1 # Depth-based multivariate descriptive statistics

2 # with hydrological applications

3 F. Chebana[*] and T.B.M.J. Ouarda

4 *Canada Research Chair on the Estimation of Hydrometeorological Variables,*

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

6 *Canada G1K 9A9*

7

8 [*]**Corresponding author:**     Tel: (418) 654-2542

9                                 Fax: (418) 654-2600

10                                Email: fateh.chebana@ete.inrs.ca

11                                November 8[th] 2010

12

## Abstract

Hydrological events are often described through various characteristics which are generally correlated. To be realistic, these characteristics are required to be considered jointly. In multivariate hydrological frequency analysis, the focus has been made on modelling multivariate samples using copulas. However, prior to this step data should be visualized and analyzed in a descriptive manner. This preliminary step is essential for all of the remaining analysis. It allows to obtain information concerning the location, scale, skewness and kurtosis of the sample as well as outlier detection. These features are useful to exclude some unusual data, to make different comparisons and to guide the selection of the appropriate model. In the present paper we introduce methods measuring these features, and which are mainly based on the notion of depth function. The application of these techniques is illustrated on two real-world streamflow data sets from Canada. In the Ashuapmushuan case study, there are no outliers and the bivariate data are likely to be elliptically symmetric and heavy–tailed. The Magpie case study contains a number of outliers, which are identified to be real observed data. These observations cannot be removed and should be accommodated by considering robust methods for further analysis. The presented depth-based techniques can be adapted to a variety of hydrological variables.

## 1. Introduction

Extreme hydrological events, such as floods, storms and droughts may have serious economic and social consequences. Frequency analysis (FA) procedures are commonly used for the analysis and prediction of such extreme events. Relating the magnitude of extreme events to their frequency of occurrence, through the use of probability distributions, is the principal aim of FA [*Chow et al., 1988*].

Generally, several correlated characteristics are required to correctly describe hydrological events. For instance, floods are described by their volume, peak and duration (e.g., Yue et al. [1999]; Ouarda et al. [2000]; Shiau [2003]; Zhang and Singh [2006] and Chebana and Ouarda [2010]). All aspects of univariate FA have already been studied extensively, see e.g. Cunnane [1987] and Rao and Hamed [2000]. On the other hand, multivariate FA has recently attracted increasing attention and the importance of jointly considering all variables characterizing an event was clearly pointed out. Justifications for adopting the multivariate framework to treat extreme events were discussed in several studies (see Chebana and Ouarda [2010] for a summary). For instance, single-variable hydrological FA can only provide limited assessment of extreme events whereas the joint study of the probabilistic characteristics leads to a better understanding of the phenomenon.

In the multivariate hydrological FA literature, the following issues have been addressed: (1) showing the importance and the usefulness of the multivariate framework, (2) selecting the appropriate copula and the marginal distributions and estimating their parameters, (3) defining and studying bivariate return periods, and (4) introducing multivariate quantiles. However, with any statistical analysis, the first stage of the study should be a close inspection of the data. If the data are found appropriate, further analysis of the issues listed above can be undertaken. Hence, exploratory

55    analysis of the data is often the initial stage of any modelling effort that uses that data. It allows to

56    understand the nature of the phenomena that generate the data. It is also useful for model selection

57    and sample comparison. This step of the study is often completely neglected in the multivariate

58    hydrological FA literature. The reason could be the unavailability of the required appropriate tools to

59    carry out this step in a clear and practical manner. Nevertheless, this step is commonly carried out in

60    practice in any univariate hydrological FA study as pointed out by Helsel, et al. [2002]. The

61    development of equivalent tools for the multivariate framework should help promote the use of

62    multivariate FA in hydrological practice.

63    Exploratory descriptive analysis consists in quantifying and summarizing the properties of the

64    samples and the distributions. Exploratory analysis is useful to guide the selection of the distribution

65    shape and summary statistics are required to characterize the sample or to judge whether the sample

66    is similar enough to some known distribution [Warner, 2008]. For instance, the location, scale,

67    skewness and kurtosis indicate respectively the centrality, dispersion, symmetry and peakedness of

68    the sample. Location and scale are summary statistics of the data whereas the shape of the data can

69    be captured by skewness and kurtosis [*Bickel and Lehmann*, 1975a; b; 1976; 1979]. On the other

70    hand, outliers, as gross errors and inconsistencies or unusual observations, can have negative

71    impacts on the selection of the appropriate distribution as well as on the estimation of the associated

72    parameters. In order to base the inference on the right data set, detection and treatment of outliers are

73    also important ([*Barnett and Lewis*, 1998] and [*Barnett*, 2004]). These concepts are well defined and

74    their computation is straightforward for univariate samples and distributions.

75    In classical multivariate analysis, several techniques were directly inspired by univariate techniques

76    and developed by analogy (multivariate normal distribution-based, component-wise and moment-

77    based). Techniques that analyse data in a component-wise manner perform badly when variables are

78 mutually dependent. Moment-based methods depend on the existence of moments. For a detailed

79 review of classical multivariate analysis techniques, the reader is referred to Anderson [1984] or

80 Schervish [1987].

81 Recently developed techniques avoid the above drawbacks by using the multivariate inward-outward

82 ranking of depth functions [*Zuo and Serfling*, 2000b]. Indeed, depth-based techniques are not

83 componentwise, and they are moment-free and affine invariant if the depth function is. These

84 advantages are useful to include distributions such as Cauchy, and also, the obtained results remain

85 the same after standardization. The depth-based ranking enables also numerous outlier detection

86 techniques, which are fundamental in FA. It is important to indicate that, unlike the univariate

87 setting, a multitude of definitions can be proposed for each sample characteristic (such as median

88 and symmetry) in the multivariate context. A key reference in the study of multivariate descriptive

89 statistics is Liu et al. [1999] where most of the above mentioned characteristic are treated. However,

90 each sample feature was subsequently studied separately by a number of authors. For instance, the

91 location was studied by Massé and Plante [2003], Zuo [2003] and Wilcox and Keselman [2004];

92 scale was treated by Li and Liu [2004], symmetry was the focus of Rousseeuw and Struyf [2004]

93 and Serfling [2006] and kurtosis was addressed by Wang and Serfling [2005]. These studies focused

94 mainly on inferential and asymptotical results. On the other hand, multivariate outlier detection, not

95 discussed in  Liu, et al. [1999], was studied recently by Dang and Serfling [2010].

96 The above features are of particular interest in hydrology since univariate data sets are generally

97 asymmetric [*Helsel et al.*, 2002] and the interest is on the tail of the distribution which is related to

98 kurtosis. In flood FA, Hosking and Wallis [1997] indicated that summary statistics, especially

99 skewness and kurtosis, are often used to judge the closeness of a sample to a target distribution.

100 Regarding outliers, in FA, we are concerned about two particular types of errors: the data may be

101     incorrect and/or the circumstances around the measurement may have changed over time [*Hosking*

102     *and Wallis*, 1997; *Rao and Hamed*, 2000].

103     The aim of the present study is to provide and to adapt recent statistical methods to the preliminary

104     analysis and exploration of multivariate hydrological data. The presented methods are mainly based

105     on the statistical notion of depth functions. Depth functions represent convenient tools for the

106     ranking of data in a multivariate context. Chebana and Ouarda [2008] presented a first application of

107     depth functions in the field of hydrology. Note that the multivariate L-moment approach represents

108     also an alternative that could be of interest for the development of multivariate descriptive statistics.

109     This approach is not treated in the present study and could be studied and compared to depth-based

110     approaches in future work. The reader is referred to Serfling and Xiao [2007] for the general

111     multivariate L-moment theory and to Chebana and Ouarda [2007] and Chebana et al. [2009] for

112     applications in hydrology.

113     The rest of the paper is organized as follows. In section 2, we present the general methodology for

114     exploratory descriptive analysis including graphical tools, measurements of location, dispersion,

115     symmetry, peakedness and outlyingness identification. We apply these concepts to real-world flood

116     data in section 3. Conclusions are presented in section 4. In the appendix, we present a brief

117     summary of the required background elements related to depth-functions.

## 2. Methodology

119     In this section, we present the general framework of the exploratory and descriptive multivariate

120     statistical tools. Let $X_1, X_2, ..., X_n \in R^d$ be a $d$-dimensional ($d \geq 1$) sample with size $n \geq d$. Using a

121     given depth function $D(.)$, we sort the sample in decreasing order of depth values to obtain

122     $X_{[1]}, X_{[2]}, ..., X_{[n]}$ and we define the "*de-class*" of $X_{[i]}$ as the set of observations with equal depth

123     values, for $i = 1,...,n$. A brief description of depth functions is given in the appendix. Note that, even

124     though, conceptually, any depth function can be used in the following visualisation and analysis

125     efforts, some combinations are not treated here because of the lack of their practical relevance and

126     since their properties are not well known. For instance, bagplots are generally based on Tukey depth

127     and are not studied using the Mahalanobis depth.

128     **2.1 Visualization**

129     Data should be visualized before any analysis can be conducted. In the 2 or 3 dimensional cases, the

130     simplest visualisation tool is the scatter plot. More useful, the bagplot is a generalization of the

131     univariate box-plot to the bivariate setting [*Rousseeuw et al.*, 1999] and is similar to the sunburst-

132     plot presented by Liu et al. [1999]. The bagplot is based on Tukey depth function whereas the

133     sunburst-plot uses either Tukey or Liu depths (given in expressions A1 and A2 respectively). The

134     bagplot is composed by a dark central bag which encircles the 50% deepest points. The Tukey

135     median, defined below in section 2.2, is indicated at the center and a light region delimited by the

136     points included in the central dark bag inflated by a factor 3 is also drawn and called the fence.

137     Points outside this region are considered as statistical outliers. We then link non-outlying points that

138     are outside the dark bag with the Tukey median. These lines have the same role as the whiskers in

139     univariate box plots [*Rousseeuw et al.*, 1999]. The bagplot generally gives indications concerning

140     the distribution of the sample, such as location, dispersion and shape. Note that the sunburst plot

141     presented by Liu et al. [1999] does not contain a fence region and hence the sunburst plot is not

142     considered as a tool to detect outliers. The points outside the fence are considered as extremes rather

143     than outliers. In the present study, we consider a more appropriate approach, given in Section 2.6, to

144     detect outliers.

145     Another way to visualise data can be obtained using the contours of the depth function. Contours can

146     reveal the shape and structure of multivariate data. Such plots enable direct comparisons of

147     geometry between bivariate data sets. The Tukey depth function is the most used and studied for

148     contour plots.

