed wavelet-based method (designated as C2 below). Specifically, these approaches are as follows.

- C1a Lomb–Scargle periodogram: the spectral density of the time series (with gaps) is estimated and the spectral slope is fit using all frequencies (Lomb, 1976; Scargle, 1982). This is a classic method for examining periodicity in irregularly sampled data, which is analogous to the more familiar fast Fourier transform method often used for regularly sampled data.
- C1b Lomb–Scargle periodogram with 5 % data: same as C1a except that the fitting of the spectral slope considers only the lowest 5 % of the frequencies (Montanari et al., 1999).
- C1c Lomb–Scargle periodogram with “binned” data: same as C1a except that the fitting of the spectral slope is performed on binned data in three steps as follows.
- The entire range of frequency is divided into 100 equal-interval bins on logarithmic scale.
 - The respective medians of frequency and power spectral density are calculated for each of the 100 bins.
 - The 100 pairs of median frequency and median spectral density are used to estimate the spectral slope on a log–log scale.
- C2 Kirchner and Neal (2013)’s wavelet method: uses a modified version of Foster’s weighted wavelet spectrum (Foster, 1996) to suppress spectral leakage from low frequencies and applies an aliasing filter (Kirchner, 2005) to remove spectral aliasing artifacts at high frequencies.

C1a was implemented using the *spec.ls* function in the *cts* R package (Wang, 2013). C2 was run in *C*, using codes modified from those in Kirchner and Neal (2013).

3.2 Evaluation of methods’ performance

Each estimation method listed above was applied to the simulated data (Sect. 2.3) to estimate β , which were then compared with the prescribed (“true”) β to quantify the performance of each method. Plots of method evaluation for all simulations are provided as Figs. S3–S12 (Supplement S2). Close inspections of these plots reveal some general patterns of the methods’ performance. For brevity, these patterns are presented with a subset of the plots, which correspond to the cases where true $\beta = 1$ and shape parameter $\lambda = 0.01, 0.1, 1, \text{ and } 10$ (Fig. 5). In general, β values estimated using the regular data (A1) are very close to 1.0, which indicates that the adopted fractional noise generation method and Whittle’s maximum likelihood estimator have small combined simulation and estimation bias. This is perhaps unsurprising, since the estimator is based on the Fourier transform and the noise generator is based on an inverse Fourier transform; thus, one method is essentially just the inverse of the other. One should also note that when fractional noises are not arbitrarily band limited at the Nyquist frequency (as they inherently are with the noise generator that is used here), spectral aliasing should lead to spectral slopes that are flatter than expected (Kirchner, 2005) and thus to underestimates of β .

For the simulated irregular data, the estimation methods differ widely in their performance. Specifically, three interpolation methods (i.e., B4–B6) consistently overestimate β , indicating that they introduce additional correlations into the time series, reducing its short-timescale variability. In contrast, the other eight interpolation methods (i.e., B1–B3 and B7–B11) generally underestimate β , indicating that the interpolated points are less correlated than the original time series, thus introducing additional variability on short timescales. As expected, results from the LOWESS methods (B7–B11) depend strongly on the size of the smoothing window; that is, β is more severely underestimated as the smoothing window becomes wider. In fact, when the smoothing window is 1.0 (i.e., method B7), LOWESS performs the interpolation using all data available and thus behaves similarly to interpolations

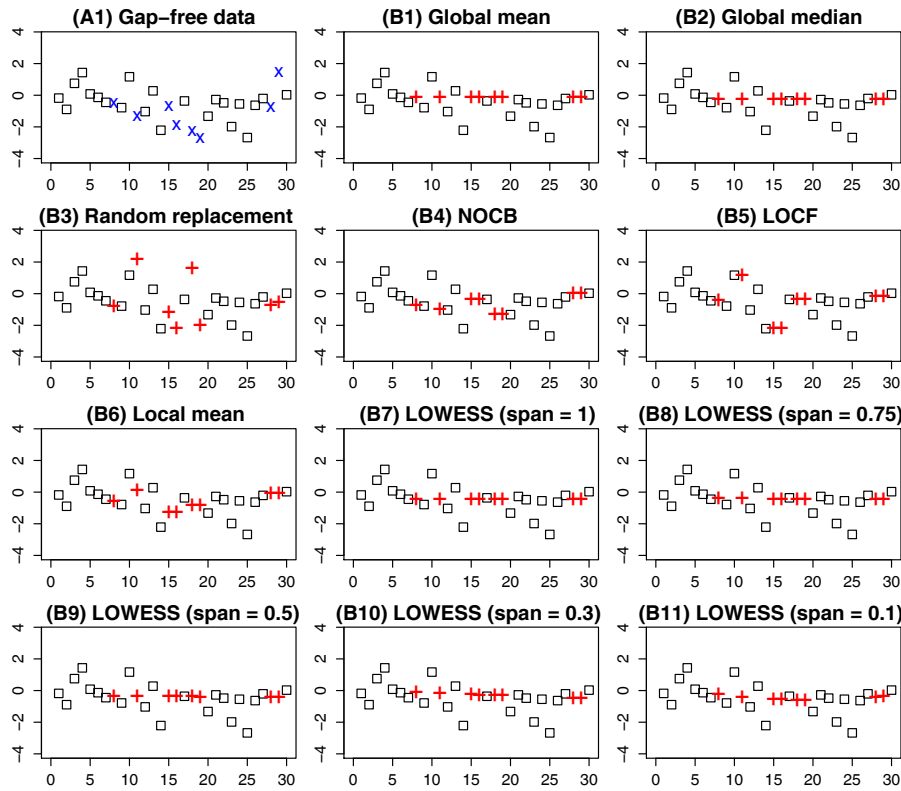


Figure 4. Illustration of the interpolation methods for gap filling. The gap-free data (A1) was simulated with a series length of 500, with the first 30 data shown. (x: omitted data for gap filling; +: interpolated data; NOCB: next observation carried backward; LOCF: last observation carried forward; LOWESS: locally weighted scatterplot smoothing.)

based on global means (B1) or global medians (B2), except that LOWESS fits a polynomial curve instead of constant values. However, whenever a sampling gap is much shorter than the smoothing window, the infilled LOWESS value will be close to the local mean or median, and the abrupt jumps produced by these infilled values will artificially increase the variance in the time series at high frequencies, leading to an artificially reduced spectral slope β and, correspondingly, an underestimate of β . This mechanism explains why LOWESS interpolation distorts β more when there are many small gaps (large λ) and therefore more jumps to, and away from, the infilled values than when there are only a few large gaps (small λ).

