

Nature's style: Naturally trendy

Timothy A. Cohn and Harry F. Lins

U.S. Geological Survey, Reston, Virginia, USA

Received 29 August 2005; revised 29 September 2005; accepted 12 October 2005; published 8 December 2005.

[1] Hydroclimatological time series often exhibit trends. While trend *magnitude* can be determined with little ambiguity, the corresponding *statistical significance*, sometimes cited to bolster scientific and political argument, is less certain because significance depends critically on the null hypothesis which in turn reflects subjective notions about what one expects to see. We consider statistical trend tests of hydroclimatological data in the presence of long-term persistence (LTP). Monte Carlo experiments employing FARIMA models indicate that trend tests which fail to consider LTP greatly overstate the statistical significance of observed trends when LTP is present. A new test is presented that avoids this problem. From a practical standpoint, however, it may be preferable to acknowledge that the concept of statistical significance is meaningless when discussing poorly understood systems. **Citation:** Cohn, T. A., and H. F. Lins (2005), Nature's style: Naturally trendy, *Geophys. Res. Lett.*, 32, L23402, doi:10.1029/2005GL024476.

1. Introduction

[2] Hydroclimatological records (henceforth "HC") such as discharge and air temperature are increasingly examined for evidence of a *structural shift* or *trend*, defined as an upward or downward tendency in the data over time. There is typically little argument about the magnitude of observed trends whether estimated by eye or statistical methods [Craigmile *et al.*, 2004] (although D. Koutsoyiannis (personal communication, 2005) has expressed doubts about the existence of a rigorous and consistent definition of trend). The *statistical significance*, or *p-value*, associated with an observed trend, however, is more difficult to assess because it depends on subjective assumptions about the underlying stochastic process [von Storch and Zwiers, 1999; Woodward and Gray, 1993; Weatherhead *et al.*, 1998]. In this paper, we consider the idea introduced by Hurst [1951] and discussed by others [Mandelbrot and Wallis, 1969a; Klemeš, 1974; Lettenmaier and Burges, 1978; Potter, 1976; Potter and Walker, 1981; Hosking, 1984; Bras and Rodriguez-Iturbe, 1985; Vogel *et al.*, 1998; Koutsoyiannis, 2000] that HC records are realizations of physical processes whose behavior exhibits long-term persistence (LTP). Such behavior is sometimes modeled as fractional Gaussian noise (fGn) or fractionally differenced ARIMA (FARIMA or *arfima*) processes. The purpose of this paper is not to evaluate claims related to LTP, but rather to explore what LTP, if present, implies about the significance of observed trends.

2. A Family of Trend Models

[3] We assume that an HC record, $\vec{Y} \equiv (Y_1, \dots, Y_N)'$, arises from a stochastic process, and that the process can be

partitioned into a deterministic linear trend component and a stochastic component [Kendall *et al.*, 1983; Craigmile *et al.*, 2004] such that

$$Y_t = \mu + \beta \cdot t + \epsilon_t \quad (1)$$

where t represents time (conveniently discretized into $(1, 2, \dots, N)$), μ is a location parameter, β is the trend coefficient (the change per unit time), and ϵ_t represents the "error."

[4] The errors are assumed to be multivariate normal with zero mean and covariance matrix Σ . The LTP, autoregressive, or moving average structure, if present, is completely characterized by Σ . To simplify the analysis, we constrain Σ to be a function of ϕ (a lag-one autoregression (AR(1)) parameter); d (the fractional differencing parameter, sometimes described by H , the Hurst coefficient, where $H = d + 0.5$); θ (a lag-one moving average (MA(1)) parameter); and σ (a scale parameter). The complete stochastic process corresponding to equation 1 is denoted by $S_{\beta, \{\phi, d, \theta\}}(t)$, where the parameters μ and σ can be omitted without loss of generality.

[5] Stationarity is an important issue if we wish to determine whether long-term "excursions" observed in the data should be attributed to ordinary process dynamics around a fixed mean versus permanent structural changes to the process. Precise conditions for stationarity of $S_{\beta, \{\phi, d, \theta\}}(t)$ are given by Kendall *et al.* [1983]; however, necessary conditions include $\beta = 0$ and $d < 0.5$.

