

ARTICLE TEMPLATE

Generalizing the normality: a novel towards different estimation methods for skewed information

Diego C Nascimento^a, Pedro Luiz Ramos^b, David Elal-Olivero^a, Milton Cortes-Araya^a, Francisco Louzada^b

^aDepartamento de Matemática, Facultad de Ingeniería, Universidad de Atacama, Copiapó 1530000, Chile

^bInstitute of Mathematical Science and Computing, University of São Paulo, São Carlos, Brazil

ARTICLE HISTORY

Compiled May 4, 2021

ABSTRACT

Normality is the most often mathematical supposition used in data modeling. Nonetheless, even based on the law of large numbers (LLN), normality is a strong presumption given that the presence of asymmetry and multi-modality in real-world problems is expected. Thus, a flexible modification in the Normal distribution proposed by Elal-Olivero [12] adds a skewness parameter, called Alpha-skew Normal (ASN) distribution, enabling bimodality and fat-tail, if needed, although sometimes not trivial to estimate this third parameter (regardless of the location and scale). This work analyzed seven different statistical inferential methods towards the ASN distribution on synthetic data and historical data of water flux from 21 rivers (channels) in the Atacama region. Moreover, the contribution of this paper is related to the probability estimation surrounding the rivers' flux level in Copiapó city neighborhood, the most important economic city of the third Chilean region, and known to be located in one of the driest areas on Earth, besides the North and the South Pole.

KEYWORDS

Alpha-skew Normal, Bimodal distribution, Asymmetry accommodation, Water monitoring

1. Introduction

We live in the Big Data Era, in which a high volume and variety of data characterization are often noticeable in a data lake, nonetheless despite its amount of observation, symmetry, and smooth-tail are not always observed. These characteristics are natural since we all live in a complex world, with nonlinear relations and outliers, describing extreme values more recurrent than easy statistical tools take into account. This new age requires flexible models and different reasoning based on data information.

An often question crossing some traditional departments worldwide, "Are Statistics methods getting old fashion?". Sir David Cox [6] explains that the focus is on the data relevance and quality, based on its coverage and representativeness, which gives confidence for the results, despite the amount of information (large volume) in the set which may hold some potentially biased estimates, with measurement errors. Efron & Hastie [11] discuss the relation across computer-related and statistical inference

as a mathematical logic system for guidance and correction, complemented by the large-scale prediction algorithms, suitable for this new century.

Therefore, complexity is intrinsic to massive data where high-dimension is often presented, and dynamic [18, 30]. Nonetheless, all the information contained in the acquired data can be extracted through an estimation method, i.e., in maximum likelihood estimation (MLE), and in a parametric version, it will be supported by a supposed distribution. Parametric approaches are easy to interpret patterns through parameters and enable association across variables, present a low computational cost, and be easier to implement in decision-making systems.

In many cases, the standard MLE may not return desirable results. Other estimation methods that return accurate estimates have been considered, such as estimators based on the least square function [32], the product of spacing [4, 28] or goodness-of-fit statistics [21]. There is no unique method that performs better for all models and it may depend on the selected parametric form [20, 27]. Thus, using an efficient estimation method jointly with a flexible parametric model that covers many data patterns is demanded, which sometimes accommodates asymmetry and multi-modality that may be contained in the data That is considered.

This is the case of meteorological data, which shows significant changes as well as a complex dynamic [2, 9]. These field demands an extra attention, on data-driven models, thus needed to incorporate spatial-time dependence [19], structural change [23], extreme value [10], but moreover in the parametric world a model supported by a probabilistic model (which deals with asymmetry and multi-modality [25, 29]). Thus, this paper was motivated by the study case of the water flux of 21 rivers (channels) from the surroundings of Copiapó city, placed in the Atacama Desert, as one of the planet's driest areas. Moreover, we are persuaded to exemplify evidence towards the probability density associated with these events' empirical distribution through seven different statistical inference approaches.

This paper is divided into four parts. Section 2 presents the motivation and details regarding the analyzed data. Section 3 provides a background about the adopted methodology towards statistical inference elements. Then, Sections 4 and 5 show the results related to the implemented methodology on synthetic data and the real-world data analysis. Finally, Section 6 discusses the conclusions based on the obtained results.

2. The Data

The adopted data set is related to Fluviometric records (average monthly flows), from the Atacama Desert region (third region of Chile), in the Copiapó city neighborhood. The historical period of these data are from the past ten years from Jan, 2011 to Dec, 2020 associated with 21 rivers (or stream channel), obtained from the Chilean government web-site called *Dirección General de Aguas (Información Oficial Hidrometeorológica y de Calidad de Aguas en Línea)*.

Historical events reveal the high periodicity of the low water flux of the region. However, cyclical events were also noticeable (such as two showers of rain, defrost of glaciers/snow in summer, amongst others), creating an expected multi-modality and large leptokurtosis.

Within the process information, a decision support system (DSS) sifts through and analyzes massive amounts of data, compiling comprehensive information that can be used to solve problems and in decision-making. Figure 1 summarizing chart flow of the knowledge discovery in databases (KDD), from the information retrieval (IR), decision

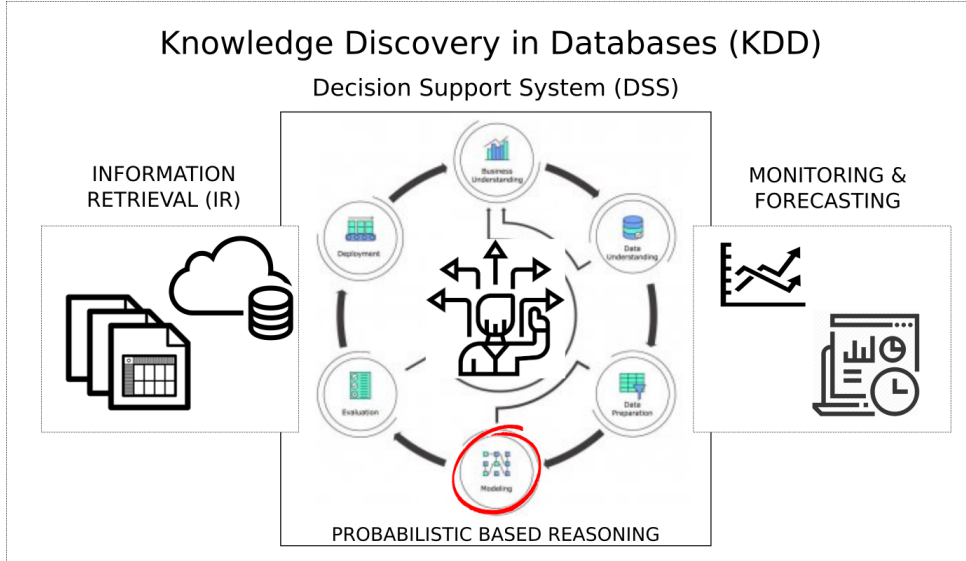


Figure 1. Visual summary of the probabilistic role in the Knowledge discovery in databases, as a cornerstone for quantification of uncertainty. Statistical inference procedures enable us to draw conclusions based on a sample, generalizing to the entire population.

support system (DSS), to monitoring & forecasting.

