

# Stochastic modelling of skewed data exhibiting long-range dependence

S.M. Papalexiou, Department of Water Resources, National Technical University of Athens;

## 1. Abstract

Time series with long-range dependence appear in many fields including hydrology, and there are several studies that have provided evidence of long autocorrelation tails. Provided that the intensity of the long-range dependence in time series of a certain process quantified by the self-similarity parameter, also known as the Hurst exponent  $H$ , could not be falsified, it is then essential that the variable of interest is modelled by a model reproducing long-range dependence. Common models of this category that have been widely used are the fractional Gaussian noise (FGN) and the fractional ARIMA (FARIMA). In case of a variable exhibiting skewness, the previous models can not be implemented in a direct manner. In order to preserve skewness in the simulated series, a normalizing transformation is typically applied in the real-life data at first. The models are then fitted to the normalized data and the produced synthetic series are finally de-normalized. In this paper, a different method is proposed, consisting of two parts. The first one regards the approximation of the long-range dependence by an autoregressive model of high order  $p$  AR( $p$ ), while the second one regards the direct calculation of the main statistical properties of the random component, that is mean, variance and skewness coefficient. The skewness coefficient calculation of the random component is done using joint sample moments. The advantage of the method is its efficiency and simplicity and the analytical solution.

## 2. Motivation

Since Hurst (1951) observed the long-term persistence phenomenon in the annual average streamflows of Nile, the same behaviour has been identified in numerous natural processes while, its importance has been underlined by scientists in many controversial disciplines. It seems that the Hurst phenomenon is ubiquitous in nature and this makes it necessary to find adequate ways to model it.

Many models have been proposed in the literature that preserve the Hurst behaviour, such as FGN (Mandelbrot, 1971), broken time averages (e.g. Diebevec, 1971), fractional ARMA (Hosking, 1981), and recently symmetric moving average models (SMA) (Koussourakis, 2000, 2002).

If the Hurst behaviour appears in a process, it needs to be modelled as it affects dramatically the time series structure. Another distinguished characteristic of hydrological processes, that needs to be modelled, is asymmetry. In this direction have been made many attempts to adapt standard models to preserve the skewness (e.g. Mandala and Wallis, 1976).

Some of the previous models are not easy to apply as the parameters are not easy to estimate, while other can preserve the skewness but not the Hurst behaviour and vice versa. Other problems are the narrow type of autocorrelation functions that those models can simulate (exception is the SMA model).

In this study a general methodology to preserve both the Hurst behaviour and skewness in the time series is proposed. The Hurst phenomenon is modelled from the innovation of the high order AR( $p$ ), which can be done easily by the easy to apply skewness coefficient of the random component of the model. The model should be able to apply and suitable for any practical purposes such as hydrologic design or water resources management.

## 3. Modelling Approach

In order to preserve the long-range dependence or the Hurst phenomenon in the simulated time series, a high order autoregressive model is implemented. The long-range dependence behaviour, is essentially the slow decay of the autocorrelation function with time. On the contrary, the AR( $p$ ) models are considered to be short-range dependence models. Nevertheless, as this study reveals, AR( $p$ ) models of high order can reproduce the Hurst phenomenon sufficiently enough for any practical modelling purposes.

In the general case of order  $p$ , the AR( $p$ ) model takes the following form:  $X_t = \epsilon_t + \sum_{i=1}^p \phi_i X_{t-i}$  where  $\epsilon_t$  is the innovation on the random component and  $\phi_i$  are coefficients. In order to fit the model to a dataset, the  $\phi_i$  coefficients and the basic statistics (mean, standard deviation) of the  $\epsilon_t$  have to be estimated.

The auto-covariance function  $\gamma_k$  of the AR( $p$ ) model for  $\log k$  and for  $k > 0$  is given by  $\gamma_k = \sum_{j=0}^{p-1} \phi_j^k$

The replacement of  $\gamma_k$  with the samples estimates and the implementation of the last equation  $p$  times gives a linear system of equations that can be solved straightforwardly, evaluating therefore the  $\phi_i$  coefficients.

Finally, the mean and the variance of the  $\epsilon_t$  can be estimated using the following two equations. 
$$\mu_{\epsilon_t} = \mu_{X_t} \left( 1 - \sum_{i=1}^p \phi_i \right)^{-1} \quad \sigma_{\epsilon_t}^2 = \gamma_0 - \sum_{i=1}^p \phi_i \gamma_i$$

## 4. Preserving the Skewness in an AR( $p$ ) Model

To preserve asymmetry in the simulated time series, it is necessary to evaluate the skewness coefficient of the innovation,  $C_{\epsilon_t}$ . It can be shown that the third central moment of the innovation of the AR( $p$ ) model is 
$$\mu_{\epsilon_t}^{(3)} = \mu_{X_t}^3 - 3 \mu_{X_t} \sigma_{X_t}^2 + \frac{3}{2} \sigma_{X_t}^4 \left( \sum_{i=1}^p \phi_i \right)^{-3} \left( X_{t-1} - \mu_{X_t} \right) \dots (X_{t-p} - \mu_{X_t}) \quad (1)$$
 it can be proven that the following equation is valid, 
$$\mu_{\epsilon_t}^{(3)} = \mu_{X_t}^3 - 3 \mu_{X_t} \sigma_{X_t}^2 + \frac{3}{2} \sigma_{X_t}^4 \sum_{i=1}^p \phi_i^2 + 6 \sum_{i=1}^{p-1} \sum_{j=i+1}^p \phi_i \phi_j \sum_{k=i+j}^p \phi_k \mu_{X_t}^{(3-k)}$$
 Replacing the multi-auto-covariance terms in the previous equation with the sample estimates, given by 
$$\hat{\mu}_{\epsilon_t}^{(3)} = \frac{1}{k} \sum_{i=1}^k \sum_{j=1}^k \sum_{l=1}^k \left( X_{t+i} - \hat{\mu}_{X_t} \right) \left( X_{t+j} - \hat{\mu}_{X_t} \right) \left( X_{t+l} - \hat{\mu}_{X_t} \right)$$
 It is then straightforward to estimate the  $\mu_{\epsilon_t}^{(3)}$  in (1) and thus the  $C_{\epsilon_t} = \frac{\mu_{\epsilon_t}^{(3)}}{\sigma_{\epsilon_t}^3}$ .

## 5. The Generalized AutoCorrelation Function (GACF)

The major criticism of a high order AR( $p$ ) model would focus on the lack of parsimony, as estimation of the autocorrelation function up to  $\log p$  is required to fit the model. Moreover, it is well known that the estimator of the ACF is highly variable and that it increases its variability with increasing lag (Bres and Rodriguez, 1983). Consequently, the uncertainty in the estimation of the ACF would lead to uncertain validation of the model parameters. To overcome this disadvantage, it is proposed to fit a generalized ACF,  $\hat{\rho}_k^{(g)}$ , to the first few empirical ACF values (where  $\alpha, \beta, \gamma$  are positive parameters and  $f$  is the lag). Subsequently, the fitted GACF can be used to extrapolate ACF values for high  $f$ .