[6] All stationary stochastic processes, $S_{0, \{\phi, 0, \theta\}}(t)$, where $d = 0$, exhibit the following property: For observations far apart in time, the correlation between $S(t)$ and $S(t + k)$ is bounded by: $\rho_k \leq c^{|k|}$ as $k \rightarrow \infty$ where c is a constant and $|c| < 1$ [Koutsoyiannis, 2000], which implies short-term persistence in the sense that the covariance structure involves exponential decay.

[7] The stochastic process $S_{0, \{\phi, d, \theta\}}(t)$, $0.5 > d > 0$, exhibits long-term persistence [Hosking, 1984]. The correlation between observations is given by [Hosking, 1984]: $\rho_k = \Gamma(1 - d)\Gamma(k + d)/(\Gamma(d)\Gamma(k + 1 - d)) \approx \Gamma(1 - d)\Gamma(d)k^{2d-1}$ where $\Gamma(\cdot)$ denotes the complete gamma function. When $0.5 > d > 0$, the correlation declines "slowly", as a power function in k . More important, as Mandelbrot and Wallis [1969b, pp. 230–231] observed, "[a] perceptually striking characteristic of fractional noises is that their sample functions exhibit an astonishing wealth of 'features' of every kind, including trends and cyclic swings of various frequencies." It is easy to imagine that LTP could be mistaken for trend.

3. Implications for Hypothesis Testing

[8] Trend assessment seeks to answer two questions:

[9] 1. What is the approximate magnitude of the trend, β ?

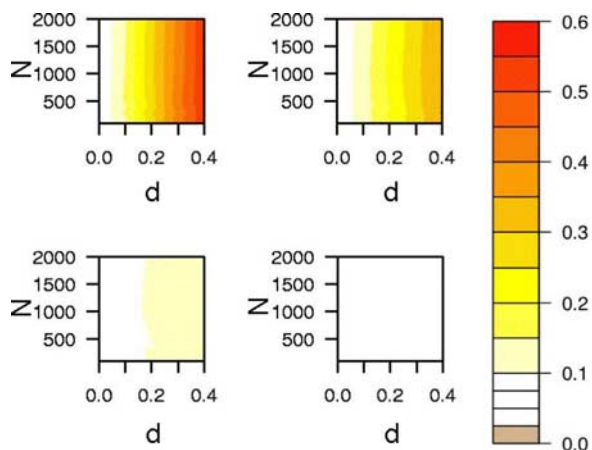


Figure 1. Observed type 1 error rate for trend tests at $\alpha = 5\%$ level, as function of sample size (N) and fractional difference parameter (d). From upper left to lower right, the four contour plots correspond to: $T_{\beta, \{\phi, 0, 0, 0\}}$; $T_{\beta, \{\phi, \phi, 0, 0\}}$; $T_{\beta, \{0, d, 0\}}$, and $T_{\beta, \{0, d, 0\}}^A$. Note that the white areas in the plots indicate type 1 error rate between 2.5% and 10%, which is about the nominal level.

[10] 2. Given what we believe about the stochastic process, how likely is it that we would have observed \bar{Y} , or something more extreme, if the true value of β is 0?

[11] The standard procedure for addressing these questions is to fit the parameters in equation 1 to the observed \bar{Y} and obtain an estimate $\hat{\beta}$. The corresponding p -value, which is the probability of observing a value at least as extreme as $\hat{\beta}$ if $\beta = 0$, is then computed. This is a straightforward exercise if we know the sampling distribution of β under the null hypothesis (H_0).

[12] The simplest case involves processes with white noise errors ($S_{0, \{0, 0, 0\}}(t)$), for which an efficient test for linear trend can be obtained by fitting an ordinary least squares (OLS) regression model with time as a predictor variable and testing to see if the fitted coefficient on time, β , differs significantly from zero. This test is denoted $T_{\beta, \{0, 0, 0\}}$, and it is the uniformly most powerful unbiased (UMPU) test if in fact the process is $S_{\beta, \{0, 0, 0\}}(t)$ [Kendall and Stuart, 1979].