149     **2.2 Location parameters**

150     A location parameter indicates where most of the data are located. This notion is useful in hydrology

151     since it appears in almost all commonly employed probability distributions. In addition, the location

152     parameter is an important constituent in the index-flood model ([*Hosking and Wallis*, 1993] and

153     [*Chebana and Ouarda*, 2009]).The concept of location is closely related to the center-outward

154     ranking of depth functions. A point maximizing a depth function can be considered as a location

155     parameter, because of the property of "maximality at the center" of depth functions given in the

156     Appendix [*Zuo and Serfling*, 2000b]. In the following, we present several location parameters some

157     of which are well-known, such as the sample mean and the component-wise median.

158     ***Sample mean:*** The simplest and common location parameter is the arithmetic mean:

159 
$$\mu_n = \frac{1}{n}\sum_{i=1}^{n} X_i \tag{1}$$

160     $\mu_n$ is d-dimensional. It corresponds simply to the component-wise arithmetic means.

161     ***α-depth-trimmed-mean*** **:** For a coefficient $0 \le \alpha \le 1$, the α-depth-trimmed-mean ([*Liu et al.*, 1999]

162     and [*Massé*, 2009]) is a generalization of the sample mean (1). Given a depth function, the α-depth-

163     trimmed-mean can be considered as the sample mean computed from the $100(1-\alpha)\%$ deepest

164     points. Formally, we first define the $R^d$- valued function $\xi_n$ on $[0,1]$ as:

165 
$$\xi_n(t) = X_{[i]} \quad \text{if} \quad \frac{i-1}{n} < t \le \frac{i}{n} \quad \text{and} \quad \xi_n(0) = X_{[1]} \tag{2}$$

166     and $\overline{\overline{\xi}}_n(t)$ as the average over the *de*-class values in which $\xi_n(t)$ is contained. We then define the

167     *DL_n-statistic* as:

168     $$DL_n = \int_0^1 \overline{\overline{\xi}}_n(t)\omega(t)\,dt \qquad (3)$$

169     where $\omega(t)$ is a non negative weight function such that $\int_0^1 \omega(t)\,dt = 1.$ The α-depth-trimmed-mean is

170     defined according to a particular function $\omega_\alpha(t)$ given by:

171     $$\omega_\alpha(t) = 1/(1-\alpha) \quad \text{if } t \in [0, 1-\alpha] \quad \text{and} \quad \omega_\alpha(t) = 0 \text{ if } t \in (1-\alpha, 1] \qquad (4)$$

172     If α = 0, the α-depth-trimmed-mean is the classical mean given in (1). For α = 1, i.e. all observations

173     are trimmed, then $DL_n$ is defined as the deepest point of the sample. For $0 < \alpha < 1$, if $n\alpha$ is an

174     integer, then $DL_n$ is simply $DL_n = \dfrac{1}{n(1-\alpha)} \sum_{i=1}^{n(1-\alpha)} X_{[i]}$.

175     The use of the $\alpha\%$ trimmed-mean as a robust estimator in the univariate framework in hydrology

176     was discussed by Ouarda and Ashkar [1998].

177     ***Component-wise median:*** As a direct extension of the univariate median and similarly to the

178     multivariate arithmetic mean, the component-wise median $CM_n$ is defined as:

179     $$CM_n = \left( med\left( X_{1,1}, X_{2,1}, ..., X_{n,1} \right); ...; med\left( X_{1,d}, X_{2,d}, ..., X_{n,d} \right) \right)' \qquad (5)$$

180     where *med* is the usual univariate median. Note that $CM_n$ is not affine equivariant.

181     ***Depth medians:*** The following three location parameters are based on depth functions. The labels of

182     these medians are directly taken from their respective depth functions, Tukey, Oja and Liu, given in

183     the appendix. Any other depth function could also be used to define a location median but the above

184     are the most studied in the literature.

185    We consider the set $E \subseteq R^d$ of points that maximize the considered depth function. The depth

186    median is the centeroïd of the polygon composed by the set of points maximizing the selected depth

187    function [*Massé and Plante*, 2003]. Generally, $E$ is a convex and compact set [*Leon and Massé*,

188    1993]. Tukey and Oja medians have suitable properties whereas those of Liu median are not studied.

189    The set $E$ corresponding to Tukey median is convex since the half-space depth function is quasi-

190    concave [*Rousseeuw and Ruts*, 1999]. Tukey median is then defined as the center of mass of $E$

191    [*Massé and Plante*, 2003]. For the Oja median, if $n$ is even, then the set $E$ is a single element

192    according to Oja and Niinimaa [1985].

193    ***Spatial median:*** The spatial median is defined as [*Massé and Plante*, 2003]:

194    $$SpMed = \frac{1}{n} \arg \min_{x \in R^d} \sum_{i=1}^{n} \|x - X_i\|$$   (6)

195    where $\|.\|$ is the Euclidean norm and $\arg \min_{t \in A} \varphi(t)$ is the minimiser of the function $\varphi(.)$ over a set $A$.

196    In the bivariate case, a numerical study by Massé and Plante [2003] compared all the above

197    mentioned location estimators. The spatial median (6) stands as the best location parameter in terms

198    of robustness and accuracy, followed by Oja and Tukey medians. In a second group, we find the Liu

199    and the component-wise medians in terms of robustness. Trimmed means (for $\alpha = 0.05, 0.10$ with

200    Tukey and Liu depths) are in a third group, followed finally by the sample mean. Overall, medians

201    were shown to be more robust location parameters than means. Note that except for the Liu and Oja

202    medians, all the above location parameters are computable in higher dimensions, though sometimes

203    under approximations.

## 2.3 Scale parameters

205    Scale parameters are useful to measure the dispersion of a distribution or a sample. The scale and

206    location parameters appear in almost all probability distributions employed in hydrology since these

207  distributions should contain at least two parameters. We present two types of multivariate scale

208  parameters: matrix-valued and scalar-valued.

209  ***α-trimmed sample dispersion matrix:*** Given a center-outward ranking of data derived from a given

210  depth function, we first define a general weighted scale matrix [*Liu et al.*, 1999]. The corresponding

211  definition is similar to that of $DL_n$ (given in (3)), except that we replace $\xi_n(t)$ by the function $S_n(t)$

212  defined on the space $\mathrm{M}_d(R)$ of $d \times d$ real-valued matrices by:

213  $$\mathbf{S}_n(t) = \left(X_{[i]} - v_n\right)\left(X_{[i]} - v_n\right)' \quad \text{if } \frac{i-1}{n} < t \le \frac{i}{n}, \quad \text{and } \mathbf{S}_n(0) = \mathbf{0}_{d \times d} \tag{7}$$

214  where $v_n$ is the sample's deepest point and $\mathbf{0}_{d \times d}$ is the $d \times d$ matrix with null elements. The *weighted*

215  *scale matrix* is defined by:

216  $$DS_n = \int_0^1 \overline{\mathbf{S}}_n(t)\,\omega(t)\,dt \tag{8}$$

217  where $\overline{\mathbf{S}}_\mathbf{n}$ indicates the average of $\mathbf{S}_\mathbf{n}$ over all *de*-classes to which $X_{[i]}$ belongs and $\omega$ is the weight

218  function as defined for the *α*-trimmed mean. The *α-trimmed sample dispersion matrix* is a particular

219  case of $DS_n$, with $\omega$ defined as in (4). Given $0 \le \alpha < 1$, if $n\alpha$ is an integer, the *α-trimmed-dispersion*

220  *matrix* is given by:

221  $$DS_n = \frac{1}{n(1-\alpha)} \sum_{i=1}^{n(1-\alpha)} \overline{\mathbf{S}}_n\left(\frac{i}{n}\right) \tag{9}$$

222  For $\alpha = 1$, we define $DS_n$ as the zeros matrix and for $\alpha = 0$, it coincides with the usual covariance

223  matrix.

224  Note that the matrix form of scale enables an easy comparison of dispersion between dimensions

225  and can reveal more information. However, a matrix is not effective for measuring the overall

226      dispersion of the distribution (e.g. [*Liu et al.*, 1999]). We can overcome this problem by taking a

227      norm of the scale matrix (9). Some known matrix norms can be found in Manly [2005].

228      ***Scalar form of scale:*** Scalar values can be seen as an information reduction regarding scales. Hence,

229      it is more appropriate to plot these values as a curve with respect to a given coefficient. We

230      introduce a graphical tool presented in Liu et al. [1999] that measures the dispersion of a

231      multivariate sample. Given a depth function, the function $Sc_n(p)$, $0 \le p \le 1$, returns the volume of

232      the central region $C_{n,p}$ composed of the $\lceil np \rceil$ deepest points, where $\lceil a \rceil$ is the smallest integer

233      larger or equal to $a$.

234      The plot of the function $Sc_n(p)$ with respect to $p$ is an evaluation of the expansion of $C_{n,p}$ with

235      respect to $p$. This kind of scale curves is a simple one-dimensional curve describing the scale. It

236      allows also to quantify the evolution of a sample. The curve $Sc_n(.)$ is interpreted as follows: "if the

237      scale curve of a distribution $G$ is consistently above the scale curve of another distribution $F$, then $G$

238      has a larger scale than $F$".

## 239      2.4 Skewness

240      Skewness can be defined as a measure of departure from symmetry. Skewness evaluation is

241      important in hydrology since generally univariate distributions are not symmetric and are one-side

242      heavily-tailed [*Helsel et al.*, 2002]. In the multivariate case, there are several types of symmetry,

243      such as: spherical, elliptical, antipodal and angular. Depth-based tools are presented in this section to

244      empirically evaluate each type of symmetry. In the following, we present the definition of each

245      symmetry as well as how it can be evaluated. The definitions are taken from Liu et al. [1999] and

246      Serfling [2006]. All the following types of symmetry have a common feature: the distribution of a

247    centered random vector $X - c$ is invariant under a given transformation and all of them reduce to the

248    usual univariate symmetry.

249    ***Spherical symmetry:*** "The distribution of the random variable $X$ is said to be spherically symmetric

250    about the point $c$ if the distributions of $(X \text{-} c)$ and $\mathbf{U}(X \text{ -} c)$ are identical, for any orthonormal

251    matrix $\mathbf{U}$." Recall that a matrix $\mathbf{U}$ is orthonormal if and only if $\mathbf{UU'} = \mathbf{U'U} = I$ where $\mathbf{U'}$ is the

252    transpose of the matrix $\mathbf{U}$ and $I$ is the identity matrix. This kind of symmetry represents a rotation of

253    $X$ about c. The probability density function of $X$, when it exists, is then of the form

254    $g\!\left((x-c)'(x-c)\right)$ for a nonnegative real-valued function $g$. Examples of this kind of distributions

255    include the multivariate versions of the standard normal, the $t$ and the logistic distributions.

