



# Statistical power estimation dataset for external validation GoF tests on EVT distribution

Federico Reghenzani, Giuseppe Massari, Luca Santinelli, William Fornaciari

## ► To cite this version:

Federico Reghenzani, Giuseppe Massari, Luca Santinelli, William Fornaciari. Statistical power estimation dataset for external validation GoF tests on EVT distribution. *Data in Brief*, 2019, 25, pp.104071. <10.1016/j.dib.2019.104071>. <hal-02650476>

**HAL Id: hal-02650476**

**<https://hal.science/hal-02650476v1>**

Submitted on 29 May 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



HAL Authorization



## Data Article

## Statistical power estimation dataset for external validation GoF tests on EVT distribution



Federico Reghenzani <sup>a,\*</sup>, Giuseppe Massari <sup>a</sup>, Luca Santinelli <sup>b</sup>,  
William Fornaciari <sup>a</sup>

<sup>a</sup> DEIB, Politecnico di Milano, Milano, Italy

<sup>b</sup> DTIS, Onera, Toulouse, France

## ARTICLE INFO

## Article history:

Received 26 March 2019

Received in revised form 17 May 2019

Accepted 20 May 2019

Available online 28 May 2019

## ABSTRACT

This paper presents the statistical power estimation of goodness-of-fit tests for Extreme Value Theory (EVT) distributions. The presented dataset provides quantitative information on the statistical power, in order to enable the sample size selection in external validation scenario. In particular, high precision estimations of the statistical power of KS, AD, and MAD goodness-of-fit tests have been computed using a Monte Carlo approach. The full raw dataset resulting from this analysis has been published as reference for future studies: <https://doi.org/10.17632/hh2byrbmf.1>.

© 2019 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

## 1. Data

The dataset described in this paper provides an estimate of the statistical power of Goodness-Of-Fit (GoF) tests. Its analytical calculation is in fact usually not easy: for most of GoF tests a closed form expression does not even exist. This estimation is necessary to properly select the sample size for testing procedures, thus reducing the type-II errors, i.e. the inability to reject the null hypothesis when

\* Corresponding author.

E-mail address: [federico.reghenzani@polimi.it](mailto:federico.reghenzani@polimi.it) (F. Reghenzani).

Specification Table

Subject Area	Statistics
More specific subject area	Extreme Value Theory
Type of data	CSV text files, tables
How data was acquired	Monte Carlo approximation via CINECA supercomputing facility (Galileo cluster). The software has been developed in C++ over MPI/OpenMP frameworks.
Data format	Raw and aggregated
Experimental factors	The statistical testing procedures have been applied to synthetic time traces generated from known distributions. The process has been repeated several times collecting the test results.
Experimental features	The statistical tests results have been aggregated to obtain the statistical power estimation.
Data source location	CINECA, Segrate (MI), Italy
Data accessibility	Full raw dataset: <a href="https://doi.org/10.17632/hh2byrbmf.1">https://doi.org/10.17632/hh2byrbmf.1</a> Data in aggregated form are presented in this paper.

**Value of the data**

- The dataset described in this paper provides an estimate of the statistical power of Goodness-Of-Fit (GoF) tests performed on Extreme Value distributions.
- Several fields can benefit from the availability of this dataset, especially where it is necessary to select a proper sample size for the execution of GoF tests.
- The statistical power data have been computed on a *case 0* scenario (also called *external validation*), i.e. when the samples used to perform the test are a different set with respect to the samples used to estimate the reference distribution.
- The Monte Carlo approximation used to compute the statistical power has been performed on a very large number of sample ( $10^9$ ) to guarantee a high level of accuracy of the results.

it is actually false. The availability of this dataset can be advantageous for several fields, where the selection of the sample size is often performed with empirical procedures and where the results are often interpreted in a too optimistic view [1]. The GoF tests aim at identifying the deviation of data samples from a given distribution. However, if the test is not able to identify such null hypothesis violation, nothing can be stated and the statistical power becomes the only quantitative value that gives us the test result reliability information. The GoF tests have not been studied in *case 0* scenario (called also *external validation*) for EVT distributions, i.e. when the samples used to perform the test are a different set w.r.t. the samples used to estimate the reference distribution. In particular, to the best of our knowledge, quantitative information of only *case 3* scenarios is available in literature [2], while no *case 0* power analysis for such distribution classes is available in literature. This dataset wants exactly to fill this gap.

