

# Measuring Extremal Clustering in Time Series <sup>†</sup>

Marta Ferreira

Centro de Matemática, Universidade do Minho, 4710-057 Braga, Portugal; msferreira@math.uminho.pt

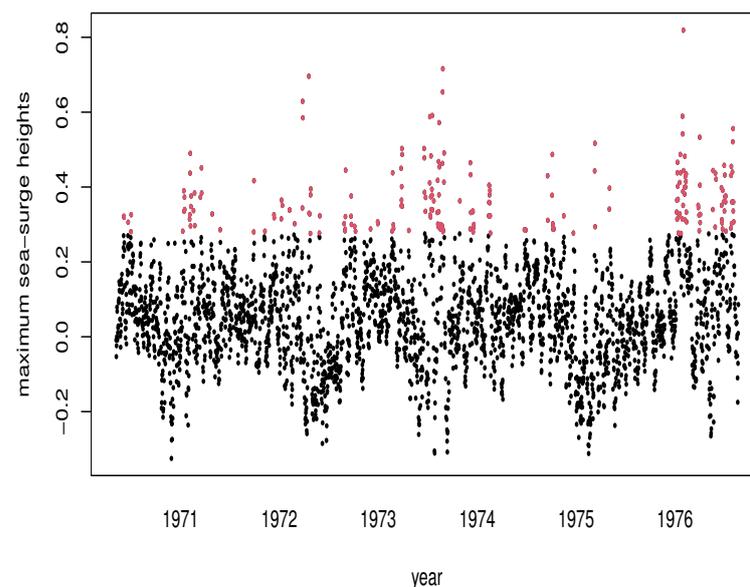
<sup>†</sup> Presented at the 9th International Conference on Time Series and Forecasting, Gran Canaria, Spain, 12–14 July 2023.

**Abstract:** The propensity of data to cluster at extreme values is important for risk assessment. For example, heavy rain over time leads to catastrophic floods. The extremal index is a measure of Extreme Values Theory that allows measurement of the degree of high-value clustering in a time series. Inference about the extremal index requires a prior choice of values for tuning parameters, which impacts the efficiency of existing estimators. In this work, we propose an algorithm that avoids these constraints. Performance is evaluated based on simulations. We also illustrate with real data.

**Keywords:** extreme values theory; stationary sequences; extremal index

## 1. Introduction

The occurrence of extreme values can lead to risky situations. Climate change, the global economic and financial crisis resulting from the COVID-19 pandemic situation, and the war in Ukraine have contributed to continuously growing attention from analysts, namely, to the risk of extreme phenomena. The duration of extreme values in time means the generation of clusters, the extension of which can aggravate the phenomenon. Extreme Values Theory (EVT) presents a set of adequate tools in this context. The extremal index is a measure of serial dependence assessing the propensity of data for the occurrence of clusters of extreme values. Figure 1 shows the maximum of sea-surge heights, where clusters of high values are visible.



**Figure 1.** Maximum hourly sea-surge heights (over contiguous 15-h time periods) in years 1971–1976 at the Newlyn Coast, Cornwall, UK.

More precisely, considering  $\mathbf{X} = \{X_n\}_{n \geq 1}$  as a stationary sequence of random variables (r.v.) with a common marginal distribution function (d.f.)  $F$  and denoting  $M_n =$



**Citation:** Ferreira, M. Measuring Extremal Clustering in Time Series. *Eng. Proc.* **2023**, *39*, 64. <https://doi.org/10.3390/engproc2023039064>

Academic Editors: Ignacio Rojas, Hector Pomares, Luis Javier Herrera, Fernando Rojas and Olga Valenzuela