256    To evaluate the spherical symmetry, we consider, for a given depth function, the smallest enclosing

257    $d$-sphere that contains the $\lceil np \rceil$ deepest points for $p \in [0,1]$. We denote $Sph(p)$ the proportion of

258    sample points falling in this sphere. The function $Sph(p)$ is increasing and $p \le Sph(p) \le 1$. The area

259    $\Delta_n$ between the curve $y = Sph(x)$ and the diagonal line $y = x$ is an indicator of spherical skewness. A

260    perfectly spherical symmetric sample would imply that the curve $Sph(.)$ is close to the diagonal (i.e.

261    $Sph(p) \approx p$ ) and hence $\Delta_n$ is close to zero.

262    ***Elliptical symmetry:*** "The distribution of the random variable $X$ is said to be elliptically symmetric

263    about a certain point $c$ if there exists a non singular matrix $\mathbf{V}$ such that $\mathbf{V}X$ is spherically symmetric

264    about  c."  The  corresponding  probability  density  function  of  $X$  is  of  the  form

265    $|V|^{-1/2} g\!\left((x-c)'V^{-1}(x-c)\right)$ which includes, for instance, the multivariate normal distribution with a

266    covariance matrix $\Sigma = V'V$. The corresponding contours of the probability density function are

267    indeed of elliptical shape.

13

268   To empirically evaluate elliptical skewness, it is suggested in Liu et al. [1999] that we first

269   standardize data using the scale matrix (see Section 2.3) of the $\lceil np \rceil$ deepest points for $p \in [0,1]$. We

270   then proceed on the basis of the transformed data as in spherical symmetry by evaluating $Sph(p)$ on

271   the transformed set. We finally plot the function $Sph(p)$ on the transformed data set. The

272   interpretation of the curves associated to the elliptical skewness is similar to that of the spherical

273   skewness.

274   ***Antipodal symmetry:*** "The distribution of the random variable $X$ is said to be antipodally symmetric

275   about the point $c$ (if such a point exists) if the distributions of $(X-c)$ and $-(X-c)$ are identical."

276   This symmetry is also called reflective or diagonal and represents the most direct extension of the

277   usual univariate symmetry. The probability density function $f$ in this case is such that

278   $f(x-c) = f(c-x)$.

279   Given a depth function and a location parameter $\mu$, we consider the reflection of the $p^{\text{th}}$ central

280   region $C_{n,p}$ about $\mu$, for $p$ in (0, 1). We denote $Ca(p)$ the proportion of the $\lceil np \rceil$ deepest points

281   falling in the intersection of $C_{n,p}$ and its reflection. By definition we have $0 \le Ca(p) \le \lceil np \rceil / n$.

282   An antipodal symmetric sample would suggest that $Ca(p) = \lceil np \rceil / n \approx p$. Thus, we can measure

283   antipodal skewness by evaluating the area between the diagonal line $y = x$ and the curve $y = Ca(x)$,

284   for $x \in [0,1]$. A larger area corresponds to a larger deviation from antipodal symmetry.

285   ***Angular symmetry:*** "The distribution of the random variable $X$ is said to be angularly symmetric

286   about the point $c$ if, conditional on $X \neq c$, the distributions of $(X-c)/\|(X-c)\|$ and

287   $-(X-c)/\|(X-c)\|$ are identical." One of the features of this symmetry is that if $c$ is a point of

288   angular symmetry, then any hyper-plane passing through $c$ divides the whole space $R^d$ into two half-

289    spaces with probability 0.5 (if the distribution is continuous). More characterizations of angular

290    symmetry can be found in [*Zuo and Serfling*, 2000a].

291    To measure angular symmetry of a given sample, we first identify the deepest point $v_n$ according to

292    a given depth function. Then, we evaluate the Tukey depth of the deepest point $v_n$ with respect to the

293    restricted data in the $p^{\text{th}}$ central region $C_{n,p}$ for each $p \in [0,1]$. The deviation of the obtained curve,

294    denoted $h(p)$, from the $x$ axes measures the degree of the antipodal symmetry. The value of Tukey

295    depth of the deepest point should be 0.5 under angular symmetry. The interpretation of the obtained

296    values and curves follows from: " […] the deviation of the half-space depth at the deepest point

297    from the value 0.5 is a measure of the departure from angular symmetry of the empirical distribution

298    determined by the sample points within each level set" [*Liu et al.*, 1999].

299    Liu et al. [1999] suggested to consider only the part of the curve with $p$ larger than 0.4 where the

300    curve stabilizes. Note that for small values of $p$, the curve is based on a small fraction of the data

301    which is not enough for the convergence of the Tukey depth function.

302    The reader may have noted that these concepts of symmetry are linked together. They can be ranked

303    from more to less restrictive:

304
$$\begin{array}{ccccccc} \text{Spherical} & & \text{Elliptical} & & \text{Antipodal} & & \text{Angular} \\ \text{symmetry} & \Rightarrow & \text{symmetry} & \Rightarrow & \text{symmetry} & \Rightarrow & \text{symmetry} \end{array} \qquad (10)$$

305    In all kinds of symmetry, a point $c$ is required. This point is generally a location parameter. Zuo and

306    Serfling [2000a] studied the performance of some location measures associated to multivariate

307    symmetry.

308    After having defined and evaluated skewness, it is important to conduct hypothesis testing for

309    symmetry. This represents a current topic of research in the multivariate setting (see for instance

310   Manzotti et al. [2002], Huffer and Park [2007], Sakhanenko [2008] and Ngatchou-Wandji [2009]).

311   The topic of hypothesis testing is beyond the scope of the present study.

## 2.5 Kurtosis

313   Peakedness and tailweight evaluations are important in hydrology as the focus is often on extreme

314   events and the tail of the distribution. These concepts are related to kurtosis which is a measure of

315   the overall spread relative to the spread in the tails. Measuring kurtosis is important in water

316   sciences since extreme events occur in the tail of the distribution (univariate or multivariate) with

317   non negligible probability. Kurtosis is generally defined as a ratio of two scale measures, i.e. scale of

318   the whole data and scale of the central part [*Bickel and Lehmann*, 1979]. We present in this section a

319   number of tools that quantify multivariate kurtosis. The reader is referred to Liu et al. [1999] and

320   Wang and Serfling [2005] for more details.

321   ***Lorenz curve of Mahalanobis distance:*** Given a non-singular scale matrix $\mathbf{S}_n$, such as the one given

322   in (8) or simply the covariance matrix, and a given depth function for which $v_n$ is the deepest point,

323   we introduce the real-valued functions:

324
$$L\left(p\right)=\frac{\sum_{i=1}^{\lceil np \rceil} Z_i}{\sum_{i=1}^{n} Z_i} \quad \text{and } L^*\left(p\right)=\frac{\sum_{i=1}^{\lceil np \rceil} Z_i \big/ \lceil np \rceil}{\sum_{i=1}^{n} Z_i \big/ n} \quad \text{for } 0 < p \leq 1 \qquad (11)$$

325   where

326
$$Z_i = \left(X_{[i]} - v_n\right)' \mathbf{S}_n^{-1} \left(X_{[i]} - v_n\right), \text{ for } i = 1, 2, ..., n \qquad (12)$$

327   We define $L(0) = L^*(0) = 0$ and we have $L(1) = L^*(1) = 1$. Note that $L^*$ is simply an adjusted

328   formulation of $L$ and each of them represents a ratio of the central variability to the total variability.

329   The functions given in (11) are then plotted and the area corresponding to the surface between the

330   curves $y = L(x)$ or $y = L^*(x)$ and the diagonal line $y = x$ is evaluated. Both areas can be

331    interpreted in the same way: a large area corresponds to a high degree of peakedness and tailweight,

332    and inversely a small area corresponds to heavy shoulders. The curves $L^*$ and $L$ have the same

333    interpretation, but the area computed from $L^*$ should be more pronounced than the one computed

334    from $L$. Consequently, sample curves can be compared more effectively using $L^*$ than $L$.

335    **Shrinkage plots:** They are based on the shrinkage of the boundary of the $p$th central region

336    $C_{n,p}$ towards its center by a given fixed coefficient $s$, $0 < s < 1$ leading to region $C_{n,p}^s$. We then plot

337    the function $a_s(p)$ of the fraction of observations in $C_{n,p}^s$ for fixed $s$. Liu et al. [1999] indicated that

338    one value of $s$ is enough to conclude and they proposed $s = 0.5$. For a fixed $s$, heavier tails

339    correspond to higher values of $a_s(p)$ especially for large $p$.

340    **Fan plots:** A fan plot is a collection of curves used to evaluate kurtosis. It consists in an arbitrary

341    number of curves, each of which is associated with a value $p \in [0,1]$. For a given $p$, we consider the

342    sub-sample $Sam(p)$ formed by the $\lceil np \rceil$ deepest points (in the central region $C_{n,p}$). For $t \in [0,1]$, we

343    denote $C_n(p,t)$ the area of the $t^{th}$ convex hull of $Sam(p)$ composed by $100t$ % of the deepest

344    observations. We define the function $b_p(t)$ for $t \in [0,1]$ by:

345 
$$b_p(t) = \frac{volume\left[C_n(p,t)\right]}{volume\left[C_n(p,1)\right]} \text{ if } C_n(p,1) \neq 0 \text{ and } b_p(t) = 0 \text{ otherwise} \qquad (13)$$

346    Intuitively, a fan plot may be regarded as a comparison of areas between central (corresponding to

347    low values of $p$), shoulder (corresponding to middle values of $p$) and tail regions (corresponding to

348    high values of $p$). A more spread out fan plot indicates that the corresponding distribution is heavy

349    tailed since $b_p(t)$ becomes smaller. This way to measure kurtosis requires a large amount of data

350    since the data size is reduced in two stages (with $p$ and then with $t$).

351 ***Quantile-based measure:*** This measure is based on the function $k_C(.)$ proposed by Wang and

352 Serfling [2005] and expressed as :

353 $$k_C(r) = \frac{V_C\left(\frac{1}{2}-\frac{r}{2}\right)+V_C\left(\frac{1}{2}+\frac{r}{2}\right)-2V_C\left(\frac{1}{2}\right)}{V_C\left(\frac{1}{2}+\frac{r}{2}\right)-V_C\left(\frac{1}{2}-\frac{r}{2}\right)} \quad \text{for} \quad 0 < r \le 1 \text{ and } k_C(0) = 0 \qquad (14)$$

354 where the function $V_C(r)$ is the volume of a central set $C(r)$. The set $C(r)$ is defined as the inner

355 set, with probability $r$, delimited by contours of a given depth function. Wang and Serfling [2005]

356 used Tukey depth function and indicated that any affine invariant depth function can be used as well.

