

# THE SCALING PROPERTIES IN THE DISTRIBUTION OF HYDROLOGICAL VARIABLES AS A RESULT OF THE MAXIMUM ENTROPY PRINCIPLE

## D. Koutsoyiannis, Department of Water Resources, National Technical University of Athens

### 1. Abstract

It is well known that the principle of maximum entropy (ME), when applied to the probability distribution of a random variable with known mean and variance, results in the normal distribution. If the variable is non-negative, as happens with hydrological variables such as rainfall and runoff, the same principle results in a truncated normal distribution. Mathematically, this distribution can have a coefficient of variation ranging from zero to unity, with the upper bound corresponding to the exponential distribution. At fine time scales, rainfall and runoff have coefficients of variations higher than one, so the classical entropy approach, constrained by known mean and variance, is not applicable. However, a generalization of entropy (specifically the use of the concept of nonextensive entropy) allows the application of the ME principle even in such cases and results in power-type distributions, which for low probabilities of exceedance have scaling properties. Thus, the ME principle can be used to infer the type of the distribution of the random variable, whether it has scaling properties or not, and to quantify the scaling exponent using a simple indicator such as the coefficient of variation. This theoretical framework is validated with several real world examples concerning rainfall, runoff and temperature data at several time scales. Given that entropy is a measure of uncertainty, the applicability of the ME principle to the distribution of hydrological variables emphasizes the dominance of uncertainty in hydrological processes.

### 2. Type-“which?” versus type-“why?” questions

The typical questions in hydrological statistics are type-“which?”: Which is the most appropriate theoretical distribution for a hydrological quantity such as rainfall and runoff? The selection is done from a repertoire that contains distributions such as normal, log-normal, Pearson, log-Pearson, Weibull, Extreme Value of maxima and minima of types I, II, and III (EV1, EV2, EV3) and is based on comparisons with empirical distributions.

Generally, hydrological statistics avoids questions of the type “why?”, which would provide an explanation of the appropriateness or inappropriateness of a certain distribution and thus would also help choose the most appropriate distribution.

The importance of type-“why?” questions becomes more evident when studying extreme events. Obviously, observations of extreme events cannot be as numerous as those of regular events. Therefore, a theoretical reasoning of the appropriateness of a statistical distribution, in addition to the empirical study of the data, would help for a more justified and correct choice.

### 6. The principle of maximum entropy (ME)

In a probabilistic context, the ME principle was introduced by Jaynes (1957) as a generalization of the “principle of insufficient reason” (PIR) attributed to Bernoulli or to Laplace.

The ME principle is used to infer unknown probabilities from known information. It is related to the homonymous physical principle that determines thermodynamical states.

It postulates that the entropy of a random variable should be at maximum, under some conditions, formulated as constraints, which incorporate the information that is given about this variable.

For a discrete  $x$  taking a finite number of possible states  $w$ , in the absence of any information, both  $\phi_q = 1/w$ , for any  $w > 0$  achieve their maximum values at equal probability ( $p_i = 1/w$ , corresponds to PIR).

For infinite number of possible states, constraints are necessary. Typical constraints used in a probabilistic or physical context are:

$$\int f(x) dx = 1, \quad E[X] = \int x f(x) dx = \mu$$

$$E[X^2] = \int x^2 f(x) dx = \sigma^2 + \mu^2, \quad x \geq 0$$

$$\text{Non-negativity}$$

$$\text{Variance/Energy}$$

$$f(x) = \exp(-\lambda_0 - \lambda_1 x - \lambda_2 x^2), \quad \phi = \lambda_0 + \lambda_1 \mu + \lambda_2 \mu^2$$

where  $\lambda_0, \lambda_1, \lambda_2$  are Lagrange multipliers. This tends to the exponential distribution as  $\sigma/\mu \rightarrow 1$ .