Thus, the Atacama Desert watershed problem is one of a multi-dimensional study related to the circular economy to be analyzed. It is essential to mention that uncertainty is always presented globally (as a measurement error, sample bias, amongst others). Nonetheless, probabilistic reasoning allows one to generalize results through statistical inference procedures.

3. Statistical Inference Elements

3.1. *Alpha-skew Normal (ASN) distribution*

Let X be a random variable following a Alpha-skew Normal (ASN) distribution then its probability density function (PDF) is given by

$$f(x|\alpha) = \frac{(1 - \alpha x)^2 + 1}{2 + \alpha^2} \phi(x)$$

where $x \in \mathbb{R}$, $\alpha \in \mathbb{R}$ and $\phi(\cdot)$ is the PDF of the standard normal distribution.

The cumulative density function (CDF) is given by

$$F(x|\alpha) = \Phi(x) + \alpha \left(\frac{2 - \alpha x}{2 + \alpha^2} \right) \phi(x)$$

Then, wrapping the ASN density $f(x|\alpha)$ with the parameters for location (μ) and scale (σ), that is, the random variable T is defined by $T = \mu + \sigma X$, for $\mu \in \mathbb{R}$ and $\sigma > 0$, given by:

$$f(t|\mu, \sigma, \alpha) = \frac{(1 - \alpha(t - \mu)\sigma^{-1})^2 + 1}{(2 + \alpha^2)\sigma} \phi\left(\frac{t - \mu}{\sigma}\right). \quad (1)$$

with CDF related to equation (1) as

$$F(t|\mu, \sigma, \alpha) = \Phi\left(\frac{t-\mu}{\sigma}\right) + \alpha\left(\frac{2\sigma - \alpha(t-\mu)}{(2+\alpha^2)\sigma}\right)\phi\left(\frac{t-\mu}{\sigma}\right). \quad (2)$$

Given its flexibility, the ASN distribution has been used in data modeling and adopted in different fields such as astronomy [34], modeling wind speed [39], and benchmark data [1]. Figure 2 presents different forms of the PDF of the ASN distribution, for instance, assuming $\mu = 0$ (location), $\sigma = 1$ (scale) and different values for α , showing the presence of asymmetry and bimodality (incorporating different heights between the modalities);

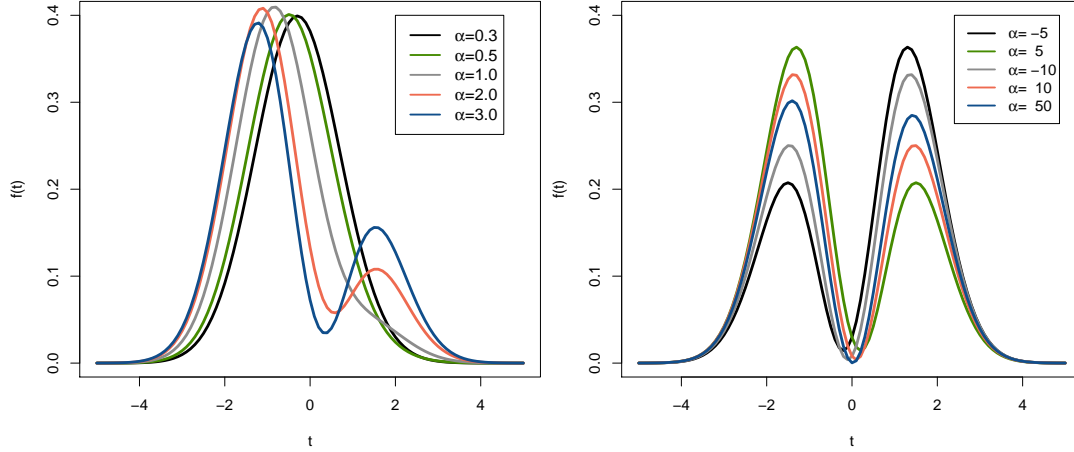


Figure 2. PDF $f(t)$ of the ASN distribution, where t is a random variable, assuming $\mu = 0$ (location), $\sigma = 1$ (scale) and different values for α .

3.2. Different estimation methods for ASN distribution

In this subsection, we will discuss seven different estimation methods (Maximum Likelihood Estimation, Ordinary and Weighted Least-Square Estimate, Method of Maximum Product of Spacings, Cramer-von Mises minimum distance estimators, Anderson-Darling and Right-tail Anderson-Darling estimators) for the parameters (μ , σ , and α) of the ASN distribution. Table 1 describe the methods used and the authors that proposed such inferential procedures. The comparison, which used the cited estimators, has been presented for other models [8, 26, 35].

Table 1. Summarizing seven inferential estimation methods.

Estimation Method	Abbreviation	Created by
Maximum Likelihood Estimation	MLE	Fisher [14]
Ordinary Least-Square Estimate	LSQ	Swain et al. [32]
Weighted Least-Square Estimate	WLQ	Swain et al. [32]
Maximum Product of Spacings	MPS	Cheng & Amin [4]
Cramer-von Mises Estimators	CME	Macdonald [22]
Anderson-Darling Estimator	ADE	Boos [3]
Right-tail Anderson-Darling Estimator	RADE	Luceno [21]

Note that while Carl Friedrich Gauss introduced the LSQ in 1822 and it is one

of the oldest estimation procedures, we have included the paper of Swain et al. [32]. The authors used such an approach for a class of non-normal models and became a standard reference when applied in different probability distributions. The additional details related to the estimation of the ASN distribution parameters are presented in the following subsections.

3.2.1. Maximum Likelihood Estimation

The Maximum Likelihood Estimation (MLE) is widely used in data analysis, where Fisher's derivation of the information inequality is seen at first for the analysis of variance, and later for estimate functions derived from Euler's Relation for homogeneous functions. Despite the fact that historical records of this technique have been widely exposed and defended by Ronald A. Fisher (maybe gained visibility because of the epic fights with Egon S. Pearson), its rationality dates back to the mid-1700s [31].

Let t_1, t_2, \dots, t_n be the sample of the random sample of size n from $F(\mathbf{t}|\mu, \sigma, \alpha)$. The maximum likelihood estimator $\hat{\mu}_{MLE}$, $\hat{\sigma}_{MLE}$ and $\hat{\alpha}_{MLE}$ can be obtained by maximizing, let's consider $z_i = \frac{(t_i - \mu)}{\sigma}$,

$$L(\mu, \sigma, \alpha) = \frac{1}{(2 + \alpha^2)^n \sigma^n} \prod_{i=1}^n ((1 - \alpha z_i)^2 + 1) \phi(z_i), \quad (3)$$

with respect to μ, σ and α . The log-likelihood function of (3) is given by

$$l(\mu, \sigma, \alpha) = \sum_{i=1}^n \log((1 - \alpha z_i)^2 + 1) - n \log(2 + \alpha^2) - n \log(\sigma) + \sum_{i=1}^n \log \phi(z_i). \quad (4)$$