Among the direct methods (i.e., C1a, C1b, C1c, and C2), the Lomb–Scargle method, with original data (C1a) or binned data (C1c) tends to underestimate β , though the underestimation by C1c is generally less severe. The modified Lomb–Scargle method (C1b), using only the lowest 5% of frequencies, yields estimates that are centered around 1.0 for large λ . However, C1b has the highest variability (i.e., least precision) in β estimates among all methods. Compared with all the above methods, the wavelet method (C2) has much better performance in terms of both accuracy and precision

when λ is 1 or 10, a slightly better or similar performance when λ is 0.1, but worse performance when λ is 0.01.

The shape parameter λ greatly affects the performance of the estimation methods. All the interpolation methods that underestimate β (i.e., B1–B3 and B7–B11) perform worse as λ increases from 0.01 to 10. This effect can be interpreted as follows: when the time series contains a large number of relatively small gaps (e.g., $\lambda = 1$ or 10), there are many jumps (which, as noted above, contain mostly high-frequency variance) between the original data and the infilled values, resulting in more severe underestimation. In contrast, when the data contain only a small number of very large gaps (e.g., $\lambda = 0.01$ or 0.1), there are fewer of these jumps, resulting in minimal underestimation. Similar effects of λ are also observed with the interpolation methods that show overestimation (i.e., B4–B6) – that is, overestimation is more severe when λ is larger. Similarly, the Lomb–Scargle method (C1a and C1c) performs worse (more serious underestimation) as λ increases. Finally, method C2 seems to perform the best when λ is large (1 or 10), but not well when λ is very small (0.01), as noted above. This result highlights the sensitivity of the wavelet method to the presence of a few large gaps in the time series. For such cases, a potentially more feasible approach is to break the whole time series into several seg-

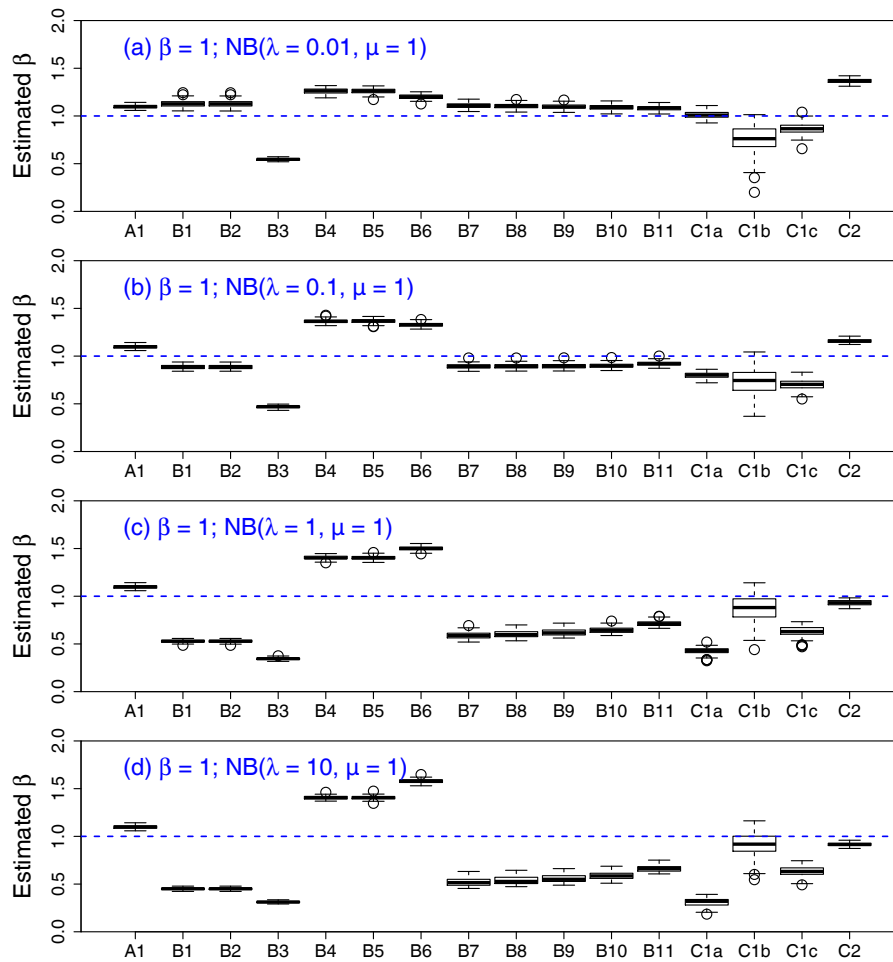


Figure 5. Comparison of bias in estimated spectral slope in irregular data that are simulated with prescribed $\beta = 1$ (100 replicates), a series length of 9125, and gap intervals simulated with (a) NB ($\lambda = 0.01$, $\mu = 1$), (b) NB ($\lambda = 0.1$, $\mu = 1$), (c) NB ($\lambda = 1$, $\mu = 1$), and (d) NB ($\lambda = 10$, $\mu = 1$). The blue dashed lines indicate the true β value.

ments (each without long gaps) and then apply the wavelet method (C2) to analyze each segment separately. If this can yield more accurate estimates, then further simulation experiments should be designed to systematically determine how long the gap needs to be to invoke such an approach.

Next, the method evaluation is extended to all the simulated spectral slopes, that is, $\beta = 0, 0.25, 0.5, 0.75, 1.0, 1.25, 1.5, 1.75$, and 2 . For ease of discussion, three quantitative criteria were proposed for evaluating performance, namely, bias (B), standard deviation (SD), and root-mean-squared error (RMSE), as defined below:

$$B_i = \bar{\beta}_i - \beta_{\text{true}}, \tag{14}$$

$$SD_i = \sqrt{\frac{1}{99} \sum_{j=1}^{100} (\beta_{i,j} - \bar{\beta}_i)^2}, \tag{15}$$

$$RMSE_i = \sqrt{B_i^2 + SD_i^2}, \tag{16}$$

where $\bar{\beta}_i$ is the mean of 100 β values estimated by method i , and β_{true} is the prescribed β value for simulation of the initial regular time series. In general, B and SD can be considered as the models' systematic error and random error, respectively, and RMSE serves as an integrated measure of both errors. For all evaluations, plots of bias and RMSE are provided in the main text. (Plots of SD are provided as Figs. S7 and S12 in the Supplement for simulations with $\mu = 1$ and $\mu = 14$, respectively.)