For  $\sigma/\mu > 1$  the Boltzmann-Gibbs-Shannon ME distribution does not exist. In this case the Tsallis entropy can be used which results in:

$$f(x) = [1 + \kappa (\lambda_0 + \lambda_1 x + \lambda_2 x^2)]^{-1/\kappa}, \quad \phi_q = (\kappa - 1)/(\lambda_0 + \lambda_1 \mu + \lambda_2 \mu^2)$$

where  $\kappa := (1 - \phi)/q$ . In the absence of any information for  $\kappa$  or  $q$ , we can set  $\lambda_2 = 0$  (Koutsoyiannis, 2005) and obtain the Pareto distribution:

$$f(x) = (1/\lambda_1)(1 + x/\lambda_1)^{-1/\kappa}$$

$$F^*(x) = (1 + \kappa x/\lambda_1)^{-1/\kappa}$$

For large  $x$ , this can be approximated by the scaling distribution in panel 3.

Maximized entropy and maximising distribution versus the coefficient of variation  $\sigma/\mu$

$$\phi := E[-\ln p(X)] = -\int f(x) \ln f(x) dx, \quad \text{where } \sum_{j=1}^w p_j \ln p_j = 1$$

For a continuous random variable  $X$  with probability density function  $f(x)$ , defined as

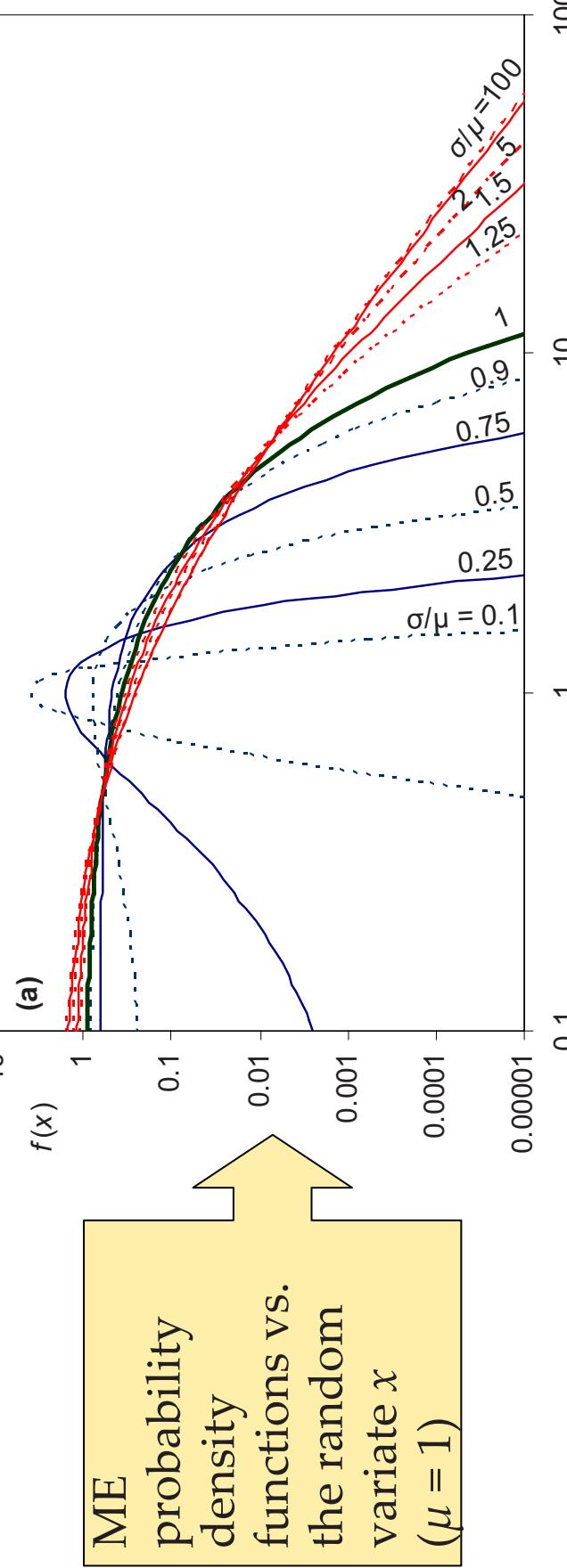
$$\phi := E[-\ln p(X)] = -\int f(x) \ln f(x) dx,$$

In both cases the entropy  $\phi$  is a measure of uncertainty about  $X$  and equals the information gained when  $X$  is observed (Papoulis, 1991).