From the expressions $\frac{\partial}{\partial \mu} l(\mu, \sigma, \alpha) = 0$, $\frac{\partial}{\partial \sigma} l(\mu, \sigma, \alpha) = 0$, $\frac{\partial}{\partial \alpha} l(\mu, \sigma, \alpha) = 0$, the likelihood equations are

$$\sum_{i=1}^n \frac{2\alpha(1 - \alpha z_i)}{(1 - \alpha z_i)^2 + 1} + \sum_{i=1}^n z_i = 0. \quad (5)$$

$$\sum_{i=1}^n \frac{2\alpha z_i(1 - \alpha z_i)}{(1 - \alpha z_i)^2 + 1} + \sum_{i=1}^n z_i^2 = 0. \quad (6)$$

$$\sum_{i=1}^n \frac{2z_i(1 - \alpha z_i)}{(1 - \alpha z_i)^2 + 1} + \frac{2n\alpha}{1 + \alpha^2} = 0. \quad (7)$$

Numerical methods such as Newton-Rapshon are required to find the solution of the nonlinear system. Under mild conditions, the MLEs are asymptotically normally distributed with a joint multivariate normal distribution given by

$$(\hat{\mu}_{MLE}, \hat{\sigma}_{MLE}, \hat{\alpha}_{MLE}) \sim N_3 [(\mu, \sigma, \alpha), I^{-1}(\mu, \sigma, \alpha)] \text{ as } n \rightarrow \infty.$$

where $I(\mu, \sigma, \alpha)$ is the Fisher information matrix given in Elal-Olivero [12].

This methodology is often adopted given the MLE sufficiency, combined with consistency (which leads to the statistical efficiency), and asymptotic normality is guaranteed. Often cases, MLE are adopted under regularity conditions and assumptions like score functions linearly approximation. Nonetheless, some more constraints on the class of estimates (whenever the increase of parameters' numbers, unbounded likelihood functions, and the possibility of local improvement) should be verified and needed to be contours of the likelihood.

Next, we will present a series of Minimum Distance Estimations (easily applied to estimate consistently unknown parameters) and designed to reflect the proposed model reproducing the probabilistic structure of the real-world phenomenon under study [38]. Minimum Distance Estimations provide consistent parameter estimates and competitive, especially when other methods did not succeed.

3.2.2. Ordinary and Weighted Least-Square Estimate

Let random sample of size n present a sequence $t_{(1)}, t_{(2)}, \dots, t_{(n)}$ then been a series of order statistics in which $F(\mathbf{t}|\mu, \sigma, \alpha)$ is a monotonic function. The least square (LSQ) estimators $\hat{\mu}_{LSE}$, $\hat{\sigma}_{LSE}$ and $\hat{\alpha}_{LSE}$ for the ASN distribution can be obtained by minimizing the parameters μ, σ and α , as

$$V(\mu, \sigma, \alpha) = \sum_{i=1}^n \left[F(t_{(i)} | \mu, \sigma, \alpha) - \frac{i}{n+1} \right]^2.$$

Thus, the LSQ equations can be obtained through solving the non-linear equations

$$\sum_{i=1}^n \left[F(t_{(i)} | \mu, \sigma, \alpha) - \frac{i}{n+1} \right] \Delta_j(t_{(i)} | \mu, \sigma, \alpha) = 0, \quad j = 1, 2, 3,$$

where

$$\begin{aligned} \Delta_1(t_{(i)} | \mu, \sigma, \alpha) &= \frac{\partial}{\partial \mu} F(t_{(i)} | \mu, \sigma, \alpha) = \frac{\phi(z_i)}{(2 + \alpha^2)\sigma} [2\alpha z_i - \alpha^2 z_i^2 - 2], \\ \Delta_2(t_{(i)} | \mu, \sigma, \alpha) &= \frac{\partial}{\partial \sigma} F(t_{(i)} | \mu, \sigma, \alpha) = \frac{z_i \phi(z_i)}{\sigma(2 + \alpha^2)} [2\alpha z_i - \alpha^2 z_i^2 - 2], \\ \Delta_3(t_{(i)} | \mu, \sigma, \alpha) &= \frac{\partial}{\partial \alpha} F(t_{(i)} | \mu, \sigma, \alpha) = \frac{\phi(z_i)}{(2 + \alpha^2)^2} [2 - 2\alpha^2 - 4\alpha z_i]. \end{aligned} \quad (8)$$

Alternative solutions are obtained through numerical approximation, with high precision, for these Δ_j for $j = 1, 2, 3$ partial derivatives.

Alternatively, the weighted least-squares (WLQ) estimates are proposed whenever efficient method is required under sets of small data, $\hat{\mu}_{WLSE}$, $\hat{\sigma}_{WLSE}$ and $\hat{\alpha}_{WLSE}$, can be obtained by minimized adopting the following equation,

$$W(\mu, \sigma, \alpha) = \sum_{i=1}^n \frac{(n+1)^2 (n+2)}{i(n-i+1)} \left[F(t_{(i)} | \mu, \sigma, \alpha) - \frac{i}{n+1} \right]^2.$$

The solutions are deviated from the non-linear equations

$$\sum_{i=1}^n \frac{(n+1)^2 (n+2)}{i(n-i+1)} \left[F(t_{(i)} | \mu, \sigma, \alpha) - \frac{i}{n+1} \right] \Delta_j(t_{(i)} | \mu, \sigma, \alpha) = 0, \quad j = 1, 2, 3,$$

where $\Delta_1(\cdot | \mu, \sigma, \alpha)$, $\Delta_2(\cdot | \mu, \sigma, \alpha)$ and $\Delta_3(\cdot | \mu, \sigma, \alpha)$ are given in (8). The WLQ estimation technique is particularly useful whenever one aims to weigh the observations proportional to the equivalence of the error variance for that observation, then overcoming the issue of non-constant variance.

3.2.3. Method of Maximum Product of Spacings

The maximum product of spacings (MPS) method is a powerful alternative to MLE for estimating unknown parameters of continuous univariate distributions, which aims to maximize the geometric mean of spacings in the data (differences between the values of the cumulative distribution function at neighborhood data points). Cheng & Amin proposed this method [4, 5], and also obtained independently by Ranney [28], as a Kullback-Leibler information approximation measurement. Some desirable properties of the MPS methods such as asymptotic efficiency, invariance, and more importantly, consistency are held broadly (under general conditions) than for MLEs [5].

Let's represented the differences between the values of the cumulative distribution function on their neighborhood data points by the function $D_i(\mu, \sigma, \alpha) = F(t_{(i)} | \mu, \sigma, \alpha) - F(t_{(i-1)} | \mu, \sigma, \alpha)$, for $i = \{1, 2, \dots, n+1, \dots\}$ as an uniform spacings of a random sample from the ASN distribution, defining by $F(t_{(0)} | \mu, \sigma, \alpha) = 0$ and $F(t_{(n+1)} | \mu, \sigma, \alpha) = 1$. The constraint of the $\sum_{i=1}^{n+1} D_i(\mu, \sigma, \alpha) = 1$ is held. Thus, the MPS estimates $\hat{\mu}_{MPS}$, $\hat{\sigma}_{MPS}$ and $\hat{\alpha}_{MPS}$ can be obtained by maximizing the geometric mean of the spacings

$$G_{ASN}(\mu, \sigma, \alpha) = \left[\prod_{i=1}^{n+1} D_i(\mu, \sigma, \alpha) \right]^{\frac{1}{n+1}} \quad (9)$$

considering the maximization of this (G_{ASN}) function by adopting its logarithm as

$$H_{ASN}(\mu, \sigma, \alpha) = \frac{1}{n+1} \sum_{i=1}^{n+1} \log D_i(\mu, \sigma, \alpha). \quad (10)$$

The estimates of the unknown parameters $\hat{\mu}_{MPS}$, $\hat{\sigma}_{MPS}$ and $\hat{\alpha}_{MPS}$ are obtained by solving the nonlinear equations

$$\frac{1}{n+1} \sum_{i=1}^{n+1} \frac{1}{D_i(\mu, \sigma, \alpha)} [\Delta_j(t_{(i)} | \mu, \sigma, \alpha) - \Delta_j(t_{(i-1)} | \mu, \sigma, \alpha)] = 0, \quad j = 1, 2, 3, \quad (11)$$

where $\Delta_1(\cdot | \mu, \sigma, \alpha)$, $\Delta_2(\cdot | \mu, \sigma, \alpha)$ and $\Delta_3(\cdot | \mu, \sigma, \alpha)$ are given respectively in (8).

