

# On the extremes of randomly sub-sampled time series

A. Hall and M. G. Scotto\*

Department of Mathematics

University of Aveiro

Portugal

## Abstract

In this paper, we investigate the extremal properties of randomly sub-sampled stationary sequences. Motivation comes from the need to account for the effect of missing values on the analysis of time series and the comparison of schemes for monitoring systems with breakdowns or systems with automatic replacement of devices in case of failures.

*MSC:* 62-02; 60G70; 60G10

*Keywords:* Extreme value theory, integer-valued stationary sequences, sub-sampling, failure, extremal index.

---

\*Manuel Scotto, Departamento de Matemática, Universidade de Aveiro. Campo Universitário de Santiago, 3810-193 Aveiro, Portugal. Email: mscotto@ua.pt. Contract/grant sponsors: PTDC/MAT/64353/2006, FCT plurianual funding.

# 1 Introduction

The extremal properties of sub-sampling stationary sequences is a rapidly developing subject and it has been a topic of active research over the last years, mainly due to its wide applicability to the analysis of environmental and financial processes. Sub-sampling may occur according to some deterministic pattern, or may occur randomly. Much of the early work on this topic paid attention on the effect of deterministic sub-sampling on the extremal properties of stationary sequences; see Scotto (2005), Hall *et al.* (2004), Martins and Ferreira (2004), Ferreira and Martins (2003), Scotto *et al.* (2003), Hall and Scotto (2003), Scotto and Ferreira (2003), Scotto and Turkman (2002) and Robinson and Tawn (2000). In contrast, the effect of random sub-sampling has not received much attention in the literature. We refer to the work of Weissman and Cohen (1995) who considered the case of i.i.d. random sub-sampling as a particular case of some mixture models. More recently, Hall and Hüsler (2006) have obtained some generalizations of Weissman and Cohen's results for sequences where the sub-sampling pattern has a weak dependence structure.

One reason for the interest in extremes observed at random sampling rates comes from the need to compare schemes for monitoring systems with breakdowns or systems with automatic replacement of devices in case of failures. Examples are encountered, for instance, in ocean engineering. The probabilistic description of the wave climate in specific sites and ocean areas is an important prerequisite for the design and assessment of coastal and offshore structures. The wave climate is commonly described from time series of sea-state parameters, such as the significant wave height and the mean zero upcrossing period. These, as well as other sea-state parameters, provide information about the sea-state that has occurred and about the way the sea-state evolves with time. Most of the early available data has been collected by *waverider* buoys (at present, however, satellite data is becoming widely available and some climate descriptions are based on this type of data). An important aspect for a correct probabilistic description of the wave climate is to work with complete records of wave

measurements. Missing values, however, are frequently encountered in time series analysis of wave measurements, mainly when *waverider* buoys are used for collecting data sets. The main reasons are damage by shipping, *freak* waves which appeared out of a calm sea and a failure on the reading device. Similar problems arise in environmental studies. For example, extreme value analysis is of particular interest in assessing the impact of high air pollution levels, because air quality guidelines are formulated in terms of the high level of permitted emissions. This methodology has been used in the analysis of levels of ozone (Smith, 1989, Nui, 1997, and Tobias and Scotto, 2005) and nitrogen dioxide (Coles and Pan, 1996). Ozone data is usually collected from sampling stations integrated within a local automatic network for the control of atmospheric pollution in a specific area. In this case, missing observations appear when the equipment is not working properly or it is out of service.

As the title of the paper suggest, the aim of this work is to extend the results known for deterministic sub-sampled processes to random-generated sub-sampling processes. In particular, we investigate the maximum limiting distribution and its corresponding extremal index, when the underlying process is represented as a moving average driven by heavy-tailed innovations and the sub-sampling process is strongly mixing. Our results both exemplify some of the findings of Hall and Hüsler (2006) and offer more precise details for this particular class of models.

The examples given in the previous paragraphs illustrate the need to account for non-i.i.d. patterns of missing-values since, in general, when an equipment is out of order its recovery time may be considerably long. In this paper we also pay special attention to discrete-valued sequences. Motivation to include discrete data models comes from the need to account for the discrete nature of certain data sets, often counts of events, objects or individuals. Examples of applications can be found in the analysis of time series of count data that are generated from stock transactions (Quoreshi, 2006), where each transaction refers to a trade

between a buyer and a seller in a volume of stocks for a given price, and also in experimental biology (Zhou and Basawa, 2005), social science (McCabe and Martin, 2005), international tourism demand (Nordström, 1996, Garcia-Ferrer and Queralt, 1997, Brännäs *et al.* 2002, and Brännäs and Nordström, 2006), and queueing systems (Ahn *et al.* 2000).