In other disciplines (statistical mechanics, thermodynamics, fluid mechanics, entropy, is regarded as a measure of disorder and complexity.

### 8. Resulting ME distributions

The ME densities are truncated normal for  $\sigma/\mu < 1$  (tending to normal as  $\sigma/\mu \rightarrow 1$ ), exponential for  $\sigma/\mu = 1$  and Pareto for  $\sigma/\mu > 1$ .



### 7. Application of the ME principle to hydrology

We use the above four constraints, which involve two parameters, the mean  $\mu$  and the standard deviation  $\sigma$ , estimated from the sample.

The non-negativity constraint is essential for hydrological variables. In the case that a variable has a lower bound  $\neq 0$ , a shift is required to make it 0. Without loss of generality, before application of ME we can standardize the variable by its mean  $\mu$ , making it have mean 1 and standard deviation  $\sigma/\mu$ , which is the coefficient of variation (CV) of the original variable.

Thus, the ME distribution depends on a single parameter, the CV =  $\sigma/\mu$ . Maximization of the Boltzmann-Gibbs-Shannon entropy with these constraints results in the truncated (for  $x \geq 0$ ) normal distribution:

$$f(x) = \exp(-\lambda_0 - \lambda_1 x - \lambda_2 x^2), \quad \phi = \lambda_0 + \lambda_1 \mu + \lambda_2 \mu^2$$

where  $\lambda_0, \lambda_1, \lambda_2$  are Lagrange multipliers. This tends to the exponential distribution as  $\sigma/\mu \rightarrow 0$  and to the exponential distribution as  $\sigma/\mu \rightarrow 1$ .

For  $\sigma/\mu > 1$  the Boltzmann-Gibbs-Shannon ME distribution does not exist. In this case the Tsallis entropy can be used which results in:

$$f(x) = [1 + \kappa (\lambda_0 + \lambda_1 x + \lambda_2 x^2)]^{-1/\kappa}, \quad \phi_q = (\kappa - 1)/(\lambda_0 + \lambda_1 \mu + \lambda_2 \mu^2)$$

where  $\kappa := (1 - \phi)/q$ . In the absence of any information for  $\kappa$  or  $q$ , we can set  $\lambda_2 = 0$  (Koutsoyiannis, 2005) and obtain the Pareto distribution:

$$f(x) = (1/\lambda_1)(1 + x/\lambda_1)^{-1/\kappa}$$

For large  $x$ , this can be approximated by the scaling distribution in panel 3.

$$F^*(x) = (1 + \kappa x/\lambda_1)^{-1/\kappa}$$

Power-type distribution (Pareto)