357 Note that the set $C(r)$ is general with a special case defined on the basis of spatial quantiles. The

358 measure $k_C(r)$ represents the difference of the volumes of two regions $A$ and $B$ divided by the sum of

359 their volumes where $A = C(1/2) - C(1/2 - r/2)$ and $B = C(1/2 + r/2) - C(1/2)$. Note that the

360 boundary associated to the region $C(1/2)$ represents the "shoulders" of the distribution and it

361 separates the "central part" from the corresponding "tail part".

362 Wang and Serfling [2005] provided indications for the interpretation of the curve $k_C(.)$. They

363 indicated that if the attention is confined to a class of distributions for which either $F$ is unimodal,

364 $F$ is uniform, or $1-F$ is unimodal, then, for any fixed $r$, a value of $k_C(r)$ near +1 suggests a

365 peakedness, a value near -1 suggests a bowl-shaped distribution, and a value near 0 suggests

366 uniformity.

367 Increasing values of $k_C(.)$ indicate that the probability mass is greater in the center than in the tails.

368 It is important to mention that, unlike kurtosis measures discussed in the above sub-sections, the

369 quantile-based measure requires some prior knowledge about the distribution of the sample to

370 interpret the obtained curves.

371 **2.6 Outlier detection**

18

372   Identifying outliers is an important statistical step to analyze data sets as indicated, for instance, by

373   Barnett and Lewis [1978] in the univariate as well as in the multivariate settings. Outlier detection in

374   hydrologic data is a common problem which has received considerable attention in the univariate

375   framework.

376   In the multivariate setting, outlyingness functions are defined and employed to detect outliers.

377   Values of these functions usually range in the interval [0, 1]. They measure outlyingness of a certain

378   point with respect to the entire sample. An outlyingness value near 1 indicates high outlyingness,

379   and inversely a value near 0 indicates centrality. In order to determine whether an observation is an

380   outlier or not, it is required to define a threshold, i.e. the minimum outlyingness value from which a

381   datum is considered to be an outlier. In the following we present the most promising and recently

382   developed outlying functions, based on depth functions and given in Dang and Serfling [2010].

383   ***Outlyingness:*** A depth outlyingness is a transformation of a depth function for a given distribution $F$

384   and $x \in R^d$. The followings are studied in Dang and Serfling [2010]:

385         Half-space: $O_{HD}(x,F) = 1 - 2HD(x,F)$ $\hspace{3cm}$ (15)

386         Mahalanobis: $O_{MD}(x,F) = d^2_{A(F)}(x,\mu(F)) / \left[ 1 + d^2_{A(F)}(x,\mu(F)) \right]$ $\hspace{1cm}$ (16)

387         Projection: $O_{PD}(x,F) = PD(x,F) / \left[ 1 + PD(x,F) \right]$ $\hspace{2cm}$ (17)

388   where $HD(.,F)$, $d^2_{A(F)}(.,\mu(F))$ and $PD(.,F)$ are given respectively in (A1), (A3) and (A6) and

389   $\mu(F)$ is a location measure and $A(F)$ is a nonsingular matrix scale measure;

390         Spatial: $O_S(x,F) = \left\| E\left( Sign(x-X) \right) \right\|$ $\hspace{3cm}$ (18)

391         Spatial Mahalanobis: $O_{SM}(x,F) = \left\| E\left[ Sign\left( \mathbf{C}^{-1/2}(x-X) \right) \right] \right\|$ $\hspace{1.5cm}$ (19)

19

392  where $\|.\|$ is the Euclidean norm, $X$ is $F$-distributed and $Sign(.)$ is the multidimensional sign function

393  given by:

394
$$Sign(x) = x/\|x\| \quad \text{if } x \neq 0 \quad \text{and} \quad Sign(0) = 0 \qquad (20)$$

395  and $\mathbf{C}$ is any affine invariant symmetric positive definite $d \times d$ matrix. The matrix $\mathbf{C}$ could be the

396  classical covariance matrix or the matrix obtained as the minimum covariance determinant

397  [*Rousseeuw and Van Driessen*, 1999].

398  ***Threshold:*** Selection of the appropriate threshold is an important step in outlier detection. It is

399  related to false positive and true positive rates. The arbitrary *false positive rate*, denoted $\alpha_n$, is the

400  proportion of non-outliers misidentified as outliers. This constant is closely related to the *true*

401  *positive rate* $\varepsilon_n$, which represents the real theoretical proportion of outliers (called also

402  contaminants). Ideally, $\alpha_n$ has to be small compared to $\varepsilon_n$. Dang and Serfling [2010] fixed a ratio of

403  false outliers $\delta = \alpha_n / \varepsilon_n$ and then used an additional coefficient $\beta = \varepsilon_n \sqrt{n}$, to define a threshold as

404  the ($1 - \alpha_n$)-quantile of the outlyingness values :

405
$$\lambda_n = F_{O(X,F)}^{-1}\left(1 - \alpha_n\right) = F_{O(X,F)}^{-1}\left(1 - \delta\varepsilon_n\right) = F_{O(X,F)}^{-1}\left(1 - \beta\delta/\sqrt{n}\right) \qquad (21)$$

406  The following example is illustrated in Dang and Serfling [2010]. By putting $\delta = 0.1$, the ratio of

407  false outliers is about 10% among the allowed ones. Assume that we allowed for $n\varepsilon_n = 15$ true

408  outliers, the constant $\beta$ takes the value $\beta = n\varepsilon_n / \sqrt{n} = 15/\sqrt{100} = 1.5$ for $n = 100$. Hence,

409  $\lambda_n = F_{O(X,F)}^{-1}\left(1 - 0.15/\sqrt{n}\right) = F_{O(X,F)}^{-1}\left(0.985\right)$ for $n = 100$ corresponds to the 0.985-quantile of the

410  outlyingness values.

411    These thresholds are given explicitly for the multivariate normal distribution. Since, these thresholds

412    are not available in general, those of the normal distribution can be employed as approximations. For

413    a multivariate sample from a multivariate normal distribution $\Phi$, the theoretical threshold $\lambda_n$ is

414    given explicitly for the Malahanobis, half-space and projection outlyingness functions of Dang and

415    Serfling [2010]. Supposing that $\beta\delta/\sqrt{n} \in [0,1]$ and that the variable $X$ is standard normally

416    distributed, then the threshold given in (21) is given more explicitly as:

$$
\begin{aligned}
\lambda_n &= \frac{T\left(d,\alpha_n\right)}{1+T\left(d,\alpha_n\right)} && \text{for the Mahalanobis outlyingness} \\
\lambda_n &= 2\Phi\left(T\left(d,\alpha_n\right)\right)-1 && \text{for the halfspace outlyingness} && (22) \\
\lambda_n &= \frac{T\left(d,\alpha_n\right)}{\Phi^{-1}\left(3/4\right)+T\left(d,\alpha_n\right)} && \text{for the projection outlyingness}
\end{aligned}
$$

417

418    where $T(d,\alpha)=\sqrt{\left(\chi_d^2\right)^{-1}\left(1-\alpha\right)}$ with $\left(\chi_d^2\right)^{-1}$ is the inverse cumulative distribution function of the

419    chi-square distribution with $d$ degrees of freedom.

420    For the spatial and spatial Mahalanobis, normal thresholds are not available. Note that it is also

421    convenient to define thresholds for each outlyingness function on the basis of the empirical quantile

422    of the outlyingness values.

423    **3. Applications**

424    In the following, the notions and methods introduced in Section 2 are applied to two real-world

425    hydrological data sets. The first one is given in details whereas in the second one, we focus on

426    outlier detection. All the methods presented in Section 2 are implemented in the Matlab environment

427    [*MathWorks*, 2008] for the bivariate setting. Few methods, such as those based on the Mahalanobis

428    distance, can be applied to higher dimensions.

429

430     **3.1 Ashuapmushuan case study**

431     The data set used in this case study is taken from Yue et al. [1999] and concerns floods in the

432     Ashuapmushuan basin located in the province of Québec, Canada. The flood annual observations of

433     flood peaks ($Q$), durations ($D$) and volumes ($V$) were extracted from a daily streamflow data set

434     from 1963 to 1995. The gauging station, with identification number 061901, is near the outlet of the

435     basin, at latitude 48.69°N and longitude 72.49°W. In this region floods are caused by high spring-

436     snowmelt.

437     To allow comparisons, we considered the study of all three combination series ($Q$, $V$), ($D$, $V$) and

438     ($Q$, $D$) by all presented methods. In all parts of the analysis, except for outlier detection, we

439     considered four depth functions: Tukey, Oja, Mahalanobis and Liu which are given respectively in

440     (A1), (A2), (A4) and (A5). Note that results were produced for the three bivariate series using the

441     four depth functions. However, the four depth functions lead to practically identical results for each

442     series. Therefore, in the following we only present results based on the Tukey depth function. A

443     sample's depth values are essential for the analysis since almost all tools presented above are depth-

444     based. The corresponding depth values for the series ($Q$, $V$) are given in Table 1 as a selected

445     example.

446     ***Displaying data***

447     Bagplots and contour plots, based on Tukey depth, are presented in Figures 1a,b respectively. The

448     Tukey depth function is the most used for bagplots and contour plots. We observe the orientation of

449     the bags which indicates the positive correlation between $Q$ and $V$. We also observe that more data

450     are concentrated in the center and that the extreme observations, with high $V$ and relatively small $Q$,

451     are located outside the fence of the ($Q$, $V$) plot. All three series are unimodal, both ($Q$, $V$) and ($D$, $V$)

452     are positive dependent whereas ($Q$, $D$) shows no clear dependence. This is in agreement with the

453    multivariate flood FA literature (e.g. [*Yue et al.*, 1999] and [*Zhang and Singh*, 2006]). The series

454    (*Q, V*) seems more concentrated and tight than the other two. The contours of (*Q, D*) are more

455    circular and more distant compared to those of the (*Q, V*) and (*D, V*) series. Note that the points

456    outside the fence of (*Q, V*) and (*D, V*) in Figure 1b (left and middle) correspond to the floods of 1994

457    and 1974 respectively. They have the smallest depth values. As indicated previously, they cannot be

458    considered as outliers at this stage of the analysis but can be seen as extremes.

459    ***Location parameters***

460    All location parameters presented in Section 2.2 are obtained in the bivariate setting. Location

461    parameters are indicated in Figure 2, both within the scatter plot and separately in a zoomed plot.

