

On Extreme Value Index Estimation under Random Censoring

Richard Minkah^{1,*} , Tertius de Wet² and Kwabena Doku-Amponsah¹

¹Department of Statistics and Actuarial Science, University of Ghana, Ghana;

²Department of Statistics and Actuarial Science, Stellenbosch University, South Africa

Received on March 01, 2018; September 30, 2018; Published Online on October 10, 2018

Copyright © 2018, African Journal of Applied Statistics (AJAS) and The Probability African Society (SPAS). All rights reserved

Abstract. Extreme value analysis in the presence of censoring is receiving much attention as it has applications in many disciplines such as survival and reliability studies. The estimation of extreme value index (EVI) is of primary importance as it is a critical parameter needed in estimating extreme events such as quantiles and exceedance probabilities. In this paper, we review several estimators of the EVI when data is subject to random censoring. In addition, we propose a reduced-bias estimator based on the exponential regression approximation of log spacings. All the estimators' performances are compared in a simulation study. The results show that no estimator is universally the best across all scenarios. However, the proposed reduced-bias estimator is found to perform well across most scenarios. Also, we present a bootstrap procedure for obtaining samples for extreme value analysis under censoring. The estimators are illustrated using a practical dataset from medical research.

Key words: Censoring, Extreme Value Index, Confidence interval; Empirical coverage probability; Confidence Interval length.

AMS 2010 Mathematics Subject Classification : 62G32; 62N02; 62F40.

Presented by Dr Tchilabola Abozou Kpanzou, University
of Kara, Togo
Corresponding Member of the Board.

*Corresponding author R. Minkah: rminkah@ug.edu.gh

T. de Wet: tdewet@sun.ac.za

K. Doku-Amponsah: KDoku-Amponsah@ug.edu.gh

Résumé (French) L'analyse de valeurs extrêmes en présence de censure fait l'objet de beaucoup d'attention car elle a des applications dans de nombreuses disciplines telles que les études de survie et de fiabilité. L'estimation de l'indice de valeur extrême (IVE) revêt une importance primordiale, car il s'agit d'un paramètre essentiel et nécessaire à l'estimation des événements extrêmes tels que les quantiles et les probabilités de dépassement. Dans cet article, nous passons en revue plusieurs estimateurs de l'IVE lorsque les données sont soumises à une censure aléatoire. En outre, nous proposons un estimateur à biais réduit basé sur l'approximation par régression exponentielle des log-espacements. Toutes les performances des estimateurs sont comparées dans une étude de simulation. Les résultats montrent qu'aucun estimateur n'est universellement le meilleur dans tous les scénarios. Cependant, l'estimateur proposé à biais réduit s'avère efficace dans la plupart des scénarios. De plus, nous présentons une procédure bootstrap pour obtenir des échantillons pour une analyse de valeur extrême sous censure. Les estimateurs sont illustrés à l'aide d'un ensemble de données pratiques issues de la recherche médicale.

1. Introduction

Statistics of extremes under random censoring is a relatively new area in extreme value analysis that has received considerable attention in the literature during the last few years. Examples of applications include estimating survival time (Einmahl *et al.*, , 2008; Ndao *et al.*, , 2014) and large insurance claims (Beirlant *et al.*, , 2017), among others.

In order to obtain estimates of parameters of extreme events, the extreme value index (EVI) is the primary parameter needed. The estimation of the EVI in the case of complete samples has been studied extensively (see e.g. Csörgő and Viharos, , 1998; Beirlant *et al.*, , 2004; Dekkers *et al.*, , 1989; Diop and Lo , 2009; Ngom and Lo, , 2016; Lo *et al.*, , 2018).

However, the same cannot be said of the estimation of the EVI when data is subject to random censoring. In this paper, we review existing estimators and propose two estimators that are aimed at reducing the bias and variance. In addition, we provide a simulation comparison of the various estimators of the Extreme Value Index (EVI).

The first work on the subject can be attributed to Beirlant and Guillou, (2001). The authors proposed an adaptation of the Hill estimator under random right censoring. The motivation for this adapted Hill estimator to censoring was the same as that of the Hill estimator obtained from the slope of the Pareto quantile plot.

However, since the censored observations have the same values (i.e the maximum), the Pareto quantile plot will be horizontal in those observations. As a result, the adaptation of the Hill estimator to censoring was based on the slope of the Pareto

quantile plot for the noncensored observations only.

In addition, by using the second order properties of the representation of log-spacings in the exponential regression model, a bias-corrected version of the adapted Hill estimator was obtained. The finite sample properties of the estimator were studied through a simulation study and the estimator was found to give credible estimates for a percentage of censoring of 5% at most. Consistency and asymptotic normality of the estimator were obtained under some restrictive conditions on the number of non-censored observations and the sample tail fraction. Delafosse and Guillou (2002) proved the almost sure convergence of the adapted Hill estimator in Beirlant and Guillou (2001) under very general conditions on the number of non-censored observations.

Delafosse and Guillou (2002) proved the almost sure convergence of the adapted Hill estimator in Beirlant and Guillou (2001) under very general conditions on the number of non-censored observations.

Also, in [Reiss and Thomas, \(2007, Section 6.1\)](#), the authors introduced an estimator of the EVI when data is randomly-or fixed censored. In the case of random right censoring, the Pareto or generalised Pareto distribution was fitted to the excesses over a given threshold. The likelihood function of the chosen distribution was adapted to censoring and maximised to obtain an estimator of the EVI. However, the authors made no attempt to study the asymptotic properties of the their proposed estimators of the EVI.

In addition, [Beirlant *et al.*, \(2007\)](#) proposed an entirely different approach by adapting the estimator of the EVI from the Peaks-Over Threshold (POT) method ([Smith, , 1987](#)) and the moment estimator ([Dekkers *et al.*, , 1989](#)) to random right censoring. The former estimator involved adapting the likelihood function to the context of censoring whereas the latter estimator was obtained by dividing the classical EVI estimator by the proportion of non-censored observations in the top order statistics selected from the sample.

Due to the difficulties in establishing the asymptotic properties of the maximum likelihood estimator of the POT method, [Beirlant *et al.*, \(2010\)](#) proposed a one-step approximation based on the Newton-Raphson algorithm. The reported simulation study showed the closeness of the approximation of the one-step estimators to the maximum likelihood estimators. The added advantage was that the asymptotic normality of the one-step estimators has been established, unlike that for the maximum likelihood estimators.

Based on the ideas of [Beirlant *et al.*, \(2007\)](#), [Einmahl *et al.*, \(2008\)](#) provided a second methodological paper which considered estimators based on the top order statistics. In addition, the authors proposed a unified method to prove the asymptotic normality of the EVI estimators.

A small scale simulation showed the superiority of the adapted Hill estimator for the Pareto domain of attraction and a slight advantage of the adapted generalised Hill for the Weibull and the Gumbel domains of attraction. [Einmahl et al., \(2008\)](#) used restrictive conditions to prove the asymptotic normality of the EVI estimators. However, these conditions were relaxed by [Brahimi et al., \(2013\)](#) to prove the asymptotic normality of the adapted Hill estimator of the EVI under random right censoring.

The estimation of the EVI has also received attention from [Gomes et al., \(2010\)](#) and [Gomes and Neves, \(2011\)](#). These papers form an overview of the EVI estimators in the context of random censoring. To the best of our knowledge, [Gomes and Neves, \(2011\)](#) made the first attempt at introducing a reduced-bias estimator of the EVI, in the form of the minimum-variance reduced-bias (MVRB) estimator ([Caeiro et al., , 2005](#)). The reported simulation study showed an overall best performance for the adapted MVRB estimator for samples generated from distributions from the Pareto domain of attraction. As in [Einmahl et al., \(2008\)](#), the generalised Hill performed better than the other adapted EVI estimators for samples whose underlying distribution functions are in the Weibull and Gumbel domains of attraction.

The Hill estimator for estimating the EVI under random censoring performs well, although in the classical case it is known to be biased, not location invariant and unstable. Efforts have been made to provide reduced-bias and minimum variance Hill-type estimators to improve on the Hill estimator for the heavy-tailed distributions (i.e. distributions in the Pareto domain of attraction).

