

Statistica Sinica Preprint No: SS-13-272R3

Title	Adaptive and minimax optimal estimation of the tail coefficient
Manuscript ID	SS-13-272R3
URL	http://www.stat.sinica.edu.tw/statistica/
DOI	10.5705/ss.2013.272
Complete List of Authors	Alexandra Carpentier and Arlene K. H. Kim
Corresponding Author	Alexandra Carpentier
E-mail	a.carpentier@statslab.cam.ac.uk
Notice: Accepted version subject to English editing.	

Adaptive and minimax optimal estimation of the tail coefficient

Alexandra Carpentier and Arlene K. H. Kim

University of Cambridge

Abstract: We consider the problem of estimating the tail index α of a distribution satisfying a (α, β) second-order Pareto-type condition, where β is the second-order coefficient. When β is available, it was previously proved that α can be estimated with the optimal rate $n^{-\frac{\beta}{2\beta+1}}$. On the contrary, when β is not available, estimating α with the optimal rate is challenging; so additional assumptions that imply the estimability of β are usually made. In this paper, we propose an adaptive estimator of α , and show that this estimator attains the rate $(n/\log \log n)^{-\frac{\beta}{2\beta+1}}$ without a priori knowledge of β and any additional assumptions. Moreover, we prove that this $(\log \log n)^{\frac{\beta}{2\beta+1}}$ factor is unavoidable by obtaining the companion lower bound.

Key words and phrases: Adaptive estimation, minimax optimal bounds, extreme value index, Pareto-type distributions.

1. Introduction

We consider the problem of estimating the tail index α of an (α, β) second-order Pareto distribution F , given n i.i.d. observations X_1, \dots, X_n . More precisely, we assume that for some $\alpha, \beta, C, C' > 0$,

$$|1 - F(x) - Cx^{-\alpha}| \leq C'x^{-\alpha(1+\beta)}. \quad (1.1)$$

We will write $\mathcal{S}(\alpha, \beta) := \mathcal{S}(\alpha, \beta, C, C')$ for the set of distributions that satisfy this property (see Definition (2)). Here the tail index α characterizes the heaviness of the tail, and β represents the proximity between F and an α -Pareto distribution $F_\alpha^P : x \in [C^{1/\alpha}, \infty) \rightarrow 1 - Cx^{-\alpha}$.

There is an abundant literature on the problem of estimating α . A very popular estimator is Hill's estimator (Hill, 1975) (see also Pickands' estimator (Pickands, 1975)). Hill (1975) considered α -Pareto distribution for the tail, and suggested an estimator $\hat{\alpha}_H(r)$ of the tail index α based on the order statistics $X_{(1)} \leq \dots \leq X_{(n)}$ where r is the fraction of order statistics from the tail,

$$\hat{\alpha}_H(r) = \left(\frac{1}{[rn]} \sum_{i=1}^{[rn]} \frac{\log(X_{(n-i+1)})}{\log(X_{(n-[rn]+1)})} \right)^{-1}. \quad (1.2)$$

For more details, see e.g. de Haan and Ferreira (2006).

Limiting distribution of Hill's estimator was first proved by Hall (1982) when β is known. Under a model that is quite similar to (1.1), he proved that if $rn^{1/(2\beta+1)} \rightarrow 0$ as $n \rightarrow \infty$, $\sqrt{nr}(\hat{\alpha}_H(r) - \alpha)$

converges in distribution to $N(0, \alpha^2)$. He also considered more restricted condition, say, the exact Hall condition,

$$|1 - F(x) - Cx^{-\alpha}| = C''x^{-\alpha(1+\beta)} + o(x^{-\alpha(1+\beta)}). \quad (1.3)$$

Under the model (1.3) with the choice of the sample fraction $r^* = Cn^{-\frac{1}{2\beta+1}}$ with some constant C , Theorem 2 of Hall (1982) states that $n^{\beta/(2\beta+1)}(\hat{\alpha}_H(r^*) - \alpha)$ converges to a Gaussian distribution with finite mean and variance, depending on the parameters of the true distribution.

The companion lower bound $n^{-\beta/(2\beta+1)}$ under the assumption (1.1) was proved by Hall and Welsh (1984). Drees (2001) improved this result by obtaining sharp asymptotic minimax bounds again when β is available. From these results, we know that the second-order parameter β is crucial to understand the behaviour of the distribution. Indeed, it determines the rate of estimation of α as well as the optimal sample fraction.

