

1 Introduction

Multivariate extremes arise when one or more of rare extremes events occur simultaneously. These events are of prime interest to assess natural hazard stemming from heavy rainfall, wind storms and earthquakes since they are driven by joint extremes of a number meteorological variables. It is well known from the classical theory that multivariate distributions can be decomposed into two distinct parts : the analysis of marginal distributions and the analysis of the dependence structure described by the copula function. Results from the extreme-value theory show that the possible dependence structure of extremes have to satisfy certain constraints. Indeed, the dependence structure may be described in various equivalent ways ([4, 14, 42]) : by the exponent measure Λ ([3]), by the Pickands dependence function A ([40]), by the stable tail dependence function L ([27]), by the madogram ([38]), by the extreme value copula C ([24]).

The dependence structure between extreme observations can be complex and characterized by different notions from the ones arises in the classical theory as given afore. For this reason, recent works bring various notions to the framework of extreme such as sparsity ([23, 37, 47]), conditional independence and graphical models ([17, 22, 46]), dimensionality reduction ([8, 15]) and unsupervised learning ([12, 29]). In this work, we are concerned about clustering as a tool for learning the dependence structure of multivariate extreme and bridge important ideas from modern statistics and machine learning to the framework of extreme-value theory.

Loosely speaking, we are able to perform clustering in two distinguish cases : by partitionning the set $\{1, \dots, n\}$ of row indices or by partitioning with respect to column indices the set $\{1, \dots, d\}$. The first problem will be designated as the data clustering problem whereas the second corresponds to the variable clustering problem discussed here so far. In data clustering, clusters are clouds of observations and corresponds to respective realizations of one of the mixture distribution, which is a distribution on the whole \mathbb{R}^d . In this framework with i.i.d. replications, [41] shows the strong consistency of k -means clustering where the result was replicated in the context of extreme by [29] for spherical k -means.

The problem of variable clustering, see e.g. [6, 16], is that of grouping similar components of a d -dimensional random vector $\mathbf{X} = (X_1, \dots, X_d)$. Those groups are referred as clusters and unknown to the statistician who wants to recover them from $\mathbf{X}_1, \dots, \mathbf{X}_n$, n independent copies of \mathbf{X} . In this problem, the effort is to define cluster models relate to subsets of components X_j , of $\mathbf{X} \in \mathbb{R}^d$. In this framework, we do not longer want to cluster independent entities but, a contrario, to cluster those who are strongly dependent. Variable clustering is of prime interest in weather extremes, with examples stemming from regionalisation, for instance [2, 5, 45], where one observes spatial phenomena at finitely many sites. An interesting special case is to cluster those sites according to their extremal dependencies. They include applications of K -means or hierarchical clustering with a dissimilarity designed for extremes. The statistical properties of those procedures have received a very limited amount of investigation. It is not currently known what probabilistic models on \mathbf{X} can

be estimated by these techniques. We will consider here model-based clustering where population-level clusters are clearly defined, offering interpretability and a benchmark to assess the performance of a peculiar clustering algorithm.

In this work, we propose the AI-block model as a model for variable clustering in extreme-value theory and show that the clusters given by this model are uniquely defined. We then motivate and develop an algorithm tailored to the model with the help of the SECO metric. We thus analyze its performance in terms of exact cluster recovery, for minimally separated clusters, under appropriately defined cluster separation metric.

We use the following notations throughout the paper. All bold letters \mathbf{x} corresponds to vector in \mathbb{R}^d . By considering $B \subseteq \{1, \dots, d\}$, we denote the $|B|$ -subvector of \mathbf{x} by $\mathbf{x}^{(B)} = (x_j)_{j \in B}$. Similarly, let G a cumulative distributive function on $[0, 1]^d$, $G^{(B)}$ is defined as

$$G^{(B)}(\mathbf{x}^{(B)}) = G(\mathbf{1}, \mathbf{x}^{(B)}, \mathbf{1}), \quad (x_j)_{j \in B} \in [0, 1]^{|B|},$$

where $(\mathbf{1}, \mathbf{x}^{(B)}, \mathbf{1})$ has j th component equals to $x_j \mathbb{1}_{\{j \in B\}} + \mathbb{1}_{\{j \notin B\}}$. In a similar way, we note $(\mathbf{0}, x^{(B)}, \mathbf{0})$ the vector in \mathbb{R}^d which equals x_j if $j \in B$ and 0 otherwise. We will denote the d consecutive integer set starting from 1 as $\llbracket d \rrbracket$. Weak convergence of processes are denoted by \rightsquigarrow . The notation δ_x corresponds to the dirac measure at x . We define by $\mathbf{X} \in \mathbb{R}^d$ a random vector with law G . Let $O = \{O_k\}_{k=1, \dots, K}$ be a partition of $\{1, \dots, d\}$ into K groups and let $s : \{1, \dots, d\} \rightarrow \{1, \dots, K\}$ be an index assignment function defined by $O_k = \{a \in \{1, \dots, d\} : s(a) = k\} = \{i_{k,1}, \dots, i_{k,d_k}\}$ with $d_1 + \dots + d_K = d$.

In the present paper, we present in Section 2 some background notions in extreme-value theory necessary for our analysis. All these notions introduced, we describe AI-block models designed for variable clustering. Central to our work is a new metric for variable clustering in AI-block models. Our metric is based on the Sum of Extremal COefficient (SECO) between several groups of variables, a measure that attempts to capture how groups of variables are independent relative to their extremes. We discuss this in detail in Section 3. In Section 4, we develop a new clustering algorithm based on the extremal correlation coefficient and a threshold. It has moderate computational complexity, polynomial in d . We prove that this algorithm can recover correctly the target partition, with high probability if the threshold is calibrated accordingly to a certain scale. Furthermore, we use the SECO metric as a tool to calibrate the threshold in our algorithm in Section 5. All proofs are deferred to the Appendix.

2 A model for variable clustering

2.1 Extreme value theory

Consider $\mathbf{Z} = (Z_1, \dots, Z_d)$ a d -dimensional random vector with law F and $\mathbf{Z}_i = (Z_{i,1}, \dots, Z_{i,d})$, $i = 1, \dots, n$ be independent copies of \mathbf{Z} . By denoting the component-wise maxima by $\mathbf{M}_n = (\max_{i=1, \dots, n} Z_{i,1}, \dots, \max_{i=1, \dots, n} Z_{i,d})$. The vector \mathbf{Z} is said to be in max-domain of attraction of the random vector $\mathbf{X} = (X_1, \dots, X_d)$, denoted as $F \in D(G)$, if for any $\mathbf{x} = (x_1, \dots, x_d)$,

$$\lim_{n \rightarrow \infty} \{F(\mathbf{a}_n \mathbf{x} + \mathbf{b}_n)\}^n = G(\mathbf{x}), \quad (1)$$

where $\mathbf{a}_n > \mathbf{0}$ (which means that $a_{i,n} > 0$ for every i and n) and $\mathbf{b}_n \in \mathbb{R}^d$. In this case, \mathbf{X} is max-stable with Generalized Extreme Value (GEV) margins and we may write

$$\mathbb{P}\{\mathbf{X} \leq \mathbf{x}\} = \exp\{-\Lambda(E \setminus [0, \mathbf{x}])\},$$

where Λ is a Radon measure on the cone $E = [0, \infty)^d \setminus \{\mathbf{0}\}$. This condition is equivalent to the notion of regular variation, that is there exists a sequence $0 < \mathbf{a}_n \rightarrow \infty$ and a limit measure such that

$$n\mathbb{P}\{\mathbf{a}_n^{-1}\mathbf{X} \in \cdot\} \xrightarrow[n \rightarrow \infty]{v} \Lambda(\cdot),$$

with \xrightarrow{v} denotes the vague convergence.

Those notions can be translated, as in the classical theory, in terms of copula. A d -dimensional random vector $\mathbf{X} = (X_1, \dots, X_d)$ follows the law of a multivariate extreme-value distribution if its one dimensional marginal distributions $G_j(x) = \mathbb{P}\{X_j \leq x\}$ for all $x \in \mathbb{R}$ and $j \in \{1, \dots, d\}$ are GEV distributions, and the joint distribution can be written, for all $\mathbf{x} \in \mathbb{R}^d$, in the form

$$G(\mathbf{x}) = C(G_1(x_1), \dots, G_d(x_d)), \quad (2)$$

where C is an extreme value copula, i.e., for all $\mathbf{u} \in (0, 1]^d$

$$C(\mathbf{u}) = \exp\{-L(-\ln(u_1), \dots, -\ln(u_d))\},$$

with L is known as the stable tail dependence function (see [24] for an overview of extreme value copulae). As it is a homogeneous function of order 1, i.e. $L(a\mathbf{z}) = aL(\mathbf{z})$ for all $a > 0$, we have, for all $\mathbf{z} \in [0, \infty)^d$,

$$L(\mathbf{z}) = (z_1 + \dots + z_d)A(\mathbf{t}),$$

with $t_j = z_j/(z_1 + \dots + z_d)$ for $j \in \{2, \dots, d\}$, $t_1 = 1 - t_2 - \dots - t_d$, and A is the restriction of L into the d -dimensional unit simplex, viz.

$$\Delta_{d-1} = \{(v_1, \dots, v_d) \in [0, 1]^d : v_1 + \dots + v_d = 1\}.$$

The function A is known as the Pickands dependence function and is often used to quantify the extremal dependence among the element of X . Indeed, A satisfies the constraints $1/d \leq \max(t_1, \dots, t_d) \leq A(\mathbf{t}) \leq 1$ for all $\mathbf{t} \in \Delta_{d-1}$, with lower and upper bounds corresponding to the complete dependence and independence cases.

2.2 AI-block models

Motivated by a rich set of applications, we consider variable clustering as the initial dimension reduction step applied to the observed vector $\mathbf{X} = (X_1, \dots, X_d)$. These models are build on the assumption that the observed variables $\mathbf{X} = (X_1, \dots, X_d)$ can be partitionned into K unknown clusters $O = \{O_1, \dots, O_K\}$ such that variables in the same cluster are as dependent as possible and the clusters are mutually independent. This structure of dependence has been observed to hold, empirically, recent studies have shown that in many applications such as spatial precipitation (see [30, 32]) or water discharges in river network ([20]). We define here a population-level cluster as a group of variables that shares the same extremal dependence structure within the cluster and is independent from the other clusters.

To keep the presentation focus, let us consider $\mathbf{X}^{(O_1)}, \dots, \mathbf{X}^{(O_K)}$ be extreme value random vectors with extreme value copulae $C^{(O_1)}, \dots, C^{(O_K)}$ respectively. We suppose that $\mathbf{X}^{(O_k)}$ and $\mathbf{X}^{(O_j)}$ are mutually asymptotically independent if $k \neq j$, i.e., the corresponding components of the limit vector in (1) are mutually independent. Let us define the following function :

$$\begin{aligned} C_{\Pi} : [0, 1]^d &\longrightarrow [0, 1] \\ \mathbf{u} &\longmapsto \prod_{k=1}^K C^{(O_k)}(u_{i_{k,1}}, \dots, u_{i_{k,d_k}}). \end{aligned}$$

We want to show that C_{Π} is an extreme value copula associated to $\mathbf{X} = (\mathbf{X}^{(O_1)}, \dots, \mathbf{X}^{(O_K)})$ in order that all objects we use below are well-defined, in particular the existence of a *stable tail dependence function* associated to C_{Π} .

Lemma 1. *C_{Π} is an extreme value copula associated to the random vector \mathbf{X} .*

A direct consequence of this lemma is that if \mathbf{X} admits a copula C_{Π} , it is an extreme-value one. In particular, there exists a stable tail dependence function L say, and one can show that L can be expressed in a convenient way such as :

$$L(z_1, \dots, z_d) = \sum_{k=1}^K L^{(O_k)}(\mathbf{z}^{(O_k)}), \quad \mathbf{z} \in [0, \infty)^d, \quad (3)$$

where $L^{(O_1)}, \dots, L^{(O_K)}$ are the stable tail dependence function associated to the extreme value copulae $C^{(O_1)}, \dots, C^{(O_K)}$ respectively. Furthermore, this model is a particular form of the nested extreme value copula and is the object of the remark below.

Remark 1. Equation (3) can restated as

$$L(\mathbf{z}) = L^{(0)} \left(L^{(O_1)} \left(z^{(O_1)} \right), \dots, L^{(O_K)} \left(z^{(O_K)} \right) \right),$$

where $L^{(0)}(z_1, \dots, z_K) = \sum_{k=1}^K z_k$ is a stable tail dependence function corresponding to the asymptotic independence. Since C is an extreme value copula by Lemma 1, we obtain that C is also a nested extreme value copula as formulated in [26].