In this regard, [Worms and Worms, \(2014\)](#) provided another methodological paper for the estimation of the EVI in the case of censoring. They provided two sets of Hill-type estimators based on the Kaplan-Meier estimation of the survival function (see [Kaplan and Meier, , 1958](#)) and the synthetic data approach of [Leurgans, \(1987\)](#).

In addition, the authors presented a small scale simulation that compared the performance of the two proposed estimators to the adapted Hill and MVRB estimators. The results showed that the two proposed estimators are superior to the Hill estimator, in particular, the estimator based on the ideas of [Leurgans, \(1987\)](#).

On the other hand, MVRB performed better than the authors' proposed estimators. However, the EVI estimator based on the synthetic data approach of [Leurgans, \(1987\)](#) compared favourably in the strong censoring framework with the MVRB estimator. The consistency of these estimators was proved under mild censoring. However, the asymptotic normality of these two estimators remains an open problem.

Furthermore, the estimation of the EVI for the Pareto domain of attraction has also been obtained from the Bayesian perspective by [Ameraoui et al., \(2016\)](#). They constructed a maximum a posteriori and mean posterior estimators for various prior distributions of the EVI, namely Jeffrey's, Maximal Data Information (MDI) and a conjugate Gamma. The asymptotic properties, namely consistency and normality of the estimators, were established. A small simulation study was used to examine the finite sample properties and the performance of the estimators. The reported simulation result showed the superiority of the maximum a posteriori estimator under maximal data information prior.

We aim to achieve two objectives in this paper. Firstly, we propose some estimators of the EVI including a reduced-bias estimator based on the exponential regression model of [Beirlant et al., \(1999\)](#).

Secondly, the above researchers compared their proposed estimators under different simulation conditions. In addition, some of the estimators' asymptotic distributions remain an open problem, and hence, theoretical comparison is not possible. Therefore, the second objective of this paper is to compare several of the existing estimators with the proposed ones in a simulation study under identical conditions.

The rest of the paper is organised as follows. In Section 2, we present the framework of extreme value analysis when data is censored. In Section 3, a simulation comparison of the various estimators is presented. In Section 4, we present a practical application of the estimators to estimate the extreme value index for a medical data set on the survival of AIDS patients. Lastly, concluding remarks are presented in Section 5.

2. Framework

Let X_1, X_2, \dots, X_n be a sequence of independent and identically distributed (*i.i.d*) random variables with distribution function F , and $X_{1,n} \leq X_{2,n} \leq \dots \leq X_{n,n}$ the associated order statistics. Therefore, the sample maximum is denoted by $X_{n,n}$. Extreme value theory attempts to solve the problem of the possible limit distributions of $X_{n,n}$. It is well-known that the distribution of the sample maximum can be obtained from the underlying distribution of X as

$$F_{X_{n,n}}(x) = F^n(x). \quad (1)$$

However, F is usually unknown and, hence, EVT focuses on the search for an approximate family of models for F^n as $n \rightarrow \infty$.

Limiting results for F^n in EVT have been addressed in the papers by [Fisher and Tippett, \(1928\)](#) and [Gnedenko, \(1943\)](#). Specifically, the results can be stated as follows: if there exist sequences of constants b_n and $a_n > 0$ ($n = 1, 2, \dots$), such that

$$\lim_{x \rightarrow \infty} P \left(\frac{X_{n,n} - b_n}{a_n} \leq x \right) \rightarrow \Psi(x), \quad (2)$$

where Ψ is a nondegenerate distribution function, then Ψ belongs to the family of distributions,

$$\Psi_\gamma(x) = \begin{cases} \exp \left(- \left(1 + \gamma \frac{x-\mu}{\sigma} \right)^{-1/\gamma} \right), & 1 + \gamma \frac{x-\mu}{\sigma} > 0, \gamma \neq 0, \\ \exp \left(- \exp \left(\frac{x-\mu}{\sigma} \right) \right), & x \in \mathbb{R}, \gamma = 0, \end{cases} \quad (3)$$

where $\mu \in \mathbb{R}$ and $\sigma > 0$. The quantity $\gamma \in \mathbb{R}$, is the *Extreme Value Index* (EVI) or the *tail index*: it determines the tail heaviness of the extreme value distributions. The EVI is classified into three groups, each representing one of the three families of distributions, Gumbel (exponential tails), Pareto (heavy-tailed) and Weibull (light-tailed). The group of families have $\gamma = 0$, $\gamma > 0$ and $\gamma < 0$ corresponding to the Gumbel, Pareto and Weibull families respectively. A distribution function F satisfying (3) is said to be in the maximum domain of attraction of Ψ_γ written as $F \in D(\Psi_\gamma)$.

In addition to (3), Balkema and de Haan, (1974) and Pickands III, (1975) showed the generalised Pareto distribution (GPD) as the limit distribution of scaled excesses over a sufficiently large threshold. The GPD can be written as

$$\Lambda_\gamma(x) = 1 + \ln \Psi_\gamma(x) = \begin{cases} 1 - \left(1 + \gamma \frac{x-\mu}{\sigma} \right)^{-1/\gamma}, & 1 + \gamma \frac{x-\mu}{\sigma} > 0, \gamma \neq 0, \\ 1 - \exp \left(\frac{x-\mu}{\sigma} \right), & x \in \mathbb{R}, \gamma = 0, \end{cases} \quad (4)$$

where Ψ_γ is given in (3).

In this paper, our interest is in the Pareto domain of attraction i.e. the case $\gamma > 0$. This family consists of distribution functions F whose tails are regularly varying with a negative index of variation. That is

$$1 - F(x) = x^{-1/\gamma} \ell_F(x), \quad x \rightarrow \infty, \quad (5)$$

where ℓ_F is the slowly varying function associated with F . A slowly varying function, ℓ , is of the form $\ell(xt)/\ell(x) \rightarrow 1$ for $x \rightarrow \infty$. Relation (5) can be stated equivalently in terms of the associated upper tail quantile function U as

$$U(x) = F^{-1} \left(1 - \frac{1}{x} \right) = x^\gamma \ell_U(x), \quad x \rightarrow \infty, \quad (6)$$

where ℓ_U is the slowly varying function associated with U .

2.1. EVT Conditions

The conditions underlying domain of attraction are presented in this section. These conditions are needed in defining estimators of tail parameters and to study their asymptotic properties.

de Haan, (1984) gave the following well-known necessary and sufficient condition for $F \in D(\Psi_\gamma)$, known as the first-order condition or extended regular variation:

$$\lim_{u \rightarrow \infty} \frac{U(ux) - U(u)}{a(u)} = h_\gamma(x) := \begin{cases} \frac{x^\gamma - 1}{\gamma} & \text{if } \gamma \neq 0 \\ \log x & \text{if } \gamma = 0, \end{cases} \quad (7)$$

where a is a positive measurable function, $x > 0$.

In addition, to study the asymptotic properties of the estimators of tail parameters, the first-order condition is generally not sufficient; a second-order condition specifying the rate of convergence of (7) is also required.

In the literature, the second-order condition can be stated in terms of U (see e.g. de Haan and Ferreira, , 2006; Gomes *et al.*, , 2008), or, equivalently, also in terms of the rate of convergence of the slowly varying function, ℓ , in (6). Beirlant *et al.*, (1999, page 602) state it as follows:

there exists a real constant $\rho < 0$ and a rate function b satisfying $b(x) \rightarrow 0$ as $x \rightarrow \infty$, such that for all $\lambda \geq 1$,

$$\lim_{x \rightarrow \infty} \frac{\log \ell(\lambda x) - \log \ell(x)}{b(x)} = \kappa_\rho(\lambda), \quad (8)$$

where $\kappa_\rho(\lambda) = \int_1^\lambda u^{\rho-1} du$.