[13] For more complicated stochastic processes, such as $S_{\beta, \{\phi, 0, 0\}}(t)$ and $S_{\beta, \{0, d, 0\}}$, the trend slope is computed by maximum likelihood using an approximation to the likelihood function [Hosking, 1984]. Statistical significance is computed using likelihood ratio tests [Kendall and Stuart, 1979], which are discussed in the online auxiliary materials¹. The tests are denoted $T_{\beta, \{\phi, 0, 0\}}$, $T_{\beta, \{0, d, 0\}}$, etc. Craigmile *et al.* [2004, 2005] consider a wavelet-based fitting method for essentially the same model.

[14] It happens that the standard likelihood ratio test (LRT) has less than ideal statistical properties. In particular, for large values of d the LRT does not come close to achieving its nominal $\alpha = 5\%$ level when H_0 is true, even for very large sample sizes. It is easy to “adjust” the LRT [Kendall and Stuart, 1979], however, to ensure that the approximate type I error rate is achieved. The adjusted test

(ALRT), denoted $T_{\beta, \{0, d, 0\}}^A$, is discussed in detail in the online auxiliary materials.

4. Trend Test Performance

[15] Monte Carlo experiments were conducted using the R programming language and the *fracdiff* package. The *fracdiff.sim* routine permits generation of simulated $S_{0, \{\phi, d, \theta\}}(t)$ trend-free time series. Linear trends can be superimposed on the $S_{0, \{\phi, d, \theta\}}(t)$ series to generate $S_{\beta, \{\phi, d, \theta\}}(t)$ series for arbitrary β .

[16] The *fracdiff* package includes a routine, *fracdiff*, for fitting the parameters of an $S_{0, \{\phi, d, \theta\}}(t)$ process to data. This routine was embedded in a loop (using the R *optimize* routine) to enable fitting the trend coefficient, β , by maximizing the value of the likelihood function (which is computed by *fracdiff*).

[17] The Monte Carlo approach used here requires simulating an approximation of the natural processes, and this requires some assumptions. In particular, the value of d , or at least a range of reasonable values for d , must be specified. Beran and Feng’s [2002] 663 year flow record for the Nile River exhibits $d = 0.39$. Hurst [1951] found that $d = H - 0.5 \approx 0.23$ for a variety of geophysical time series. Vogel *et al.* [1998] looked at the Hurst coefficient corresponding to the USGS’s Hydroclimatic Data Network (HCDN) data set [Slack and Landwehr, 1992], and found that, given the shortness of the records, the correlation structure could be explained either by LTP or by non-LTP Box-Jenkins ARMA processes. Assuming that the correlation structure was due exclusively to LTP, Vogel reported that the interquartile range of d for streamflow records in the United States was approximately (0.3–0.4). To ensure that the range of simulated populations could represent the range of characteristics of data observed in HC data, 35 separate experiments were run comprising all combinations of samples of size $N = \{100, 200, 300, 400, 500, 1000, 2000\}$ and fractional differencing values of $d = \{0, 0.1, 0.2, 0.3, 0.4\}$.

4.1. Type I Error Rates

[18] The first set of experiments was designed to determine the true type I error rate (for a nominal 5% test) for each of the trend tests as a function of the true value of d and N . Time series were generated without trend (i.e., $\beta = 0$). Four trend tests were used to determine if a trend was present at the $\alpha = 5\%$ level: $T_{\beta, \{0, 0, 0\}}$ (white noise); $T_{\beta, \{\phi, 0, 0\}}$ (autoregressive); $T_{\beta, \{0, d, 0\}}$ (LRT with fractional differencing); $T_{\beta, \{0, d, 0\}}^A$ (ALRT with fractional differencing).

[19] Figure 1 depicts the actual type I error rates ($\alpha = 5\%$ test), for each of the four tests, as a function of the true value of d and the sample size N . For small d , all of the tests exhibit type 1 error rates in the “white” contour – in the range 2.5% to 10% – reasonably close to the nominal level of 5%. For $d \geq 0.3$ (a plausible level for HC processes), however, the type I error rates exceed 10% for all but the ALRT test, and generally exceed 50% for $T_{\beta, \{0, 0, 0\}}$ regardless of sample size. This is indicated by contours depicted in increasingly dark colors.