462    The corresponding values are given in Table 2. Generally, all location parameters are located in the

463    center of the sample. We observe that locations based on the mean are slightly influenced by the

464    extreme values of the sample, for instance, in the series (*Q, D*). This result is in agreement with the

465    study by Massé and Plante [2003] where the authors recommend, on the basis of accuracy and

466    robustness, the use of spatial median followed by Oja and Tukey medians.

467    ***Scale parameters***

468    The α-trimmed dispersion matrix, given in (8), is easily computed for any multivariate setting.

469    Corresponding values associated to each series are presented in Table 3 for α = 0.00, 0.05 and 0.10.

470    For a given series, all matrices are in the same order of magnitude with a slight decrease with respect

471    to α. Values in the matrices corresponding to (*Q, V*) are larger than those of (*D, V*) and the smallest

472    are those of (*Q, D*). All values in the dispersion matrices are positive except for those representing

473    the covariance between *Q* and *D*. This was already indicated when displaying data and is again in

474    agreement with the hydrological literature (e.g. [*Yue et al.*, 1999] and [*Zhang and Singh*, 2006]).

475    In addition, Figure 3 presents, for each series, the function $Sc_n(p)$ with respect to $p$ of the volume

476    of the $p$th central region $C_{n,p}$. We observe that $(Q, V)$ is more dispersed than both $(D, V)$ and $(Q, D)$

477    since $Sc_n(p)$ corresponding to $(Q, V)$ is larger for any fixed $p$. This can be partially explained by

478    comparing the magnitudes of volumes ($\approx 10^4$), flood peaks ($\approx 10^3$) and durations ($\approx 10^1$). Moreover,

479    the variances of the marginal variables differ greatly: the variance of $V$ ($\sigma^2 = 1.55\text{e}+008$) is larger

480    than the variance of $Q$ ($\sigma^2 = 1.29\text{e}+005$) and the variance of $D$ ($\sigma^2 = 211.30$). The variability induced

481    by $D$ is included in both $Q$ and $V$ because they are evaluated on $D$. This is in concordance with

482    matrix dispersion given in Table 3. These findings, both with matrices and scalars, confirm what was

483    previously revealed from bagplots and contour plots in Figure 1.

484    ***Skewness measures***

485    The measures of the four kinds of symmetry, presented in Section 2.4, are applied on each one of the

486    three series. Figure 4 illustrates the curves of the four skewness measures. We notice that the

487    $(D, V)$ sample is the closest to spherical symmetry with a small volume $\Delta_n = 0.09$ (Figure 4a).

488    Results from Figure 4b suggest that the $(Q, V)$, $(D, V)$ and $(Q, D)$ distributions are likely to be

489    elliptically symmetric, since $Sph_n(p)$ is very close to the diagonal with a very small $\Delta_n$. This can be

490    confirmed with the bagplots and contour plots of Figures 1a and 1b respectively. Regarding

491    antipodal skewness, Figure 4c shows that all the considered series seem to be symmetric since the

492    obtained curves are similar to those in Figure 14 in Liu et al. [1999] and are already elliptically

493    symmetric. Among the three series, $(Q, D)$ is the closest to angular symmetry since the function

494    $h(p)$ converges to 0.5 for $p$ larger than 0.4 (Figure 4d). Hence, the three series seem to be

495    elliptically symmetric. Note that the procedure treats the whole distribution including copula and

496    margins. The univariate skewness coefficient values are 0.978, 0.522 and 0.286 respectively for $D$, $V$

497  and $Q$. Since these values are significantly non null, the corresponding marginal distributions are

498  positively skewed. In contexts similar to the present one, the so-called meta-elliptical distributions

499  could be a reasonable model to consider. In the statistical literature, meta-elliptical copulas are

500  studied by Abdous et al. [2005] and applied in hydrology by Wang et al. [2010]. Meta-elliptical

501  distributions allow margin variables to follow different distributions. It is advisable to check the

502  significance of this symmetry by using statistical tests given in the references provided in Section

503  2.4. These findings are useful to guide the selection of the appropriate distribution for further

504  analysis.

505  ***Kurtosis parameters***

506  For all three series $(Q, V)$, $(D, V)$ and $(Q, D)$, the curves to evaluate kurtosis are presented in Figure

507  5. The functions $L$ and $L^*$ defined in (11) are presented in Figures 5a,b respectively. Clearly, as

508  expected, $L^*$ is more distinctive than $L$. Hence, the series $(Q, V)$ represents the most peaked sample,

509  followed by $(Q, D)$ and then by $(D, V)$ according to $L^*$.

510  Shrinkage plots, in Figure 5c, are very similar and do not allow to compare the various series, apart

511  that all the three series are heavy-tailed. However, fan plots indicate again that $(Q, V)$ is the most

512  peaked series (Figure 5d). As explained in Section 2.5, quantile-based curves, provided in Figure 5e,

513  do not reveal indications concerning kurtosis for the studied series since they require some

514  information regarding the generating distribution.

515  Overall, we conclude that $(Q, V)$ is the most peaked series and that the $L^*$-based kurtosis measure

516  seems to be the best option since it is simple, distribution-free and able to distinguish between

517  kurtosis of distributions. Therefore, the appropriate distribution candidates should be heavy-tailed as

518  expected.

519

520     *Outlier detection*

521     We evaluated spatial and both depth-based Mahalanobis and Tukey outlyingness functions for the

522     three series. The results are presented in Table 4. The corresponding thresholds are obtained by

523     selecting the values discussed in Section 2.6: that is the ratio of false outliers $\delta = 0.1$, the true

524     number of outliers $n\varepsilon_n = 5$ corresponding to approximately 15% of the sample, and the constant $\beta$

525     $\beta = n\varepsilon_n / \sqrt{n} = 0.8704$ for $n = 33$. Hence, from expression (21), $\lambda_n = F_{o(X,F)}^{-1}(0.985)$ which

526     corresponds to the 0.985-quantile of the outlyingness values.

527     Table 4 illustrates the normal and empirical thresholds for each series and each outlyingness

528     function as well as the corresponding detected outliers as years. The results show that there is no

529     outlier for the three series on the basis of the empirical thresholds using the three kinds of

530     outlyingness. However, the normal thresholds are not convenient in the present case. They lead to

531     very small thresholds for Mahalanobis and very high thresholds for Tukey. The reason could be the

532     short sample size of the series which does not allow for appropriate approximations. Furthermore, it

533     is well documented that flood series are not normally distributed. Note that the $(Q, V)$ and $(D, V)$ of

534     the years 1974 and 1994 are not detected as outliers even by relaxing the coefficients $\delta$ and $n\varepsilon_n$.

535     **3.2 Magpie case study**

536     The data series related to the second case study consists in daily natural streamflow measurements

537     from the Magpie station (reference number 073503). This station is located at the discharge of the

538     Magpie Lake in the Côte-Nord region in the province of Québec, Canada. Data are available from

539     1979 to 2004. In this case study we focus on outlier detection for the flood peak $Q$ and the flood

540     volume $V$ series. The corresponding Tukey depth and the outlyingness values are reported in Table

541     5.

542   To obtain the threshold that the outlyingness of an outlier exceeds, we considered $\delta = 0.15$ as the

543   ratio of false outliers and $n\varepsilon_n = 5$ as the number of true outliers. Therefore, from expression (21),

544   the threshold corresponds to the empirical 97%-quantile of the outlyingness values. Numerically, the

545   obtained thresholds are respectively 0.9231, 0.8676 and 0.9462 for $O_{HD}$, $O_{MD}$ and $O_S$. Consequently,

546   the flood of 1981 is detected by all the measures as outlier, whereas 1987 is detected only by $O_{HD}$

547   and has the second highest outlying value by both $O_{MD}$ and $O_S$. The measure $O_{HD}$ detects several

548   other outliers, such as 1999 and 2002, with the same outlyingness value (equal to the threshold).

549   However, if a quantile of order higher than 97% is considered, by modifying the parameters related

550   to the threshold, then $O_{HD}$ will not detect any outliers. Note that according to Dang and Serfling

551   [2010], the $O_{HD}$ measure is not recommended.

552   To explain these outliers, hydrological characteristics were derived and the corresponding

553   meteorological data were examined. These data were extracted from Environment Canada's Web

554   site (www.climat.meteo.gc.ca/climateData/canada\_f.html). The hydrograph of the year 1981 is

555   characterized by very high $V$ and $Q$ whereas 1987 seems to correspond to a dry year since the flow

556   was the lowest during the spring season and has the lowest $V$ and $Q$ values in the series. For 1981

557   there was an important amount of snow in early winter (October to January) followed by thaw and

558   rain during February-March. In comparison to the previous and following years, 1987 was

559   characterised by a warm end of winter and a very cold and less rainy fall. Hence, snow melted

560   earlier compared to other years. The flood of 1999 is characterised by a high $V$, although lower than

561   the one corresponding to 1981. The year 1999 was characterised by an important quantity of snow

562   on the ground with high temperatures in March. The observed hydrograph of 2002 contains two

563   peaks: the first one is characterised by a high magnitude while the second one is smaller and occurs

564   later in the summer. This year was particular with a very cold winter and a large amount of snow on

565   the ground until early May. In conclusion, the flows of the above detected years seem unusual but

566   are actually observed and do not correspond to incorrect measurements or changes over time in the

567   circumstances under which the data were collected. Hence, these observations should be kept and

568   employed for further analysis. However, it is recommended to use robust statistical methods to avoid

569   sensitivity of the obtained results to outliers.

570   The Tukey median and the arithmetic mean are evaluated. We observe that the median corresponds

571   to the year 1980 with $Q = 847.72$ and $V = 2216.22$. The bivariate mean vector is ($Q = 859.15$, $V =$

572   2138.70). After removing any of the above outliers, the mean changes significantly whereas the

573   median remains the same. For instance, the mean becomes (835.25, 2067.89) after removing the

574   1981 outlier. This result illustrates the effect of the detected outliers on the mean which is not the

575   case for the median. Since the detected outliers represent actual observations, it is not advised to

576   remove them. In that case, the median is recommended as a location measure. For further analysis,

577   robust methods and measures are recommended for this data set.

578   **4. Conclusions**

579   The techniques and methods presented in the present paper constitute the first step in a multivariate

580   frequency analysis. In the present paper, several features of the sample are treated, such as location,

581   scale, skewness, kurtosis and outlier detection. The methods discussed in the present paper are

582   superior to the classical multivariate methods based on moments, the assumption of normality, and

583   componentwise techniques. These recent methods, mainly based on the notion of depth function, are

584   moment-free, not normally-based and affine invariant (if the depth function is). This preliminary

585   step of the analysis is useful for the modeling of hydrological variables and for risk evaluation. It

586   allows to screen the data, to guide the selection of the appropriate model and to make comparisons

587   of multivariate samples. The methods discussed in the present paper were applied to flood data from

588    the Ashuapmushuan and Magpie data sets in the province of Québec, Canada. These methods can

589    also be adapted and applied to other hydrometeorelogical variables such as storms, heat waves and

590    draughts.

