# STATISTICS OF EXTREMES IN ATHLETICS

Authors:    Lígia Henriques-Rodrigues
– Instituto Politécnico de Tomar and C.E.A.U.L., Portugal
Ligia.Rodrigues@aim.estt.ipt.pt

M. Ivette Gomes
– Universidade de Lisboa, F.C.U.L. (D.E.I.O.) and C.E.A.U.L., Portugal
ivette.gomes@fc.ul.pt

Dinis Pestana
– Universidade de Lisboa, F.C.U.L. (D.E.I.O.) and C.E.A.U.L., Portugal
dinis.pestana@fc.ul.pt

Abstract:

• TV shows on any athletic event make clear that those who want *gold medals* cannot dispense *statistics*. And the statistics more appealing to champions and coachers are the *extreme order statistics*, and in particular *maximum* (or *minimum*) *values* and *records*. The models in *statistics of extremes* are usually semi-parametric or even non-parametric in nature, with the imposition of a few *regularity conditions* in the appropriate tail of the unknown model underlying the available data. The primordial parameter is the *extreme value index*, the shape parameter in the (unified) *extreme value* distribution. The estimation of the *extreme value index* is one of the basis for the estimation of other parameters of rare events, like the *right endpoint* of the model underlying the data, a *high quantile*, the *return period* and the *probability of exceedance* of a high level. In this paper, we are interested in an application of *statistics of extremes* to the best personal marks in a few athletic events. Due to the way data are collected, we begin with a parametric data analysis, but we pay special attention to the semi-parametric estimation of the extreme value index and the right endpoint whenever finite, the possible *world record*, given the actual conditions. In order to achieve a better decision we consider a few alternative semi-parametric estimators available in the literature, and heuristic rules for the choice of thresholds.

## 1. INTRODUCTION AND OUTLINE OF THE PAPER

Statistical facts are quite commonly used by sports commentators. We all have listened to programs on different athletic events, showing that *statistics* is an instrument that champions instructors need to use. It is without doubt a subject which cannot be dispensed by those who want *gold medals*, and the statistics more appealing to the champions are the *extreme order statistics* (o.s.'s), and in particular *maximum* (or *minimum*) *values* and *records*.

The models in *statistics of extremes* are usually semi-parametric or even non-parametric in nature, with the imposition of a few "regularity conditions" in the right-tail (or left-tail), $\overline{F}(x) := 1 - F(x)$, as $x \to +\infty$ $\big($or $F(x)$, as $x \to -\infty\big)$, of an unknown model $F$ underlying the available data, whenever we are interested in large (or small) values. The primordial parameter is the *extreme value index*. For large values, the extreme value index is the shape parameter $\gamma$ in the distribution function (d.f.)

$$(1.1) \qquad G_\gamma(x) = \begin{cases} \exp\big(-(1+\gamma x)^{-1/\gamma}\big), & 1+\gamma x > 0 , & \text{if } \gamma \neq 0 , \\ \exp\big(-\exp(-x)\big), & x \in \mathbb{R} , & \text{if } \gamma = 0 , \end{cases}$$

the (unified) *extreme value* distribution. The extreme value index needs to be estimated in a "precise" way, because such an estimation plays a major role in the estimation of other parameters of extreme and large events, like the *right endpoint* of the model $F$ underlying the data,

$$(1.2) \qquad\qquad x^* := \sup\Big\{x \colon F(x) < 1\Big\} ,$$

a *high quantile* with probability $1 - p$, $p$ small, i.e., $\chi_{1-p} := \inf\big\{x \colon F(x) \geq 1 - p\big\}$, $p < 1/n$, with $n$ the available sample size, the *return period* and the *probability of exceedance* of a *high level*.

In this paper, we shall be interested in an application of *statistics of extremes* to the best personal marks attained at a few athletic events, in a context similar to the one used in Einmahl and Magnus (2008). We shall pay special attention to the estimation of $\gamma$, in (1.1), as well as of the right endpoint $x^*$, in (1.2), whenever finite, and of an indicator of the "excellence" of the level $x_{n:n}$, the maximum of the $n$ available observations. The right endpoint provides an estimate of the possible "world record" given the actual conditions, and the closer to one the "excellence" indicator of the level $x_{n:n}$ is, the better is the actual world record. In Section 2, we present some preliminary results in *extreme value theory*. In Section 3, we refer a few details on the semi-parametric estimation of a few parameters of extreme events. In Section 4, we provide heuristic choices of the thresholds for an adaptive semi-parametric estimation of the parameters of interest. Such heuristic choices take essentially into account the similarities of a few simple and alternative estimators of those parameters. In Section 4, we analyze data related to six athletic events and draw some final comments.

## 2.     PRELIMINARY RESULTS IN EXTREME VALUE THEORY

Let us think on any athletic event, like for instance the women marathon. Let us denote the best personal marks of $n$ different athletes by $X_1, X_2, ..., X_n$ and by

$$X_{1:n} \leq X_{2:n} \leq \cdots \leq X_{n:n}$$

the associated ascending o.s.'s. Under this set-up, $(X_1, X_2, ..., X_n)$ can be considered as independent, identically distributed (i.i.d.) observations from an underlying model $F$, obviously unknown. Let us also assume that, if necessary, data are transformed so that we can speak of maximum values (and not of minimum values). Indeed, any result for maxima can be easily reformulated for minima, due to the fact that $\min(X_1, X_2, ..., X_n) = -\max(-X_1, -X_2, ..., -X_n)$. However, this is not the transformation used later on in this paper for the analysis of data in athletics, where we shall convert running times in speeds, so that the higher the speed, the better. We shall thus work with upper o.s.'s.

One of the main results in extreme value theory is related to the possible limiting laws of the sequence $X_{n:n} := \max(X_1, X_2, ..., X_n)$, of maximum values, as $n \to \infty$. Since

$$\mathbb{P}\big(X_{n:n} \leq x\big) \;=\; \mathbb{P}\left(\bigcap_{i=1}^{n} \{X_i \leq x\}\right) = F^n(x) \;\underset{n\to\infty}{\longrightarrow}\; \begin{cases} 0 & \text{if } F(x) < 1 \;, \\ 1 & \text{if } F(x) = 1 \;, \end{cases}$$

we obviously have

$$X_{n:n} \;\xrightarrow[n\to\infty]{p}\; x^* \;,$$

with $x^*$ given in (1.2).

In order to obtain a possible non-degenerate behaviour for $X_{n:n}$, we thus need to normalize it. Similarly to the *central limit theorem* for sums or means, we know that if the maximum $X_{n:n}$, linearly normalized, converges to a non-degenerate random variable (r.v.), then there exist real constants $\{a_n\}_{n\geq 1}$ $(a_n > 0)$ and $\{b_n\}_{n\geq 1}$, the so-called *attraction coefficients* of $F$ to $G_\gamma$, in (1.1), such that

$$(2.1) \qquad \lim_{n\to\infty} \mathbb{P}\left(\frac{X_{n:n} - b_n}{a_n} \leq x\right) = \lim_{n\to\infty} F^n(a_n x + b_n) \;=\; G_\gamma(x) \;,$$

for some $\gamma \in \mathbb{R}$ (Gnedenko, 1943; de Haan, 1970). We then say that $F$ is in the *domain of attraction* (for maxima) of $G_\gamma$ and we use the notation $F \in \mathcal{D}_{\mathcal{M}}(G_\gamma)$.

The extreme value index $\gamma$, in (1.1), measures essentially the weight of the right-tail $\overline{F} = 1 - F$. If $\gamma < 0$, the right-tail is *light*, i.e., $F$ has a finite right endpoint $(x^* < +\infty)$. If $\gamma > 0$, the right-tail is *heavy*, of a negative polynomial type, i.e., $F$ has an infinite right endpoint. If $\gamma = 0$, the right-tail is of an *exponential* type and the right endpoint can be either finite or infinite. In Figure 1,

we represent graphically the probability density function (p.d.f.) associated with the extreme value d.f., in (1.1), i.e. $g_\gamma(x) = dG_\gamma(x)/dx$, for $\gamma = -0.5$, 0 and 0.5. We also picture the standard normal p.d.f., $\varphi(x) = \exp(-x^2/2)/\sqrt{2\pi}$, $x \in \mathbb{R}$, as well as a "zoom" of the right-tails of these four models. It is clear the lightness of the right-tail of $G_\gamma$ for $\gamma < 0$ (finite right endpoint), followed by the normal tail and next the Gumbel tail ($\gamma = 0$). It is also clear the heaviness of the right-tail of $G_\gamma$ for $\gamma > 0$.
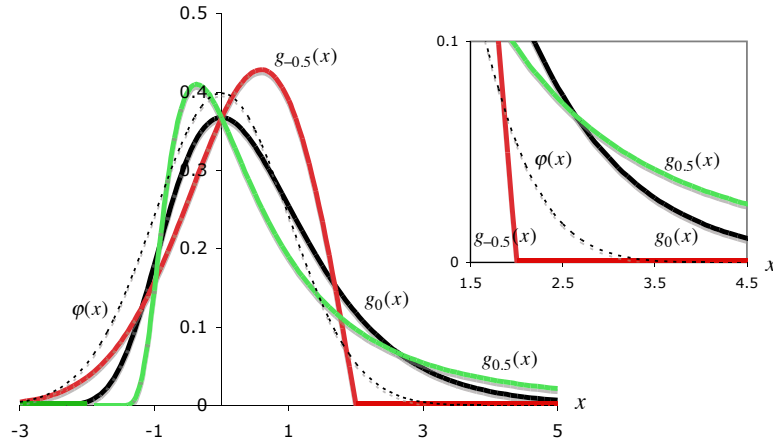


**Figure 1**: Extreme value p.d.f.'s, $g_\gamma(x)$, with $\gamma = -0.5, 0$ and 0.5, and normal p.d.f., $\varphi(\cdot)$.

**Remark 2.1.** Note that to say that $F \in \mathcal{D}_\mathcal{M}(G_\gamma)$ is equivalent to saying that for all $x$ real and such that $0 < G_\gamma(x) < 1$, $\lim_{n\to\infty} n \ln F(a_n x + b_n) = \ln G_\gamma(x) = -(1+\gamma x)^{-1/\gamma}$. Consequently, $F(a_n x + b_n) \to 1$ for those values of $x$. Since $\lim_{n\to\infty} \left(-\ln F(a_n x + b_n)\right)/\left(1 - F(a_n x + b_n)\right) = 1$, we equivalently have

$$(2.2) \qquad \lim_{n\to\infty} n\left(1 - F(a_n x + b_n)\right) = -\ln G_\gamma(x) = (1+\gamma x)^{-1/\gamma}.$$

Let us define

$$(2.3) \qquad U(t) := F^\leftarrow(1 - 1/t) \ \ (t > 1), \qquad F^\leftarrow(x) := \inf\{y \colon F(y) \geq x\},$$

with $F^\leftarrow$ denoting thus the generalized inverse function of $F$. It is reasonably easy to prove (de Haan, 1984) that, with $G_\gamma^{-1}$ denoting the inverse function of the extreme value d.f. $G_\gamma$ in (1.1),

$$\lim_{t\to\infty} \frac{U(tx) - b_t}{a_t} = G_\gamma^{-1}\left(\exp(-1/x)\right) = \frac{x^\gamma - 1}{\gamma},$$

for all $x > 0$, with $a_t \equiv a(t) \equiv a_{[t]}$, $[t] = $ integer part of $t$ and $a_n$ the scale attraction coefficient in (2.1). Also $b_t \equiv b(t) \equiv b_{[t]}$, with $b_n$ the location attraction

coefficient, also in (2.1). Moreover, we can choose $b_t = U(t)$, with $U(\cdot)$ defined in (2.3) (see Theorem 1.1.2 of de Haan and Ferreira, 2006). More generally,

$$(2.4) \qquad F \in \mathcal{D}_{\mathcal{M}}(G_\gamma) \quad \Longleftrightarrow \quad \lim_{t\to\infty} \frac{U(tx) - U(t)}{a(t)} = \frac{x^\gamma - 1}{\gamma} \,,$$

for all $x > 0$, with $U(\cdot)$ defined in (2.3).