The statistical power computation has been performed with Monte Carlo approximations on a very large number of samples ( $10^9$ ). This guaranteed a high level of accuracy of the results. This, together with the external validation scenario, is an interesting feature for recent applications of the EVT. One of the possible use-case of this dataset is *probabilistic real-time computing* [3], where EVT is used to estimate the probabilistic Worst-Case Execution Time (WCET) of the computer tasks. In this scenario, the confidence level of the statistical test is critical. A false-negative result may indeed lead to an under-estimation of the WCET, which may be unacceptable for the production system [4]. This is the reason why we decided to build the statistical power dataset with the highest possible accuracy, enabling the selection of suitable sample size and ensuring a sufficient test result reliability [5].

1.1. Hypothesis testing and statistical power

In hypothesis testing, the null hypothesis ( $H_0$ ) is rejected when the observed data strongly suggest that it is false, in favour to an alternative hypothesis ( $H_1$ ). On the contrary, if the null hypothesis cannot be rejected, nothing can be inferred about the truthfulness of any hypothesis. The **statistical power** is defined as the probability to incur in a *Type II error*, i.e. the failure to reject the null hypothesis when it is

actually false. This concept can be expressed with the following conditional probability:  $P(\text{not reject } H_0 | H_0 \text{ is false})$ . This work presents the estimated statistical power of three Goodness-of-Fit (GoF) tests: Kolmogorov-Smirnov (KS) [6], Anderson-Darling (AD) [7], and Modified Anderson-Darling (MAD) [8] for EVT distributions. Other common tests have been excluded, for example the Chi-Squared (CS) and Cramer-von Mises (CvM) test, because state-of-the-art works already showed that they have a lower statistical power with respect to KS or AD [9,10].

Regarding the specific EVT case, the work of Heo et al. [2] estimated the AD and MAD test critical values and power, by using a Montecarlo approach for GoF test of EVT distributions. The critical values were computed for a scenario where the model parameters to be tested were estimated from the same data used for the test. This scenario is commonly referred to as *Case 3*, i.e., the assumed distribution parameters are unknown. The *a priori* knowledge of the distribution parameters (*Case 0*) in fact, is not usually available for most of classical EVT applications. However, in some cases, e.g. the probabilistic real-time computing previously mentioned, we can easily increase the sample size, because getting new samples requires a low effort. For this reason, the *Case 0* can be applied, by drawing different independent samples for model parameter estimation and for model validation. This enables the possibility to perform the *external validation* that leads, in general, to the most stringent and unbiased test [11].

Generally, statistical power estimations for *Case 0* are not representative of *Case 3* and vice versa. This makes the data provided with this paper extremely valuable, because they represent a highly accurate estimation of the GoF statistical power for the external validation scenario and EVT distributions.

## 1.2. Statistical power estimation

The EVT distributions can be grouped under the Generalized Extreme Value distribution:  $GEV(\mu, \sigma, \xi)$ , where  $\mu$  is the *location* parameter,  $\sigma$  is the *scale* parameter, and  $\xi$  is the *shape* parameter. The *location* and *scale* parameters determine the linear transformation of the standard GEV, while the *shape* parameter determines the distribution class. In this work, we explored all the three GEV classes as distribution references: a Gumbel distribution  $GEV(0, 1, 0)$ , a Weibull distribution  $GEV(0, 1, -0.5)$  and a Fréchet distribution  $GEV(0, 1, 0.5)$ . For each of these distributions, the Goodness-of-Fit tests have been run on samples drawn from the other two GEV and from: a normal  $N(0, 1)$ , a t-student  $t(10)$ , and a uniform distribution  $U(-2, 3)$ . The results for KS are shown in Table 1, for AD in Table 2, and for MAD in Table 3.

## 1.3. Sensitivity analysis

Given the statistical power results of the representative test cases, we performed a sensitivity analysis on the sample size and the shape parameter  $\xi$  of the GEV distribution. The results are depicted in Fig. 1, while the raw data are available in the dataset.