[20] It is possible to condense Figure 1 into a single graph because the results do not vary substantially with sample size (N). Figure 2 depicts the case where $N = 100$, and shows clearly that the standard trend tests, particularly

¹Auxiliary material is available at <ftp://ftp.agu.org/apend/gl/2005GL024476>.

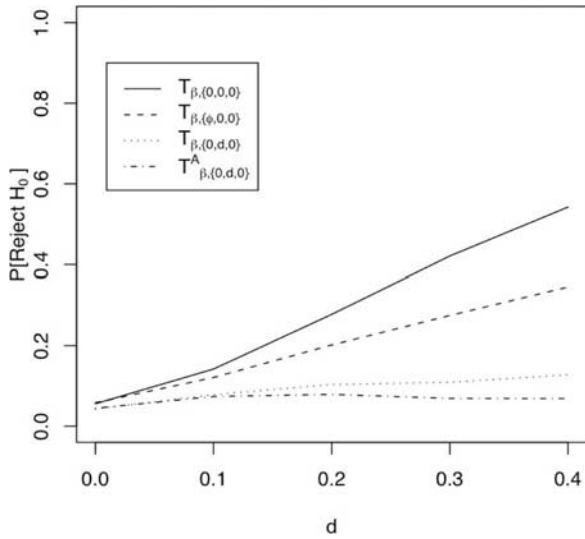


Figure 2. The probability of rejecting H_0 when H_0 is true as a function of d for the case $N = 100$.

$T_{\beta, \{0,0,0\}}$, become increasingly likely to find statistical significance as d increases. It is worrisome that when LTP is present, and well within the range observed for many natural phenomena ($d \approx 0.35$), the commonly used OLS trend test ($T_{\beta, \{0,0,0\}}$) is likely to report significant trends about half the time when we know there is no trend in the stochastic process. The AR(1) trend test ($T_{\beta, \{\phi,0,0\}}$), though better, has a type 1 error rate 5 times larger than the nominal value when $d = 0.4$. The LRT ($T_{\beta, \{0,d,0\}}$) is better still, although its type I error rates are close to 15% in some cases. Among these tests, only the ALRT test ($T_{\beta, \{0,d,0\}}^A$) comes close to achieving the nominal $\alpha = 5\%$ type 1 error rate when substantial LTP is present.

4.2. Power and Type II Error Rates

[21] Figure 3 shows power curves for each of the four trend tests when no LTP is present ($d = 0$). These curves indicate the probability of rejecting H_0 when it is false or,

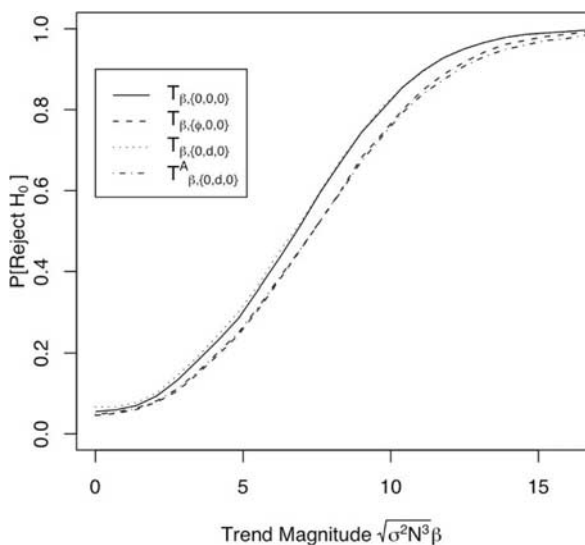


Figure 3. Power curves indicating the probability of rejecting H_0 when H_0 is false, as function of true value of trend magnitude $b \equiv \sqrt{\sigma^2 N^3} \beta$ for the case $d = 0$ and $N = 100$.