Published: 6 July 2023



**Copyright:** © 2023 by the author. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

$\max(X_1, \dots, X_n)$ , then  $\mathbf{X}$  has extremal index  $\theta \in (0, 1]$  if for each real  $\tau > 0$  there exists a sequence of normalized levels  $u_n$ , i.e., satisfying  $n(1 - F(u_n)) \rightarrow \tau$ , as  $n \rightarrow \infty$ , such that  $P(M_n \leq u_n) \rightarrow \exp(-\theta\tau)$ . In the independent and identically distributed (i.i.d.) case, we have  $P(M_n \leq u_n) \rightarrow \exp(-\tau)$  and thus  $\theta = 1$ . On the other hand, if  $\theta = 1$ , then the tail behavior of  $\mathbf{X}$  resembles an i.i.d. sequence. Clustering of extreme values takes place whenever  $\theta < 1$ , and the smaller the  $\theta$  is, the larger is the propensity for clusters to appear. Under some dependence conditions,  $\theta$  is stated as the arithmetic inverse of the mean cluster size (Hsing et al. [1] 1988).

Assuming  $F$  is continuous, we have  $U_i = F(X_i)$ ,  $i = 1, \dots, n$  standard uniform r.v. and  $P(-n \log(F(M_n)) \geq \tau) \approx P(n(1 - F(M_n)) \geq \tau) = P(M_n \leq u_n) \rightarrow \exp(-\theta\tau)$ , with  $F(M_n) = \max(U_1, \dots, U_n)$ . Thus,  $Y_n = -n \log(F(M_n))$  and  $Z_n = n(1 - F(M_n))$  follow asymptotically an exponential distribution with parameter  $\theta$ . The maximum likelihood estimator was considered by Northrop ([2] 2015) based on  $Y_n$ . More precisely, dividing the time series  $X_1, \dots, X_n$  into  $k_n$  blocks of length  $b_n$ , with  $n = b_n k_n$ , and considering  $M_{ni} = M_{((i-1)b_n+1):(ib_n)} = \max(X_{(i-1)b_n+1}, \dots, X_{ib_n})$ ,  $i = 1, \dots, k_n$ , the maximum of the  $i$ -th block in the disjoint blocks case, and  $M_{ni} = M_{((i-1):(i+b_n-1))} = \max(X_{i-1}, \dots, X_{i+b_n-1})$ ,  $i = 1, \dots, n - b_n + 1$ , the maximum of the  $i$ -th block in the sliding blocks case, the Northrop estimator is given by

$$\tilde{\theta}^N = \left( \frac{1}{t_n} \sum_{i=1}^{t_n} \hat{Y}_{ni} \right)^{-1}, \tag{1}$$

where  $\hat{Y}_{ni} = -b_n \log(\hat{F}(M_{ni}))$  and  $\hat{F}$  denotes the empirical d.f. estimating the usually unknown  $F$ , with  $t_n = k_n$  or  $t_n = n - b_n + 1$  depending on whether we are using disjoint or sliding blocks, respectively. Berghaus and Bücher ([3] 2018) considered

$$\tilde{\theta}^B = \left( \frac{1}{t_n} \sum_{i=1}^{t_n} \hat{Z}_{ni} \right)^{-1}, \tag{2}$$

with  $Z_{ni} = b_n(1 - \hat{F}(M_{ni}))$ , a more amenable formulation to derive the asymptotic properties. Here, we consider the Berghaus and Bücher estimator with bias adjustment given by

