

Introduction

0.1 Context

Management of environmental resources often requires the analysis of multivariate extreme values. In the classical theory, one is often interested in the behavior of the mean or average of a random variable X defined on a probability space $(\Omega, \mathcal{A}, \mathbb{P})$. This average will then be described through the expected value $\mathbb{E}[X]$ of the distribution. The central limit theorem yields, under mild conditions, the asymptotic behavior of the sample mean \bar{X} . This result can be used to provide a confidence interval for $\mathbb{E}[X]$ for a level $\alpha \in [0, 1]$. But in case of extreme events, it can be just as important to estimate tails probabilities. Furthermore, what if the second moment $\mathbb{E}[X^2]$ or even the mean is not finite? Then the central limit theorem does not apply and the classical theory, carried by the normal distribution, is no longer relevant [Beirlant et al., 2004].

Some extreme events, such as heavy precipitation or wind speed has spatial characteristics and geostatisticians are striving to better understand the physical processes in hand. In geostatistics, we often consider $(\Omega, \mathcal{F}, \mathbb{P})$ a probability space, S a set of locations and (E, \mathcal{E}) a measurable state space. We define on this probability space a stochastic process $X = \{X(s), s \in S\}$ with values on (E, \mathcal{E}) . It is classical to define the following second-order statistic as the variogram (see [Gaetan and Guyon, 2008] Chapter 1.3 for definition and basic properties) :

$$2\gamma(h) = \mathbb{E}[|X(s+h) - X(s)|^2],$$

where $\{X(s), s \in S\}$ represents a spatial and stationary process with a well-defined covariance function. The function $\gamma(\cdot)$ is called the semi-variogram of X . With respect to extremes, this definition is not well adapted because a second order statistic is difficult to interpret inside the framework of extreme value theory or may not even be defined. To ensure that we always work with finite moments quantities, the following type of first-order variogram is introduced by [Cooley et al., 2006]

$$\nu(h) = \frac{1}{2} \mathbb{E}[|F(X(s+h)) - F(X(s))|],$$

where $F(u) = \mathbb{P}(X(s) \leq u)$ is named as the FMadogram. His link to the pairwise extremal dependence function (Section 4.3 of [Coles et al., 1999]) or the Pickands dependence function ([Pickands, 1981]) make him an interesting quantity to capture the dependence between the extremas of stochastic processes or random variables. Indeed, this quantity may be seen as a

dissimilarity measure among bivariate maxima to be used for clustering time-series as shown by [Bernard et al., 2013] or [Bador et al., 2015].

The main drawback of this quantity is that she only focus on the value of the diagonal section of the pairwise extremal dependence function. In the bivariate case, the FMadogram characterize solely the extremal dependence coefficient for random variables X and Y (see Section 8.2.7 of [Beirlant et al., 2004]). To overpass this drawback, [Naveau et al., 2009] introduce the λ -FMadogram defined as,

$$\nu(h, \lambda) = \frac{1}{2} \mathbb{E} [|F^\lambda(X(s+h)) - F^{1-\lambda}(X(s))|],$$

for every $\lambda \in [0, 1]$.

This quantity characterize the pairwise extremal dependence function outside the diagonal section but also the whole Pickands dependence function ([Marcon et al., 2017]) and contribute to the vast litterature of the estimation of the Pickands dependence function for bivariate extreme value copulas (see for example [Pickands, 1981], [Deheuvels, 1991], [Hall and Tajvidi, 2000] or [Capéraà et al., 1997]). Statisticians may estimate this quantity but the classical results may applied only if the data in hands are clean as possible. This induce that the process of data collection has not been corrupted such as the data table is complete and that the implicit law of the observations is still the same.

Nevertheless, as the volume of data expands, the problem of missing or contaminated data has been increasingly present in many fields of statistical applications. It frequently happens that all of the individuals of a sample of statistical data from a multivariate population are not observed. If a sample be represented in matrix form by allowing the rows to represent the individuals and the columns the variables, then the matrix of the type of sample with which we are concerned is incomplete in that some elements are not present. In dealing with fragmentary samples, it is important to have at hand techniques which will enable the statistician to extract as much information as possible from the data. A useful reference for general parametric statistical inferences with missing data was provided by [Little R.J.A., 1987].