It is important to mention that if $t_{(i+k)} = t_{(i+k-1)} = \dots = t_{(i)}$ then $D_{i+k}(\mu, \sigma, \alpha) = D_{i+k-1}(\mu, \sigma, \alpha) = \dots = D_i(\mu, \sigma, \alpha) = 0$. Therefore, the MPS estimators are sensitive to closely spaced observations, especially ties. When the ties are due to multiple observations, $D_i(\mu, \sigma, \alpha)$ should be replaced by the corresponding likelihood $f(t_{(i)}, \mu, \sigma, \alpha)$ since $t_{(i)} = t_{(i-1)}$.

Under mild conditions for the ASN distribution, the MPS estimators are asymptotically normally distributed with a joint trivariate normal distribution given by

$$(\hat{\mu}_{MPS}, \hat{\sigma}_{MPS}, \hat{\alpha}_{MPS}) \sim N_3 [(\mu, \sigma, \alpha), I^{-1}(\mu, \sigma, \alpha)] \text{ as } n \rightarrow \infty.$$

3.2.4. The Cramer-von Mises minimum distance estimators

Alternatively, an estimator that requires no assumption about the distributions' parametric form, the Cramer-von Mises estimator (CME), is based on the difference between the estimate of the cumulative distribution function and the empirical distribution function [7, 37]. These estimators operate based on the minimum distance across the "true" distribution (observed) and the "modeled" distribution (adjusted) through the maximum goodness-of-fit.

Macdonald [22] showed that the bias of the estimator, from the CME, presents smaller distances alternatively to other minimum distance estimators. The Cramer-von Mises estimates $\hat{\mu}_{CME}$, $\hat{\sigma}_{CME}$ and $\hat{\alpha}_{CME}$ of the parameters μ , σ and α are obtained by minimizing through

$$C(\mu, \sigma, \alpha) = \frac{1}{12n} + \sum_{i=1}^n \left(F(t_{(i)} | \mu, \sigma, \alpha) - \frac{2i-1}{2n} \right)^2. \quad (12)$$

Thus, these estimates are also obtained by solving the non-linear equations:

$$\sum_{i=1}^n \left(F(t_{(i)} | \mu, \sigma, \alpha) - \frac{2i-1}{2n} \right) \Delta_j(t_{(i)} | \mu, \sigma, \alpha) = 0, \quad j = 1, 2, 3,$$

where $\Delta_1(\cdot | \mu, \sigma, \alpha)$, $\Delta_2(\cdot | \mu, \sigma, \alpha)$ and $\Delta_3(\cdot | \mu, \sigma, \alpha)$ are given respectively in (8).

3.2.5. The Anderson-Darling and Right-tail Anderson-Darling estimators

Another type of minimum distance estimator is based on Anderson-Darling statistics, often called Anderson-Darling estimator (ADE). This estimator is based on the minimum distance estimator obtained from sampling a sort data in ascending order of observed set (Y), then $X = \text{Sort}(Y)$, and also combined with the permutation of $\{1, 2, \dots, n\}$ which makes the X series sorted. Thus, this process is associated with the cumulative distribution function $F(\cdot)$ and the survival function $S(\cdot) = 1 - F(\cdot)$ for any PDF. In contrast, samples are drawn from a uniform distribution only if Y (and X) are samples from the PDF distribution.

The Anderson-Darling estimates $\hat{\mu}_{ADE}$, $\hat{\sigma}_{ADE}$ and $\hat{\alpha}_{ADE}$ of the parameters μ , σ and α are obtained by minimizing, with respect to μ , σ and α , the function

$$A(\mu, \sigma, \alpha) = -n - \frac{1}{n} \sum_{i=1}^n (2i-1) (\log F(t_{(i)} | \mu, \sigma, \alpha) + \log S(t_{(n+1-i)} | \mu, \sigma, \alpha)). \quad (13)$$

These estimates can also be obtained by solving the non-linear equations

$$\sum_{i=1}^n (2i-1) \left[\frac{\Delta_j(t_{(i)} | \mu, \sigma, \alpha)}{F(t_{(i)} | \mu, \sigma, \alpha)} - \frac{\Delta_j(t_{(n+1-i)} | \mu, \sigma, \alpha)}{S(t_{(n+1-i)} | \mu, \sigma, \alpha)} \right] = 0, \quad j = 1, 2, 3.$$

Alternatively, one can improve the ADE performance by taking into account the information held on the non-symmetrical differences between theoretical CDF and empirical CDF [40]. Thus, the Right-tail Anderson-Darling estimates (RADE) is an alternative though $\hat{\mu}_{RADE}$, $\hat{\sigma}_{RADE}$ and $\hat{\alpha}_{RADE}$ of the parameters μ , σ and α are obtained by minimizing the function

$$R(\mu, \sigma, \alpha) = \frac{n}{2} - 2 \sum_{i=1}^n F(t_{i:n} | \mu, \sigma, \alpha) - \frac{1}{n} \sum_{i=1}^n (2i - 1) \log S(t_{n+1-i:n} | \mu, \sigma, \alpha). \quad (14)$$

These estimates can also be obtained by solving the non-linear equations:

$$-2 \sum_{i=1}^n \Delta_j(t_{i:n} | \mu, \sigma, \alpha) + \frac{1}{n} \sum_{i=1}^n (2i - 1) \frac{\Delta_j(t_{n+1-i:n} | \mu, \sigma, \alpha)}{S(t_{n+1-i:n} | \mu, \sigma, \alpha)} = 0, \quad j = 1, 2, 3.$$

where $\Delta_1(\cdot | \mu, \sigma, \alpha)$, $\Delta_2(\cdot | \mu, \sigma, \alpha)$ and $\Delta_3(\cdot | \mu, \sigma, \alpha)$ are given respectively in (8).