However, β is unknown in general. To cope with this problem, Hall and Welsh (1985) proved that under condition (1.3), it is possible to estimate β in a consistent way, and thus also to estimate the sample fraction r^* consistently by \hat{r} (see Theorem 4.2 in their paper). Theorem 4.1 of Hall and Welsh (1985) deduces from these results that the estimate $\hat{\alpha}_H(\hat{r})$ is asymptotically as efficient as $\hat{\alpha}_H(r^*)$, that is, $n^{\beta/(2\beta+1)}(\hat{\alpha}_H(\hat{r}) - \alpha)$ converges to a Gaussian distribution with the same mean and variance as the one resulting from the choice r^* . Their result is pointwise, but not uniform under the model (1.3), as opposed to the uniform convergence when β is known.

This first result on adaptive estimation was extended in several ways. For instance, Gomes, et. al. (2008) provided more precise ways to reduce the bias of the estimate of α using the estimate of β by supposing the third order condition. The adaptive estimates of α under the third order condition was considered in Gomes, et. al. (2012). In addition, several other methods for estimating r^* have been proposed, e.g. bootstrap (e.g. Danielsson, et. al. (2001)) or regression (e.g. Beirlant, et. al. (1996)). In particular, Drees and Kaufmann (1998) considered a method that is related to Lepski's method (see Lepski (1992) for more details in a functional estimation setting) by choosing the sample fraction that balances the squared bias and the variance of the resulting estimate. They proved that Hill's estimate computed with this sample fraction is asymptotically as efficient as the oracle estimate if F satisfies a condition that is slightly more restrictive than the condition (1.3). Finally, Grama and Spokoiny (2008) consider a more general setting than (1.1). However, when they apply their results to the exact Hall model (without little o), their estimator obtains the optimal rate up to a $\log(n)$ factor, which is clearly sub-optimal as proven in Hall and Welsh (1985).

In this paper, we focus on deriving results for the setting (1.1). Indeed, many common dis-

tributions (in particular some distributions with change points in the tail) belong to it, and the construction of the lower bound in Hall and Welsh (1984) was proved in this model. However, to the best of our knowledge, either the existing results that we mentioned previously hold in a more restrictive setting than the model (1.1), typically in a model that is close to the model (1.3) (see e.g. Hall and Welsh (1985); Beirlant, et. al. (1996); Drees and Kaufmann (1998); Danielsson, et. al. (2001); Gomes, et. al. (2008, 2012)), or the convergence rates for the setting (1.1) in the previous results are worse than one could expect (see e.g. Grama and Spokoiny (2008)). It is important to note here that the set of distributions described in Equation (1.1) is significantly larger than the set of distributions that satisfy the restricted condition (1.3). As will be explained later, the adaptive estimation in our setting (i.e. condition (1.1)) is more involved since the second-order parameter β is not always estimable (even a consistent estimator does not exist for all distributions in this model), and the adaptive procedures based on estimating β or the oracle sample fraction r^* as in the papers (Hall and Welsh (1985); Gomes, et. al. (2008, 2012)) might not work on all the functions satisfying (1.1).

The contributions of this paper are the following. We construct an adaptive estimator $\hat{\alpha}$ of α in the setting (1.1) and prove that $\hat{\alpha}$ converges to α with the rate $(n/\log\log(n))^{-\beta/(2\beta+1)}$. More precisely, for an arbitrarily small $\epsilon > 0$, and some arbitrarily large range I_1 for α and $[\beta_1, \infty)$ for β , there exist large constants $D, E > 0$ such that for any $n > D \log(\log(n)/\epsilon)$

$$\sup_{\alpha \in I_1, \beta > \beta_1} \sup_{F \in \mathcal{S}(\alpha, \beta)} \mathbb{P}_F \left(|\hat{\alpha} - \alpha| \geq E \left(\frac{n}{\log(\log(n)/\epsilon)} \right)^{-\frac{\beta}{2\beta+1}} \right) \leq \epsilon. \quad (1.4)$$