**Remark 2.2.** When $\gamma = 0$, and by continuity arguments, the functions $-\ln G_\gamma(x) = (1 + \gamma x)^{-1/\gamma}$ and $G_\gamma^{-1}(\exp(-1/x)) = (x^\gamma - 1)/\gamma$ should be interpreted as $\exp(-x)$ and $\ln x$, respectively.

## 3. SEMI-PARAMETRIC ESTIMATION OF A FEW RELEVANT PARAMETERS OF EXTREME EVENTS

On the basis of the available random sample, $(X_1, X_2, ..., X_n)$, let us see how to estimate the extreme value index $\gamma$, the primordial parameter in statistics of extremes, the scale $a$, the location $b$, the right endpoint $x^*$ and the return period of a high level $x_{\mathrm{H}}$, usually defined as the expected number of exceedances of such a level.

### 3.1. Estimation of the extreme value index

For any integer $j \geq 1$, let us denote

$$(3.1) \qquad L_{k,n}^{(j)} := \frac{1}{k} \sum_{i=1}^{k} \left\{ 1 - \frac{X_{n-k:n}}{X_{n-i+1:n}} \right\}^j$$

and

$$(3.2) \qquad M_{k,n}^{(j)} := \frac{1}{k} \sum_{i=1}^{k} \left\{ \ln X_{n-i+1:n} - \ln X_{n-k:n} \right\}^j.$$

These statistics have revealed to be fundamental in *statistics of extremes*. For the estimation of $\gamma$, we shall first refer three estimators, valid, i.e. consistent, for all $\gamma \in \mathbb{R}$:

**1.** The *moment* ($M$) estimator (Dekkers *et al.*, 1989), with the functional form

$$(3.3) \qquad \hat{\gamma}_{k,n}^M \equiv M_{k,n} := M_{k,n}^{(1)} + \frac{1}{2} \left\{ 1 - \left( \frac{M_{k,n}^{(2)}}{[M_{k,n}^{(1)}]^2} - 1 \right)^{-1} \right\},$$

$M_{k,n}^{(j)}$, $j = 1, 2$, defined in (3.2).

**2.** The *generalized Hill* (*GH*) estimator introduced in Beirlant *et al.* (1996), further studied in Beirlant *et al.* (2005), and based on the Hill estimator (Hill, 1975), the statistic $M_{k,n}^{(1)}$, in (3.2), also denoted

$$(3.4) \qquad \hat{\gamma}_{k,n}^{H} \equiv H_{k,n} := \frac{1}{k} \sum_{i=1}^{k} \left\{ \ln X_{n-i+1:n} - \ln X_{n-k:n} \right\},$$

and valid only for $\gamma \geq 0$. The *GH*-estimator, valid for all $\gamma \in \mathbb{R}$, and with $\hat{\gamma}_{k,n}^{H}$ given in (3.4), has the functional form

$$(3.5) \qquad \hat{\gamma}_{k,n}^{GH} \equiv GH_{k,n} := \hat{\gamma}_{k,n}^{H} + \frac{1}{k} \sum_{i=1}^{k} \left\{ \ln \hat{\gamma}_{i,n}^{H} - \ln \hat{\gamma}_{k,n}^{H} \right\}.$$

**3.** The *mixed moment* (*MM*) estimator (Fraga Alves *et al.*, 2009), with the functional form

$$(3.6) \qquad \hat{\gamma}_{k,n}^{MM} \equiv MM_{k,n} := \frac{\hat{\varphi}_{k,n} - 1}{1 + 2 \min(\hat{\varphi}_{k,n} - 1, 0)}, \qquad \hat{\varphi}_{k,n} := \frac{M_{k,n}^{(1)} - L_{k,n}^{(1)}}{\left( L_{k,n}^{(1)} \right)^2},$$

$L_{k,n}^{(1)}$ and $M_{k,n}^{(1)}$ defined in (3.1) and (3.2), respectively.

The three estimators in (3.3), (3.5) and (3.6) are consistent in $\mathcal{D}_{\mathcal{M}}(G_\gamma)$, $\gamma \in \mathbb{R}$, if $k = k_n$ is an intermediate sequence, i.e., a sequence of integers such that

$$(3.7) \qquad k = k_n \to \infty \quad \text{and} \quad k_n = o(n), \qquad \text{as} \quad n \to \infty.$$

Due to the specificity of the data, we shall also consider another simple estimator:

**4.** The *location invariant* estimator (*F*) introduced in Falk (1995),

$$(3.8) \qquad \hat{\gamma}_{k,n}^{F} \equiv F_{k,n} := \frac{1}{k} \sum_{i=1}^{k-1} \ln \frac{X_{n:n} - X_{n-i:n}}{X_{n:n} - X_{n-k:n}},$$

valid only for a negative extreme value index smaller than $-0.5$.

We still would like to refer the so-called "*maximum likelihood*" (*ML*) estimator, introduced in Smith (1987) and further studied in Drees *et al.* (2004). Such an estimator is valid and asymptotically normal for all $\gamma > -1$ (see Zhou, 2009, 2010, for details in the region $-1 < \gamma \leq -1/2$). The extreme value index *ML*-estimator is based on the application of the maximum likelihood methodology to the excesses $X_{n-i+1:n} - X_{n-k:n}$, $1 \leq i \leq k$. These excesses are approximately the $k$ top o.s.'s of a sample of size $k$ from a *generalized Pareto* model, strongly related to the *extreme value* d.f. $G_\gamma$ in (1.1), through the relation

$$(3.9) \qquad GP(x; \gamma, \alpha) = 1 + \ln G_\gamma(\alpha x / \gamma) = 1 - (1 + \alpha x)^{-1/\gamma}, \quad 1 + \alpha x > 0, \ x > 0,$$

with $\alpha, \gamma \in \mathbb{R}$. This is a re-parametrization due to Davison (Davison, 1984).

Then, with such a re-parametrization, the *ML*-estimator of $\gamma$ has an explicit expression as a function of the *ML*-estimator $\widehat{\alpha} = \widehat{\alpha}_{ML}$ of $\alpha$ and the sample of the excesses. We have

$$(3.10) \qquad \hat{\gamma}_{k,n}^{ML} \; = \; \widehat{\gamma}_{k,n,\hat{\alpha}}^{ML} \; \equiv \; ML_{k,n} := \frac{1}{k} \sum_{i=1}^{k} \ln\big(1 + \hat{\alpha}(X_{n-i+1:n} - X_{n-k:n})\big) \,.$$

The estimates $\widehat{\alpha} = \widehat{\alpha}_{ML}$ are obtained through numerical iterative methods, usually computationally time-consuming. This is the reason why we shall not consider these estimators in the Monte Carlo simulation in Section 4, related to heuristic choices of thresholds. We shall however consider the *ML*-estimators in the data analysis provided in Section 5.3.2, due to their nice asymptotic properties for $-1/2 < \gamma < 0$ (see, for instance, Gomes and Neves, 2008, among others).

For a large variety of models, and under mild second-order conditions, we can guarantee the asymptotic normality of all the above mentioned estimators and can build approximate confidence intervals (CI's) for $\gamma$, as well as for all other parameters of extreme events, like the ones discussed next in Section 3.2. We merely need to assume the existence of a function $A(t)$, converging to 0, as $t \to \infty$, which measures the rate of convergence of the sequence of maximum values to a non-degenerate limit r.v. and that "measures" also the bias of the estimators in a great variety of situations (see de Haan and Ferreira, 2006, for details). Such a second-order condition can be written as

$$(3.11) \qquad \lim_{t\to\infty} \frac{\frac{U(tx)-U(t)}{a(t)} - \frac{x^{\gamma}-1}{\gamma}}{A(t)} \; = \; H_{\gamma,\rho}(x) := \frac{1}{\rho}\left(\frac{x^{\gamma+\rho}-1}{\gamma+\rho} - \frac{x^{\gamma}-1}{\gamma}\right),$$

for all $x > 0$, where $\rho \le 0$ is a second order parameter controlling the speed of convergence in the first-order condition, (2.4), and $|A(t)| \in RV_{\rho}$, with $RV_a$ standing for the class of regularly varying functions at infinity with an index of regular variation $a$, i.e. positive measurable functions $g$ such that $\lim_{t\to\infty} g(tx)/g(t) = x^a$, for all $x > 0$. Note that for the *extreme value* d.f., in (1.1), condition (3.11) holds, with $\rho = -1$ if $\gamma \neq 1$ and $\rho = -2$ if $\gamma = 1$. For the Generalized Pareto d.f., in (3.9), $U(t) = (t^{\gamma} - 1)/\gamma$, and we can say that condition (3.11) holds with $A(t) \equiv 0$ $(\rho = -\infty)$.

## 3.2.  Semi-parametric estimation of other parameters of interest

### 3.2.1. Estimation of location and scale

As mentioned before, we have $b_t = U(t)$, with $U(\cdot)$ defined in (2.3). On another side, the universal uniform transformation enables us to guarantee that $\forall F$, unknown and underlying the r.v. $X$, $X \stackrel{d}{=} U(Y)$, with $Y$ a unit Pareto r.v.,

i.e. a r.v. with d.f. $F_Y(y) = 1 - 1/y$, $y \geq 1$. Consequently,

$$X_{n-k:n} \overset{d}{=} U(Y_{n-k:n}), \quad \text{and since} \quad Y_{n-k:n} \overset{p}{\sim} n/k, \quad \text{as } n \to \infty ,$$

where $X_n \overset{p}{\sim} Y_n$ means that $X_n/Y_n$ converges in probability to one, as $n \to \infty$, it is sensible to consider

$$\hat{b} = \hat{b}_{k,n} = \widehat{U}(n/k) = X_{n-k:n} .$$

And for any extreme value index estimator, $\hat{\gamma}^\bullet \equiv \hat{\gamma}_{k,n}^\bullet$, we can consider (de Haan and Ferreira, 2006)

$$\hat{a}^\bullet = \hat{a}_{k,n}^\bullet = X_{n-k:n} M_{k,n}^{(1)}\big(1 - \min(0, \hat{\gamma}^\bullet)\big) ,$$

with $M_{k,n}^{(1)}$ given in (3.2).