4. Numerical Analysis

In this section, we investigated the behavior of the ASN distribution based on artificial (synthetic) data, and its parameters modification conditioning on the estimation methodology. Thus, a Monte Carlo simulation was carried out, seven frequentist estimation methods were considered for the parameters, and comparing their efficiency. The following approach was adopted. The procedure is:

- (1) Generate N samples of size n given a set of parameters from the $ASN(\mu, \sigma, \alpha)$ distribution;
- (2) For each generated set, based on the estimation methods (MLEs, LSQs, WLQs, MPSs, CMEs, ADEs and RTADEs, estimates of the parameters (μ , σ and α) were calculated;
- (3) Then, considering $\hat{\boldsymbol{\theta}} = (\hat{\mu}, \hat{\sigma}, \hat{\alpha})$ and $\boldsymbol{\theta} = (\mu, \sigma, \alpha)$ it was computed the Bias and Mean Squared Error (MSE) of $\hat{\boldsymbol{\theta}}$, which are given, respectively, by $\frac{1}{N} \sum_{k=1}^N (\hat{\theta}_j^{(k)} - \theta_j)$ and $\frac{1}{N} \sum_{k=1}^N (\hat{\theta}_j^{(k)} - \theta_j)^2$, for $j = \{1, 2, 3\}$ (each parameter). Whereas, $\hat{\theta}_j^{(k)}$ denotes the estimate of θ_j obtained from sample k , for $k = 1, 2, \dots, N$.

This simulated study's results shall return the most expected efficient estimation method conditioning on their estimations both Bias and MSE closer to zero. For this simulation study, we adopted the R software (R Core Team 2014), and for the maximization method used package *maxLik* and *stats4* (Henningsen and Toomet, 2011). The chosen values of the simulation parameters were: $N = 10,000$ and $n = \{40, 60, 80, \dots, 300\}$. Due to lack of space, we will present the results only for $\{\mu = 0.5, \sigma = 0.5, \alpha = 3\}$ and $\{\mu = 0, \sigma = 1, \alpha = 5\}$. Nonetheless, the following results are generalized by other choices of the vector of parameters $\boldsymbol{\theta}$. The estimation methods are considered under the same conditions in terms of samples, limit iterations numbers, and initial values. Here, we considered the true values as initial values. However, we provide a simple approach discussed in the next section to deal with real cases where good initial values are not available.

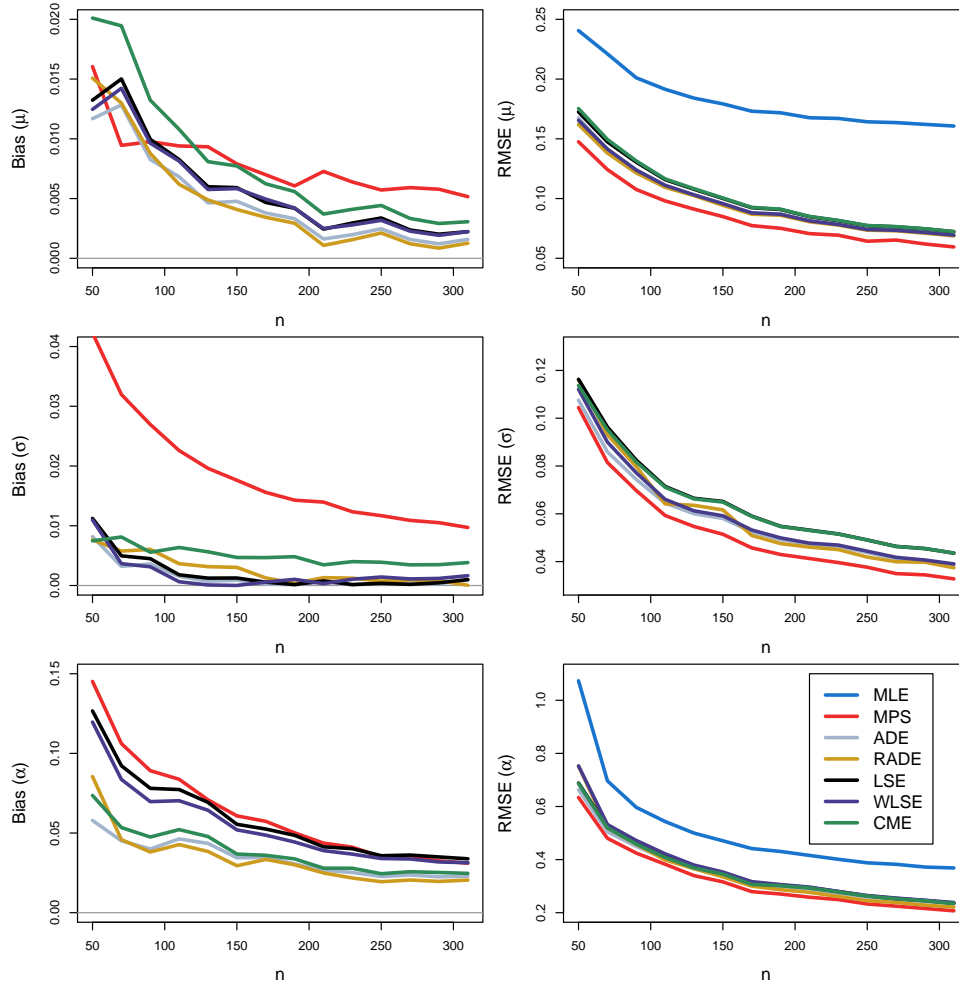


Figure 3. Bias and MSE of the estimates of $\mu = 0.5$, $\sigma = 0.5$ and $\alpha = 3$, for $N = 10,000$ simulated samples of size n , using the following methods: 1- MLE, 2- MPS, 3- ADE, 4- RADE, 5- LSE, 6- WLSE, 7 - CME.

Figures 3 and 4 present the performance of the estimators in terms of Bias and MSE for the parameters μ , σ , and α using the MLEs, LSQs, WLQs, MPSs, CMEs, ADEs and RTADEs, with $N = 10.000$ simulated samples, and different values of n . It can be observed that the MLEs do not return adequate estimates for some parameter values and only converges for large samples of size. These results show a drawback in the current approach used to obtain the parameter estimates of ASN distribution. Although there is no uniform method that returns better estimates for all parameters and different parameter values, we observed that the ADEs obtained the best results in terms of minimum Bias and MSE. Additionally, obtaining an estimate for α is quite challenging to estimate given the influence from the parameter of location and scale (μ and σ). In this approach, we faced fewer computational issues to obtain such estimates, and therefore we recommend using the ADEs to achieve the estimates for all practical purposes.