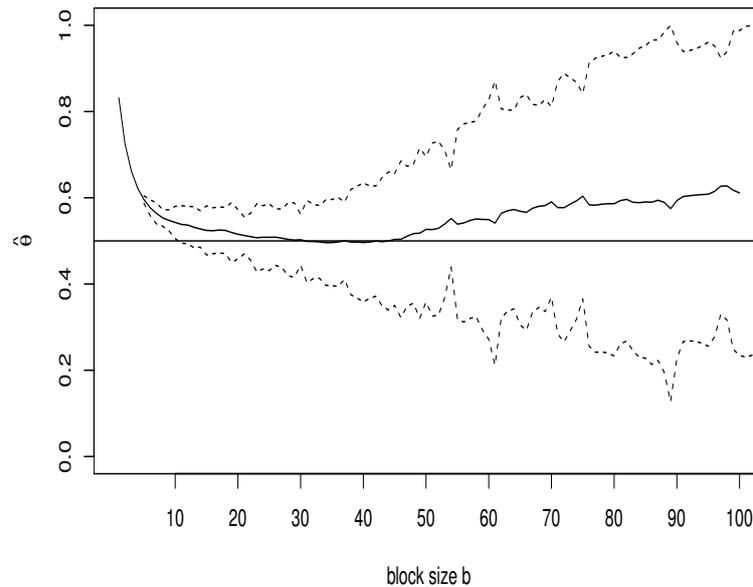
$$\hat{\theta} = \tilde{\theta}^B - 1/b_n. \tag{3}$$

We also consider the sliding blocks version since it usually performs better (Northrop [2] 2015, Berghaus and Bücher [3] 2018).

Observe that the estimators above only depend on a tuning parameter: the block length  $b \equiv b_n$ . This is an advantage of these methods since most estimators of  $\theta$  presented in the literature have two sources of uncertainty and thus two parameters to be defined in advance: the clustering generation of high values and the choice of a high threshold above which the clusters occur. To mention the best known ones, there are the Nandagopalan ([4] 1990), Runs and Blocks (Weissman and Novak, [5] 1998 and references there in),  $K$ -gaps (Süveges and Davison, [6] 2010), censored/truncated (Holěšovský and Fusek, [7,8] 2020/22), and cycles estimator (Ferreira and Ferreira, [9] 2018). We also refer to other estimators that require a single tuning parameter, such as the intervals estimator, which needs to fix a high threshold (Ferro and Segers, [10] 2003), and, similar to the Northrop estimator above, where we only choose the block length for maxima, we cite Gomes ([11] 1993), Ancona-Navarrete and Tawn ([12] 2000), and Ferreira and Ferreira ([13] 2022).

As already highlighted in the literature, there is no simple optimal methodology for the best choice of block length and a single estimate for  $\theta$ . In EVT, we have a typical bias-variance trade-off observed in sample path estimates of rare event parameters. For block estimators, the bias decreases with  $b$  while the variance increases. A recurrent method is to plot the estimates obtained for successive block size values and visually identify case-by-case plateau zones of these estimates. The stability around a value is an indicator

of a reasonable estimate, and this stability region, in general, should have neither too small nor too large a value of  $b$  due to the trade-off between bias and variance already mentioned. Figure 2 is a plot of the trajectory of estimates (full line) along with 95% confidence intervals (CI) (dashed line) obtained for each block length  $b$  from 1 to 100 in a random sample of dimension 1000 generated from a moving maximum model with standard Fréchet margins. We can see a plateau region in the estimates around the true value (horizontal line)  $\theta = 0.5$  for the block sizes between 25 and 45. Observe the large variability occurring for large values of  $b$  and the higher bias for small values of  $b$ .



**Figure 2.** Estimates of  $\hat{\theta}$  given in (3) for successive values of block size  $b = 1, \dots, 100$  (full line) obtained for a sample simulated from a moving maxima Fréchet model with  $\theta = 0.5$  (horizontal line). The dashed lines correspond to 95% CI.

Some methods have been proposed in the literature to help in the choice of tuning parameters based on the stability regions of the estimates graph: see, e.g., Frahm et al. ([14] 2005), Gomes and Neves ([15] 2020), and their references. In particular, the algorithm proposed in Frahm et al. ([14] 2005) was implemented in the context of estimating the bivariate tail dependence, and in Ferreira ([16] 2018), it was applied to extremal index estimators requiring the choice of a high threshold. In this work, our objective is to propose an adaptation of the algorithm developed in Frahm et al. ([14] 2005) applied to estimator (3) in order to find a suitable plateau of estimates taking into account the bias–variance trade-off. As a byproduct, this will allow us to circumvent the unique tuning parameter selection corresponding to the block size of where the sequence of maximums will be extracted, as described above. The method will be detailed in Section 2 and analyzed through simulation in Section 3. We end with an application to real data.