For simulations with $\mu = 1$, results of estimation bias and RMSE are summarized in Figs. 6 and 7, respectively. (More details are provided in Figs. S3–S6 in the Supplement.) For brevity, we focus on three direct methods (C1a, C1b, and C2) and three representative interpolation methods. (Specifically, B1 represents B1–B3 and B7, B6 represents B4–B6, and B8 represents B8–B11.) Overall, these six methods show mixed performances. In terms of bias (Fig. 6), B1 (global mean) and B8 (LOWESS with a smoothing span of 0.75) tend to have negative bias, particularly for time series with

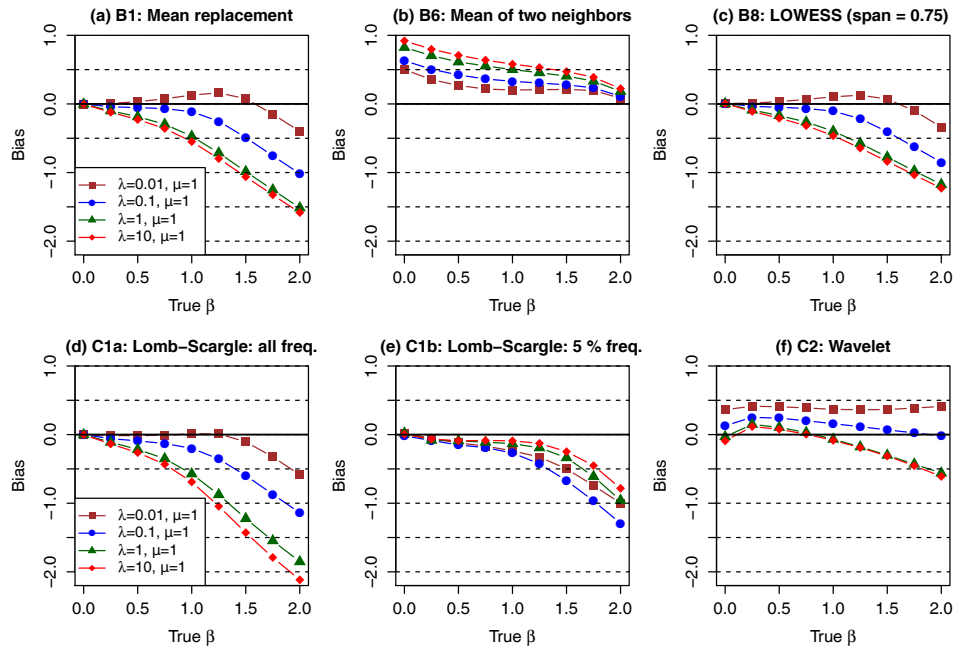


Figure 6. Comparison of bias in estimated spectral slope in irregular data that are simulated with varying prescribed β values (100 replicates), a series length of 9125, and a mean gap interval of 2 (i.e., $\mu = 1$).

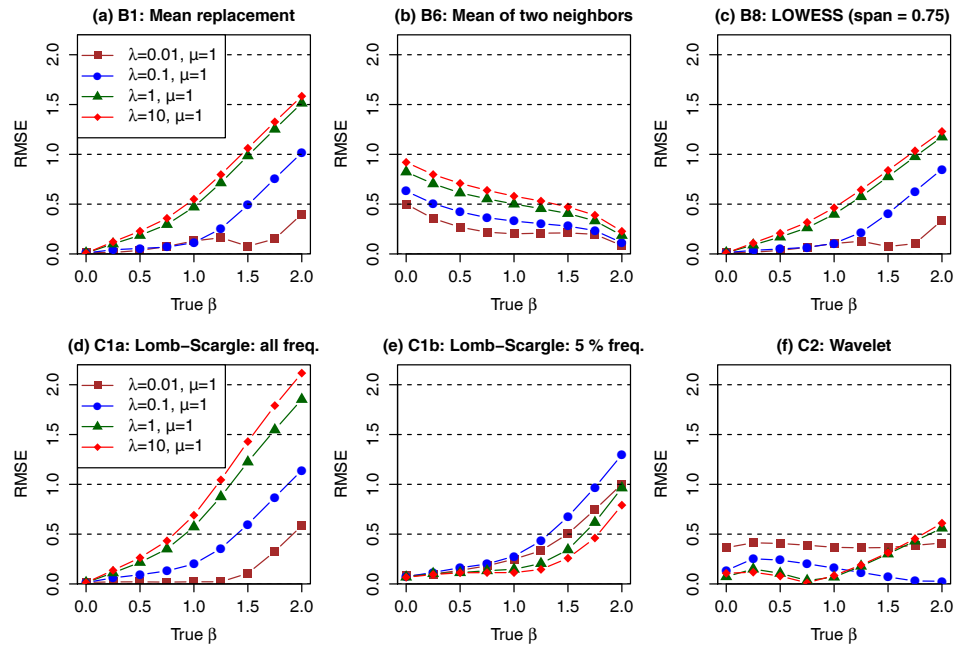


Figure 7. Comparison of root-mean-squared error (RMSE) in estimated spectral slope in irregular data that are simulated with varying prescribed β values (100 replicates), a series length of 9125, and a mean gap interval of 2 (i.e., $\mu = 1$).