The rest of the paper is organized as follows: Section 2 provides a background description of basic theoretical results related to conventional and non-negative integer-valued moving averages with regularly varying tails. Moreover, a suitable representation for the randomly sub-sampled process is described. In Section 3 we obtain the limiting distribution of the maximum term of the sub-sampled moving average sequence and the expression of its extremal index. Finally, in Section 4 the results are applied to conventional and discrete autoregressive processes.

## 2 Preliminaries

For the purpose of this work we shall consider stationary sequences  $\mathbf{X} = (X_n)_{n \in \mathbb{N}_0}$  of the form

$$X_n = \sum_{j=0}^{\infty} \beta_j * Z_{n-j}, \quad (1)$$

where  $\mathbf{Z} = (Z_n)_{n \in \mathbb{Z}}$  is an i.i.d. sequence of random variables (rv's) with distribution function  $F_Z$  belonging to the domain of attraction of the Fréchet distribution with parameter  $\alpha > 0$ , (hereafter  $F_Z \in D(\Phi_\alpha)$ ):

$$P(|Z_1| > x) = x^{-\alpha} L(x), \quad x > 0, \quad (2)$$

where  $L$  is slowly varying at infinity and

$$\lim_{x \rightarrow \infty} \frac{P(Z_1 > x)}{P(|Z_1| > x)} = p, \quad \lim_{x \rightarrow \infty} \frac{P(Z_1 < -x)}{P(|Z_1| > x)} = q, \quad (3)$$

for some  $p + q = 1$  with  $0 \leq p \leq 1$ . We further assume that the coefficients  $(\beta_j)_{j \in \mathbb{N}_0}$  are such that

$$\sum_{j=0}^{\infty} |\beta_j|^\delta < \infty, \quad \delta < \min(\alpha, 1). \quad (4)$$

Throughout the paper we consider two different cases:

(a) The  $*$ -operator denotes multiplication and  $\mathbf{Z}$  is an i.i.d. sequence of continuous rv's. In this case  $\mathbf{X}$  represents a conventional (i.e., continuous-valued) moving average model.

(b) The  $*$ -operator denotes *binomial thinning*, say  $\circ$ , and  $\mathbf{Z}$  represents an i.i.d. sequence of non-negative integer-valued rv's; that is

$$\beta \circ Z = \sum_{s=1}^Z B_s(\beta), \quad \beta \in [0, 1],$$

where  $(B_s(\beta))$  forms an i.i.d. sequence of Bernoulli rv's satisfying  $P[B_s(\beta) = 1] = \beta$ . In this case  $\mathbf{X}$  represents a discrete analogue of case (a). It is important to stress the fact that discreteness of the process  $\mathbf{X}$  is ensured by the  $\circ$ -operator since this operator incorporates the discrete nature of the variates and acts as the analogue of the standard multiplication used in the continuous-valued moving average model. Note that thinning is a random operation which reflects the behavior of many natural phenomena. For instance, if  $Z_n$  represents the number of individuals of a certain specie at time  $n$ ,  $\beta \circ Z_n$  will represent the number of survivors at the next time instant with  $\beta$  representing the probability of surviving. The concept of thinning is well known in classical probability theory and has been in use in the Bienaymé-Galton-Watson branching processes literature as well as in the theory of stopped-sum distributions.

We further consider within the discrete case the general class of models consisting of all stationary sequences defined by (1) in which all thinning operations involved are independent, for each  $n$ . Nevertheless, dependence is allowed to occur between the thinning operators

$\beta_j \circ Z_n$  and  $\beta_i \circ Z_n$ ,  $j \neq i$  (which belong to  $X_{n+j}$  and  $X_{n+i}$  respectively). We therefore obtain a rich class of discrete models which share some properties with the conventional case. For particular examples and estimation procedures see Brännäs and Hall (2001).

The tail properties of  $X_n$  have been studied by Davis and Resnick (1985) for the conventional case and by Hall (2001) for the discrete case. The result below summarises the tail behavior of the random variables  $W = \beta * Z$  and  $X_n$ , when  $F_Z \in D(\Phi_\alpha)$ .