The following remark states the form of the stable tail dependence function in terms of regular variation. This corresponding statement might be found in the literature, see Lemma 7.2 of [43] or Theorem 1.28 of [33] for example.

Remark 2. Let $K = 2$, consider $\mathbf{X}^{(O_1)} \in \mathbb{R}^{d_1}$ and $\mathbf{X}^{(O_2)} \in \mathbb{R}^{d_2}$ defined in the same probability space, independent, and satisfy the regular variation assumption

$$n\mathbb{P} \left\{ a_n^{-1} \mathbf{X}^{(O_1)} \in \cdot \right\} \xrightarrow[n \rightarrow \infty]{v} \nu_1(\cdot), \quad n\mathbb{P} \left\{ a_n^{-1} \mathbf{X}^{(O_2)} \in \cdot \right\} \xrightarrow[n \rightarrow \infty]{v} \nu_2(\cdot),$$

with the same sequence $0 < a_n \rightarrow \infty$. then the distribution tail of $\mathbf{X} = (\mathbf{X}^{(O_1)}, \mathbf{X}^{(O_2)})$ is also regularly varying with

$$n\mathbb{P} \left\{ a_n^{-1} \mathbf{X} \in \cdot \right\} \xrightarrow[n \rightarrow \infty]{v} \nu(\cdot),$$

where

$$\nu(d\mathbf{x}^{(O_1)}, d\mathbf{x}^{(O_2)}) = \nu_1(d\mathbf{x}^{(O_1)})\delta_0(d\mathbf{x}^{(O_2)}) + \delta_0(d\mathbf{x}^{(O_1)})\nu_2(d\mathbf{x}^{(O_2)}).$$

With all notations and definitions previously introduced, we now are able to state the definition of the considered model here.

Definition 1 (Asymptotic Independence-block model). The random vector $\mathbf{X} = (\mathbf{X}^{(O_1)}, \dots, \mathbf{X}^{(O_K)})$ follows an AI-block model if $\mathbf{X}^{(O_k)} = (X_{i_{k,1}}, \dots, X_{i_{k,d_k}})$ are extreme value random vectors for $k \in \{1, \dots, K\}$ and are mutually independent if $k \neq j$.

Notice that, when $K = 1$, the definition of the AI-block model thus reduces to $\mathbf{X} = (X_1, \dots, X_d)$ is an extreme value random vector, which is trivially obtained by the definition.

Let $\mathcal{L}(\mathbf{X}) = \{O : \mathbf{X} \text{ is AI-block model}\}$ which is nonempty and finite so it does have maximal elements, we note $\mathbf{X} \sim O$ if \mathbf{X} follows an AI-block model. We introduce the following partial order on sets, let $O = \{O_k\}_k$, $S = \{S_{k'}\}_{k'}$ be two partitions of $\{1, \dots, d\}$. We say that S is a sub-partition of O if for each k' there exists k such that $S_{k'} \subseteq O_k$. We define the partial order \leq between two partitions O, S of $\{1, \dots, d\}$ by $O \leq S$ if S is a sub-partition of O . For any partition $O = \{O_k\}_{1 \leq k \leq K}$, we write $j \overset{O}{\sim} k$ if there exist $k \in \{1, \dots, K\}$ such that $j, k \in O_k$.

Definition 2. For any two partitions O, S of $\{1, \dots, d\}$, we define $O \cap S$ as the partition induced by the equivalence relation $j \overset{O \cap S}{\sim} k$ iff $j \overset{O}{\sim} k$ and $j \overset{S}{\sim} k$.

Checking that $j \stackrel{O \cap S}{\sim} k$ is an equivalence relation is straightforward. With this definition, we have the following interesting properties that lead to the desired result, the identifiability of the introduced AI-block models.

Theorem 1. *Let \mathbf{X} be an extreme value random vector, then the following properties hold :*

1. *Consider $O \leq S$. Then $\mathbf{X} \sim S$ implies $\mathbf{X} \sim O$.*
2. *$O \leq O \cap S$ and $S \leq O \cap S$*
3. *If $\mathbf{X} \sim O$ and $\mathbf{X} \sim S$, then $\mathbf{X} \sim O \cap S$.*
4. *The set $\mathcal{L}(\mathbf{X})$ has a unique maximum $\bar{O}(\mathbf{X})$, with respect to the partition partial order.*

We refer to Appendix for a proof of this Theorem. It shows that the maximum partition for which \mathbf{X} is an AI-block model always exists and its structure is intrinsic of the definition of \mathbf{X} . The maximal partition of $\bar{O}(\mathbf{X})$ matches our intuition regarding what would constitute a reasonable clustering target for these models, *i.e.* the finest partition $\bar{O}(\mathbf{X})$ where \mathbf{X} is an AI-block model as stated in the introduction of [44]. With a slight abuse of notation, we will write $\bar{O}(\mathbf{X})$ as \bar{O} .

In an AI-block model, we may restrict Equation (3) to the simplex and thus, expressing the stable tail dependence function in terms of the Pickands dependence function. Furthermore, its expression is obtained as a convex combination of the Pickands dependence function $A^{(O_1)}, \dots, A^{(O_K)}$ as follows

$$\begin{aligned} A(t_1, \dots, t_d) &= \frac{1}{z_1 + \dots + z_d} \left[\sum_{k=1}^K (z_{i_{k,1}} + \dots + z_{i_{k,d_k}}) A^{(O_k)}(\mathbf{t}^{(O_k)}) \right] \\ &= \sum_{k=1}^K w^{(O_k)}(\mathbf{t}) A^{(O_k)}(\mathbf{t}^{(O_k)}) =: A_O \end{aligned}$$

with $t_j = z_j / (z_1 + \dots + z_d)$ for $j \in 2, \dots, d$ and $t_1 = 1 - t_2 - \dots - t_d$, $w^{(O_k)}(\mathbf{t}) = z_{i_{k,1}} + \dots + z_{i_{k,d_k}} / (z_1 + \dots + z_d)$ for $k \in \{2, \dots, K\}$ and $w_1 = 1 - w^{(O_2)}(\mathbf{t}) - \dots - w^{(O_K)}(\mathbf{t})$ and $t_{i_{k,\ell}} = z_{i_{k,\ell}} / (z_{i_{k,1}} + \dots + z_{i_{k,d_k}})$ for $k \in \{1, \dots, K\}$ and $\ell \in \{1, \dots, d_k\}$. The function A is still a Pickands dependence function as a convex combination of Pickands dependence function (see p. 123 of [19]).

When considering independence between random variables, we know that $A(\mathbf{t}) \leq 1$ for $\mathbf{t} \in \Delta_{d-1}$ where inequality stands for asymptotic independence between all random variables. A more general statement can also be formulated considering random vectors, where the former case directly comes down by taking $d_1 = \dots = d_K = 1$.

Proposition 1. *In general case, we have for every $\mathbf{t} \in \Delta_{d-1}$,*

$$(A_O - A)(\mathbf{t}) \geq 0,$$

with equality if and only if $\mathbf{X} \sim O$.

Remark 3. A beautiful way to state asymptotic independence between random vectors is by the mean of the exponent measure. Taking notations of Remark 2, $\mathbf{X}^{(O_1)}$ is independent of $\mathbf{X}^{(O_2)}$ if and only if, for $\mathbf{y} > \mathbf{0}$

$$\nu \left\{ \mathbf{x} \in E, \mathbf{x}^{(O_1)} > \mathbf{y}^{(O_1)}, \mathbf{x}^{(O_2)} > \mathbf{y}^{(O_2)} \right\} = 0.$$

Thus, the exponent measure ν concentrates on

$$]0, \infty[^{d_1} \times \{\mathbf{0}\}^{d_2} \cup \{\mathbf{0}\}^{d_1} \times]0, \infty[^{d_2}.$$

These conditions generalize straightforwardly those stated in Proposition 5.24 of [42].

3 The SECO similarity metric

Let $\mathbf{X} \sim O$ has a block structure with associated Pickands A which can be expressed as a convex combination of the K Pickands $A^{(O_k)}$. One expect a consistent clustering could be possible when

$$A_O - A \tag{4}$$

is small enough. To motivate this metric, notice that in AI-block model, if a given partition $O \in \mathcal{L}(\mathbf{X})$, then (4) is equal to 0 while it is strictly greater than 0 if $O \notin \mathcal{L}(\mathbf{X})$. This statement is precised by the following proposition.

Proposition 2. *Let $\mathbf{X} \sim O$ and $S_{k'}$ an arbitrary partition of $\llbracket d \rrbracket$. We thus have for $\mathbf{t} \in \Delta_{d-1}$*

$$0 \leq \sum_{l=1}^L \sum_{k=1}^K w^{(O_k \cap S_{k'})}(\mathbf{t}) A^{(O_k)}(\mathbf{0}, \mathbf{t}^{(O_k \cap S_{k'})}, \mathbf{0}) - A(\mathbf{t}), \quad \forall \mathbf{t} \in \Delta_{d-1}, \tag{5}$$

where

$$w^{(O_k \cap S_{k'})}(\mathbf{t}) = \sum_{j \in S_{k'} \cap O_k} t_j, \quad \forall k \in \{1, \dots, K\}, k' \in \{1, \dots, K'\}.$$

Furthermore, the right hand side of the inequality is equal to zero if and only if $S \leq O$.

Proposition 2 gives the theoretical value if $\mathbf{X} \sim O$ for an arbitrary partition which is stricly greater than 0 if the given partition does not belong to $\mathcal{L}(\mathbf{X})$. Now, using this result, one can think of an algorithm to recover the maximal element \bar{O} . Indeed, if the value of the Pickands dependence function is known for every $\mathbf{t} \in \Delta_{d-1}$, one has to evaluate the metric given in (4) with \hat{C} and \hat{R} a given partition of $\llbracket d \rrbracket$. Iterate this process for every possible partition (\hat{C}, \hat{R}) of $\llbracket d \rrbracket$ and stop when it is equal to 0 to obtain an independent partition of $\llbracket d \rrbracket$. However, this process is computationally burdensome as the number of partition (\hat{C}, \hat{R}) to consider is given by the Stirling number of the second kind and will grows drastically as d increases. We thus need a metric that avoids to call for

value for each partition of $\{1, \dots, d\}$. Indeed, by Proposition 1, one has :

$$A(\mathbf{t}) \leq w^{(\llbracket d-1 \rrbracket)}(\mathbf{t}) A\left(\frac{t_1}{w^{(\llbracket d-1 \rrbracket)}(\mathbf{t})}, \frac{t_2}{w^{(\llbracket d-1 \rrbracket)}(\mathbf{t})}, \dots, \frac{t_{d-1}}{w^{(\llbracket d-1 \rrbracket)}(\mathbf{t})}, 0\right) + t_d, \quad \mathbf{t} \in \Delta_{d-1}, \quad (6)$$

where inequalities holds if and only if $X_d \perp\!\!\!\perp (X_1, \dots, X_{d-1})$. This statement holds whatever the chosen order but we choose the trivial one for notational convenience. This result gives us a metric to introduce an index in a candidate cluster. Indeed, take as candidate cluster $\hat{C} = \{1, \dots, d-1\}$ and we ask if the element $\{d\}$ belongs to \hat{C} . Thus, with knowledge of the above equation, $\{d\}$ do not belongs to \hat{C} if (6) holds as an equality for every $\mathbf{t} \in \Delta_{d-1}$.

Recall that \hat{C} and \hat{R} are an independent partition of $\llbracket d \rrbracket$ if and only if the criteria given in (4) is equal to 0. Another difficulty is this has to hold for every $\mathbf{t} \in \Delta_{d-1}$ which is not countable and implies computationally infeasibility for every considered dimension $d \in \mathbb{N}_*$. Indeed, one shall discretize the simplex of \mathbb{R}^d to have an approximation of (4). Take $\mathbf{t} \in \Delta_{d-1} \subset [0, 1]^d$, we want to have at least one observation at distance less than 1 say, from \mathbf{t} , then we must increase the number N of observations as d increases. We thus need at least

$$N \geq \frac{\Gamma(d/2 + 1)}{\pi^{d/2}} \stackrel{d \rightarrow \infty}{\sim} \left(\frac{d}{2\pi e}\right)^{d/2} \sqrt{d\pi}$$

points in order to fill the hypercube $[0, 1]^d$, thus Δ_{d-1} . This number of points grows exponentially fast with d . So, even with discretization, where we might obtain theoretical concentration bounds for the sup-norm over the simplex (using chaining methods), the criteria would not be computable in practice.