(1) moderate-to-large β_{true} values and (2) large λ values (i.e., less skewed gap intervals). By contrast, B1 and B8 generally have minimal bias when (1) β_{true} is close to zero (i.e., when the simulated time series is close to white noise) and (2) λ is small (e.g., 0.01), since interpolating a few large gaps cannot

significantly affect the overall correlation structure. In addition, LOWESS interpolation with a larger smoothing window tends to yield more negatively biased estimates (data not shown). The other interpolation method, B6 (mean of the two nearest neighbors) tends to overestimate β , particularly for

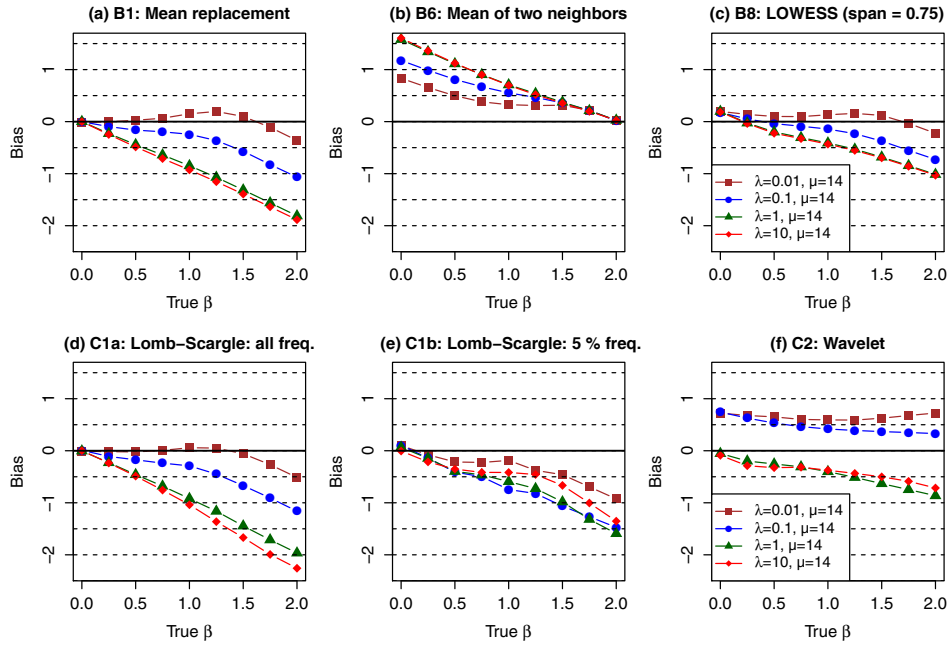


Figure 8. Comparison of bias in estimated spectral slope in irregular data that are simulated with varying prescribed β values (100 replicates), a series length of 9125, and a mean gap interval of 15 (i.e., $\mu = 14$).

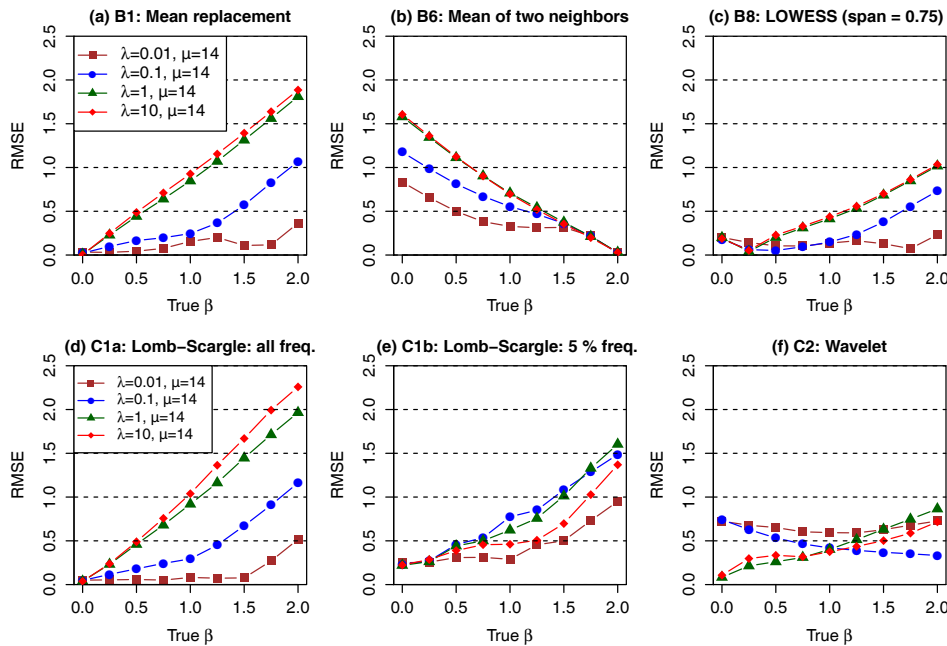


Figure 9. Comparison of root-mean-squared error (RMSE) in estimated spectral slope in irregular data that are simulated with varying prescribed β values (100 replicates), a series length of 9125, and a mean gap interval of 15 (i.e., $\mu = 14$).

time series with (1) small β_{true} values and (2) large λ values. At large β_{true} values (e.g., 2.0), the autocorrelation is already very strong such that taking the mean of two neighbors for gap filling does not introduce much additional correlation, as opposed to the case of small β_{true} values. The Lomb–Scargle methods (C1a and C1b) generally have negative bias, partic-

ularly for time series with (1) moderate-to-large β_{true} values (for both methods) and (2) large λ values (for C1a), which is similar to B1 and B8. However, C1b overall shows less severe bias than C1a. Finally, the wavelet method (C2) shows generally the smallest bias among all methods. However, its performance advantage is not as great when the time series

has small λ values (i.e., very skewed gap intervals), as noted above, which may be due to the fact that the aliasing filter was designed for regular time series. In terms of SD (Fig. S7 in the Supplement), method C1b performs the worst among all methods (as noted above), method B6 and B8 perform poorly for large β_{true} values, and method C2 performs poorly for $\beta_{\text{true}} = 0$. In terms of RMSE (Fig. 7), methods B1, B8, C1a, and C1b perform well for small β_{true} values and small λ values, whereas method B6 performs well for large β_{true} values and small λ values. In comparison, method C2 generally has the smallest RMSEs among all methods, and its RMSEs are similarly small for the wide range of β_{true} and λ values. In general, the wavelet method can be considered the best among all the tested methods.

For simulations with $\mu = 14$, results of estimation bias and RMSE are summarized in Figs. 8 and 9, respectively. (More details are provided in Figs. S8–S11 in the Supplement.) Overall, these methods show mixed performances that are generally similar to the cases when $\mu = 1$, as discussed above. These results highlight the generality of these methods' performances, which applies at least to the range of $\mu = [1, 14]$. In addition, all methods show generally larger RMSE for $\mu = 14$ than $\mu = 1$, indicating their dependence on the mean gap interval (Fig. 9). Perhaps the most notable difference is observed with method C2, which in this case shows positive bias for small λ values (0.01 and 0.1) and negative bias for large λ values (1 and 10) (Fig. 8f). It nonetheless generally shows the smallest RMSEs among all the tested methods as in the cases of $\mu = 1$ above.