## 2. Estimation Method

Our proposed estimation of  $\theta$  is based on the bias-corrected estimator  $\hat{\theta}$  in (3) by considering sliding blocks and on the heuristic plateau-finding algorithm of Frahm et al. ([14] 2005). The algorithm is described in the following steps:

- Step 1. Calculate estimates  $\hat{\theta}_b$  from estimator (3) for  $1 \leq b \leq t < n$ ;
- Step 2. Smooth the results of the previous step by taking means of  $2w + 1$  successive estimates; we consider bandwidth  $w = \lfloor 0.02t \rfloor$ ;
- Step 3. Define plateaus of length  $m = \lfloor \sqrt{t - 2w} \rfloor$ , i.e.,  $p_j = (\hat{\theta}_j, \dots, \hat{\theta}_{j+m-1})$ ,  $j = 1, \dots, t - 2w - m + 1$ ;

Step 4. Compute the standard deviation  $s$  of  $\bar{\theta}_1, \dots, \bar{\theta}_{t-2w}$  and choose the first plateau  $p_j$  satisfying  $\sum_{i=j+1}^{j+m-1} |\bar{\theta}_i - \bar{\theta}_j| \leq 2s$ ;

Step 5. The extremal index is estimated through  $\frac{1}{m} \sum_{i=1}^m \bar{\theta}_{j+i-1}$ , i.e., taking the average of the estimates that constitute the plateau chosen in the previous step. This is denoted the plateau estimator.

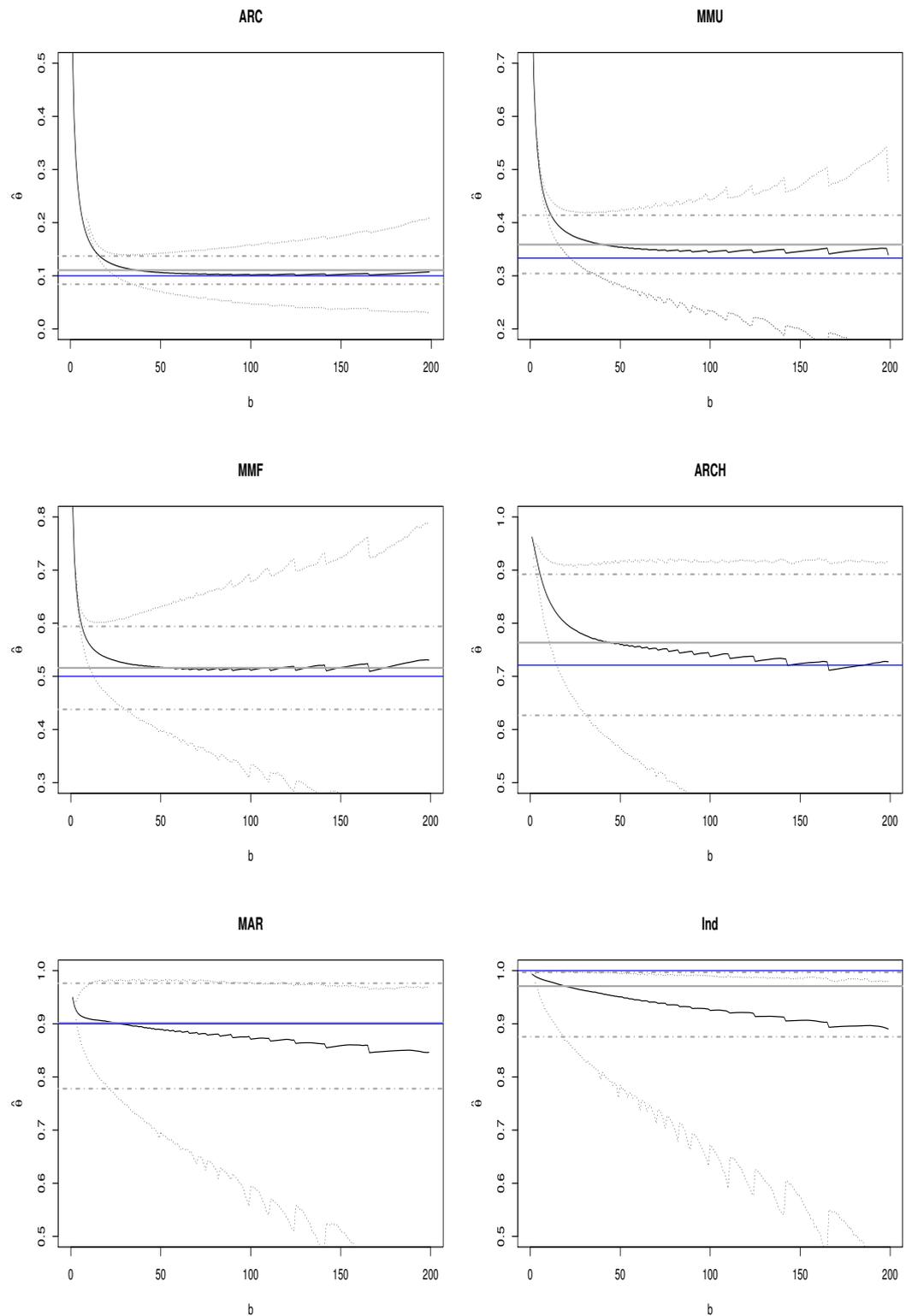
The estimators (1), (2), and (3) are already implemented in package *exdex* of software R (Northrop and Christodoulides [17] 2019) with the respective CIs. We use package *exdex* to compute estimator (3) under sliding blocks and the respective upper and lower 95% CI bounds. We also apply Steps 1, 2, and 3 to the lower and upper bounds of the CIs. Once the plateau of *theta* estimates is chosen in Step 4, we pick the corresponding plateau in the CI limits, and in Step 5, we apply the average of the plateau values of the lower limit of the CI as well as the average of the plateau values of the upper limit of the CI.