However, in the framework of extremes, independence between the components X_1, \dots, X_d of an extreme-value random vector $\mathbf{X} \in \mathbb{R}^d$ can be stated in a useful manner as stated in the Remark 3 with help of the exponent measure. One simple necessary and sufficient conditions is those stated in [49, 50]. It is shown that

$$G(\mathbf{x}) = G_1(x_1) \dots G_d(x_d), \quad \mathbf{x} \in \mathbb{R}^d, \quad (7)$$

that is \mathbf{X} is totally independent, if and only if there exists $\mathbf{p} = (p_1, \dots, p_d) \in \mathbb{R}^d$ such that (7) holds. Roughly speaking, in extreme-value theory if for at least one $\mathbf{p} \in \mathbb{R}^d$ Equation (7) holds, then it extends for every $\mathbf{x} \in \mathbb{R}^d$ using max-stability. An analogy of this statement in the non-extreme world could be similar to X_1 and X_2 are independent if

$$\text{cov}(f_1(X_1), f_2(X_2)) = 0, \quad (8)$$

holds for every continuous bounded functions f_1 and f_2 is equivalent to there exists f, g continuous and bounded functions such that (8) is equal to 0. This statement is known to be false in general but stand as true for the Gaussian case (take $f = g = id$).

The result of [50] is extended to our framework by $G(\mathbf{x}) = G^{(O_1)}(\mathbf{x}^{(O_1)})G^{(O_2)}(\mathbf{x}^{(O_2)})$ for any $\mathbf{x} = (\mathbf{x}^{(O_1)}, \mathbf{x}^{(O_2)}) \in \mathbb{R}^d$ if and only if the equation holds for at least one $\mathbf{p} = (\mathbf{p}^{(O_1)}, \mathbf{p}^{(O_2)}) \in \mathbb{R}^d$. The formal statement and its proof are given in Appendix B.3. Moreover, this result belongs to those who were stated for single components of max-stable random vector \mathbf{X} that can be extended for subvectors of max-stable random vector (see Lemma 6 of [39] or Exercise 5.5.1 of [42] for instance).

One direct application of this result is that $\mathbf{X}^{(O_1)}$ and $\mathbf{X}^{(O_2)}$ are independent if and only one has :

$$A\left(\frac{1}{d}, \dots, \frac{1}{d}\right) = \frac{d_1}{d} A^{(O_1)}\left(\frac{1}{d_1}, \dots, \frac{1}{d_1}\right) + \frac{d_2}{d} A^{(O_2)}\left(\frac{1}{d_2}, \dots, \frac{1}{d_2}\right).$$

By denoting $\theta = d A(d^{-1}, \dots, d^{-1})$ the so-called extremal coefficient, we restate the equation above as the SECO (Sum of Extremal COefficients) metric

$$\text{SECO}(O_1, O_2) = \theta^{(O_1)} + \theta^{(O_2)} - \theta. \quad (9)$$

Loosely speaking, if $\mathbf{X}^{(O_1)}$ and $\mathbf{X}^{(O_2)}$ are mutually independent, thus a characterizing property is that the extremal coefficient of \mathbf{X} is written as the sum of the extremal coefficients $\theta^{(O_1)}$ and $\theta^{(O_2)}$ of $\mathbf{X}^{(O_1)}$ and $\mathbf{X}^{(O_2)}$ respectively.

We now state an assumption in order to recover our hidden clusters and guarantee the positivity of the MSECO metric.

Assumption A. For every $k \in \{1, \dots, K\}$, $\mathbf{X}^{(\bar{O}_k)}$ exhibits asymptotic dependence between all components.

A sufficiency condition in order that Assumption A is satisfied is to suppose that each exponent measure of the extreme value random vectors $\mathbf{X}^{(\bar{O}_k)}$ has a nonnegative Lebesgue density on the non negative orthant $[0, \infty)^{d_k} \setminus \{\mathbf{0}^{(\bar{O}_k)}\}$ for every $k \in \{1, \dots, K\}$ (see [17] and Kirstin Strokorb's discussion contribution). Various classes of tractable extreme value distributions satisfy Assumption A. Popular models that are commonly used for statistical inference include the asymmetric logistic model ([51]), the asymmetric Dirichlet model ([9]), the pairwise Beta model ([10]) or the Hüsler Reiss model ([28]).

Remark 4. In its seminal work, [44] proposes a conditional independence test to decide whether an element j belongs to a cluster \hat{C} , say. Formally, one may ask if $\mathbf{X}^{(\hat{C})} \perp\!\!\!\perp X_j | \llbracket d \rrbracket \setminus (\hat{C} \cup \{j\})$. However, defining a suitable conditional independence for extremes is relatively new and mainly designed for tree inference (see e.g., [1, 17, 18, 46]). In those models, as in Gaussian graphical model (we refer to [31] for an overview), one supposes that the extreme-value random vector admits a Hüsler Reiss density which can be seen as a Gaussian extremal model with variogram. Interesting properties are also obtained such that conditional dependencies are encoded in the precision matrix of the Hüsler Reiss distribution.

Highlighted by Remark 4, if we suppose, as in the litterature of extremal graphs, that \mathbf{X} is absolutely continous with respect to the Lebesgue measure and admits an Hüsler Reiss density

in order to study conditional dependencies (defined in a peculiar way, see [17] for details), then \mathbf{X} shall exhibit extremal dependence between all its components. Thus, the resulting maximal element of the AI-block model is trivial and given by $\bar{O} = \{\{1\}, \dots, \{d\}\}$ or $\bar{O} = \{1, \dots, d\}$, namely each component is independent or none are. In order to drop Assumption A, further works are needed and left for future investigations.

Remark 5. In this remark we go outside the extreme framework, we suppose that \mathbf{X} has a Gaussian copula distribution with zero mean, and copula function with parameters $\mu = 0$ and Σ , a correlation matrix. Recall that this implies that

$$\mathbf{Y} := (Y_1, \dots, Y_d) := (h_1(X_1), \dots, h_d(X_d)) =: h(\mathbf{X}) \stackrel{d}{\sim} \mathcal{N}_d(0, \Sigma),$$

with $h_j = \Phi^{-1} \circ G_j$ for each $j \in \{1, \dots, d\}$, where Φ is the cumulative distribution function of a standard Gaussian random variable. One can show that

$$\mathbf{X}^{(O_1)} \perp\!\!\!\perp \mathbf{X}^{(O_2)} \iff |\Sigma^{(O_1)}| |\Sigma^{(O_2)}| = |\Sigma|,$$

where $\Sigma^{(O_k)}$ is a sub-matrix of Σ where we only kept the i th rows and columns with $i \in O_k$, $k \in \{1, 2\}$. Also, conditional independencies are encoded in the correlation matrix such as $\mathbf{X}^{(O_1)} \perp\!\!\!\perp X_j | (S \setminus R)$, where $R = S \setminus (O_1 \cup \{j\})$ is equivalent to

$$|\Sigma^{(R)}| |\Sigma| = |\Sigma^{(R \setminus \{j\})}| |\Sigma^{(S \setminus \{j\})}|.$$

One can hope that these properties might find an equivalent statement in the case of Husler-Reiss for extreme frameworks. However, such a distributional assumption will lead to a trivial model as stated right after Remark 4.

4 Consistent estimation of minimally separated clusters via SECO

In this section, we propose an estimation approach that utilizes nonparametric estimation of the Pickands dependence function and we use it to recover clusters. Let $\mathbf{X}_1, \dots, \mathbf{X}_n$ be i.i.d. copies of \mathbf{X} . The estimator that we present is based on the madogram concept ([11, 34]), a notion borrowed from geostatistics in order to capture the spatial dependence structure. Our estimator is defined as

$$\hat{A}_n(\mathbf{t}) = \frac{\hat{\nu}_n(\mathbf{t}) + c(\mathbf{t})}{1 - \hat{\nu}_n(\mathbf{t}) - c(\mathbf{t})}, \quad (10)$$

where

$$\begin{aligned}\hat{\nu}_n(\mathbf{t}) &= \frac{1}{n} \sum_{j=1}^n \left[\bigvee_{j=1}^d \{G_{n,j}(X_{i,j})\}^{1/t_j} - \frac{1}{d} \sum_{j=1}^d \{G_{n,j}(X_{i,j})\}^{1/t_j} \right], \\ c(\mathbf{t}) &= \frac{1}{d} \sum_{j=1}^d \frac{t_j}{1+t_j}.\end{aligned}$$

by convention, here $u^{1/0} = 0$ for $u \in (0, 1)$. One can use our results presented to design a new test of stochastic vectorial independence. Let $k \in \{1, \dots, K\}$ and let $\hat{A}_n^{(O_k)}(\mathbf{t}^{(O_k)}) = \hat{A}_n(\mathbf{0}, \mathbf{t}^{(O_k)}, \mathbf{0})$ denotes the empirical Pickands dependence function associated to the k -th subvector of X . Then consider the empirical process

$$\mathcal{E}_{nK}(\mathbf{t}) = \sqrt{n} \left(\hat{A}_n(\mathbf{t}) - \hat{A}_{n\Sigma}(\mathbf{t}) \right), \quad (11)$$

where $\hat{A}_{n\Sigma} = \sum_{k=1}^K w^{(O_k)}(\mathbf{t}) \hat{A}_n^{(O_k)}(\mathbf{t}^{(O_k)})$. We define the set of hypotheses :

$$\mathcal{H}_0 : O_1, \dots, O_K \text{ are mutually independent}, \quad \mathcal{H}_1 : \exists j \neq k \text{ } O_j, O_k \text{ are mutually dependent}.$$

Theorem below states the asymptotic behaviour of our statistic test \mathcal{E}_{nK} .

Theorem 2. *Under \mathcal{H}_0 , the empirical process \mathcal{E}_{nK} converges weakly in $\ell^\infty(\Delta_{d-1})$ to a tight Gaussian process having representation*

$$\begin{aligned}\mathcal{E}_K(\mathbf{t}) &= -(1 + A(\mathbf{t}))^2 \int_{[0,1]} N_C(u^{t_1}, \dots, u^{t_d}) du \\ &\quad + \sum_{k=1}^K w^{(O_k)}(\mathbf{t}) \left(1 + A^{(O_k)}(\mathbf{t}^{(O_k)}) \right)^2 \int_{[0,1]} N_C(\mathbf{1}, u^{t_{i,1}}, \dots, u^{t_{i,d_k}}, \mathbf{1}) du,\end{aligned}$$

where N_C is a continuous tight Gaussian process with representation

$$N_C(u_1, \dots, u_d) = B_C(u_1, \dots, u_d) - \sum_{j=1}^d \dot{C}_j(u_1, \dots, u_d) B_C(\mathbf{1}, u_i, \mathbf{1}),$$

and B_C is a continuous tight Gaussian process with covariance function

$$\text{cov}(B_C(\mathbf{u}), B_C(\mathbf{v})) = C(\mathbf{u} \wedge \mathbf{v}) - C(\mathbf{u})C(\mathbf{v}) \stackrel{X \sim O}{=} C_\Pi(\mathbf{u} \wedge \mathbf{v}) - C_\Pi(\mathbf{u})C_\Pi(\mathbf{v})$$

Using previous results stated in the present document, the test is equivalent to test the hypothesis under $\mathbf{t} = (d^{-1}, \dots, d^{-1})$, that is, whether the extremal coefficient of \mathbf{X} is equal to the sum of the other extremal coefficients of $\mathbf{X}^{(O_1)}, \dots, \mathbf{X}^{(O_K)}$. Furthermore, the asymptotic variance is easily (but technical) computable under this specific point.

We can estimate \bar{O} by applying a multiple test procedure by testing if

$$\mathbf{X}^{(\hat{C})} \perp\!\!\!\perp X_j, \quad j \in \llbracket d \rrbracket \setminus \hat{C}, \quad (12)$$

where $\hat{C} \subset \bar{O}_k$ is a candidate cluster for $k \in \{1, \dots, K\}$. This procedure makes perfect sense when n is large, much larger than d . But when d is large, the procedure will leads to poor results and induce multiple testing procedure.