## 2. Experimental design, materials, and methods

The analytical computation of the statistical power, and consequently the selection of the appropriate sample size, is usually not possible, due to the frequent lack of the *effect size knowledge*, i.e. the real characterization of the population's distribution from which the samples have been collected. Consequently, Munthen et al. [12] studied the usage of Monte Carlo methods to select the sample size and determine the testing power. To this purpose, we need to define a set of tuples representing the test conditions. In particular, the Monte Carlo sampling is executed for every tuple  $(D, n, \alpha, \mathcal{S}_1, \mathcal{S}_2)$ , where  $D$  is the statistic of the test under analysis,  $n$  is the sample size,  $\alpha$  the level of significance,  $\mathcal{S}_1, \mathcal{S}_2$  are respectively the reference distribution with cumulative distribution function  $F(x)$  and the empirical distribution with cumulative distribution function  $F_n(x)$ .

The statistics  $D$  for KS, AD and MAD test can be computed using their discretized forms [13–15]:

**Table 1**  
Statistical powers of Kolmogorov-Smirnov (KS) test.

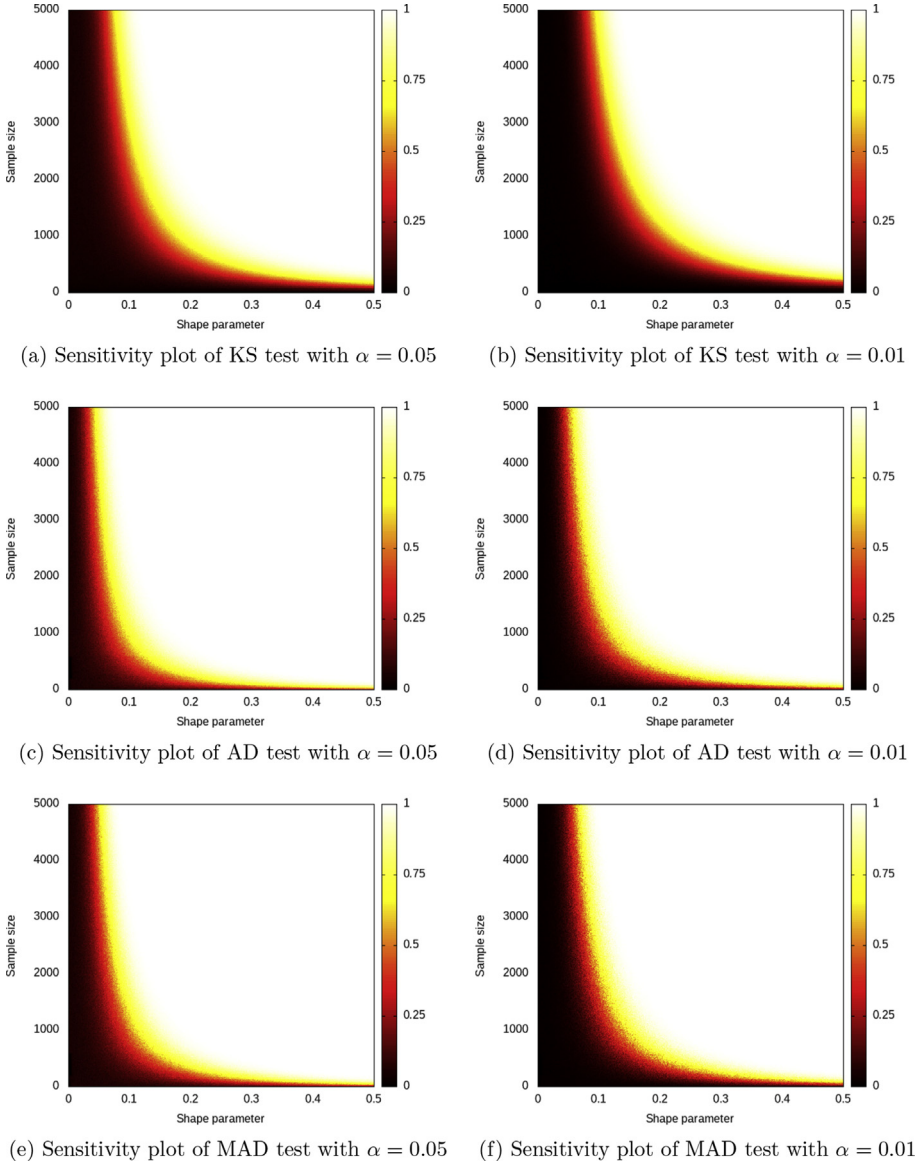
$G_0$	$G_1$	$\alpha$	Sample size ( $n$ )									
			50	100	150	200	300	400	500	750	1000	2500
GEV (0, 1, 0)	$N(0, 1)$	0.05	0.433100925	0.883347765	0.991603951	0.999615010	0.999999874	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.173221448	0.643452223	0.926426259	0.990668281	0.999969746	0.999999977	1.000000000	1.000000000	1.000000000	1.000000000
	$t(10)$	0.05	0.402221446	0.827773062	0.975704624	0.997320621	0.999988383	0.999999963	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.164124225	0.581214439	0.872639669	0.972233002	0.999500289	0.999996110	0.999999985	1.000000000	1.000000000	1.000000000
	$U(-2, 3)$	0.05	0.286787074	0.754349990	0.962845802	0.996346524	0.999991250	0.999999992	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.092007617	0.442246778	0.782790629	0.944778594	0.999292443	0.999997245	0.999999994	1.000000000	1.000000000	1.000000000
	GEV (0, 1, -0.5)	0.05	0.061924052	0.847865820	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.005506233	0.173621553	0.914148452	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	GEV (0, 1, 0.5)	0.05	0.109873147	0.293153608	0.549437525	0.740308878	0.943703112	0.991959149	0.999253784	0.999999586	1.000000000	1.000000000