The figure depicts the empirical autocorrelation function of the Nilometer dataset (analysis follows) and the fitted GACF. The GACF has been fitted to the first 10 empirical values of the ACF by minimizing the square error.

$$\hat{\rho}_k^{(g)} = (\alpha \rho_k^{\beta} + \gamma)^{\frac{1}{\gamma}}$$

## 6. The Generalized Lambda Distribution (GLD)

In order to preserve the skewness in the simulated series, the innovation  $\epsilon_t$  must be sampled from a flexible family of distributions, such as the Generalized Lambda Distribution (GLD), proposed by John Tucker (1960) and generalized for Monte Carlo simulation purposes by John Ramirez and Bruce Schmeiser (1974). Although the GLD has been applied in many fields since the early 1970s (Karim and Dudevec, 2000), it has never been used in hydrology.

The GLD family with parameters  $\lambda_1, \lambda_2, \lambda_3$  and  $\lambda_4$ , is defined in terms of its percentile function, where  $0 < \gamma < 1$ . The parameters  $\lambda_1$  and  $\lambda_2$  are, respectively, location and scale parameters, while  $\lambda_3$  and  $\lambda_4$  are the generalized lambda density function is  $f(x) = \frac{\lambda_2}{\lambda_1} \gamma^{1-\lambda_3} (1-\gamma)^{\lambda_3-1} \lambda_4$  at  $x = Q(\gamma)$

The restrictions on  $\lambda_1, \lambda_2, \lambda_3$  and  $\lambda_4$  that yield a valid GLD distribution, the parameter space and the corresponding regions in the parameter space are given by Ramirez and Dudevec (2000). In the next figure GLD pdfs are plotted with mean = 0, variance = 1 and skewness coefficient,  $C_{\epsilon_t}$  ranging from 0 to -4.5

## 7. Fitting the GLD and Sampling

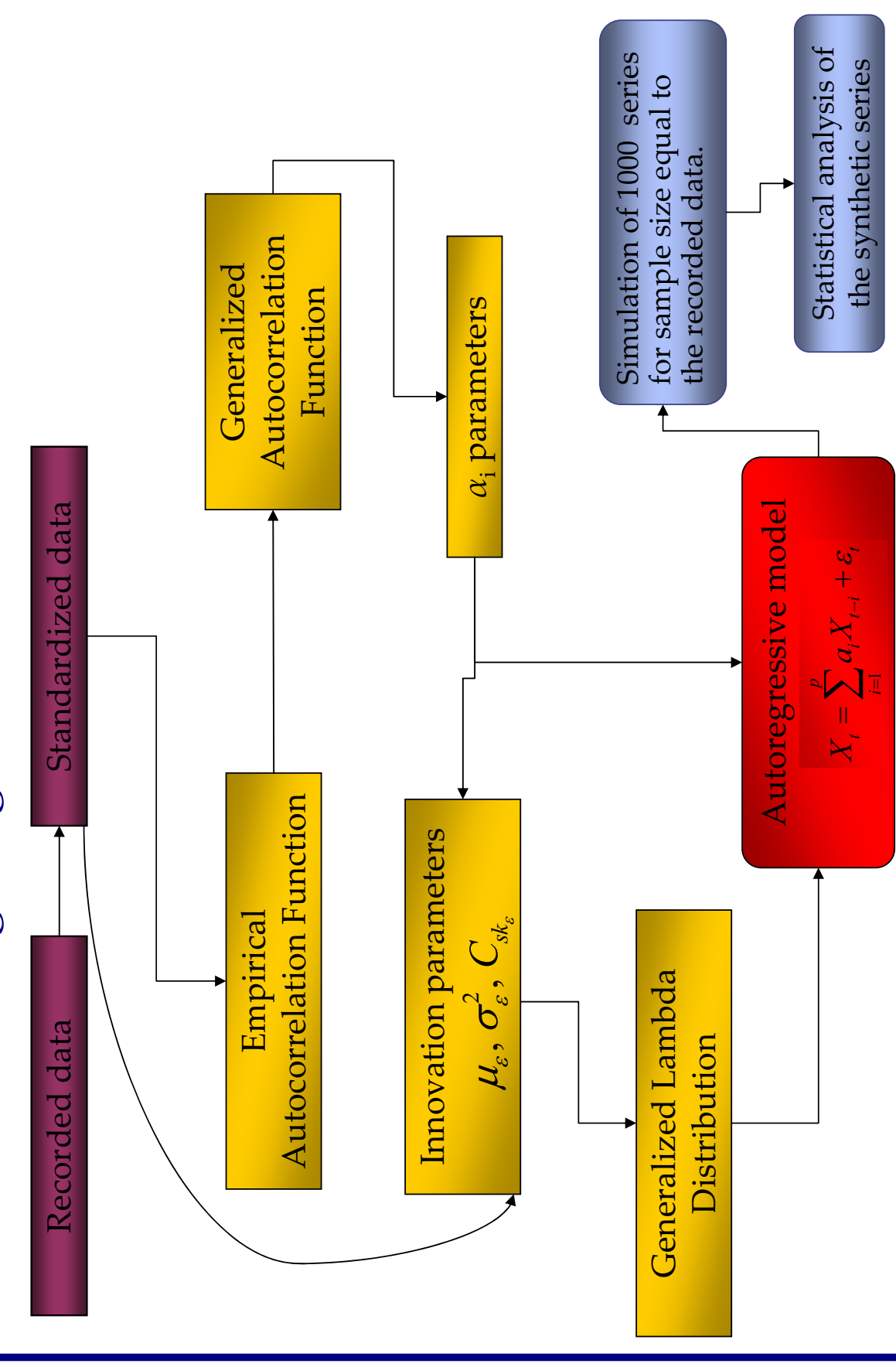
If  $X$  is GLD( $\lambda_1, \lambda_2, \lambda_3, \lambda_4$ ) with  $\lambda_2 > -1/A$  and  $\lambda_3 > -1/A$ , then its first four moments (Ramirez et al., 1970),  $\mu, \mu_2, \mu_3, \mu_4$  (mean, variance, skewness coefficient, and kurtosis coefficient), are given by 
$$\mu = \lambda_1 + \lambda_2 \quad \mu_2 = \lambda_2^2 \frac{1-\lambda_3}{\lambda_3} \quad \mu_3 = 2\lambda_2^3 \frac{1-\lambda_3}{\lambda_3^2} \quad \mu_4 = \frac{3}{2} \lambda_2^4 \frac{1-\lambda_3}{\lambda_3^3} \quad C_{\epsilon_t} = \frac{\mu_3 - 3\mu\mu_2}{\sigma_{\epsilon_t}^3} \quad \kappa_{\epsilon_t} = \frac{\mu_4 - 3\mu\mu_3 + 2\lambda_2\mu_2\mu_3}{\sigma_{\epsilon_t}^4}$$
 where  $A = \frac{1}{\lambda_1 + \lambda_2} + \frac{1}{\lambda_3}$

$$B = 3\lambda_2(\lambda_3 + 1) \left[ \lambda_4 - \frac{1}{2\lambda_3} + \frac{1}{2\lambda_3^2} \right] \quad C = 3\lambda_2(\lambda_3 + 1)(2\lambda_4 + 1) + \frac{3}{\lambda_3^2} \left[ \frac{1}{\lambda_3} - \frac{1}{\lambda_3^2} \right] \quad D = -4\lambda_2(\lambda_3 + 1)(3\lambda_4 + 1) + 6\lambda_2(\lambda_3 + 1)(2\lambda_4 + 1) - 4\lambda_2(\lambda_3 + 1) \left[ \frac{1}{\lambda_3} + \frac{1}{\lambda_3^2} \right]$$