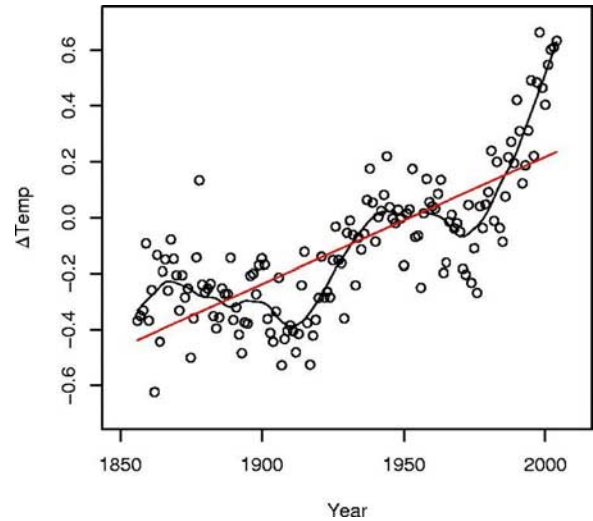


Figure 4. Annual departures from the period-of-record mean northern hemisphere temperature in degrees C, 1856–2002, with least squares fit (red line) and loess smooth (black line).

stated differently, the probability of correctly identifying a trend when a trend is, in fact, present. The standard “OLS” test, $T_{\beta, \{0,0,0\}}$, is known to be optimal in this case. The power curves are plotted as a function of the real trend (expressed in terms of the nearly invariant $b = \sqrt{\sigma^2 N^3} \beta$) for a sample size of $N = 100$ and no fractional differencing ($d = 0$). The most powerful test in this case is $T_{\beta, \{0,0,0\}}$, and as long as ϕ , d , and θ are all known to be zero, there is not much doubt about which test to use. It is noteworthy, however, that all of the tests, and particularly $T_{\beta, \{0,d,0\}}^A$, are only slightly less powerful than $T_{\beta, \{0,0,0\}}$ at detecting a real trend in the absence of LTP. Thus, the penalty for using the alternative tests is small.

5. An HC Example

[22] Figure 4 presents annual departures from the long-term mean in northern hemisphere surface air temperature (“NHT”) during the last century and a half [Jones et al., 1999]. To gain some perspective on the two issues discussed in sections 3 and 4, we can apply trend tests to this well known data set whose LTP and trend properties have been considered by Smith [1993], Smith and Chen [1996], Beran and Feng [2002], and Craigmile et al. [2005].

[23] Table 1 contains estimates of the trends and corresponding p-values for the tests discussed in section 3

Table 1. Estimates of Trend Magnitudes and p-values Corresponding to Various Models Fitted to the Annual Northern Hemisphere Temperature Departure Data, 1856–2002

H_0 Process	Test	$\hat{\beta}^a$	p-Value
White noise	$T_{\beta, \{0,0,0\}}$	0.0045	1.8e-27
MA(1)	$T_{\beta, \{0,0,0\}}$	0.0046	1.9e-21
AR(1)	$T_{\beta, \{\phi,0,0\}}$	0.0047	5.2e-11
LTP	$T_{\beta, \{0,d,0\}}$	0.0050	4.8e-3
LTP	$T_{\beta, \{0,d,0\}}^A$	0.0050	9.4e-3
ARMA(1,1)	$T_{\beta, \{\phi,0,0\}}$	0.0053	1.7e-4
LTP + MA(1)	$T_{\beta, \{0,d,0\}}$	0.0045	7.2e-2
LTP + AR(1)	$T_{\beta, \{\phi,d,0\}}$	0.0045	7.1e-2

^aTrend magnitude, $\hat{\beta}$, is expressed in units of $^\circ\text{C}/\text{year}$.

applied to the $N = 149$ temperature observations. All of the tests report nearly the same estimated trend magnitude ($\hat{\beta}$), which ranges from 0.0045 to 0.0053 °C/year. As far as the magnitude is concerned, it makes little difference which test is used. Choice of trend test, however, does matter when computing trend *significance*. The simplest test, $T_{\beta,\{0,0,0\}}$ (which assumes no LTP), finds strong evidence of trend, a p-value of 1.8×10^{-27} . $T_{\beta,\{\phi,0,0\}}$ (which allows for short-term persistence) yields a p-value of 5.2×10^{-11} , 16 orders of magnitude larger and still highly significant. The p-value corresponding to either $T_{\beta,\{0,d,0\}}$ or $T_{\beta,\{\phi,d,0\}}$, an unadjusted LRT trend test that considers both short-term and long-term persistence, is about 7%, which is not significant under the null hypothesis. In changing from one test to another, 25 orders of magnitude of significance vanished. This result is somewhat troubling given the uncertainty about the stochastic process and consequently about which test to rely on.