Considering a random sample of incomplete data,

$$(X_t, Y_t, \delta_t), \quad t \in \{1, \dots, T\}, \tag{1}$$

where all the X_t 's are observed and $\delta_t = 0$ if Y_t is missing, otherwise $\delta_t = 1$. The simple missing data pattern describe by (1) is basically created by the double sampling or two phase sampling (see chapter 12 of [Cochran, 2007]). Samples like (1) may arise in survival analysis : The study of the duration time preceding an event of interest is considered with series of random censors, which might prevent the capture of the whole survival time. This is known as the censoring mechanism and it arises from restrictions depending from the nature of the study. Typically, they may occur in medicine, with studies of the survival times before the recovery / decease from a specific disease. Another important example is often realized in comparing treament effects of

two educational programs. Individuals with lower scores on a preliminary test are more likely to receive the experimental treatment (*i.e.*, a composatory study program), whereas those with higher preliminary scores are more inclined to take the standard control. This phenomenon is well-known as the selection problem and we refer to Chapter 2 of [Angrist and Pischke, 2008] for more details. Beside of missing observations, the process of data might be disturb in a way that innerly deteriorate the quality of some data and one may ask that the estimation process should be robust.

The topic of Robustness in estimation has known an important research activity developed in the 60's and 70's resulting in a large number of publications. For a summary, the interested reader is referred to [Huber, 2011]. Robustness can be seen as an estimation procedure in which both stochastic and approximation error are low (see Section 1.1 from [Baraud et al., 2016]). In other words, an estimator is said to be robust if our model provides a reasonable approximation of the true one and derive an estimator which remains close to the true distribution. In this report, we mean by *robust* as *robust against outlier*, *e.g* the ϵ -contamination model (see [Huber, 1964]), or *robust again heavy-tailed data* where only low-order moments are assumed to be finite for the data distribution. There is no simple relation between the two definitions and the first framework of robustness that we have depicted. We want to propose a robust estimator of the Madogram. In our perspective, we only know [Escobar-Bach et al., 2018] that include the contamination framework in their estimation of the Pickands dependence function in the extreme value theory. To achieve our goal, we leverage the idea of Median-Of-Means (MoN). Intuitively, we replace the linear operator of expectation with the median of averages taken over non-overlapping blocks of the data, in order to get a robust estimate thanks to the median step (see [Lerasle et al., 2019] for a similar idea applied to Kernel).

0.2 Definitions and Notation

Let $(\Omega, \mathcal{A}, \mathbb{P})$ be a probability space and (X, Y) be a bivariate random vector with values in $(\mathbb{R}^2, \mathcal{B}(\mathbb{R}^2))$. This random vector has a joint distribution function H and marginal distribution function F and G . A function $C : [0, 1]^2 \rightarrow [0, 1]$ is called a bivariate copula if it is the restriction to $[0, 1]^2$ of a bivariate distribution function whose marginals are given by the uniform distribution on the interval $[0, 1]$. Since the work of [Sklar, 1959], it is well known that every distribution function H can be decomposed as $H(x, y) = C(F(x), G(y))$, for all $(x, y) \in \mathbb{R}^2$. This function C characterizes the dependence between X and Y and is called an extreme value copula if and if it admits a representation of the form [Gudendorf and Segers, 2009]

$$C(u, v) = (uv)^{A(\log(v)/\log(uv))}, \quad (2)$$

for all $u, v \in [0, 1]$ and where $A(\cdot)$ is the Pickands dependence function, *i.e.*, $A : [0, 1] \rightarrow [1/2, 1]$ is convex and satisfies $t \vee (1 - t) \leq A(t) \leq 1$, $\forall t \in [0, 1]$. The upper and lower bound of A has special meanings, the upper bound $A(t) = 1$ corresponds to independence, whereas the lower bound $A(t) = t \vee 1 - t$ corresponds to the perfect dependence (comonotonicity). Notice that,