2.2. General Estimation under Censored Data

Let the random variable of interest be X with distribution function, F . Since samples on X may not be fully observed, we introduce another positive random variable C , which is independent of X , with distribution function G . In this setting, we then observe $(Z_i, \delta_i), i = 1, \dots, n$ with

$$Z_i = \min(X_i, C_i) \quad (9)$$

and

$$\delta_i = \begin{cases} 1 & \text{if } X_i \leq C_i; \\ 0 & \text{if } X_i > C_i. \end{cases} \quad (10)$$

Here, δ_i is a variable indicating whether Z_i is censored or not. Let H be the distribution function of Z defined in (9). Thus, by the independent assumption of the random variables Y and C , we have $1 - H = (1 - F)(1 - G)$.

In addition, let $\vartheta_F = \sup\{F(x) < 1\}$ be the corresponding right endpoint of the underlying distribution function, F . Similarly, let ϑ_G and ϑ_H be the right endpoints of the underlying distribution functions of C and Z respectively. If we assume $F \in D(\Psi_{\gamma_1})$ and $G \in D(\Psi_{\gamma_2})$ for some real numbers, γ_1 and γ_2 , then $H \in D(\Psi_\gamma)$ where $\gamma \in \mathbb{R}$. Einmahl *et al.*, (2008) considered these three combinations of γ_1 and

γ_2 :

Case 1.

(1) F and G are Pareto types: $\gamma_1 > 0, \gamma_2 > 0 \rightarrow \gamma = \frac{\gamma_1 \gamma_2}{\gamma_1 + \gamma_2}$

(2) F and G are Gumbel types: $\gamma_1 = 0, \gamma_2 = 0, \vartheta_F = \vartheta_G \rightarrow \gamma = 0$

(3) F and G are Weibull types: $\gamma_1 < 0, \gamma_2 < 0, \vartheta_F = \vartheta_G = \infty \rightarrow \gamma = \frac{\gamma_1 \gamma_2}{\gamma_1 + \gamma_2}$.

The other two possibilities, $\{\gamma_1 > 0, \gamma_2 < 0\}$ and $\{\gamma_1 < 0, \gamma_2 > 0\}$, correspond closely to the completely noncensored case which has been studied widely whereas the latter corresponds closely to the completely censored case where estimation is impossible.

2.3. Extreme Value Index Estimation Methods

The estimation of the extreme value index (EVI) when observations are censored needs some modification from that of the complete sample. This is because the observed sample is $(Z_i, \delta_i), i = 1, \dots, n$, and hence, the application of the classical EVI estimation methods will yield estimators that converge to γ , the EVI of the underlying distribution of the random variable Z . However, our interest is in γ_1 , the EVI of the underlying distribution of the random variable X . Therefore, some modification is needed to adapt the estimation of γ from the Z sample to estimate γ_1 .

The existing methodologies for estimating the EVI under right censoring can be grouped into four categories:

Case 2.

(1) adapting a classical EVI by dividing it by the proportion of noncensored observations (Beirlant *et al.*, , 2007; Einmahl *et al.*, , 2008; Gomes and Neves, , 2011);

(2) adapting the likelihood function of an extreme value distribution (Beirlant *et al.*, , 2010);

(3) Censored regression (Worms and Worms, , 2014).

(4) Bayesian estimation (Ameraoui *et al.*, , 2016; Beirlant *et al.*, , 2017)

In this paper, we consider the frequentist methods only i.e. the first three cases. These methods and the resulting estimators are grouped into three categories and presented in the three sub-sections that follow. Following that, we propose a reduce-bias estimator based on exponential regression model and adapted to the censored case.

2.3.1. First Method

The first method was introduced in [Beirlant et al., \(2007\)](#) and further developed by [Einmahl et al., \(2008\)](#). In this method a classical estimator of the EVI is obtained from the Z sample and then adapted to censoring. Among these estimators are: the maximum likelihood estimator from the Peaks-Over Threshold (POT) method and the moment estimator ([Beirlant et al., \(2007\)](#)); Hill, Moment, Generalised Hill and the maximum likelihood estimator from the POT method ([Einmahl et al., \(2008\)](#)); and Hill, moment, mixed moment and generalised Hill ([Gomes and Neves, \(2011\)](#)).

In addition, [Einmahl et al., \(2008\)](#) provides a uniform way to establish the asymptotic normality of the proposed estimators of the EVI (i.e. Hill, Moment, Generalised Hill and the maximum likelihood estimators). These estimators are reviewed below in terms of the random variable Z , and thus estimates γ the EVI of Z .

The Hill Estimator: The Hill estimator ([Hill, \(1975\)](#)) is arguably the most common estimator of γ in the Pareto case i.e. $\gamma > 0$. The Hill estimator is defined for the $(k + 1)$ -largest order statistics as

$$\hat{\gamma}_{Z,k,n}^{(Hill)} = \frac{1}{k} \sum_{j=1}^k \log Z_{n-j+1,n} - \log Z_{n-k,n}. \quad (11)$$

The properties of the Hill estimator have been studied widely and its attractive properties include consistency ([Mason, \(1982\)](#)) and asymptotic normality ([Hall, \(1982\)](#); [de Haan and Peng, \(1998\)](#)).

The Generalised Hill Estimator: [Beirlant et al., \(1996\)](#) proposed the generalised Hill (UH) estimator in a bid to extend the Hill estimator to the case where $\gamma \in \mathbb{R}$. The UH estimator is obtained as the slope of the ultimately linear part of the generalised Pareto quantile plot,

$$\left(-\log \left(\frac{j+1}{n+1} \right), \log (Z_{n-j,n} H_{Z,j,n}) \right), j = 1, 2, \dots, n-1. \quad (12)$$

It is given by

$$\hat{\gamma}_{Z,k,n}^{(UH)} = \frac{1}{k} \sum_{j=1}^k \log UH_{Z,j,n} - \log UH_{Z,k+1,n}, \quad (13)$$

where $UH_{Z,j,n} = Z_{n-j,n} \left(\frac{1}{j} \sum_{i=1}^j \log Z_{n-i+1,n} - \log Z_{n-j,n} \right)$.

The Minimum-Variance Reduced Bias Estimator: [Caeiro et al., \(2005\)](#) proposed the Minimum-Variance Reduced Bias (MVRB) estimator for heavy-tailed distributions belonging to the Hall class ([Hall, \(1982\)](#)) of models. The estimator is a direct modification of the Hill estimator using the second order parameters to reduce bias.

It has the added advantage of having the same asymptotic variance as the Hill estimator. The MVRB estimator is obtained by using the second-order condition (8) with $b(u) = \gamma\beta u^\rho$. It is given by

$$\hat{\gamma}_{Z,k,n}^{(MVRB)} = \hat{\gamma}_{Z,k,n}^{(Hill)} \left(1 - \frac{\hat{\beta}}{1 - \hat{\rho}} \left(\frac{k}{n} \right)^{-\hat{\rho}} \right), \quad (14)$$

where, $\hat{\gamma}_{Z,k,n}^{(Hill)}$ is the Hill estimator in (11) and the pair $(\hat{\beta}, \hat{\rho})$ is the estimator for the pair of parameters (β, ρ) of the second-order auxiliary function b .

The Moment Estimator: Dekkers *et al.*, (1989) introduced another estimator known as the moment estimator as an adaptation of the Hill estimator valid for all domains of attraction. The moment estimator is defined for $k \in \{2, \dots, n - 1\}$ and it is given by

$$\hat{\gamma}_{Z,k,n}^{(MOM)} = M_{Z,k,n}^{(1)} + 1 - \frac{1}{2} \left(1 - \frac{(M_{Z,k,n}^{(1)})^2}{M_{Z,k,n}^{(2)}} \right)^{-1}, \quad (15)$$

where

$$M_{Z,k,n}^{(j)} = \frac{1}{k} \sum_{i=1}^k (\log Z_{n-i+1,n} - \log Z_{n-k,n})^j, \quad j = 1, 2.$$

Adapting EVI Estimators

Beirlant *et al.*, (2007) and Einmahl *et al.*, (2008) proposed that the EVIs for the complete sample, $\gamma_{Z,k,n}^{(\cdot)}$, (i.e. (11) - (15)) can be adapted to censoring by dividing each estimator by the proportion of noncensored observations, $\hat{\phi}$, in the k largest Z observations. Thus, the estimator of γ_1 is given by

$$\hat{\gamma}_1 = \hat{\gamma}_{Z,k,n}^{(c,\cdot)} = \frac{\hat{\gamma}_{Z,k,n}}{\hat{\phi}}. \quad (16)$$

Here, $\hat{\phi}$ is given by

$$\hat{\phi} = \frac{1}{k} \sum_{i=1}^k \delta_{n-i+1,n}, \quad (17)$$

where $\delta_{i,n}$, $i = 1, \dots, n$ are the δ -values corresponding to $Z_{i,n}$, $i = 1, \dots, n$ respectively. In the literature, (16) has primarily been used to adapt the EVI estimators to censoring.

