

1 **Multivariate extreme value identification using depth functions**

2

3

4

Fateh Chebana* and Taha B.M.J. Ouarda

5

Canada Research Chair on the Estimation of Hydrometeorological Variables,

6

INRS-ETE, 490 rue de la Couronne, Quebec (QC),

7

Canada G1K 9A9

8

9

10

11

12

13

***Corresponding author:** Tel: +1 (418) 654-2542

14

Fax: +1 (418) 654-2600

15

Email: fateh.chebana@ete.inrs.ca

16

17

October 13th 2010

18

Revised version submitted to Environmetrics

19

20 **Abstract**

21 Extreme value theory (EVT) is commonly applied in several fields such as finance, hydrology
22 and environmental modeling. It is extensively developed in the univariate setting. A number of
23 studies have focused on the extension of EVT to the multivariate context. However, most of these
24 studies are based on a direct extension of univariate extremes. In the present paper, we present a
25 procedure to identify the extremes in a multivariate sample. The present procedure is based on
26 the statistical notion of depth function combined with the orientation of the observations. The
27 extreme identification itself is important and it can also serve as basis for the modeling and the
28 asymptotic studies. The proposed procedure is also employed to detect peaks-over-thresholds in
29 the multivariate setting. This method is general and includes several special cases. Furthermore, it
30 is flexible and can be applied to several situations depending on the degree of extreme event risk.
31 The procedure is mainly motivated by application considerations. A simulation study is carried
32 out to evaluate the performance of the procedure. An application, based on air quality data, is
33 presented to show the various elements of the procedure. The procedure is also shown to be
34 useful in other statistical areas.

35

36

37 **1. Introduction**

38 Extreme value theory (EVT) plays an important role in several fields such as finance, hydrology,
39 insurance and Internet traffic, see e.g., de Haan and Ferreira [2006], Reiss and Thomas [2007].

40 EVT is extensively studied in the univariate setting when the extreme event is only described by
41 one characteristic (e.g., Leadbetter et al. [1983] and Coles [2001]). In reality, extreme events are
42 often described through a number of dependent variables. For instance, floods are described by
43 their peak, volume and duration (e.g., Yue et al. [1999]) and air quality is monitored
44 simultaneously through several variables such as the levels of ozone and nitrogen dioxide (e.g.,
45 Heffernan and Tawn [2004]). In the multivariate context, EVT is also developing increasingly to
46 treat these events (e.g., Coles and Tawn [1991; 1994]; Mikosch [2005]; Heffernan and Tawn
47 [2004]; Boldi and Davison [2007]; Li [2009] and the references therein).

48 In the univariate case, extremes are directly identified and the focus is then on modeling efforts.
49 This scheme is commonly extended to the multivariate setting despite the fact that the extreme
50 identification step is not similar and much more complex. In the multivariate literature, the focus
51 is most often on the modeling of extremes, especially on describing the dependence of extreme
52 observations, and also providing asymptotic results (e.g. Li [2009]). However, before developing
53 these important issues (modeling and asymptotic behavior), it is important to correctly identify
54 the notion of extreme in the multivariate context. That is, it is important to specify with respect to
55 which characteristic an observation is extreme, and to identify and quantify the impacts (social,
56 economic) of such extremes. In some fields, such as hydrology and car manufacturing, it is not
57 appropriate to employ asymptotic results since sample sizes are generally small.

58 In the univariate case, the limiting distribution of the block-maxima is shown to belong to the
59 class of the generalized extreme value (GEV) distributions. For practical considerations, even in

60 the univariate setting, where the extremes are simply and clearly identified, the GEV distributions
61 do not represent a systematic choice to fit extremes as shown in El Adlouni et al. [2010]. Several
62 other distributions are appropriate for extreme modeling and should be considered as well, such
63 as the Halphen family, Pearson, log-Pearson and Gamma (see e.g., Hosking and Wallis [1997]).
64 The usual extension of EVT to the multivariate setting is based merely on the *component-wise*
65 maxima or minima of the vector sample. This extension is not appropriate since the obtained
66 point does not necessarily belong to the sample and this extension is based on mathematical and
67 theoretical justifications. In the recent literature, the component-wise approach was criticized
68 from both the theoretical and the practical points of view (see e.g., Smith [2004] and Salvadori et
69 al. [2007]). The multivariate EVT approaches neglected the identification step and focused on the
70 modeling and asymptotic aspects.

71 One of the drawbacks of the multivariate EVT approaches, especially the component-wise
72 extension, is the absence of a convenient notion of ordering in the multivariate context. The
73 notion of order can be statistically extended to the multivariate context using depth functions, see
74 e.g. Zuo and Serfling [2000]. Multivariate EVT does not take advantage of the potential of depth
75 functions. In addition, in the multivariate setting the notion of the median is not employed to
76 define extremes. Usually, in the univariate framework, an extreme observation is the one for
77 which the deviation from the median is the highest. A drawback of the block-maxima approach
78 identification is that it could select extremes in blocks even if all values in a particular block are
79 low and could identify only one extreme in blocks where several high values should be identified
80 in a particular block.

81 The aim of the present paper is to propose a procedure to identify extreme values in a
82 multivariate sample. Then, modeling and studying the asymptotic properties will be based on the
83 appropriate extremes. The present study has an exploratory objective, rather than modeling or

84 inferential. As pointed out in Liu et al. [1999], depth values do not provide enough information
85 from the multivariate sample to define extremes in the present context. The proposed procedure is
86 based on a combination of depth functions with orientations of the observations with respect to
87 the median. This combination (depth and orientation) is also employed by Serfling [2002] to
88 introduce multivariate median-oriented quantiles. The proposed extreme identification procedure
89 is motivated and defined by practical considerations.

90 In the univariate setting, extremes in a sample are selected as the minimum, the maximum or
91 both. The choice is based on the underlying variable, the associated risk as well as the case study
92 in hand. By extension, in the multivariate context, it is also of interest to focus on a *part* of the
93 extreme observations by reducing the orientation space to a convenient part.

94 When dealing with extremes, an alternative of interest is also known as the peaks-over-threshold
95 (POT) approach. In this situation we are interested in identifying all values “over” a given
96 threshold. A more detailed description of the POT approach can be found for instance in Lang et
97 al. [1999] and the references therein. Again as in the extremes, the multivariate POT theory
98 focused on modeling and asymptotic results, see e.g. Reiss and Thomas [2007]. The present
99 depth-based approach is generalized to identify POT observations. By analogy with extremes, the
100 multivariate POT can also be focused on a part of the orientation space.

101 For reference and clarity of presentation, depth functions are presented briefly in section 2. In the
102 following sections, we present a description and a general algorithm of the proposed method
103 (section 3), we evaluate the consistency of the procedure in a simulation study (section 4), we
104 illustrate it on a real-world air quality dataset (section 5), and we present a discussion (section 6).
105 Conclusions and perspectives are presented in the last section.

106

107 **2. Background: Depth functions**

108 In the univariate setting, one of the most important notions related to extremes is the *ordering* of
109 observations. In the multivariate setting several extensions of the order are developed and
110 employed. Depth functions are introduced by Tukey [1975] to provide an outward ordering in a
111 multivariate sample. A detailed description of the theoretical background of depth functions can
112 be found in Zuo and Serfling [2000]. Depth functions are developed for several multivariate
113 statistical applications, e.g. in Mizera and Müller, [2004] and Ghosh and Chaudhuri [2005] and
114 are applied in several areas such as economic and social sciences by Caplin and Nalebuff
115 [1991a; b], industrial quality control by Liu and Singh [1993] and in water sciences by Chebana
116 and Ouarda [2008].

117 A statistical depth function $D(\cdot; F)$, or simply $D(\cdot)$, for a given cumulative distribution function
118 F on R^d ($d \geq 1$) is bounded and nonnegative which provides a F -based center-outward ordering
119 of points x in R^d that satisfies the following properties:

- 120 i. *Affine invariance*
- 121 ii. *Maximality at center*
- 122 iii. *Monotonicity relative to the deepest point*
- 123 iv. *Vanishing at infinity*

124 For a formal definition of depth functions, the reader is referred to Zuo and Serfling [2000]. In
125 the literature several kinds of depth functions are introduced and studied. Here we present a non
126 exhaustive list of some of the key ones:

127 1. The *Mahalanobis depth* is given by:

$$128 \quad MD(x; F) = \frac{1}{1 + d_A^2(x, \mu)} \quad (1)$$

129 where $d_A^2(x, y) = (x - y)' A^{-1} (x - y)$ is the Mahalanobis distance between two points
 130 $x, y \in R^d$ with respect to a positive definite matrix A , F is a given distribution and μ and A are
 131 any corresponding location and covariance measures, respectively.

132 2. The *Simplicial depth* whose expression is given by:

$$133 \quad SD(x; P) = P \{x \in S[X_1, \dots, X_{d+1}]\} \quad (2)$$

134 where $S[X_1, \dots, X_{d+1}]$ is the random d -dimensional simplex with vertices X_1, \dots, X_{d+1} which is a
 135 random sample from the distribution P .

136 3. The *Simplicial volume depth* is given through the expression :

$$137 \quad SVD^\alpha(x; F) = \left(1 + E \left[\left(\frac{\Delta(S[x, X_1, \dots, X_d])}{\sqrt{\det(\Sigma)}} \right)^\alpha \right] \right)^{-1} \quad \text{for } x \in R^d \quad (3)$$

138 where $\Delta(S[x, X_1, \dots, X_d])$ denotes the volume of the d -dimensional simplex $S[x, X_1, \dots, X_d]$, Σ
 139 is the covariance matrix of F and $\alpha > 0$.

140 4. The *Halfspace depth* is defined for $x \in R^d$ with respect to a probability P on R^d as:

$$141 \quad HD(x; P) = \inf \{P(H) : H \text{ a closed halfspace that contains } x\} \quad (4)$$

142 A corresponding sample version of a statistical depth function $D(x; F)$ may be defined by
 143 replacing F with a suitable empirical function \hat{F}_n and denoted by $D_n(x) = D(x; \hat{F}_n)$. The
 144 asymptotic properties of $D_n(x)$ are studied in several papers including Liu [1990], Massé [2002;
 145 2004] and Lin and Chen [2006]. The evaluation of some depth functions is complex and requires
 146 approximations and specific algorithms. For instance, Miller et al. [2003] developed an algorithm
 147 for the computation of the halfspace depth. Recently, Massé and Plante [2009] provided a
 148 package in the R software to evaluate several depth functions.

149 3. Methodology

150 In this section we present a description of the proposed procedure followed by a general
151 algorithm for practical implementation. Probabilistic formulations of the approach as well as a
152 diagnostic of its use are presented. The main notations used throughout the paper are summarized
153 in the notation list and are illustrated in the case study section.