In multivariate extreme-value theory, several summary statistics have been developped to measure the strength of dependence between the extreme of different variables. The most popular one is the extremal correlation, which for $a, b \in \llbracket d \rrbracket$ is defined as

$$\chi(a, b) = \lim_{q \rightarrow 0} \mathbb{P} \{F_a(X_a) > 1 - q | F_b(X_b) > 1 - q\}, \quad (13)$$

whenever the limit exists. It ranges between 0 and 1 when the boundary cases are asymptotic independence and complete extremal dependence, respectively. In particular, if \mathbf{X} is a multivariate extreme-value distribution, then the extremal correlation always exists and $\chi(a, b) = 2 - \theta(a, b)$ where $\theta(a, b)$ is the bivariate extremal dependence coefficient between X_a and X_b . In an AI-block model, the statement

$$\mathbf{X}^{(O_k)} \perp\!\!\!\perp \mathbf{X}^{(O_j)}, \quad k \neq j$$

is equivalent to

$$\chi(a, b) = \chi(b, a) = 0, \quad \forall a \in O_k, \forall b \in O_j, \quad k \neq j \quad (14)$$

(see Appendix for a proof of this statement). Let us denote by $\mathbb{X} = (\chi(a, b))_{a, b \in \llbracket d \rrbracket}$, i.e. the matrix composed of all extremal correlation. We will present an algorithm that recovers clusters in AI-block model using a dissimilarity based on the extremal correlation. Indeed, by Equation (14), two variables X_a and X_b belong to the same cluster of an AI-block model if and only if

$$\chi(a, b) > 0.$$

Then the difficulty of clustering in AI-block model can be assessed via the size of the Minimal Extremal CORrelation (MECO) separation between two variables in concomittant groups :

$$\text{MECO}(\mathcal{X}) = \min_{a \overset{\bar{O}}{\sim} b} \chi(a, b). \quad (15)$$

In AI-block models, with Assumption A, we always have $\text{MECO}(\mathcal{X}) > \eta$ for every η with $\eta = 0$. However, a large value of η will be needed for retrieving consistently the partition \bar{O} from identical and independent observation. Let $\eta > 0$ and define

$$\mathbb{X}(\eta) = \{\mathcal{X}, \text{MECO}(\mathcal{X}) > \eta\}.$$

We further recall that, under the extreme value model on \mathbf{X} ,

$$\chi(a, b) = 2 - \theta(a, b). \quad (16)$$

Therefore, using block maxima approach, we can estimate $\theta(a, b)$ by

$$\hat{\theta}(a, b) = \frac{\hat{\nu}(a, b) + 1/2}{1/2 - \hat{\nu}(a, b)},$$

where $\hat{\nu}$ is the so-called madogram ([11]) defined as

$$\hat{\nu}(a, b) = \frac{1}{n} \sum_{i=1}^n \left| \frac{\hat{R}_a}{n} - \frac{\hat{R}_b}{n} \right|,$$

and \hat{R}_a (resp. \hat{R}_b) corresponds to the empirical ranks of the sample $(X_{i,a})_{i=1}^n$ (resp. $(X_{i,b})_{i=1}^n$). For a given estimator $\hat{\mathbf{R}} = (\hat{\mathbf{R}}_1, \dots, \hat{\mathbf{R}}_d)$, where $\hat{\mathbf{R}}_j$ is n -dimensional vector composed out of the ranks of the corresponding sample from X_j , we associate the estimator

$$\hat{\chi}(a, b) = 2 - \hat{\theta}(a, b), \quad a, b \in \{1, \dots, d\},$$

of the extremal correlation. Using these, we define $\hat{\mathcal{X}} = (\hat{\chi}(a, b))_{a,b \in \{1, \dots, d\}}$. We then estimate the partition \hat{O} according to the extremal correlation as explained in Algorithm 1. We emphasize that this algorithm does not require as input the specification of the number of K groups which is automatically estimated by our procedure.

The algorithm complexity for computing $\hat{\mathbf{R}}$ is of order $O(dn \ln(n))$ ([13], Section 2). Given the empirical ranks, computing $\hat{\mathcal{X}}$ and performing the algorithm require $O(d^3 n \ln(n))$ and $O(d^3)$ computations, respectively. So the overall complexity of the estimation procedure is $O(d^3 n \ln(n))$.

In the following, we provide conditions ensuring that our algorithm is consistent, that is $\hat{O} = \bar{O}$. For notational convenience, we will write

$$|\hat{\theta} - \theta|_\infty = \sup_{a,b \in \llbracket d \rrbracket} |\hat{\theta}(a, b) - \theta(a, b)|.$$

Proposition 3. *Consider the AI-block model with $d \geq 2$ and $K \leq d$. Let us define $\tau = |\hat{\theta} - \theta|_\infty$ and consider parameters (α, η) such that*

$$\alpha \geq \tau, \quad \eta \geq \tau + \alpha. \quad (17)$$

Then if $\mathcal{X} \in \mathbb{X}(\eta)$ and under Assumption A, our algorithm yields $\hat{O} = \bar{O}$.

We below state a concentration inequality for $|\hat{\theta}(a, b) - \theta(a, b)|$ which is of a prime interest to study the consistency of the algorithm with high probability.

Algorithm 1 Split procedure with A unknown

```

1: procedure CLUST( $S, \alpha, \hat{\mathcal{X}}$ )
2:   Initialize :  $S = \{1, \dots, d\}$ ,  $\hat{\chi}(a, b) = 2 - \hat{\theta}(a, b)$  for  $a, b \in \{1, \dots, d\}$  and  $l = 0$ 
3:   while  $S \neq \emptyset$  do
4:      $l = l + 1$ 
5:     if  $|S| = 1$  then
6:        $\hat{O}_l = S$ 
7:     if  $|S| > 1$  then
8:        $(a_l, b_l) = \arg \max_{a, b \in S} \hat{\chi}(a, b)$ 
9:       if  $\hat{\chi}(a_l, b_l) < \alpha$  then
10:         $\hat{O}_l = \{a_l\}$ 
11:       if  $\hat{\chi}(a_l, b_l) \geq \alpha$  then
12:         $\hat{O}_l = \{s \in S : \hat{\chi}(a_l, s) \wedge \hat{\chi}(b_l, s) \geq \alpha\}$ 
13:        $S = S \setminus \hat{O}_l$ 
14:   return  $\hat{O} = (\hat{O}_l)_{l=1, \dots, K}$ 

```

Proposition 4. For $t > 0$, we have

$$\mathbb{P} \left\{ |\hat{\theta}(a, b) - \theta(a, b)| \geq t \right\} \leq 8 \exp \left\{ -\frac{nt^2}{8 \cdot 9^2} \right\}.$$

Corollary 1. Let us consider parameters fulfilling

$$\alpha \geq 18 \sqrt{\frac{4(1+\gamma)}{n} \ln(d)}, \quad \eta \geq 18 \sqrt{\frac{4(1+\gamma)}{n} \ln(d)} + \alpha,$$

for some $\gamma > 0$. If $\mathbf{X} \sim O$ and for $k \in \{1, \dots, K\}$, $\mathcal{X} \in \mathbb{X}(\eta)$, and under Assumption A, then the output of our algorithm applied to the estimator $\hat{\mathbf{R}}_n$, is consistent : $\hat{O} = \bar{O}$, with probability higher than $1 - 8d^{-2\gamma}$.

Remark 6. Therefore, Corollary 1 say that τ is of order $\sqrt{\ln(d)n^{-1}}$, thresholding $\hat{\chi}(a, b)$ at level τ guarantees exact recovery, whenever the separation MECO is at least 2τ . Such results are similar to the hard thresholding estimator in Gaussian sequence model, see e.g. Section 4.1 of [52] where similar bounds are found.

Note that above, independently of any distributional assumption, the χ metric and its related estimator $\hat{\chi}$ is sufficient to recover the original clusters. This results is quite interesting, since it is stronger than what is known in the classical, non-extremal independence clustering [44]. Indeed, in classical theory, one ask the existence of density in order to compute the mutual information between two vector $\mathbf{X}^{(\hat{C})}$ and $\mathbf{X}^{(\hat{R})}$. Here, we just use a simple statistic, which can be estimated non parametrically and without the existence of densities to recover our clusters. A similar statement could be obtained for Gaussian models where the correlation coefficient between random variables gives knowledge about their dependences as the extremal coefficient gives for extreme value distributions.

Such stronger results obtained in extreme value theory but not in the classical theory are also valid for trees, we refer to [18] for further details.

Some comments on the implications of the above result are in order. On a high level, larger dimension d leads to an higher bound of the threshold α . This is also implicated for the bound of the MECO metric. Thus, whereas the dimension d increases, the dependence between each component should be stronger in order to distinguish between noise and a signal. In other words, for alternatives that are sufficiently separated from the asymptotic independence case at the $\sqrt{\ln(d)n^{-1}}$ scale, the algorithm will be able to distinguish between the case of asymptotic independence or asymptotic dependence.

5 Numerical results

5.1 Exact recovery in AI-block models

In this Section, we verify numerically our theoretical findings. We consider a number of type of AI-block models of increasing complexity. We introduce a extreme-value distribution which will be used for our simulations. The symmetric logistic, or Gumbel model [25] is defined by the following Pickands" dependence function

$$A(w_1, \dots, w_d) = \left(\sum_{j=1}^d w_j^{1/\theta} \right)^\theta,$$

with $\theta \in [1, \infty)$. Algorithms developped by [48] is of prime interest to sample from the Logistic distribution as we can increase the dimension of the considered random vector \mathbf{X} without being limited by the computation time of generating the sample, which is out of interests.

E1 \mathbf{X} is made of two blocks O_1 and O_2 , of equal lengths where $\mathbf{X}^{(O_1)}$ and $\mathbf{X}^{(O_2)}$ are extreme-valued random vectors with a Logistic distribution and $\theta = 0.7$.

E2 \mathbf{X} is composed out of 5 blocks of random sizes d_1, \dots, d_5 sample form multinomial distribution with $\mathbf{p} = (p_k)_k$ where $p_k = 0.5^k$ for $k \in \{1, \dots, K-1\}$ and $p_K = \sum_{k=1}^{K-1} p_k$. Each random vectors are distributed according to a Logistic distribution where parameters θ_k are sampled uniformly under the segment $[0.65, 0.75]$ for $k \in \{1, \dots, K\}$.

E3 We consider the same model as E2 where we add 5 singletons, resulting in $K+5$ clusters.

The goal of our algorithm is to create sub-groups from a d random vector where n copies are observed, using the extremal correlation χ . This task can be viewed as that of clustering d objects in \mathbb{R}^n . In our simulations, the algorithm use estimation of the extremal correlation under the extreme-value setting using an estimator of the madogram. In the following simulation study, we

use the fixed threshold $\alpha = 2\sqrt{\ln(d)/n}$ since our theoretical results suggest the usage of a threshold proportional $\sqrt{\ln(d)/n}$.

We study the performance of our algorithm in terms exact true group recovery. In Figure 1 is shown the performance of Algorithm 1 where $d \in \{200, 1600\}$ and for $n \in \{100, 200, \dots, 1000\}$. The case $d = 1600$ is made to assert the efficiency of the algorithm in the high dimensional setting, i.e. where $d > n$. We saw that the algorithm for a suitable sample length n recover all the clusters with a exact recovery rate of 1. Even in Experiment E3 where, traditionnally, classical clustering methods struggle to recover singletons, our algorithm is still performant.

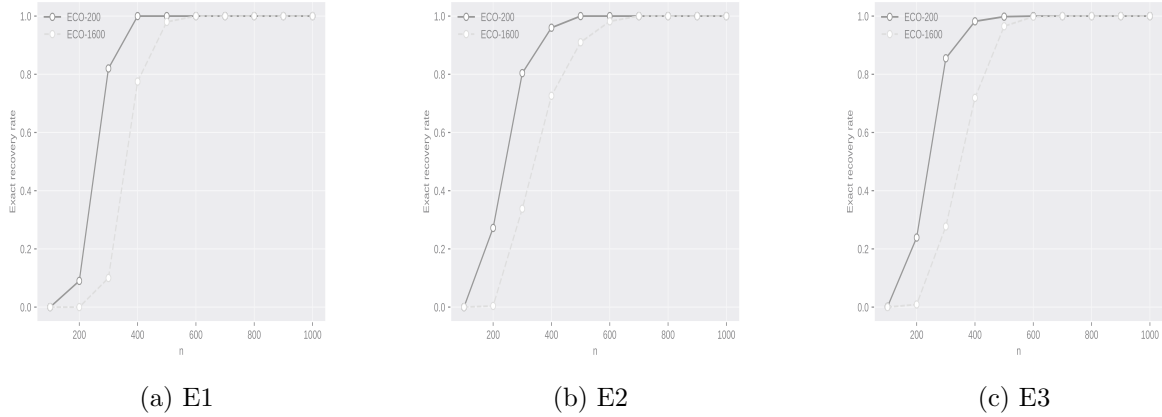


Figure 1: Exact recovery rate for Experiments E1, E2 and E3.