## 3.2.2. Estimation of the right endpoint for $\gamma < 0$

For large values of $t$ and $\gamma \neq 0$, if we take into account the validity of condition (2.4), we can write the approximation $U(tx) \approx U(t) + a(t)(x^\gamma - 1)/\gamma$. But $x^* = U(\infty)$ and for all $\gamma < 0$, $(x^\gamma - 1)/\gamma \to -1/\gamma$, as $x \to \infty$. If we consider $t = n/k$, with $k$ intermediate, we can thus guarantee that, whenever $\hat{\gamma}^\bullet < 0$,

$$x^* \approx U(n/k) - a(n/k)/\gamma \quad \Longrightarrow \quad \hat{x}_\bullet^* := \hat{b} - \hat{a}^\bullet/\hat{\gamma}^\bullet .$$

As we have the obvious restriction $x_{n:n} \leq x^*$, we shall instead consider the right endpoint estimator

$$(3.12) \qquad \hat{x}_{k,n|\bullet}^* = \max\left(X_{n:n}, X_{n-k:n}\Big(1 - M_{k,n}^{(1)}(1 - \min(0, \hat{\gamma}_{k,n}^\bullet))/\hat{\gamma}_{k,n}^\bullet\Big)\right) .$$

## 3.2.3. Estimation of the return period of a high level $x_{\mathrm{H}}$ and similar indicators

In a pure framework of i.i.d. observations, if we think on the number of observations $N_{\mathrm{H}}$ needed to reach a value higher than $x_{\mathrm{H}}$, such a r.v. has support $\{1, 2, ...\}$ and $\mathbb{P}(N_{\mathrm{H}} = r) = p_{\mathrm{H}}(1 - p_{\mathrm{H}})^{r-1}$, $r \geq 1$, with $p_{\mathrm{H}} = \mathbb{P}(X > x_{\mathrm{H}}) = 1 - F(x_{\mathrm{H}})$, i.e. $N_{\mathrm{H}}$ is a *geometric* r.v.. The *return period* of the high level $x_{\mathrm{H}}$ is usually defined as the mean value of $N_{\mathrm{H}}$, being given by

$$R(x_{\mathrm{H}}) := \frac{1}{p_{\mathrm{H}}} = \frac{1}{1 - F(x_{\mathrm{H}})} .$$

In the framework of this paper, it is perhaps sensible to think on the $n$ athletes under consideration, and to define an indicator associated with a high level $x_{\mathrm{H}}$

as the mean number of athletes, among the $n$, who will have in the future a personal mark larger than $x_H$. We thus have the mean value of a Binomial$\big(n, p_H = 1 - F(x_H)\big)$ r.v., given by

$$MN(x_H) := n\big(1 - F(x_H)\big) \, ,$$

with $MN$ standing for *mean number*.

On the basis of the limiting relation in (2.2), we can then consider the estimators

$$\widehat{R}^{\bullet}(x_H) \, := \, \frac{n}{k} \left( \min\left( +\infty, \ 1 + \hat{\gamma}^{\bullet}\left(\frac{x_H - \hat{b}}{\hat{a}^{\bullet}}\right) \right) \right)^{1/\hat{\gamma}^{\bullet}}$$

and

$$\widehat{MN}^{\bullet}(x_H) \, \equiv \, n\,\widehat{p}_H^{\bullet} \, := \, k \left( \max\left( 0, \ 1 + \hat{\gamma}^{\bullet}\left(\frac{x_H - \hat{b}}{\hat{a}^{\bullet}}\right) \right) \right)^{-1/\hat{\gamma}^{\bullet}}$$

of $R(x_H)$ and $MN(x_H)$, respectively.

Note that for $x_H = x_{n:n}$, $F$ absolutely continuous, and denoting $(U_1, ..., U_n)$ a random sample from a uniform d.f. in $(0,1)$, we have

$$MN(X_{n:n}) \, = \, n\big(1 - F(X_{n:n})\big) \, \overset{d}{=} \, n\,U_{1:n} \, ,$$

which converges weakly towards a unit exponential r.v., as $n \to \infty$. Consequently, the sequence of r.v.'s $\exp\big(-MN(X_{n:n})\big)$ converges weakly towards a uniform r.v. in $(0,1)$. In the data analysis provided in Section, 5.3.2 we shall thus consider

(3.13) $$\widehat{E}_n^{\bullet} \, \equiv \, \widehat{E}_{k,n}^{\bullet} \, := \, \exp\big(-\widehat{MN}^{\bullet}(X_{n:n})\big)$$

as an estimator of an indicator of the "excellence" of the world record $X_{n:n}$, given by $E_n := \exp\big(-MN(X_{n:n})\big)$. Note that the $E$-indicator was chosen merely because it lies in the finite support [0,1]. The closer to 1 this indicator is, the better is the actual world record. Such an indicator is strongly related to the *quality of the current world record*'s indicator $Q := -\ln E_n = n\big(1 - F(X_{n:n})\big)$ of Einmahl and Magnus (2008), the expected number of exceedances of the current world record, $X_{n:n}$, conditional on this world record.

For further details on most of the subjects of this Section, see Chapters 1 and 4 of de Haan and Ferreira (2006).

# 4. HEURISTIC CHOICES OF THRESHOLDS IN THE SEMI-PARAMETRIC EXTREME VALUE INDEX, RIGHT END-POINT AND EXCEEDANCE PROBABILITY ESTIMATION: A MONTE-CARLO STUDY

For any arbitrary estimator, $\hat{\gamma}^{\bullet}_{k,n}$, of $\gamma$, like the ones in (3.3), (3.4), (3.5), (3.6), (3.8) and (3.10), and under the validity of a second-order condition like the one in (3.11), we get an asymptotic distributional representation of the type

$$(4.1) \qquad \hat{\gamma}^{\bullet}_{k,n} \overset{d}{=} \gamma + \frac{\sigma_{\bullet} P^{\bullet}_k}{\sqrt{k}} + v_{\bullet} A(n/k)\left(1 + o_p(1)\right) ,$$

with $P^{\bullet}_k \overset{a}{\sim} \text{Normal}(0,1)$. Consequently, for intermediate levels $k$, i.e., levels such that (3.7) holds, and also such that $\sqrt{k}\, A(n/k) \to \lambda$, finite, $\exists\, v_{\bullet} \in \mathbb{R}$ and $\sigma_{\bullet} \in \mathbb{R}^+$ such that

$$(4.2) \qquad \sqrt{k}\,(\hat{\gamma}^{\bullet}_{k,n} - \gamma) \xrightarrow[n\to\infty]{d} \text{Normal}(\lambda v_{\bullet}, \sigma^2_{\bullet}) .$$

The "*asymptotic mean squared error*" (*AMSE*) is defined as

$$AMSE\big(\hat{\gamma}^{\bullet}_{k,n}\big) := \frac{\sigma^2_{\bullet}}{k} + v^2_{\bullet}\, A^2(n/k) ,$$

i.e. we get asymptotic bias and variance given by $BIAS_{\infty}\big(\hat{\gamma}^{\bullet}_{k,n}\big) := v_{\bullet}\, A(n/k)$ and $\text{Var}_{\infty}\big(\hat{\gamma}^{\bullet}_{k,n}\big) := \sigma^2_{\bullet}/k$, respectively. If $\lambda = 0$, the mean value of the limiting normal law in (4.2) is equal to zero.

Let us define $k^{\bullet}_0 = k^{\bullet}_0(n) := \arg\min_k MSE\big(\hat{\gamma}^{\bullet}_{k,n}\big) \sim \arg\min_k AMSE\big(\hat{\gamma}^{\bullet}_{k,n}\big)$, the level associated with a minimal *AMSE*, as the optimal level for the estimation of $\gamma$ through $\hat{\gamma}^{\bullet}_{k,n}$, and let us denote $\hat{\gamma}^{\bullet}_{n0} := \hat{\gamma}^{\bullet}_{k^{\bullet}_0,n}$, the estimator computed at its optimal level. With the notation $A(t) = \beta t^{\rho}$, $\rho < 0$, the value $\sigma_{\bullet}$ is a function of $\gamma$ and $v_{\bullet}$ is usually a function of $\beta$ and $\rho$ (possibly also of $\gamma$). We then get

$$(4.3) \qquad k^{\bullet}_0 = \big(\sigma^2_{\bullet}/(-2\,\rho\, v^2_{\bullet}\beta^2)\big)^{1/(1-2\rho)}\, n^{-2\rho/(1-2\rho)} .$$

In order to estimate $k^{\bullet}_0$ in (4.3), in a simple and precise way, we thus need to have "nice" estimates of the second-order parameters $(\beta, \rho)$. However, whereas such an estimation is reliable for $\gamma > 0$ (see, for instance, Caeiro *et al.*, 2005; Gomes and Pestana, 2007; Gomes *et al.*, 2007, 2008, among others), this is not the case for $\gamma \le 0$. Notice however that we can estimate $\rho$, for a general $\gamma \in \mathbb{R}$, through the estimators in Fraga Alves *et al.* (2003). Even so, the optimal level, in (4.3), depends often not only on $\beta$ but also on $\gamma$. The estimation of $k_0$ can then be made recursively, but it induces a high volatility in the estimates and a drastic loss of efficiency. Alternatively, we could also use, for instance, bootstrap

methods (Draisma *et al.*, 1999; Danielson *et al.*, 2001; Gomes and Oliveira, 2001) for an optimal adaptive choice of $k$. Here, after deciding on a negative value for $\gamma$, as will be the case in Section 5, we propose the following heuristic choice of the threshold $k$. Let us denote $\hat{\gamma}_{k,n}^{(i)}$, $i \in \mathcal{K} = \{1, 2, 3, 4\}$, the set of alternative (and computationally simple to obtain) *EVI*-estimators in (3.3), (3.5), (3.6) and (3.8). Then, consider

$$(4.4) \qquad k_{\min}^* := \arg\min_k \sum_{(i,j)\in\mathcal{K},\, i\neq j} \big(\hat{\gamma}_{k,n}^{(i)} - \hat{\gamma}_{k,n}^{(j)}\big)^2$$

and

$$(4.5) \qquad T^* := T_{k_{min}^*,n}\,, \qquad \text{with} \quad T = M \text{ or } GH \text{ or } MM \text{ or } F\,,$$

$M_{k,n}$, $GH_{k,n}$, $MM_{k,n}$, $F_{k,n}$ and $k_{\min}^*$ given in (3.3), (3.5), (3.6), (3.8) and (4.4), respectively. We cannot claim any kind of asymptotic optimality for the choice $k_{\min}^*$, in (4.4), in the sense that we would like to have $k_{\min}^*/k_0^{\bullet} \to 1$, as $n \to \infty$. However, if $b_{\bullet} \neq 0$, we can guarantee that the value $k_{\min}^*$ in (4.4) is of the order of $n^{-2\rho/(1-2\rho)}$, i.e., of the same order of the optimal value $k_0^{\bullet}$ in (4.3). Consequently, (4.2) holds whenever we there replace $k$ by $k_{\min}^*$. Moreover, the value in (4.4) seems to be heuristically appealing whenever we want to take into account a set of alternative semi-parametric estimators of a parameter of extreme events. It is expected that there will be a region where all estimators work, and in such a region we surely get close values for all estimates and the smallest possible value for the indicator in (4.4). If we enlarge the set $\mathcal{K}$, in order to include the extreme value index *ML*-estimator, in (3.10), as we shall do in the data analysis performed in Section 5.3.2, we shall use the notations $k_{\min}^{**}$ and $T^{**}$ for the entities equivalent to the ones in (4.4) and (4.5), respectively.