**Theorem 2.1** *Let  $Z$  be a random variable with  $F_Z \in D(\Phi_\alpha)$ ,  $\alpha > 0$ .*

1. *For both meanings of the  $*$ -operator,  $F_W \in D(\Phi_\alpha)$  and*

(a) *for the conventional case*

$$\lim_{n \rightarrow \infty} \frac{1 - F_W(n)}{1 - F_Z(n)} = p(\beta^+)^{\alpha} + q(\beta^-)^{\alpha},$$

*with  $\beta^+ = \max(\beta, 0)$  and  $\beta^- = \max(-\beta, 0)$ ;*

(b) *for the discrete case*

$$\lim_{n \rightarrow \infty} \frac{1 - F_W(n)}{1 - F_Z(n)} = \beta^{\alpha}.$$

2. *If  $F_Z \in D(\Phi_\alpha)$  then, for both meanings of the  $*$ -operator,  $F_X \in D(\Phi_\alpha)$ , and for all  $\tau > 0$  and some sequence of constants  $(u_n)$*

$$\lim_{n \rightarrow \infty} n(1 - F_Z(u_n)) = \tau' \Rightarrow \lim_{n \rightarrow \infty} n(1 - F_X(u_n)) = \tau,$$

*with*

$$\tau' = \frac{\tau}{\sum_{j=0}^{\infty} p(\beta_j^+)^{\alpha} + q(\beta_j^-)^{\alpha}}. \quad (5)$$

*for the conventional case and*

$$\tau' = \frac{\tau}{\sum_{j=0}^{\infty} \beta_j^{\alpha}}, \quad (6)$$

*for the discrete case.*

The result above implies that every random variables  $Z_n$  contributes to the tail  $P(X > x)$ . This contribution depends on the size of the weight  $\beta_j$  for both meanings of the  $*$ -operator , as well as on the sign of the weight  $\beta_j$  in the conventional case.

Now we define the randomly sub-sampled sequence  $\mathbf{Y} = (Y_n)_{n \in \mathbb{N}_0}$  obtained from  $\mathbf{X}$  and induced through a strictly increasing function  $g(n) : \mathbb{N}_0 \rightarrow \mathbb{N}_0$  as follows:

$$Y_n = X_{g(n)}, \quad n \geq 0.$$

In addition, let  $\mathbf{U} = (U_n)_{n \in \mathbb{N}_0}$  be a Bernoulli stationary sequence independent of  $\mathbf{X}$  having marginal distribution with parameter  $\gamma$  ( $0 \leq \gamma \leq 1$ ). The  $U_n$ s are used as indicator variables that signal which observations are sampled whereas the  $g(\cdot)$  function gives the sampled time, that is the increasing sequence of  $ns$  for which  $U_n = 1$ . As an example take

$$U_1 = 1, U_2 = 0, U_3 = 1, U_4 = 0, U_5 = 0, U_6 = 1, U_7 = 1, \dots,$$

providing

$$g(1) = 1, g(2) = 3, g(3) = 6, g(4) = 7, \dots$$

The sequences  $\mathbf{U}$  considered in this paper will either be i.i.d. or strongly mixing.

The study of the extremal properties of stationary sequences is frequently based on the verification of appropriate dependence conditions which assure that the limiting distribution of the maximum term is of the same type as the limiting distribution of the maximum of i.i.d. rv's with the same marginal distribution  $F$ . For stationary sequences, usual conditions used in the literature are Leadbetter's  $D(u_n)$  condition (Leadbetter *et al.* 1983) and condition  $D^{(k)}(u_n)$ ,  $k \in \mathbb{N}$ , (Chernick *et al.* 1991). For completeness and reader's convenience the definition of conditions  $D(u_n)$  and  $D^{(k)}(u_n)$  are given below.

**Definition 2.1** *The condition  $D(u_n)$  is said to hold for a stationary sequence  $(X_n)_{n \in \mathbb{N}}$  with marginal distribution  $F$ , if for any integers  $i_1 < \dots < i_p < j_1 < \dots < j_q < n$  such that*

$j_1 - i_p \geq l_n$  we have

$$|F_{i_1, \dots, i_p, j_1, \dots, j_q}(u_n, \dots, u_n) - F_{i_1, \dots, i_p}(u_n, \dots, u_n)F_{j_1, \dots, j_q}(u_n, \dots, u_n)| \leq \alpha_{n, l_n}$$

with  $\alpha_{n, l_n} \rightarrow 0$  for some sequence  $(l_n)$ ,  $l_n = o(n)$ .