on sections, the extreme value copula is of the form

$$C(u^t, u^{1-t}) = u^{A(t)}. \quad (3)$$

Let $(X_t, Y_t)_{t=1, \dots, T}$ be an *i.i.d.* sample of a bivariate random vector whose underlying copula is denoted by C and whose margins by F, G . For $x, y \in \mathbb{R}$, let $x \wedge y = \min(x, y)$ and $x \vee y = \max(x, y)$. Let $(b_{t,j})_{t \geq 1, j \in \{1,2\}}$ and $(a_{t,j})_{t \geq 1, j \in \{1,2\}}$ be respectively a sequence of numbers and a sequence of positive numbers. We say that the sequence $(a_{t,1}^{-1}(\bigvee_{t=1}^T X_t - b_{t,1}), a_{t,2}^{-1}(\bigvee_{t=1}^T Y_t - b_{t,2}))$ belongs to the domain of attraction of H , if for all real values x, y (at which the limit is continuous)

$$\mathbb{P} \left(\frac{\bigvee_{t=1}^T X_t - b_{t,1}}{a_{t,1}} \leq x, \frac{\bigvee_{t=1}^T Y_t - b_{t,2}}{a_{t,2}} \leq y \right) \xrightarrow{T \rightarrow \infty} H(x, y).$$

If this relationship hold, H is said to be a multivariate extreme value distribution. We will call by FMadogram the following quantity

$$\nu = \frac{1}{2} \mathbb{E} [|F(X) - G(Y)|], \quad (4)$$

and the λ -FMadogram by the expression

$$\nu(\lambda) = \frac{1}{2} \mathbb{E} [|F^\lambda(X) - G^{1-\lambda}(Y)|]. \quad (5)$$

A classical estimator of the λ -FMadogram when the margins F, G are unknown is

$$\hat{\nu}(\lambda) = \frac{1}{2T} \sum_{t=1}^T |\hat{F}_T^\lambda(X_t) - \hat{G}_T^{1-\lambda}(Y_t)| \quad (6)$$

with \hat{F}_T (resp. \hat{G}_T) the empirical cumulative distribution function of X (resp. Y). We suppose that we observe sequentially a quadruple defined by

$$(I_t X_t, J_t Y_t, I_t, J_t), \quad t \in \{1, \dots, T\}, \quad (7)$$

where $I_t = 0$ (resp. $J_t = 0$) if X_t (resp. Y_t) is missing, otherwise $I_t = 1$ (resp. $J_t = 1$), *i.e.* at each $t \in \{1, \dots, T\}$, one of both entries may be missing. The probability of observing a realisation partially or completely is denoted by $p_X = \mathbb{P}(I_t = 1) > 0$, $p_Y = \mathbb{P}(J_t = 1) > 0$ and $p_{XY} = \mathbb{P}(I_t = 1, J_t = 1) > 0$. Let us now define the empirical cumulative distribution of X (resp. Y and (X, Y)) in case of missing data,

$$\hat{F}_T(u) = \frac{\sum_{t=1}^T 1_{\{X_t \leq u\}} I_t}{\sum_{t=1}^T I_t}, \quad \hat{G}_T(v) = \frac{\sum_{t=1}^T 1_{\{Y_t \leq v\}} J_t}{\sum_{t=1}^T J_t}, \quad \hat{H}_T(u, v) = \frac{\sum_{t=1}^T 1_{\{X_t \leq u, Y_t \leq v\}} I_t J_t}{\sum_{t=1}^T I_t J_t}. \quad (8)$$

Here, we weight the estimator by the number of observed data which is a natural estimator (if divided by T) of the probabilities of missing. We have all tools in hand to define the *hybrid*

copula estimator introduced by [Segers, 2014],

$$\hat{C}_T^{\mathcal{H}}(u, v) = \hat{H}_T(\hat{F}_T(u), \hat{G}_T(v)). \quad (9)$$

Given a rate $r_T > 0$ and $r_T \rightarrow \infty$ as $T \rightarrow \infty$, the normalized estimation error of the hybrid copula estimator is :

$$\mathbb{C}_T^{\mathcal{H}}(u, v) = r_T \left(\hat{C}_T^{\mathcal{H}}(u, v) - C(u, v) \right). \quad (10)$$