		0.01	0.029933482	0.121708898	0.280459902	0.438755733	0.781235667	0.939418508	0.988237286	0.999932671	0.999999921	1.000000000
GEV (0, 1, 0.5)	$N(0, 1)$	0.05	0.869488165	0.999999315	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.454885837	0.998873700	0.999999998	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	$t(10)$	0.05	0.766801312	0.999765759	0.999999991	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.359259013	0.983640502	0.999988015	0.999999997	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	$U(-2, 3)$	0.05	0.367774702	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.139806008	0.744566527	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	GEV (0, 1, -0.5)	0.05	0.987657414	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.632639500	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	GEV (0, 1, 0)	0.05	0.032299787	0.231060818	0.576355557	0.852314165	0.995633829	0.999963685	0.999999914	1.000000000	1.000000000	1.000000000
		0.01	0.003184289	0.034069354	0.173820625	0.443183696	0.892633450	0.993613818	0.999882386	1.000000000	1.000000000	1.000000000
GEV (0, 1, -0.5)	$N(0, 1)$	0.05	0.284370451	0.629365918	0.862809163	0.953616572	0.996003172	0.999709102	0.999984209	0.999999990	1.000000000	1.000000000
		0.01	0.115952998	0.409541419	0.678458448	0.854765234	0.979383833	0.997758518	0.999791031	0.999999723	1.000000000	1.000000000
	$t(10)$	0.05	0.283091343	0.616658436	0.853660417	0.948713438	0.995402308	0.999673102	0.999984703	0.999999999	1.000000000	1.000000000
		0.01	0.116120515	0.399945876	0.664716864	0.844156028	0.976087763	0.997213016	0.999726448	0.999999688	0.999999999	1.000000000
	$U(-2, 3)$	0.05	0.826074656	0.998836102	0.999998678	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.564242941	0.981799855	0.999913887	0.999999780	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	GEV (0, 1, 0.5)	0.05	0.726868690	0.993646689	0.999951322	0.999999872	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
		0.01	0.466562747	0.954538684	0.999232900	0.999988813	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000	1.000000000
	GEV (0, 1, 0)	0.05	0.200916378	0.658118197	0.907265288	0.983782467	0.999851447	0.999999511	0.999999999	1.000000000	1.000000000	1.000000000
		0.01	0.062618618	0.341107470	0.726574304	0.899661698	0.997334290	0.999960032	0.999999740	1.000000000	1.000000000	1.000000000

F. Reghenzani et al. / Data in brief 25 (2019) 104071

5

Statistical powers of Modified Anderson-Darling (MAD) test.