591    The findings related to the first case study of the Ashuapmushuan basin show that there are no

592    outliers and the data are likely to be elliptically symmetric and heavy–tailed. Therefore, the

593    appropriate multivariate distribution should be in a class with similar features. The second case

594    study of the Magpie station contains a number of outliers which are checked to be real observed

595    data. Therefore, they cannot be removed from the sample and robust methods should be adopted for

596    further analysis.

## Acknowledgments

602

# Bibliography

Abdous, B., C. Genest, and B. Rémillard (2005), Dependence Properties of Meta-Elliptical Distributions, in *Statistical Modeling and Analysis for Complex Data Problems*, edited by P. Duchesne and B. RÉMillard, pp. 1-15, Springer US.

Aloupis, G., C. Cortés, F. Gómez, M. Soss, and G. Toussaint (2002), Lower bounds for computing statistical depth, *Computational Statistics and Data Analysis*, *40*(2), 223-229.

Anderson, T. W. (1984), *An introduction to multivariate statistical analysis*, 2nd ed., 675 pp., Wiley, New York.

Barnett, V., and T. Lewis (1978), *Outliers in statistical data*, Reprint ed., 365 pp., Wiley, Chichester.

Barnett, V., and T. Lewis (1998), *Outliers in statistical data*, 3rd ed., 584 pp., Wiley, Chichester [etc.].

Barnett, V. (2004), *Environmental statistics : methods and applications*, xi, 293 p. pp., J. Wiley, Chichester, England ; Hoboken, N.J.

Bickel, P. J., and E. L. Lehmann (1975a), Descriptive statistics for nonparametric models. II. Location, *Ann. Statist.*, *3*(5), 1045--1069.

Bickel, P. J., and E. L. Lehmann (1975b), Descriptive statistics for nonparametric models. I. Introduction, *Ann. Statist.*, *3*(5), 1038-1044.

Bickel, P. J., and E. L. Lehmann (1976), Descriptive statistics for nonparametric models. III. Dispersion, *Ann. Statist.*, *4*(6), 1139-1158.

Bickel, P. J., and E. L. Lehmann (1979), Descriptive statistics for nonparametric models. IV. Spread, in *Contributions to statistics*, edited by Reidel, pp. 33-40, Dordrecht.

Caplin, A., and B. Nalebuff (1991a), Aggregation and social choice - a mean voter theorem, *Econometrica*, *59*(1), 1-23.

Caplin, A., and B. Nalebuff (1991b), Aggregation and imperfect competition - on the existence of equilibrium, *Econometrica*, *59*(1), 25-59.

Caplin, A., Nalebuff, B. (1988), On 64%-majority rule, *Econometrica*, *56*, 787-814.

Chebana, F., and T. B. M. J. Ouarda (2007), Multivariate L-moment homogeneity test, *Water Resources Research*, *43*(8).

Chebana, F., and T. B. M. J. Ouarda (2008), Depth and homogeneity in regional flood frequency analysis, *Water Resources Research*, *44*(11).

Chebana, F., and T. B. M. J. Ouarda (2009), Index flood-based multivariate regional frequency analysis, *Water Resources Research*, *45*(10).

Chebana, F., T. B. M. J. Ouarda, P. Bruneau, M. Barbet, S. El Adlouni, and M. Latraverse (2009), Multivariate homogeneity testing in a northern case study in the province of Quebec, Canada, *Hydrological Processes*, *23*(12), 1690-1700.

Chebana, F., and T. B. M. J. Ouarda (2010), Multivariate quantiles in hydrological frequency analysis, *Environmetrics*, *in press*.

Chow, V. T., D. R. Maidment, and L. R. Mays (1988), *Applied Hydrology*, 572 pp., McGraw-Hill, New York.

Cunnane, C. (1987), Review of statistical models for flood frequency estimation.

Dang, X., and R. Serfling Nonparametric depth-based multivariate outlier identifiers, and masking robustness properties, *Journal of Statistical Planning and Inference*.

Dang, X., and R. Serfling (2010), Nonparametric depth-based multivariate outlier identifiers, and masking robustness properties, *Journal of Statistical Planning and Inference*, *140*(1), 198-213.

648     Ghosh, A. K., and P. Chaudhuri (2005), On maximum depth and related classifiers, *Scandinavian*
649     *Journal of Statistics*, *32*(2), 327-350.
650     Helsel, D. R., R. M. Hirsch, and Geological Survey (U.S.) (2002), Statistical methods in water
651     resources, edited, U.S. Geological Survey, [Reston, Va.].
652     Hosking, J. R. M., and J. R. Wallis (1993), Some statistics useful in regional frequency analysis,
653     *Water Resources Research*, *29*(2), 271-281.
654     Hosking, J. R. M., and J. R. Wallis (1997), *Regional Frequency Analysis: An Approach Based on L-*
655     *Moments*, 240 pp., Cambridge University Press, Cambridge.
656     Huffer, F. W., and C. Park (2007), A test for elliptical symmetry, *Journal of Multivariate Analysis*,
657     *98*(2), 256-281.
658     Leon, C. A., and J. C. Massé (1993), A simplex Oja median - existence, uniqunes, stability,
659     *Canadian Journal of Statistics*, *21*(4), 397-408.
660     Li, J., and R. Y. Liu (2004), New nonparametric tests of multivariate locations and scales using data
661     depth, *Statistical Science*, *19*(4), 686-696.
662     Liu, R. Y. (1990), On a notion of data depth based on random simplices, *Annals of Statistics*, *18*(1),
663     405-414.
664     Liu, R. Y., and K. Singh (1993), A quality index based on data depth and multivariate rank-tests,
665     *Journal of the American Statistical Association*, *88*(421), 252-260.
666     Liu, R. Y. (1995), Control charts for multivariate processes, *Journal of the American Statistical*
667     *Association*, *90*(432), 1380-1387.
668     Liu, R. Y., J. M. Parelius, and K. Singh (1999), Multivariate analysis by data depth: Descriptive
669     statistics, graphics and inference, *Annals of Statistics*, *27*(3), 783-858.
670     Manly, B. F. J. (2005), *Multivariate statistical methods, a primer*, 3rd ed., 214 pp., Chapman &
671     Hall/CRC, Boca Raton.
672     Manzotti, A., F. J. Prérez, and A. J. Quiroz (2002), A statistic for testing the null hypothesis of
673     elliptical symmetry, *Journal of Multivariate Analysis*, *81*(2), 274-285.
674     Massé, J. C., and J. F. Plante (2003), A Monte Carlo study of the accuracy and robustness of ten
675     bivariate location estimators, *Computational Statistics and Data Analysis*, *42*(1-2), 1-26.
676     Massé, J. C. (2009), Multivariate trimmed means based on the Tukey depth, *Journal of Statistical*
677     *Planning and Inference*, *139*(2), 366-384.
678     MathWorks (2008), MATLAB Version 7.6.0.324, edited, MathWorks, Inc., , Natick, MA.
679     Mizera, I., and C. H. Müller (2004), Location-scale depth, *Journal of the American Statistical*
680     *Association*, *99*(468), 949-966.
681     Ngatchou-Wandji, J. (2009), Testing for symmetry in multivariate distributions, *Statistical*
682     *Methodology*, *6*(3), 230-250.
683     Oja, H. (1983), Descriptive statistics for multivariate distributions, *Statistics and Probability Letters*,
684     *1*(6), 327-332.
685     Oja, H., and A. Niinimaa (1985), Asymptotic Properties of the Generalized Median in the Case of
686     Multivariate Normality, *Journal of the Royal Statistical Society. Series B (Methodological)*, *47*(2),
687     372-377.
688     Ouarda, T. B. M. J., and F. Ashkar (1998), Effect of trimming on LP III flood quantile estimates,
689     *Journal of Hydrologic Engineering*, *3*(1), 33-42.
690     Ouarda, T. B. M. J., M. Hache, P. Bruneau, and B. Bobee (2000), Regional flood peak and volume
691     estimation in northern Canadian basin, *Journal of Cold Regions Engineering*, *14*(4), 176-191.
692     Rao, A. R., and K. H. Hamed (2000), *Flood Frequency Analysis*, 376 pp., CRC Press, Boca Raton.

693 Rousseeuw, P. J., and I. Ruts (1996), Bivariate location depth, *Applied Statistics-Journal of the*
694 *Royal Statistical Society Series C*, *45*(4), 516-526.
695 Rousseeuw, P. J., and I. Ruts (1999), The depth function of a population distribution, *Metrika*, *49*(3),
696 213-244.
697 Rousseeuw, P. J., I. Ruts, and J. W. Tukey (1999), The bagplot: A bivariate boxplot, *American*
698 *Statistician*, *53*(4), 382-387.
699 Rousseeuw, P. J., and K. Van Driessen (1999), Fast algorithm for the minimum covariance
700 determinant estimator, *Technometrics*, *41*(3), 212-223.
701 Rousseeuw, P. J., and A. Struyf (2004), Characterizing angular symmetry and regression symmetry,
702 *Journal of Statistical Planning and Inference*, *122*(1-2), 161-173.
703 Sakhanenko, L. (2008), Testing for ellipsoidal symmetry: A comparison study, *Computational*
704 *Statistics and Data Analysis*, *53*(2), 565-581.
705 Schervish, M. J. (1987), A Review of Multivariate Analysis, *Statistical Science*, *2*(4), 396-413.
706 Serfling, R. (2006), Multivariate symmetry and asymmetry, in *Encyclopedia of Statistical Sciences*,
707 edited by S. Kotz, N. Balakrishnan, C. B. Read and B. Vidakovic, pp. 5338-5345, Wiley.
708 Serfling, R., and P. Xiao (2007), A contribution to multivariate L-moments: L-comoment matrices,
709 *Journal of Multivariate Analysis*, *98*(9), 1765-1781.
710 Shiau, J. T. (2003), Return period of bivariate distributed extreme hydrological events, *Stochastic*
711 *Environmental Research and Risk Assessment*, *17*(1-2), 42-57.
712 Tukey, J. W. (1975), Mathematics and the picturing of data, paper presented at Proceedings of the
713 International Congress of Mathematicians, SIAM, Philadelphia, Vancouver
714 Wang, J., and R. Serfling (2005), Nonparametric multivariate kurtosis and tailweight measures,
715 *Journal of Nonparametric Statistics*, *17*(4), 441-456.
716 Wang, X., M. Gebremichael, and J. Yan (2010), Weighted likelihood copula modeling of extreme
717 rainfall events in Connecticut, *Journal of Hydrology*, *390*(1-2), 108-115.
718 Warner, R. M. (2008), *Applied statistics : from bivariate through multivariate techniques*, xxvi,
719 1101 p. pp., SAGE Publications, Thousand Oaks, Calif.
720 Wilcox, R. R., and H. J. Keselman (2004), Multivariate location: Robust estimators and inference,
721 *Journal of Modern Applied Statistical Methods*, *3*(1), 2-12.
722 Yue, S., T. B. M. J. Ouarda, B. Bobée, P. Legendre, and P. Bruneau (1999), The Gumbel mixed
723 model for flood frequency analysis, *Journal of Hydrology*, *226*(1-2), 88-100.
724 Zhang, L., and V. P. Singh (2006), Bivariate flood frequency analysis using the copula method,
725 *Journal of Hydrologic Engineering*, *11*(2), 150-164.
726 Zuo, Y., and R. Serfling (2000a), On the performance of some robust nonparametric location
727 measures relative to a general notion of multivariate symmetry, *Journal of Statistical Planning and*
728 *Inference*, *84*(1-2), 55-79.
729 Zuo, Y., and R. Serfling (2000b), General notions of statistical depth function, *Annals of Statistics*,
730 *28*(2), 461-482.
731 Zuo, Y. J. (2003), Finite sample tail behavior of multivariate location estimators, *Journal of*
732 *Multivariate Analysis*, *85*(1), 91-105.
733
734