There is an additional $(\log\log(n))^{\frac{\beta}{2\beta+1}}$ factor in the rate with respect to the oracle rate, which comes from the fact that we adapt over β on a set of distributions where β is not estimable. Although we obtain worse rates of convergence than the oracle rate, we actually prove the optimality of our adaptive estimator by obtaining a matching lower bound. Indeed, there exists a small enough constant $E' > 0$ such that for any n large enough, and for any estimator $\tilde{\alpha}$,

$$\sup_{\alpha \in I_1, \beta > \beta_1} \sup_{F \in \mathcal{S}(\alpha, \beta)} \mathbb{P}_F \left(|\tilde{\alpha} - \alpha| \geq E' \left(\frac{n}{\log(\log(n))} \right)^{-\frac{\beta}{2\beta+1}} \right) \geq \frac{1}{4}.$$

Both lower and upper bounds containing the $(\log\log(n))^{\beta/(2\beta+1)}$ factor are new to the best of our knowledge (we do not provide a tight scaling factor as in the paper by Novak (2013), but the setting in this paper is different and their rate does not involve this additional $(\log\log(n))^{\beta/(2\beta+1)}$ factor). The presence of the $\log\log n$ factor is not unusual in adaptive estimation (see Spokoiny (1996) in a signal detection setting). This issue is also discussed in the paper (Drees and Kaufmann, 1998).

The adaptive estimator $\hat{\alpha}$ we propose in this paper is based on a sequence of estimates $\hat{\alpha}(k)$ defined in (3.1), where the parameter $k \in \mathbb{N}$ plays a role similar to the sample fraction in Hill's estimator (see Subsection 3.1 for more details). These estimates $\hat{\alpha}(k)$ are not based on order statistics, but on probabilities of tail events. We first prove that for an appropriate choice of this threshold k (independent of α or β), $\hat{\alpha}(k)$ is consistent. We then prove that for an oracle choice of k (as a function of β), this estimate is minimax-optimal for distributions satisfying (1.1) with the rate $n^{-\frac{\beta}{2\beta+1}}$. Finally an adaptive version of this estimate, where the parameter k is chosen in a data-driven way without knowing β in advance, is proved to satisfy Equation (1.4).

2. Definitions of distribution classes

In this section, we introduce two sets of distributions of interest, namely the class of approximately α -Pareto distributions, and the class of approximately (α, β) second-order Pareto distributions. We let \mathcal{D} be the class of distribution functions on $[0, \infty)$.

Definition 1. Let $\alpha > 0$, $C > 0$. We denote by $\mathcal{A}(\alpha, C)$ the class of approximately α -Pareto distributions:

$$\mathcal{A}(\alpha, C) = \left\{ F \in \mathcal{D} : \lim_{x \rightarrow \infty} (1 - F(x))x^\alpha = C \right\}.$$

Distributions in $\mathcal{A}(\alpha, C)$ converge to Pareto distributions for large x , and these distributions have been used as a first attempt to understand heavy tail behavior (see Hill (1975); de Haan and Ferreira (2006)). The first-order parameter α characterizes the tail behavior such that distributions with smaller α correspond to heavier tails.

In order to provide rates of convergence (of an estimator of α), we define the set of second-order Pareto distributions.

Definition 2. Let $\alpha > 0$, $C > 0$, $\beta > 0$ and $C' > 0$. We denote by $\mathcal{S}(\alpha, \beta, C, C')$ the class of approximately (α, β) second-order Pareto distributions:

$$\mathcal{S}(\alpha, \beta, C, C') = \left\{ F \in \mathcal{D} : \forall x \text{ s.t. } F(x) \in (0, 1], |1 - F(x) - Cx^{-\alpha}| \leq C'x^{-\alpha(1+\beta)} \right\}. \quad (2.1)$$

From Definition 2, we know that not only are the distributions in $\mathcal{S}(\alpha, \beta, C, C')$ approximately α -Pareto, but we additionally have a bound on the rate at which they approximate Pareto distributions. This rate of approximation is linked to the second-order parameter β —a large β corresponds to a distribution that is very close to a Pareto distribution (in particular, when $\beta = \infty$, it becomes exactly Pareto), and a small β corresponds to a distribution that is well approximated by a Pareto distribution only for a very large x . From now, if there is no confusion, we call the distributions in

$\mathcal{S}(\alpha, \beta, C, C')$ second-order Pareto distributions, and we use the notation \mathcal{A} and \mathcal{S} without writing parameters explicitly.