We shall also consider the same type of heuristic procedure for the estimation of the right endpoint $x^*$, in (1.2), done through similar adaptive right endpoint estimators,

$$(4.6) \qquad \hat{x}_T^* := \hat{x}_{k_{\min}^{*x},n|T}^*\,, \qquad \text{again with} \quad T = M \text{ or } GH \text{ or } MM \text{ or } F\,,$$

$\hat{x}_{k,n|\bullet}^*$ given in (3.12), and where, for the same set $\mathcal{K}$ and the same notation as before,

$$(4.7) \qquad k_{\min}^{*x} := \arg\min_k \sum_{(i,j)\in\mathcal{K},\, i\neq j} \big(\hat{x}_{k,n|(i)}^* - \hat{x}_{k,n|(j)}^*\big)^2 =: k^{*x}\,.$$

Similarly, we shall use the notations $k_{\min}^{**x}$ and $\hat{x}_T^{**}$, whenever we include in $\mathcal{K}$ the *ML*-estimator, in (3.10), for the estimation of the right endpoint. A similar method was also applied to the estimators of the "excellence" indicators, in (3.13) (or equivalently to the exceedance probability of $X_{n:n}$). We shall use the obvious similar notations $\widehat{E}_{\bullet}^*$, $\widehat{E}_{\bullet}^{**}$ for those adaptive estimators and $k_{\min}^{*E} \equiv k^{*E}$, $k_{\min}^{**E} \equiv k^{**E}$ for the adaptive choices of $k$.

In order to obtain distributional properties of the adaptive estimators under consideration, we have performed simulation studies of size $5000 \times 10$ for sample sizes $n = 100, 200, 300, 400, 500, 1000, 2000$ and $5000$, from a reasonably large variety of models. Due to characteristics of the data, which are maxima of a certain number of marks, and should consequently be associated with an underlying d.f. quite close to the *extreme value* ($EV$) model, we shall uniquely present, as an illustration, the results associated with an underlying model $F(x) = G_\gamma(x)$, with $G_\gamma(x)$ given in (1.1), $\gamma = -0.1$ and $-0.3$.

For each value of $n$, we have simulated not only the mean values and root mean squared errors of the four estimators in (4.5), but also of the similar adaptive right endpoint estimators in (4.6). A similar method was also applied to the estimators of the "excellence" indicators, in (3.13) (or equivalently to the exceedance probability of $X_{n:n}$). As mentioned before, we shall use the obvious notation $\widehat{E}_\bullet^*$ for those adaptive estimators and $k_{\min}^{*E} \equiv k^{*E}$ for the adaptive choice of $k$. Due to the stability of the sample paths of the estimators in (3.8), even when we cannot guarantee their consistency, the results do not depend on either the inclusion or the non-inclusion of such an estimator.

For underlying $EV$ models, with $\gamma = -0.1$ and $-0.3$, the estimates of the absolute bias ($|BIAS|$) and root mean squared error ($RMSE$) of the adaptive $EVI$-estimators are presented in Figures 2 and 3, respectively. We also present in these figures the corresponding values of at least one of the estimators at its simulated optimal level, denoted $T_0$, with $T = M$, $GH$, $MM$ or $F$. For the bias structure, we present only one $T_0$, the one with the lowest absolute bias for large values of $n$. The introduction of the $ML$-estimator, in (3.10), does not lead to very different conclusions, but increases drastically the time of computation and consequently the loss of precision associated with the $REFF$ indicators. Similar patterns have been obtained for underlying $GP$ and reversed-Burr parents, and we see no need to present those extra results.



**Figure 2**: Absolute values of bias (*left*) and root mean squared errors (*right*) of the adaptive extreme value index estimators in (4.5), for an *extreme value* model with $\gamma = -0.1$.
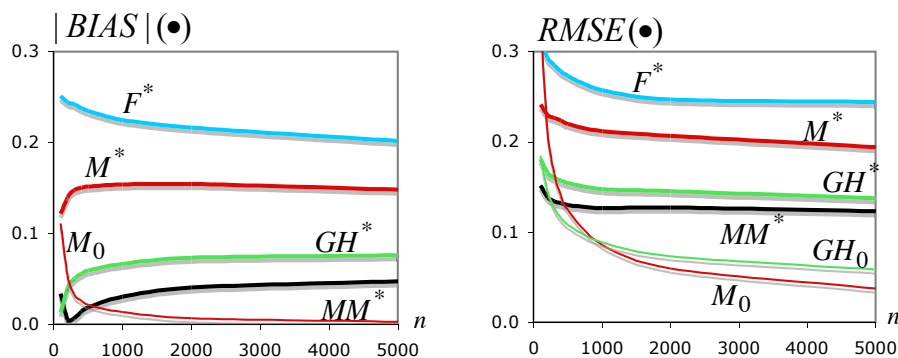
**Figure 3**:   Absolute values of bias (*left*) and root mean squared errors (*right*)
of the adaptive extreme value index estimators in (4.5), for an
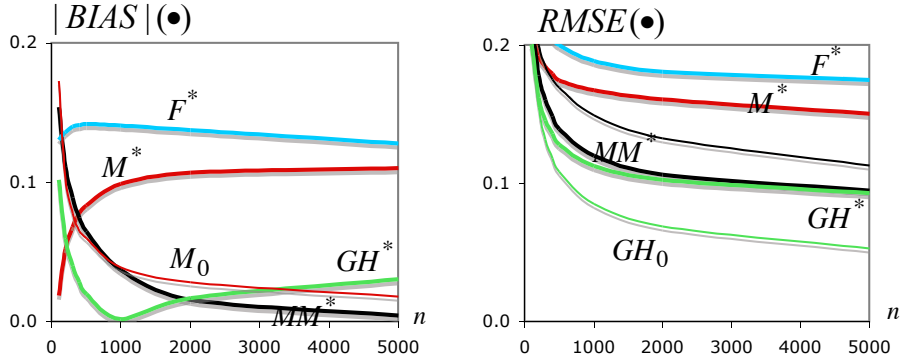*extreme value* model with $\gamma = -0.3$.

A few remarks related with the adaptive estimators in (4.5) for underlying
*EV* models:

- For $\gamma = -0.1$, the absolute bias of $MM^*$ is the smallest one, except for
  $n = 100$. For this sample size, and regarding absolute bias, $GH^*$ beats
  $MM^*$. Regarding *MSE*, the best of the adaptive estimates is $MM^*$, for
  all $n$. As $\gamma$ decreases, and regarding bias, $MM^*$ is replaced by $GH^*$ for
  moderate $n$ and by $M^*$ for small $n$.

- For $\gamma = -0.3$, the absolute bias of $MM^*$ is the smallest one, only for
  $n \geq 2000$. For $300 \leq n \leq 1000$, $GH^*$ beats the other estimators. For
  $n \leq 200$, the smallest absolute bias is the one of $M^*$. Regarding *MSE*,
  the best of the estimates is $GH^*$, quite close to $MM^*$ for all $n$.

- Notice the overall worst performance of the estimator $F^*$, essentially
  due to the region of $\gamma$-values under consideration.

Regarding the "potential" estimators $T_0$ at simulated optimal levels, with
$T = M, GH, MM$ or $F$, we draw the following comments:

- At simulated optimal levels, $GH_0$ achieves the minimum *MSE* for all $n$,
  if $\gamma = -0.3$. For the other values of $\gamma$, $GH_0$ is the best one for small $n$,
  but $M_0$ becomes the best for large $n$ ($n \geq 1000$ for $\gamma = -0.1$).

- Regarding smallest absolute bias at simulated optimal levels, $M_0$ is the
  best for all $n$, if $\gamma = -0.1$. For the other values of $\gamma$, $M_0$ is the best for
  $n \geq 200$. For $n = 100$, $GH_0$ overpasses all other ones.

In Table 1, for *EV* underlying parents, for a few values of $n$, and for $T =
M, GH, MM$ and $F$, we present two relative efficiency indicators of $T^*$, in (4.5),
relatively to $T_0$, and to $S_0$, the best $T_0$-estimator, i.e. the one with smallest *MSE*

at optimal level. With the notation $MSE_s(S_0) = \min\big(MSE_s(M_0), MSE_s(GH_0),$ $MSE_s(MM_0), MSE_s(F_0)\big)$, we have simulated

$$REFF_1 := \sqrt{\frac{MSE_s(T_0)}{MSE_s(T^*)}} \qquad \text{and} \qquad REFF_2 := \sqrt{\frac{MSE_s(S_0)}{MSE_s(T^*)}} \ ,$$

the values placed in the first and second line, respectively, of any entrance $T^*$. Note that the higher than one these indicators are, the better $T^*$ performs. Moreover, we obviously have $REFF_1 \geq REFF_2$. For all $n$, the highest $REFF_1$-indicator is written in **bold** and the highest $REFF_2$ indicator is written in *italic*. The results obtained are consistent with the remarks made above.

**Table 1**: Simulated *REFF*'s of the adaptive *EVI*-estimators under study, together with associated 95% CI's, for *extreme value* underlying parents.

| | | | | | |
|---|---|---|---|---|---|
| ***EV* parent**, $\gamma = -0.1$ | | | | | |
| $n$ | 100 | 200 | 500 | 1000 | 2000 |
| $M^*$ | **1.403** ± 0.015 | 0.921 ± 0.011 | 0.508 ± 0.130 | 0.404 ± 0.002 | 0.289 ± 0.003 |
| | 0.762 ± 0.008 | 0.632 ± 0.006 | 0.491 ± 0.006 | 0.404 ± 0.002 | 0.289 ± 0.003 |
| $MM^*$ | 1.212 ± 0.009 | **1.116** ± 0.016 | **0.973** ± 0.016 | **0.874** ± 0.013 | **0.771** ± 0.014 |
| | *1.212* ± *0.009* | *1.050* ± *0.014* | *0.833* ± *0.014* | *0.677* ± *0.010* | *0.469* ± *0.007* |
| $GH^*$ | 1.021 ± 0.010 | 0.877 ± 0.010 | 0.700 ± 0.009 | 0.599 ± 0.009 | 0.506 ± 0.006 |
| | 1.015 ± 0.010 | 0.877 ± 0.010 | 0.700 ± 0.009 | 0.579 ± 0.009 | 0.410 ± 0.004 |
| $F^*$ | 1.023 ± 0.009 | 0.898 ± 0.002 | 0.768 ± 0.003 | 0.695 ± 0.004 | 0.628 ± 0.003 |
| | 0.587 ± 0.006 | 0.496 ± 0.004 | 0.395 ± 0.002 | 0.332 ± 0.003 | 0.243 ± 0.001 |

| | | | | | |
|---|---|---|---|---|---|
| ***EV* parent**, $\gamma = -0.3$ | | | | | |
| $n$ | 100 | 200 | 500 | 1000 | 2000 |
| $M^*$ | **1.932** ± 0.020 | **1.355** ± 0.020 | 0.896 ± 0.013 | 0.657 ± 0.006 | 0.489 ± 0.003 |
| | 0.978 ± 0.011 | 0.807 ± 0.005 | 0.621 ± 0.006 | 0.509 ± 0.008 | 0.428 ± 0.003 |
| $MM^*$ | 1.060 ± 0.006 | 1.121 ± 0.009 | **1.214** ± 0.008 | **1.245** ± 0.014 | **1.246** ± 0.012 |
| | 0.962 ± 0.008 | 0.871 ± 0.005 | 0.780 ± 0.006 | 0.711 ± 0.007 | 0.645 ± 0.006 |
| $GH^*$ | 1.038 ± 0.010 | 0.959 ± 0.006 | 0.847 ± 0.010 | 0.757 ± 0.007 | 0.670 ± 0.008 |
| | *1.038* ± *0.010* | *0.959* ± *0.006* | *0.847* ± *0.010* | *0.757* ± *0.007* | *0.670* ± *0.008* |
| $F^*$ | 1.097 ± 0.010 | 0.959 ± 0.005 | 0.788 ± 0.007 | 0.693 ± 0.145 | 0.597 ± 0.006 |
| | 0.895 ± 0.009 | 0.718 ± 0.006 | 0.544 ± 0.005 | 0.452 ± 0.005 | 0.388 ± 0.004 |

The behaviour of the right endpoint semi-parametric estimators is quite erratic, even when we consider equation (3.12), to make them coherent with the data. Such a behaviour is even more catastrophic when we do not make them coherent with the data, and the most usual estimators in the literature are in fact "raw", in the sense that they have not been modified in order to be larger than the maximum in the sample, as needed. Indeed, alternative semi-parametric estimators of the right endpoint are urgently needed. The bias and the *RMSE* of the estimators in (4.6) almost overlap, and we see no reason to present figures similar to the ones drawn for the adaptive *EVI*-estimators in (4.5). A similar comment applies to the adaptive estimators of the "excellence indicator".