**Definition 2.2** The condition  $D^{(k)}(u_n)$ ,  $k \geq 1$ , holds for a stationary sequence  $(X_n)_{n \in \mathbb{N}}$  if there exist sequences  $(s_n)$  and  $(l_n)$  of integers, and  $(u_n)$  of reals, such that  $s_n \rightarrow \infty$ ,  $s_n \alpha_{n, l_n} \rightarrow 0$ ,  $\frac{s_n l_n}{n} \rightarrow 0$ , and

$$\lim_{n \rightarrow \infty} nP(X_1 > u_n \geq M_{2,k}, M_{k+1, r_n} > u_n) = 0, \quad (7)$$

where  $r_n = \left\lceil \frac{n}{s_n} \right\rceil$ ,  $M_{i,j} = \begin{cases} -\infty & \text{if } i > j \\ \max_{i \leq t \leq j} X_t & \text{if } i \leq j \end{cases}$ .

The main result is due to Chernick *et al.* (1991), in which the extremal index is computed by knowledge of the joint distribution of  $k$  consecutive terms.

**Theorem 2.2 (Chernick *et al.* 1991)** Suppose that for some  $k \geq 1$  the conditions  $D(u_n)$  and  $D^{(k)}(u_n)$  hold for  $u_n = u_n(\tau), \forall \tau > 0$ . Then, the extremal index of  $(X_n)_{n \in \mathbb{N}}$  exists and is equal to  $\theta$  iff

$$P(M_{2,k} \leq u_n | X_1 > u_n) \rightarrow \theta, \quad \text{as } n \rightarrow \infty, \quad \forall \tau > 0.$$

A convenient way to apply the above result may be through the following:

**Theorem 2.3 (Chernick *et al.* 1991)** Suppose  $(X_n)_{n \in \mathbb{N}}$  and  $(X_n^{(m)})_{n \in \mathbb{N}}$ ,  $m \geq 1$ , are stationary sequences defined on the same probability space such that for some sequence of constants  $\{u_n\}$

$$\lim_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} nP((1 - \epsilon)u_n < X_1 \leq (1 + \epsilon)u_n) = 0,$$

$$\lim_{m \rightarrow \infty} \limsup_{n \rightarrow \infty} nP(|X_1 - X_1^{(m)}| > \epsilon u) = 0, \quad \epsilon > 0.$$

Then



1. If condition  $D(u_n)$  holds for  $(X_n^{(m)})_{n \in \mathbf{N}}$ , for each  $m$ , then it holds for  $(X_n)_{n \in \mathbf{N}}$  as well.
2. If  $(X_n^{(m)})_{n \in \mathbf{N}}$  has extremal index  $\theta^{(m)}$ ,  $(X_k)_{n \in \mathbf{N}}$  has extremal index  $\theta$  iff

$$\lim_{m \rightarrow \infty} \theta^{(m)} = \theta.$$

### 3 Extremal behavior

The main task of this section is to derive the extremal behavior of the sub-sampled  $\mathbf{Y}$  process. The main result is formalized through the following theorem.

**Theorem 3.1** *Let  $\mathbf{X}$  be a moving average process defined as previously. Assume that  $(|\beta_j|)_{j \geq 0}$  forms a decreasing sequence. Consider the sub-sampled sequence  $\mathbf{Y}$  obtained by random sub-sampling according to an auxiliary stationary sequence  $\mathbf{U}$ . Furthermore, assume that  $F_Z \in D(\Phi_\alpha)$  satisfying  $\lim_{n \rightarrow \infty} n(1 - F_Z(u_n)) = \tau'$  with  $\tau'$  defined as in (5) for the conventional case, and defined as in (6) for the discrete case. Then, the distribution of  $Y_k$  satisfies*

$$\lim_{n \rightarrow \infty} n(1 - F_Y(u_n)) = \tau',$$

and it holds that:

1. The sequence  $\mathbf{Y}$  has extremal index

$$\theta_C = \frac{\sum_{j=1}^{\infty} P(g(2) - g(1) = j) (\sum_{i=0}^{j-1} p(\beta_i^+)^\alpha + q(\beta_i^-)^\alpha)}{\sum_{j=0}^{\infty} p(\beta_j^+)^\alpha + q(\beta_j^-)^\alpha}, \quad (8)$$

for the conventional case, with  $\beta_j^+$  and  $\beta_j^-$  defined as in Theorem 2.1, and

$$\theta_D = \frac{\sum_{j=1}^{\infty} P(g(2) - g(1) = j) \sum_{i=0}^{j-1} \beta_i^\alpha}{\sum_{j=0}^{\infty} \beta_j^\alpha}, \quad (9)$$

for the discrete case.