154 3.1. Description of the methodology

155 Let X_1, \dots, X_n be a R^d vector sample of size n , denoted $\Lambda_{n,d}$, where $X_i = (X_{i1}, \dots, X_{id})$ and d is
156 a positive integer ($d \geq 1$). Let M_n represent the multivariate median of the sample. It corresponds,
157 in the present study, to the maximum depth value in the sample. It is natural to assume that the
158 median is the “center” of the sample and an extreme is so with respect to the median. Hence, the
159 median is considered as the origin of the multivariate space and then data are median-centered by
160 translation. The orientation set of the observations is the unit hyper-sphere centered at the median
161 M_n and denoted $\Omega^{d-1}(M_n)$. The space $\Omega^{d-1}(M_n)$ will be denoted Ω^{d-1} after translation of the
162 data to be centered at M_n . In the bivariate case ($d = 2$), the unit sphere Ω^{d-1} reduces to the interval
163 $[-1, 1]$, and it becomes $\{-1, +1\}$ in the univariate case.

164 For each observation i from the sample, we assign a depth value $D_i = D(X_i)$ and an orientation u_i
165 $= u(X_i)$ from Ω^{d-1} . Note that in the bivariate case, the presentation is analogous to the polar
166 coordinates but with depths instead of Euclidian distances. In higher dimensions, the expressions
167 of an orientation u with respect to the Cartesian coordinates x are more complex and can be
168 found, for instance, in Stanley [1990]. For a fixed orientation $u_0 \in \Omega^{d-1}$, we identify an extreme
169 as the observation corresponding to the smallest depth value in that orientation.

170 Since the orientation space Ω^{d-1} is continuous, it is convenient to proceed to its discretization. To
 171 this end, let $\lambda \in (0,1]$ be a coefficient that defines a “partition” of Ω^{d-1} . The obtained partition of
 172 Ω^{d-1} is composed by the subsets $\Pi_{k,\lambda}$, $k = 1, \dots, n_e$ where n_e is related to λ as shown bellow.
 173 The coefficient λ represents the volume of each portion $\Pi_{k,\lambda}$. For each portion $\Pi_{k,\lambda}$, we select
 174 the observation corresponding to the smallest depth value in this portion as the extreme one. On
 175 the other hand, the coefficient λ indicates the size of the sub-sample composed by extreme
 176 observations, say, $n_e = \lfloor 1/\lambda^{d-1} \rfloor$ where $\lfloor \cdot \rfloor$ represents the integer part of a real number. The
 177 condition $1 \leq n_e \leq n$ leads to the constraint on λ : $1/n \leq \lambda^{d-1} \leq 1$. We define the set of the
 178 identified extreme observations $\Sigma_n(\lambda, D)$ as:

$$179 \quad \Sigma_n(\lambda, D) = \{x_k \in \Lambda_{n,d} : D(x_k) \text{ is the samllest in } \Pi_{k,\lambda}, k = 1, \dots, n_e\} \quad (5)$$

180 As special cases for λ , we identify one observation ($n_e = 1$) as extreme over the whole sample
 181 for $\lambda = 1$. When λ is close to zero (which requires n to be large), the number of extreme
 182 observations is high. In the case where $\lambda^{d-1} = 1/n$, all the observations are identified as extremes,
 183 unless located on the same orientation.

184 The selection of the coefficient λ is important to define the extremes. In general, its selection
 185 depends on the context of the case study. At this stage, an acceptable general option could be
 186 related to the number of blocks (n_b) in the traditional block-maxima approach. Indeed, λ can be
 187 selected such that n_e is the same (or approximately) as n_b , that is $\lambda^{d-1} = 1/n_b$. For instance,
 188 blocks are generally related to the length of the temporal series, such as daily, weekly, monthly,
 189 seasonally or annually. This characteristic depends on the application field such as air quality
 190 modeling, hydrology or climatology. The coefficient λ should be small enough to lead to a

191 reasonable number of extremes with which one can perform a statistical analysis and at the same
192 time λ should be large enough to avoid obtaining a large number of extremes which would
193 contradict the rarity principle of extremes. However, an optimal and automatic selection
194 procedure of λ would be useful and should be developed in a future work.

195 By relating the extreme observations in $\Sigma_n(\lambda, D)$ we obtain the hyper-surface $\mathbb{C}_n(\lambda, D)$ which
196 may be convex or not. The coefficient λ controls, in an inversely proportional way, the
197 regularity (or smoothness) of $\mathbb{C}_n(\lambda, D)$. In terms of risk, a small value of λ indicates that the
198 decision is hard and it should be taken with care whereas a large value of λ is associated to safer
199 situations. Indeed, a small value of λ leads to several extreme combinations that should be
200 considered to prevent the associated risk whereas large values of λ are representative for
201 situations that require less attention. The coefficient λ can be interpreted as a confidence degree
202 against the corresponding risk. Risk is often defined as the probability of occurrence of an
203 extreme event (see Niwa [1989] and Ouarda and Labadie [2001]). Hence, a small value of λ
204 indicates that a large number of extreme events have occurred and therefore the probability of
205 occurrence of similar events is high. It is important to mention that, in contrast to the univariate
206 case where we have one extreme observation (minimum or maximum), the existence of several
207 extreme observations in the multivariate context is natural. It can be justified by the fact that
208 several combinations of variable values lead to the same risk.

209 In the univariate case, extreme value refers to the maximum *or* the minimum of a sample.
210 According to the problem to be treated, the focus is made on the minimum *or* the maximum *or*
211 both. The above extreme identification procedure allows to generalize this aspect to the
212 multivariate setting. Indeed, we consider the identification of extremes on a *part* T of Ω^{d-1} of

213 orientations. Hence, the subdivision in portion of volume λ can be limited only to the part T
 214 instead of the whole set Ω^{d-1} . The corresponding set of extreme observations is given by:

$$215 \quad \Sigma_n(\lambda, T, D) = \Sigma_n(\lambda, D) \cap T \quad (6)$$

216 For instance, for $d = 2$ where $\Omega^1 = [-1, 1]$, if the focus is on simultaneous non-exceedence events
 217 $(X \leq x, Y \leq y)$, it is convenient to choose $T = [0, 0.5]$ which corresponds to the first quadrant.

218 In the univariate case where $\Omega^0 = \{-1, 1\}$, the maximum is associated to $T = \{1\}$ whereas the
 219 minimum is associated to $T = \{-1\}$. Note that the volume of the range T should be larger than λ .

220 In the equality case $\lambda = \text{volume}(T)$, we have $e_n = \Sigma_n(\text{volume}(T), T, D)$. By analogy with
 221 $\Sigma_n(\lambda, T, D)$, the hyper-surface $\mathbb{C}_n(\lambda, D)$ can be restricted to a given part T as:

$$222 \quad \mathbb{C}_n(\lambda, T, D) = \mathbb{C}_n(\lambda, D) \cap T \quad (7)$$

223 The present approach can be generalized for the identification of POTs for a given multivariate
 224 sample. In the univariate POT, one of the criteria used to define the threshold is based on a given
 225 percentile of the sample. Hence, in the multivariate framework we select in each λ -portion $\Pi_{k,\lambda}$

226 the observations for which the depth deviations from that of the median do not exceed a given
 227 proportion s ($0 \leq s \leq 1$) of the deviation of the minimum depth in $\Pi_{k,\lambda}$. That is the depth value is

228 smaller than $D_{s,k} = D_{\max} - (D_{\max} - D_{\min,k})s$ and, therefore, the set of POTs is given by:

$$229 \quad \Sigma_n(\lambda, s, D) = \{x_j \in \Lambda_{n,d} \cap \Pi_{k,\lambda} : D(x_j) < D_{s,k}, \Pi_{k,\lambda} \subset \Omega^{d-1}\} \quad (8)$$

230 where D_{\max} is the depth value of the median and $D_{\min,k}$ is the smallest depth value over $\Pi_{k,\lambda}$.

231 Note that for a fixed value of the threshold s , the value $D_{s,k}$ is not necessarily the same for all

232 $\Pi_{k,\lambda}$ since $D_{\min,k}$ depends on k . Clearly, the special case $s = 0$ leads to the selection of all data as

233 POTs whereas $s = 1$ identifies the extreme observations. As it is the case for extremes, the POT
 234 may also be of interest in a part $T \subset \Omega^{d-1}$. The corresponding set is:

$$235 \quad \Sigma_n(\lambda, s, T, D) = \Sigma_n(\lambda, s, D) \cap T \quad (9)$$

236 The hyper-surfaces $\mathbb{C}_n(\lambda, s, D)$ and $\mathbb{C}_n(\lambda, s, T, D)$ can be defined, for instance, by connecting
 237 the observations with the largest depth value and smaller than $D_{s,k}$ in each portion $\Pi_{k,\lambda}$. Other
 238 options are presented in Sections 5 and 6.

239 In the depth-based approach, the case with $\lambda = 1$ corresponds to the usual POT approach in the
 240 whole data set. However, smaller values of λ are useful to adapt the approach in the presence of
 241 trend or seasonality in the data.

242 Now that the descriptive presentation of the depth-based identification approach is complete, we
 243 provide a brief probabilistic formulation of the approach. Assume that the original random
 244 vectors X_1, X_2, \dots have a multivariate distribution F on R^d , then the random vector of the depth-
 245 based extreme values is given by:

$$246 \quad B_{k,\lambda} = \arg \min_{X_i \in \Lambda_{n,d} \text{ with } u(X_i) \in \Pi_{k,\lambda}} D(X_i) \quad (10)$$

247 The exact or asymptotic distribution of the random vector $B_{k,\lambda}$ is related to F as well as to the
 248 distribution of $D(X)$. Basically, the problem can be seen as a minimization of a special
 249 transformation $D(\cdot)$ of the original random vectors X . In addition, in the present context, the study
 250 of $B_{k,\lambda}$ can be conducted by considering previous work such as Massé [2004; 2009], Arcones et
 251 al. [2006] and Zu and He [2006] where asymptotic results are obtained for $D(X)$. Further
 252 developments in this direction are outside the scope of the present study and are the subject of
 253 future work.

254 In a similar way the extremes on a range T can be defined as:

255
$$B_{k,\lambda,T} = \arg \min_{X_i \in \Lambda_{n,d} \text{ with } u(X_i) \in \Pi_{k,\lambda} \cap T} D(X_i) \quad (11)$$

256 The corresponding POTs on Ω^{d-1} and on a range T are defined respectively as follows:

257
$$K_{k,\lambda,s} = X_i * I\{D(X_i) \leq D_{s,k}, u(X_i) \in \Pi_{k,\lambda}\} \quad (12)$$

258
$$K_{k,\lambda,s,T} = X_i * I\{D(X_i) \leq D_{s,k}, u(X_i) \in \Pi_{k,\lambda} \cap T\} \quad (13)$$

259 where $I\{A\}$ stands for the indicator function of a set A , that is $I\{A\} = 1$ if A holds and 0 if not.

260 As indicated in the introduction, the block-maxima approach imposes a uniform repartition of the
 261 extremes over time (one extreme per time block). The proposed depth-based approach avoids this
 262 constraint since it is based on the magnitudes of the values and not on their time of occurrence. In
 263 situations where it is necessary to define time blocks, an intermediate option could be to combine
 264 both approaches by employing the depth-based approach in each large time block. A large time
 265 block is composed of a number of the usual blocks. For instance, large and usual blocks could be
 266 respectively season and month or year and season. An illustration is given in the case study where
 267 each season has four months and hence $\lambda = 1/4 = 0.25$ for each season.