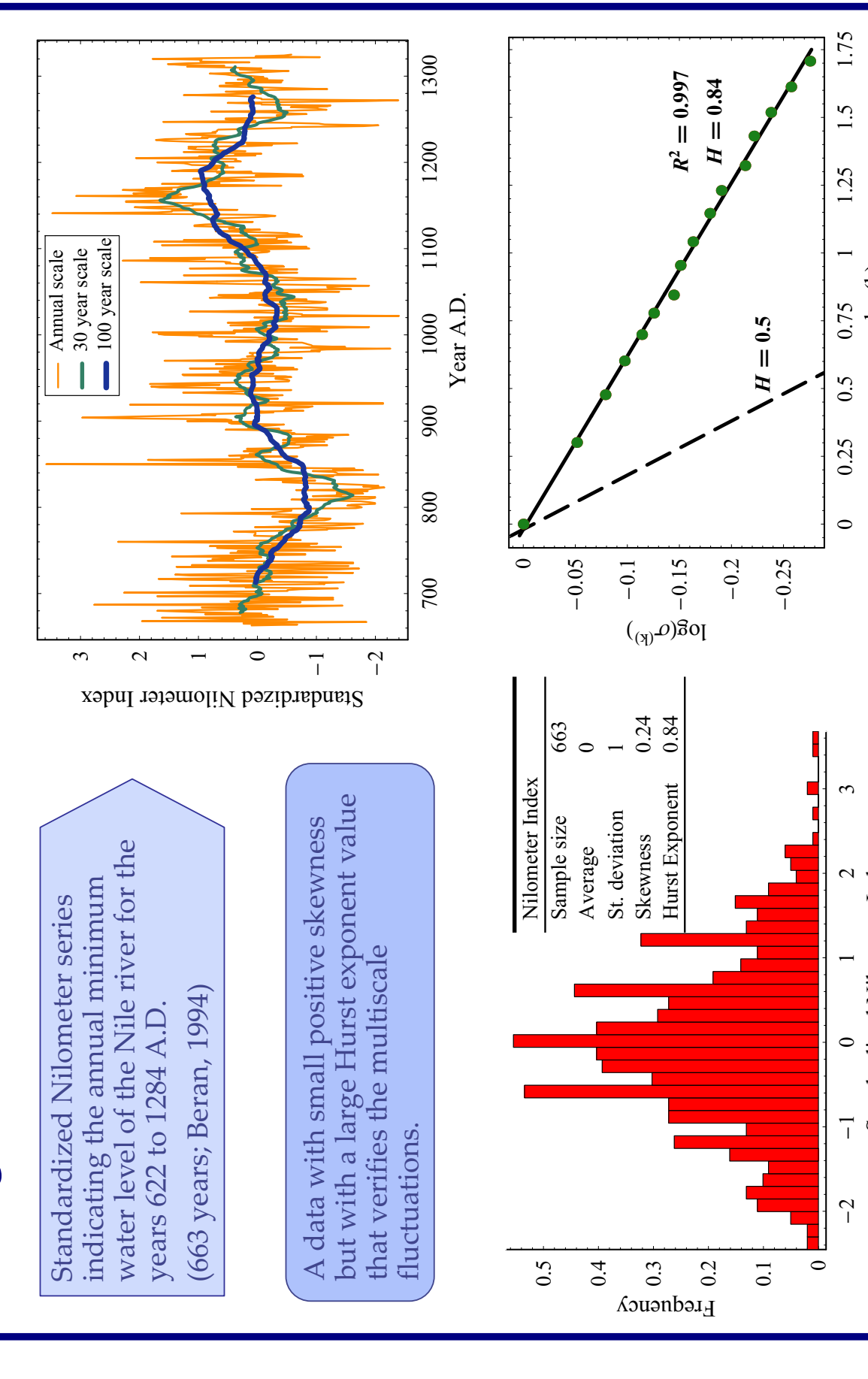
If  $B$  is the Beta function defined as  $B(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt$  and if we consider the innovation  $\epsilon_t$  as a random variable with known estimation of the mean, the variance, the skewness and the kurtosis coefficient, a GLD distribution can be fitted by solving numerically the previous nonlinear system. The mean, the variance and the skewness coefficient of  $\epsilon_t$  can be analytically estimated as described in slide four. At the moment there is no analytical way to estimate the kurtosis coefficient of  $\epsilon_t$ , but heuristically for this study was taken the minimum so as  $\lambda_2 > 0$  and  $\lambda_3 < 0$ , which implies that the fitted GLD ranges from  $-\infty$  to  $\infty$  (Karim and Dudevec, 2000).

Once the parameters  $\lambda_1, \lambda_2, \lambda_3$  and  $\lambda_4$  of the GLD are estimated, the sampling is very easy as the percentile function has a simple and analytical formulae.