2. Moreover the limiting distribution of the maximum  $M_n(Y) = \max_{1 \leq g(k) \leq n} \{Y_k\}$  is given by

$$\lim_{n \rightarrow \infty} P(M_n(Y) \leq u_n) = \exp\{-\theta^* x^{-\alpha}\},$$

where  $\theta^*$  equals  $\theta_C$  for the conventional case and  $\theta_D$  for the discrete case.

*Proof.* By Theorem 2.3, to prove (8) we first obtain the extremal index of the the auxiliary finite-order sub-sampled moving average sequence

$$Y_k^{(m)} = \sum_{j=0}^m \beta_j * Z_{g(k)-j},$$

for fixed  $m > 0$ . We also temporarily take  $\beta_j = 0$  for  $j > m$ . Note that the local dependence  $D^{(m+1)}(u_n)$  condition trivially holds for  $\mathbf{Y}^{(m)} = (Y_k^{(m)})$ . For simplicity in notation we define  $M_{2,m+1}^{(m)} = \max_{2 \leq k \leq m+1} Y_k^{(m)}$ , and

$$\mu_{m+1}^{(m)}(u_n) = P(Y_1^{(m)} > u_n \geq M_{2,m+1}^{(m)}).$$

By Theorem 2.2 we have that the extremal index of the sequence  $\mathbf{Y}^{(m)}$ , for both meanings of the \*-operator, is given by

$$\theta^{(m)} = \lim_{n \rightarrow \infty} \frac{n \mu_{m+1}^{(m)}(u_n)}{nP(Y_1^{(m)} > u_n)}.$$

Moreover by arguments as in Chernick *et al.* (1991, Prop. 2.1)

$$\begin{aligned} \lim_{n \rightarrow \infty} n \mu_{m+1}^{(m)}(u_n) &= \lim_{n \rightarrow \infty} n \sum_{j=0}^m P(M_{2,m+1}^{(m)} \leq u_n, \beta_j * Z_{g(1)-j} > u_n) \\ &= \lim_{n \rightarrow \infty} n \sum_{j=0}^m [P(\beta_j * Z_{g(1)-j} > u_n) - P(M_{2,m+1}^{(m)} > u_n, \beta_j * Z_{g(1)-j} > u_n)]. \end{aligned}$$

Now

$$\begin{aligned} \lim_{n \rightarrow \infty} nP(M_{2,m+1}^{(m)} > u_n, \beta_j * Z_{g(1)-j} > u_n) &= \\ \lim_{n \rightarrow \infty} \left\{ nP(M_{2,m+1}^{(m)} > u_n, \beta_j * Z_{g(1)-j} > u_n, \bigvee_{\substack{0 \leq i' \leq m \\ 2 \leq t \leq m+1}} \beta_{i'} * Z_{g(t)-i'} > u_n) \right. & \end{aligned}$$

$$\begin{aligned}
& + \left. nP(M_{2,m+1}^{(m)} > u_n, \beta_j * Z_{g(1)-j} > u_n, \bigvee_{\substack{0 \leq i' \leq m \\ 2 \leq t \leq m+1}} \beta_{i'} * Z_{g(t)-i'} \leq u_n) \right\} \\
& = \lim_{n \rightarrow \infty} nP(\beta_j * Z_{g(1)-j} > u_n, \bigvee_{\substack{0 \leq i' \leq m \\ 2 \leq t \leq m+1}} \beta_{i'} * Z_{g(t)-i'} > u_n),
\end{aligned}$$

since as in Chernick *et al.* (1991, p. 842) and with the convention that  $\beta_j = 0$  for  $j > m$  it follows that

$$\lim_{n \rightarrow \infty} nP(M_{2,m+1}^{(m)} \leq u_n, \beta_j * Z_{g(1)-j} > u_n, \bigvee_{\substack{0 \leq i' \leq m \\ 2 \leq t \leq m+1}} \beta_{i'} * Z_{g(t)-i'} > u_n) = 0.$$