268 Before presenting the procedure steps, we state a number of simple properties of the above
 269 concepts. For a given sample, using the same depth function D , we have:

270
$$\text{If } \lambda_1 < \lambda_2 \text{ and } \lambda_2/\lambda_1 \text{ is an integer, then } \Sigma_n(\lambda_2, D) \subset \Sigma_n(\lambda_1, D) \quad (14)$$

271 The condition λ_2/λ_1 is an integer insures that for each k , there exists k' such that $\Pi_{k,\lambda_2} \subset \Pi_{k',\lambda_1}$.

272 A counter example is given in the case study section when λ_2/λ_1 is not an integer.

273 For a given sample, on the same part T , using the same depth function D , we have:

274
$$\text{If } s_1 < s_2, \quad \text{then } \Sigma_n(\lambda, s_2, T, D) \subset \Sigma_n(\lambda, s_1, T, D) \text{ for a fixed } \lambda \quad (15)$$

275
$$\text{For } s < 1, \quad \text{we have } \Sigma_n(\lambda, s, T, D) \subset \Sigma_n(\lambda, T, D) \quad (16)$$

276 **3.2. Procedure steps**

277 In the following we present the proposed procedure to identify extremes and POTs for a given
278 multivariate sample. Identification of multivariate extremes requires a depth function D , a
279 coefficient λ and, if necessary, a range $T \subseteq \Omega^{d-1}$. A threshold s in $(0, 1]$ is also to be specified
280 for POTs. On the basis of the description and notations introduced in Section 3.1, the POTs and
281 the extremes are identified through the following steps:

- 282 1. Find the median M_n of the multivariate sample. In the present study, it corresponds to the
283 largest depth value D_{\max} ;
- 284 2. Standardize data, especially when variables are not of the same nature;
- 285 3. Evaluate, in the range $T \subseteq \Omega^{d-1}$, the depth D_i of the observation i using the selected depth
286 function D , $i = 1, \dots, n$;
- 287 4. Evaluate, in the range $T \subseteq \Omega^{d-1}$, the orientation u_i of the observation i , $i = 1, \dots, n$;
- 288 5. Select, in the range $T \subseteq \Omega^{d-1}$, the observations for which the depth values are smaller than a
289 threshold s of $D_{\min,k}$ in each λ -portion $\Pi_{k,\lambda}$, i.e., with depth smaller than
290 $D_{s,k} = (1-s)D_{\max} + sD_{\min,k}$.

291 The above procedure is general and covers all possible scenarios depending on the special cases
292 of T and s . Indeed, the range T is taken to be $T \subseteq \Omega^{d-1}$ and the threshold s is inclusively between
293 0 and 1. The case where T represents the whole space Ω^{d-1} is hence a special case. When $s = 1$,
294 the identified observations are extremes whereas when $s < 1$, the identified observations are
295 POTs. The observations corresponding to the depth values obtained in step 5 constitute,
296 depending on s and T , one of the sets $\Sigma_n(\lambda, D)$, $\Sigma_n(\lambda, T, D)$, $\Sigma_n(\lambda, s, D)$ or $\Sigma_n(\lambda, s, T, D)$.

297 In step 1, various options are available in the literature to obtain the multivariate median, see e.g.,
298 Small [1990]; Liu et al. [1999] and Zuo and Serfling [2000]. The selection of a depth function,
299 among the various options presented in the literature, depends on its convenience for the specific
300 data in hand as well as the simplicity of its evaluation algorithm. Note that a depth function is
301 more general than a simple transformed distance and it combines both geometry and statistics.
302 Generally depth functions are affine invariant, i.e., depth values remain the same after
303 standardization of data. Therefore, step 2 is not required when the depth function is affine
304 invariant. Note that some depth functions, such as the simplicial volume depth, meet this property
305 only under some assumptions on the parent distribution. The reader is referred to Zuo and
306 Serfling [2000] for more details.

307 The choice of a depth function may affect the identification of the extremes. As a first criterion to
308 select a depth function, we propose to consider depth functions which are evaluated with respect
309 to a given centre of the data (such as the median). This is the case for the Mahalanobis and the
310 projection depth functions. Note that, in general, depth functions cannot be directly and
311 analytically evaluated. To this end, numerical algorithms are required and a package in the R
312 software is provided by Massé and Plante [2009] for several depth functions.

313 Since the distribution of the identified extremes is not developed in the present study, the
314 following diagnostic strategy is proposed. First, if the data are highly correlated, the component-
315 wise approach can be adopted to identify extremes and the corresponding models can be selected
316 for further studies by considering the existing literature. This can be checked, for instance, from a
317 scatter plot per block and an evaluation of different dependence association parameters such as
318 the correlation coefficient, Spearman's rho or Kendall's tau (Joe [1997]). The scatter plot should
319 have an elliptical shape in the first diagonal line and the values of the dependence parameters
320 should be high to insure that the identified component-wise extremes are part or closely part of

321 the data (a component-wise extreme is not necessarily an observation). Second, if the above
 322 situation does not occur, which can be the case for several multivariate data sets, we propose to
 323 consider the depth-based extreme identification. Even though this is an exploratory study and
 324 inferential concerns are a subject of future efforts, one can consider the identified sub-sample of
 325 extremes to select the appropriate distribution among those existing in the literature. This can be
 326 done on the basis of goodness-of-fit tests in the multivariate context. For instance, the univariate
 327 GEV and GPD distributions can be employed as marginal distributions and combined to a copula
 328 to obtain the whole multivariate distribution according to Sklar's theorem (Sklar [1959]). Among
 329 the available and convenient copulas, one can consider the extreme value copula or the
 330 Archimedean copula families.

331 Finally, even though the identified depth-based extremes are extreme observations by definition,
 332 it is advised to check them on the basis of the different types of plots such as scatter plots with
 333 the partitions $\Pi_{k,\lambda}$, depth-orientation plots and multiple chronological plots.

334 **4. Procedure evaluation**

335 In the present section, we evaluate the performance of the proposed procedure on the basis of
 336 simulations. To generate samples, we consider a bivariate distribution commonly used in EVT.
 337 The margins of this distribution are the Gumbel distribution given by:

$$338 \quad F_X(x) = \exp\left\{-\exp\left(-\frac{x-\beta_X}{\alpha_X}\right)\right\}, \quad x \text{ real, } \alpha_X > 0 \text{ and } \beta_X \text{ real} \quad (17)$$

339 and the dependence structure is the Gumbel logistic copula expressed as:

$$340 \quad C_\gamma(u,v) = \exp\left\{-\left[(-\log u)^\gamma + (-\log v)^\gamma\right]^{1/\gamma}\right\}, \quad \gamma \geq 1 \text{ and } 0 \leq u, v \leq 1 \quad (18)$$

341 The Gumbel logistic copula C_γ is at the same time an Archimedean copula and an extreme value
 342 copula. The considered parameters of the marginal distributions are $\alpha_X = 300.22$, $\beta_X = 1239.80$

343 and $\alpha_Y = 15.85$, $\beta_Y = 51.85$. The parameter of the Gumbel logistic copula is considered to take
 344 each one of the values $\gamma = 1, 1.414, 3.162$ which correspond respectively to the correlation
 345 coefficient values $\rho = 0, 0.5, 0.9$. The parameters of the margins, with $\gamma = 1.414$, represent a
 346 real-world flood data set studied in Yue et al. [1999] corresponding to the Ashuapmushuan river
 347 basin in the province of Quebec, Canada. The sample generation is based on the algorithm
 348 developed by Ghoudi et al. [1998]. The number of the generated samples is taken to be $M = 2000$
 349 samples (higher values lead to similar results).

350 For the evaluation of the procedure, we select values of $\lambda = 0.05$ and 0.10 and a value of $s = 0.90$
 351 for the POTs. The procedure is judged consistent if for different samples of the same nature, all
 352 identified extremes are similar. Hence, we evaluate the consistency on the basis of the volume of
 353 the polygon composed by the identified extremes. For the k^{th} generated sample, let $V_e(k)$ be the
 354 volume of the polygon $\Pi_n^{(k)}(\lambda, D)$ composed by the set of the identified extreme observations
 355 $\Sigma_n^{(k)}(\lambda, D)$.

356 In a similar manner, the consistency of the POTs identification is evaluated by the volume of the
 357 area between the polygons composed by the extremes and the deepest POTs in each portion. Let
 358 $\tilde{\Sigma}_n^{(k)}(\lambda, s, D)$ be the subset of the identified POTs corresponding to the observations with the
 359 largest depth value among the POTs in each portion. Similarly, define $\tilde{\Pi}_n^{(k)}(\lambda, s, D)$ as the
 360 polygon composed by $\tilde{\Sigma}_n^{(k)}(\lambda, s, D)$ and let $V_{POT}(k)$ be the difference between the volumes of
 361 $\Pi_n^{(k)}(\lambda, D)$ and $\tilde{\Pi}_n^{(k)}(\lambda, s, D)$. Then, the evaluation is based on the mean and the standard-
 362 deviation over the M generated samples of $V_e(\cdot)$ for the extremes and of $V_{POT}(\cdot)$ for the POTs
 363 given respectively by :

$$364 \quad M_e = \frac{1}{M} \sum_{k=1}^M V_e(k) \quad \text{and} \quad STD_e = \sqrt{\frac{1}{M-1} \sum_{k=1}^M (V_e(k) - M_e)^2} \quad (19)$$

$$365 \quad M_{POT} = \frac{1}{M} \sum_{k=1}^M V_{POT}(k) \quad \text{and} \quad STD_{POT} = \sqrt{\frac{1}{M-1} \sum_{k=1}^M (V_{POT}(k) - M_{POT})^2} \quad (20)$$

366 Table 1 presents the evaluation results obtained from the simulations. Results show the general
367 consistency of the identified extremes and POTs for each one of the considered cases. Generally
368 we observe that M_e values have a slight variation (e.g. between 0.85 and 0.97 for
369 $\gamma = 1.414$ and $\lambda = 0.05$) with respect to n whereas STD_e decreases slightly (e.g. from 0.24 to 0.19
370 for the same case). On the other hand, M_{POT} increases with respect to n (e.g. from 0.16 to 0.59 for
371 the same case) and STD_{POT} remains almost constant. Both values of M_{POT} and STD_{POT} are
372 smaller than those corresponding to the extremes. This is mainly due to the definition of $V_e(\cdot)$ and
373 $V_{POT}(\cdot)$ where $V_{POT} \leq V_e$. Hence, we always have $M_{POT} \leq M_e$ and we can also write
374 $V_{POT}(k) \approx aV_e(k)$ for some constant $a < 1$ and independent of the index k since the values of s is
375 constant ($s = 0.9$). Therefore, we have $STD_{POT} \approx aSTD_e \leq STD_e$. Furthermore, no significant
376 differences are observed between the independence and moderate dependence cases ($\gamma = 1$ and
377 $\gamma = 1.414$ respectively) whereas the case of higher dependence ($\gamma = 3.162$) produces clearly
378 smaller values of both mean and standard-deviation. The reason could be related to the shape of
379 the scatter plot of the last case which is more elliptical and concentrated as illustrated in Figure 1.