The condition in (2.1) is related to the condition (1.3), but is weaker. Indeed, the condition (1.3) implies

$$\lim_{x \rightarrow \infty} \frac{1 - F(x) - Cx^{-\alpha}}{x^{-\alpha(1+\beta)}} = C',$$

whereas our condition imposes only an upper bound,

$$\limsup_{x \rightarrow \infty} \left| \frac{1 - F(x) - Cx^{-\alpha}}{x^{-\alpha(1+\beta)}} \right| \leq C'.$$

This difference is essential in the estimation problem. For instance, in the setting (1.3), it is possible to estimate β consistently (see e.g. Hall and Welsh (1985)), whereas in our setting (2.1), it is not possible to estimate β consistently over the set \mathcal{S} of distributions for $\beta \in [\beta_1, \beta_2]$ with $0 < \beta_1 < \beta_2$. Adaptive estimation of α is thus likely to be more involved in our setting than in the more restricted model (1.3). For instance, many adaptive techniques rely on estimating β or the sample fraction as a function of β , which is not directly applicable in our setting (see e.g. Hall and Welsh (1985); Danielsson, et. al. (2001); Gomes, et. al. (2012)).

Remark 1. *The difference between the functions satisfying the condition in Definition 2 and the condition (1.3) is related to the difference between Hölder functions that actually attain their Hölder exponent and Hölder functions that are in a given Hölder ball but do not attain their Hölder exponent (see e.g. Giné and Nickl (2010) for a comparison of these two sets, and the problem for estimation when the second set is considered).*

3. Main results

Most estimates in the literature are based on order statistics (as Hill's estimate or Pickands' estimate), which causes a difficulty for one to analyse them in a non-asymptotic way. In contrast, the estimate we will present in Section 3.1 verifies large deviation inequalities in a simple way. This estimate is based on probabilities of well chosen tail events.

3.1. A new estimate

Let X_1, \dots, X_n be an i.i.d. random sample from a distribution $F \in \mathcal{A}$. We write, for any $k \in \mathbb{N}$,

$$p_k := \mathbb{P}(X > e^k) = 1 - F(e^k),$$

and its empirical estimate

$$\hat{p}_k := \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i > e^k\}.$$

We define the following estimate of α for any $k \in \mathbb{N}$

$$\hat{\alpha}(k) := \log(\hat{p}_k) - \log(\hat{p}_{k+1}). \quad (3.1)$$

This estimate gives the following large deviation inequalities, which is crucial for proving consistency and convergence rates of $\hat{\alpha}(k)$.

Lemma 1 (Large deviation inequality). *Let X_1, \dots, X_n be an i.i.d. sample from F .*

A. *Suppose $F \in \mathcal{A}$ and let $\delta > 0$. For any k such that $p_{k+1} \geq \frac{16 \log(2/\delta)}{n}$, with probability larger than $1 - 2\delta$,*

$$|\hat{\alpha}(k) - (\log(p_k) - \log(p_{k+1}))| \leq 6\sqrt{\frac{\log(2/\delta)}{np_{k+1}}}. \quad (3.2)$$

B. *Assume now that $F \in \mathcal{S}$ and let $\delta > 0$. For any k such that $p_{k+1} \geq \frac{16 \log(2/\delta)}{n}$ and $e^{-k\alpha\beta} \leq C/(2C')$, with probability larger than $1 - 2\delta$,*

$$|\hat{\alpha}(k) - \alpha| \leq 6\sqrt{\frac{\log(2/\delta)}{np_{k+1}}} + \frac{3C'}{C}e^{-k\alpha\beta} \quad (3.3)$$

$$\leq 6\sqrt{\frac{e^{(k+1)\alpha+1} \log(2/\delta)}{Cn}} + \frac{3C'}{C}e^{-k\alpha\beta}. \quad (3.4)$$

For this new estimate $\hat{\alpha}(k)$, k plays a similar role as the sample fraction in Hill's estimate (1.2). The bias-variance trade-off should be solved by choosing k in an appropriate way as a function of β (we will explain this more in details later). Choosing a too large k leads to using a small sample fraction, and the resulting estimate has a large variance and a small bias. On the other hand, choosing a too small k yields a large bias and a small variance for the estimate. The optimal k equalises the bias term and the standard deviation.

3.2. Rates of convergence

We first consider the set of approximately Pareto distributions, and prove that the estimate $\hat{\alpha}(k_n)$ is consistent if we choose k_n such that it diverges to ∞ but not too fast.