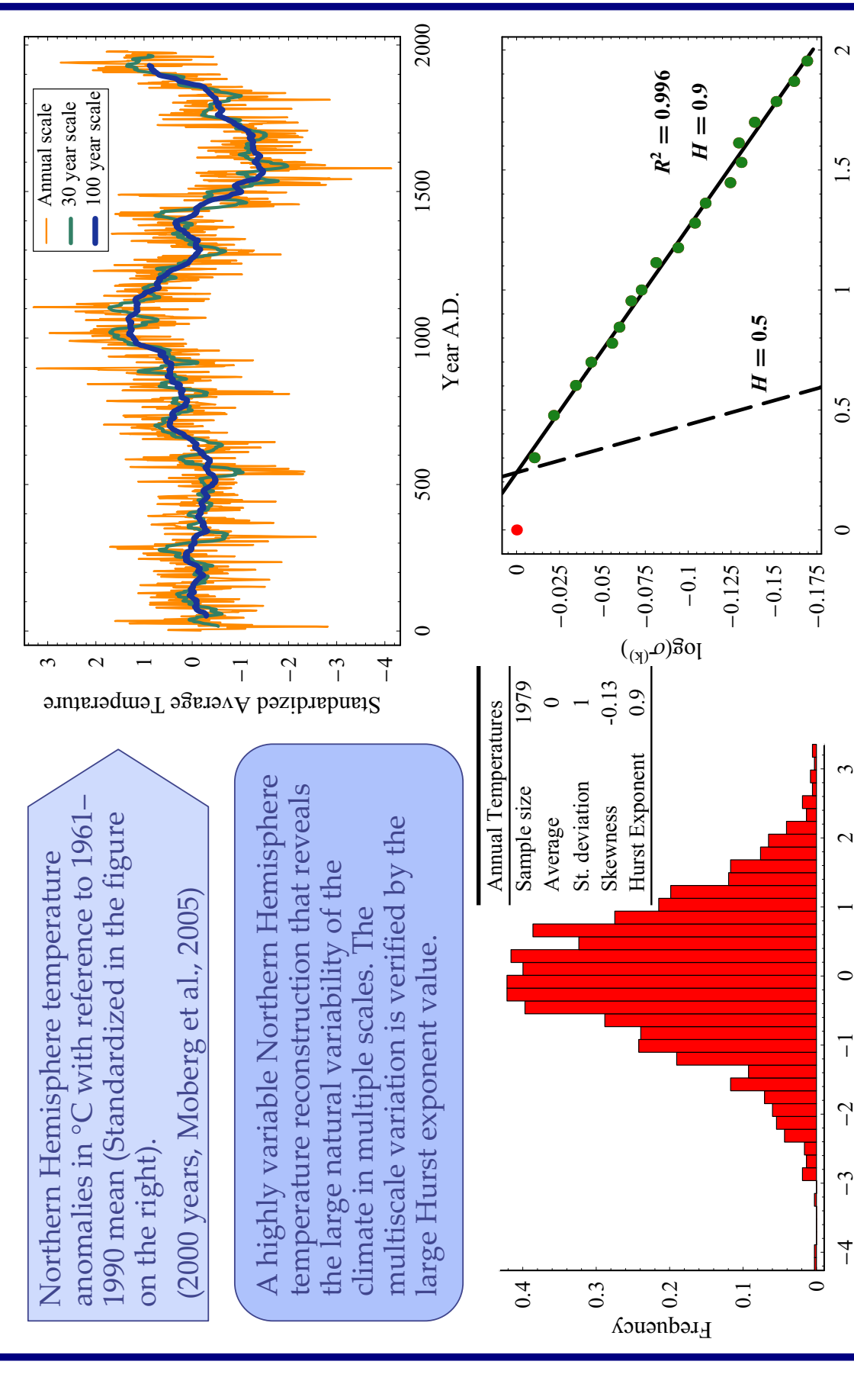
## 8. Simulation Organogram



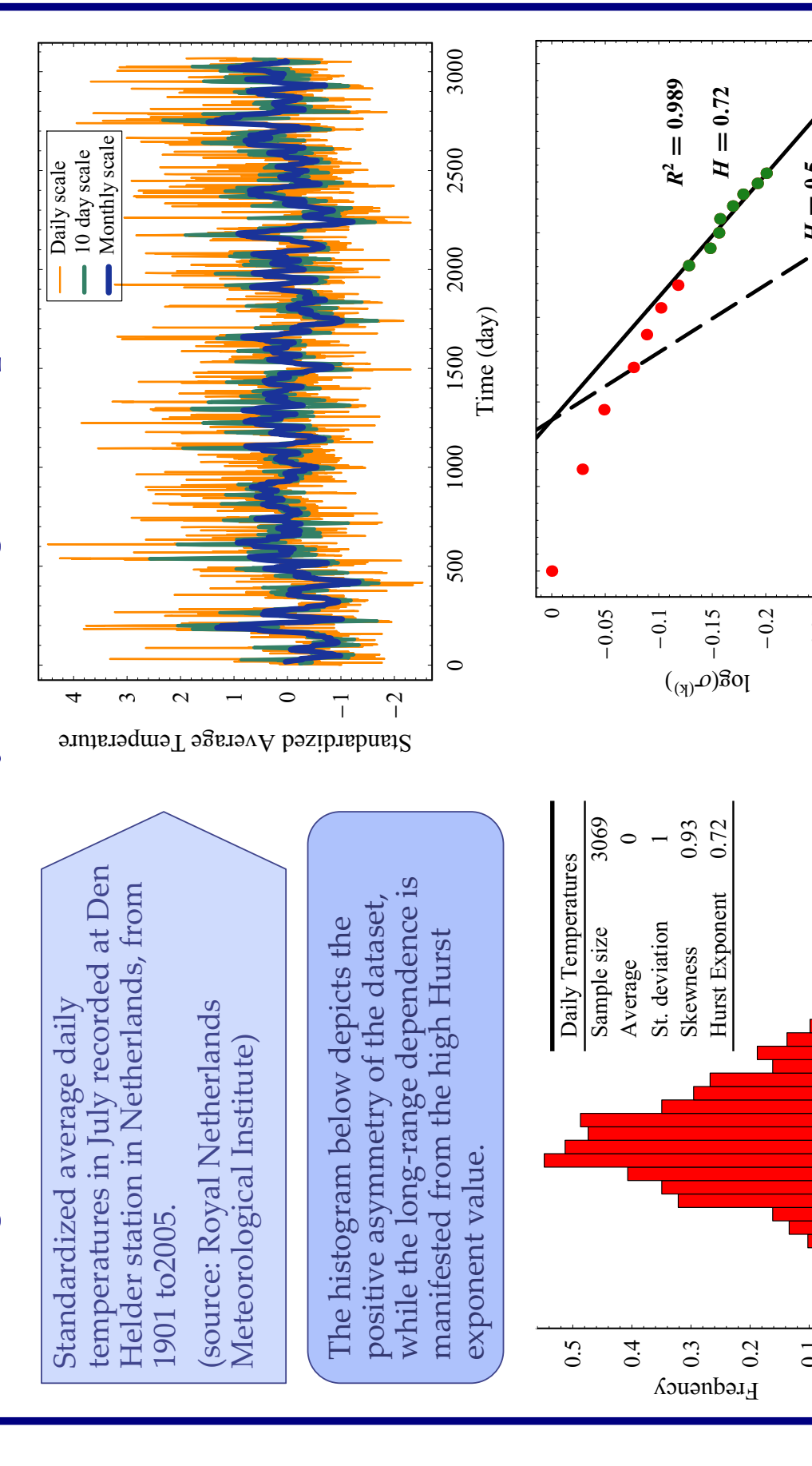
## 9. Original Data I: Nilometer Index



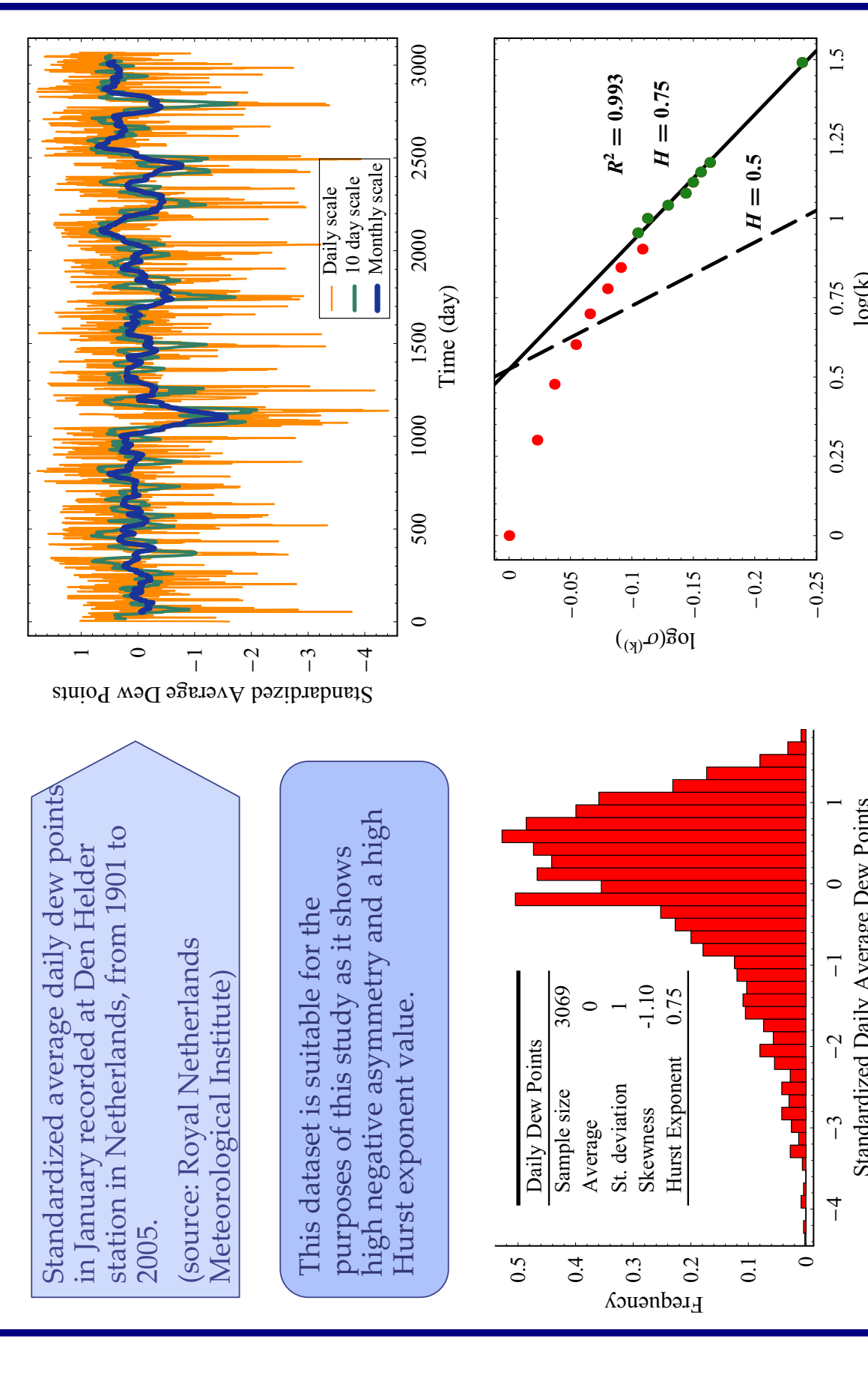
## 10. Original Data II: Annual Temperatures