380 5. Case study

381 In this section we present a case study to illustrate the different aspects of the proposed
382 methodology. We consider air quality monitoring data employed by Heffernan and Tawn [2004].
383 The data is represented by a series of summer daily maximum measurements (in parts per billion)
384 of ground level ozone (O_3) and nitrogen dioxide (NO_2) in Leeds city centre, UK, during the years

385 1994-1998 inclusively with a sample size $n = 578$ measurements (with some missing data). The
 386 summer data set corresponds to observations during the months of April - July. The outliers
 387 mentioned in Heffernan and Tawn [2004] were excluded from these data sets.
 388 The Mahalanobis depth (MD) function given in (1) is considered for its simplicity and convenient
 389 properties (Zuo and Serfling [2000]). The bivariate median is obtained as the observation that
 390 maximizes the MD function. For each observation, we obtain the corresponding MD value as
 391 well as the orientation u . Note that the MD function is affine invariant, i.e., depth values are the
 392 same for the original and the standardized data. The orientation space in the present case is the
 393 interval $[-1, 1]$ since $d = 2$. Figure 2 illustrates the main employed notations for a selected value
 394 of λ ($\lambda = 0.05$). The corresponding sets of extremes $\Sigma_n(\lambda, D)$ and the associated curves
 395 $\mathbb{C}_n(\lambda, D)$ are presented in Figure 3 for each value of $\lambda = 0.0625$ and 0.05 . The regularity of the
 396 curve $\mathbb{C}_n(\lambda, D)$ as well as the number of extreme observations depend on λ . As it can be seen,
 397 the number of extreme observations in the present case is $n_e = 16$ and 20 for $\lambda = 0.0625$ and 0.05
 398 respectively. From Figures 3.a and 3.b, one can see that even though $\lambda = 0.05$ is smaller than $\lambda =$
 399 0.0625 , the set $\Sigma_n(0.0625, D)$ is not included in the set $\Sigma_n(0.05, D)$, since the ratio $0.0625/0.05$
 400 $= 1.25$ does not meet the condition of being an integer as specified in (15).
 401 Since in air quality, high values of both variables O_3 and NO_2 are considered as extreme cases of
 402 air pollution, it is more realistic to restrict attention to the part $T = [0, 0.5]$ representing the first
 403 quadrant as illustrated in Figure 2. The identified extremes and corresponding depth values are
 404 given in Table 2 for $\lambda=0.05$. Table 2 and Figure 4 indicate that depth values of the extremes vary
 405 between 0.0287 and 0.2590 which are associated to the most “outer” and “interior” observations.
 406 For comparison purposes, we present in Table 3 and Figure 4 the component-wise extremes per
 407 month as well as those obtained by considering the depth-based approach on each season

408 composed of four months. Therefore, to identify 20 extreme observations (four extremes in each
409 season) we set $\lambda=0.25$ per season. Table 3 and Figure 4 indicate that out of 20 component-wise
410 extremes only 4 correspond to observations. In addition, the component-wise extreme (68, 105),
411 associated to the smallest depth value 0.0246, is very far from the observations. The two largest
412 depth values 0.8198 and 0.3262 correspond respectively to the component-wise extremes (34, 41)
413 and (45, 45). These two component-wise extremes are very close to the median and are unlikely
414 to be true extremes. On the other hand, the depth-based approach per season identified 3 unusual
415 extremes out of 20 which are relatively close to the median with large depth values 0.2836,
416 0.3550 and 0.4937.

417 In order to check the identified extremes by the different approaches, Figure 5 presents a multiple
418 chronological plot of the series. Figure 5 shows that the component-wise extremes occur at
419 different dates for each variable and in some situations these dates are very distant within the
420 block such as in May 1994 and May 1997. On the one hand, the component-wise approach
421 identified only one extreme in months with high O_3 and NO_2 (e.g. May 1995 and June 1996). On
422 the other hand, it identified an extreme for ordinary months such as July 1998 where during the
423 whole month the levels of O_3 and NO_2 are low. The extremes identified by the depth-based per
424 season approach are generally clustered such as in June-July 1994 and April-May 1998.

425 It is important to point out that the usual scatter plot (O_3 , NO_2) may be misleading since it is
426 based on the Euclidian distance which represents more the geometric aspect of the data whereas
427 the (u, D) plot is more appropriate since it exhibits the probabilistic aspects. For the entire data
428 set, the (u, D) plot is presented in Figure 6 where the extremes are clearly shown. We observe
429 that the range $T = [0, 0.5]$ contains the observations with the lowest depth values. The above
430 elements indicate that the depth-based approach seems more appropriate especially if we take

431 into account the fact that the identified extremes are observations and are generally “far” from the
432 median.

433 Figures 7a,b show the identified POTs for the studied data set where the threshold is taken to be s
434 = 90% with $\lambda = 1$ and 0.05 on the whole data set. It is easier to visualize the threshold in the
435 space (u, D) than in the space (O_3, NO_2) as it is shown in Figures 7c,d for the present case study.
436 Again, we observe that almost all the identified POTs are found to be in the first
437 quadrant $T = [0, 0.5]$.

438 As indicated in Section 3, it is of interest to consider other depth functions. In the following, we
439 considered four depth functions (Mahalanobis, simplicial volume with $\alpha = 1$, halfspace and
440 simplicial). Figure 8 illustrates the histograms for each one of the considered depth functions. It
441 can be seen that the Mahalanobis depth is convenient for the current study. Indeed, the
442 Mahalanobis histogram shows that the majority of data are in the centre (depth values between
443 0.1 and 0.9). However, a smaller portion of the observations is found to be very close to the
444 median (depth values between 0.9 and 1) or is at the boarder (depth values approximately
445 between 0 and 0.1). This distribution of depth values is natural and reflects the distribution of the
446 data, especially the fact that extreme observations are rare. The simplicial volume depth (SVD)
447 function, given in equation (3), could be a good choice as well.

448 **6. Discussion**

449 The identification of extremes treated in the present paper presents some similarities with a
450 number of commonly used statistical techniques, either in their aims or in their concepts. The
451 illustrations in this section are based on the above case study.

452 The hyper-surface $\mathbb{C}_n(\lambda, D)$ obtained by connecting the extreme observations can be employed
453 to define the range of the multivariate sample. The theoretical version of $\mathbb{C}_n(\lambda, D)$ can be seen

454 as the support of the corresponding multivariate distribution. In addition, the volume of the “inner
455 set” with boundary $\mathbb{C}_n(\lambda, D)$ can be employed to measure and compare the spread of
456 multivariate samples. This spread measure can be compared, for instance, to the measures
457 proposed by Liu et al. [1999] which are also based on depth functions.

458 On the other hand, the hyper-surface $\mathbb{C}_n(\lambda, D)$ can also be viewed as the contour that includes
459 the entire sample. In the bivariate case, it is possible to present this curve in the space (u, D) as
460 shown in Figure 6. When presented in this manner, the contours can be associated to frontier
461 estimation. Frontier estimation is a statistical technique useful and commonly employed in
462 econometrics. The reader can refer, for instance, to Simar and Wilson [2000] for a review. Hence,
463 the elements of the present procedure can be useful to frontier estimation problems. In addition,
464 based on the presentation (u_i, D_i) , the estimation of the curve $\mathbb{C}_n(\lambda, D)$ can be considered as in a
465 regression analysis. However, in the current regression estimation we are dealing with minimum
466 values of D whereas in the usual regression analysis the focus is on the mean values and on the
467 global trend of the series. Furthermore, the presentation of the curve $\mathbb{C}_n(\lambda, D)$ in the space
468 (u, D) as an open curve (function) is the opposite of the presentation used in time series analysis
469 to illustrate seasonality trends (see e.g. Cunderlik et al., [2004] and Ouarda et al., [2006]).

470 An analogy can be established between the present procedure and the generalized additive model
471 (GAM) estimation using spline functions [Wood, 2006]. Indeed, the coefficient λ in the present
472 procedure has a similar role to the penalizing coefficient that controls the regularity of the
473 estimated function in GAM inference. The estimated function using GAM is similar to the
474 contour $\mathbb{C}_n(\lambda, D)$ shown in Figure 6. To be more general, and smoother, the contour
475 $\mathbb{C}_n(\lambda, D)$ can be obtained by connecting the extreme observations by functions similar to splines
476 instead of straight lines. An illustration is given in Figure 9 with $\lambda = 0.10$ and 0.05 .

477 Consequently, the developed tools in GAM inference can be adapted to the present context. One
478 of these important tools is the generalized cross-validation technique which can be adapted to
479 select the coefficient λ .

480 Note that one of the criteria to be imposed to the curve $\mathbb{C}_n(\lambda, D)$ is to include all data as well as
481 to be the closest to the data. In other words, $\mathbb{C}_n(\lambda, D)$ should be the convex-hull on each portion
482 $\Pi_{k, \lambda}$. The present procedure leads to a convex-hull that is not only geometrically-based but also
483 statistically-based through depth functions.

484 **7. Conclusions and future work**

485 In the present paper a new procedure is proposed to identify extremes in a multivariate sample.
486 The proposed procedure, as a natural extension of the univariate setting, is based on depth
487 functions and the orientation of the observations toward the median. It can also be used to
488 identify multivariate extremes in a POT framework. From a simulation study, the procedure is
489 shown to be generally consistent. The procedure is applied to a case study representing
490 environmental data. In addition, the component-wise maxima are shown not to represent realistic
491 scenarios in several situations.

492 The identification of extremes is directly useful to build warning environmental or health
493 systems. However, in numerous situations, the identification is not an end in itself. Indeed, the
494 identification is an important step for the study of the asymptotic properties and the modeling of
495 the identified extremes. It is also shown that the obtained extreme sets and curves are related to
496 other statistical topics such as multivariate spread measures, frontier estimation and generalized
497 additive modeling. The proposed procedure is general and offers a large degree of flexibility
498 through the coefficient λ , the range T and the threshold s . It is useful to practitioners as well as
499 to methodologists.

500 Even though a major part of the elements related to the procedure are treated in the present paper,
501 others are worth developing in future work. An important issue to be developed is related to the
502 inferential aspects of the approach including the modeling of the identified extremes. It is also of
503 interest to optimize and automate the selection of the coefficient λ for a given data set. In
504 addition, the coefficient λ does not need to be constant. This issue is analogue to the smoothing
505 window in the kernel density nonparametric estimation. The impact of the choice of the depth
506 function should also be studied thoroughly by considering different depth functions. More
507 precisely, the consistency of the identified extremes according to the depth functions can be
508 considered. Finally, it is of interest to associate the obtained extreme curve with a confidence
509 band representing the identification errors.

510 **Acknowledgments**

511 Financial support for this study was graciously provided by the Natural Sciences and Engineering
512 Research Council (NSERC) of Canada, and the Canada Research Chair Program. The authors
513 would like to thank Barbara Martel for her assistance. The authors wish to thank the Editor,
514 Associate Editor, and two anonymous reviewers whose comments helped considerably improve
515 the quality of the paper.