3.3 Quantification of spectral slopes in real water-quality data

In this section, the proposed estimation approaches were applied to quantify β in real water-quality data from the two monitoring programs presented in Sect. 2.2 (Table 1). As noted in Sect. 1.3, such real data are typically much more complex than our simulated time series, because of (1) strong deviations from normal distributions and (2) effects of flow dependence, seasonality, and temporal trends (Hirsch et al., 1991; Helsel and Hirsch, 2002). In this regard, future research may simulate time series with these important characteristics and evaluate the performance of various estimation approaches, perhaps following the modeling framework described here. Alternatively, one may quantify β in transformed time series after accounting for the above aspects. In this work, we have taken the latter approach for a preliminary investigation. Specifically, we have used the published weighted regressions on time, discharge, and season (WRTDS) method (Hirsch et al., 2010) to transform the original time series. This widely accepted method estimates daily concentrations based on discretely collected concentration samples using time, season, and discharge as explanatory

variables, i.e.,

$$\ln(C) = \beta_0 + \beta_1 t + \beta_2 \ln(Q) + \beta_3 \sin(2\pi t) + \beta_4 \cos(2\pi t) + \varepsilon, \quad (17)$$

where C is concentration, Q is daily discharge, t is time in decimal years, β_i are fitted coefficients, and ε is the error term. The second and third terms on the right represent time and discharge effects, respectively, whereas the fourth and fifth terms collectively represent cyclical seasonal effects. For a full description of this method, see Hirsch et al. (2010). In this work, WRTDS was applied to obtain time series of estimated daily concentrations for each constituent at each site. The difference between observed concentration (C_{obs}) and estimated concentration (C_{est}) was calculated in logarithmic space to obtain the concentration residuals,

$$\text{residuals} = \ln(C_{\text{obs}}) - \ln(C_{\text{est}}). \quad (18)$$

For our data sets, histograms of concentration residuals (expressed in natural log concentration units) are shown in Figs. S13–S16 in the Supplement. Compared with the original concentration data, these model residuals are much more nearly normal and homoscedastic. Moreover, the model residuals are less susceptible to the issues of temporal, seasonal, and discharge-driven variations than the original concentrations. Therefore, the model residuals are more appropriate than the original concentrations for β estimation using the simulation framework adopted in this work.

The estimated β values for the concentration residuals are summarized in Fig. 10. Clearly, the estimated β varies considerably with the estimation method. In addition, the estimated β varies with site and constituent (i.e., TP, TN, or NO_x). Our discussion below focuses on the wavelet method (C2), because it is established above that this method performs better than the other estimation methods under a wide range of gap conditions. We emphasize that it is beyond our current scope to precisely quantify β in these water-quality data sets, but our simulation results presented above (Sect. 3.2) can be used as references to qualitatively evaluate the reliability of C2 and/or other methods for these data sets.

For TN and TP concentration data at the Chesapeake River input monitoring sites (Table 1), μ varies between 9.5 and 24.4, whereas λ is ~ 1.0 . Thus, the simulated gap scenario of NB($\mu = 14$, $\lambda = 1$) can be used as a reasonable reference to assess methods' reliability (Fig. 8). Based on method C2, the estimated β ranges between $\beta = 0.36$ and $\beta = 0.61$ for TN and between $\beta = 0.30$ and $\beta = 0.58$ for TP at these sites (Fig. 10). For such ranges, the simulation results indicate that method C2 tends to moderately underestimate β under this gap scenario (Fig. 8), and hence spectral slopes for TN and TP at these Chesapeake sites are probably slightly higher than those presented here (Fig. 10).

For NO_x and TP concentration data at the Lake Erie and Ohio sites (Table 1), μ varies between 0.06 and 0.22, whereas λ is ~ 0.01 . Thus, the simulated gap scenario of NB($\mu = 1$,

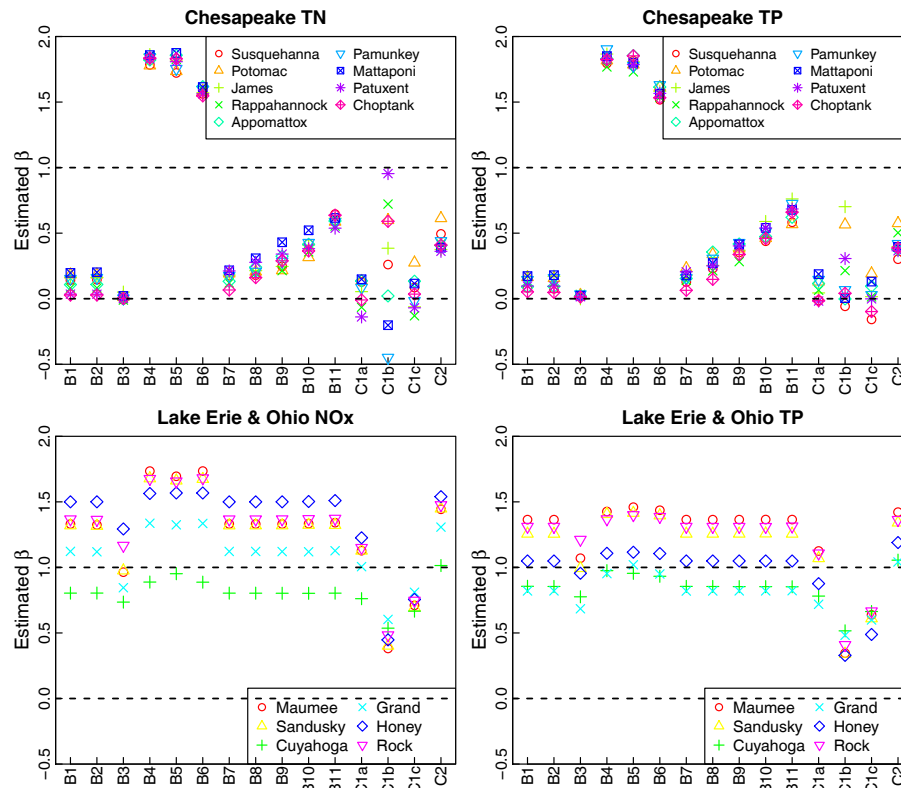


Figure 10. Quantification of spectral slope in real water-quality data from the two regional monitoring networks, as estimated using the set of examined methods. All estimations were performed on concentration residuals (in natural log concentration units) after accounting for effects of time, discharge, and season. The two dashed lines in each panel indicate white noise ($\beta = 0$) and pink (flicker) noise ($\beta = 1$), respectively. See Table 1 for site and data details.