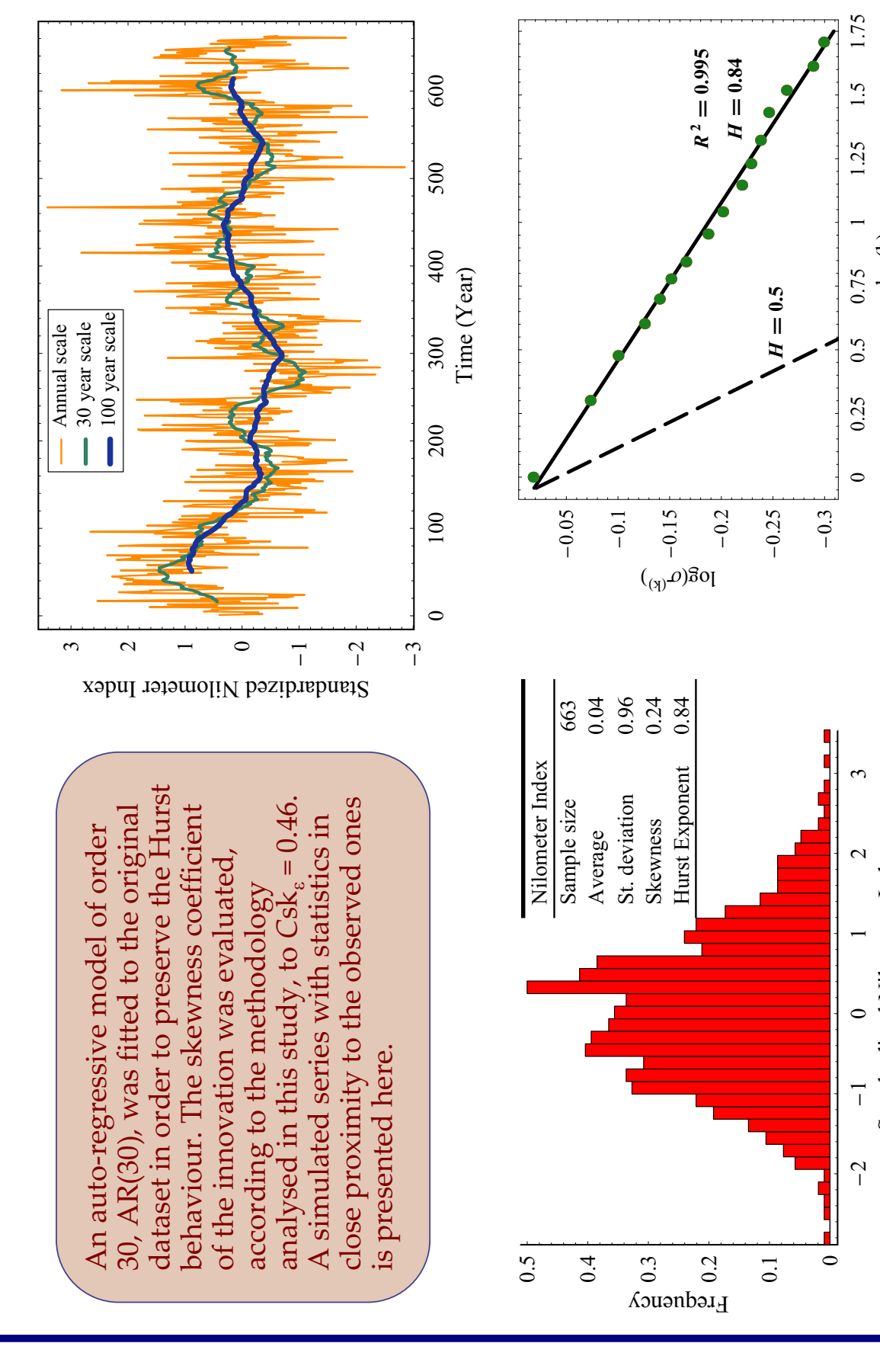
## 11. Original Data III: Daily Average Temperatures