516

517 **Main notation list:**

d	Dimension of data vector
D	Depth function
$\lambda \in (0, 1]$	Coefficient that defines a “grid” of the sphere Ω^{d-1} and represents the volume of a portion of the unit sphere Ω^{d-1}
MD	Mahalanobis depth function
M_n	Multivariate median of the sample
n	Sample size
$\Omega^{d-1}(M_n)$	Unit sphere centered at the median M_n . It represents the space of orientations
u	Orientation
$\Lambda_{n,d}$	d -dimension sample with size n
$\Pi_{k,\lambda}$	λ -portion from Ω^{d-1}
$\Sigma_n(\lambda, D)$	Set of extreme observations of the sample with coefficient λ using a depth function D
$B_{k,\lambda}$	The corresponding random vector of $\Sigma_n(\lambda, D)$
$\Sigma_n(\lambda, T, D)$	Restriction of the set $\Sigma_n(\lambda, D)$ on the range T
$B_{k,\lambda,T}$	The corresponding random vector of $\Sigma_n(\lambda, T, D)$
$\Sigma_n(\lambda, s, D)$	Set of POT observations of the sample over a threshold s with coefficient λ using a depth function D
$K_{k,\lambda,s}$	The corresponding random vector of $\Sigma_n(\lambda, s, D)$
$\Sigma_n(\lambda, s, T, D)$	Restriction of the set $\Sigma_n(\lambda, s, D)$ on the range T
$K_{k,\lambda,s,T}$	The corresponding random vector of $\Sigma_n(\lambda, s, T, D)$
e_n	Extreme observation when $\lambda = volume(T)$
$\mathbb{C}_n(\lambda, D)$ and $\mathbb{C}_n(\lambda, T, D)$	Extreme hyper-surfaces obtained by connecting the observations of respectively the sets $\Sigma_n(\lambda, D)$ and $\Sigma_n(\lambda, T, D)$
$\mathbb{C}_n(\lambda, s, D)$ and $\mathbb{C}_n(\lambda, s, T, D)$	POT hyper-surfaces obtained by connecting the observations of respectively the sets $\Sigma_n(\lambda, s, D)$ and $\Sigma_n(\lambda, s, T, D)$ in the space (u, D)

519 **Bibliography**

- 520 Arcones, M. A., H. Cui, and Y. Zuo (2006), Empirical depth processes, *Test*, 15(1), 151-177.
521 Boldi, M. O., and A. C. Davison (2007), A mixture model for multivariate extremes, *Journal of*
522 *the Royal Statistical Society. Series B: Statistical Methodology*, 69(2), 217-229.
523 Caplin, A., and B. Nalebuff (1991a), Aggregation and imperfect competition - on the existence of
524 equilibrium, *Econometrica*, 59(1), 25-59.
525 Caplin, A., and B. Nalebuff (1991b), Aggregation and social choice - a mean voter theorem,
526 *Econometrica*, 59(1), 1-23.
527 Chebana, F., and T. B. M. J. Ouarda (2008), Depth and homogeneity in regional flood frequency
528 analysis, *Water Resources Research*, 44(11).
529 Coles, S. (2001), *An introduction to statistical modeling of extreme values*, xiv+208 pp.,
530 Springer-Verlag London Ltd., London.
531 Coles, S. G., and J. A. Tawn (1991), Modeling extreme multivariate events, *Journal of the Royal*
532 *Statistical Society Series B-Methodological*, 53(2), 377-392.
533 Coles, S. G., and J. A. Tawn (1994), Statistical-methods for multivariate extremes - an
534 application to structural design, *Applied Statistics-Journal of the Royal Statistical Society Series*
535 *C*, 43(1), 1-48.
536 Cunderlik, J. M., T. B. M. J. Ouarda, and B. Bobée (2004), Determination of flood seasonality
537 from hydrological records, *Hydrological Sciences Journal*, 49(3), 511-526.
538 de Haan, L., and A. Ferreira (2006), *Extreme value theory, An introduction*, xviii+417 pp.,
539 Springer, New York.
540 El Adlouni, S., F. Chebana, and B. Bobee (2010), Generalized Extreme Value versus Halphen
541 System: Exploratory Study, *Journal of Hydrologic Engineering*, 15(2), 79-89.
542 Ghosh, A. K., and P. Chaudhuri (2005), On maximum depth and related classifiers, *Scandinavian*
543 *Journal of Statistics*, 32(2), 327-350.
544 Ghoudi, K., A. Khoudraji, and L. P. Rivest (1998), Propriétés statistiques des copules de valeurs
545 extrêmes bidimensionnelles, *Canadian Journal of Statistics*, 26(1), 187-197.
546 Heffernan, J. E., and J. A. Tawn (2004), A conditional approach for multivariate extreme values,
547 *Journal of the Royal Statistical Society. Series B: Statistical Methodology*, 66(3), 497-530.
548 Hosking, J. R. M., and J. R. Wallis (1997), *Regional Frequency Analysis: An Approach Based on*
549 *L-Moments*, 240 pp., Cambridge University Press, Cambridge.
550 Joe, H. (1997), *Multivariate models and dependence concepts*, 1st ed., xviii, 399 p. pp., Chapman
551 & Hall, London ; New York.
552 Lang, M., T. B. M. J. Ouarda, and B. Bobée (1999), Towards operational guidelines for over-
553 threshold modeling, *Journal of Hydrology*, 225(3-4), 103-117.
554 Leadbetter, M. R., G. Lindgren, and H. Rootzén (1983), *Extremes and Related Properties of*
555 *Random Sequences and Series.*, Springer, New York.
556 Li, H. (2009), Orthant tail dependence of multivariate extreme value distributions, *Journal of*
557 *Multivariate Analysis*, 100(1), 243-256.
558 Lin, L., and M. H. Chen (2006), Robust estimating equation based on statistical depth, *Statistical*
559 *Papers*, 47(2), 263-278.
560 Liu, R. Y. (1990), On a notion of data depth based on random simplices, *Annals of Statistics*,
561 18(1), 405-414.
562 Liu, R. Y., and K. Singh (1993), A quality index based on data depth and multivariate rank tests,
563 *J. Amer. Statist. Assoc.*, 8, 252-260.

564 Liu, R. Y., J. M. Parelius, and K. Singh (1999), Multivariate analysis by data depth: Descriptive
565 statistics, graphics and inference, *Annals of Statistics*, 27(3), 783-858.

566 Massé, J.-C., and J.-F. Plante (2009), *Depth, an R package for depth functions in multivariate*
567 *analysis* (www.cran.r-project.org).

568 Massé, J. C. (2002), Asymptotics for the Tukey median, *Journal of Multivariate Analysis*, 81(2),
569 286-300.

570 Massé, J. C. (2004), Asymptotics for the Tukey depth process, with an application to a
571 multivariate trimmed mean, *Bernoulli*, 10(3), 397-419.

572 Massé, J. C. (2009), Multivariate trimmed means based on the Tukey depth, *Journal of Statistical*
573 *Planning and Inference*, 139(2), 366-384.

574 Mikosch, T. (2005), How to model multivariate extremes if one must?, *Statistica Neerlandica*,
575 59(3), 324-338.

576 Miller, K., S. Ramaswami, P. Rousseeuw, J. A. Sellares, D. Souvaine, I. Streinu, and A. Struyf
577 (2003), Efficient computation of location depth contours by methods of computational geometry,
578 *Statistics and Computing*, 13(2), 153-162.

579 Mizera, I., and C. H. Müller (2004), Location-scale depth, *Journal of the American Statistical*
580 *Association*, 99(468), 949-966.

581 Niwa, K. (1989), *Knowledge-based risk management in engineering : a case study in human-*
582 *computer cooperative systems*, xi, 132 p. pp., Wiley, New York.

583 Ouarda, T. B. M. J., and J. W. Labadie (2001), Chance-constrained optimal control for
584 multireservoir system optimization and risk analysis, *Stochastic Environmental Research and*
585 *Risk Assessment*, 15(3), 185-204.

586 Ouarda, T. B. M. J., J. M. Cunderlik, A. St-Hilaire, M. Barbet, P. Bruneau, and B. Bobée
587 (2006), Data-based comparison of seasonality-based regional flood frequency methods, *Journal*
588 *of Hydrology*, 330(1-2), 329-339.

589 Reiss, R.-D., and M. Thomas (2007), *Statistical analysis of extreme values with applications to*
590 *insurance, finance, hydrology and other fields.*, Third edition ed., Birkhäuser Verlag, Basel.

591 Salvadori, G., C. De Michele, N. T. Kottegoda, and R. Rosso (2007), *Extremes in Nature: An*
592 *Approach Using Copulas*, Springer

593 Serfling, R. (2002), Quantile functions for multivariate analysis: Approaches and applications,
594 *Statistica Neerlandica*, 56(2), 214-232.

595 Simar, L., and P. W. Wilson (2000), Statistical Inference in Nonparametric Frontier Models: The
596 State of the Art, *Journal of Productivity Analysis*, 13(1), 49-78.

597 Sklar, M. (1959), Fonctions de répartition à n dimensions et leurs marges, *Publ. Inst. Statist.*
598 *Univ. Paris*, 8, 229--231.

599 Small, C. G. (1990), A survey of multidimensional medians, *International Statistical Review*,
600 58(3), 263-277.

601 Smith, R. L. (2004), Discussion of "A conditional approach for multivariate extreme value" by
602 J.E. Heffernan and J.A. Tawn., *Journal of the Royal Statistical Society. Series B: Statistical*
603 *Methodology*, 66(3), 530-532.

604 Stanley, C. R. (1990), Descriptive statistics for N-dimensional closed arrays: A spherical
605 coordinate approach, *Mathematical Geology*, 22(8), 933-956.

606 Tukey, J. W. (1975), Mathematics and picturing data, paper presented at Proceedings of the
607 International Congress on Mathematics, Canadian Math. Congress.

608 Wood, S. N. (2006), *Generalized Additive Models, An introduction with R.*, Chapman & Hall/
609 CRC.

610 Yue, S., T. B. M. J. Ouarda, B. Bobée, P. Legendre, and P. Bruneau (1999), The Gumbel mixed
611 model for flood frequency analysis, *Journal of Hydrology*, 226(1-2), 88-100.
612 Zu, Y. J., and X. M. He (2006), On the limiting distributions of multivariate depth-based rank
613 sum statistics and related tests, *Annals of Statistics*, 34(6), 2879-2896.
614 Zuo, Y., and R. Serfling (2000), General notions of statistical depth function, *Annals of Statistics*,
615 28(2), 461-482.
616
617

618 List of Tables and Figure Captions

619 Table 1: Evaluation of the consistency of the extreme and POT identification procedures based on the
620 polygon volume.

621 Table 2: Original and standardized values of (O_3 , NO_2) of the identified extreme observations
622 corresponding to $\lambda=0.05$ as well as their MD depth values in the first quadrant.