F. Reghenzani et al. / Data in brief 25 (2019) 104077



**Fig. 1.** Sensitivity plots for  $\mathcal{G}_0 \sim \text{GEV}(0, 1, 0)$ ,  $\mathcal{G}_1 \sim \text{GEV}(0, 1, 0.5)$

$$D_{KS} = \sup_x \left| F_n(x) - F(x) \right|$$

$$D_{A^2} = -n - \frac{1}{n} \sum_{i=1}^n (2i-1) \log(F(x_i)) - \frac{1}{n} \sum_{i=1}^n (2n-2i+1) \log(F(1-x_i))$$



$$D_{AU^2} = \frac{n}{2} - 2 \sum_{i=1}^n F(x_i) - \sum_{i=1}^n \frac{2n - 2i + 1}{n} \log(F(1 - x_i))$$

The critical values (line 8) are computed with the following closed form – valid for  $n > 30$  – for KS test [16]:

$$\text{critical\_value}_{KS} = \frac{\sqrt{-\frac{1}{2} \log \frac{\alpha}{2}}}{\sqrt{n}}$$

Instead, for (M)AD test no closed form is available because the critical value computation procedure strongly depends on  $\mathcal{G}_0$ . We performed a dedicated Monte Carlo estimation similar to the method used by Heo et al. [2] to get (M)AD critical values. To double check, the resulting values have been used in the statistic comparison against data coming from  $\mathcal{G}_0$  (i.e. when  $H_0$  is true) and the tests failed to not reject  $H_0$  with  $\alpha$  probability, as expected by the definition of significance level.

**Algorithm 1.** Power estimation with Monte Carlo simulations.

---

**Input:**  $N$  (number iterations),  $D$  (test statistic),  $n$  (sample size),  $\alpha$  (significance level),  $\mathcal{G}_0, \mathcal{G}_1$  (null and alternative distributions)

**Output:**  $P_w$  (test power)

```

1 reject = not_reject = 0;
2 critical_value = get_critical_value( $D, \mathcal{G}_0, n, \alpha$ );
3 for  $i \in [1; N]$  do
4    $X = \text{collect\_sample}(\mathcal{G}_1, n)$ ;
5   if  $D(F_{\mathcal{G}_0}(\cdot), X) > \text{critical\_value}$  then
6     reject++;
7   else
8     not_reject++;
9   end
10 end
11  $P_w = \text{reject} / (\text{reject} + \text{not\_reject})$ ;

```

---

The estimation algorithm is shown in Listing 2. For each scenario, the critical value is computed (line 2) and a large number of explorations  $N$  is performed (lines 3–10). Each time, we draw a sample from the reference distribution (line 4) and we check if the statistic  $D$  of the ecdf matches or not with the drawn sample, comparing it with the critical value (line 5). If the statistic value is higher than the critical value, then the sample is rejected (line 6), otherwise not (line 8). Finally, the ratio rejection over total samples provide us the statistical power (line 11). If the test is able to detect the differences between  $\mathcal{G}_1$  and  $\mathcal{G}_2$  we expect to get a value near 1 for this ratio. In this specific Monte Carlo simulation, the standard error of *power* can be computed as [17]:

$$\sqrt{\frac{R(N - R)}{N^3}} \quad (1)$$

where  $R \leq N$  is the number of rejects (the accumulation variable of line 12). The standard error is decreasing when  $N \rightarrow \infty$  and when  $R \rightarrow N$ , i.e. when statistical power approaches the maximum value 1.

The selected values for parameters of each Monte Carlo estimation are:

- $N = 10^9$ : number of Monte Carlo iterations;
- $D$ : the test statistics previously described;

- $n$ : the sample size. Exploring all the possible values would have increased in a non-sustainable way the computational effort required by the Monte Carlo simulations. Since the power test function is a non-decreasing function of  $n$ , we explored them easily selecting the following values:  $n = (50, 100, 150, 200, 300, 400, 500, 750, 1000, 2500)$ ;
- $\alpha$ : the significance level. We studied the traditional values of 0.05 and 0.01.

The simulations ran on 4 nodes of CINECA supercomputing facility (GALILEO-A1 cluster, 2 x Intel Xeon E5-2697v4@2.3GHz per node) for a total of 144 CPU cores. It took  $\approx 13$ h for KS tests,  $\approx 17.5$ h for AD test,  $\approx 16$ h for MAD test.

Given the statistical power results of the representative test cases, we performed a sensitivity analysis on sample size  $n$  and shape parameter  $\xi$ . The power was obtained by using the same procedure of Algorithm 2, but reducing considerably the number of iterations  $N$ , in order to enable a fine-grain analysis with a sustainable computational effort. By exploring the integer sample size space and the real shape parameter space, the Monte Carlo simulations carry out a power matrix of sizes  $\bar{\xi} \times \bar{n}$  (where  $\bar{\cdot}$  is the cardinality of the set of all the possible values of  $\cdot$ ).