Theorem 1 (Consistency in \mathcal{A}). *Let $F \in \mathcal{A}$. Let $k_n \in \mathbb{N}$ be such that $k_n \rightarrow \infty$ and $(\log(n)/n)e^{k_n\alpha} \rightarrow 0$ as $n \rightarrow \infty$. Then*

$$\hat{\alpha}(k_n) \rightarrow \alpha \text{ a.s.}$$

Choosing (for instance) $k_n = (\log \log(n))$ ensures almost sure convergence.

The estimate $\hat{\alpha}(\log \log(n))$ converges to α almost surely under the rather weak assumption that F belongs to \mathcal{A} . But on such sets, no uniform rate of convergence exists, and this is the reason why the restricted set \mathcal{S} is introduced.

Let $\alpha, \beta, C, C' > 0$. Consider now the set $\mathcal{S} := \mathcal{S}(\alpha, \beta, C, C')$ of second-order Pareto distributions. We assume in a first instance that, although we do not have access to α , we know the parameter $\alpha(2\beta+1)$. It is not very realistic assumption, but we will explain soon how we can modify the estimate so that it is minimax optimal on the class of second-order Pareto distributions.

Theorem 2 (Rate of convergence when $\alpha(2\beta+1)$ is known). *Let n be such that (4.7) is satisfied. Let $k_n^* = \lfloor \log(n^{\frac{1}{\alpha(2\beta+1)}}) + 1 \rfloor$. Then for any $\delta > 0$, we have*

$$\sup_{F \in \mathcal{S}} \mathbb{P}_F \left(|\hat{\alpha}(k_n^*) - \alpha| \geq \left(B_1 + \frac{3C'}{C} \right) n^{-\frac{\beta}{2\beta+1}} \right) \leq 2\delta,$$

where $B_1 = 6\sqrt{e^{2\alpha+1} \frac{\log(2/\delta)}{C}}$.

Theorem 2 states that, uniformly on the class of second-order Pareto distributions, the estimate $\hat{\alpha}(k_n^*)$ converges to α with the minimax optimal rate $n^{-\frac{\beta}{2\beta+1}}$ (see Hall and Welsh (1984) for the matching lower bound).

Remark 2. *Theorem 2 can be used to prove the convergence rate of our estimator by modifying the choice of k_n^* , when $\alpha(2\beta+1)$ is unknown but only β is known. For instance, we can plug a rough estimate $\tilde{\alpha} := \hat{\alpha}((\log \log(n))^2)$ of α into k_n^* . The idea behind this choice is that with sufficiently large n , we have with high probability,*

$$|\hat{\alpha}((\log \log(n))^2) - \alpha| = O\left(\frac{1}{\log n}\right).$$

Then \hat{k}_n^1 is defined as $\lfloor \log(n^{\frac{1}{\tilde{\alpha}(2\beta+1)}}) + 1 \rfloor$. Finally, the rate of convergence of $\hat{\alpha}(\hat{k}_n^1)$ can be shown as $n^{-\beta/(2\beta+1)}$ by proving $\exp(\hat{k}_n^1) = O(n^{1/(\alpha(2\beta+1))})$ with high probability.

However, the previous optimal choice of k (k_n^* or \hat{k}_n^1) still depends on β , which is unavailable in general. To deal with this problem, we construct an adaptive estimate of α that does not depend on β but still attains a rate that is quite close to the minimax optimal rate $n^{-\frac{\beta}{2\beta+1}}$ on the class of β second-order Pareto distributions.

The adaptive estimator is obtained by considering a kind of bias and variance trade-off based on the large deviation inequality (3.2). Suppose we know the optimal choice of k^* . Then this k^* will optimize the squared error by making bias and standard error (of the estimate with respect to its expectation) equal. Since the bias is decreasing while the standard error is increasing as k

increases, for all k' larger than this optimal k^* , the bias will be smaller than the standard error. Based on this heuristic (originally proposed by Lepski (1992)), we pick the smallest k which satisfies for all k' larger than k , the proxy for the bias is smaller than the proxy for the standard error $O(\sqrt{1/(n\hat{p}_{k'+1})})$ as in (3.2). For the proxy for the bias, we use $|\hat{\alpha}(k') - \hat{\alpha}(k)|$ by treating $\hat{\alpha}(k)$ as the true α based on the idea that $\hat{\alpha}(k)$ would be very close in terms of the rate to the true α (if k is selected in an optimal way).