$\lambda = 0.01$) can be used as a reasonable reference to assess the methods' reliability (Fig. 6). For such small λ (i.e., a few gaps that are very dissimilar from others), C2 is not reliable for β estimation, as reflected by the generally positive bias in the simulation results. By contrast, methods B1 (interpolation with global mean) and B8 (LOWESS with span 0.75) both perform quite well under this gap scenario (Fig. 6). These two methods provide almost identical β estimates for each site–constituent combination, ranging from $\beta = 0.8$ to $\beta = 1.5$ for NO_x and TP (Fig. 10).

Overall, the above analysis of real water-quality data has illustrated the wide variability in β estimates, with different choices of estimation methods yielding very different results. To our knowledge, these water-quality data have not previously been analyzed in this context. As illustrated above, our simulation experiments (Sect. 3.2) can be used as references to coarsely evaluate the reliability of each method under specific gap scenarios, thereby considerably narrowing the likely range of the estimated spectral slopes. Nonetheless, our results demonstrate that the analyzed water-quality time series can exhibit strong fractal scaling, particularly at the Lake Erie and Ohio tributary sites. Thus, an important implication is that researchers and analysts should be cautious when

applying standard statistical methods to identify temporal trends in such water-quality data sets (Kirchner and Neal, 2013). In future work, one may consider applying Bayesian statistical analysis or other approaches to more accurately quantify the spectral slope and associated uncertainty for real water-quality data analysis. In addition, the modeling framework presented here (including both gap simulation and β estimation) may be extended to simulations of irregular time series that have prescribed spectral slopes and also superimposed temporal trends, which can then be used to evaluate the validity of various statistical methods for identifying trends and their associated statistical significance.

4 Conclusions

River water-quality time series often exhibit fractal scaling behavior, which presents challenges to the identification of deterministic trends. Because traditional spectral estimation methods are generally not applicable to irregularly sampled time series, we have examined two broad types of estimation approaches and evaluated their performances against synthetic data with a wide range of prescribed β values and gap

intervals that are representative of the sampling irregularity of real water-quality data.

The results of this work suggest several important messages. First, the results remind us of the risks in using interpolation for gap filling when examining autocorrelation, as the interpolation methods consistently underestimate or overestimate β under a wide range of prescribed β values and gap distributions. Second, the widely used Lomb–Scargle spectral method also consistently underestimates β . Its modified form, using the 5 % lowest frequencies for spectral slope estimation, has very poor precision, although the overall bias is small. Third, the wavelet method, coupled with an aliasing filter, has the smallest bias and root-mean-squared error among all methods for a wide range of prescribed β values and gap distributions, except for cases with small prescribed β values (i.e., close to white noise) or small λ values (i.e., very skewed gap distributions). Thus, the wavelet method is recommended for estimating spectral slopes in irregular time series until improved methods are developed. In this regard, future research should aim to develop an aliasing filter that is more applicable to irregular time series with very skewed gap intervals. Finally, all methods' performances depend strongly on the sampling irregularity in terms of both the skewness and mean of gap-interval lengths, highlighting that the accuracy and precision of each method are data specific.

Overall, these results provide new contributions in terms of better understanding and quantification of the proposed methods' performances for estimating the strength of fractal scaling in irregularly sampled water-quality data. In addition, the work has provided an innovative and general approach for modeling sampling irregularity in water-quality records. Moreover, this work has proposed and demonstrated a generalizable framework for data simulation (with gaps) and β estimation, which can be readily applied to evaluate other methods that are not covered in this work. More generally, the findings and approaches may also be broadly applicable to irregularly sampled data in other scientific disciplines. Last but not least, we note that accurate quantification of fractal scaling in irregular water-quality time series remains an unresolved challenge for the hydrologic community and for many other disciplines that must grapple with irregular sampling.

Data availability. River monitoring data used in this study are available through the US Geological Survey National Water Information System (<https://doi.org/10.5066/F7P55KJN>) and Heidelberg University's National Center for Water Quality Research.

The Supplement related to this article is available online at <https://doi.org/10.5194/hess-22-1175-2018-supplement>.

Competing interests. The authors declare that they have no conflict of interest.

Acknowledgements. Zhang was supported by the Maryland Sea Grant through awards NA10OAR4170072 and NA14OAR1470090 and by the Maryland Water Resources Research Center through a graduate fellowship while he was a doctoral student at Johns Hopkins University. Subsequent support to Zhang was provided by the US EPA under grant “EPA/CBP Technical Support 2017” (no. 07-5-230480). Harman's contribution to this work was supported by the National Science Foundation through grants CBET-1360415 and EAR-1344664. We thank Bill Ball (Johns Hopkins University) and Bob Hirsch (US Geological Survey) for many useful discussions. We are very grateful to the Editor and two anonymous reviewers for their comments and suggestions. This is contribution no. 5449 of the University of Maryland Center for Environmental Science.