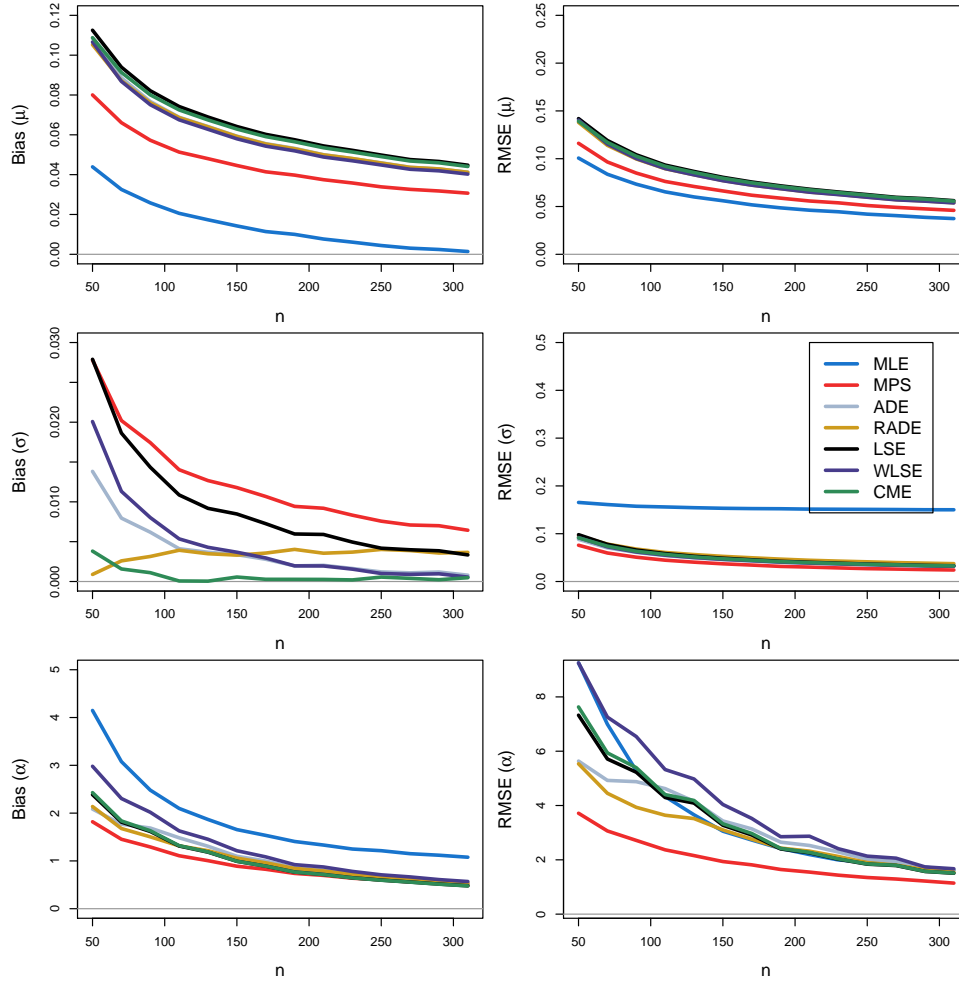


Figure 4. Bias and MSE of the estimates of $\mu = 0$, $\sigma = 1$ and $\alpha = 5$, for $N = 10,000$ simulated samples of size n , using the following methods: 1- MLE, 2- MPS, 3- ADE, 4- RADE, 5- LSE, 6- WLSE, 7 - CME.

5. Results

As we presented in Section 2, the motivation of this paper is driven by the Atacama Desert water flux. More precisely in the Copiapó neighborhood. Figure 5 shows the empirical density of this phenomenon, whereas a high concentration of low-values is presented (near to zero) although important events also captured in this 10-year time window, such as a big rain retrieving a large leptokurtosis.

Table 2 presents the statistical summaries (minimum, 1st Quartile, Median, Mean, 3rd Quartile, maximum) of the water flux per month. Since the weather is very constant in the region, the seasonality across years is to be ignored since low flux is common (close values through the minimum of the months). Nonetheless, the cycle per month is essential given events like defrost at the end of spring/beginning of summer (higher values in the 3rd quartile in NOV and DEC), it is expected to receive more water in the system.

The use of the logarithm transformation has been used for a long time [13], though often obtained a normal distribution, some other situations, not the case, for instance, the dynamic of the water flux through the observed period in Figure 7, showing the presence of bimodality in this data transformation, and its monthly representation by

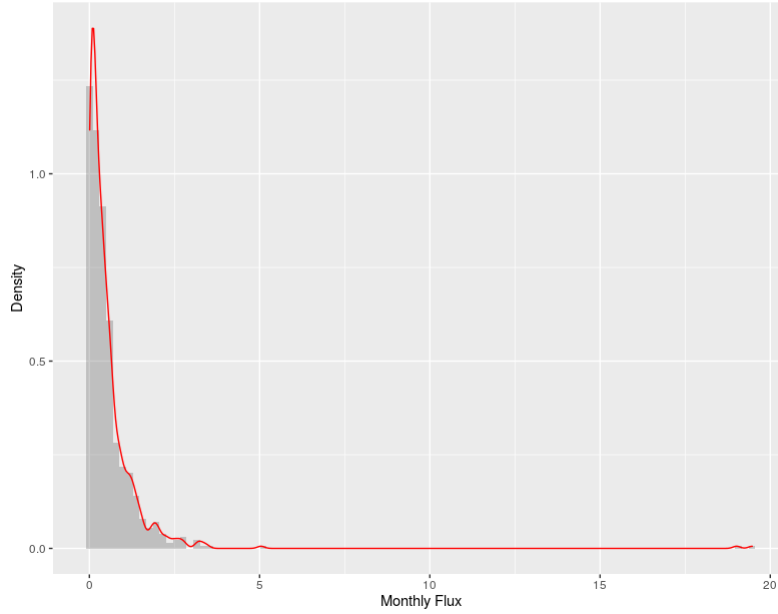


Figure 5. Empirical density function of the water flux from the 21 river/channels from the Copiapó neighborhood. The gray solid-shade represents the density (frequency) of each numerical records of the water flux, and the red solid-line a smooth adjusted function.

Figure 6. It is essential to mention that the maximum historical values were in MAY and JUN (of 2017), related to the heavy rains that occurred in the region, also notable in the previous Table 2.

The empirical distribution of this phenomenon is visualized in Figure 7, whereas the dashed-line represents the adjusted density functions, in black adopting the MLE and in red adopting the ADE. The initial values used to start the iteration procedures were obtained from

$$\tilde{\mu} = \sum_{i=1}^n \frac{x_i}{n} \quad \text{and} \quad \tilde{\sigma} = \sum_{i=1}^n \frac{(x_i - \tilde{\mu})^2}{n}.$$

Table 2. Summary Statistics of the Water Flux per month.

Month	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.	NA's
JAN	0.02	0.06	0.31	0.5374	0.68	3.45	39
FEB	0.01	0.065	0.2	0.5165	0.6875	3.15	40
MAR	0.01	0.06	0.31	0.5449	0.85	3.24	37
APR	0.03	0.08	0.27	0.4494	0.5275	2.25	36
MAY	0.03	0.12	0.29	0.7859	0.55	19.47	47
JUN	0.02	0.12	0.35	0.9106	0.62	19.01	51
JUL	0.01	0.14	0.46	0.5636	0.64	2.58	53
AUG	0.01	0.1125	0.33	0.4692	0.6175	2.23	50
SEP	0.01	0.2025	0.45	0.5356	0.6775	2.66	52
OCT	0.02	0.1	0.37	0.5229	0.77	2.46	49
NOV	0.01	0.0775	0.365	0.5536	0.855	3.36	48
DEC	0.01	0.055	0.265	0.5639	0.7975	5.04	46