We are going to analyze the estimation method described above through simulation. The models that will be used are the following:

- First-order auto-regressive model with Cauchy standard marginals (ARC),  $X_i = \rho X_{i-1} + \epsilon_i$ ,  $\{\epsilon_i\}$  i.i.d. having Cauchy d.f. with mean 0 and scale  $1 - |\rho|$  and  $\theta = 1 - \rho$  if  $\rho > 0$  (Chernick et al. [18], 1991); we consider  $\rho = 0.9$  and  $\theta = 0.1$ ;
- An  $m$ -dependent model (MMU),  $X_i = \max(U_i, U_{i+1}, \dots, U_{i+m-1})$ ,  $i \geq 1$ , where  $\{U_i\}$  is an i.i.d. sequence of r.v. (Newell [19] 1964) with  $\theta = 1/m$ ; we consider  $U_i$ ,  $i \geq 1$ , standard uniform r.v., and  $m = 3$ , and thus,  $\theta = 1/3$ ;
- Moving maxima Fréchet model (MMF),  $X_i = \max_{j=0, \dots, d} a_j Z_{i-j}$  with  $a_j \geq 0$ ,  $\sum_{j=0}^d a_j = 1$  and  $\{Z_i\}$  i.i.d. standard Fréchet where  $\theta = \max_{j=0, \dots, d} a_j$  (Weissman and Cohen [20] 1995); we consider  $d = 2$  and parameters  $a_0 = 1/3$ ,  $a_1 = 1/6$ , and  $a_2 = 1/2$ , and thus,  $\theta = 1/2$ ;
- ARCH(1) process,  $X_i = (\beta + \alpha X_{i-1}^2)^{1/2} \epsilon_i$ , with i.i.d. Gaussian innovations  $\{\epsilon_i\}$ ,  $\alpha = 0.7$ , and  $\beta = 2 \cdot 10^{-5}$ , where  $\theta = 0.721$  (Cai, [21] 2019);
- First-order max auto-regressive (MAR),  $X_i = \max(\phi X_{i-1}, \epsilon_i)$ ,  $i \geq 1$ ,  $X_0 = \epsilon_1 / (1 - \phi)$ ,  $\{\epsilon_i\}$  i.i.d. with standard Fréchet marginals and  $\theta = 1 - \phi$  (Davis and Resnick [22] 1989); we consider  $\phi = 0.1$  and  $\theta = 0.9$ ;
- An i.i.d. sequence (Ind) of Fréchet r.v. where  $\theta = 1$ .