Edited by: Erwin Zehe

Reviewed by: two anonymous referees

References

- Aubert, A. H., Kirchner, J. W., Gascuel-Oudou, C., Fauchoux, M., Gruau, G., and Mérot, P.: Fractal water quality fluctuations spanning the periodic table in an intensively farmed watershed, *Environ. Sci. Technol.*, 48, 930–937, <https://doi.org/10.1021/es403723r>, 2014.
- Beran, J.: Long-range dependence, *Wiley Interdiscip. Rev. Comput. Stat.*, 2, 26–35, <https://doi.org/10.1002/wics.52>, 2010.
- Beran, J., Feng, Y., Ghosh, S., and Kulik, R.: *Long-Memory Processes: Probabilistic Properties and Statistical Methods*, Berlin, Heidelberg, Springer Berlin Heidelberg, 884 pp., 2013.
- Boutahar, M., Marimoutou, V., and Noura, L.: Estimation Methods of the Long Memory Parameter: Monte Carlo Analysis and Application, *J. Appl. Stat.*, 34, 261–301, <https://doi.org/10.1080/02664760601004874>, 2007.
- Box, G. E. P., Jenkins, G. M., and Reinsel, G. C.: *Time Series Analysis*, Fourth Edition. Hoboken, NJ, John Wiley & Sons, Inc., 47–92, 2008.
- Clarke, R. T.: Calculating uncertainty in regional estimates of trend in streamflow with both serial and spatial correlations, *Water Resour. Res.*, 49, 7120–7125, <https://doi.org/10.1002/wrcr.20465>, 2013.
- Cleveland, W. S.: LOWESS: A program for smoothing scatterplots by robust locally weighted regression, *Am. Stat.*, 35, 54, <https://doi.org/10.2307/2683591>, 1981.
- Cohn, T. A. and Lins, H. F.: Nature's style: Naturally trendy, *Geophys. Res. Lett.*, 32, L23402, <https://doi.org/10.1029/2005GL024476>, 2005.
- Constantine, W. and Percival, D.: *fractal: Fractal Time Series Modeling and Analysis*, available at: <https://cran.r-project.org/web/packages/fractal> (last access: 6 April 2015.), 2014.
- Darken, P. F., Zipper, C. E., Holtzman, G. I., and Smith, E. P.: Serial correlation in water quality variables: Estimation and implications for trend analysis, *Water Resour. Res.*, 38, 1117, <https://doi.org/10.1029/2001WR001065>, 2002.

- Delignette-Muller, M. L. and Dutang, C.: *fitdistrplus: An R Package for Fitting Distributions*, *J. Stat. Softw.*, 64, 1–34, 2015.
- Ehsanzadeh, E. and Adamowski, K.: Trends in timing of low stream flows in Canada: impact of autocorrelation and long-term persistence, *Hydrol. Process.*, 24, 970–980, <https://doi.org/10.1002/hyp.7533>, 2010.
- Fatichi, S., Barbosa, S. M., Caporali, E., and Silva, M. E.: Deterministic versus stochastic trends: Detection and challenges, *J. Geophys. Res.*, 114, D18121, <https://doi.org/10.1029/2009JD011960>, 2009.
- Foster, G.: Wavelets for period analysis of unevenly sampled time series, *Astron. J.*, 112, 1709–1729, 1996.
- Franzke, C.: Nonlinear Trends, Long-Range Dependence, and Climate Noise Properties of Surface Temperature, *J. Clim.*, 25, 4172–4183, <https://doi.org/10.1175/JCLI-D-11-00293.1>, 2012a.
- Franzke, C.: On the statistical significance of surface air temperature trends in the Eurasian Arctic region, *Geophys. Res. Lett.*, 39, L23705, <https://doi.org/10.1029/2012GL054244>, 2012b.
- Godsey, S. E., Aas, W., Clair, T. A., de Wit, H. A., Fernandez, I. J., Kahl, J. S., Malcolm, I. A., Neal, C., Neal, M., Nelson, S. J., Norton, S. A., Palucis, M. C., Skjelkvåle, B. L., Soulsby, C., Tetzlaff, D., and Kirchner, J. W.: Generality of fractal 1/f scaling in catchment tracer time series, and its implications for catchment travel time distributions, *Hydrol. Process.*, 24, 1660–1671, <https://doi.org/10.1002/hyp.7677>, 2010.
- Graham, J.: Missing Data Analysis: Making It Work in the Real World, *Annu. Rev. Psychol.*, 60, 549–576, <https://doi.org/10.1146/annurev.psych.58.110405.085530>, 2009.
- Helsel, D. R. and Hirsch, R. M.: *Statistical Methods in Water Resources*, US Geological Survey Techniques of Water-Resources Investigations Book 4, Chapter A3, US Geological Survey, Reston, VA, p. 522, <http://pubs.usgs.gov/twri/twri4a3/> (last access: 11 June 2016.), 2002.
- Hirsch, R. M., Alexander, R. B., and Smith, R. A.: Selection of methods for the detection and estimation of trends in water quality, *Water Resour. Res.*, 27, 803–813, <https://doi.org/10.1029/91WR00259>, 1991.
- Hirsch, R. M., Moyer, D. L., and Archfield, S. A.: Weighted regressions on time, discharge, and season (WRTDS), with an application to Chesapeake Bay river inputs, *J. Am. Water Resour. Assoc.*, 46, 857–880, <https://doi.org/10.1111/j.1752-1688.2010.00482.x>, 2010.
- Hurst, H. E.: Long-term storage capacity of reservoirs, *Trans. Amer. Soc. Civil Eng.*, 116, 770–808, 1951.
- Kasmarek, M. C. and Ramage, J. K.: Water-Level Measurement Data Collected during 2015-2016 and Approximate Long-term Water-Level Altitude Changes of Wells Screened in the Chicot, Evangeline, and Jasper Aquifers, Houston-Galveston Region, Texas: US Geological Survey data release, <https://doi.org/10.5066/F77H1GP3>, 2016.
- Khalik, M. N., Ouarda, T. B. M. J., and Gachon, P.: Identification of temporal trends in annual and seasonal low flows occurring in Canadian rivers: The effect of short- and long-term persistence, *J. Hydrol.*, 369, 183–197, <https://doi.org/10.1016/j.jhydrol.2009.02.045>, 2009.
- Khalik, M. N., Ouarda, T. B. M. J., Gachon, P., and Sushama, L.: Temporal evolution of low-flow regimes in Canadian rivers, *Water Resour. Res.*, 44, W08436, <https://doi.org/10.1029/2007WR006132>, 2008.
- Kirchner, J.: Aliasing in $1/f^{\alpha}$ noise spectra: Origins, consequences, and remedies, *Phys. Rev. E*, 71, 066110, <https://doi.org/10.1103/PhysRevE.71.066110>, 2005.
- Kirchner, J. W. and Neal, C.: Universal fractal scaling in stream chemistry and its implications for solute transport and water quality trend detection, *P. Natl. Acad. Sci. USA*, 110, 12213–12218, <https://doi.org/10.1073/pnas.1304328110>, 2013.
- Kirchner, J. W. and Weil, A.: No fractals in fossil extinction statistics, *Nature*, 395, 337–338, <https://doi.org/10.1038/26384>, 1998.
- Kirchner, J. W., Feng, X., and Neal, C.: Fractal stream chemistry and its implications for contaminant transport in catchments, *Nature*, 403, 524–527, <https://doi.org/10.1038/35000537>, 2000.
- Kirchner, J. W., Feng, X., and Neal, C.: Catchment-scale advection and dispersion as a mechanism for fractal scaling in stream tracer concentrations, *J. Hydrol.*, 254, 82–101, [https://doi.org/10.1016/S0022-1694\(01\)00487-5](https://doi.org/10.1016/S0022-1694(01)00487-5), 2001.
- Langland, M. J., Blomquist, J. D., Moyer, D. L., and Hyer, K. E.: Nutrient and suspended-sediment trends, loads, and yields and development of an indicator of streamwater quality at nontidal sites in the Chesapeake Bay watershed, 1985–2010, US Geological Survey Scientific Investigations Report 2012-5093, Reston, VA, p. 26., available at: <http://pubs.usgs.gov/sir/2012/5093/pdf/sir2012-5093.pdf> (last access: 6 April 2015), 2012.
- Lennartz, S. and Bunde, A.: Trend evaluation in records with long-term memory: Application to global warming, *Geophys. Res. Lett.*, 36, L16706, <https://doi.org/10.1029/2009GL039516>, 2009.
- Lomb, N. R.: Least-squares frequency analysis of unequally spaced data, *Astrophys. Space Sci.*, 39, 447–462, <https://doi.org/10.1007/BF00648343>, 1976.
- Montanari, A., Rosso, R., and Taquq, M. S.: Fractionally differenced ARIMA models applied to hydrologic time series: Identification, estimation, and simulation, *Water Resour. Res.*, 33, 1035–1044, <https://doi.org/10.1029/97WR00043>, 1997.
- Montanari, A., Taquq, M. S., and Teverovsky, V.: Estimating long-range dependence in the presence of periodicity: An empirical study, *Math. Comput. Model.*, 29, 217–228, [https://doi.org/10.1016/S0895-7177\(99\)00104-1](https://doi.org/10.1016/S0895-7177(99)00104-1), 1999.
- Montanari, A., Rosso, R., and Taquq, M. S.: A seasonal fractional ARIMA Model applied to the Nile River monthly flows at Aswan, *Water Resour. Res.*, 36, 1249–1259, <https://doi.org/10.1029/2000WR900012>, 2000.
- National Center for Water Quality Research: Tributary Data Download, <https://ncwqr.org/monitoring/data/> (last access: 23 July 2015), 2015.
- Noguchi, K., Gel, Y. R., and Duguay, C. R.: Bootstrap-based tests for trends in hydrological time series, with application to ice phenology data, *J. Hydro.*, 410, 150–161, <https://doi.org/10.1016/j.jhydrol.2011.09.008>, 2011.
- R Development Core Team: R: A language and environment for statistical computing, R Foundation for Statistical Computing, Vienna, Austria, <http://www.r-project.org> (last access: 6 April 2015), , 2014.
- Rea, W., Oxley, L., Reale, M., and Brown, J.: Estimators for Long Range Dependence: An Empirical Study, *Electron. J. Stat.*, <http://arxiv.org/abs/0901.0762> (last access: 6 April 2015), 2009.
- Sang, Y.-F., Wang, Z., and Liu, C.: Comparison of the MK test and EMD method for trend identification in