2.3.2. Second Method

The second method introduced by Beirlant *et al.*, (2010) involves using the POT method and adapting the log-likelihood function for censoring. We know from (4) that given a high threshold, u , the limit distribution of excesses

$V_j = Z_i - u, j = 1, \dots, k$ given $Z_i > u, i = 1, \dots, n$ can be approximated by the generalised Pareto (GP) distribution. In Beirlant *et al.*, (2007) and Einmahl *et al.*, (2008), the maximum likelihood estimator, $\hat{\gamma}^{(c,POT)}_{Z,k,n}$, is obtained from the GP approximation of the distribution of the V_j 's and is adapted to censoring using (16).

An alternative approach in Beirlant *et al.*, (2010) involves adapting the likelihood function of the random variable $V_j, j = 1, \dots, k$,

$$L(\gamma_1, \sigma_{1,k}) = \prod_{j=1}^k [\lambda(V_j)]^{\delta_j} [1 - \Lambda(V_j)]^{1-\delta_j} \tag{18}$$

where Λ is the GP distribution and λ the corresponding density function of the GP distribution. However, there are difficulties with obtaining explicit expressions for the maximum likelihood estimators of γ_1 and $\sigma_{1,k}$. In addition, their asymptotic properties remain an open problem. As a result, Beirlant *et al.*, (2010) proposed solving the maximum likelihood equations using one-step approximations based on the Newton-Raphson algorithm. The resulting estimator of the parameters is given by

$$\begin{pmatrix} \hat{\gamma}_{Z,k,n}^{(c,POT.L)} \\ \frac{\hat{\sigma}_{Z,k}^{(c,POT.L)}}{\sigma_{1,k}} \end{pmatrix} = \begin{pmatrix} \hat{\gamma}_{Z,k,n}^{(c,I)} \\ \frac{\hat{\sigma}_{Z,k}^{(c,I)}}{\sigma_{1,k}} \end{pmatrix} - \begin{pmatrix} L''_{11} & \sigma_{1,k} L''_{12} \\ \sigma_{1,k} L''_{12} & \sigma_{1,k}^2 L''_{22} \end{pmatrix} \begin{pmatrix} L'_1 \\ \sigma_{1,k} L'_2 \end{pmatrix} \tag{19}$$

where L'_i and $L''_{ij}, i = 1, 2, j = 1, 2$ are the first and second derivatives of $\log L(\gamma_1, \sigma_{1,k})$, evaluated at $(\hat{\gamma}_{Z,k,n}^{(c,I)}, \hat{\sigma}_{Z,k}^{(c,I)})$. The estimators, $\hat{\gamma}_{Z,k,n}^{(c,I)}$ and $\hat{\sigma}_{Z,k}^{(c,I)}$, are the initial estimators and must be asymptotically normal. The authors state that the moment estimator provides a good example of the initial estimators. The performance of the estimators, $\hat{\gamma}_{Z,k,n}^{(c,POT.L)}$ and $\hat{\sigma}_{Z,k}^{(c,POT.L)}$, were found to be close to the maximum likelihood estimators obtained from (18). In addition, the asymptotic normality of the one-step Newton-Raphson estimators obtained in (19) has been established in that paper.

2.3.3. Third Method

The third method introduced by Worms and Worms, (2014) is based on censored regression method of Koul *et al.*, (1981). The estimators are valid for estimating the EVI for distributions in the Pareto domain of attraction. From the well known result of deriving the Hill estimator from the mean excess function, they define an adaptation of the classical Hill estimator valid for case 1 as,

$$\hat{\gamma}_{Z,k,n}^{(c,WW.KM)} := \frac{1}{n(1 - \hat{F}(Z_{n-k,n}))} \sum_{j=1}^k \frac{\delta_{n-j+1,n}}{1 - \hat{G}(Z_{n-j+1,n}^-)} \log \left(\frac{Z_{n-j+1,n}}{Z_{n-k,n}} \right), \tag{20}$$

where \hat{F} and \hat{G} are the Kaplan-Meier estimators for F and G respectively. Here, the Kaplan-Meier estimators of the survival functions are defined for $b < Z_{n,n}$ as

$$1 - \hat{F}(b) = \Pi_{Z_{j,n} \leq b} \left(\frac{n-j}{n-j+1} \right)^{\delta_{j,n}} \quad (21)$$

and

$$1 - \hat{G}(b) = \Pi_{Z_{j,n} \leq b} \left(\frac{n-j}{n-j+1} \right)^{1-\delta_{j,n}}. \quad (22)$$

In practice, the estimator $1 - \hat{G}(Z_{n-j+1,n}^-)$ can be equal to zero, making (20) undefined. Therefore, Worms and Worms, (2014) defined $\hat{G}(Z_{n-j+1,n}^-)$ as a function of the form $g(z^-) = \lim_{\nu \rightarrow z} g(\nu)$.

As an alternative to the Kaplan-Meier estimators of F and G , Worms and Worms, (2014) provides a variant of (20) based on the ideas of “synthetic data” introduced by Leurgans, (1987). The estimator turns out to be a weighted version of the Hill estimator, (20), and is given by

$$\hat{\gamma}_{Z,k,n}^{(c,WW.L)} := \frac{1}{n(1 - \hat{F}(Z_{n-k,n}))} \sum_{j=1}^k \frac{1}{1 - \hat{G}(Z_{n-j+1,n}^-)} j \log \left(\frac{Z_{n-j+1,n}}{Z_{n-k,n}} \right). \quad (23)$$

The consistency of the estimators (20) and (23) were proven under some restrictive conditions. However, the asymptotic normality of the estimators (20) and (23) remains an open-problem.

2.3.4. The Proposed Estimator

We propose adapting the exponential regression method of Beirlant *et al.*, (1999) to censoring. This method yields a maximum likelihood (ML) estimator for $\gamma > 0$, and hence, for $\gamma_1 > 0$.

Beirlant *et al.*, (1999) provide an approximate representation for the log-spacings of successive order statistics:

$$R_j = j(\log Z_{n-j+1,n} - \log Z_{n-j,n}) \sim \left(\gamma + b_{n,k} \left(\frac{j}{k+1} \right)^{-\rho} \right) E_j, \quad j = 1, \dots, k, \quad (24)$$

where E_j , $j = 1, \dots, k$ are standard exponential random variables, $b_{n,k} = b((n+1)/(k+1)) \in \mathbb{R}$ (also $b_{n,k} \rightarrow 0$, as $k, n \rightarrow \infty$) and ρ are second-order parameters from (8). From the approximate distribution of log-spacings (24), a likelihood function can be formed. Maximisation of the likelihood function leads to the maximum likelihood estimators $\hat{\gamma}_{Z,k,n}^{(ERM)}$, $\hat{b}_{n,k}$ and $\hat{\rho}$ of γ , $b_{n,k}$ and ρ respectively. We note that (24) simplifies to $R_j \sim \gamma E_j$, $j = 1, \dots, k$ if $b_{n,k} = 0$. In addition, the resulting maximum likelihood estimator is the usual Hill estimator.

The maximum likelihood estimator, $\hat{\gamma}_{Z,k,n}^{(ERM)}$, of γ is adapted to censoring to obtain an estimator of γ_1 using (16). Moreover, the estimation of γ leads to concurrent estimates of the second order parameters, $\hat{b}_{n,k}$ and $\hat{\rho}$. These estimators can be adapted to censoring and used to obtain reduced-bias estimators for quantiles and exceedance probabilities.