This makes explicit the precise way in which a single large  $Z$  asymptotically dominates the behavior of the maximum of the sequence  $\mathbf{Y}^{(m)}$ . For the conventional case, it follows that

$$\begin{aligned}
\lim_{n \rightarrow \infty} n\mu_{m+1}^{(m)}(u_n) & = \lim_{n \rightarrow \infty} n \sum_{j=0}^m P(\beta_j * Z_{g(1)-j} > u_n, \bigvee_{\substack{0 \leq i' \leq m \\ 2 \leq t \leq m+1}} \beta_{i'} * Z_{g(t)-i'} \leq u_n) \\
& = \lim_{n \rightarrow \infty} n \sum_{j=0}^m P(\beta_j * Z_1 > u_n, \bigvee_{2 \leq t \leq m+1} \beta_{g(t)-g(1)+j} * Z_1 \leq u_n) \\
& = \lim_{n \rightarrow \infty} n \sum_{j=0}^m [P(\bigvee_{2 \leq t \leq m+1} \beta_{g(t)-g(1)+j}^+ * Z_1 \leq u_n) \\
& + P(\bigvee_{2 \leq t \leq m+1} \beta_{g(t)-g(1)+j}^- * Z_1 \leq u_n) - P(\bigvee_{1 \leq t \leq m+1} \beta_{g(t)-g(1)+j}^+ * Z_1 \leq u_n) \\
& - P(\bigvee_{1 \leq t \leq m+1} \beta_{g(t)-g(1)+j}^- * Z_1 \leq u_n)] \\
& = \lim_{n \rightarrow \infty} n \sum_{j=0}^m [P(\beta_{g(2)-g(1)+j}^+ * Z_1 \leq u_n) + P(\beta_{g(2)-g(1)+j}^- * Z_1 \leq u_n) \\
& - P(\beta_j^+ * Z_1 \leq u_n) - P(\beta_j^- * Z_1 \leq u_n)],
\end{aligned}$$

since  $(|\beta_j|)_{j \in \mathbf{N}_0}$  forms a decreasing sequence with  $\beta_j = 0$  for  $j \geq m+1$ . Conditioning on  $V = g(2) - g(1)$  we obtain

$$\lim_{n \rightarrow \infty} n\mu_{m+1}^{(m)}(u_n) = \sum_{j=0}^m P(g(2) - g(1) = j) \left( \sum_{i=0}^{j-1} p(\beta_i^+)^{\alpha} + q(\beta_i^-)^{\alpha} \right).$$

Following Davis and Resnick (1985) the tail behavior of  $Y_k^{(m)}$  is given as follows:

$$\lim_{n \rightarrow \infty} \frac{P(Y_k^{(m)} > u_n)}{P(Z_1 > u_n)} = \sum_{j=0}^m p(\beta_j^+)^\alpha + q(\beta_j^-)^\alpha,$$

yielding

$$\theta^{(m)} = \frac{\sum_{j=1}^m P(g(2) - g(1) = j) (\sum_{i=0}^{j-1} p(\beta_i^+)^\alpha + q(\beta_i^-)^\alpha)}{\sum_{j=0}^m p(\beta_j^+)^\alpha + q(\beta_j^-)^\alpha}.$$

Finally as an application of Lemma 3.1 in Hall and Hüsler (2006, p. 547), condition  $D(u_n)$  holds for the sub-sampled sequence  $\mathbf{Y}$ , and hence by Theorem 2.3 the extremal index  $\theta_C$  is

$$\theta_C = \lim_{m \rightarrow \infty} \theta^{(m)}.$$

The discrete case follows as an application of the results given in Hall (2001) and Hall *et al.* (2004).

## 4 Examples

We now illustrate the effect of random sub-sampling on the extremal index of an AR(1) process

$$X_k = \beta * X_{k-1} + Z_k,$$

considering two different cases: (a) the conventional case with  $\beta \in (-1, 0)$  and the sequence of innovations  $\mathbf{Z}$  satisfying (2) and (3); and (b) the discrete case with  $\mathbf{Z}$  being a sequence of non-negative integer-valued rv's. This type of autoregressive sequence is known as *INteger-valued AutoRegressive process of order one* (INAR(1) in short) process and has been considered by several authors in the literature; see Aly and Bouzar (2005) for details. It is worth noting that in the former case, Hall and Hüsler's results can not be applied since condition  $D''(u_n)$  does not hold. In contrast, the AR(1) model with  $\beta \in (0, 1)$  satisfies  $D''(u_n)$  condition.