Table 2 is equivalent to Table 1, but for the adaptive right endpoint estimation. Similarly, Table 3 is equivalent to Table 1, now for the exceedance probability estimation (or equivalently, for the "excellence" indicator).

**Table 2**: Simulated *REFF*'s of the adaptive right endpoint estimators under study, together with associated 95% CI's, for *extreme value* underlying parents.

| | | | | | |
|---|---|---|---|---|---|
| **EV parent**, $\gamma = -0.1$ | | | | | |
| $n$ | 100 | 200 | 500 | 1000 | 2000 |
| $M^*$ | $3.272 \pm 0.771$<br>*0.877 ± 0.006* | **2.786** ± 0.569<br>*0.865 ± 0.005* | **1.660** ± 0.293<br>*0.866 ± 0.007* | **1.001** ± 0.002<br>*0.863 ± 0.008* | $0.998 \pm 0.000$<br>*0.858 ± 0.005* |
| $MM^*$ | $1.819 \pm 0.483$<br>*0.877 ± 0.006* | $1.000 \pm 0.000$<br>*0.865 ± 0.005* | $1.000 \pm 0.294$<br>*0.866 ± 0.007* | $1.000 \pm 0.000$<br>*0.863 ± 0.008* | **1.000** ± 0.000<br>*0.858 ± 0.005* |
| $GH^*$ | **5.454** ± 1.139<br>*0.877 ± 0.006* | $2.106 \pm 0.698$<br>*0.865 ± 0.005* | $1.000 \pm 0.000$<br>*0.866 ± 0.007* | $1.000 \pm 0.000$<br>*0.863 ± 0.008* | **1.000** ± 0.000<br>*0.858 ± 0.005* |
| $F^*$ | $0.877 \pm 0.006$<br>*0.877 ± 0.006* | $0.865 \pm 0.005$<br>*0.865 ± 0.005* | $0.866 \pm 0.007$<br>*0.866 ± 0.007* | $0.863 \pm 0.008$<br>*0.863 ± 0.008* | $0.858 \pm 0.005$<br>*0.858 ± 0.005* |

| | | | | | |
|---|---|---|---|---|---|
| **EV parent**, $\gamma = -0.3$ | | | | | |
| $n$ | 100 | 200 | 500 | 1000 | 2000 |
| $M^*$ | $8.116 \pm 1.817$<br>$0.810 \pm 0.003$ | **3.586** ± 0.728<br>$0.793 \pm 0.005$ | $1.001 \pm 0.001$<br>$0.779 \pm 0.006$ | $0.988 \pm 0.007$<br>*0.773 ± 0.005* | $0.947 \pm 0.004$<br>*0.765 ± 0.005* |
| $MM^*$ | **8.118** ± 1.926<br>$0.844 \pm 0.005$ | $1.429 \pm 0.440$<br>$0.809 \pm 0.006$ | **1.002** ± 0.003<br>$0.779 \pm 0.006$ | **1.000** ± **0.000**<br>*0.773 ± 0.005* | **1.000** ± **0.000**<br>*0.765 ± 0.005* |
| $GH^*$ | $0.875 \pm 0.007$<br>$0.828 \pm 0.003$ | $0.822 \pm 0.009$<br>$0.804 \pm 0.006$ | $0.779 \pm 0.006$<br>$0.779 \pm 0.006$ | $0.773 \pm 0.005$<br>*0.773 ± 0.005* | $0.765 \pm 0.005$<br>*0.765 ± 0.005* |
| $F^*$ | $0.875 \pm 0.007$<br>*0.875 ± 0.007* | $0.822 \pm 0.009$<br>*0.822 ± 0.009* | $0.779 \pm 0.006$<br>*0.779 ± 0.006* | $0.773 \pm 0.005$<br>*0.773 ± 0.005* | $0.768 \pm 0.003$<br>*0.765 ± 0.005* |

**Table 3**: Simulated *REFF*'s of the adaptive estimators of exceedance probabilities of $x_{n:n}$, and associated 95% CI's, for *extreme value* underlying parents.

| | | | | | |
|---|---|---|---|---|---|
| **EV parent**, $\gamma = -0.1$ | | | | | |
| $n$ | 100 | 200 | 500 | 1000 | 2000 |
| $M^*$ | $4.149 \pm 0.899$<br>$0.698 \pm 0.005$ | **5.265** ± 0.923<br>$0.700 \pm 0.009$ | **2.257** ± 0.293<br>$0.657 \pm 0.015$ | $1.259 \pm 0.021$<br>$0.590 \pm 0.011$ | $1.048 \pm 0.021$<br>$0.537 \pm 0.004$ |
| $MM^*$ | $3.343 \pm 1.369$<br>$0.733 \pm 0.008$ | $1.518 \pm 0.043$<br>$0.702 \pm 0.014$ | $1.453 \pm 0.294$<br>$0.641 \pm 0.015$ | $1.394 \pm 0.016$<br>$0.579 \pm 0.011$ | $1.371 \pm 0.011$<br>$0.533 \pm 0.004$ |
| $GH^*$ | **7.690** ± 1.944<br>$0.702 \pm 0.007$ | $3.349 \pm 1.020$<br>$0.680 \pm 0.013$ | $1.498 \pm 0.036$<br>$0.629 \pm 0.015$ | **1.422** ± 0.017<br>$0.570 \pm 0.011$ | **1.387** ± 0.012<br>$0.525 \pm 0.004$ |
| $F^*$ | $0.858 \pm 0.010$<br>*0.858 ± 0.010* | $0.793 \pm 0.016$<br>*0.793 ± 0.016* | $0.710 \pm 0.017$<br>*0.710 ± 0.017* | $0.637 \pm 0.013$<br>*0.637 ± 0.013* | $0.580 \pm 0.005$<br>*0.580 ± 0.005* |

| | | | | | |
|---|---|---|---|---|---|
| **EV parent**, $\gamma = -0.3$ | | | | | |
| $n$ | 100 | 200 | 500 | 1000 | 2000 |
| $M^*$ | $12.177 \pm 3.181$<br>$0.587 \pm 0.006$ | $5.746 \pm 1.475$<br>$0.601 \pm 0.018$ | $1.774 \pm 0.149$<br>$0.580 \pm 0.022$ | $1.111 \pm 0.040$<br>$0.545 \pm 0.029$ | $0.920 \pm 0.039$<br>$0.540 \pm 0.022$ |
| $MM^*$ | **24.493** ± 6.833<br>$0.702 \pm 0.015$ | $8.523 \pm 3.839$<br>$0.722 \pm 0.038$ | **6.527** ± 0.426<br>$0.692 \pm 0.039$ | **6.928** ± 0.343<br>$0.633 \pm 0.048$ | **4.865** ± 0.305<br>$0.623 \pm 0.034$ |
| $GH^*$ | $8.535 \pm 3.339$<br>$0.518 \pm 0.008$ | **9.516** ± 5.754<br>$0.583 \pm 0.024$ | $1.840 \pm 0.095$<br>$0.614 \pm 0.033$ | $0.812 \pm 0.047$<br>$0.586 \pm 0.041$ | $0.601 \pm 0.034$<br>$0.592 \pm 0.029$ |
| $F^*$ | $0.865 \pm 0.012$<br>*0.865 ± 0.012* | $0.787 \pm 0.029$<br>*0.787 ± 0.029* | $0.712 \pm 0.031$<br>*0.712 ± 0.031* | $0.641 \pm 0.040$<br>*0.641 ± 0.040* | $0.618 \pm 0.029$<br>*0.618 ± 0.029* |

On the basis of the simulated results, the adaptive estimation procedure seems to provide interesting results, in the sense that we have obtained *REFF* indicators reasonably high for small $n$ and all parameters of interest. Regarding *EVI*-estimation, and despite of the fact that it is not possible to claim that $MM^*$ has, for all models in $\mathcal{D}_\mathcal{M}(G_\gamma)$, $\gamma < 0$, the best performance among the four adaptive estimators in (4.5), it is clear that if we have to elect one of these four adaptive estimators, we are inclined to the choice of $MM^*$. This is particularly if the model is not a long way from an *EV* model, and we have a light indication for this underlying parent, not only on the basis of the undertaken parametric data analysis in Section 5.1, but also due to the nature of the data. This is the reason why in Section 5.3.2, we shall compute the final estimates of $\gamma$ on the basis of $MM^*$. Note however that, for small $n$, $GH^*$ is also a serious alternative. For the right endpoint estimation all adaptive estimators in (4.6) are almost equivalent, and we thus see no reason not to use also $x^*_{MM}$. A similar comment applies to the estimators of the exceedance probability (or equivalently, of the excellence indicator).

## 5. DATA ANALYSIS OF INDOOR ATHLETIC EVENTS

The data under analysis are related to three running and three jumping events, all for men, the 60 Metres Hurdles (60 MH), 400 (400 M) and 1500 Metres (1500 M), as well as the high jump (HJ), long jump (LJ) and pole vault (PV). The *sources* were `http://www.iaaf.org/statistics/toplists/index.htmx` and `http://hem.bredband.net/athletics/athletics_all-time_best.htm`. Data was collected until the end of 2007 and for any athlete only the best mark was taken into account. As mentioned before, we are dealing with right-tails. Consequently, for all running events we have converted *running times* into *speeds*, i.e., 10.00 seconds in the 60 MH (equal to 0.06 kilometers) is transformed to a speed of $3600 \times 0.06/10 = 21.6$ km/h. Like this, the higher the speed, the better, just as the higher the jump, the better. Contrarily to what has been done in Einmahl and Smeets (2009), we have not paid attention to doping related times, and we are conscious that slightly different estimates could then be obtained, despite of the usual robustness of the methods to a few outliers in the data.