5.2 Data driven choice of the threshold

The threshold α in Corollary 1 does not involve any unknown quantity and it can be used directly. In practice, however, it is advisable to use a data-driven choice of the threshold for Algorithm 1. We propose to use the following type of cross-validation for this purpose. The idea is to use the SECO criteria presented in Section 3. Let $\mathbf{X} \sim O$, given a partition $\hat{O} = \{\hat{O}_k\}_k$, we know from Proposition 2 that the SECO similarity given by

$$\text{SECO}(\hat{O}) = \sum_k \theta^{(\hat{O}_k)} - \theta$$

is equal to 0 if and only if $\hat{O} \leq \bar{O}$. We thus construct a loss function given by the SECO where we evaluate its value over a grid of the α values. The value of α which the SECO similarity has minimum values is also the value of α for which we have consistent recovery of our communities. Our procedure ask to split the data in three independent sub-sets : on one subset, we construct a set of candidate partitions. The other two sets are used to estimate the extremal dependence coefficient for the whole distribution \mathbf{X} and the sub-vector of the candidate partition $\mathbf{X}^{(\hat{O}_k)}$ respectively. Let $\hat{\theta}_{(i)}$, $i \in \{1, 2, 3\}$ be the madogram-based estimator of θ based on three independent samples, each

of size n . The cross-validation based estimator of the SECO is thus defined as the following

$$\widehat{\text{SECO}}(\hat{O}) = \sum_k \hat{\theta}_{(1)}^{(\hat{O}_k)} - \hat{\theta}_{(2)}. \quad (18)$$

Let $\hat{\mathcal{O}}$ be a collection of partitions computed with Algorithm 1 from the third sample, by varying α around its theoretical optimal value (of order $\sqrt{\ln(d)n^{-1}}$, on a fine grid. For any $\hat{O} \in \hat{\mathcal{O}}$, we evaluate our cross-validation SECO in (18). Proposition 5 offers theoretical support for this procedure, for large n . It shows that the minimum of the proposed criterion, over a grid that is fine enough to include the target \bar{O} , is asymptotically attained at \bar{O} , in expectation.

Proposition 5. *Let $\mathbf{X} \sim O$, then we have*

$$\lim_{n \rightarrow \infty} \mathbb{E} [\widehat{\text{SECO}}(\bar{O})] < \lim_{n \rightarrow \infty} \mathbb{E} [\widehat{\text{SECO}}(\hat{O})], \quad \hat{O} \not\leq \bar{O}.$$

To illustrate how the cross-validation SECO enables selection of α numerically, we use the case $n = 800$ and $d = 1600$ to study the high-dimensional performance of our cross validation approach.

We show the relationship between the average SECO and exact recovery percentages in Figure 2 by our algorithm. This figure shows that the optimal ranges of α values for high exact recovery percentages are also associated with low average SECO values. Also, we can notice that for $n = 800, d = 1600$, the constant 2 belongs to the range where the exact recovery rate is equal to 1, consistent to what has been shown previously in Section 5.1.

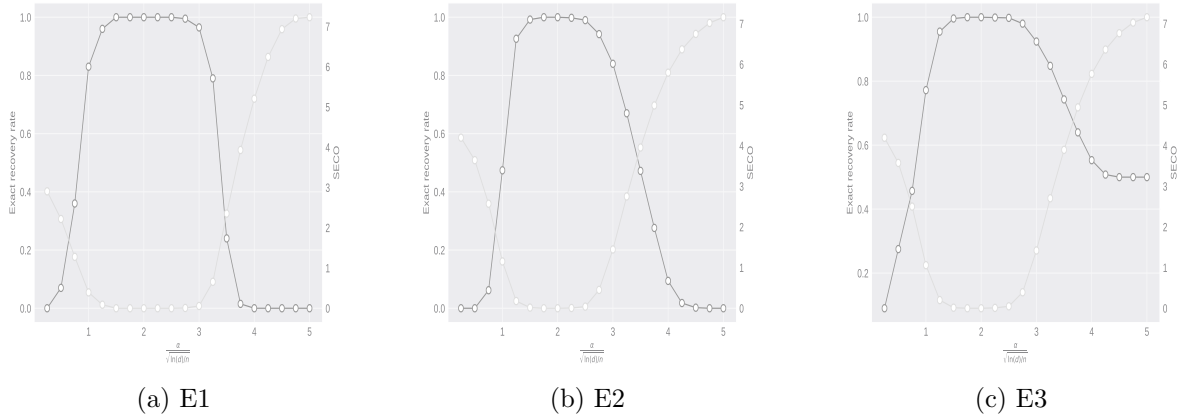


Figure 2: Average SECO (lightgrey) and exact recovery percentages (grey) across 100 runs. For better illustration, the SECO are standardized first by subtracting the minimal SECO in each figure, and the standardized SECO added by 1 are then plotted on the logarithmic scale.

5.3 Weak asymptotic dependence block model

Recall that the framework of an AI-block model is required for our theoretical guarantees. Here, we consider a scenario where the blocks have common elements which exhibits a lower dependence among extremes. To do so, we consider the same framework in Section 6.3. of [20]. Half of the observations come from the (randomly sampled) Hüsler-Reiss model with two groups $\{1, \dots, 40\}, \{41, \dots, 100\}$ and half form an analogous model with two groups $\{21, \dots, 60\}, \{1, \dots, 20, 61, \dots, 100\}$. The four block to be identified is thus $\{1, \dots, 20\}, \{21, \dots, 40\}, \{41, \dots, 60\}$ and $\{61, \dots, 100\}$. The estimated matrix of tail dependence $\chi(a, b)$ where $a, b \in \llbracket d \rrbracket$ is presented in Figure 3a. We indeedly observe 4 blocks which exhibits a strong asymptotic dependence where some blocks have common elements with a weaker dependence. This image may give a false feeling that identification of groups is easy, and so we reorder the indices according where the results might be found in Figure 3b. Clearly, no clear pattern emerges from it.

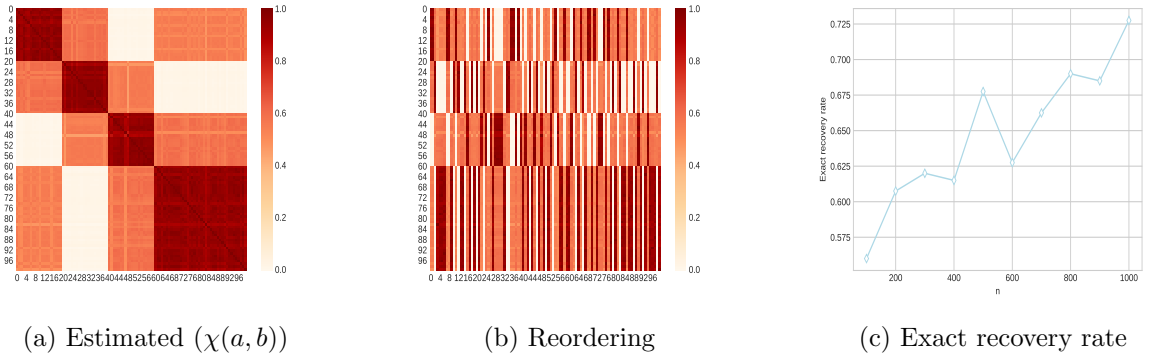


Figure 3: Block models with weak asymptotic dependence.

References

- [1] S. Asenova, G. Mazo, and J. Segers. Inference on extremal dependence in the domain of attraction of a structured hüsler-reiss distribution motivated by a markov tree with latent variables. *Extremes*, 24(3):461–500, 2021.
- [2] M. Bador, P. Naveau, E. Gilleland, M. Castellà, and T. Arivelo. Spatial clustering of summer temperature maxima from the cnrm-cm5 climate model ensembles & e-obs over europe. *Weather and climate extremes*, 9:17–24, 2015.
- [3] A. A. Balkema and S. I. Resnick. Max-infinite divisibility. *Journal of Applied Probability*, 14(2):309–319, 1977.
- [4] J. Beirlant, Y. Goegebeur, J. Segers, and J. L. Teugels. *Statistics of extremes: theory and applications*, volume 558. John Wiley & Sons, 2004.

- [5] E. Bernard, P. Naveau, M. Vrac, and O. Mestre. Clustering of maxima: Spatial dependencies among heavy rainfall in france. *Journal of climate*, 26(20):7929–7937, 2013.
- [6] F. Bunea, C. Giraud, X. Luo, M. Royer, and N. Verzelen. Model assisted variable clustering: minimax-optimal recovery and algorithms. *The Annals of Statistics*, 48(1):111–137, 2020.
- [7] F. Bunea, C. Giraud, X. Luo, M. Royer, and N. Verzelen. Model assisted variable clustering: Minimax-optimal recovery and algorithms. *The Annals of Statistics*, 48(1):111 – 137, 2020.
- [8] E. Chautru. Dimension reduction in multivariate extreme value analysis. *Electronic journal of statistics*, 9(1):383–418, 2015.
- [9] S. G. Coles and J. A. Tawn. Modelling extreme multivariate events. *Journal of the Royal Statistical Society. Series B (Methodological)*, 53(2):377–392, 1991.
- [10] D. Cooley, R. A. Davis, and P. Naveau. The pairwise beta distribution: A flexible parametric multivariate model for extremes. *Journal of Multivariate Analysis*, 101(9):2103–2117, 2010.
- [11] D. Cooley, P. Naveau, and P. Poncet. Variograms for spatial max-stable random fields. In *Dependence in probability and statistics*, pages 373–390. Springer, 2006.
- [12] D. Cooley and E. Thibaud. Decompositions of dependence for high-dimensional extremes. *Biometrika*, 106(3):587–604, 2019.
- [13] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to algorithms*. MIT press, 2022.
- [14] L. De Haan, A. Ferreira, and A. Ferreira. *Extreme value theory: an introduction*, volume 21. Springer, 2006.
- [15] H. Drees and A. Sabourin. Principal component analysis for multivariate extremes. *Electronic Journal of Statistics*, 15(1):908–943, 2021.
- [16] C. Eisenach, F. Bunea, Y. Ning, and C. Dinicu. High-dimensional inference for cluster-based graphical models. *Journal of machine learning research*, 21(53), 2020.
- [17] S. Engelke and A. S. Hitz. Graphical models for extremes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 82(4):871–932, 2020.
- [18] S. Engelke and S. Volgushev. Structure learning for extremal tree models. *arXiv preprint arXiv:2012.06179*, 2020.
- [19] M. Falk, J. Hüsler, and R. Reiss. *Laws of Small Numbers: Extremes and Rare Events*. Springer Basel, 2010.
- [20] V. Fomichov and J. Ivanovs. Spherical clustering in detection of groups of concomitant extremes. *Biometrika*, 2022.

- [21] J. Galambos. The asymptotic theory of extreme order statistics. Technical report, 1978.
- [22] N. Gissibl and C. Klüppelberg. Max-linear models on directed acyclic graphs. *Bernoulli*, 24(4A):2693–2720, 2018.
- [23] N. Goix, A. Sabourin, S. Clémen, et al. Learning the dependence structure of rare events: a non-asymptotic study. In *Conference on Learning Theory*, pages 843–860. PMLR, 2015.
- [24] G. Gudendorf and J. Segers. Extreme-value copulas. In P. Jaworski, F. Durante, W. K. Härdle, and T. Rychlik, editors, *Copula Theory and Its Applications*, pages 127–145, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.
- [25] E. Gumbel. Distributions de valeurs extrêmes en plusieurs dimensions. *Publications de l’institut de Statistique de l’Université de Paris*, 9:171–173, 1960.
- [26] M. Hofert, R. Huser, and A. Prasad. Hierarchical archimax copulas. *Journal of Multivariate Analysis*, 167:195–211, 2018.
- [27] X. Huang. *Statistics of bivariate extreme values*. Thesis Publishers Amsterdam, 1992.
- [28] J. Hüsler and R.-D. Reiss. Maxima of normal random vectors: Between independence and complete dependence. *Statistics & Probability Letters*, 7(4):283–286, 1989.
- [29] A. Janßen and P. Wan. k -means clustering of extremes. *Electronic Journal of Statistics*, 14(1):1211–1233, 2020.
- [30] M. Lalancette, S. Engelke, and S. Volgushev. Rank-based estimation under asymptotic dependence and independence, with applications to spatial extremes. *The Annals of Statistics*, 49(5):2552–2576, 2021.
- [31] S. L. Lauritzen. *Graphical models*, volume 17. Clarendon Press, 1996.
- [32] P. D. Le, A. C. Davison, S. Engelke, M. Leonard, and S. Westra. Dependence properties of spatial rainfall extremes and areal reduction factors. *Journal of Hydrology*, 565:711–719, 2018.
- [33] F. Lindskog. *Multivariate extremes and regular variation for stochastic processes*. PhD thesis, ETH Zurich, 2004.
- [34] G. Marcon, S. Padoan, P. Naveau, P. Muliere, and J. Segers. Multivariate nonparametric estimation of the pickands dependence function using bernstein polynomials. *Journal of Statistical Planning and Inference*, 183:1–17, 2017.
- [35] A. W. Marshall and I. Olkin. Domains of Attraction of Multivariate Extreme Value Distributions. *The Annals of Probability*, 11(1):168 – 177, 1983.
- [36] A. W. Marshall and I. Olkin. Domains of Attraction of Multivariate Extreme Value Distributions. *The Annals of Probability*, 11(1):168 – 177, 1983.