- hydrological time series, *J. Hydrol.*, 510, 293–298, <https://doi.org/10.1016/j.jhydrol.2013.12.039>, 2014.
- Scargle, J. D.: Studies in Astronomical Time-Series Analysis. II. Statistical Aspects of Spectral-Analysis of Unevenly Spaced Data, *Astrophys. J.*, 263, 835–853, <https://doi.org/10.1086/160554>, 1982.
- Stroe-Kunold, E., Stadnytska, T., Werner, J., and Braun, S.: Estimating long-range dependence in time series: an evaluation of estimators implemented in R, *Behav. Res. Meth.*, 41, 909–923, <https://doi.org/10.3758/BRM.41.3.909>, 2009.
- Taqqu, M. S., Teverovsky, V., and Willinger, W.: Estimators for long-range dependence: an empirical study, *Fractals*, 3, 785–798, <https://doi.org/10.1142/S0218348X95000692>, 1995.
- Wang, Z.: cts: An R Package for Continuous Time Autoregressive Models via Kalman Filter, *J. Stat. Softw.*, 53, 1–19, 2013.
- Witt, A. and Malamud, B. D.: Quantification of Long-Range Persistence in Geophysical Time Series: Conventional and Benchmark-Based Improvement Techniques, *Surv. Geophys.*, 34, 541–651, <https://doi.org/10.1007/s10712-012-9217-8>, 2013.
- Yue, S., Pilon, P., Phinney, B., and Cavadias, G.: The influence of autocorrelation on the ability to detect trend in hydrological series, *Hydrol. Process.*, 16, 1807–1829, <https://doi.org/10.1002/hyp.1095>, 2002.
- Zeileis, A. and Grothendieck, G.: zoo: S3 Infrastructure for Regular and Irregular Time Series, *J. Stat. Softw.*, 14, 1–27, 2005.
- Zetterqvist, L.: Statistical Estimation and Interpretation of Trends in Water Quality Time Series, *Water Resour. Res.*, 27, 1637–1648, <https://doi.org/10.1029/91wr00478>, 1991.
- Zhang, Q. and Ball, W. P.: Improving Riverine Constituent Concentration and Flux Estimation by Accounting for Antecedent Discharge Conditions, *J. Hydrol.*, 547, 387–402, <https://doi.org/10.1016/j.jhydrol.2016.12.052>, 2017.
- Zhang, Q., Brady, D. C., Boynton, W. R., and Ball, W. P.: Long-Term Trends of Nutrients and Sediment from the Nontidal Chesapeake Watershed: An Assessment of Progress by River and Season, *J. Am. Water Resour. Assoc.*, 51, 1534–1555, <https://doi.org/10.1111/1752-1688.12327>, 2015.