### 3. Simulation Study and Application

The simulation study is based on random generation of samples with size 1000 replicated 1000 times within each of the models described above. We consider different models with different values of  $\theta$ . We apply the estimation plateau method of Section 2 both to estimate  $\theta$  and the respective 95% CI lower and upper bounds. Table 1 contains the estimation global results of the plateau method. See also the simulation results of  $\hat{\theta}$  given in (3) for each block size  $b$  in Figure 3 as well as the results of the plateau method. We can observe in each model that the plateau estimate (thicker gray horizontal full line) is located in a plateau zone of the sample path of estimates plotted as a function of block size  $b$  (full black line), as expected. We can also see that the plateau estimate is close to the true value (blue horizontal full line). In all cases, it is verified that the limits of the 95% CIs estimated by the plateau method (thicker gray horizontal dotted–dashed lines) include the true value of  $\theta$ . In the ARCH case, the estimates closest to the true value of  $\theta$  occur for large values of  $b$  where the variability is very high, which makes it difficult to apply the plateau methodology. Even so, the root mean squared error (rmse) of 0.1126 is not very expressive. The independent model (Ind) has  $\theta = 1$  and, therefore, constitutes a frontier value of the parameter support, which typically leads to difficulties in statistical estimation. Still, the plateau method shows relatively low bias and rmse. Observe also that in the MAR model, we have  $\theta = 0.9$ , which is quite near to the boundary value of 1, and the plateau method does a very good job.



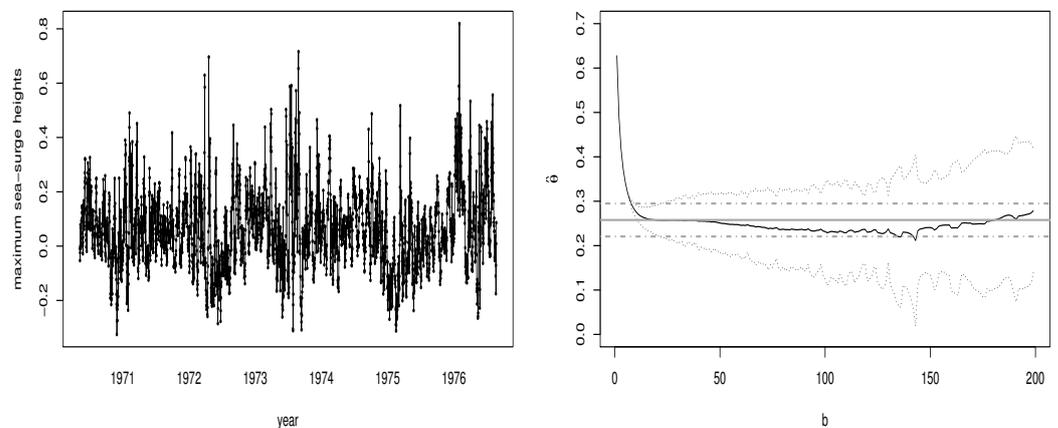
**Figure 3.** Simulation results: average of estimates of  $\theta$  for each block size  $b = 2, \dots, 200$  using  $\hat{\theta}$  in (3) (full black line) and average of respective 95% CI upper and lower bounds (dotted lines); plateau estimation of  $\theta$  (thicker gray horizontal full line) and respective plateau estimates of 95% CI upper and lower bounds (thicker gray horizontal dotted–dashed lines). The true value of  $\theta$  corresponds to the blue horizontal full line.