- [37] N. Meyer and O. Wintenberger. Sparse regular variation. *Advances in Applied Probability*, 53(4):1115–1148, 2021.
- [38] P. Naveau, A. Guillou, D. Cooley, and J. Diebolt. Modelling pairwise dependence of maxima in space. *Biometrika*, 96(1):1–17, 2009.
- [39] I. Papastathopoulos and K. Strokorb. Conditional independence among max-stable laws. *Statistics & Probability Letters*, 108:9–15, 2016.
- [40] J. Pickands. Multivariate extreme value distribution. *Proceedings 43th, Session of International Statistical Institution, 1981*, 1981.
- [41] D. Pollard. Strong consistency of k-means clustering. *The Annals of Statistics*, pages 135–140, 1981.
- [42] S. Resnick. *Extreme Values, Regular Variation, and Point Processes*. Applied probability. Springer, 2008.
- [43] S. I. Resnick. *Heavy-tail phenomena: probabilistic and statistical modeling*. Springer Science & Business Media, 2007.
- [44] D. Ryabko. Independence clustering (without a matrix). In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [45] K. Saunders, A. Stephenson, and D. Karoly. A regionalisation approach for rainfall based on extremal dependence. *Extremes*, 24(2):215–240, 2021.
- [46] J. Segers. One-versus multi-component regular variation and extremes of markov trees. *Advances in Applied Probability*, 52(3):855–878, 2020.
- [47] E. S. Simpson, J. L. Wadsworth, and J. A. Tawn. Determining the dependence structure of multivariate extremes. *Biometrika*, 107(3):513–532, 2020.
- [48] A. Stephenson. Simulating multivariate extreme value distributions of logistic type. *Extremes*, 6(1):49–59, 2003.
- [49] R. Takahashi. Some properties of multivariate extreme value distributions and multivariate tail equivalence. *Annals of the Institute of Statistical Mathematics*, 39:637–647, 1987.
- [50] R. Takahashi. Asymptotic independence and perfect dependence of vector components of multivariate extreme statistics. *Statistics & Probability Letters*, 19(1):19–26, 1994.
- [51] J. A. Tawn. Modelling multivariate extreme value distributions. *Biometrika*, 77(2):245–253, 06 1990.

- [52] A. B. Tsybakov. Aggregation and high-dimensional statistics (preliminary notes of saint-flour lectures, july 8-20, 2013). 2014.
- [53] A. van der Vaart, A. van der Vaart, A. van der Vaart, and J. Wellner. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer Series in Statistics. Springer, 1996.

A Proofs

A.1 Proofs of main results

We show here that there exists a unique maximal element \bar{O} in AI-block models. To do that, we show several items concerning the introduced the partition partial order and the partition induced by the equivalence relation $\overset{O \cap S}{\sim}$. Using these properties, we construct an explicit unique maximal element of $\mathcal{L}(\mathbf{X})$.

Proof of Theorem 1 For the first item, if $\mathbf{X} \sim O$, then the partition O induces mutually independent random vectors. As S is a sub-partition of O , it also generates a partition where vectors are mutually independent.

Now, take $k \in \{1, \dots, K\}$ and $j, k \in (O \cap S)_k$, we thus have in particular $j \overset{O}{\sim} k$, thus there exists $k' \in \{1, \dots, K'\}$ such that $a, b \in O'_k$. Thus $(O \cap S)_k \subseteq O_{k'}$. Hence the second statement.

The third result comes down from 1. and 2. We conclude the proof of this proposition by proving the fourth claim. The set $\mathcal{L}(\mathbf{X})$ is non-empty since the trivial partition $O = \{1, \dots, d\}$ belongs to $\mathcal{L}(\mathbf{X})$. It is also a finite set, and we can enumerate it $\mathcal{L}(\mathbf{X}) = \{O_1, \dots, O_m\}$. Define the sequence O'_1, \dots, O'_M recursively according to

- $O'_1 = O_1$,
- $O'_k = O_k \cap O'_{k-1}$ for $k = 2, \dots, M$.

According to Theorem 1, we have that by induction $O'_1, \dots, O'_M \in \mathcal{L}(\mathbf{X})$. In addition, we have both $O'_{k-1} \leq O'_k$ and $O_k \leq O'_k$, so by induction $O_1, \dots, O_k \leq O'_k$. Hence the partition $O^*(\mathbf{X}) := O'_M = O_1 \cap \dots \cap O_{M-1}$ is the maximum of $\mathcal{L}(\mathbf{X})$. \square

We extend the well-known inequality $A \leq 1$ where inequality holds if and only if \mathbf{X} is totally independent to the case of mutual independence, that is $A \leq A_O$ using notations of Section 2. Two ways are given to obtain this statement. The first use the convexity and homogeneity of order one of the stable tail dependence function. The associativity of extreme-value random vectors are of a prime interest to obtain the statement in a second sketch.

Proof of Proposition 1 As L is an homogeneous and convex function under a cone, it is also subadditive, *i.e.*

$$L(\mathbf{x} + \mathbf{y}) \leq L(\mathbf{x}) + L(\mathbf{y}),$$

for every $\mathbf{x}, \mathbf{y} \in [0, \infty)^d$. In particular, we obtain

$$L\left(\sum_{k=1}^K \mathbf{z}_k\right) \leq \sum_{k=1}^K L(\mathbf{z}_k),$$

where $\mathbf{x}_k \in [0, \infty)^d$ with $k \in \{1, \dots, K\}$. Consider now $\mathbf{x}_k = (\mathbf{0}, x_{i_{k,1}}, \dots, x_{i_{k,d_k}}, \mathbf{0})$, we directly obtain

$$L(\mathbf{z}) = L\left(\sum_{k=1}^K \mathbf{z}_k\right) \leq \sum_{k=1}^K L(\mathbf{z}_k) = \sum_{k=1}^K L^{(k)}(z_{i_{k,1}}, \dots, z_{i_{k,d_k}})$$

Expression this equation in terms of Pickands gives

$$A(\mathbf{t}) \leq \sum_{k=1}^K \frac{1}{z_1 + \dots + z_d} L^{(O_k)}(z_{i_{k,1}}, \dots, z_{i_{k,d_k}}) = \sum_{k=1}^K \frac{z_{i_{k,1}} + \dots + z_{i_{k,d_k}}}{z_1 + \dots + z_d} A^{(O_k)}(t_{i_{k,1}}, \dots, t_{i_{k,d_k}}),$$

where $t_i = z_i / (z_1 + \dots + z_d)$. Hence the result.

Another proof is obtained using that extreme-value distributions are associated (see Proposition 5.1 of [35] or Section 5.4.1 of [42]), *i.e.*

$$\mathbb{E}[f(\mathbf{X})g(\mathbf{X})] \geq \mathbb{E}[f(\mathbf{X})]\mathbb{E}[g(\mathbf{X})],$$

for every increasing (or decreasing) functions f, g . By induction, we can prove that, $K \in \mathbb{N}_*$,

$$\mathbb{E}[\Pi_{k=1}^K f_k(\mathbf{X})] \geq \Pi_{k=1}^K \mathbb{E}[f_k(\mathbf{X})] \quad (19)$$

Take $f_k(\mathbf{x}) = \mathbb{1}_{\{-\infty, \mathbf{x}_k\}}$ for each $k \in \{1, \dots, K\}$, thus Equation (19) gives

$$C(G_1(x_1), \dots, G_d(x_d)) \geq \Pi_{k=1}^K C^{(O_k)}\left(G^{(O_k)}\left(\mathbf{x}^{(O_k)}\right)\right),$$

which can be restated in terms of stable tail dependence function

$$L(\mathbf{z}) \leq \sum_{k=1}^K L^{(O_k)}(\mathbf{z}^{(O_k)}).$$

We obtain the statement expressing this inequality with Pickands dependence function.

For the rest of the proof, notice that (19) with $f_k(\mathbf{x}) = \mathbb{1}_{\{-\infty, \mathbf{x}^{(O_k)}\}}$ for each $k \in \{1, \dots, K\}$ holds as an inequality if and only if $\mathbf{X}^{(O_1)} \perp\!\!\!\perp \dots \perp\!\!\!\perp \mathbf{X}^{(O_K)}$. \square

We thus state the theoretical value taken by $\mathbf{X} \sim O$ for a given arbitrary partition called $\bigsqcup_{l=1}^L S_l$.

We show that this value is strictly greater than 0 if and only if the considered partition does not induce mutually independent random vectors. Again, the proof makes use of the associativity of extreme random vectors. To obtain the theoretical value, we decompose each element of the partition S_l , $l \in \{1, \dots, L\}$ as a subset of \bar{O}_k . Using mutual independence, we thus obtain the closed form of the theoretical value through the Pickands dependence function.

Proof of Proposition 2 Note that $\mathbf{X}^{(S_l)} = (\mathbf{X}_{j_{l,1}}, \dots, \mathbf{X}_{j_{l,D_l}})$ and $\mathbf{X}^{(S_m)} = (\mathbf{X}_{j_{m,1}}, \dots, \mathbf{X}_{j_{m,D_m}})$ are not, in general independent where $l, m \in \{1, \dots, L\}$ because the partition $\bigsqcup_{l=1}^L S_l$ is arbitrarily given. Consider for some $l \in \{1, \dots, L\}$,

$$f_l = \mathbb{1}_{\{-\infty, \mathbf{x}^{(S_l)}\}}, \quad \mathbf{x}^{(S_l)} = (x_{j_{l,1}}, \dots, x_{j_{l,D_l}}).$$

We thus have in one hand

$$\mathbb{E}[\Pi_{l=1}^L f_l(\mathbf{X})] = \mathbb{P}\{X_1 \leq x_1, \dots, X_d \leq x_d\} = \Pi_{k=1}^K C^{(O_k)}\left(G^{(O_k)}(\mathbf{x}^{(O_k)})\right),$$

on the other hand

$$\Pi_{l=1}^L \mathbb{E}[f_l(\mathbf{X})] = \Pi_{l=1}^L \mathbb{P}\left\{X_{j_{l,1}} \leq x_{j_{l,1}}, \dots, X_{j_{l,D_l}} \leq x_{j_{l,D_l}}\right\}.$$

It is worth noticing that the above set can be rewritten as

$$\left\{X_{j_{l,1}} \leq x_{j_{l,1}}, \dots, X_{j_{l,D_l}} \leq x_{j_{l,D_l}}\right\} = \bigcap_{k=1}^K \bigcap_{i \in S_l \cap O_k} \{X_i \leq x_i\}.$$

Notice that $\bigcap_{i \in \emptyset} \{X_i \leq x_i\} = \Omega$. As a subvector of $\mathbf{X}^{(k)}$, $\mathbf{X}_{kl} = (X_i)_{i \in S_l \cap O_k}$ is independent of $\mathbf{X}^{(j)}$ for $j \neq k$. Using this block-independence, we thus obtain :

$$\mathbb{P}\left\{X_{j_{l,1}} \leq x_{j_{l,1}}, \dots, X_{j_{l,D_l}} \leq x_{j_{l,D_l}}\right\} = \Pi_{k=1}^K \mathbb{P}\left\{\bigcap_{i \in S_l \cap O_k} \{X_i \leq x_i\}\right\}.$$

Now using positive association, we obtain that

$$\Pi_{k=1}^K C^{(O_k)}\left(G^{(O_k)}(x^{(O_k)})\right) \geq \Pi_{l=1}^L \Pi_{k=1}^K C^{(O_k)}(\mathbf{x}_{S_l \cap O_k}),$$

where $\mathbf{x}_{kl} = (x_i)_{i \in S_l \cap O_k}$. Reexpressing the whole in terms of stable tail dependence function leads to :

$$\sum_{k=1}^K L^{(O_k)}(z_{i_{k,1}}, \dots, z_{k,d_k}) \leq \sum_{l=1}^L \sum_{k=1}^K L^{(O_k)}(\mathbf{0}, \mathbf{z}^{(S_l \cap O_k)}, \mathbf{0}),$$

with $z_i = -\ln(G_i(x_i))$ for every $i \in \{1, \dots, d\}$. Dividing by $z_1 + \dots + z_d$ gives

$$A_O(\mathbf{t}) = \sum_{k=1}^K w^{(O_k)}(\mathbf{t}) A^{(O_k)}(\mathbf{t}^{(O_k)}) \leq \sum_{l=1}^L \sum_{k=1}^K \frac{1}{z_1 + \dots + z_d} L^{(O_k)}(\mathbf{0}, \mathbf{z}^{(S_l \cap O_k)}, \mathbf{0}).$$

The right hand side of the equation can be written as $\forall k \in \{1, \dots, K\}, \forall l \in \{1, \dots, L\}$

$$\begin{aligned} \frac{1}{z_1 + \dots + z_d} L^{(O_k)}(\mathbf{0}, \mathbf{z}^{(S_l \cap O_k)}, \mathbf{0}) &= \frac{\sum_{j \in S_l \cap O_k} z_j}{z_1 + \dots + z_d} L^{(O_k)}\left(\mathbf{0}, \frac{\mathbf{z}^{(S_l \cap O_k)}}{\sum_{j \in S_l \cap O_k} z_j}, \mathbf{0}\right) \\ &\triangleq w^{(S_l \cap O_k)}(\mathbf{t}) A^{(O_k)}(\mathbf{0}, \mathbf{t}^{(S_l \cap O_k)}, \mathbf{0}). \end{aligned}$$

We thus obtain the statement.