In order to propose a robust estimator we will assume that the sample is partitioned into K disjoint subsets B_1, \dots, B_K of cardinalities $n_j := \text{card}(B_j)$ respectively, where the partitioning scheme is independent of the data. Let f be a measurable function from $(\mathbb{R}^2, \mathcal{B}(\mathbb{R}^2))$ to $(\mathbb{R}^2, \mathcal{B}(\mathbb{R}^2))$, we define the following estimator of $\mathbb{E}[f(X, Y)]$ by

$$\bar{\mathbb{P}}_{n_j} f = \frac{1}{n_j} \sum_{j \in B_j} f(X_j, Y_j).$$

We define the MoN estimator of f as solutions of the optimization problem

$$\hat{f}_{MoN} = \underset{z \in \mathbb{R}}{\operatorname{argmin}} \sum_{j=1}^K |\bar{\mathbb{P}}_{n_j} f - z|, \quad (11)$$

which, if we note $\text{med}(\cdot)$ the usual univariate median

$$\hat{f}_{MoN} = \text{med}(\bar{\mathbb{P}}_{n_1} f, \dots, \bar{\mathbb{P}}_{n_K} f), \quad (12)$$

is a solution of Equation (11).

We will write the generalized inverse function of F (respectively G) as $F^{\leftarrow}(u) = \inf\{v \in \mathbb{R} | F(v) \geq u\}$ (respectively $G^{\leftarrow}(u) = \inf\{v \in \mathbb{R} | G(v) \geq u\}$) where $0 < u, v < 1$. Given $\mathcal{X} \subset \mathbb{R}^2$, let $l^\infty(\mathcal{X})$ denote the spaces of bounded real-valued function on \mathcal{X} . For $f : \mathcal{X} \rightarrow \mathbb{R}$, let $\|f\|_\infty = \sup_{x \in \mathcal{X}} |f(x)|$. Here, we use the abbreviation $Qf = \int f dQ$ for a given measurable function f and signed measure Q . The arrows $\xrightarrow{a.s.}$, \xrightarrow{d} and \rightsquigarrow denote almost sure convergence, convergence in distribution of random vectors and weak convergence in the sense of J.Hoffman-Jørgensen (see Part 1 in the monograph by [van der Vaart and Wellner, 1996]).

This work is organized as follows : In Chapter 1 we state some results on the weak convergence of the estimator of the λ -FMadogram with missing data. To propose a robust estimator of the λ -FMadogram, we leverage the idea of Median-Of-Means (MoN) and state a concentration inequality that this estimator does verify. We also propose a closed formula for the asymptotic variance of the λ -FMadogram for a fixed $\lambda \in [0, 1]$.

Chapter 2 will present our results in a finite-sample framework. The asymptotic variance of the normalized estimation error of several models would be drawn with their empirical counterpart obtained through simulation. We also propose a reproduction of the experiment of the λ -FMadogram with a Smith's process as found in [Naveau et al., 2009] and we will explain the

augmentation of the Mean Squared Error while h is close to zero. This phenomenon would be also thoroughly explained through simulation and a counterexample.

In Chapter 3 we will present in details the mathematical proof of our statement.