3. Simulation Study

To investigate and compare the performance of different EVI estimators, we shall make use of simulation. The simulation study is grouped into two categories: point and confidence interval estimation. The former involves assessing the performance of the estimators in terms of Median Absolute Deviation (MAD) and median bias. The latter case consists of diagnostic checks on 95% confidence intervals based on the coverage probabilities and interval lengths.

We consider the following combination of factors in the simulation: distributions, sample sizes, threshold levels, proportions of censoring. Several samples sizes, $n = 500, 1000, 2000$ and 5000 , and number of top order statistics, taken as 10%, 20% and 30% of the sample size. However, the result did not differ so much and hence, for brevity and ease of presentation, we consider samples of size, $n = 1000$ and the number of top order statistics taken as 10% of the sample size.

Data were generated from the three distributions presented in Table 1.

Table 1. Distributions

Distribution	$1 - F(z)$	γ
Burr (η, τ, λ)	$(\eta/(\eta + z^\tau))^\lambda, \quad z > 0; \eta, \lambda, \tau > 0$	$\frac{1}{\tau\lambda}$
Pareto (α)	$z^{-\alpha}, \quad z > 1; \alpha > 0$	$\frac{1}{\alpha}$
Fréchet (α)	$1 - \exp(-z^{-\alpha}), \quad z > 1; \alpha > 0$	$\frac{1}{\alpha}$

With regard to the proportion of censoring in the right tail, we consider three values: 0.10 (small), 0.35 (medium) and 0.65 (large). This allows us to study the performance of the estimators as censoring increases or decreases.

3.1. Simulation Design

In this section, we examine the procedure for measuring the performance of point and interval estimators of the EVI. In the case of point estimators, the median of R ($R = 1000$) repetitions was used as the point estimate of γ_1 , and MAD and median bias are obtained as the performance measures.

On the other hand, the comparison of the confidence intervals are based on two properties: interval length and coverage probability. Before, we introduce the simulation algorithm to compute the diagnostics of the confidence interval, we present a procedure known as the conditional block bootstrap for obtaining samples for extreme value analysis in the case of censoring.

3.1.1. Conditional Block Bootstrap for Censored Data

In order to obtain the performance measures, coverage probability and average interval length, the bootstrap samples are required. However, as stated in Section 2.2, two scenarios in EVT in the case of censoring are to be avoided in this study. Firstly, if none of the observations are censored (i.e. as can happen in cases where $\gamma_1 > 0$ and $\gamma_2 < 0$), then the classical EVT estimation techniques apply. This has been widely studied in the literature and is not of interest in this paper. Secondly, for a completely censored case (which can occur when $\gamma_1 < 0$ and $\gamma_2 > 0$) the estimation of the EVI and the other extreme events are impossible.

Therefore, any bootstrap procedure implemented for the estimation of parameters of extreme events for censored data must be constrained to exclude the above scenarios, particularly where the estimation is impossible. However, the bootstrap sampling Efron and Tibshirani, (1993) and bootstrap for censored data Efron, (1981) do not guarantee the exclusion of these two scenarios.

We present here a bootstrap procedure, termed the “conditional block bootstrap”, for selecting bootstrap samples that exclude the two scenarios in statistics of extremes when data is subject to random censoring. The conditional block bootstrap is a combination of ideas from the moving block bootstrap (Efron and Tibshirani, , 1993) and the bootstrap for censored data (Efron, , 1981).

In this procedure, the censored data is grouped into randomly chosen blocks and it is crucial that each block must contain at least one censored observation. This ensures that the second case is eliminated from each generated bootstrap sample. The bootstrap observations are obtained by repeatedly sampling with replacement from these blocks and placing them together to form the bootstrap sample. Enough blocks must be sampled to obtain approximately the same sample size as the original censored sample.

Given a sample of size, n , a proportion of censoring in the right tail, φ , and assuming $\varphi \leq 0.5$, the conditional block bootstrap procedure is as follows:

(1) Group the n observations into two groups namely, censored and noncensored with sample sizes n_c and $n_{\bar{c}}$ respectively. Thus, $\varphi = n_c/n$.

(2) Let d ($d \geq 1$) denote the number of censored observations to be included in each block. The size of each block, s , is obtained as $(n \times d)/n_c$. If s is not an integer, then let $s = \lceil (n \times d)/n_c \rceil$.

(3) The number of blocks, m , is chosen such that $n \cong m \times s$. In the case, $n = m \times s$, the blocks will have the same number of observations. Otherwise, if $n \approx m \times s$, then m is taken as $\lceil n/s \rceil$, in which case the first $m - 1$ blocks are allocated s observations each and the remaining $n - s(m - 1)$ observations, allocated to the m th block.

(4) Let b_i , $i = 1, \dots, m$ denote the m blocks. Assign observations to each block by randomly sampling, $s - d$ observations without replacement from the noncensored group. In addition, randomly sample d observations without replacement from the censored-group and assign to each block b_i , $i = 1, \dots, m$. Thus, each block would contain d and $s - d$ observations that are censored and noncensored respectively.

(5) Sample m times with replacement from b_1, b_2, \dots, b_m and place them together to form the bootstrap sample. Note that, more than m blocks may be sampled, in the case, $n \approx m \times s$, for the bootstrap sample to be approximately equal to the original sample size, n .

(6) Repeat (5) a large number of times, B , to obtain B bootstrap samples.

In the case, $\varphi < 0.5$, the above procedure can be used to constitute the blocks. However, the allocations should be done such that each block contains at least one noncensored observation.

3.1.2. Simulation algorithm

The following algorithm is used to obtain performance measures of the estimators of γ_1 :

A1.]

(1) Generate n observations from Y and C respectively, and hence, obtain $Z^{(1)} = \{Z_1, \dots, Z_n\}$ and $\delta^{(1)} = \{\delta_1, \dots, \delta_n\}$. Repeat a large number of times $R - 1$ ($R = 1000$) to obtain R pairs of $(Z^{(i)}, \delta^{(i)})$, $i = 1, \dots, R$ samples.

(2) Select the pair of samples, $(Z^{(1)}, \delta^{(1)})$. Draw B ($B = 1000$) bootstrap samples each of size n using the conditional block bootstrap procedure in Section 3.1.1.

(3) Compute the bootstrap replicates, $\hat{\gamma}_{1,1}^{*(c,\cdot)}, \dots, \hat{\gamma}_{1,B}^{*(c,\cdot)}$, using the estimators of γ_1 .

(4) Compute the $100(1 - \alpha)\%$, bootstrap confidence interval.

(5) Repeat A3.1.2 through to A3.1.2 for the remainder of the pairs of samples, $(Z^{(j)}, \delta^{(j)})$, $j = 2, \dots, R$ to obtain R confidence intervals for γ_1 .

(6) Compute the properties of confidence intervals i.e. coverage probability and average interval length using the R confidence intervals in A3.1.2.

3.2. Results and Discussions

In this section, we discuss the results of the simulation study for each distribution. General comments across the various distribution are presented in the last section. The simulation results for the Burr, Pareto and Fréchet distributions are presented in Appendices A, B and C respectively. In most cases, estimators having small values of MAD and median bias generally give better coverage probability and interval length. Therefore, our performance measuring criterion focuses on the coverage probability (CP) and interval lengths. Generally, we regard a good estimator as having a coverage probability of at least 0.90 and a reasonable interval length among these estimators.

3.2.1. Burr Distribution

– **For** $\gamma_1 = 0.1$:

The ERM estimator is undoubtedly the best confidence interval estimator of $\gamma_1 = 0.1$ as it has small bias, MAD, CP approximately equal to the nominal level and shorter average confidence interval length. For percentage of censoring in the right tail, $\varphi = 10\%$ (or more generally $\varphi \leq 10\%$), other estimators of $\gamma_1 = 0.1$ including MOM and occasionally POT.L, have good CP values. However, these estimators have wider average interval lengths compared to the ERM estimator.