**Table 1.** Simulation results of plateau method: average of  $\theta$  estimates (mean), average of lower and upper 95% CI bound estimates, bias, root mean squared error (rmse), and standard deviation of  $\theta$  estimates (sd).

	mean	lower	upper	bias	rmse	sd
ARC ( $\theta = 0.1$ )	0.1106	0.0841	0.1372	0.0106	0.0218	0.0190
MMU ( $\theta = 1/3$ )	0.3587	0.3042	0.4139	0.0254	0.0494	0.0424
MMF ( $\theta = 0.5$ )	0.5160	0.4379	0.5940	0.0160	0.0636	0.0616
ARCH ( $\theta = 0.721$ )	0.7634	0.6267	0.8920	0.0424	0.1126	0.1044
MAR ( $\theta = 0.9$ )	0.9017	0.7779	0.9763	0.0017	0.0827	0.0827
Ind ( $\theta = 1$ )	0.9709	0.8756	0.9969	−0.0291	0.0643	0.0573

#### Application to Real Data

We illustrate the method with the real data *newlyn* available in the R package *exdex* consisting of 2894 sea-surge heights measured at the coast of Newlyn, Cornwall, UK, over years 1971–1976. The observations correspond to the maximum hourly surge heights during periods of 15 h. See the left plot in Figure 4. Previous analysis of this data can be seen in Northrop ([2] 2015) and references therein. The sample path of estimates from (3) as a function of block size  $b$  and respective 95% confidence limits are plotted on the right graph of Figure 4. The horizontal full line corresponds to the plateau estimate of  $\theta$  given by 0.2577, and the horizontal dotted–dashed lines correspond to the plateau 95% CI estimate (0.2206, 0.2948).



**Figure 4.** (Left) Maximum hourly (within successive 15-hour periods) surge height time series at Newlyn Coast, Cornwall, UK, in years 1971–1976; (Right) Sample path estimates obtained from estimator in (3) (full line) and respective 95% CI limits (dotted lines) for successive values of block size  $b$ , plateau estimate of  $\theta$  (horizontal full line), and respective 95% CI plateau estimate limits (horizontal dotted–dashed lines).

#### 4. Conclusions

This work addresses the estimation of the extremal index  $\theta$ . This is an important measure in time series, namely in assessing risky phenomena, as it measures the propensity for the occurrence of clusters of extreme values. The estimation of  $\theta$  requires a prior setting of tuning parameter values that impacts the precision of estimates. In this work, we presented an algorithm that allows estimation of  $\theta$  free of tuning parameters. We applied this methodology to diverse models, and the results are encouraging in several cases. In

the future, it is intended to continue the study of this methodology and develop it in order to improve its applicability to different types of models.

**Funding:** The research at CMAT was partially financed by Portuguese Funds through FCT - Fundação para a Ciência e a Tecnologia within the Projects UIDB/00013/2020 and UIDP/00013/2020.