## Appendix: brief presentation of depth functions

The main aim of introducing depth functions was to define multivariate extensions of the rank and order notions. Tukey [1975] presented pioneering work in this direction by proposing the half-space depth function. Several types of depth functions were defined later then standardized and classified by Zuo and Serfling [2000b]. A depth function $D(x; F)$, defined for a given cumulative distribution function $F$ on $R^d$ $(d \geq 1)$ and $x$ in $R^d$, is any bounded and nonnegative function that meets the following properties:

i. *Affine invariance*: the depth of a point $x \in R^d$ should not depend on the underlying coordinate system or, in particular, on the scales of the underlying measurements. That is, $D(Ax + b; F_{AX+b}) = D(x; F_X)$ holds for any random vector $X$ in $R^d$, any $d \times d$ nonsingular matrix $A$ and any $d$-vector $b$;

ii. *Maximality at center*: for a distribution having a uniquely defined center, the depth function should attain its maximum value at this center;

iii. *Monotonicity relative to deepest point*: as a point $x \in R^d$ moves away from the deepest point along any fixed ray through the center, the depth at $x$ should decrease monotonically;

iv. *Vanishing at infinity*: the depth of a point $x$ should be close to zero as the corresponding norm $\|x\|$ approaches infinity.

The following depth functions have received more attention in the literature [*Zuo and Serfling*, 2000b]:

1. **Tukey depth (called also the Half-space depth):** Given a probability $P$ on $R^d$ and $x \in R^d$, the *Half-space depth* [*Tukey*, 1975], noted *HD*, is given by:

756 $$HD(x;P) = \inf \{P(H) : H \text{ a closed halfspace that contains } x\} \qquad \text{(A1)}$$

757 The empirical *half-space* depth function is defined by replacing the probability function $P(H)$ by the

758 proportion of sample observations falling into a half-space $H$. An illustration based on a simple

759 example is given in Figure 6. The depth value of $\theta$ is the minimum number of observations falling

760 in the half-spaces (here 2) divided by the sample size 9. Note that $\theta$ does not belong to the sample.

761 2. **Oja depth (called also the Simplicial volume depth):** The *Simplicial volume depth* [*Oja,*

762 *1983*], noted *SVD*, is given through the expression:

763 $$SVD(x,F) = \left(1 + E\left[\Delta\left(S_n\left[x, X_1,..., X_d\right]\right)\right]\right)^{-1} \quad \text{for } x \in R^d \qquad \text{(A2)}$$

764 where $\Delta\left(S_n\left[x, x_1,..., x_d\right]\right)$ is the volume of the closed *d-simplex* $S_n\left[x, x_1,..., x_d\right]$ formed by the points

765 $x, x_1..., x_d \in R^d$. A *d-simplex* is defined as the convex hull of these points. This is a *d-dimensional*

766 generalization of triangles.

767 3. **Mahalanobis depth**: We introduce the *Mahalanobis* distance:

768 $$d_A^2(x,y) = (x-y)' A^{-1}(x-y) \qquad \text{(A3)}$$

769 where $x, y \in R^d$ are column vectors and $A$ is any semi-definite-positive matrix. Given a distribution

770 $F$, a scatter measure $A(F)$ and a location parameter $\mu(F)$, the Mahalanobis depth, noted *MD*, is:

771 $$MD(x,F) = \left(1 + d_{A(F)}^2\left(x, \mu(F)\right)\right)^{-1} \qquad \text{(A4)}$$

772 4. **Liu depth (called also the Simplicial depth) :** The *Simplicial depth* [*Liu, 1990*], noted *SD*, of

773 $x \in R^d$ with respect to a distribution $F$ is given by:

774 $$SD(x,F) = P_F \left\{x \in S_n\left[X_1,..., X_{d+1}\right]\right\} \qquad \text{(A5)}$$

775 where $S_n$ is as defined above and $X_i \sim F, \ i = 1,...,d+1$.

776 5. **Projection depth**: For a given distribution $F$ of a variable $X$, we define $F_{u'X}$ as the univariate

777 distribution of the variable $u'X$. Then, given a location and a scatter parameters $\mu(.)$ and $\sigma(.)$, the

778 *projection depth PD*(.) is defined as:

779
$$PD(x,F) = \sup_{\|u\|=1} \left| \left( u'x' - \mu(F_{u'X}) \right) \sigma^{-1}(F_{u'X}) \right| \qquad (A6)$$

780 where $\|.\|$ is the Euclidian norm. The empirical version of $PD$ is obtained by substituting the location

781 and scale measures $\mu(.)$ and $\sigma(.)$ with their estimations, and $F_{u'X}$ by the empirical distribution of the

782 sample $\{u'X_1, u'X_2,..., u'X_n\}$.

783 The computation of depth functions is generally not straightforward and requires specific

784 algorithms. For instance, Rousseeuw and Ruts [1996] and Aloupis et al. [2002] developed

785 algorithms for the computation of the half-space and the simplicial depth functions. The

786 Mahalanobis depth is among the simplest ones to evaluate. However, computational algorithms for

787 the projection depth are not available yet.

788 Depth functions are applied in several fields such as in econometric and social studies [*Caplin and*

789 *Nalebuff*, 1991a; b; 1988]. Liu and Singh [1993] and Liu [1995] employed depth functions in

790 industrial quality control. Recently, the depth-based approach proposed by Chebana and Ouarda

791 [2008] improved the performance of Canonical Correlation Analysis in the context of regional flood

792 frequency analysis. Depth functions were also investigated in nonparametric discriminant analysis

793 by Ghosh and Chaudhuri [2005]. Mizera and Müller [2004] defined and studied the location-scale

794 depth and gave some statistical applications.

795

**List of tables and figures**

796

797   Table 1: Depth values and *de*-classes for the flood peak-volume data set (Ashuapmushuan)

798   Table 2: Location parameters (Ashuapmushuan)

799   Table 3: Dispersion matrices (Ashuapmushuan)

800   Table 4: Outlier detection for the three considered bivariate series using Mahalanobis, Spatial and Tukey
801   outlyingness with normal and empirical thresholds (Ashuapmushuan)

802   Table 5: Tukey depth and outlyingness values for the flood peak-volume series (Magpie)

803   Figure 1a: Bagplots using Tukey depth (Ashuapmushuan)

804   Figure 1b: Contour plots using Tukey depth (Ashuapmushuan)

805   Figure 2: Location parameters: $(Q, V)$ left, $(D, V)$ middle and $(Q, D)$ right. Top figures present the location
806   parameters within the data and in the bottom figures a zoom is made to show the different location parameters
807   (Ashuapmushuan)

808   Figure 3: Scales using Tukey depth (Ashuapmushuan)

809   Figure 4a: Spherical skewness using Tukey depth (Ashuapmushuan)

810   Figure 4b: Elliptical skewness using Tukey depth (Ashuapmushuan)

811   Figure 4c: Antipodal skewness using Tukey depth (Ashuapmushuan)

812   Figure 4d: Angular skewness using Tukey depth (Ashuapmushuan)

813   Figure 5a: Kurtosis measure with L(p) using Tukey depth (Ashuapmushuan)

814   Figure 5b: Kurtosis measure with L*(p) using Tukey depth (Ashuapmushuan)

815   Figure 5c: Kurtosis measure with shrinkage using Tukey depth (Ashuapmushuan)

816   Figure 5d: Kurtosis measure with fan plots using Tukey depth (Ashuapmushuan)

817   Figure 5e: Kurtosis measure with quantile using Tukey depth (Ashuapmushuan)

818   Figure 6: Half-space depth evaluation for the point $\theta$ in an arbitrary generated sample. The numbers in boxes represent
819   the number of points in the associated half-space. The minimum value is 2 which gives the depth value of $\theta$ which is
820   equal to as 2 divided by the sample size.