### 5.1. Parametric data analysis

Prior to a semi-parametric analysis of the data, the most common framework of *statistics of extremes*, we shall proceed to a parametric data analysis, in the lines of Robinson and Tawn (1995) and Barão and Tawn (1999), who considered the annual best times in the women's 3000 m event. Also Smith (1988) has

proposed a maximum likelihood method of fitting models to a series of records, and applied his method to athletics records for the mile and the marathon. The attempts made in these papers to predict an ultimate world record are based on the development of top performances over time. This is not the case in this paper. Here, as in Einmahl and Magnus (2008), as well as in Einmahl and Smeets (2009), we are not interested in predicting the world record in the future. We are using only the top performances associated with a set of $n$ athletes, and consequently, our estimated ultimate record tells us what, in principle, is possible at this moment, given today's knowledge and material.

We first illustrate in Figures 4 and 5, the Gumbel QQ-plots associated with all data sets under analysis. In all figures we have thus plotted the points $\left(x_{i:n}, \; p_i^\Lambda = -\ln(-\ln(i/(n+1)))\right)$, $1 \le i \le n$, and proceeded to the fitting of a least-squares line.
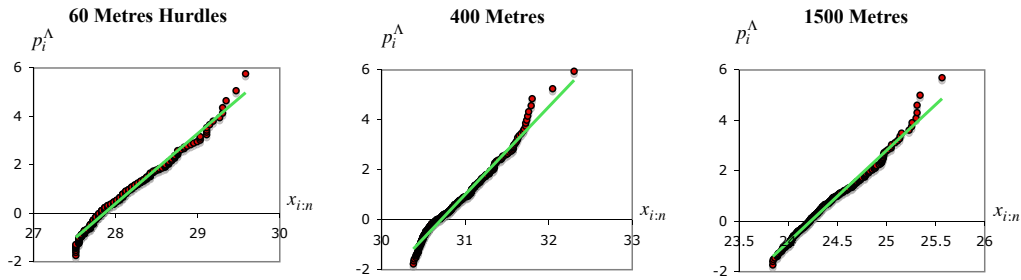


**Figure 4**: Gumbel QQ-plot related to the *running events* under analysis — 60 Metres Hurdles, 400 and 1500 Metres.
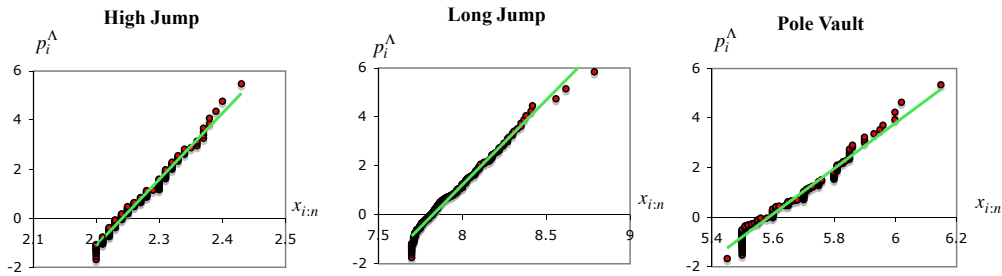


**Figure 5**: Gumbel QQ-plot related to the *jumping events* under analysis — High Jump, Long Jump and Pole Vault.

Apart from the Long Jump event, where $\gamma = 0$ can perhaps provide a reasonable fit to the right-tail, despite of a slight deviation of top o.s.'s smaller than the third largest value, all other events exhibit a light right-tail, i.e. a negative extreme value index, and consequently a finite right endpoint $x^*$.

Due to the fact that the observed data considered are already maxima, possibly of a small and dependent number of marks associated with any of the $n$ athletes, but the *extreme value* limiting law, in (1.1), is "robust" to changes of the i.i.d. assumption, we have first tried the fitting, through maximum likelihood, of an extreme value model $F(x; \lambda, \delta, \gamma) = G_\gamma\big((x-\lambda)/\delta\big)$, with $G_\gamma(x)$ given in (1.1). We have used the EVIR package in the $R$-software. The estimate of the right endpoint is then provided by $\hat{x}^* = \max(x_{n:n}, \hat{\lambda} - \hat{\delta}/\hat{\gamma})$, with $(\hat{\lambda}, \hat{\delta}, \hat{\gamma})$ the maximum likelihood estimates of the unknown parameters, $(\lambda, \delta, \gamma)$. The results obtained are presented in Table 4.

**Table 4**:  Maximum likelihood estimates of $(\lambda, \delta, \gamma, x^*)$ for an underlying model $G_\gamma\big((x-\lambda)/\delta\big)$, with $G_\gamma(x)$ given in (1.1): $'$– Km/h, $''$– metres.

| **Event** | $n$ | $(x_{1:n}, x_{n:n})$ | $\hat{\lambda}$ | $\hat{\delta}$ | $\hat{\gamma}$  (95% CI) | $\hat{x}^*$ |
|---|---|---|---|---|---|---|
| 60 MH | 312 | $(27.52, 29.59)'$ | 27.84 | 0.28 | $-0.21$  $(-0.328, -0.090)$ | 29.59 |
| 400 M | 380 | $(30.38, 32.31)'$ | 30.70 | 0.25 | $-0.15$  $(-0.277, -0.024)$ | 32.36 |
| 1500 M | 296 | $(23.84, 25.57)'$ | 24.23 | 0.26 | $-0.06$  $(-0.166, +0.042)$ | 28.33 |
| HJ | 235 | $(2.20, 2.43)''$ | 2.24 | 0.03 | $-0.09$  $(-0.223, +0.040)$ | 2.61 |
| LJ | 340 | $(7.70, 8.79)''$ | 7.81 | 0.11 | $-0.26$  $(-0.392, -0.130)$ | 8.79 |
| PV | 205 | $(5.45, 6.15)''$ | 5.58 | 0.09 | $-0.15$  $(-0.342, +0.041)$ | 6.21 |

As expected, all estimates of $\gamma$ are negative. But for the 1500 Metres, High Jump and Pole Vault, the upper limits of the associated 95% CI's are positive, suggesting that the value $\gamma = 0$ could possibly be adequate. The estimation of the right endpoint, which provides estimates equal to the maximum value in the data, the value $x_{n:n}$, for two of the athletic events, 60 Metres Hurdles and Long Jump, can be considered slightly problematic.

## 5.2. Fitting the extreme value model

In Figure 6, we picture in *light grey* the asymptotic 95% critical values (CV), $1.36/\sqrt{n}$, of the Kolmogorov–Smirnov statistic for testing a model without unknown parameters. The observed values of the Kolmogorov–Smirnov statistic, $KS_n := \max_{1 \leq i \leq n}\Big(\big|G_{\hat{\gamma}}\big((x_{i:n} - \hat{\lambda})/\hat{\delta}\big) - i/n\big|, \big|G_{\hat{\gamma}}\big((x_{i:n} - \hat{\lambda})/\hat{\delta}\big) - (i-1)/n\big|\Big)$  are pictured in *black*. The simulated 95% critical points of the Kolmogorov–Smirnov statistic, for testing an extreme value model $G_{\hat{\gamma}}\big((x - \hat{\lambda})/\hat{\delta}\big)$, have been based on 1000 runs, and are pictured in *grey*, showing again the "conservative property" of the Kolmogorov–Smirnov test — if we are led to rejection of a model without taking into account the maximum likelihood estimation of the parameters, we are *a fortiori* led to a rejection of the same model whenever we appropriately estimate the unknown parameters through the maximum likelihood approach.
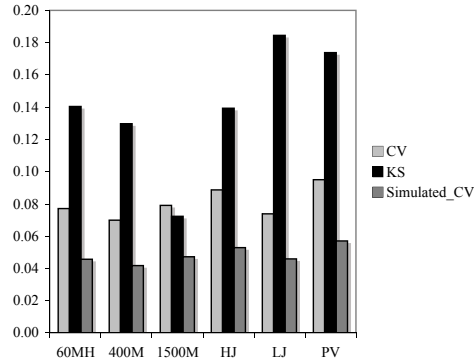
**Figure 6**:  Asymptotic critical values CV$= 1.36/\sqrt{n}$ (*light grey*), simulated critical values (*grey*) and observed values (*black*) of the Kolmogorov–Smirnov statistic, for all athletic events under analysis.

At the significance level $\alpha = 0.05$, the hypothesis of a (unified) *extreme value* model has thus been rejected by the Kolmogorov–Smirnov test for all data sets, as could also have been inferred graphically from Figure 7 and Figure 8, where we picture the empirical d.f., in grey, and the fitted extreme value d.f., in black.
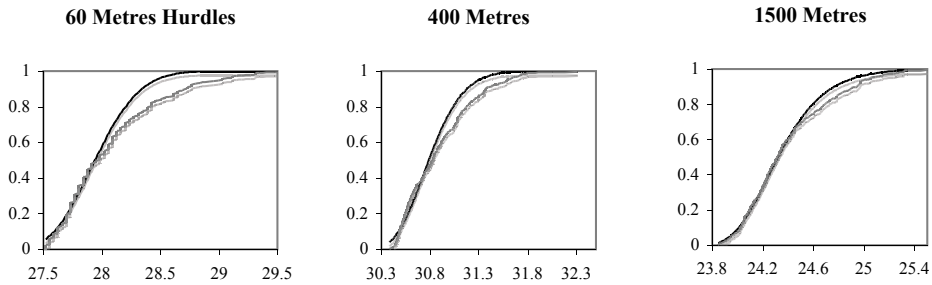


**Figure 7**:  Empirical d.f. (*grey*) and fitted extreme value d.f. (*black*) for the running events under analysis.
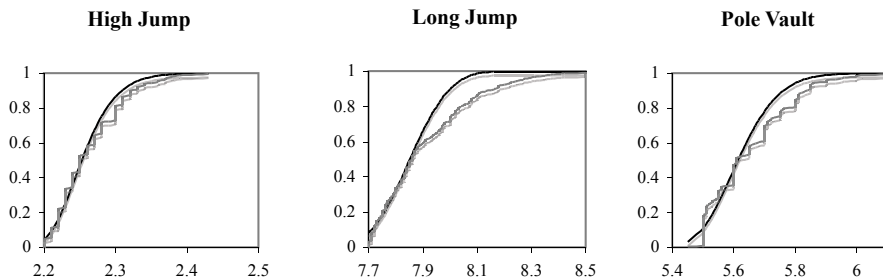


**Figure 8**:  Empirical d.f. (*grey*) and fitted extreme value d.f. (*black*) for the jumping events under analysis.

Alternative parametric models have even provided worse fitting results. There is thus a claim for the need of a semi-parametric data analysis, to be developed next, in Section 5.3.

## 5.3. A semi-parametric data analysis

### 5.3.1. Testing the extreme value index sign

As mentioned before, whenever we place ourselves under a semi-parametric framework, we assume only that (2.4) holds, or equivalently, that $F \in \mathcal{D}_{\mathcal{M}}(G_\gamma)$, for a certain $\gamma$, being $\gamma$ the primordial parameter of extreme events.