## Acknowledgment

This research was partially funded by EU project RECIPE H2020 (grant no. 801137 [18]). We thank CINECA supercomputing facility for the availability of high performance computing resources and support.

## Transparency document

Transparency document associated with this article can be found in the online version at <https://doi.org/10.1016/j.dib.2019.104071>.

## References

- [1] R.L. Lieber, Statistical significance and statistical power in hypothesis testing, *J. Orthop. Res.* 8 (1990) 304–309.
- [2] J.-H. Heo, H. Shin, W. Nam, J. Om, C. Jeong, Approximation of modified anderson–darling test statistics for extreme value distributions with unknown shape parameter, *J. Hydrol* 499 (2013) 41–49.
- [3] L. Santinelli, F. Guet, J. Morio, Revising measurement-based probabilistic timing analysis, in: *IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS)*, IEEE, 2017, pp. 199–208.
- [4] F. Reghenzani, G. Massari, W. Fornaciari, The misconception of exponential tail upper-bounding in probabilistic real-time, *IEEE Embedded Systems Letters* (2018), <https://doi.org/10.1109/LES.2018.2889114>, 1–1.
- [5] F. Reghenzani, G. Massari, W. Fornaciari, A. Galimberti, Probabilistic-WCET reliability: on the experimental validation of EVT hypotheses, in: *Proceedings of the 1st International Conference on Omni-Layer Intelligent Systems, COINS '19*, ACM, 2019.
- [6] F.J. Massey Jr., The Kolmogorov-smirnov test for goodness of fit, *J. Am. Stat. Assoc.* 46 (1951) 68–78.
- [7] T.W. Anderson, D.A. Darling, Asymptotic theory of certain goodness of fit criteria based on stochastic processes, *Ann. Math. Stat.* 23 (1952) 193–212.
- [8] M.I. Ahmad, C.D. Sinclair, B.D. Spurr, Assessment of flood frequency models using empirical distribution function statistics, *Water Resour. Res.* 24 (1988) 1323–1328.
- [9] A. Zemplén, Goodness-of-fit Test in Extreme Value Applications, Technical Report, Ludwig-Maximilians-Universität München, 2004.
- [10] R. Alpin, L. Fattorini, Empirical performance of some goodness-of-fit tests for the weibull and type i extreme value distributions, *Statistica Applicata* 5 (1993).
- [11] R.A. Giancristofaro, L. Salmaso, Model performance analysis and model validation in logistic regression, *Statistica* 63 (2007) 375–396.
- [12] L.K. Muthn, B.O. Muthn, How to use a Monte Carlo study to decide on sample size and determine power, *Struct. Equ. Model.: A Multidisciplinary J.* 9 (2002) 599–620.
- [13] F.J. M Jr., The Kolmogorov-smirnov test for goodness of fit, *J. Am. Stat. Assoc.* 46 (1951) 68–78.
- [14] T.W. Anderson, D.A. Darling, A test of goodness of fit, *J. Am. Stat. Assoc.* 49 (1954) 765–769.
- [15] C. Sinclair, B. Spurr, M. Ahmad, Modified anderson darling test, *Commun. Stat. Theor. Methods* 19 (1990) 3677–3686.
- [16] L. Sachs, *Angewandte Statistik*, Springer-Verlag Berlin Heidelberg (1997).
- [17] Z. Zhang, Monte Carlo based statistical power analysis for mediation models: methods and software, *Behav. Res. Methods* 46 (2014) 1184–1198.

- [18] W. Fornaciari, G. Agosta, D. Atienza, C. Brandolese, L. Cammoun, L. Cremona, A. Cilaro, A. Farres, J. Flich, C. Hernandez, M. Kulchewski, S. Libutti, J.M. Martínez, G. Massari, A. Oleksiak, A. Pupykina, F. Reghenzani, R. Tornero, M. Zanella, M. Zapater, D. Zoni, Reliable power and time-constraints-aware predictive management of heterogeneous exascale systems, in: *Proceedings of the 18th International Conference on Embedded Computer Systems: Architectures, Modeling, and Simulation, SAMOS '18*, ACM, New York, NY, USA, 2018, pp. 187–194.