Chapter 1

On the variance of the Madogram

Chapter 2

Numerical results

Chapter 3

Mathematical section

Bibliography

- [Angrist and Pischke, 2008] Angrist, J. D. and Pischke, J.-S. (2008). *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press.
- [Bador et al., 2015] Bador, M., Naveau, P., Gilleland, E., Castellà, M., and Arivelo, T. (2015). Spatial clustering of summer temperature maxima from the cnrm-cm5 climate model ensembles & e-obs over europe. *Weather and Climate Extremes*, 9:17–24. The World Climate Research Program Grand Challenge on Extremes – WCRP-ICTP Summer School on Attribution and Prediction of Extreme Events.
- [Baraud et al., 2016] Baraud, Y., Birgé, L., and Sart, M. (2016). A new method for estimation and model selection: ρ -estimation. *Inventiones Mathematicae*.
- [Beirlant et al., 2004] Beirlant, J., Goegebeur, Y., Segers, J., and Teugels, J. (2004). *Statistics of Extremes: Theory and Applications*. Wiley. Pagination: 522.
- [Bernard et al., 2013] Bernard, E., Naveau, P., Vrac, M., and Mestre, O. (2013). Clustering of Maxima: Spatial Dependencies among Heavy Rainfall in France. *Journal of Climate*, 26(20):7929–7937.
- [Capéraà et al., 1997] Capéraà, P., Fougères, A.-L., and Genest, C. (1997). A nonparametric estimation procedure for bivariate extreme value copulas. *Biometrika*, 84:567–577.
- [Cochran, 2007] Cochran, W. (2007). *Sampling Techniques, 3Rd Edition*. A Wiley publication in applied statistics. Wiley India Pvt. Limited.
- [Coles et al., 1999] Coles, S., Heffernan, J., and Tawn, J. (1999). Dependence measures for extreme value analyses. *Extremes*, 2:339 – 365.
- [Cooley et al., 2006] Cooley, D., Naveau, P., and Poncet, P. (2006). *Variograms for spatial max-stable random fields*, pages 373–390. Springer New York, New York, NY.
- [Deheuvels, 1991] Deheuvels, P. (1991). On the limiting behavior of the pickands estimator for bivariate extreme-value distributions. *Statistics & Probability Letters*, 12(5):429–439.
- [Escobar-Bach et al., 2018] Escobar-Bach, M., Goegebeur, Y., and Guillou, A. (2018). Local robust estimation of the Pickands dependence function. *The Annals of Statistics*, 46(6A):2806 – 2843.
- [Gaetan and Guyon, 2008] Gaetan, C. and Guyon, X. (2008). *Modélisation et statistique spatiales*. Mathématiques & applications. Springer, Berlin Heidelberg New York.
- [Gudendorf and Segers, 2009] Gudendorf, G. and Segers, J. (2009). Extreme-value copulas.

- [Hall and Tajvidi, 2000] Hall, P. and Tajvidi, N. (2000). Distribution and dependence-function estimation for bivariate extreme-value distributions. *Bernoulli*, 6(6):835–844.
- [Huber, 1964] Huber, P. J. (1964). Robust Estimation of a Location Parameter. *The Annals of Mathematical Statistics*, 35(1):73 – 101.
- [Huber, 2011] Huber, P. J. (2011). *Robust Statistics*, pages 1248–1251. Springer Berlin Heidelberg, Berlin, Heidelberg.
- [Lerasle et al., 2019] Lerasle, M., Szabó, Z., Mathieu, T., and Lecué, G. (2019). MONK – Outlier-Robust Mean Embedding Estimation by Median-of-Means. In *ICML 2019 - 36th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Long Beach, United States.
- [Little R.J.A., 1987] Little R.J.A., R. D. (1987). *Statistical analysis with missing data*. Wiley Series in Probability and Statistics.
- [Marcon et al., 2017] Marcon, G., Padoan, S., Naveau, P., Muliere, P., and Segers, J. (2017). Multivariate nonparametric estimation of the pickands dependence function using bernstein polynomials. *Journal of Statistical Planning and Inference*, 183:1–17.
- [Naveau et al., 2009] Naveau, P., Guillou, A., Cooley, D., and Diebolt, J. (2009). Modeling pairwise dependence of maxima in space. *Biometrika*, 96(1):1–17.
- [Pickands, 1981] Pickands, J. (1981). Multivariate extreme value distribution. *Proceedings 43th, Session of International Statistical Institution, 1981*.
- [Segers, 2014] Segers, J. (2014). Hybrid copula estimators.
- [Sklar, 1959] Sklar, A. (1959). Fonctions de répartition à n dimensions et leurs marges. *Publications de l’Institut de Statistique de l’Université de Paris*, 8:229–231.
- [van der Vaart and Wellner, 1996] van der Vaart, A. W. and Wellner, J. A. (1996). *Weak Convergence and Empirical Process: With Applications to Statistics*. Springer.