In many areas where extremes are relevant, the simplest case $\gamma = 0$ is often considered. Moreover, if we clearly think that $\gamma < 0$ or that $\gamma > 0$, we have specific procedures for the estimation of $\gamma$, possibly more reliable than the procedures valid for a general $\gamma \in \mathbb{R}$. Prior to a deeper semi-parametric analysis of the tail associated with this type of data, it thus seems sensible to test

$$
\begin{aligned}
& H_0\colon\ F \in \mathcal{D}_{\mathcal{M}}(G_\gamma)_{\gamma=0}\ \ \big(\text{or}\ \ F \in \mathcal{D}_{\mathcal{M}}(G_\gamma)_{\gamma\geq0}\big) \\
& \text{versus} \\
& H_1\colon\ F \in \mathcal{D}_{\mathcal{M}}(G_\gamma)_{\gamma<0}\ ,
\end{aligned}
$$

(5.1)

through the use of any semi-parametric test statistic.

We shall consider here two test statistics of a similar type, i.e. both based on the excesses over a high random threshold $X_{n-k:n}$, with $k$ satisfying (3.7). The first one was introduced by Greenwood (1946) and the second one by Hasofer and Wang (1992). These two statistics were further studied, under a semi-parametric framework, by Neves and Fraga Alves (2007). They are given by

$$
G_{k,n} := \frac{\frac{1}{k}\sum\limits_{i=1}^{k}\big(X_{n-i+1:n}-X_{n-k:n}\big)^2}{\left(\frac{1}{k}\sum\limits_{i=1}^{k}X_{n-i+1:n}-X_{n-k:n}\right)^2}
$$

and

$$
W_{k,n} := \frac{1}{k(G_{k,n}-1)}\ .
$$

Under the null hypothesis $H_0$ in (5.1) and extra mild conditions on the right-tail of $F$ and on the growth of $k = k_n$, they both have an asymptotic normal behaviour. More specifically,

(5.2) $\qquad G_{k,n}^* := \sqrt{k/4}\ \big(G_{k,n}-2\big)\big|_{F\in\mathcal{D}_{\mathcal{M}}(G_0)} \xrightarrow[n\to\infty]{d} N(0,1)$

and

(5.3) $\qquad W_{k,n}^* := \sqrt{k/4}\ \big(k\,W_{k,n}-1\big)\big|_{F\in\mathcal{D}_{\mathcal{M}}(G_0)} \xrightarrow[n\to\infty]{d} N(0,1)\ .$

Motivated by the important contribution of the maximum to the sum of the $k$ excesses, $X_{n-i+1:n} - X_{n-k:n}$, $1 \leq i \leq k$, Neves *et al.* (2006) introduced the following complimentary statistic,

$$R_{k,n} := \frac{X_{n:n} - X_{n-k:n}}{\frac{1}{k}\sum\limits_{i=1}^{k} X_{n-i+1:n} - X_{n-k:n}} \ ,$$

also considered in the analysis of the data under study. The asymptotic behaviour of $R_{k,n}$ is provided by the Gumbel d.f., $\Lambda = G_0$, with $G_\gamma$ given in (1.1). More specifically,

(5.4) $$R_{k,n}^* := R_{k,n} - \ln k \Big|_{F \in \mathcal{D}_\mathcal{M}(G_0)} \ \xrightarrow[n\to\infty]{d} \ Z \frown G_0 \ .$$

As a function of $k$ both $G_{k,n}^*$ and $R_{k,n}^*$ tend to have a slope with the sign of $\gamma$. The statistic $W_{k,n}^*$ works the other way round.

As an illustration, we present, in Figure 9, the sample paths of the three test statistics $G_{k,n}^*$, $W_{k,n}^*$ and $R_{k,n}^*$ in (5.2), (5.3) and (5.4), respectively, associated with the Long Jump. In this figure we also picture the quantiles $\left(\chi_{0.025}^\bullet, \chi_{0.975}^\bullet\right)$ of the standard normal $\Phi$, equal to $(-1.96, +1.96)$, and of the standard Gumbel $\Lambda \equiv G_0$, equal to $(-1.31, +3.68)$.
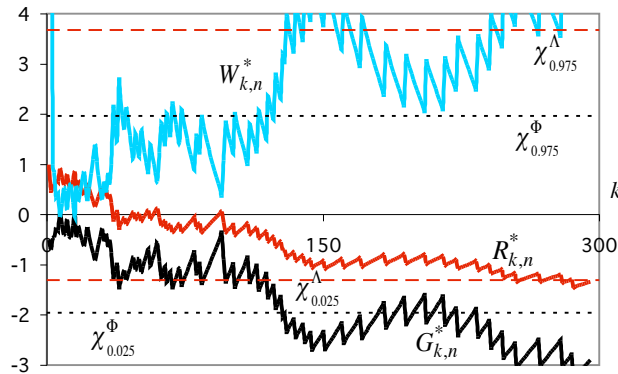


**Figure 9**:  Sample paths of the test statistics for the Long Jump event.

For all other data sets under analysis the graphs are quite similar, showing clearly a decreasing trend of $R_{k,n}^*$ and $G_{k,n}^*$ (with $G_{k,n}^*$ below $\chi_{0.025}^\Phi$ for a large number of $k$-values), as well as an increasing trend of $W_{k,n}^*$ (above $\chi_{0.975}^\Phi$ for moderate up to large values of $k$). Such a trend is mainly related to bias, but bias is strongly related to the extreme value index sign. These results provide a strong suggestion of a negative extreme value index, as expected. Despite of that, notice that, in Figure 9, the sample path of $R_{k,n}^*$ is within the 95% CI for almost all $k$-values. This was also expected, because it is well known (see, for instance, Neves and Fraga Alves, 2008) that $R_{k,n}^*$ tends to be a conservative test and the true value of $\gamma$ is for sure close to zero.

### 5.3.2. Semi-parametric estimates of the extreme value index and the right endpoint

In Table 5 we present a summary of the performed data analysis, with estimates and 95% CI's for the extreme value index $\gamma$. These estimates of $\gamma$ were obtained through the mixed moment ($MM$) estimates, computed at the value $k^*_{\min}$, in (4.4), i.e. they are the adaptive estimate $MM^*$ in (4.5).

**Table 5**: Estimates of the extreme value index, based on $M^*$: $'$ – Km/h, $''$ – metres.

| **Event** | $n$ | $(x_{1:n}, x_{n:n})$ | $MM^*$ (95% CI) | $k^*_{\min}$ |
|---|---|---|---|---|
| 60 MH | 312 | $(27.52, 29.59)'$ | $-0.34$ $(-0.469, -0.214)$ | 305 |
| 400 M | 380 | $(30.38, 32.31)'$ | $-0.26$ $(-0.445, -0.080)$ | 128 |
| 1500 M | 296 | $(23.84, 25.57)'$ | $-0.38$ $(-0.520, -0.241)$ | 275 |
| HJ | 235 | $(2.20, 2.43)''$ | $-0.32$ $(-0.468, -0.173)$ | 219 |
| LJ | 340 | $(7.70, 8.79)''$ | $-0.20$ $(-0.315, -0.087)$ | 296 |
| PV | 205 | $(5.45, 6.15)''$ | $-0.31$ $(-0.472, -0.151)$ | 182 |

In this semi-parametric data analysis, we have also considered the adaptive estimators $MM^{**}$ and $ML^{**}$, the estimators in (3.6) and (3.10), respectively, computed at the value $k^{**}_{\min}$, obtained through a minimization procedure of the type of the one in (4.4), but including also the $ML$-estimator. The reason for the consideration of the $ML$-estimator lies on the fact that in the region $-1/2 < \gamma < 0$, where the estimates lie, $\sigma^2_{ML} = (1 + \gamma)^2$ is smaller than $\sigma^2_{MM} = \sigma^2_M = (1 - \gamma)^2 (1 - 2\gamma)(1 - \gamma + 6\gamma^2)/\big((1 - 3\gamma)(1 - 4\gamma)\big)$ for all $\gamma$, with $\sigma_\bullet$ the asymptotic standard deviation in the asymptotic representation (4.1) (see Gomes and Neves, 2008, for further details). These estimates are presented in Table 6. For the LJ athletic event $k^*_{\min} = k^{**}_{\min}$. Then, the estimates $MM^* = MM^{**}$ and associated CI's are written in *italic*.

**Table 6**: Estimates of the extreme value index, based on $ML^{**}$ and $MM^{**}$.

| **Event** | $n$ | $ML^{**}$ (95% CI) | $MM^{**}$ (95% CI) | $k^{**}_{\min}$ |
|---|---|---|---|---|
| 60 MH | 312 | $-0.30$ $(-0.377, -0.216)$ | $-0.31$ $(-0.438, -0.186)$ | 294 |
| 400 M | 380 | $-0.22$ $(-0.351, -0.085)$ | $-0.26$ $(-0.434, -0.077)$ | 133 |
| 1500 M | 296 | $-0.32$ $(-0.400, -0.239)$ | $-0.31$ $(-0.440, -0.180)$ | 273 |
| HJ | 235 | $-0.29$ $(-0.387, -0.201)$ | $-0.31$ $(-0.456, -0.165)$ | 220 |
| LJ | 340 | $-0.17$ $(-0.262, -0.073)$ | *$-0.20$ $(-0.315, -0.087)$* | 296 |
| PV | 205 | $-0.29$ $(-0.396, -0.191)$ | $-0.30$ $(-0.458, -0.142)$ | 183 |

As it can be seen from Tables 5 and 6, there is only a small difference between $k^*_{\min}$ and $k^{**}_{\min}$, as expected. All semi-parametric $\gamma$-estimates at $k = k^{**}_{\min}$ are within the CI's provided in Table 5 and based on $MM^*$. Similarly, all estimates in Table 5 are within the CI's provided in Table 6. However, apart from the parametric estimates of $\gamma$ associated with the 400 Metres and Long Jump events, the parametric estimates in Table 4 are outside the CI's provided in Table 5, as well as the other way round. The parametric estimates are above the semi-parametric estimates for the six events considered. Note also that, contrarily to what generally happens, the values $k^*_{\min}$ and $k^{**}_{\min}$ are quite large, comparatively with the sample size $n$. This is essentially due to the fact that, for large $k$, the samples paths of the different estimators are reasonably stable as functions of $k$ and close to each other (a small bias, contrarily to the most common situations in practice) and volatile for small $k$ (large variance for small $k$, as usual).