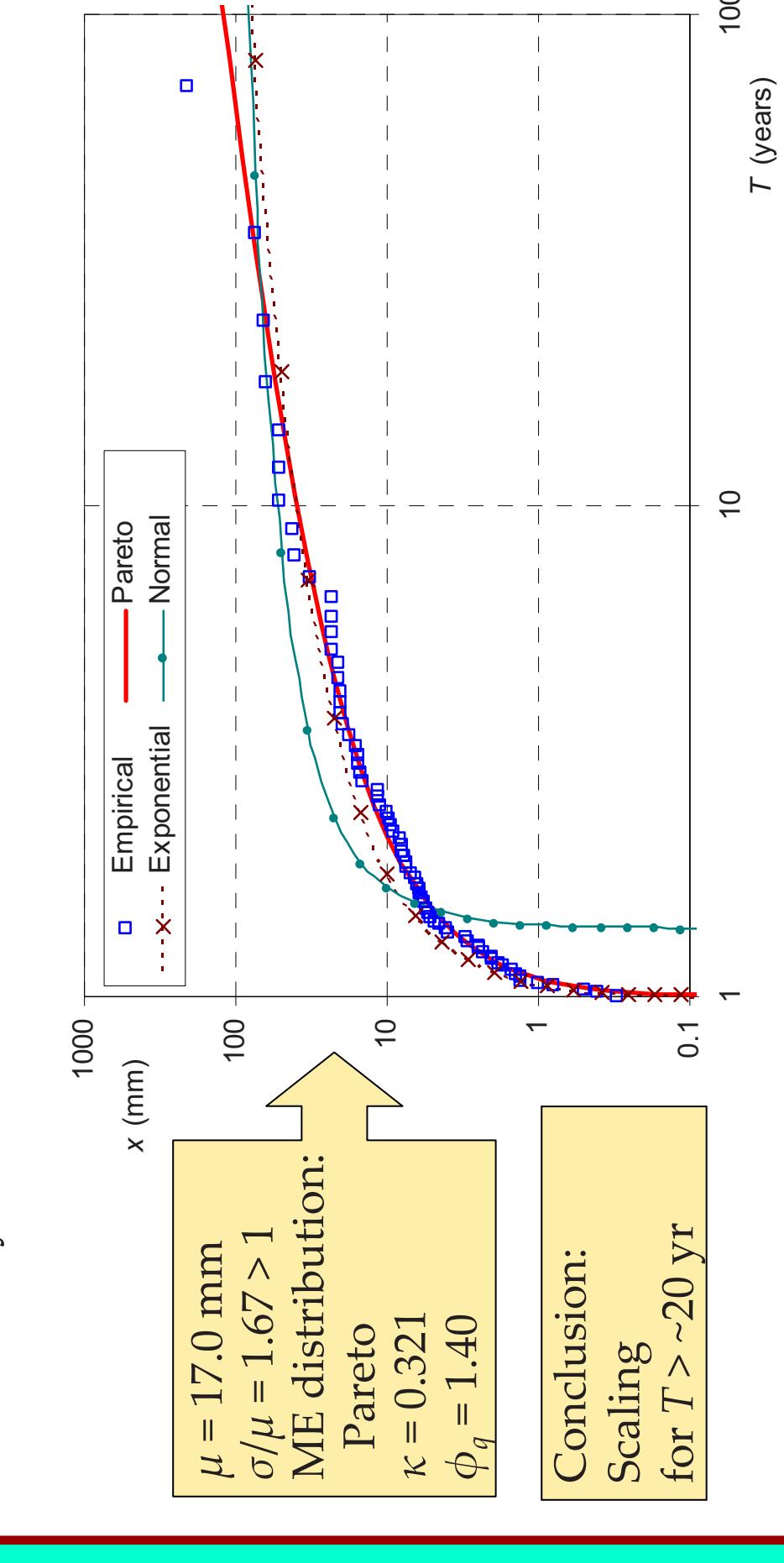
$$\text{Exponential-L shaped density}$$

$$\text{Truncated normal}$$

$$\text{Beta-shaped density}$$

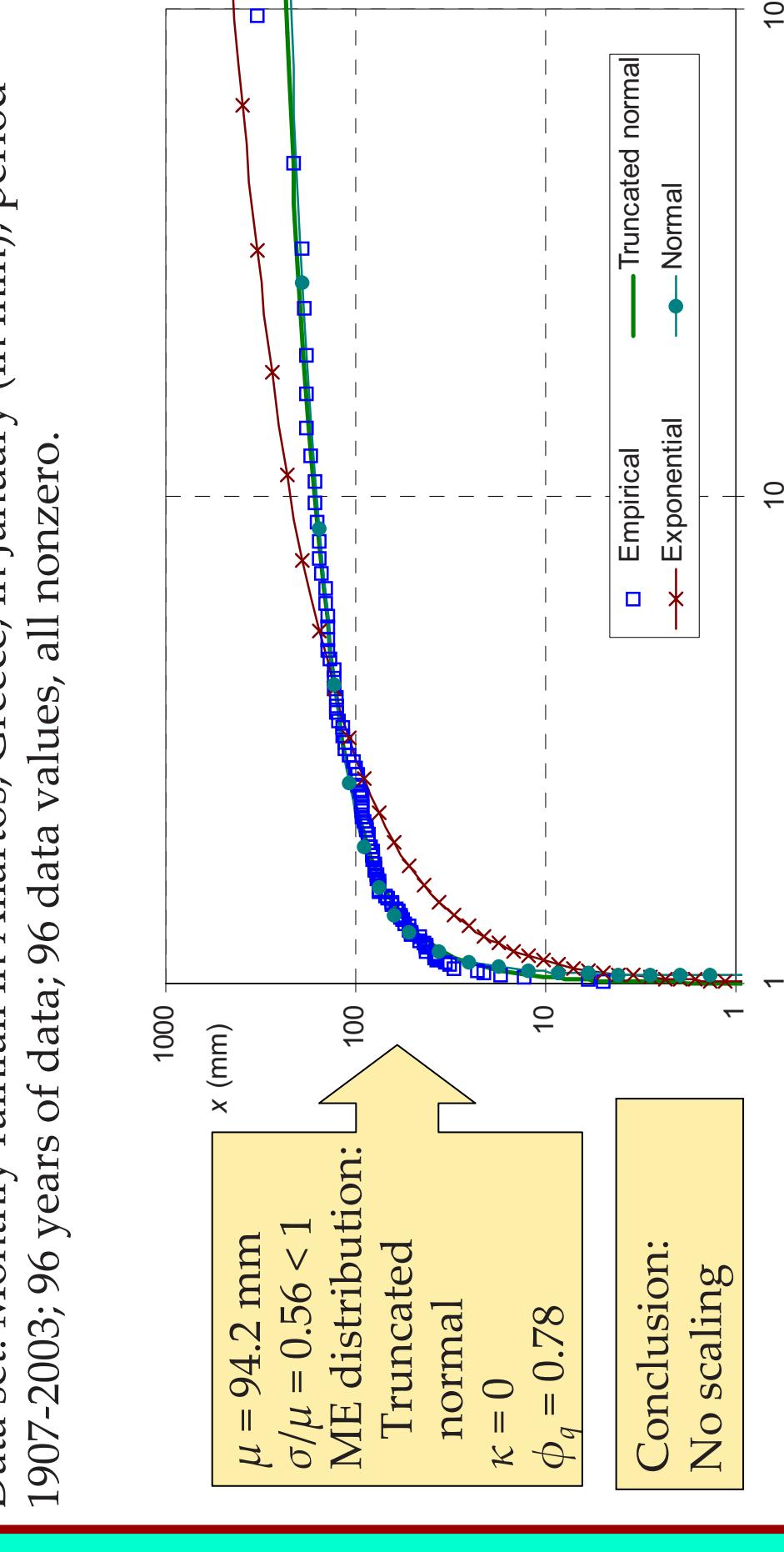
### 11. Application to monthly rainfall in a dry month

Data set: Monthly rainfall in Aliartos, Greece, in August (in mm); period 1907-2003; 96 years of data; 96 data values, 71 of which are nonzero.