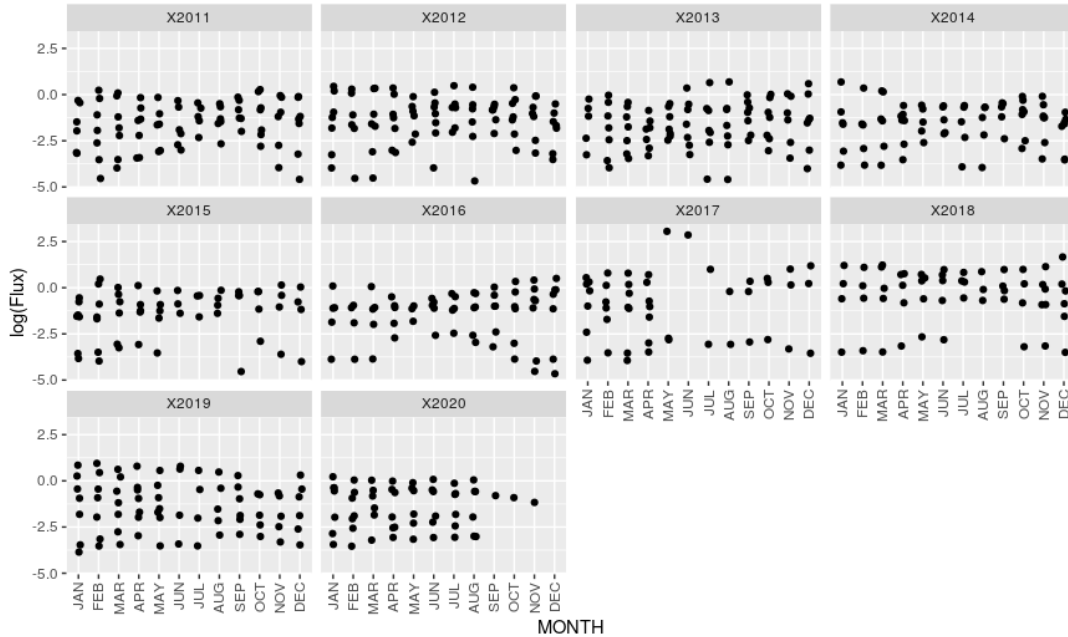


Figure 6. The logarithm of the water flux dispersion records (y-coordinates), per year (per panel), through the months (x-coordinates).

while $\tilde{\alpha}$ is obtained from a grid search considering the range $(-10, -9.5, \dots, 9, 10)$. The initial values of $\tilde{\mu}$ and $\tilde{\sigma}$ are clearly biased as they are the standar MLE for the normal distribution assuming that $\alpha = 1$, on the other hand, the obtained values are not so far from the true value of the ASN, therefore, they can be used as good initial values only. Additionally, the Kolmogorov-Smirnov test of the MLE showed an statistic test $D = 0.094$, presenting a p-value of 0.5, whereas the ADE had a $D = 0.12$, with the associated p-value of 0.2 (suggesting the adequacy of the methods).

After confirming the ASN distribution's good-of-fitness, event occurrence can be associated with its density (or cumulative) probability function. For instance, extreme values are to be seen as Table 3 shows some exemplifications taking into account the 1%, 10%, 50%, 99%, and 99.99%.

Table 3. Cumulative event probability based on the adjusted ASN distribution (using ADE).

CUM Prob.	1%	10%	50%	99%	99.99%
Flux	0.0059	0.0174	0.3396	1.5068	16.281

6. Conclusion

Uncertainty reveals a wide variety of processes and experiences, which may follow different rules, although different attributions of uncertainty such as external (disposition) versus internal (ignorance) are assessed by statistical inference, given philosophical interpretations of probability [17]. The utilities of each possible outcome lead to choosing rational actions regardless of the observed results' uncertainty.

The example brought by this paper is the water flux modeling related to an essential element stressed by the significant population growth and the increase in the demand

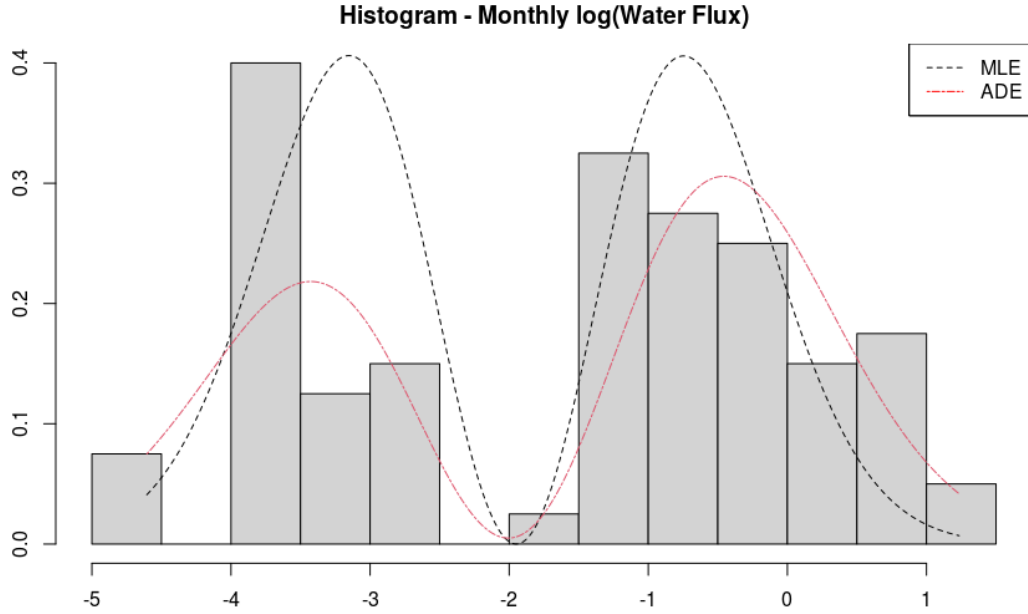


Figure 7. The empirical distribution of the log of the water flux, and its frequency represented by gray blocks. The black dashed line represents the adjusted ASN distribution based on the MLE ($\mu = -1.949, \sigma = 0.85, \alpha = 4147.07$), and the red dot-dashed line based on the ADE ($\mu = -1.879, \sigma = 1.05, \alpha = -8.36$).

for water supply (for agriculture, industrial process, mineral extraction, human consumption) [33]. Therefore, planning the logistic of these excessive water decisions, upon the probabilistic distribution, helps in unraveling this complex task [16, 36], initiated by the analysis of the water flux, especially when environmental factors present limit sources towards the water level.

Thus, this work proposed the investigation towards comparing different inference methods towards the ASN probabilistic distribution, which shows a promising and flexible distribution. Bimodality was noticeable and skewed information observed in the historical series, nevertheless accommodated by the adopted probabilistic approach.

Big-data solution may now be implemented, using real-time analysis, once the process's probabilistic function was estimated and evidence towards the good-of-fitness was shown. Monitoring charts and other statistical process control (SPC) tools can also explore since parametric distribution is often adopted, and here shown elements onto the ASN distribution. Future works should expand into the reasoning towards quantile estimations associated with this problem with explainable features (in a regression structure), and forecasting may also be a further research motivation later.

Disclosure statement

No potential conflict of interest was reported by the author(s)

Acknowledgments

Diego Nascimento acknowledges the support from the São Paulo State Research Foundation (FAPESP process 2020/09174-5). Pedro L. Ramos acknowledge the support

from the São Paulo State Research Foundation (FAPESP process 2017/25971-0). Francisco Louzada acknowledges support from the São Paulo State Research Foundation (FAPESP Processes 2013/07375-0) and CNPq (grant no. 301976/2017-1).