623 Table 3: Values of (O_3 , NO_2) of the identified extremes corresponding to $\lambda = 0.25$ within each season as
624 well as their MD depth values in the first quadrant (left); similar values using the component-wise
625 approach within monthly blocks (right).

626 Figure 1: Scatter plot illustration for samples generated with $n = 300$ and a) $\gamma = 1$ b) $\gamma = 1.414$ and c)
627 $\gamma = 3.162$

628 Figure 2: Illustration of a λ -portion $\Pi_{k,\lambda}$, the orientation interval $[-1, 1]$, the set of extreme observations
629 $\Sigma_n(\lambda, D)$, the corresponding curve $\mathbb{C}_n(\lambda, D)$, and a part T as the first quadrant of (O_3 , NO_2) for $\lambda =$
630 0.05 where D is MD

631 Figure 3: Identified extreme observations set $\Sigma_n(\lambda, D)$ of (O_3 , NO_2) and the corresponding curve
632 $\mathbb{C}_n(\lambda, D)$ for a) $\lambda = 0.0625$ and b) $\lambda = 0.05$ where D is MD

633 Figure 4: Extremes identified as component-wise, depth-based ($\lambda=0.05$) and depth-based per season
634 ($\lambda=0.25$) in the first quadrant

635 Figure 5: Chronological (O_3 , NO_2) series and the extremes identified as component-wise, depth-based
636 ($\lambda=0.05$) and depth based per season ($\lambda=0.25$) in the first quadrant. The vertical lines indicate month limits
637 in each season

638 Figure 6: Identified extreme observation set $\Sigma_n(\lambda, D)$ of (O_3 , NO_2) and the corresponding curve
639 $\mathbb{C}_n(\lambda, D)$ in the space (u, D) for a) $\lambda = 0.10$ and b) $\lambda = 0.05$ where D is MD

640 Figure 7: Identified POT observation set $\Sigma_n(\lambda, s, D)$ of (O_3 , NO_2) and the corresponding curve
641 $\mathbb{C}_n(\lambda, s, D)$ for a) $\lambda = 1$ and b) $\lambda = 0.05$ where D is MD and $s = 0.90$; and for c) $\lambda = 1$ and d) $\lambda = 0.05$ in
642 the space (u, D)

643 Figure 8: Histograms of depth values of the data set (O_3 , NO_2) for different depth functions

644 Figure 9: Identified extreme observation set $\Sigma_n(\lambda, D)$ of (O_3 , NO_2) and the corresponding smooth curve
645 $\mathbb{C}_n(\lambda, D)$ in the space (u, D) for a) $\lambda = 0.10$ and b) $\lambda = 0.05$ where D is MD

646

647 Table 1: Evaluation of the consistency of the extreme and POT identification procedures based
 648 on the polygon volume.

649

		$\lambda = 0.05$				$\lambda = 0.1$			
		Extremes		POT with $s = 0.9$		Extremes		POT with $s = 0.9$	
		M_e	STD_e	M_{POT}	STD_{POT}	M_e	STD_e	M_{POT}	STD_{POT}
$\gamma = 1$	n = 100	0.87	0.23	0.14	0.07	1.09	0.26	0.29	0.11
	n = 300	0.95	0.21	0.36	0.09	1.08	0.23	0.52	0.12
	n = 500	0.96	0.20	0.45	0.09	1.06	0.20	0.57	0.11
	n = 1000	0.96	0.18	0.54	0.10	1.03	0.18	0.63	0.11
$\gamma = 1.414$	n = 100	0.85	0.24	0.16	0.07	1.05	0.28	0.31	0.11
	n = 300	0.95	0.22	0.39	0.09	1.05	0.23	0.54	0.11
	n = 500	0.95	0.20	0.48	0.09	1.05	0.20	0.61	0.11
	n = 1000	0.97	0.19	0.59	0.10	1.04	0.19	0.68	0.12
$\gamma = 3.162$	n = 100	0.44	0.16	0.11	0.05	0.45	0.17	0.13	0.07
	n = 300	0.46	0.13	0.19	0.05	0.45	0.14	0.21	0.07
	n = 500	0.46	0.12	0.22	0.05	0.46	0.13	0.24	0.07
	n = 1000	0.46	0.11	0.27	0.06	0.45	0.11	0.28	0.07

650

651

652 Table 2: Original and standardized values of (O₃, NO₂) of the identified extreme observations
 653 corresponding to $\lambda=0.05$ as well as their *MD* depth values in the first quadrant.
 654

O3	NO2	Standardized O3	Standardized NO2	Depth
74	37	0.7167	0.0286	0.0528
80	40	0.8167	0.0714	0.0418
64	44	0.5500	0.1286	0.0894
84	53	0.8833	0.2571	0.0365
71	52	0.6667	0.2429	0.0615
53	46	0.3667	0.1571	0.1759
71	61	0.6667	0.3714	0.0565
65	60	0.5667	0.3571	0.0733
58	59	0.4500	0.3429	0.1018
64	70	0.5500	0.5000	0.0617
69	86	0.6333	0.7286	0.0372
63	79	0.5333	0.6286	0.0505
58	85	0.4500	0.7143	0.0471
40	55	0.1500	0.2857	0.2590
42	61	0.1833	0.3714	0.1712
38	60	0.1167	0.3571	0.1938
46	105	0.2500	1.0000	0.0287
36	62	0.0833	0.3857	0.1722
37	82	0.1000	0.6714	0.0621
32	58	0.0167	0.3286	0.2206

655
 656 Bold character indicates the (O₃, NO₂) corresponding to the largest and smallest depth values.

657
 658
 659
 660
 661

Table 3: Values of (O₃, NO₂) of the identified extremes corresponding to $\lambda = 0.25$ within each season as well as their *MD* depth values in the first quadrant (left); similar values using the component-wise approach within monthly blocks (right).

Depth-based per season			Component-wise		
O3	NO2	Depth	O3	NO2	Depth
55	41	0.1329	45	58	0.1850
71	61	0.0456	46	62	0.1466
48	62	0.1339	53	78	0.0644
34	78	0.0806	71	71	0.0483
74	37	0.0640	46	60	0.1616
64	70	0.0945	71	86	0.0357
69	86	0.0566	63	79	0.0505
37	82	0.0678	74	60	0.0509
84	53	0.0453	52	62	0.1194
46	54	0.2595	41	55	0.2514
43	55	0.2836	68	105	0.0246
46	105	0.0292	84	53	0.0365
46	46	0.2100	39	60	0.1916
57	81	0.0439	58	66	0.0840
39	66	0.1138	53	51	0.1608
28	47	0.4937	57	81	0.0542
59	42	0.0694	47	63	0.1357
40	41	0.3550	59	62	0.0903
42	61	0.0989	45	45	0.3262
36	62	0.1114	34	41	0.8198

662
 663
 664
 665

Bold character indicates the (O₃, NO₂) corresponding to the largest and smallest depth values.
 Shaded character indicates component-wise extremes that coincide with observations.

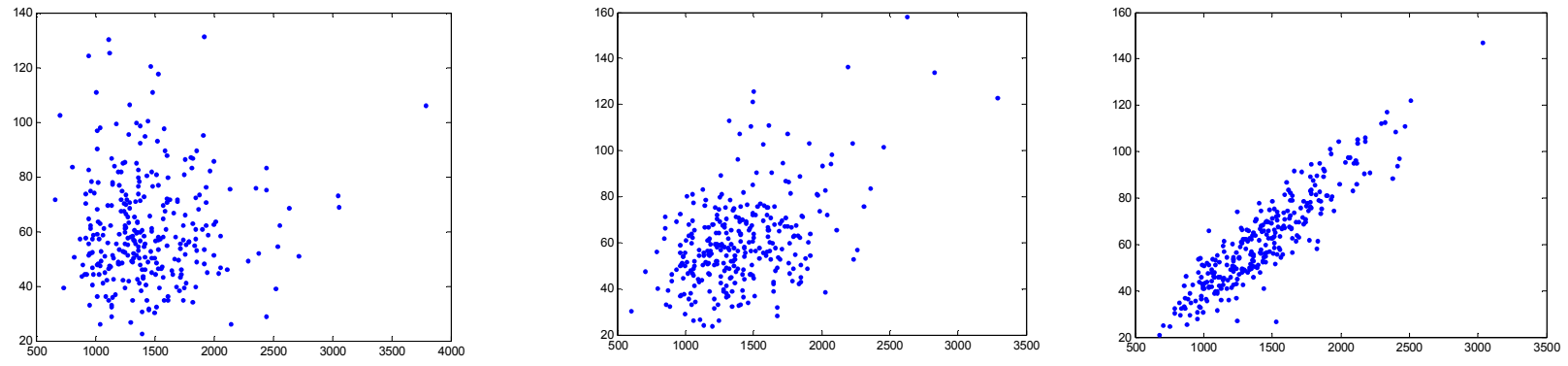


Figure 1: Scatter plot illustration for samples generated with $n = 300$ and a) $\gamma = 1$ b) $\gamma = 1.414$ and c) $\gamma = 3.162$

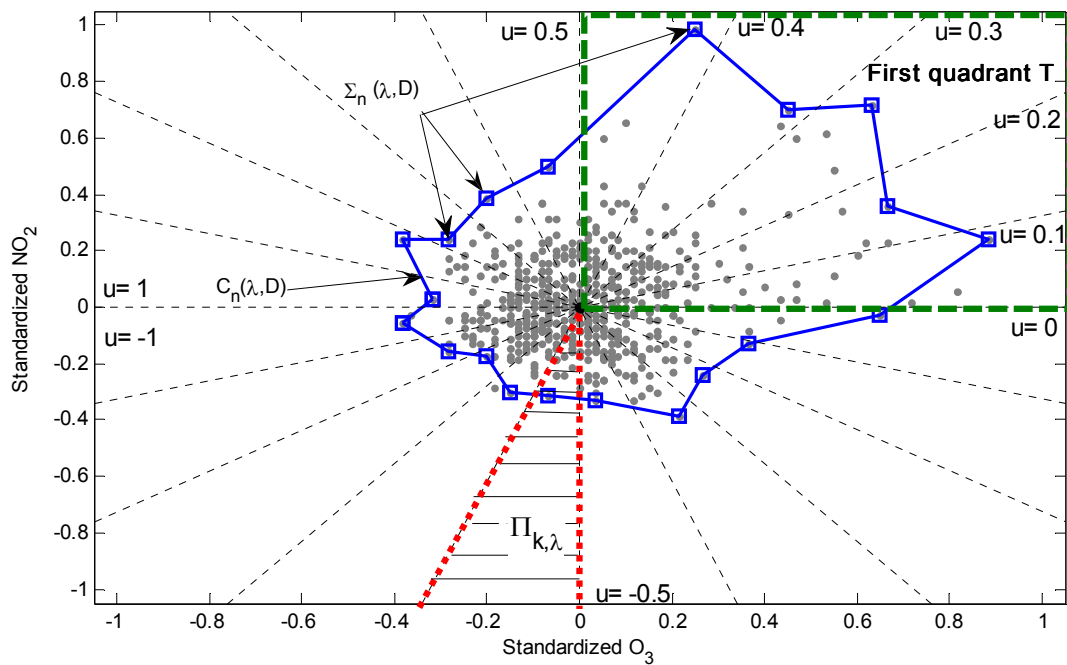
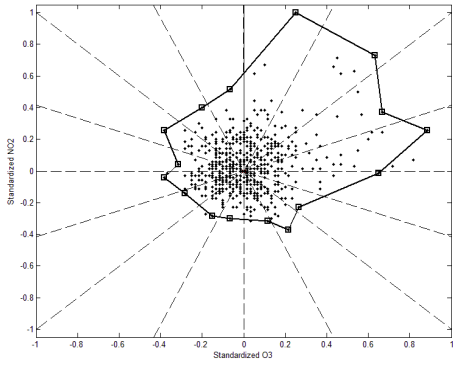
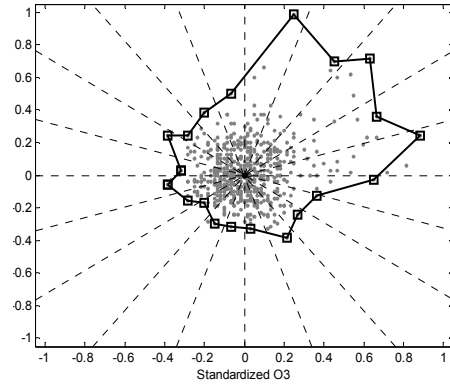