More precisely, we choose k as follows, for $1/4 > \delta > 0$

$$\hat{k}_n = \inf \left\{ k \in \mathbb{N} : \hat{p}_{k+1} > \frac{24 \log(2/\delta)}{n} \text{ and } \forall k' > k \text{ s.t. } \hat{p}_{k'+1} > \frac{24 \log(2/\delta)}{n}, |\hat{\alpha}(k') - \hat{\alpha}(k)| \leq A(\delta) \sqrt{\frac{1}{n\hat{p}_{k'+1}}} \right\}, \quad (3.5)$$

where $A(\delta)$ satisfies the condition (3.6) in the following theorem.

Theorem 3 (Rates of convergence with unknown β). *Let $1/4 > \delta > 0$ and let n be such that (4.9) is satisfied. Consider the adaptive estimator $\hat{\alpha}(k_n)$ where k_n is chosen as described in (3.5) where $A(\delta)$ satisfies the following condition*

$$A(\delta) \geq 6\sqrt{2(C + C') \log(2/\delta)} \left(2\sqrt{\frac{e^{2\alpha+1}}{C}} + \frac{C'}{C} \right). \quad (3.6)$$

Then we have

$$\sup_{F \in \mathcal{S}} \mathbb{P}_F \left(|\hat{\alpha}(\hat{k}_n) - \alpha| \geq \left(B_2 + \frac{3C'}{C} \right) \left(\frac{n}{\log(2/\delta)} \right)^{-\frac{\beta}{2\beta+1}} \right) \leq \left(1 + \frac{1}{\alpha} \log \left(\frac{(C + C')n}{16} \right) \right) \delta.$$

where $B_2 = \left(B_1 + 2A(\delta) \sqrt{\frac{e^{2\alpha}}{C}} \right) \frac{1}{\sqrt{\log(2/\delta)}}$ and B_1 is defined in Theorem 2.

Theorem 3 holds for any (α, β) provided that n and $A(\delta)$ are larger than some constants depending on α, β, C, C' , and on the probability δ . The advantage of our adaptive estimator is that since the threshold \hat{k}_n is chosen adaptively to the samples, the second-order parameter β does not need to be known in the procedure in order to obtain the convergence rate of $\hat{\alpha}(\hat{k}_n)$. Theorem 3 gives immediately the following corollary.

Corollary 1. *Let $\epsilon \in (0, 1)$ and $C' > 0$ and let $0 < \alpha_1 < \alpha_2$ and $0 < C_1 < C_2$. We use \hat{k}_n as in*

(3.5) where $A(\delta) = A(\delta(\epsilon)) =: A(\epsilon)$ is chosen as in Equation (4.19). If n satisfies (4.21), then

$$\sup_{\substack{\alpha \in [\alpha_1, \alpha_2], \beta \in [\beta_1, \infty) \\ C \in [C_1, C_2]}} \sup_{F \in \mathcal{S}(\alpha, \beta, C, C')} \mathbb{P}_F \left(|\hat{\alpha}(\hat{k}_n) - \alpha| \geq B_3 \left(\frac{n}{\log \left(\frac{2}{\epsilon} \left(1 + \frac{\log((C_2 + C')n)}{\alpha_1} \right) \right)} \right)^{-\frac{\beta}{2\beta+1}} \right) \leq \epsilon,$$

where B_3 is a constant explicitly expressed in (4.20), which only depends on α_2, C_1, C_2 , and C' .

In other words, if we fix the range of the α and C and a lower bound on β to which we wish to adapt, we can tune the parameters of the adaptive choice of \hat{k}_n so that we adapt to the maximal β such that F is β second-order Pareto. Moreover, this adaptive procedure works uniformly well over the set of second-order Pareto distributions satisfying (1.1) (for $\alpha \in [\alpha_1, \alpha_2], \beta \in [\beta_1, \infty], C \in [C_1, C_2]$), which is much larger than the class of distributions that verify the condition (1.3). Then this gives *non-asymptotic guarantees with explicit bounds*.