References

- [1] Ara, A. and F. Louzada (2019). The multivariate alpha skew gaussian distribution. *Bulletin of the Brazilian Mathematical Society, New Series* 50(4), 823–843.
- [2] Bonnail, E., R. C. Lima, and G. M. Turrieta (2018). Trapping fresh sea breeze in desert? health status of camanchaca, atacama’s fog. *Environmental Science and Pollution Research* 25(18), 18204–18212.
- [3] Boos, D. D. (1982). Minimum anderson-darling estimation. *Communications in Statistics-Theory and Methods* 11(24), 2747–2774.
- [4] Cheng, R. and N. Amin (1979). Maximum product of spacings estimation with application to the lognormal distribution. *Math. Report*, 79–1.
- [5] Cheng, R. and N. Amin (1983). Estimating parameters in continuous univariate distributions with a shifted origin. *Journal of the Royal Statistical Society. Series B (Methodological)*, 394–403.
- [6] Cox, D., C. Kartsonaki, and R. H. Keogh (2018). Big data: Some statistical issues. *Statistics & probability letters* 136, 111–115.
- [7] Cramér, H. (1928). On the composition of elementary errors: First paper: Mathematical deductions. *Scandinavian Actuarial Journal* 1928(1), 13–74.
- [8] Dey, S., D. Kumar, P. L. Ramos, and F. Louzada (2017). Exponentiated chen distribution: Properties and estimation. *Communications in Statistics - Simulation and Computation* 46(10), 8118–8139.
- [9] Du, H., L. V. Alexander, M. G. Donat, T. Lippmann, A. Srivastava, J. Salinger, A. Kruger, G. Choi, H. S. He, F. Fujibe, et al. (2019). Precipitation from persistent extremes is increasing in most regions and globally. *Geophysical Research Letters* 46(11), 6041–6049.
- [10] Dutfoy, A., S. Parey, and N. Roche (2014). Multivariate extreme value theory—a tutorial with applications to hydrology and meteorology. *Dependence Modeling* 2(1).
- [11] Efron, B. and T. Hastie (2016). *Computer age statistical inference*, Volume 5. Cambridge University Press.
- [12] Elal-Olivero, D. (2010). Alpha-skew-normal distribution. *Proyecciones (Antofagasta)* 29(3), 224–240.
- [13] Finney, D. (1941). On the distribution of a variate whose logarithm is normally distributed. *Supplement to the Journal of the Royal Statistical Society* 7(2), 155–161.
- [14] Fisher, R. A. (1922). On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society of London. Series A, Containing Papers of a Mathematical or Physical Character* 222(594-604), 309–368.
- [15] Henningsen, A. and O. Toomet (2011). maxlik: A package for maximum likelihood estimation in r. *Computational Statistics* 26(3), 443–458.
- [16] Jain, A. and L. E. Ormsbee (2002). Short-term water demand forecast modeling techniques—conventional methods versus ai. *Journal-American Water Works Association* 94(7), 64–72.
- [17] Kahneman, D. and A. Tversky (1982). Variants of uncertainty. *Cognition* 11(2), 143–157.
- [18] Leonelli, M., E. Riccomagno, and J. Q. Smith (2020). Coherent combination of probabilistic outputs for group decision making: an algebraic approach. *OR Spectrum* 42(2), 499–528.
- [19] Lopes, H. F., E. Salazar, D. Gamerman, et al. (2008). Spatial dynamic factor analysis. *Bayesian Analysis* 3(4), 759–792.
- [20] Louzada, F., P. Luiz Ramos, and P. Henrique Ferreira (2020). Exponential-poisson distribution: estimation and applications to rainfall and aircraft data with zero occurrence.

- Communications in Statistics-Simulation and Computation* 49(4), 1024–1043.
- [21] Luceño, A. (2006). Fitting the generalized pareto distribution to data using maximum goodness-of-fit estimators. *Computational Statistics & Data Analysis* 51(2), 904–917.
- [22] Macdonald, P. (1971). An estimation procedure for mixtures of distribution. *Journal of the Royal Statistical Society. Series B (Methodological)* 33, 326–329.
- [23] Mutti, P. R., P. S. Lúcio, V. Dubreuil, and B. G. Bezerra (2020). Ndvi time series stochastic models for the forecast of vegetation dynamics over desertification hotspots. *International Journal of Remote Sensing* 41(7), 2759–2788.
- [24] R Core Team (2014). *R: A Language and Environment for Statistical Computing. (Version 3.3.1)*. Vienna, Austria: R Foundation for Statistical Computing.
- [25] Ramos, P. and F. Louzada (2016). The generalized weighted lindley distribution: Properties, estimation and applications. *Cogent Mathematics* 3(1), 1256022.
- [26] Ramos, P. L., F. Louzada, T. K. Shimizu, and A. O. Luiz (2018). The inverse weighted lindley distribution: Properties, estimation and an application on a failure time data. *Communications in Statistics - Theory and Methods* 99, 1–20.
- [27] Ramos, P. L., D. C. Nascimento, P. H. Ferreira, K. T. Weber, T. E. Santos, and F. Louzada (2019). Modeling traumatic brain injury lifetime data: Improved estimators for the generalized gamma distribution under small samples. *PLoS one* 14(8), e0221332.
- [28] Ranneby, B. (1984). The maximum spacing method. an estimation method related to the maximum likelihood method. *Scandinavian Journal of Statistics* 11, 93–112.
- [29] Rodrigues, G. C., F. Louzada, and P. L. Ramos (2018). Poisson–exponential distribution: different methods of estimation. *Journal of Applied Statistics* 45(1), 128–144.
- [30] Smith, J. Q. (1987). *Decision analysis: a Bayesian approach*. Chapman & Hall, Ltd.
- [31] Stigler, S. M. et al. (2007). The epic story of maximum likelihood. *Statistical Science* 22(4), 598–620.
- [32] Swain, J. J., S. Venkatraman, and J. R. Wilson (1988). Least-squares estimation of distribution functions in johnson’s translation system. *Journal of Statistical Computation and Simulation* 29(4), 271–297.
- [33] Södergren, K. and J. Palm (2021). How organization models impact the governing of industrial symbiosis in public wastewater management. an explorative study in sweden. *Water* 13(6).
- [34] Tarnopolski, M. (2016). Analysis of gamma-ray burst duration distribution using mixtures of skewed distributions. *Monthly Notices of the Royal Astronomical Society* 458(2), 2024–2031.
- [35] Teimouri, M., S. M. Hoseini, and S. Nadarajah (2013). Comparison of estimation methods for the Weibull distribution. *Statistics* 47(1), 93–109.
- [36] Tu, Z., X. Gao, J. Xu, W. Sun, Y. Sun, and D. Su (2021). A novel method for regional short-term forecasting of water level. *Water* 13(6).
- [37] Von Mises, R. (1928). Statistik und wahrheit. *Julius Springer* 20.
- [38] Wolfowitz, J. (1957). The minimum distance method. *The Annals of Mathematical Statistics*, 75–88.
- [39] Yang, K. and M. Aziz (2018). Modeling wind speed distributions using skewed probability functions: A monte carlo simulation with applications to real wind speed data.
- [40] Ye, Y., G. Lu, Y. Li, and M. Jin (2017). Unilateral right-tail anderson-darling test based spectrum sensing for cognitive radio. *Electronics Letters* 53(18), 1256–1258.