Also as an illustration, we present, in Figure 10, the estimates $M \equiv M_{k,n}$, $GH \equiv GH_{k,n}$, $MM \equiv MM_{k,n}$, and $F \equiv F_{k,n}$ of $\gamma$, defined in (3.3), (3.5), (3.6) and (3.8), respectively, again for the Long Jump athletic event. We also picture the sample paths of the $\gamma$-estimator $ML \equiv ML_{k,n}$, in (3.10).
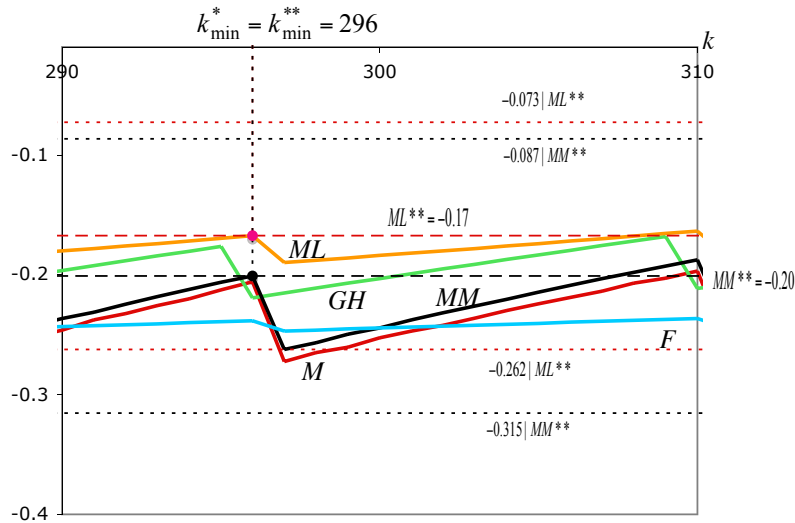


**Figure 10**: Sample paths of the extreme value index estimates under consideration, for the Long Jump event.

Analogously, and for the estimation of the right endpoint, apart from the adaptive estimators $\hat{x}^*_{MM}$, the estimators $\hat{x}^*_{k,n|\bullet}$, in (3.12), for $\bullet \equiv MM$, computed at the value $k^{*x}$, in (4.7), we have also considered the adaptive estimators $\hat{x}^{**}_{MM}$ and $\hat{x}^{**}_{ML}$, the estimators $\hat{x}^*_{k,n|\bullet}$ in (3.12) for $\bullet \equiv MM$ and $ML$, computed at the value $k^{**x}_{\min} \equiv k^{**x}$, obtained through a minimization procedure of the type of the one used for the adaptive endpoint estimators in (4.6), but including also the

*ML*-estimator. Similarly, and as mentioned before, for the estimation of the excellence indicator, we use the notation $\widehat{E}_\bullet^{**}$ for the estimator $\widehat{E}_{k,n}^\bullet$, in (3.13), computed at the value $k_{\min}^{**E} \equiv k^{**E}$. Next, in Table 7 we present the estimates of the right endpoints of the models underlying the different data sets under study.

**Table 7**: Estimates of the right endpoint: $'$ – km/h, $\bullet$ – minutes, $''$ – metres.

| **Event** | $x_{n:n}$ | $k^{*x}$ | $\hat{x}_{MM}^*$ | $k^{**x}$ | $\hat{x}_{ML}^{**}$ |
|---|---|---|---|---|---|
| 60 MH | $29.59'$ $(00:07.30)^\bullet$ | 53 | $29.81'$ $(00:07.25)^\bullet$ | 53 | $29.71'$ $(00:07.27)^\bullet$ |
| 400 M | $32.31'$ $(00:44.57)^\bullet$ | 32 | $32.45'$ $(00:44.37)^\bullet$ | 128 | $32.68'$ $(00:44.06)^\bullet$ |
| 1500 M | $25.57'$ $(03:31.18)^\bullet$ | 119 | $25.63'$ $(03:30.69)^\bullet$ | 119 | $25.70'$ $(03:3012)^\bullet$ |
| HJ | $2.43''$ | 219 | $2.44''$ | 219 | $2.46''$ |
| LJ | $8.79''$ | 144 | $8.84''$ | 281 | $9.12''$ |
| PV | $6.15''$ | 82 | $6.16''$ | 182 | $6.22''$ |

In Table 8 we present the estimates of the associated "excellence" indicators of the levels $x_{\mathrm{H}} = x_{n:n}$, provided in (3.13). Note that for all data sets we got $k^{*E} = k^{**E}$, smaller than expected for some of the data sets (60 MH, 1500 M and LJ).

**Table 8**: Estimates of an "excellence" indicator of the level $x_{n:n}$.

| **Event** | $k^{*x}$ | $\widehat{E}_{k^{*x},n}^{MM}$ | $k^{**x}$ | $\widehat{E}_{k^{**x},n}^{ML}$ | $k_{\min}^{*E} = k_{\min}^{**E}$ | $\widehat{E}_{MM}^* \mid \widehat{E}_{ML}^*$ |
|---|---|---|---|---|---|---|
| 60 MH | 53 | 0.66 | 53 | 0.67 | 11 | 0.62 $\mid$ 0.72 |
| 400 M | 32 | 0.98 | 128 | 0.88 | 148 | 0.99 $\mid$ 0.90 |
| 1500 M | 119 | 0.95 | 119 | 0.82 | 36 | 0.89 $\mid$ 0.81 |
| HJ | 219 | 0.98 | 219 | 0.90 | 222 | 0.91 $\mid$ 0.89 |
| LJ | 144 | 0.99 | 281 | 0.92 | 39 | 0.80 $\mid$ 0.78 |
| PV | 82 | 0.98 | 182 | 0.94 | 132 | 0.99 $\mid$ 0.94 |

Despite of slight discrepancies of the different estimates of the relevant parameters of extreme events, the results in Tables 5, 6, 7 and 8 mean that, under the present conditions, there are finite upper limits for all jumping events under analysis, as well as finite lower limits in the times associated with all running events under analysis. From the "excellence" indicators of the world records, we can say that the current 400 Metres, High Jump and Pole Vault world records are very good (indicators above 89%). The lowest "excellence" indicator, around 65%, corresponds to the 60 Metres Hurdles.

## ACKNOWLEDGMENTS

## REFERENCES

[1]    BARÃO, M.I. and TAWN, J. (1999). Extremal analysis of short series with out-
       liers: sea-levels and athletic records, *Applied Statistics*, **48**, 469–487.

[2]    BEIRLANT, J.; VYNCKIER, P. and TEUGELS, J. (1996). Excess functions and
       estimation of the extreme-value index, *Bernoulli*, **2**, 293–318.

[3]    BEIRLANT, J.; DIERCKX, G. and GUILLOU, A. (2005). Estimation of the
       extreme-value index and generalized quantile plots, *Bernoulli*, **11**(6), 949–970.

[4]    CAEIRO, C.; GOMES, M.I. and PESTANA, D. (2005). Direct reduction of bias
       of the classical Hill estimator, *Revstat*, **3**(2), 113–136.

[5]    DANIELSSON, J.; HAAN, L. DE; PENG, L. and DE VRIES, C.G. (2001). Using
       a bootstrap method to choose the sample fraction in the tail index estimation,
       *J. Multivariate Analysis*, **76**, 226–248.

[6]    DAVISON, A.C. (1984). *Modelling excesses over high thresholds*. In "Statistical
       Extremes and Applications" (J. Tiago de Oliveira, Ed.), D. Reidel, Dordrecht,
       Holland, 461–482.

[7]    DEKKERS, A.; EINMAHL, J. and DE HAAN, L. (1989). A moment estimator for
       the index of an extreme-value distribution, *Annals of Statistics*, **17**, 1833–1855.

[8]    DRAISMA, G.; HAAN, L. DE; PENG, L. and PEREIRA, T.T. (1999). A bootstrap-
       based method to achieve optimality in estimating the extreme-value index,
       *Extremes*, **2**, 367–404.

[9]    DREES, H.; FERREIRA, A. and DE HAAN, L. (2004). On maximum likelihood
       estimation of the extreme value index, *Annals of Applied Probability*, **14**, 1179–
       1201.

[10]   EINMAHL, J. and MAGNUS, J.R. (2008). Records in athletics through extreme-
       value theory, *J. American Statistical Association*, **103**, 1382–1391.

[11]   EINMAHL, J. and SMEETS, S.G.W.R. (2011). Ultimate 100-m world records
       through extreme-value theory, *Statistica Neerlandica*, **65**(1), 32–42.

[12]   FALK, M. (1995). Some best parameter estimates for distributions with finite
       endpoint, *Statistics*, **27**(1)–(2), 115–125.

[13]   FRAGA ALVES, M.I.; DE HAAN, L. and LIN, T. (2003). Estimation of the param-
       eter controlling the speed of convergence in extreme value theory, *Mathematical
       Methods of Statistics*, **12**(2), 155–176.

[14]   FRAGA ALVES, M.I.; GOMES, M.I.; DE HAAN, L. and NEVES, C. (2009).
       The mixed moment estimator and location invariant alternatives, *Extremes*, **12**,
       149–185.

[15]   GNEDENKO, B.V. (1943). Sur la distribution limite du terme maximum d'une série aléatoire, *Annals of Mathematics*, **44**(6), 423–453.

[16]   GOMES, M.I. and NEVES, C. (2008). Asymptotic comparison of the mixed moment and classical extreme value index estimators, *Statistics and Probability Letters*, **78**(6), 643–653.

[17]   GOMES, M.I. and OLIVEIRA, O. (2001). The bootstrap methodology in Statistics of Extremes: choice of the optimal sample fraction, *Extremes*, **4**(4), 331–358.

[18]   GOMES, M.I. and PESTANA, D. (2007). A sturdy reduced bias extreme quantile (VaR) estimator, *J. American Statistical Association*, **102**(477), 280–292.

[19]   GOMES, M.I.; MARTINS, M.J. and NEVES, M. (2007). Improving second order reduced bias extreme value index estimation, *Revstat*, **5**(2), 177–207.

[20]   GOMES, M.I.; DE HAAN, L. and HENRIQUES-RODRIGUES, L. (2008). Tail Index estimation for heavy-tailed models: accommodation of bias in weighted log-excesses, *J. Royal Statistical Society*, **B70**(1), 31–52.

[21]   GREENWOOD, M. (1946). The statistical study of infectious diseases, *J. Royal Statistical Society*, **A109**, 85–109.

[22]   HAAN, L. DE (1970). *On Regular Variation and its Application to the Weak Convergence of Sample Extremes*, Mathematical Centre Tract 32, Amsterdam.

[23]   HAAN, L. DE (1984). *Slow variation and characterization of domains of attraction*. In "Statistical Extremes and Applications" (J. Tiago de Oliveira, Ed.), D. Reidel, Dordrecht, Holland, 31–48.

[24]   HAAN, L. DE and FERREIRA, A. (2006). *Extreme Value Theory: an Introduction*, Springer, USA.

[25]   HASOFER, A. and WANG, J.Z. (1992). A test for extreme value domain of attraction, *J. American Statistical Association*, **87**, 171–177.

[26]   HILL, B.M. (1975). A simple general approach to inference about the tail of a distribution, *Annals of Statistics*, **3**(5), 1163–1174.

[27]   NEVES, C. and FRAGA ALVES, M.I. (2007). Semi-parametric approach to Hasofer-Wang and Greenwood statistics in extremes, *Test*, **16**, 297–313.

[28]   NEVES, C. and FRAGA ALVES, M.I. (2008). Testing extreme value conditions — an overview and recent approaches, *Revstat*, **6**(1), 83–100.

[29]   NEVES, C.; PICEK, J. and FRAGA ALVES, M.I. (2006). The contribution of the maximum to the sum of excesses for testing max-domains of attraction, *J. Statistical Planning and Inference*, **136**(4), 1281–1301.

[30]   ROBINSON, M.E. and TAWN, J. (1995). Statistics for exceptional athletic records, *Applied Statistics*, **44**, 499–511.

[31]   SMITH, R.L. (1987). Estimating tails of probability distributions, *Annals of Statistics*, **15**(3), 1174–1207.

[32]   SMITH, R.L. (1988). Forecasting records by maximum likelihood, *J. American Statistical Association*, **83**, 331–338.

[33]   ZHOU, C. (2009). Existence and consistency of the maximum likelihood estimator for the extreme value index, *J. Multivariate Analysis*, **100**(4), 794–815.

[34]   ZHOU, C. (2010). The extent of the maximum likelihood estimator for the extreme value index, *J. Multivariate Analysis*, **101**(4), 971–983.