Moreover, in the case of $\varphi > 10\%$, ERM is the only estimator that has coverage probability close to the nominal level and has a shorter confidence interval length. Also, POT.L has good CP values but larger interval lengths, and hence, not recommended for estimating $\gamma_1 = 0.1$. The apparent poor performance of most of the estimators of γ_1 may be due to the second-order parameter $\rho \rightarrow 0$.

– **For** $\gamma_1 = 0.5$:

Hill, MVRB, WW.KM and WW.L are the best estimators of γ_1 for a small percentage of censoring less than or equal to 10%. These estimators have CP values close to the nominal level and small average interval lengths. As the percentage of censoring increases, the ERM and POT.L estimators have the best CP values: the other estimators have poor coverage probabilities. Overall, ERM and POT.L are the estimators which have good CP values and can be considered for estimating $\gamma_1 = 0.5$. In addition, ERM has shorter interval lengths, and hence, can be considered as the most appropriate for estimating $\gamma_1 = 0.5$.

– **For** $\gamma_1 = 0.9$:

Most of the confidence interval estimators perform very well for the estimation of $\gamma_1 = 0.9$ compared with $\gamma_1 \leq 0.5$. The Hill, MVRB, WW.L, ERM and POT.L estimators generally give CP values close to the desired level of 0.95 regardless of the percentage of censoring in the right tail. Among these estimators, POT.L has the largest average interval length followed by ERM. Overall, the Hill, MVRB and WW.L possesses the best attributes in terms of MAD, CP and interval lengths, and hence, are the most appropriate estimators of $\gamma_1 = 0.9$.

3.2.2. Pareto Distribution

– **For** $\gamma_1 = 0.1$:

In this case, regardless of the percentage of censoring in the right tail, few estimators of γ_1 have CP values close to the nominal level and moderate interval lengths. These include MOM and POT. The rest of the estimators have poor CP values close to zero except ERM and POT.L. However, POT.L has larger interval length, and hence, may not be appropriate an appropriate estimator of γ_1 . Thus, MOM and POT are the most robust to censoring when estimating $\gamma_1 = 0.1$.

– **For** $\gamma_1 = 0.5$:

In the case of the estimation of $\gamma_1 = 0.5$, more estimators satisfy the CP-Interval length criterion when compared to $\gamma_1 = 0.1$. Estimators such as ERM, MOM, POT and POT.L mostly have high CP values close to 0.95 regardless of the value of φ . Again, the POT.L estimator has the largest interval length. Overall, the MOM and ERM are the preferred estimators as they have better CP values and moderate interval lengths compared with the others.

– **For** $\gamma_1 = 0.9$:

For a small percentage of censoring in the right tail, $\varphi = 10\%$, most of the estimators have good CP values. The exceptions to this include WW.KM and WW.L. Also, when $\varphi = 0.35$ and 0.65 the WW.KM, MOM, POT, POT.L and ERM estimators have good CP values and relatively moderate interval lengths. However, POT.L always has the largest interval length of at least twice the estimator with the shortest interval length. Therefore, ERM, MOM and POT can be considered as more robust for the estimation of $\gamma_1 = 0.9$, as φ increases.

3.2.3. Fréchet Distribution

– **For** $\gamma_1 = 0.1$:

In the estimation of $\gamma_1 = 0.1$, for a small percentage of censoring, $\varphi \leq 10\%$, several confidence interval estimators with the exception of POT.L provide good coverage probabilities and reasonable interval lengths. Among these estimators, Hill, MVRB, WW.L, WW.KM and ERM have CP values close to 0.95. In addition, for $\varphi \geq 0.35$, similar performance is observed as with $\varphi \leq 10$. Here, we noticed a better performance in CP values of WW.L compared with WW.KM. This is in conformity to the simulation results reported in [Worms and Worms, \(2014\)](#). Generally, the Hill, MVRB and ERM are the most appropriate for estimating $\gamma_1 = 0.1$ for various levels of censoring in the right tail.

– **For** $\gamma_1 = 0.5$:

At 10% censoring in the right tail, the ERM, POT.L and POT estimators provide good coverage probabilities. In terms of interval length, Zipf provide approximately half of the average interval lengths of the other estimators. Thus, these two estimators are the most appropriate estimators of $\gamma_1 = 0.5$. However, as the percentage of censoring in the right tail increases, the ERM, POT.L and

MOM estimators provide the best CP values. Moreover, the POT.L estimator has larger interval lengths, and hence, the ERM estimator is regarded as the most appropriate for estimating $\gamma_1 = 0.5$.

– **For** $\gamma_1 = 0.9$:

In the case of $\varphi = 10\%$, most of the estimators of γ_1 performed well with CP values close to the nominal level of 0.95 except Hill, MVRB and WW.KM. The ERM, POT.L and POT estimators consistently have CP values close to 0.95 and relatively good interval lengths. In addition, as with the case $\varphi = 10\%$, the estimators of $\gamma_1 = 0.9$ exhibited similar performance when φ was increased to 35% or 65%. Overall, ERM and POT can be used as estimators of $\gamma_1 = 0.9$ that are more robust to censoring.

3.2.4. General Comments

As may be expected, no single estimator is universally the best for estimating the EVI across distributions, size of the EVI and percentage of censoring in the right tail. However, some common underlying behaviours exist. In what follows, we present some general comments on the estimators in all the distributions considered.

In the first place, we found that the estimators' performance diminish with increasing levels of the percentage of censoring. In this regard, we noticed either a decline in the values of the coverage probability or a wider confidence interval lengths as the percentage of censoring in the right tail increases.

Secondly, most estimators exhibit large bias when estimating small values of γ_1 , especially in the Burr and Pareto distributions. However, the proposed ERM estimator is an exception to this as it exhibits high coverage even for the Burr distribution.

Thirdly, in the case of specific distributions, the following observations were made. In the Burr distribution, ERM and MOM are generally the best estimators of the EVI. For samples from the Fréchet distribution, ERM and MOM are universally good for estimating various sizes of the EVI and most robust to censoring whereas in the case of samples from the Pareto distribution, ERM and POT estimators of the EVI appear to be the best.

Lastly, we found the two estimators, ERM and MOM as the most appropriate for the estimation of the EVI across all the distributions. In addition, these estimators are the most robust to censoring and the size of the EVI. More importantly, the proposed ERM estimator was observed to be consistently robust for the estimation of the EVI regardless of latter's size and the percentage of censoring. Moreover, the estimation of from the exponential regression, the basis of the ERM estimator, leads to estimators of the second order parameters. These second-order parame-

ters can be used to obtain reduced-bias estimators of quantiles and exceedance probabilities.

4. Practical Application

In this section, we present an application of the estimators of the EVI discussed in the previous section to study the tails of the distribution of the survival time of AIDS patients. Data was obtained from [Venables and Ripley, \(2002\)](#) based on a study by Dr. P. J. Solomon and the Australian National Centre in HIV Epidemiology and Clinical Research.

The data consists of 2,843 patients of which 1,761 patients died while the remaining were right censored. Out of the total number of patients, 2,754 were males, of which 1,708 died and the remaining 1,046 were right censored. In this study, we consider the male patients only.

This data has been studied in the extreme value theory literature in [Einmahl et al., \(2008\)](#) and [Ndao et al., \(2014\)](#). In the former, the EVI is used to assess the tail heaviness of the right tail of the survival function, $1 - F$, and extreme quantiles are estimated to obtain an indication of how long a healthy man can survive AIDS. The latter uses survival time as a response variable with the age of the patient at diagnosis as covariate to obtain conditional EVI (or tail index) and extreme quantiles. Thus, the tails of the distribution of the survival time of male AIDS patients is studied conditional on the age at diagnosis.

Figure 1 shows the scatter plot and histogram of the Australian AIDS survival data. The scatter plot indicates that most of the males who survive longer are censored and the histogram indicates that there is a lower chance of survival after 7 years of diagnosis with AIDS.