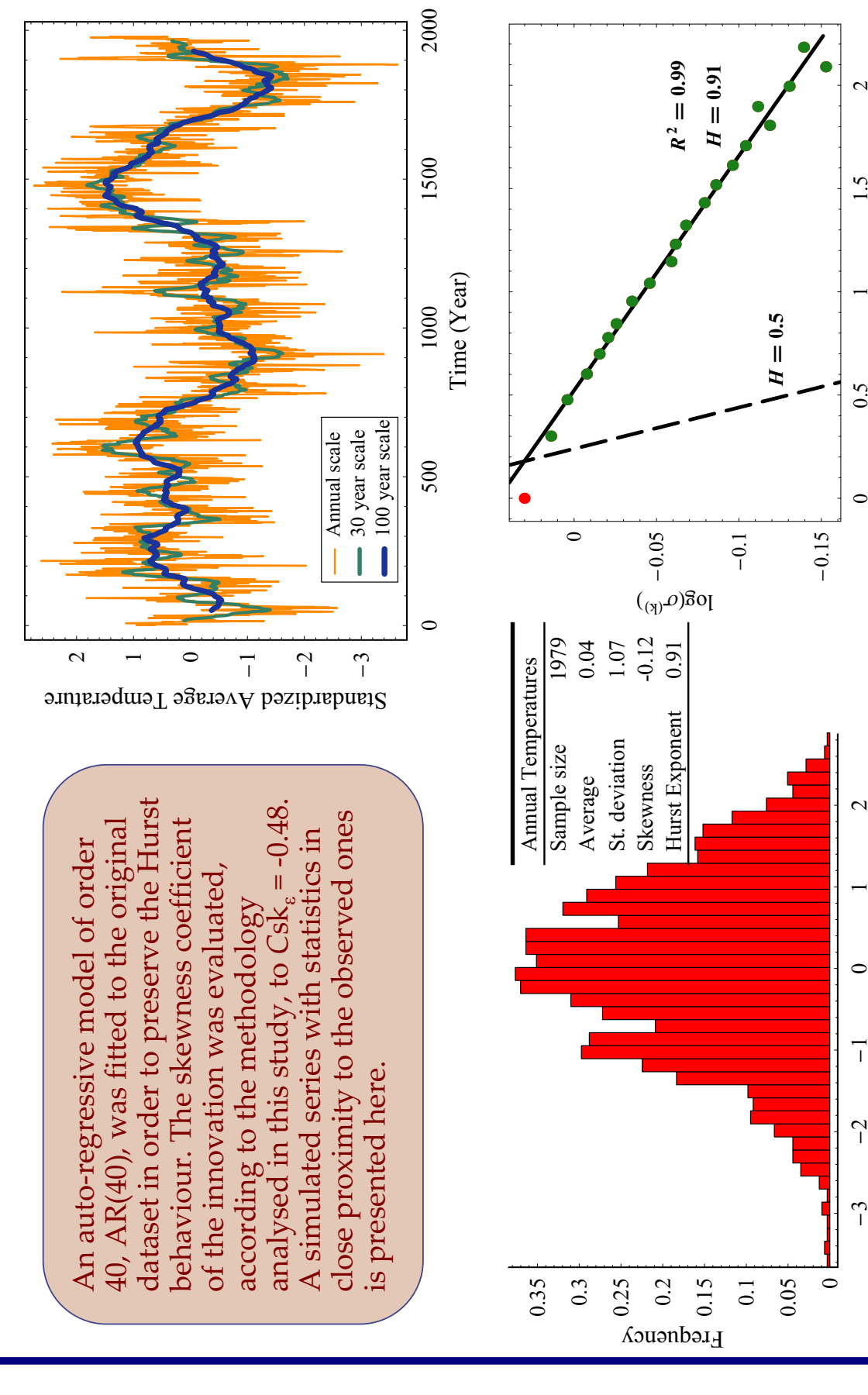
## 12. Original Data IV: Daily Average Dew Points



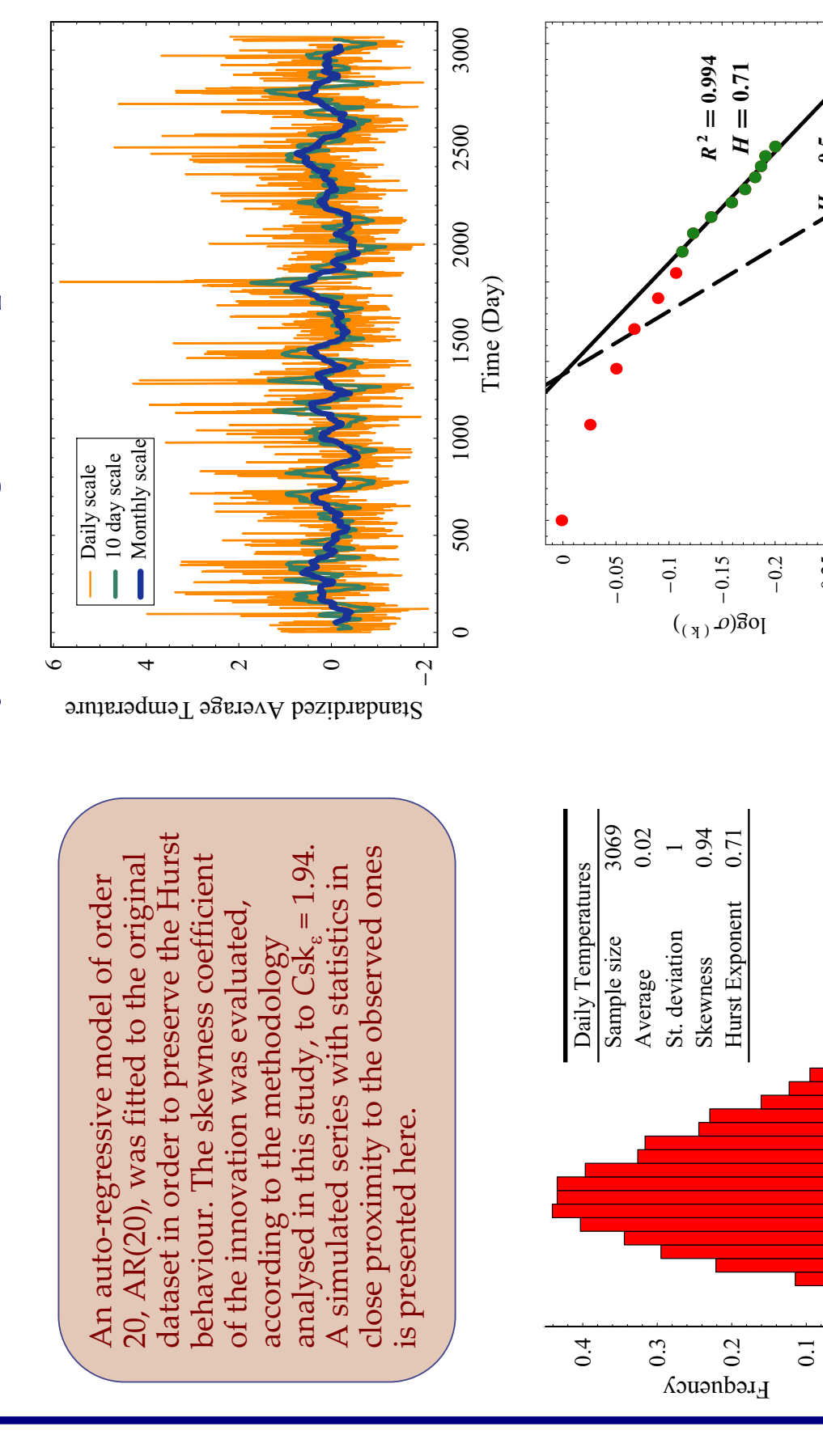
## 13. Simulated Data I: Nilometer Index