Figure 2: Illustration of a λ -portion $\Pi_{k,\lambda}$, the orientation interval $[-1, 1]$, the set of extreme observations $\Sigma_n(\lambda, D)$, the corresponding curve $\mathbb{C}_n(\lambda, D)$, and a part T as the first quadrant of (O_3, NO_2) for $\lambda = 0.05$ where D is MD



a)



b)

Figure 3: Identified extreme observations set $\Sigma_n(\lambda, D)$ of (O_3, NO_2) and the corresponding curve $\mathbb{C}_n(\lambda, D)$ for a) $\lambda = 0.0625$ and b) $\lambda = 0.05$ where D is MD

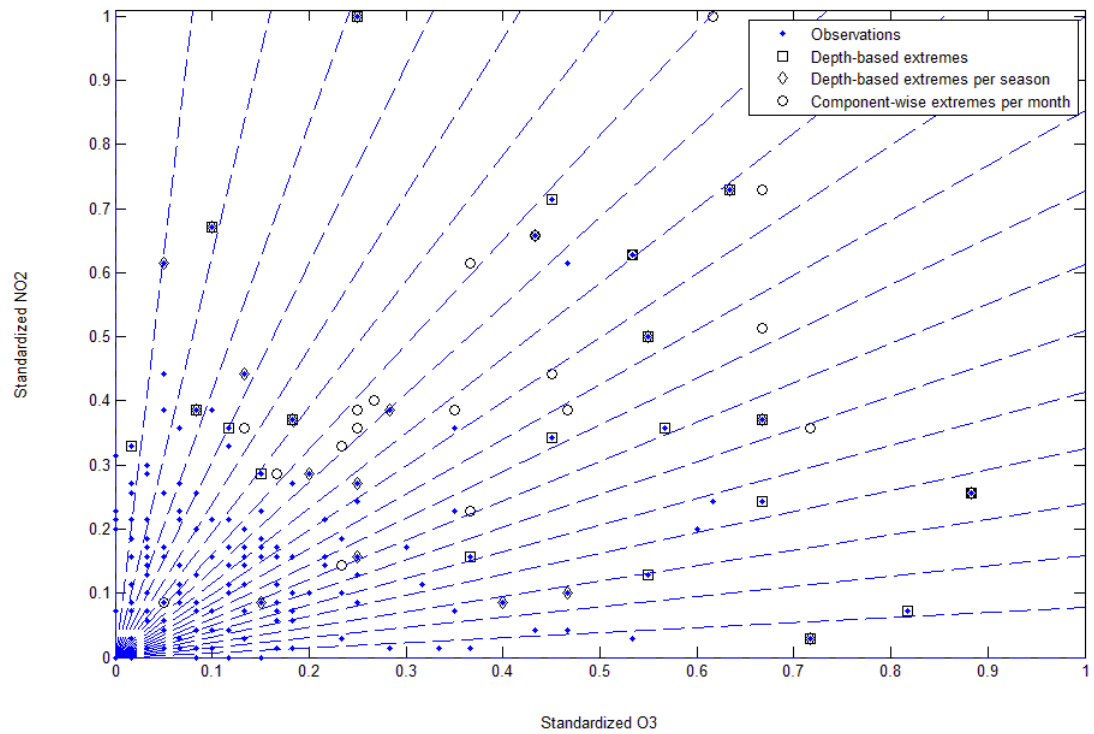


Figure 4: Extremes identified as component-wise, depth-based ($\lambda=0.05$) and depth-based per season ($\lambda=0.25$) in the first quadrant.

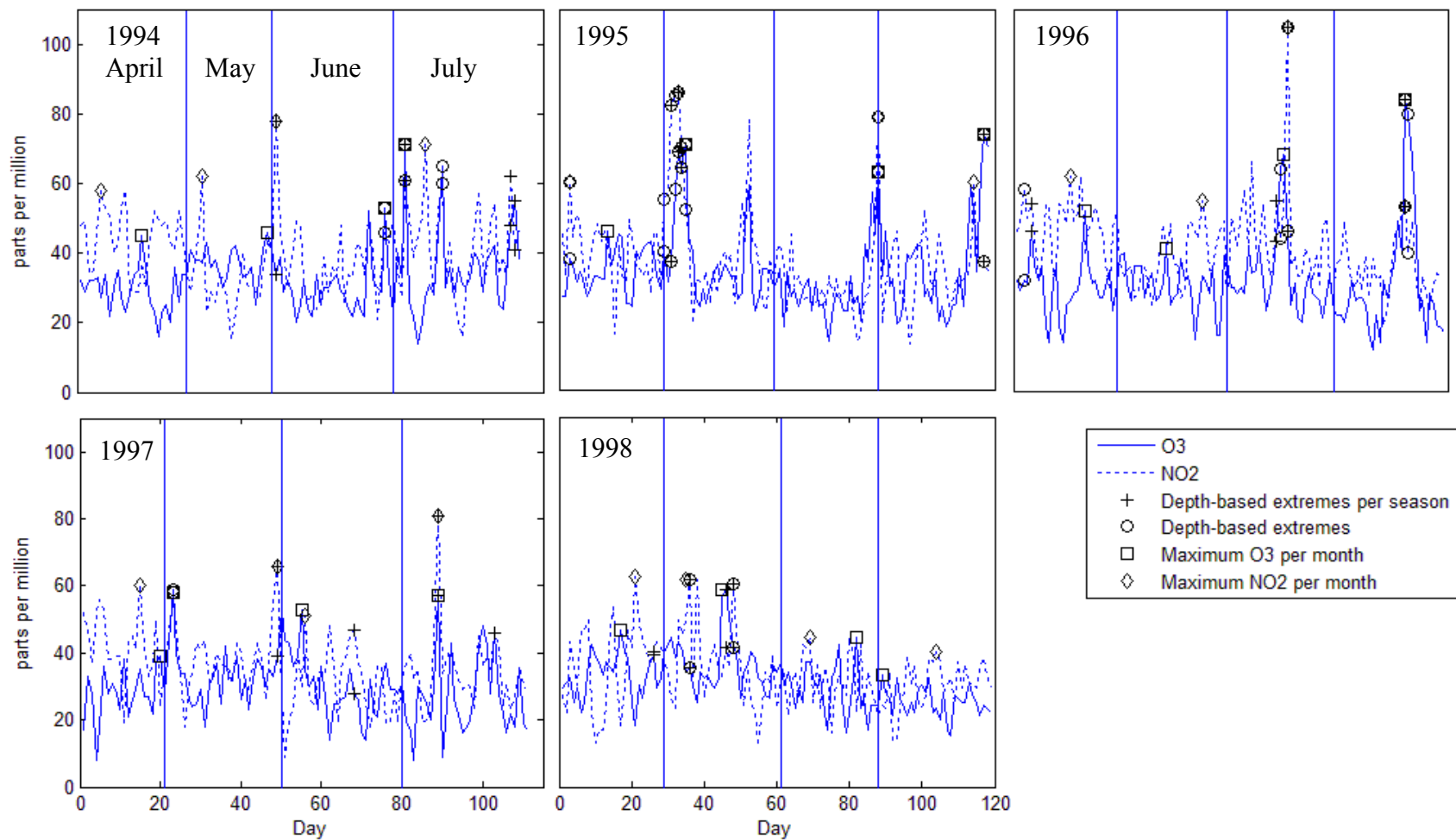
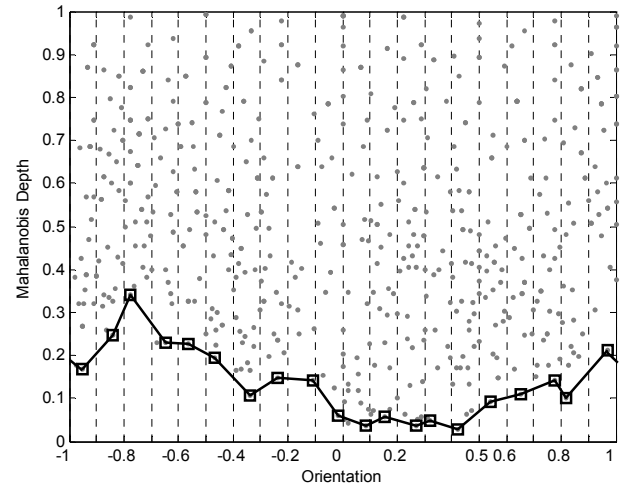
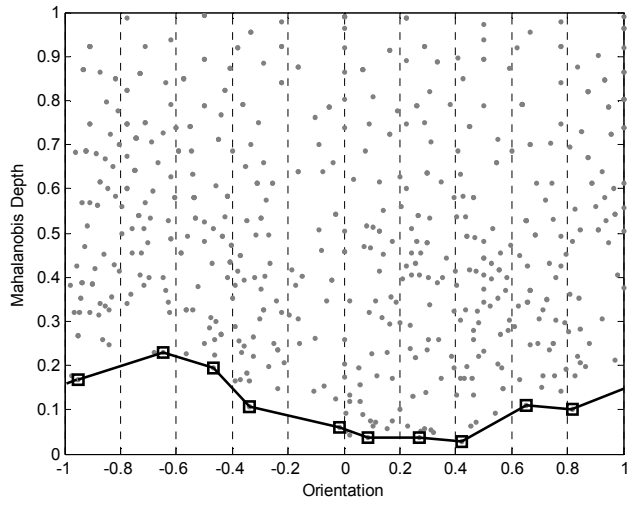