6. Discussion and Conclusions

[24] The problems with significance testing are well-documented [McCloskey, 1995; Nicholls, 2000], and significance testing for HC trends is particularly problematical because we do not know what null hypothesis to use. Because statistical tests are proofs by contradiction, any inconsistency between the null hypothesis and the natural system can itself lead to rejection of the null hypothesis. As demonstrated in section 4.1 above, rejection of H_0 can occur because $\beta \neq 0$ (the hoped for explanation) or because $d \neq 0$ and the trend test does not recognize the possibility of LTP. In short, the presence of LTP in a stochastic process can induce a significant trend result when no trend is present, if an inappropriate trend test is used.

[25] The question remains whether natural HC processes in fact possess LTP. The idea was introduced more than 50 years ago by Hurst [1951], and has been debated ever since [Mandelbrot and Wallis, 1968; Klemeš, 1974; Potter and Walker, 1981; Hosking, 1984; Loucks et al., 1981; Koutsoyiannis, 2000, 2003]. Hurst's fundamental finding has neither been discredited nor universally embraced, but persuasive arguments have been presented (for discussion and additional references, see Koutsoyiannis [2003]). Given the LTP-like patterns we see in longer HC records, however, such as the periods of multidecadal drought that occurred during the past millennium and our planet's geologic history of ice ages and sea level changes, it might be prudent to assume that HC processes could possess LTP.

[26] In any case, powerful trend tests are available that can accommodate LTP [Hosking, 1984; Craigmile et al., 2005]. In particular, Hosking [1984] developed a unified approach for modeling fractional Gaussian noise as a generalization of ARIMA models [Box et al., 1994] and provided a practical technique for fitting data exhibiting LTP. Moreover, the ALRT test presented here, which is based on Hosking's approach, is both accurate (in the sense that it comes close to achieving its nominal α -level), and nearly as powerful as the commonly used OLS procedure when applied to processes with little or no persistence. It is therefore surprising that nearly every assessment of trend significance in geophysical variables published during the past few decades has failed to account properly for long-term persistence.

[27] These findings have implications for both science and public policy. For example, with respect to temperature

data there is overwhelming evidence that the planet has warmed during the past century. But could this warming be due to natural dynamics? Given what we know about the complexity, long-term persistence, and non-linearity of the climate system, it seems the answer might be yes. Finally, that reported trends are real yet insignificant indicates a worrisome possibility: natural climatic excursions may be much larger than we imagine. So large, perhaps, that they render insignificant the changes, human-induced or otherwise, observed during the past century.

[28] **Acknowledgments.** This paper benefited significantly from comments by M. Beran, W. Kirby, D. Koutsoyiannis, K. Potter, G. Schwarz, and S. Vecchia.