Now, if $\forall k \in \{1, \dots, K\}, \exists l \in \{1, \dots, L\}$ such that $O_k \subseteq S_l$, then the random vectors $\mathbf{X}^{(S_l)} = (\mathbf{X}_{j_{l,1}}, \dots, \mathbf{X}_{j_{l,D_l}})$ and $\mathbf{X}^{(S_m)} = (\mathbf{X}_{j_{m,1}}, \dots, \mathbf{X}_{j_{m,D_m}})$ are now independent. If we suppose $L \geq K$, we thus have $S_1 = O_1, \dots, S_K = O_K$ and $S_l = \emptyset$ for $l > K$. The equality comes down from Lemma 1. Now, if $L < K$, we have with the same notations in the proof

$$\mathbb{E} [\Pi_{l=1}^L f_l(\mathbf{X})] = \Pi_{l=1}^L \mathbb{P} \left\{ X_{j_{l,1}} \leq x_{j_{l,1}}, \dots, X_{j_{l,D_l}} \leq x_{j_{l,D_l}} \right\} = \Pi_{k=1}^K \mathbb{P} \left\{ \mathbf{X}^{(k)} \leq \mathbf{x}_k \right\}.$$

Expressing this two terms as before, we obtain that $A = A_O$. For the converse, suppose that the right hand side of the inequality in (5) is equal to zero. Applying Lemma 1 gives that the arbitrary partition has the same value as A_O . That is saying that the random vectors inside S_1, \dots, S_L are mutually independent. If $L \geq K$, thus, apart for the K first clusters say for which $S_k = O_k$, the others are empty set. Now, if $L < K$, we group one or more O_k in a given cluster S_l without one overlaps to an other cluster S_j say (if it does, the value could not be equal as zero by Lemma 1). Hence the statement. \square

To prove the weak convergence of our process given in Theorem 2, we make of use of empirical processes as stated in [53].

Proof The proof is straightforward, notice that (see Figure 4)

$$\mathcal{E}_{nK} = \psi \circ \phi \left(\sqrt{n}(\hat{\theta}_n - \theta) \right),$$

where ϕ is detailed as

$$\begin{aligned} \phi : \ell^\infty(\Delta_{d-1}) &\rightarrow \ell^\infty(\Delta_{d-1}) \otimes (\ell^\infty(\Delta_{d-1}), \dots, \ell^\infty(\Delta_{d-1})) \\ x &\mapsto (x, \phi_1(x), \dots, \phi_K(x)), \end{aligned}$$

$$\begin{array}{ccc}
& \sqrt{n}(\hat{A}_n - A) & \rightarrow \mathcal{E}_{nK} \\
& \downarrow \phi & \nearrow \psi \\
& \left(\sqrt{n}(\hat{A}_n - A); w_1 \sqrt{n}(\hat{A}_{n1} - A_1), \dots, w_k(\mathbf{t}) \sqrt{n}(\hat{A}_{nk} - A^{(O_k)}) \right)
\end{array}$$

Figure 4: Diagram of composition of function.

with for every $k \in \{1, \dots, K\}$

$$\begin{array}{ccc}
\phi_k & : \ell^\infty(\Delta_{d-1}) & \rightarrow \ell^\infty(S_d) \\
x & \mapsto & x(\mathbf{0}, t_{i_{k,1}}, \dots, t_{i_{k,d_k}}, \mathbf{0}).
\end{array}$$

Thus ϕ_k is a linear and bounded function hence continuous, it follows that ϕ is continuous since each coordinate functions are continuous. Using that (see [34])

$$\sqrt{n}(\hat{A}_n(\mathbf{t}) - A(\mathbf{t})) \rightsquigarrow -(1 + A(\mathbf{t}))^2 \int_{[0,1]} N_C(u^{t_1}, \dots, u^{t_d}) du,$$

and applying the continuous mapping theorem for the weak convergence in $\ell^\infty(\Delta_{d-1})$ (Theorem 1.3.6 of [53]) leads the result. \square

We now present all tools used obtain the exact recovery of our algorithm, that is $\hat{O} = \bar{O}$. The result is obtained by induction on step l while we assume that the algorithm remains consistent at this step. We show that under conditions (17) and the cluster separation condition, that is $\Theta \in \mathcal{X}(\eta)$, we have that :

$$a \stackrel{\bar{O}}{\sim} b \iff \hat{\chi}(a, b) > \alpha.$$

And thus the algorithm, under this equivalence, the algorithm remains consistent at step l . The arguments of using induction reasoning to prove the consistency is the same as used in [7] while the equivalence stated afore is proven using tool of Theorem 7 (iii) of [52].

Proof of Proposition 3 If $a \stackrel{\bar{O}}{\not\sim} b$, then $\chi(a, b) = 0$ and

$$\hat{\chi}(a, b) = \hat{\chi}(a, b) - \chi(a, b) \leq \tau \leq \alpha.$$

Now, if $a \stackrel{\bar{O}}{\sim} b$, then we have, under $\Theta \in \mathcal{X}(\eta)$, so $\chi(a, b) > \tau + \alpha$ and

$$\tau + \alpha < \chi(a, b) - \hat{\chi}(a, b) + \hat{\chi}(a, b),$$

and thus

$$\hat{\chi}(a, b) > \alpha.$$

In particular, under (17) and the separation condition $\Theta \in \mathcal{X}(\eta)$, we have

$$a \stackrel{\bar{O}}{\sim} b \iff \hat{\chi}(a, b) > \alpha. \quad (20)$$

Let us prove the proposition by induction on l . We consider the algorithm at some step $l - 1$ and assume that the algorithm was consistent up to this step, *i.e.* $\hat{O}_j = \bar{O}_j$ for $j = 1, \dots, l - 1$.

If $\hat{\chi}(a_l, b_l) \leq \alpha$, then according to (20), no $b \in S$ is in the same group of a_l . Since the algorithm has been consistent up to this step l , it means that a_l is a singleton and $\hat{O}_l = \{a_l\}$.

If $\hat{\chi}(a_l, b_l) > \alpha$, then $a_l \stackrel{\bar{O}}{\sim} b$ according to (20). Furthermore, the equivalence implies that $\hat{O}_l = S \cap \bar{O}_l$. Since the algorithm has been consistent up to this step, we have $\hat{O}_l = \bar{O}_l$. To conclude, the algorithm remains consistent at the step l and the proposition follows by induction. \square

We will now prove that Proposition 3 holds with high probability for given parameters. We need to specify some threshold τ such that $|\hat{\theta} - \theta|_\infty \leq \tau$ with high probability. To do so, we make use of concentration inequality stated in Proposition 4 that gives concentration of the extremal coefficient using madogram estimator.

Proof of Corollary 1 We have that for $t > 0$:

$$\mathbb{P} \left\{ |\hat{\theta} - \theta|_\infty \geq t \right\} \leq d^2 \mathbb{P} \left\{ |\hat{\theta}(a, b) - \theta(a, b)| \geq t \right\}.$$

Using Proposition 4, one has

$$\mathbb{P} \left\{ |\hat{\theta}(a, b) - \theta(a, b)| \geq t \right\} \leq 8 \exp \left\{ -\frac{nt^2}{8 \cdot 9^2} \right\},$$

By considering $\delta \in]0, 1[$ and solve the following equation

$$\frac{\delta}{d^2} = 8 \exp \left\{ -\frac{nt^2}{8 \cdot 9^2} \right\},$$

with respect to t gives that :

$$\mathbb{P} \left\{ |\hat{\theta} - \theta|_\infty \geq 18 \sqrt{\frac{2}{n} \ln \left(\frac{8d^2}{\delta} \right)} \right\} \leq \delta.$$

Now, taking $\delta = 8d^{-2\gamma}$, with $\gamma > 0$, we have

$$|\hat{\theta} - \theta|_\infty \leq 18 \sqrt{\frac{4(1 + \gamma)}{n} \ln(d)},$$

with probability higher than $1 - 8d^{-2\gamma}$. The result then follows from Proposition 3, since $|\hat{\theta} - \theta|_\infty \leq$

$18\sqrt{\frac{4(1+\gamma)}{n}} \ln(d)$ with probability higher than $1 - 8d^{-2\gamma}$. □

Proof of Proposition 5

As $\hat{\theta}_{(i)}$ is strongly consistent for $i \in \{1, 2, 3\}$, we obtain that

$$\widehat{\text{SECO}}(\bar{O}) \xrightarrow[n \rightarrow \infty]{\text{a.s.}} \text{SECO}(\bar{O}) = 0$$

Using Lebesgue's theorem, we obtain

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\widehat{\text{SECO}}(\bar{O}) \right] = 0.$$

Furthermore, applying again Lebesgue theorem

$$0 < \mathbb{E} \left[\text{SECO}(\hat{O}) \right] = \lim_{n \rightarrow \infty} \mathbb{E} \left[\widehat{\text{SECO}}(\hat{O}) \right].$$

To conclude,

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\widehat{\text{SECO}}(\bar{O}) \right] < \lim_{n \rightarrow \infty} \mathbb{E} \left[\widehat{\text{SECO}}(\hat{O}) \right].$$

What had to be proven. □

B Proofs of auxiliary results

B.1 extreme-value copula

In this first lemma, we prove that the function introduced in Paragraph 2.2 is indeed an extreme-value copula. For the ease of reading, we recall here its definition

$$\begin{aligned} C_{\Pi} : [0, 1]^d &\longrightarrow [0, 1] \\ \mathbf{u} &\longmapsto \prod_{k=1}^K C^{(k)}(u_{i_{k,1}}, \dots, u_{i_{k,d_k}}). \end{aligned}$$

To prove this statement, we show that each margins is indeed distributed uniformly on the unit segment $[0, 1]$. Hence C is a copula function. In order to prove that C is an extreme-value copula, we show that C is max-stable as it is a characterizing property of extreme-value copula or, more generally, of extreme-value distribution.

Proof of Lemma 1 We first show that C is a copula function. It is clear that $C(\mathbf{u}) \in [0, 1]$ for every $\mathbf{u} \in [0, 1]^d$. We check that its univariate marginals are uniformly distributed on $[0, 1]$. Without loss of generality, take $u_{i_{1,1}} \in [0, 1]$ and let us compute

$$C(1, \dots, u_{i_{1,1}}, \dots, 1) = C^{(1)}(u_{i_{1,1}}, 1, \dots, 1) \prod_{k=1}^K C^{(k)}(1, \dots, 1) = C^{(1)}(u_{i_{1,1}}, 1, \dots, 1) = u_{i_{1,1}}.$$

So C is a copula function. We now have to prove that C is an extreme-value copula. We recall that C is an extreme-value copula if and only if C is max-stable, that is for every $m \geq 1$

$$C(u_1, \dots, u_d) = C(u_1^{1/m}, \dots, u_d^{1/m})^m.$$

By definition, we have

$$C(u_1^{1/m}, \dots, u_d^{1/m})^m = \left(\prod_{k=1}^K C^{(k)} \left(u_{i_k,1}^{1/m}, \dots, u_{i_k,d_k}^{1/m} \right) \right)^m = \prod_{k=1}^K \left\{ C^{(k)} \left(u_{i_k,1}^{1/m}, \dots, u_{i_k,d_k}^{1/m} \right) \right\}^m.$$

Using that $C^{(1)}, \dots, C^{(K)}$ are extreme-value copulae, thus max stable, we obtain

$$C(u_1^{1/m}, \dots, u_d^{1/m})^m = \prod_{k=1}^K C^{(k)} \left(u_{i_k,1}, \dots, u_{i_k,d_k} \right) = C(u_1, \dots, u_d).$$

Thus C is an extreme-value copula. We end the proof by proving that C is associated to the random vector $\mathbf{X} = (\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(K)})$, that is

$$\mathbb{P} \{ \mathbf{X} \leq \mathbf{x} \} = C(G_1(x_1), \dots, G_d(x_d)), \quad \mathbf{x} \in \mathbb{R}^d.$$

Using mutual independence between random vectors, we have

$$\begin{aligned} \mathbb{P} \{ \mathbf{X} \leq \mathbf{x} \} &= \prod_{k=1}^K \mathbb{P} \left\{ X_{i_k,1} \leq x_{i_k,1}, \dots, X_{i_k,d_k} \leq x_{i_k,d_k} \right\} \\ &= \prod_{k=1}^K C^{(k)} \left(G_{i_k,1}(x_{i_k,1}), \dots, G_{i_k,d_k}(x_{i_k,d_k}) \right) \\ &= C(G_1(x_1), \dots, G_d(x_d)). \end{aligned}$$

Hence the result. □

B.2 Proof of proposition 4

We will prove this result for a given dimension $d \geq 2$ and use the result in the core of the paper for $d = 2$. Using the same tools introduced in [34], a generalization of the bivariate madogram is given by

$$\nu = \mathbb{E} \left[\bigvee_{j=1}^d G_j(X_j) - \frac{1}{d} \sum_{j=1}^d G_j(X_j) \right].$$

This quantity can be expressed in terms of the extremal dependence coefficient θ with the following relationship

$$\theta = \frac{1/2 + \nu}{1/2 - \nu}.$$

Thus, a plug-in estimation process for θ is directly obtained using

$$\hat{\theta}_n = \frac{1/2 + \hat{\nu}_n}{1/2 - \hat{\nu}_n},$$

where

$$\hat{\nu}_n = \frac{1}{n} \sum_{i=1}^n \left[\bigvee_{j=1}^d \hat{G}_{n,j}(X_{i,j}) - \frac{1}{d} \sum_{j=1}^d \hat{G}_{n,j}(X_{i,j}) \right].$$

Technical details for proving the concentration bound of $|\hat{\theta}_n - \theta|$ will be subdivided in some lemmas of which the combined use will give the statement of Proposition 4. The first lemma gives an upper bound of $|\hat{\theta}_n - \theta|$ with respect to $|\hat{\nu}_n - \nu|$. This follows from the link the extremal coefficient and the madogram.