The estimation of the EVI has been shown in the simulation to be sensitive to the value of φ . The values of φ must be reasonably moderate in the top order statistics to enable the application of the estimators of the EVI. Therefore, it is necessary in applications to assess the percentage (or proportion) of censoring in the right tail. The left panel of Figure 2 shows a plot of the proportion of censoring as a function of k . [Einmahl et al., \(2008\)](#) chose the proportion of censoring as $\varphi = 0.28$ and justified the selection as corresponding to the most stable part of the graph i.e. $60 \leq k \leq 200$. However, owing to the sensitivity of the estimators of γ_1 to φ , we compute our estimates using the actual φ values in the data.

From the conclusions drawn from the simulation study and in order to make it less cumbersome, we selected five estimators for illustration. These estimators are ERM, POT, MOM, WW.KM and Hill. The estimators of the EVI, γ_1 , are presented in the right panel of Figure 2. As with the UH estimator used in [Einmahl et al., \(2008\)](#), the estimators of γ_1 are relatively constant for $k \geq 200$.

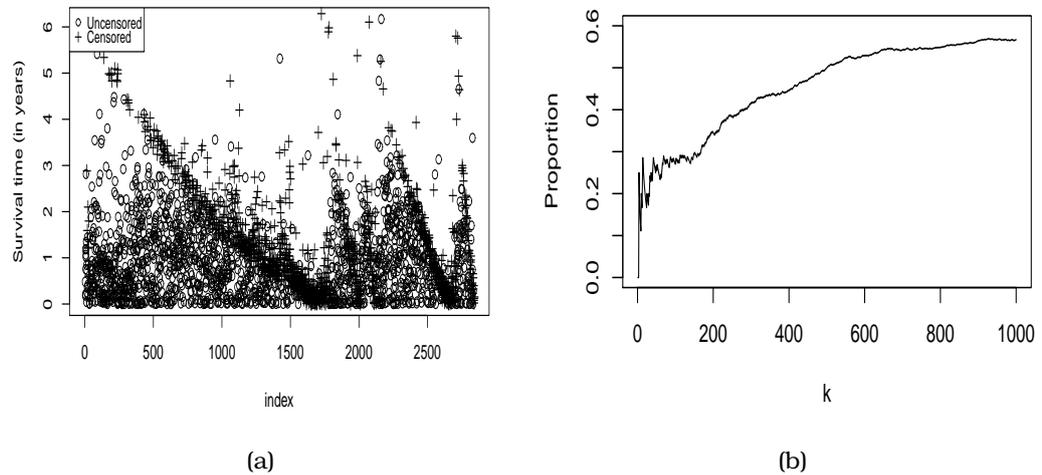


Fig. 1. (a) Survival time of AIDS patients. (b) Estimates of the φ .

Also, in practice, when a set of EVI estimators are to be taken into account, [Henriques-Rodrigues et al., \(2011\)](#) provide a simple heuristic approach to aid in selecting an appropriate threshold.

We follow a modification of the heuristic approach of selecting an optimal k instead of a percentage of the sample size as used in Section 3. Let $\gamma_1^{(i)}$, $i \in \Omega$ be the list of estimators under consideration where $\Omega = \{\text{Hill, WW.KM, ERM, MOM, POT}\}$. The optimum value of k , is chosen as

$$k_{\text{opt}} = \underset{k}{\operatorname{argmin}} \sqrt{\sum_{(i,j) \in \Omega, i \neq j} \left(\hat{\gamma}_1^{(i)} - \hat{\gamma}_1^{(j)} \right)^2}. \quad (25)$$

We apply (25) to the EVI estimators for the AIDS survival data and the results are presented in the right panel of 2. A closer look at the graph shows a stable region between 200 and 600: we choose $k_{\text{opt}} = 339$ (which is equal to 12% of the sample size and close to the 10% used in the simulation study) for the estimation of γ_1 .

The EVI estimates at $k_{\text{opt}} = 339$ are shown in Table 2. In [Einmahl et al., \(2008\)](#), only the generalised Hill estimator, $\hat{\gamma}_1^{(c,UH)}$ was used for the estimation of the EVI. The estimate of $\hat{\gamma}_1^{(c,UH)}$ was found to be 0.14. In addition, [Ndao et al., \(2014\)](#) estimates γ_1 as 0.304, 0.340 and 0.323 for males diagnosed with AIDS at ages 27, 37 and 47 years respectively.

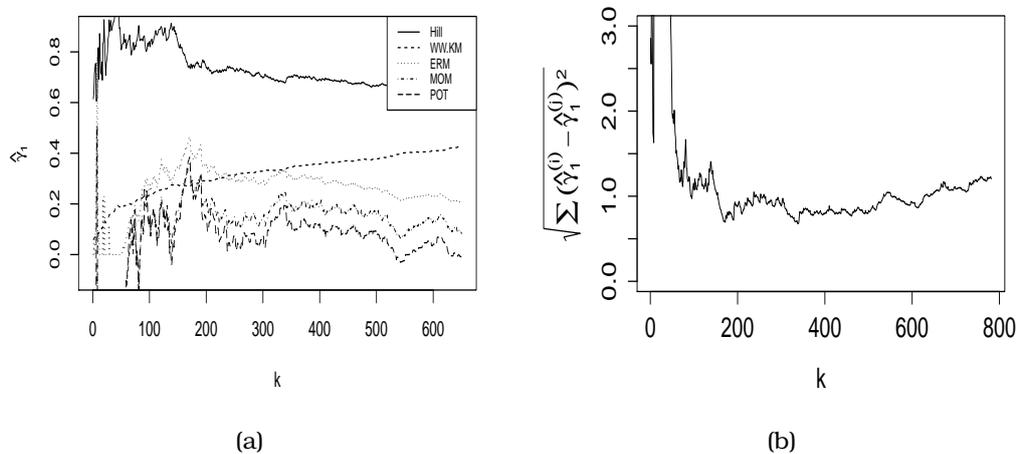


Fig. 2. Estimates of γ_1 , left panel; Heuristic choice of the threshold k , right panel

Therefore, with the exception of the Zipf estimator, all the other estimators considered give estimates within the range of the values provided by Einmahl *et al.*, (2008) and Ndao *et al.*, (2014). In particular, our ERM estimator of γ_1 and the WW.KM give estimates close to that of Ndao *et al.*, (2014), although age was not considered as a factor. Moreover, the ERM estimator is quite stable for most part of the values of k .

Table 2. Estimates of the EVI and the corresponding extreme quantile at k_{opt}

Estimator	EVI				
	WW.KM	Zipf	MOM	POT	ERM
Estimate	0.334	0.587	0.244	0.193	0.334

5. Conclusions

This paper reviews various estimators of extreme value index when observations are subject to right random censoring. In addition, an estimator based on exponential regression model was proposed. Since the asymptotic distributions are not known for all the estimators, theoretical comparison was not possible. Therefore, a simulation study was conducted to compare the performance of the various estimators under different distributions, size of the EVI and percentage of censoring in the right tail.

The performance criterion used were bias, MAD, confidence interval length and coverage probability. The simulation results show that the performance of

the estimators differ, depending on: the underlying distribution; EVI size; and percentage of censoring in the right tail. Therefore, no estimator was shown to be universally the best across all these scenarios.

However, certain estimators perform reasonably well across most distributions. These are the estimators that we recommend as appropriate for the estimation of the EVI. In this regard, if a practitioner is interested in estimators that perform well across distributions in the sense of having good coverage and small interval size, then we recommend the proposed ERM and MOM estimators. The estimators that performed well in the simulation study were illustrated using real data on the survival of AIDS patients.

Generally, we recommend that practitioners should assess the distribution of a dataset, size of γ_1 and proportion of censoring using other external information. This includes graphical plots to assist in knowing the tail behaviour of the underlying distribution and plot of the proportion of censoring at different values of k . In addition, several estimators can be used to compute estimates of γ_1 to assess the possible size of γ_1 , and hence, the selection of an appropriate estimator. We believe that the findings from this simulation will help practitioners in the selection of estimators of EVI when data is subject to right random censoring.