Furthermore, for the sequence  $\mathbf{U}$  two different cases will be considered:

- Independent and identically distributed failure instants: in this case  $\mathbf{U}$  forms an i.i.d sequence with  $P(U_k = 1) = \gamma = 1 - P(U_k = 0)$ , providing

$$P(g(2) - g(1) = j) = \gamma(1 - \gamma)^{j-1}, \quad j = 1, 2, \dots;$$

- Failures via a Markov Chain: within this framework  $\mathbf{U}$  forms an stationary Markov sequence defined by

$$\begin{cases} P(U_k = 1|U_{k-1} = 1) = \eta \\ P(U_k = 1|U_{k-1} = 0) = \nu \end{cases}.$$

This model defines a system where the probability of failure depends only on whether there occurred or not a failure just before. Given any values of  $\eta, \nu \in [0, 1]$  it is easy to obtain that

$$P(U_1 = 1) = \frac{\nu}{1 - \eta + \nu}.$$

Note that for a fixed value of  $\kappa = \frac{\nu}{1 - \eta + \nu} \in [0, 1]$ , the parameters  $\nu$  and  $\eta$  are not entirely arbitrary since if  $\kappa > 1/2$  then  $\eta \in [2 - 1/\kappa, 1]$ . The sequence  $\mathbf{U}$  is regenerative with finite mean duration of renewal epochs and hence it is strongly mixing. Moreover

$$P(g(2) - g(1) = j) = \begin{cases} \eta & j = 1 \\ (1 - \nu)^{j-2}(1 - \eta)\nu & j \geq 2 \end{cases}.$$

#### 4.1 Conventional case with negative parameter

In this case, the sub-sampled sequence  $\mathbf{Y}$  generated through the i.i.d sequence  $\mathbf{U}$  has extremal index

$$\theta_C = \frac{1 - \beta^{2\alpha}}{1 - (1 - \gamma)\beta^{2\alpha}}. \quad (10)$$

When  $\gamma = 1$ , (i.e., no sub-sampling), the extremal index in (10) becomes  $\theta_C = 1 - \beta^{2\alpha}$  which may be derived from the results given in Davis and Resnick (1985). Moreover, if the sub-sampled sequence  $\mathbf{Y}$  is generated through the stationary Markov sequence  $\mathbf{U}$ , the extremal

index becomes

$$\theta_C = \frac{1 - \beta^{2\alpha}[1 - (\nu - \eta)(1 - \beta^{2\alpha})]}{1 - (1 - \nu)\beta^{2\alpha}}.$$

## 4.2 Discrete case

In the discrete case, the extremal index of the sub-sampled sequence  $\mathbf{Y}$  generated through the i.i.d sequence  $\mathbf{U}$ , takes the form

$$\theta_D = \frac{1 - \beta^\alpha}{1 - (1 - \gamma)\beta^\alpha};$$

whereas for the stationary Markov sequence, the extremal index is given by

$$\theta_C = \frac{1 - \beta^\alpha[1 - (\nu - \eta)(1 - \beta^\alpha)]}{1 - (1 - \nu)\beta^\alpha}.$$

## References

- [1] AHN, S.; GYEMIN, L. and JONGWOO, J. (2000). Analysis of the M/D/1-type queue based on an integer-valued autoregressive process. *Oper. Res. Lett.* **27**, 235–241.
- [2] ALY, E.-E. and BOUZAR, N. (2005). Stationary solutions for integer-valued autoregressive processes. *Int. J. Math. Math. Sci.* **1**, 1-18.
- [3] BRÄNNÄS, K. and HALL, A. (2001). Estimation in integer-valued moving average models. *Appl. Stochastic Models Bus. Ind.* **17**, 277-291.
- [4] BRÄNNÄS, K.; HELLSTRÖM, and J. NORDSTRÖM, J. (2002). A new approach to modelling and forecasting monthly guest nights in hotels. *Int. J. Forecast.* **18**, 9–30.
- [5] BRÄNNÄS, K. and NORDSTRÖM, J. (2006). Tourist accomodation effects of festivals. *Tourism Economics* **12**, 291–302.
- [6] CHERNICK, M.; HSING, T. and MCCORMICK, W. (1991). Calculating the extremal index for a class of stationary sequences. *Adv. Appl. Prob.* **23**, 835-850.