Remark 3. The parameter C' plays a role in the definition of the second order Pareto class that is slightly different than the one of C or α, β . Unlike α or C , C' is not uniquely defined: if $F \in \mathcal{S}(\alpha, \beta, C, \tilde{C}')$, then $F \in \mathcal{S}(\alpha, \beta, C, C')$ with $C' \geq \tilde{C}'$. This implies in particular that the results of Corollary 1 could have been rewritten, fixing a constant $C' > 0$ and writing \tilde{C}' for a constant that fits more closely F , by taking supremum over $F \in \mathcal{S}(\alpha, \beta, C, \tilde{C}')$ where $\tilde{C}' \leq C'$. Being non-adaptive over \tilde{C}' and choosing a loose constant C' instead of \tilde{C}' will only worsen the bound by a constant factor, unlike making a mistake on β which will worsen the exponent of the bound.

It seems that we lose a $(\log \log(n))^{\frac{\beta}{2\beta+1}}$ factor with respect to the optimal rate, due to adaptivity to β . However, the lower bound below implies that this $(\log \log(n))^{\frac{\beta}{2\beta+1}}$ loss is inevitable; hence the rate provided in Theorem 3 is sharp.

Theorem 4 (Lower bound). Let $\alpha_1, \beta_1, C_1, C_2, C' > 0$ be such that $C_1 \leq \exp(-\frac{1}{2\alpha_1(2\beta_1+1)})$, $C_2 \geq 1$ and $C' \geq \frac{1}{2\alpha_1\beta_1}$. Let n be sufficiently large. Then for any estimate $\tilde{\alpha}$ of α ,

$$\sup_{\substack{\alpha \in [\alpha_1, 2\alpha_1], \beta \in [\beta_1, \infty) \\ C \in [C_1, C_2]}} \sup_{F \in \mathcal{S}(\alpha, \beta, C, C')} \mathbb{P}_F \left(|\tilde{\alpha} - \alpha| \geq B_4 \left(\frac{n}{\log(\log(n)/2)} \right)^{-\frac{\beta}{2\beta+1}} \right) \geq \frac{1}{4},$$

where B_4 is a constant depending on α_1 and β_1 , which is provided in (4.30).

The lower bound result is proved with specific ranges of the parameters (e.g. restrictions on C_1, C_2, C' in the statement of Theorem 4), but it can be modified by considering different ranges

(see Remark 4).

3.3. Additional remarks on our estimate

In the definition of our estimate, we use exponential spacings (i.e. we estimate the probability that the random variable is larger than e^k), but we can generalize our estimate by considering the probability of other tail events. For some parameters $u > v \geq 1$, define

$$\hat{q}_u = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i > u\}, \text{ and } \hat{q}_v = \frac{1}{n} \sum_{i=1}^n \mathbf{1}\{X_i > v\}.$$

We define the following estimate of α as

$$\hat{\alpha}(u, v) = \frac{\log(\hat{q}_v) - \log(\hat{q}_u)}{\log(u) - \log(v)}. \quad (3.7)$$

If we fix $v \sim O(n^{1/(\alpha(2\beta+1))})$ and $u/v \sim O(1)$, then we will also obtain the oracle rate for estimating α with $\hat{\alpha}(u, v)$. However, the choice of u/v will have an impact on the constants. In practice, these parameters are important to tune well (in particular for the exact Pareto case, or for distributions satisfying Equation (1.3)). However, a precise analysis of the best choices for u and v (in terms of constants) is beyond the scope of this paper.

Another point we want to address is the relation between our estimate and usual estimates based on order statistics. To estimate the tail index α , it is natural to consider the quantiles associated with the tail probabilities. For the estimates based on order statistics, one fixes some tail-probabilities and then observes the order statistics in order to estimate the quantiles. On the other hand, we fix some values corresponding to the quantiles, and estimate the associated tail probabilities. Based on such a link, one could relate any existing method based on order statistics to the method based on tail probabilities.

In particular, the estimator based on order statistics corresponding to our estimator would be of the form, for some parameters $1 \geq q_v > q_u \geq 0$,

$$\tilde{\alpha}(q_u, q_v) = \frac{\log(q_v) - \log(q_u)}{\log(\hat{u}) - \log(\hat{v})}, \quad (3.8)$$

where $\hat{u} = X_{(n-[q_u n])}$ and $\hat{v} = X_{(n-[q_v n])}$. This estimate can be interpreted as the inverse of some generalized Pickands' estimate (see Pickands (1975), it is however *not* Pickands' estimate). There is actually a duality between these two estimators: for any couple (q_u, q_v) in the definition (3.8), it is possible to find (u, v) in the definition (3.7) such that these two estimates exactly match (see Figure 3.1 for an illustration). However, there is no analytical transformation from one estimate to the other since such a transformation will be data dependent.