References

- Beran, J., and Y. Feng (2002), SEMIFAR models—A semiparametric approach to modelling trends, *Comput. Stat. Data Anal.*, 40(2), 393–419.
- Box, G. E. P., G. M. Jenkins, and G. C. Reinsel (1994), *Time Series Analysis: Forecasting and Control*, 3rd ed., Prentice Hall, Upper Saddle River, N. J.
- Bras, R. L., and I. Rodriguez-Iturbe (1985), *Random Functions in Hydrology*, 1st ed., Addison-Wesley, Boston, Mass.
- Craigmile, P. F., P. Guttorp, and D. B. Percival (2004), Trend assessment in a long memory dependence model using the discrete wavelet transform, *Environmetrics*, 15(4), 35–313.
- Craigmile, P. F., D. B. Percival, and P. Guttorp (2005), Wavelet based estimation for polynomial contaminated fractionally differenced processes, *IEEE Trans. Signal Process.*, 53(8), 3151–3161.
- Hosking, J. (1984), Modeling persistence in hydrological time series using fractional differencing, *Water Resour. Res.*, 20(12), 1898–1908.
- Hurst, H. E. (1951), Long term storage capacities of reservoirs, *Trans. Am. Soc. Civ. Eng.*, 116, 776–808.
- Jones, P. D., M. New, D. E. Parker, S. Martin, and I. G. Rigor (1999), Surface air temperature and its changes over the past 150 years, *Rev. Geophys.*, 37, 173–199. (Data available at <http://www.cru.uea.ac.uk/ftpdata/tavenh2v.dat>)
- Kendall, M., and A. Stuart (1979), *The Advanced Theory of Statistics*, vol. 2, *Inference and Relationship*, 4th ed., 748 pp., Oxford Univ. Press, New York.
- Kendall, M., A. Stuart, and J. K. Ord (1983), *The Advanced Theory of Statistics*, vol. 3, *Design and Analysis, and Time Series*, 4th ed., 780 pp., Oxford Univ. Press, New York.
- Klemeš, V. (1974), The Hurst phenomenon: A puzzle?, *Water Resour. Res.*, 10(4), 675–688.
- Koutsoyiannis, D. (2000), A generalized mathematical framework for stochastic simulation and forecast of hydrologic time series, *Water Resour. Res.*, 36(6), 1519–1533.
- Koutsoyiannis, D. (2003), Climate change, the Hurst phenomenon, and hydrologic statistics, *Hydrol. Sci.*, 48(1), 3–24.
- Lettenmaier, D. P., and S. J. Burges (1978), Climate change: Detection and its impact on hydrologic design, *Water Resour. Res.*, 14(4), 679–687.
- Loucks, D. P., J. Stedinger, and D. Haiht (1981), *Water Resource Systems Planning and Analysis*, 559 pp., Prentice Hall, Upper Saddle River, N. J.
- Mandelbrot, B. B., and J. R. Wallis (1968), Noah, Joseph, and operational hydrology, *Water Resour. Res.*, 4(5), 909–918.
- Mandelbrot, B. B., and J. R. Wallis (1969a), Some long-run properties of geophysical records, *Water Resour. Res.*, 5(2), 321–340.
- Mandelbrot, B. B., and J. R. Wallis (1969b), Computer experiments with fractional Gaussian noises: 1. Averages and variances, *Water Resour. Res.*, 5(1), 228–241.
- McCloskey, D. (1995), The insignificance of statistical significance, *Sci. Am.*, 272(4), 32–33.
- Nicholls, N. (2000), The insignificance of significance testing, *Bull. Am. Meteorol. Soc.*, 81(5), 981–986.
- Potter, K. W. (1976), Evidence of nonstationarity as a physical explanation of the Hurst phenomenon, *Water Resour. Res.*, 12(5), 1047–1052.
- Potter, K. W., and J. F. Walker (1981), A model of discontinuous measurement error and its effects on the probability distribution of flood discharge measurements, *Water Resour. Res.*, 21, 1505–1509.
- Slack, J., and J. Landwehr (1992), Hydro-climatic data network (HCDN)—A U.S. Geological Survey streamflow data set for the United States for the study of climate fluctuations, 1874–1988, *U.S. Geol. Surv. Open File Rep.* 92-129.

- Smith, R. (1993), Long-range dependence and global warming, in *Statistics for the Environment*, edited by V. Barnett and F. Turkman, pp. 141–161, John Wiley, Hoboken, N. J.
- Smith, R., and F.-L. Chen (1996), Regression in long-memory time series, in *Athens Conference on Applied Probability and Time Series*, vol. 2, *Time Series Analysis in Memory of E. J. Hannan*, edited by P. Robinson and M. Rosenblatt, *Springer Lecture Notes Stat.*, 115, 378–391.
- Vogel, R. M., Y. Tsai, and J. F. Limbrunner (1998), The regional persistence and variability of annual streamflow in the United States, *Water Resour. Res.*, 34(12), 3445–3459.
- von Storch, H., and F. W. Zwiers (1999), *Statistical Analysis in Climate Research*, 484 pp., Cambridge Univ. Press, New York.
- Weatherhead, E., et al. (1998), Factors affecting the detection of trends: Statistical considerations and applications to environmental data, *J. Geophys. Res.*, 103(D14), 17,149–17,161.
- Woodward, W., and H. Gray (1993), Global warming and the problem of testing for trend in time series data, *J. Clim.*, 6(5), 953–962.

T. A. Cohn and H. F. Lins, U.S. Geological Survey, MS 415, Reston, VA 20192, USA. (tacohn@usgs.gov; hlins@usgs.gov)