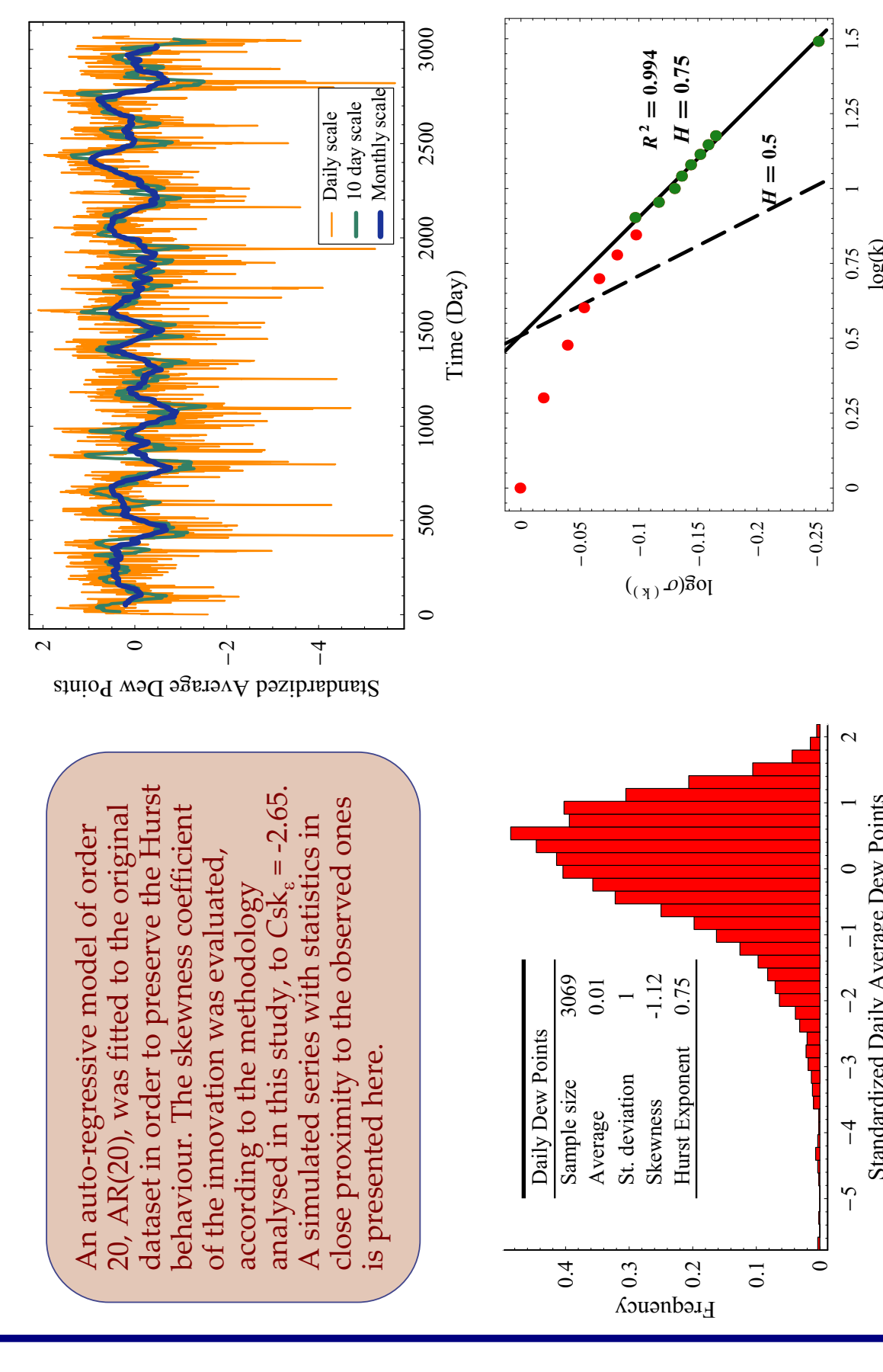
## 14. Simulated Data II: Global Temperatures



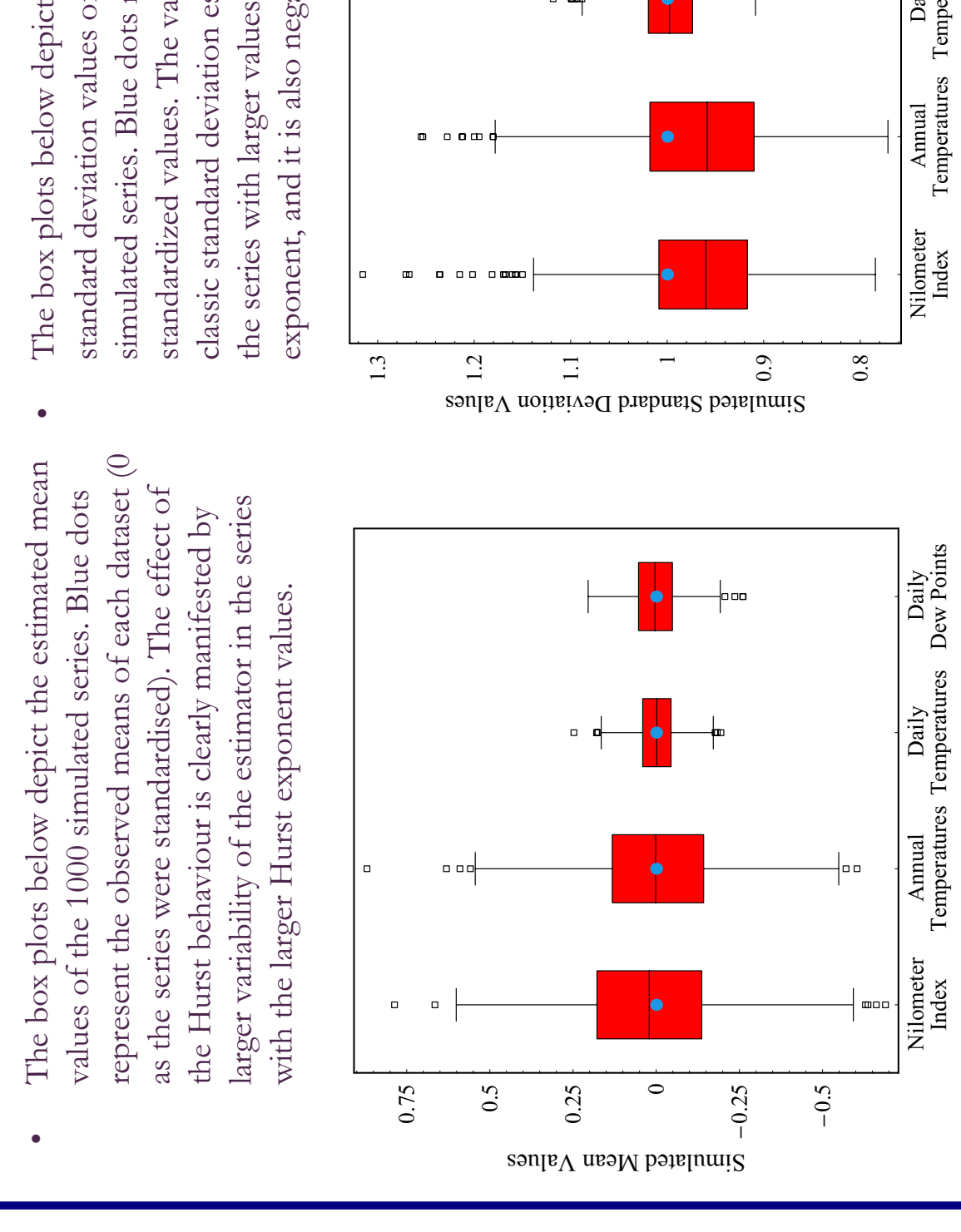
## 15. Simulated Data III: Daily Average Temperatures