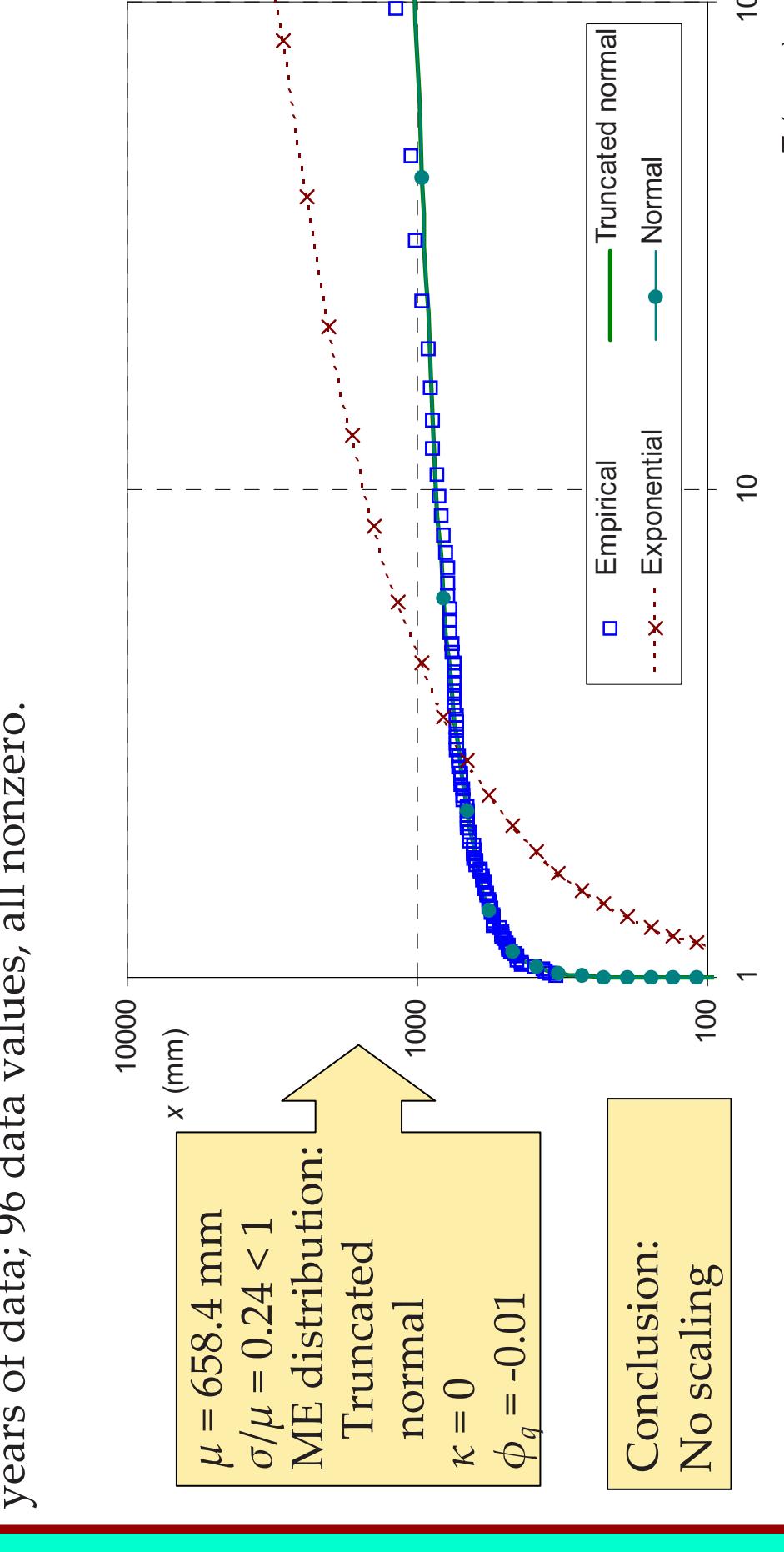
### 12. Application to monthly rainfall in a wet month

Data set: Monthly rainfall in Aliartos, Greece, in January (in mm); period 1907-2003; 96 years of data; 96 data values, all nonzero.