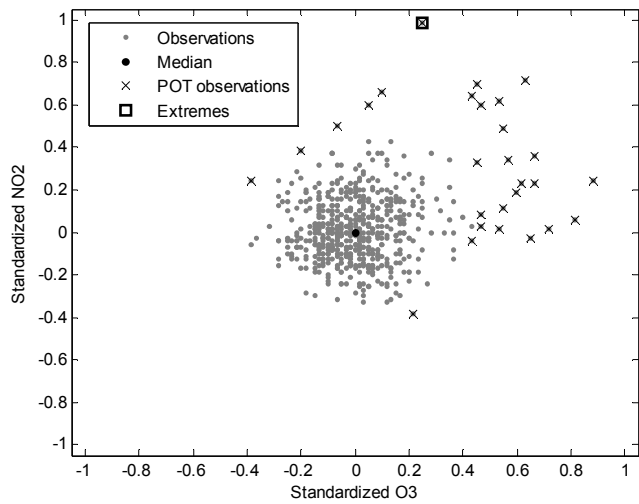
Figure 5: Chronological (O3, NO2) series and the extremes identified as component-wise, depth-based ($\lambda=0.05$) and depth based per season ($\lambda=0.25$) in the first quadrant. The vertical lines indicate month limits in each season



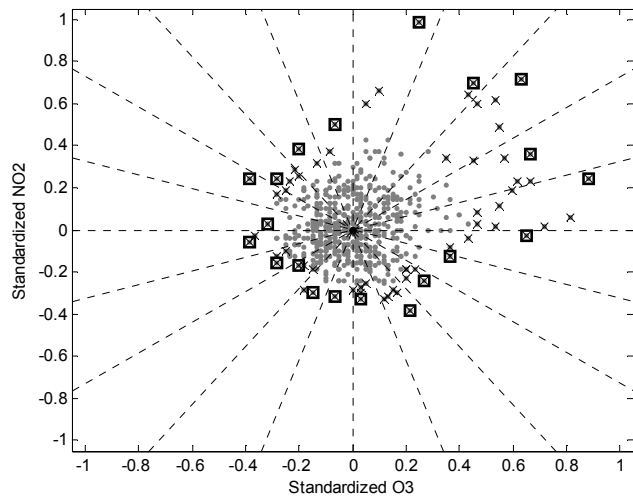
a)

b)

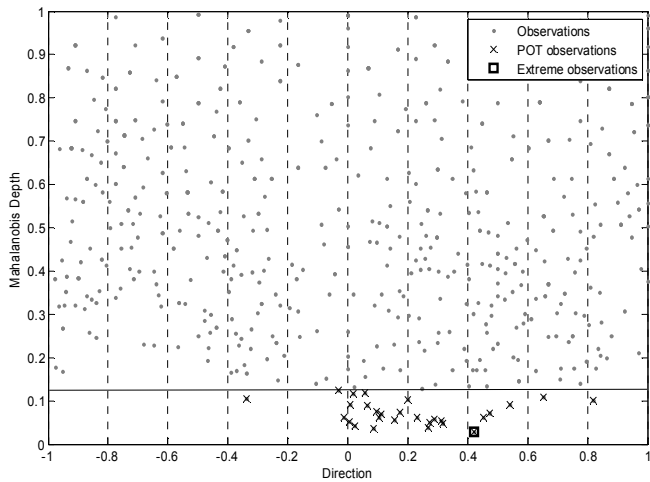
Figure 6: Identified extreme observation set $\Sigma_n(\lambda, D)$ of (O_3, NO_2) and the corresponding curve $\mathbb{C}_n(\lambda, D)$ in the space (u, D) for a) $\lambda = 0.10$ and b) $\lambda = 0.05$ where D is MD



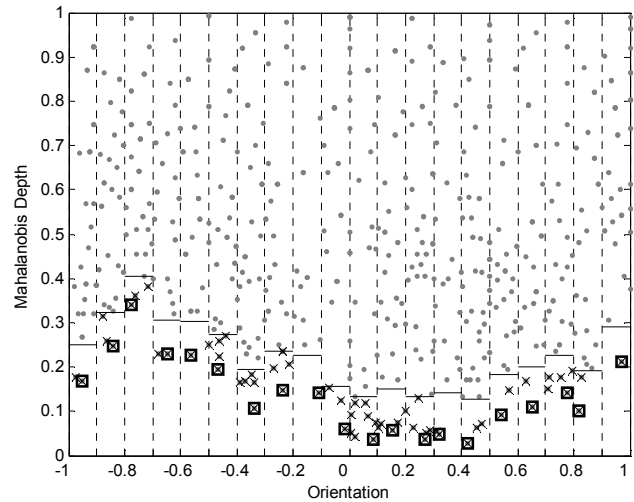
a)



b)

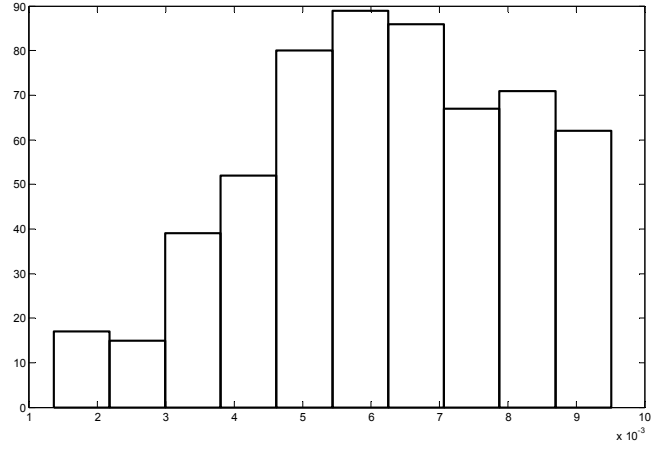
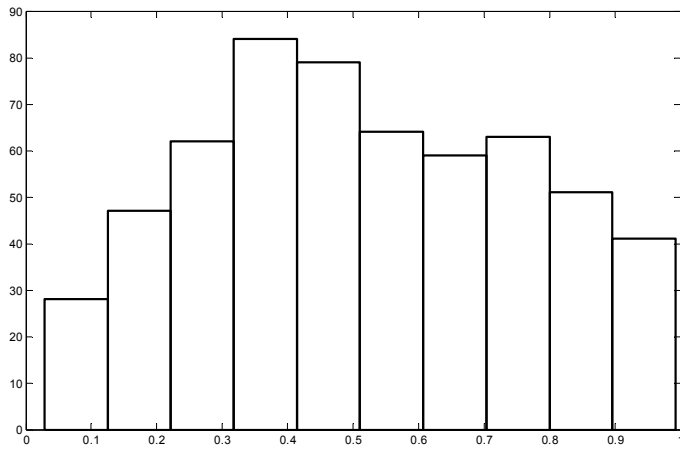


c)



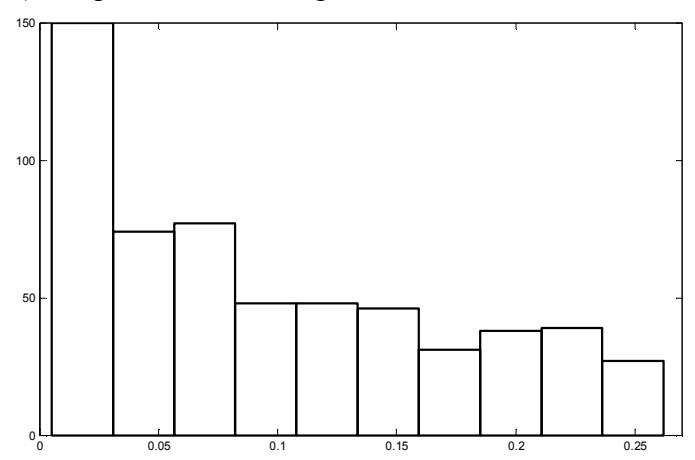
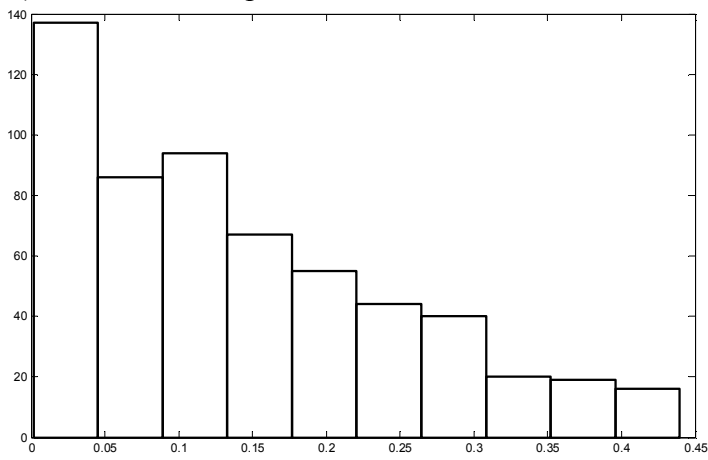
d)

Figure 7: Identified POT observation set $\Sigma_n(\lambda, s, D)$ of (O_3, NO_2) and the corresponding curve $\mathbb{C}_n(\lambda, s, D)$ for a) $\lambda = 1$ and b) $\lambda = 0.05$ where D is MD and $s = 0.90$; and for c) $\lambda = 1$ and d) $\lambda = 0.05$ in the space (u, D)



a) Mahalanobis depth

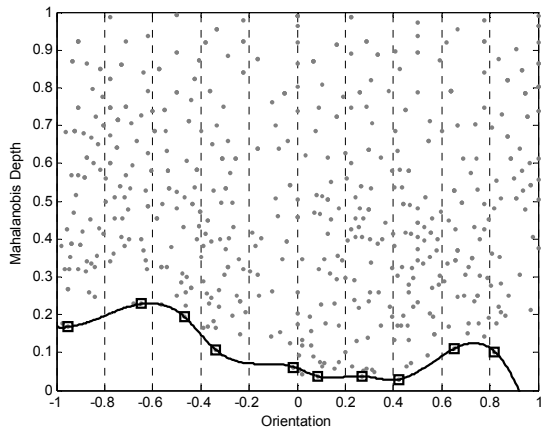
b) Simplicial volume depth



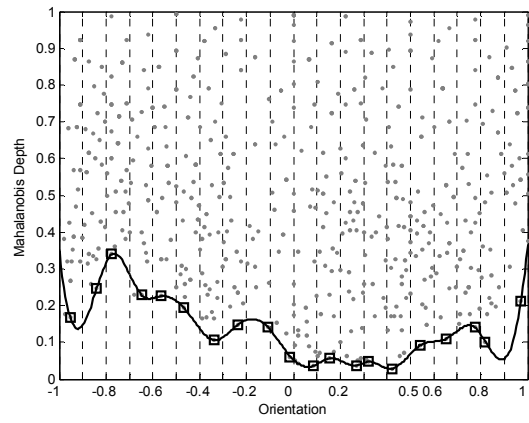
c) Halfspace depth

d) Simplicial depth

Figure 8: Histograms of depth values of the data set (O_3 , NO_2) for different depth functions



a)



b)

Figure 9: Identified extreme observation set $\Sigma_n(\lambda, D)$ of (O_3, NO_2) and the corresponding smooth curve $\mathbb{C}_n(\lambda, D)$ in the space (u, D) for a) $\lambda = 0.10$ and b) $\lambda = 0.05$ where D is MD