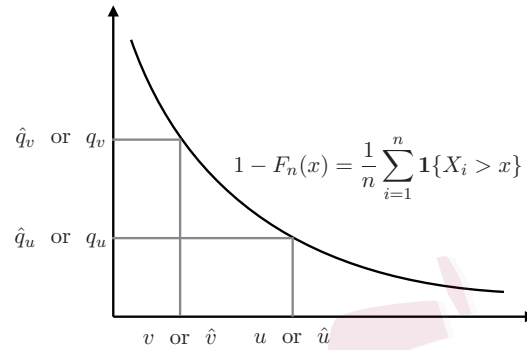


FIG 3.1. Duality between the estimate (3.7) and the estimate (3.8).

Acknowledgment

We are grateful to Richard J. Samworth and Richard Nickl for their comments and advice. We are also grateful to the anonymous Referees, the Associate Editor and the Editor for very insightful comments that were helpful in enhancing the quality of the paper.

References

- Beirlant, J., and Vynckier, P., and Teugels, J. (1996). Tail index estimation, Pareto quantile plots and regression. *Journal of American Statistical Association*, **70**, 1659–1667.
- Cover, T. M. and Thomas, J. A. (2012). *Elements of Information Theory*. Wiley-interscience.
- Danielsson, J. and de Haan, L. and Peng, L. and de Vries, C.G. (2001). Using a Bootstrap Method to Choose the Sample Fraction in Tail Index Estimation. *Journal of Multivariate Analysis*, **2**, 226–248
- Drees, H. (2001). Minimax risk bounds in extreme value theory. *The Annals of Statistics*, **29**, (1) 266–294.
- Drees, H. and Kaufmann, E. (1998). Selecting the optimal sample fraction in univariate extreme value estimation. *Stochastic Processes and Their Applications*, **75**, 149–172.
- Giné, E. and Nickl, R. (2010). Confidence bands in density estimation. *The Annals of Statistics*, **38**, (2) 1122–1170.
- Gomes, M. I. and F. Figueiredo, and Neves, M. (2012). Adaptive estimation of heavy right tails: resampling-based methods in action. *Extremes*, **15**. 463–489
- Gomes, I. M and De Haan, L., and Rodrigues, Ligia Henriques (2008). Tail index estimation for

- heavy-tailed models: accommodation of bias in weighted log-excesses. *Journal of the Royal Statistical Society: Series B*, **91**, 31–52.
- Grams, I. and Spokoiny, V. (2008). Statistics of extremes by oracle estimation. *The Annals of Statistics*, **36**, (4) 1619–1648.
- de Haan, L. and Ferreira, A. (2006). *Extreme Value Theory: An Introduction*. Springer series in operations research und financial engineering.
- Hall, P. (1982). On some simple estimates of an exponent of regular variation. *Journal of the Royal Statistical Society: Series B.*, **44**, (1) 37–42.
- Hall, P. and Welsh, A. H. (1984). Best attainable rates of convergence for estimates of parameters of regular variation. *The Annals of Statistics*, **12**, (3) 1079–1084.
- Hall, P. and Welsh, A. H. (1985). Adaptive estimates of parameters of regular variation. *The Annals of Statistics*, **75**, (1) 331–341.
- Hill, B. M. (1975). A simple general approach to inference about the tail of a distribution. *The Annals of Statistics*, **3**, (5) 1163–1174.
- Lepski, O. V. (1992). On problems of adaptive estimation in white gaussian noise. *Topics in nonparametric estimation*, **12**, 87–106.
- Novak, S. Y. (2013). Lower bounds to the accuracy of inference on heavy tails. *Bernoulli*(to appear)
- Pickands, J. (1975). Statistical inference using extreme order statistics. *The Annals of Statistics*, 119–131.
- Spokoiny, V. G. (1996). Adaptive hypothesis testing using wavelets. *The Annals of Statistics*, **24**, (6) 2477–2498.
- Tsybakov, A. B.(2008). *Introduction to Nonparametric Estimation*. Springer.
- Van der Vaart, A. W.(2000). *Asymptotic Statistics*. Cambridge University Press.

University of Cambridge

E-mail: a.carpentier@statslab.cam.ac.uk

University of Cambridge

E-mail: a.kim@statslab.cam.ac.uk