- [7] COLES, S. G. and PAN, F. (1996). The analysis of extreme pollution levels: a case study. *J. Appl. Stats.* **23**, 333-348.
- [8] DAVIS, R. A. and RESNICK, S. I. (1985). Limit theory for moving averages of random variables with regularly varying tail probabilities. *Ann. Probab.* **13**, 179-195.
- [9] FERREIRA, H. and MARTINS, A. P. (2003). The extremal index of sub-sampled periodic sequences with strong local dependence. *Revstat - Statistical Journal*, **1** 15-24.
- [10] GARCIA-FERRER, A. and QUERALT, R. A. (1997). A note on forecasting international tourism demand in Spain. *Int. J. Forecast.* **13**, 539-549.
- [11] HALL, A. (2001). Extremes of integer-valued moving averages models with regularly varying tails. *Extremes* **4**, 219-239.
- [12] HALL, A. and HÜSLER, J. (2006). Extremes of stationary sequences with failures. *Stochastic Models* **22**, 537-557.
- [13] HALL, A. and SCOTTO, M. G. (2003). Extremes of sub-sampled integer-valued moving average models with heavy-tailed innovations. *Statist. Probab. Lett.* **63**, 97-105.
- [14] HALL, A.; SCOTTO, M. G. and FERREIRA, H. (2004). On the extremal behaviour of generalised periodic sub-sampled moving average models with regularly varying tails. *Extremes* **7**, 149-160.
- [15] LEADBETTER, M. R.; LINDGREN, G. and ROOTZÉN, H. (1983). *Extremes and Related Properties of Random Sequences and Processes*, Springer-Verlag, New York.
- [16] MARTINS, A. P. and FERREIRA, H. (2004). The extremal index of sub-sampled processes. *J. Statist. Plann. Inference* **1**, 145-152.
- [17] MCCABE, B. P. M. and MARTIN, G. M. (2005). Bayesian prediction of low count time series. *Int. J. Forecast.* **21**, 315-330.
- [18] NUI, X. F. (1997). Extreme value for a class of non-stationary time series with applications. *Ann. Appl. Probab.* **7**, 508-522.
- [19] NORDSTRÖM, J. (1996). Tourism satellite account for Sweden 1992-93. *Tourism Economics* **2**, 13-42.
- [20] QUORESHI, A. M. M. S. (2006). Bivariate time series modelling of financial count data. *Comm. Statist. Theory Methods* **35**, 1343-1358.

- [21] ROBINSON, M. E. and TAWN, J. A. (2000). Extremal analysis of processes sampled at different frequencies. *J. Roy. Statist. Soc. B* **62**, 117-135.
- [22] SCOTTO, M. G. (2005). Extremes of a class of deterministic sub-sampled processes with applications to stochastic difference equations. *Stochastic Process. Appl.* **115**, 417-434.
- [23] SCOTTO, M. G.; TURKMAN, K. F. and ANDERSON, C. W. (2003). Extremes of some sub-sampled time series. *J. Time Ser. Anal.* **24**, 579-590.
- [24] SCOTTO, M. G. and FERREIRA, H. (2003). Extremes of deterministic sub-sampled moving averages with heavy-tailed innovations. *Appl. Stochastic Models Bus. Ind.* **19**, 303-313.
- [25] SCOTTO, M. G. and TURKMAN, K. F. (2002). On the extremal behavior of sub-sampled solutions of stochastic difference equations. *Portugal. Math.* **59**, 267-282.
- [26] SMITH, R. L. (1989). Extreme value analysis of environmental time series: an application to trend detection in ground-level ozone. *Statist. Science* **4**, 367-393.
- [27] TOBIAS, A. and SCOTTO, M. G. (2005). Prediction of extreme ozone levels in Barcelona, Spain. *Environmental Monitoring and Assessment* **100**, 23-32.
- [28] WEISSMAN, I. and COHEN, U. (1995). The extremal index and clustering of high values for derived stationary sequences. *J. Appl. Prob.* **32**, 972-981.
- [29] ZHOU, J. and BASAWA, I. V. (2005). Least-squared estimation for bifurcation autoregressive processes. *Statist. Probab. Lett.* **74**, 77-88.