References

- Ameraoui, A., Boukhetala, K., and Dupuy, J.-F. (2016). Bayesian estimation of the tail index of a heavy tailed distribution under random censoring. *Computational Statistics & Data Analysis*, 104:148–168.
- Balkema, A. A. and de Haan, L. (1974). Residual life time at great age. *Annals of Probability*, 2(5):792–804.
- Beirlant, J., Bardoutsos, A., de Wet, T., and Gijbels, I. (2016). Bias reduced tail estimation for censored Pareto type distributions. *Statistics and Probability Letters*, 109:78–88.
- Beirlant, J., Dierckx, G., Goegebeur, Y., and Matthys, G. (1999). Tail index estimation and an exponential regression model. *Extremes*, 2:177–200.
- Beirlant, J., Goegebeur, Y., Segers, J., and Teugels, J. L. (2004). *Statistics of extremes: Theory and applications*. Wiley, England.
- Beirlant, J. and Guillou, A. (2001). Pareto index estimation under moderate right censoring. *Scandinavian Actuarial Journal*, 2:111–125.
- Beirlant, J., Guillou, A., Dierckx, G., and Fils-Villetard, A. (2007). Estimation of the extreme value index and extreme quantiles under random censoring. *Extremes*, 10:151–174.
- Beirlant, J., Guillou, A., and Toulemonde, G. (2010). Peaks-Over-Threshold modeling under random censoring. *Communications in Statistics - Theory and Methods*, 39(7):1158–1179.
- Beirlant, J., Maribe, G., and Verster, A. (2017). Penalized bias reduction in extreme value estimation for censored Pareto-type data , and long-tailed insurance

- applications. *arXiv1705.0663v1*, pages 1–21.
- Beirlant, J., Vynckier, P., and Teugels, J. L. (1996). Excess functions and estimation of the extreme-value index. *Bernoulli*, 2(4):293–318.
- Brahimi, B., Meraghi, D., and Necir, A. (2013). On the asymptotic normality of Hill's estimator of the tail index under random censoring. *arXiv*, pages 1–11.
- Caeiro, F., Gomes, M. I., and Pestana, D. (2005). Direct reduction of bias of the classical Hill estimator. *REVSTAT*, 3(2):113–136.
- Coles, S. (2001). *An introduction to statistical modelling of extreme values*. Springer, London.
- Csörgő, S. and Viharos, L. (1998). Estimating the tail index. In Szyszkowicz, B., editor, *Asymptotic methods in probability and statistics*, pages 833–881. North Holland.
- de Haan, L. (1970). *On regular variation and its application to the weak convergence of sample extremes*. PhD thesis, University of Amsterdam.
- de Haan, L. (1984). Slow Variation and Characterization of Domains of Attraction. In Reidel, D., editor, *Statistical Extremes and Applications*, pages 31–48.
- de Haan, L. and Ferreira, A. (2006). *Extreme value theory: An introduction*. Springer, New York, NY.
- de Haan, L. and Peng, L. (1998). Comparison of tail index estimators. *Statistica Neerlandica*, 52:60–70.
- Deheuvels, P., de Haan, L., Peng, L., and Pereira, T. T. (1997). Comparison of extreme value index estimators. Technical Report T400:EUR-09, The Erasmus University Rotterdam, Rotterdam.
- Dekkers, A. L. M., Einmahl, J. H. J., and de Haan, L. (1989). A moment estimator for the index of an extreme-value distribution. *Annals of Statistics*, 17(4):1833–1855.
- Delafosse, E. and Guillo, A. (2002). Almost sure convergence of a tail index estimator in the presence of censoring. *Comptes Rendus Mathématique*, 335(4):375–380.
- Diop, A. and Lo, G. S. (2009). Ratio of Generalized Hill's Estimator and its asymptotic normality theory. *Math. Method. Statist.*, 18(2): 117–133.
- Efron, B. (1981). Censored data and the bootstrap. *Journal of the American Statistical Association*, 76(374):312–319.
- Efron, B. and Tibshirani, R. J. (1993). *An introduction to the bootstrap*. Chapman and Hall, London.
- Einmahl, J. H. J., Elie, A. M., and Guillo, A. (2008). Statistics of extremes under random censoring. *Bernoulli*, 14(1):207–227.
- Embrechts, P., Klüppelberg, C., and Mikosch, T. (1997). *Modelling extremal events: For insurance and finance*. Springer, Berlin, Heidelberg.
- Fisher, R. and Tippett, L. (1928). On the estimation of the frequency distributions of the largest or smallest member of a sample. *Proceedings of the Cambridge Philosophical Society*, 24:80–190.
- Gnedenko, B. (1943). Sur la distribution limite du terme maximum d'une série aléatoire. *Annals of Mathematics*, 44(3):423–453.
- Gomes, M. I., Bloco, C., and Grande, C. (2010). A note on statistics of extremes for censoring schemes on a heavy right tail. In *International Conference on Informa-*

- tion Technology Interfaces*, pages 539–544, Cavtat.
- Gomes, M. I., Luísa, C. e. C., Fraga Alves, M. I., and Pestana, D. (2008). Statistics of extremes for IID data and breakthroughs in the estimation of the extreme value index: Laurens de Haan leading contributions. *Extremes*, 11(1):3–34.
- Gomes, M. I. and Neves, C. (2011). Estimation of the extreme value index for randomly censored data. *Biometrical Letters*, 48(1):1–22.
- Hall, P. (1982). On some simple estimates of an exponent of regular variation. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 44(1):37–42.
- Henriques-Rodrigues, L., Gomes, M. I., and Pestana, D. (2011). Statistics of extremes in athletics. *REVSTAT*, 9(2):127–153.
- Hill, B. (1975). A simple general approach to inference about the tail of a distribution. *Annals of Statistics*, 3:1163–1174.
- Hosking, J. R. M. and Wallis, J. R. (1987). Parameter and quantile estimation for the generalized Pareto distribution. *Technometrics*, 29:339–349.
- Kaplan, E. L. and Meier, P. (1958). Non-parametric estimation from incomplete observations. *Journal of American Statistical Association*, 53(282):457–481.
- Koul, H., Susarla, V., and Van Ryzin, J. (1981). Regression analysis with randomly right-censored data. *Annals of Statistics*, 9(6):1276–1288.
- Lo G.S., Ngom M., Kpanzou T.A., Diallo M.(2018) Weak Convergence (IIA) - Functional and Random Aspects of the Univariate Extreme Value Theory. ArXiv:1810.01625
- Leurgans, S. (1987). Linear models, random censoring and synthetic data. *Biometrika*, 74:301–309.
- Mason, D. M. . (1982). Laws of large numbers for sums of extreme values. *Annals of Probability*, 10:754–764.
- Matthys, G. and Beirlant, J. (2003). Estimating the extreme value index and high quantiles with exponential regression models. *Statistica Sinica*, 13:853–880.
- Ndao, P., Diop, A., and Dupuy, J. F. (2014). Nonparametric estimation of the conditional tail index and extreme quantiles under random censoring. *Computational Statistics and Data Analysis*, 79:63–79.
- Ngom, M., Lo, G. S. (2016) A Double-indexed Functional Hill Process and Applications. *Journal of Mathematics Research*: 8(4): 144-165.
- Pickands III, J. (1975). Statistical inference using extreme order statistics. *Annals of Statistics*, 3(1):119–131.
- Reiss, R.-D. and Thomas, M. (2007). *Statistical analysis of extreme values*. Birkhäuser, Basel, 2nd edition.
- Ripley, B. D. and Solomon, P. J. (1994). A note on Australian AIDS survival. Technical Report 94/3, Department of Statistics, University of Adelaide, Adelaide.
- Smith, R. L. (1987). Estimating tails of probability distributions. *Annals of Statistics*, 15(3):1174–1207.
- Smith, R. L. (1989). Extreme value analysis of environmental time series: An application to trend detection in ground-level ozone. *Statistical Science*, 4(4):367–377.
- Venables, W. N. and Ripley, B. D. (2002). *Modern Applied Statistics with S*. Springer, New York, NY, 4 edition.
- Worms, J. and Worms, R. (2014). New estimators of the extreme value index under random right censoring, for heavy-tailed distributions. *Extremes*, 17(2):337–358.