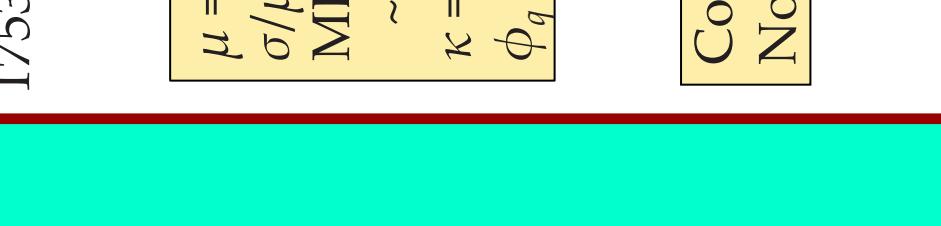
### 13. Application to annual rainfall

Data set: Annual rainfall in Aliartos, Greece (in mm); period 1907-2003; 96 years of data; 96 data values, all nonzero.



### 17. Application to mean annual temperature

Data set: Mean annual temperature in Geneva, Switzerland (in K); period 1753-1980 (one of the longest worldwide); 228 data values.



### 18. Conclusions

- Maximum entropy + Low variation  $\rightarrow$  Exponential-type (truncated normal) distribution.
- Maximum entropy + High variation  $\rightarrow$  Power-type (Pareto) distribution.
- Maximum entropy + High variation + High return periods  $\rightarrow$  State scaling.
- State scaling is only an approximation, good for high return periods and variables with high variation.
- Real world data series validate the applicability of the principle of maximum entropy in hydroeteorological processes.
- This can be interpreted as dominance of uncertainty in nature.

**References**

Jaynes, E.T. (1957) Information theory and statistical mechanics. Physical Review 106(4), 620-630.

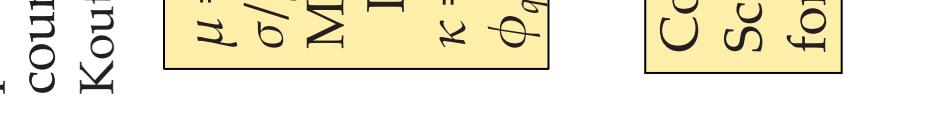
Koutsoyiannis, D. (2004) Statistics of extremes and estimation of long rainfall records. Hydrolog. Sci. J., 49(4), 591-610.

Koutsoyiannis, D. (2005) Uncertainty, entropy and hydrological processes. 1. Marginal distributional properties of hydrologic processes. A stochastic process approach. Random Variables and Stochastic Processes, (third ed.), 57 McGraw-Hill, New York, USA.

Tenmori, C. (1989) Revised generalization of Boltzmann's entropy. Statistical Science, 4(2), 175-179.

### 15. Application to daily runoff

Data set: Daily runoff in Boeotian Kephisos river basin, Greece (in mm); period 1978-2003; 5402 data values, out of which 6637 are nonzero. (Note: see details in Koutsoyiannis, 2005)



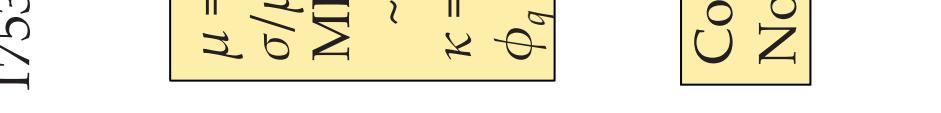
### 16. Application to mean daily temperature

Data set: Mean daily temperature in Athens, Greece, in three months (in K); period 1930-2003; 74 years of data; 2294/2294 data values for January/August/April. (Note: temperatures are expressed in K to have lower bound 0)



### 19. Application to hourly rainfall

Data set: Hourly rainfall in Athens, Greece, in January (in mm); period 1927-1996; 70 years of data (65 equivalent years if missing data are not counted); 48397 data values, out of which 2919 are nonzero. (Note: only nonzero values are modelled.)



### 14. Application to extreme daily rainfall worldwide

Data set: Daily rainfall from 168 stations worldwide (Koutsoyiannis 2004, 2005) each having at least 100 years of measurements; series above threshold, standardized by mean and unified; period 1822-2002; 17922 station-years of data.