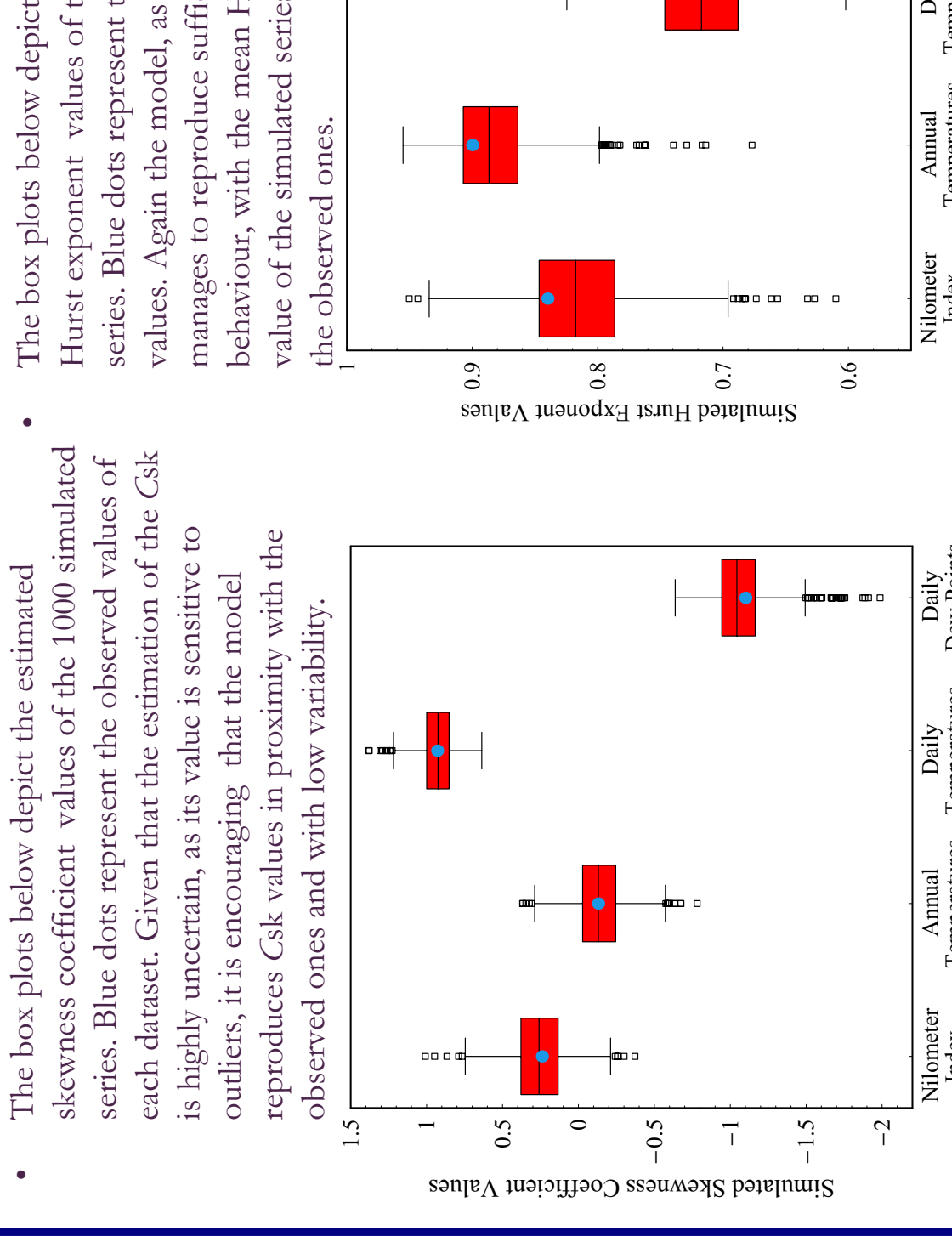
## 16. Simulated Data IV: Daily Average Dew Points



## 17. Simulated Mean and Standard Deviation



## 18. Simulated Skewness and Hurst Exponent



## 19. Conclusions

- While the autoregressive models are considered to be short range persistence models, it is concluded in this study that a higher order AR model preserves adequately the Hurst behaviour, for Hurst exponent values as high as 0.9. It seems that the model can preserve even more intense long-term persistence but this needs to be further examined.
- To preserve the asymmetry, an analytical expression for the estimation of the skewness coefficient of the innovation is given. Subsequently, the innovation sequence is sampled from a flexible skewed distribution, the so-called Generalized Lambda Distribution. The model manages to preserve sufficiently the skewness as the mean skewness coefficient of the simulated series is in proximity with the observed ones.
- As the simulated series are in accordance with the observed ones, the model can be used for any practical modelling purposes.
- Overall, the proposed methodology is simple and robust.

## 20. References

Brenys, J. Statistics for Long Memory processes, Volume 61 of Monographs on Statistics and Applied Probability, New York, 1985.

Diebevec, J. M. and Rodriguez Iturbe, Random Functions and Hydrology, 589 pp, Addison-Wesley Longman, Reading, Mass., 1971.

Hosking, J. R. M. Fractional Differentiation, Biometrika, vol. 68, pp. 163-174, 1981.

Hurst, H. H. Long term storage capacity of reservoirs, Transactions of the American Society of Civil Engineers 116, 770-808, 1957.

Karim, Z., Dudevec, E., Fring, Statistical Distributions: The Generalized Lambda Distribution and Generalized Bootstrap Methods, Boreas, CRC Press, 2000.

Koussourakis, S. M. An analytical framework for stochastic simulation and forecast of hydrologic time series. *Water Resources Research*, 36(6), 1919-1933, 2000.

Mandelbrot, B. B. The Hurst phenomenon and fractional Gaussian noise: some comments. *Hydrological Science Journal*, 4(1), 573-580, 1968.

Mandala, H. B. Use class processes with long-range dependence for simulation: an application to the hydrologic de H. E. *Journal of Hydrology*, 35(3), 251-263, 1978.

Ramirez, J. P., and Wallis, J. R. Wills, Computer experiments with fractional Gaussian noises, 1, Averages and variances. *Water Resources Research*, 12(4), 1033-1043, 1976.

Ramirez, J. P., and Wallis, J. R. Wills, Computer experiments with fractional Gaussian noises, 2, Rescaled ranges and spectra, *Water Resources Research*, 12(4), 1045-1056, 1976.

Gaussian noises, 3, Multivariate approach. *Water Resources Res.*, 12(4), 267-276, 1976.

Hosking, J. R. M., and Wallis, J. R. Wills, Computer experiments with fractional Gaussian noises, 4, Rescaled ranges and spectra. *Water Resources Research*, 12(4), 1057-1066, 1976.

Hosking, J. R. M., Wallis, J. R. Wills, and Wood, E. W., Methods for the analysis of extreme events, Communications of the Institute of Mathematical Statistics, 24(4), 1059-1072, 1970.

Samko, S. S., Kilbas, A. A., Marichev, O. I., Fractional Integrals and Derivatives. Theory and Applications, G. Y. T. Yung, Ed., Heldermann Verlag, Berlin, 1993.

Tucker, J., *Journal of the Royal Statistical Society*, vol. 22, 1960.

Wallis, J. R. Wills, and Wood, E. W., An approximation method for generating asymmetric random variables, Communications of the Institute of Mathematical Statistics, 24(4), 1073-1077, 1970.

Ramirez, J. P., S. Dudevec, E., J. T. Koutoukoulas, P. R. Molykides, E. A. probability distribution and its use in fitting data. *Communications in Statistics—Theory and Methods*, 21(23), 2319-2374, 1979.

Ramirez, J. P., S. Dudevec, E., J. T. Koutoukoulas, P. R. Molykides, E. A. Probability distribution and its use in fitting data. Technical Report 36, Statistical Techniques Research Group, Princeton University, 1980.