**Institutional Review Board Statement:** Not applicable.

**Informed Consent Statement:** Not applicable.

**Data Availability Statement:** Not applicable.

**Conflicts of Interest:** The author declares no conflict of interest.

## References

1. Hsing, T.; Hüsler, J.; Leadbetter, M.R. On the exceedance point process for a stationary sequence. *Probab. Theory Relat. Fields* **1988**, *78*, 97–112. [[CrossRef](#)]
2. Northrop, P.J. An efficient semiparametric maxima estimator of the extremal index. *Extremes* **2015**, *18*, 585–603. [[CrossRef](#)]
3. Berghaus, B.; Bücher, A. Weak convergence of a pseudo maximum likelihood estimator for the extremal index. *Ann. Stat.* **2018**, *46*, 2307–2335. [[CrossRef](#)]
4. Nandagopalan, S. Multivariate Extremes and Estimation of the Extremal Index. Ph.D. Thesis, University of North Carolina, Chapel Hill, NC, USA, 1990.
5. Weissman, I.; Novak, S.Y. On blocks and runs estimators of the extremal index. *J. Stat. Plan. Inference* **1998**, *66*, 281–288. [[CrossRef](#)]
6. Süveges, M.; Davison, A.C. Model misspecification in peaks over threshold analysis. *Ann. Appl. Stat.* **2010**, *4*, 203–221. [[CrossRef](#)]
7. Holěšovský, J.; Fusek, M. Estimation of the extremal index using censored distributions. *Extremes* **2020**, *23*, 197–213. [[CrossRef](#)]
8. Holěšovský, J.; Fusek, M. Improved interexceedance-times-based estimator of the extremal index using truncated distribution. *Extremes* **2022**, *25*, 695–720. [[CrossRef](#)]
9. Ferreira, H.; Ferreira, M. Estimating the extremal index through local dependence. *Ann. L'Institut Henri-Poincaré-Probab. Stat.* **2018**, *54*, 587–605. [[CrossRef](#)]
10. Ferro, C.A.T.; Segers, J. Inference for clusters of extreme values. *J. R. Stat. Soc. Ser. B* **2003**, *65*, 545–556. [[CrossRef](#)]
11. Gomes, M. On the estimation of parameters of rare events in environmental time series. In *Statistics for the Environment 2: Water Related Issues*; Barnett, V., Turkman, K., Eds.; Wiley: Hoboken, NJ, USA, 1993; pp. 225–241.
12. Ancona-Navarrete, M.A.; Tawn, J.A. A comparison of methods for estimating the extremal index. *Extremes* **2000**, *3*, 5–38. [[CrossRef](#)]
13. Ferreira, H.; Ferreira, M. A new blocks estimator for the extremal index. *Commun.-Stat.-Theory Methods* **2022**, in press. [[CrossRef](#)]
14. Frahm, G.; Junker, M.; Schmidt, R. Estimating the tail-dependence coefficient: Properties and pitfalls. *Insur. Math. Econ.* **2005**, *37*, 80–100. [[CrossRef](#)]
15. Gomes, D.P.; Neves, M.M. Extremal index blocks estimator: The threshold and the block size choice. *J. Appl. Stat.* **2020**, *47*, 2846–2861. [[CrossRef](#)] [[PubMed](#)]
16. Ferreira, M. Heuristic Tools for the Estimation of The Extremal Index: A Comparison of Methods. *Revstat-Stat. J.* **2018**, *16*, 115–136.
17. Northrop, P.J.; Christodoulides, C. Exdex: Estimation of the Extremal Index. R Package Version 1.0.1. 2019. Available online: <https://CRAN.R-project.org/package=exdex> (accessed on 10 January 2023).
18. Chernick, M.R.; Hsing, T.; McCormick, W.P. Calculating the extremal index for a class of stationary sequences. *Adv. Appl. Probab.* **1991**, *23*, 835–850. [[CrossRef](#)]
19. Newell, G.F. Asymptotic Extremes for  $m$ -Dependent Random Variables. *Ann. Math. Stat.* **1964**, *35*, 1322–1325. [[CrossRef](#)]
20. Weissman, I.; Cohen, U. The extremal index and clustering of high values for derived stationary sequences. *J. Appl. Prob.* **1995**, *32*, 972–981. [[CrossRef](#)]
21. Cai, J.J. Statistical inference on  $D^{(d)}(u_n)$  condition and estimation of the Extremal Index. *arXiv* **2019**, arXiv:1911.06674.
22. Davis, R.; Resnick, S. Basic properties and prediction of max-ARMA processes. *Adv. Appl. Probab.* **1989**, *21*, 781–803. [[CrossRef](#)]

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.