Lemma 2. *We have,*

$$|\hat{\theta}_n - \theta| \leq (d+1)^2 |\hat{\nu}_n - \nu|.$$

Proof We have respectively $\theta = f(\nu)$ and $\hat{\theta}_n = f(\hat{\nu}_n)$ where

$$\begin{aligned} f &: \mathbb{R}_+ \rightarrow \mathbb{R}_+ \\ x &\mapsto \frac{1/2+\nu}{1/2-\nu}. \end{aligned}$$

Now, using that $\theta \leq d$, we have that

$$\nu + \frac{1}{2} \leq \frac{d}{2} - d\nu,$$

which implies that

$$\nu \leq \frac{d-1}{2(d+1)} < \frac{1}{2}.$$

In particular $2^{-1} - \nu > (d+1)^{-1} > 0$. Now taking derivation, we directly have

$$|f'(x)| = \frac{1}{(1/2 - \nu)^2} \leq (d+1)^2, \quad x \in \left[0, \frac{d-1}{2(d+1)}\right].$$

Thus f is $(d+1)^2$ -Lipschitz continuous and we have

$$|\hat{\theta}_n - \theta| \leq (d+1)^2 |\hat{\nu}_n - \nu|.$$

Hence the statement. □

Now, we state a concentration inequality for the madogram estimator. This inequality is obtained through two main arguments, that are Hoeffding's inequality and the DKW inequality bound.

Lemma 3. *For $t > 0$ one has*

$$\mathbb{P} \{ |\nu_n - \nu| > t \} \leq 4d \exp \left\{ -\frac{nt^2}{8} \right\}.$$

Proof observe that

$$|\hat{\nu}_n - \nu| \leq |\hat{\nu}_n - \nu_n| + |\nu_n - \nu|,$$

where

$$\nu_n := \sum_{i=1}^n Y_i = \sum_{i=1}^n \frac{1}{n} \left[\bigvee_{j=1}^d F_j(X_{i,j}) - \frac{1}{d} \sum_{j=1}^d F_j(X_{i,j}) \right].$$

As the following inequalities holds for every $i \in \{1, \dots, n\}$

$$Y_i \leq \frac{(d-1)}{dn}.$$

Hoeffding's inequality applies and we obtain that

$$\mathbb{P} \{ |\nu_n - \nu| > t/2 \} \leq 2 \exp \left(-\frac{nd^2t^2}{4(d-1)^2} \right).$$

Furthermore, we have

$$|\hat{\nu}_n - \nu_n| \leq 2 \sup_{j \in \{1, \dots, d\}} \sup_{i \in \{1, \dots, n\}} \left| \hat{F}_{n,j}(X_{i,j}) - F_j(X_{i,j}) \right|$$

Applying DKW inequality, we obtain

$$\mathbb{P} \left\{ \sup_{j \in \{1, \dots, d\}} \sup_{i \in \{1, \dots, n\}} \left| \hat{F}_{n,j}(X_{i,j}) - F_j(X_{i,j}) \right| > \frac{t}{4} \right\} \leq 2d \exp \left(-\frac{nt^2}{8} \right).$$

We thus have for $d \geq 2$

$$\mathbb{P} \{ |\hat{\nu}_n - \nu| > t \} \leq 4d \exp \left(-\frac{nt^2}{8} \right).$$

Hence the statement. □

Combine Lemma 2 and Lemma 3 gives Proposition 4.

B.3 Asymptotic independence between multivariate extreme distribution

In this subsection, we extend the result given in Theorem 2.1 of [50] for asymptotic independence between extreme random vector. The used arguments are similar of those used in the proof in [50]. We make extensive use of the following result (see, for example [36] and the proof of Theorem 5.3.1 of [21]) *i.e.* $F \in D(G)$ is equivalent to

$$\lim_{n \rightarrow \infty} n \{1 - F(\mathbf{a}_n \mathbf{x} + \mathbf{b}_n)\} = -\ln G(\mathbf{x}) \quad (21)$$

for all \mathbf{x} such that $0 < G(\mathbf{x}) < 1$. In this section, we denote by \bar{F} the survival function of F .

Theorem 3. *Let F be a d -distribution function and let $G^{(O_i)}$ be a d_i -extreme-value distribution for*

$i = 1, 2$. Then for $\mathbf{a}_n > 0$ and $\mathbf{b}_n \in \mathbb{R}^d$

$$\{F(\mathbf{a}_n \mathbf{x} + \mathbf{b}_n)\}^n \xrightarrow{n \rightarrow \infty} G^{(O_1)}(\mathbf{x}^{(O_1)})G^{(O_2)}(\mathbf{x}^{(O_2)}) \quad (22)$$

if and only if

$$\left\{F^{(O_i)}(\mathbf{a}_n^{(O_i)} \mathbf{x}^{(O_i)} + b_n^{(O_i)})\right\}^n \xrightarrow{n \rightarrow \infty} G^{(O_i)}(\mathbf{x}^{(O_i)}), \quad (23)$$

and there exists a $\mathbf{p} = (\mathbf{p}^{(O_1)}, \mathbf{p}^{(O_2)}) \in \mathbb{R}^d$ such that $0 < H^{(O_1)}(\mathbf{x}^{(O_1)}), H^{(O_2)}(\mathbf{x}^{(O_2)}) < 1$ and

$$\{F(\mathbf{a}_n \mathbf{x} + \mathbf{b}_n)\}^n \xrightarrow{n \rightarrow \infty} G^{(O_1)}(\mathbf{p}^{(O_1)})G^{(O_2)}(\mathbf{p}^{(O_2)}) \quad (24)$$

Proof The proof follows exactly the same lines as in Theorem 2.1 of [50]. One substantial difference is emphasizes in Remark. For any $\mathbf{x} \in \mathbb{R}^d$, $0 < H^{(O_1)}(\mathbf{x}^{(O_1)}), H^{(O_2)}(\mathbf{x}^{(O_2)}) < 1$, there exists $s > 0$ such that $\{H^{(O_1)}(\mathbf{x}^{(O_1)})\}^{1/s} > H^{(O_1)}(\mathbf{p}^{(O_1)})$, $\{H^{(O_2)}(\mathbf{x}^{(O_2)})\}^{1/s} > H^{(O_2)}(\mathbf{p}^{(O_2)})$. By Equation (23)

$$\left\{F^{(O_i)}\left(\mathbf{a}_{[sn]}^{(O_i)} \mathbf{x}^{(O_i)} + b_{[sn]}^{(O_i)}\right)\right\}^{sn} \xrightarrow{n \rightarrow \infty} H^{(O_i)}(\mathbf{x}^{(O_i)})$$

thus

$$\left\{F^{(O_i)}\left(\mathbf{a}_{[sn]}^{(O_i)} \mathbf{x}^{(O_i)} + b_{[sn]}^{(O_i)}\right)\right\}^n \xrightarrow{n \rightarrow \infty} \left\{H^{(O_i)}(\mathbf{x}^{(O_i)})\right\}^{1/s}.$$

Notice that

$$\begin{aligned} \mathbb{P}\left\{\left[\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}, \mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right]^c\right\} &= \mathbb{P}\left\{\left[\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}\right]^c \cup \left[\mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right]^c\right\} \\ &= \mathbb{P}\left\{\left[\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}\right]^c\right\} + \mathbb{P}\left\{\left[\mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right]^c\right\} \\ &\quad - \mathbb{P}\left\{\left[\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}\right]^c \cap \left[\mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right]^c\right\}. \end{aligned}$$

We thus obtain

$$\begin{aligned} \mathbb{P}\left\{\left[\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}\right]^c \cap \left[\mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right]^c\right\} &= \left[1 - \mathbb{P}\left\{\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}\right\}\right] + \left[1 - \mathbb{P}\left\{\mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right\}\right] \\ &\quad - \left[1 - \mathbb{P}\left\{\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}, \mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}\right\}\right]. \end{aligned}$$

Using Equations (23) and (24) we have, in joint hands, that

$$\begin{aligned} n \left[1 - \mathbb{P}\left\{\mathbf{X}^{(O_i)} \leq \mathbf{a}_n^{(O_i)} \mathbf{p}^{(O_i)} + \mathbf{b}_n^{(O_i)}\right\}\right] &\xrightarrow{n \rightarrow \infty} -\ln G^{(O_i)}(\mathbf{p}^{(O_i)}), \quad i = 1, 2 \\ n \left[1 - \mathbb{P}\left\{\mathbf{X} \leq \mathbf{a}_n \mathbf{p} + \mathbf{b}_n\right\}\right] &\xrightarrow{n \rightarrow \infty} -\ln G^{(O_1)}(\mathbf{p}^{(O_1)})G^{(O_2)}(\mathbf{p}^{(O_2)}). \end{aligned}$$

Thus,

$$\mathbb{P}\left\{\left[\mathbf{X}^{(O_1)} \leq \mathbf{a}^{(O_1)} \mathbf{p}^{(O_1)} + \mathbf{b}_n^{(O_1)}\right]^c \cap \left[\mathbf{X}^{(O_2)} \leq \mathbf{a}^{(O_2)} \mathbf{p}^{(O_2)} + \mathbf{b}_n^{(O_2)}\right]^c\right\} \xrightarrow{n \rightarrow \infty} 0.$$

Using now that $\left\{\mathbf{X}^{(O_1)} > \mathbf{x}^{(O_1)}, \mathbf{X}^{(O_2)} > \mathbf{x}^{(O_2)}\right\} \subset \left\{[\mathbf{X}^{(O_1)} \leq \mathbf{x}^{(O_1)}]^c \cap [\mathbf{X}^{(O_2)} \leq \mathbf{x}^{(O_2)}]^c\right\}$, we obtain

$$n\bar{F}(\mathbf{a}_n\mathbf{p} + \mathbf{b}_n) \xrightarrow{n \rightarrow \infty} 0.$$

By $\mathbf{q} \geq \mathbf{p}$, we now have

$$0 \leq n\bar{F}(\mathbf{a}_n\mathbf{q} + \mathbf{b}_n) \leq n\bar{F}(\mathbf{a}_n\mathbf{p} + \mathbf{b}_n) \xrightarrow{n \rightarrow \infty} 0.$$

The rest of the proof is similar to [50]. □

Remark 7. When $d_1 = d_2 = 1$, we immediately have that :

$$\mathbb{P}\{X_1 > x_1, X_2 > x_2\} = [1 - \mathbb{P}\{X_1 \leq x_1\}] + [1 - \mathbb{P}\{X_2 \leq x_2\}] - [1 - \mathbb{P}\{X_1 \leq x_1, X_2 \leq x_2\}].$$

We immediately obtain that, under the same hypotheses of Theorem 3 that

$$n\bar{F}(\mathbf{a}_n\mathbf{p} + \mathbf{b}_n) \xrightarrow{n \rightarrow \infty} 0. \tag{25}$$

The arguments exposed in this remark are those used in the proof of Theorem 2.1 [50]. In our work framework, we do not directly obtain (25) but we can upper bound this quantity with respect to an other which indeedly converges to 0 as $n \rightarrow \infty$ in the framework of Theorem 3.