821

822    **Table 1:** Depth values and *de*-classes for the flood peak-volume data set (Ashuapmushuan)

| Year | Q (m³/s) | V (day m³/s) | Oja Depth | Oja de-class | Tukey Depth | Tukey de-class | Liu Depth | Liu de-class | Mahalanobis Depth | Mahalanobis de-class |
|------|------|------|------|------|------|------|------|------|------|------|
| 1969 | 1380 | 50895 | 2.84E-07 | 1 | 0.3939 | 1 | 0.3310 | 1 | 0.9824 | 1 |
| 1973 | 1470 | 55766 | 2.75E-07 | 2 | 0.3636 | 2 | 0.3248 | 2 | 0.9211 | 2 |
| 1975 | 1260 | 48790 | 2.62E-07 | 3 | 0.3030 | 4 | 0.2980 | 4 | 0.8236 | 4 |
| 1984 | 1460 | 57769 | 2.58E-07 | 4 | 0.3333 | 3 | 0.3021 | 3 | 0.8023 | 5 |
| 1995 | 1550 | 51853 | 2.56E-07 | 5 | 0.3030 | 4 | 0.2892 | 5 | 0.8319 | 3 |
| 1993 | 1360 | 45263 | 2.50E-07 | 6 | 0.3030 | 4 | 0.2757 | 6 | 0.7436 | 6 |
| 1985 | 1210 | 47627 | 2.46E-07 | 7 | 0.2424 | 5 | 0.2515 | 7 | 0.7341 | 7 |
| 1976 | 1490 | 60767 | 2.31E-07 | 8 | 0.2121 | 6 | 0.2482 | 8 | 0.6420 | 9 |
| 1966 | 1650 | 54139 | 2.27E-07 | 9 | 0.2121 | 6 | 0.2368 | 9 | 0.6860 | 8 |
| 1972 | 1160 | 42497 | 2.22E-07 | 10 | 0.1818 | 7 | 0.2346 | 10 | 0.5794 | 10 |
| 1991 | 1130 | 49226 | 2.04E-07 | 11 | 0.1212 | 8 | 0.1683 | 12 | 0.5625 | 11 |
| 1978 | 1530 | 63663 | 2.03E-07 | 12 | 0.1818 | 7 | 0.1877 | 11 | 0.5121 | 12 |
| 1977 | 1370 | 60824 | 1.95E-07 | 13 | 0.0909 | 9 | 0.1602 | 14 | 0.5043 | 13 |
| 1981 | 1500 | 64631 | 1.88E-07 | 14 | 0.0909 | 9 | 0.1290 | 19 | 0.4478 | 14 |
| 1989 | 1490 | 41943 | 1.85E-07 | 15 | 0.1212 | 8 | 0.1606 | 13 | 0.4216 | 15 |
| 1965 | 1330 | 38682 | 1.81E-07 | 16 | 0.0909 | 9 | 0.1290 | 19 | 0.4161 | 16 |
| 1968 | 1100 | 37213 | 1.80E-07 | 17 | 0.0909 | 9 | 0.1345 | 17 | 0.3991 | 17 |
| 1983 | 1590 | 67223 | 1.72E-07 | 18 | 0.0909 | 9 | 0.1158 | 22 | 0.3900 | 19 |
| 1988 | 993 | 36882 | 1.69E-07 | 19 | 0.0909 | 9 | 0.1246 | 20 | 0.3498 | 23 |
| 1970 | 1780 | 66879 | 1.69E-07 | 20 | 0.0909 | 9 | 0.1437 | 15 | 0.3983 | 18 |
| 1986 | 1690 | 46735 | 1.68E-07 | 21 | 0.0909 | 9 | 0.1290 | 19 | 0.3667 | 20 |
| 1971 | 1420 | 38634 | 1.67E-07 | 22 | 0.0606 | 10 | 0.1107 | 23 | 0.3562 | 21 |
| 1967 | 934 | 39744 | 1.63E-07 | 23 | 0.0606 | 10 | 0.1294 | 18 | 0.3417 | 24 |
| 1992 | 1820 | 51752 | 1.59E-07 | 24 | 0.0909 | 9 | 0.1426 | 16 | 0.3411 | 25 |
| 1964 | 1780 | 68828 | 1.59E-07 | 25 | 0.0606 | 10 | 0.1184 | 21 | 0.3525 | 22 |
| 1980 | 949 | 33010 | 1.47E-07 | 26 | 0.0303 | 11 | 0.0909 | 25 | 0.2751 | 26 |
| 1990 | 1570 | 38568 | 1.39E-07 | 27 | 0.0303 | 11 | 0.0909 | 25 | 0.2553 | 27 |
| 1882 | 1920 | 50525 | 1.30E-07 | 28 | 0.0303 | 11 | 0.0909 | 25 | 0.2331 | 29 |
| 1979 | 2040 | 59254 | 1.25E-07 | 29 | 0.0606 | 10 | 0.0964 | 24 | 0.2368 | 28 |
| 1963 | 968 | 58538 | 1.12E-07 | 30 | 0.0606 | 10 | 0.0964 | 24 | 0.1953 | 30 |
| 1987 | 610 | 35600 | 1.07E-07 | 31 | 0.0303 | 11 | 0.0909 | 25 | 0.1626 | 31 |
| 1994 | 1170 | 74840 | 8.50E-08 | 32 | 0.0303 | 11 | 0.0909 | 25 | 0.1073 | 32 |
| 1974 | 2400 | 84198 | 8.10E-08 | 33 | 0.0303 | 11 | 0.0909 | 25 | 0.1027 | 33 |

823

824

825

826

827 **Table 2:** Location parameters (Ashuapmushuan)

828

| | | (Q, V) | | (D, V) | | (Q, D) | |
|---|---|---|---|---|---|---|---|
| **Mean** | | 1.43E+03 | 5.22E+04 | 84.3 | 5.22E+04 | 1.43E+03 | 84.3 |
| **Trimmed mean 5%** | **Tukey** | 1.43E+03 | 5.22E+04 | 84.2 | 5.22E+04 | 1.43E+03 | 84.0 |
| | **Oja** | 1.41E+03 | 5.08E+04 | 83.3 | 5.08E+04 | 1.41E+03 | 83.3 |
| | **Mahalanobis** | 1.41E+03 | 5.08E+04 | 83.3 | 5.08E+04 | 1.41E+03 | 83.3 |
| | **Liu** | 1.43E+03 | 5.22E+04 | 84.2 | 5.22E+04 | 1.43E+03 | 84.0 |
| **Trimmed mean 10%** | **Tukey** | 1.43E+03 | 5.21E+04 | 84.1 | 5.21E+04 | 1.43E+03 | 83.6 |
| | **Oja** | 1.43E+03 | 5.09E+04 | 81.8 | 4.99E+04 | 1.43E+03 | 82.4 |
| | **Mahalanobis** | 1.43E+03 | 5.09E+04 | 82.4 | 5.08E+04 | 1.43E+03 | 82.4 |
| | **Liu** | 1.43E+03 | 5.21E+04 | 84.1 | 5.21E+04 | 1.43E+03 | 83.6 |
| **Median** | **Componentwise** | 1.46E+03 | 5.09E+04 | 80.0 | 5.09E+04 | 1.46E+03 | 80.0 |
| | **Tukey** | 1.41E+03 | 5.13E+04 | 81.0 | 5.03E+04 | 1.40E+03 | 81.0 |
| | **Oja** | 1.40E+03 | 5.15E+04 | 80.0 | 5.03E+04 | 1.43E+03 | 81.0 |
| | **Mahalanobis** | 1.38E+03 | 5.09E+04 | 83.0 | 4.88E+04 | 1.49E+03 | 84.0 |
| | **Liu** | 1.38E+03 | 5.09E+04 | 80.0 | 5.09E+04 | 1.38E+03 | 80.0 |
| | **Spacial** | 1.50E+03 | 5.09E+04 | 80.0 | 5.09E+04 | 1.46E+03 | 84.0 |

829

830    **Table 3:** Dispersion matrices (Ashuapmushuan)

831

| | | (Q, V) | | (D, V) | | (Q, D) | |
|---|---|---|---|---|---|---|---|
| **Dispersion (0%)** | | 1.29E+05 | 2.67E+06 | 2.11E+02 | 1.03E+05 | 1.29E+05 | -9.56E+02 |
| | | 2.67E+06 | 1.55E+08 | 1.03E+05 | 1.55E+08 | -9.56E+02 | 2.11E+02 |
| **Trimmed dispersion 5%** | **Tukey** | 1.25E+05 | 2.75E+06 | 1.64E+02 | 7.39E+04 | 9.60E+04 | -6.32E+02 |
| | | 2.75E+06 | 1.36E+08 | 7.39E+04 | 1.35E+08 | -6.32E+02 | 2.09E+02 |
| | **Oja** | 1.00E+05 | 1.81E+06 | 1.68E+02 | 7.33E+04 | 1.13E+05 | -5.60E+02 |
| | | 1.81E+06 | 1.13E+08 | 7.33E+04 | 1.20E+08 | -5.60E+02 | 1.58E+02 |
| | **Mahalanobis** | 1.00E+05 | 1.81E+06 | 1.79E+02 | 9.03E+04 | 1.16E+05 | -4.92E+02 |
| | | 1.81E+06 | 1.13E+08 | 9.03E+04 | 1.16E+08 | -4.92E+02 | 1.59E+02 |
| | **Liu** | 1.18E+05 | 2.73E+06 | 2.29E+02 | 1.08E+05 | 9.47E+04 | -5.97E+02 |
| | | 2.73E+06 | 1.52E+08 | 1.08E+05 | 1.36E+08 | -5.97E+02 | 2.28E+02 |
| **Trimmed dispersion 10%** | **Tukey** | 1.13E+05 | 2.47E+06 | 1.62E+02 | 7.78E+04 | 8.66E+04 | -3.82E+02 |
| | | 2.47E+06 | 1.30E+08 | 7.78E+04 | 1.22E+08 | -3.82E+02 | 1.98E+02 |
| | **Oja** | 8.37E+04 | 1.61E+06 | 1.28E+02 | 6.07E+04 | 8.62E+04 | -1.15E+02 |
| | | 1.61E+06 | 1.04E+08 | 6.07E+04 | 1.08E+08 | -1.15E+02 | 1.53E+02 |
| | **Mahalanobis** | 8.37E+04 | 1.61E+06 | 1.52E+02 | 8.29E+04 | 8.65E+04 | -1.24E+02 |
| | | 1.61E+06 | 1.04E+08 | 8.29E+04 | 1.06E+08 | -1.24E+02 | 1.54E+02 |
| | **Liu** | 1.04E+05 | 2.51E+06 | 2.23E+02 | 1.12E+05 | 8.46E+04 | -3.04E+02 |
| | | 2.51E+06 | 1.34E+08 | 1.12E+05 | 1.09E+08 | -3.04E+02 | 2.18E+02 |

**Table 4:** Outlier detection for the three considered bivariate series using Mahalanobis, Spatial and Tukey outlyingness with normal and empirical thresholds (Ashuapmushuan)

|  |  |  | **Mahalanobis** | **Spatial** | **Tukey** |
|---|---|---|---|---|---|
| $(Q, V)$ | Normal | Threshold | 0.7297 | --- | 0.9931 |
|  |  | Outliers (years) | 1989-1995 | --- | None |
|  | Empirical | Threshold | 0.8973 | 0.9695 | 0.9394 |
|  |  | Outliers (years) | None | None | None |
| $(D, V)$ | Normal | Threshold | 0.7297 | --- | 0.9931 |
|  |  | Outliers (years) | 1982;1988; 1990-1995 | --- | None |
|  | Empirical | Threshold | 0.9181 | 0.9697 | 0.9394 |
|  |  | Outliers (years) | None | None | None |
| $(Q, D)$ | Normal | Threshold | 0.7297 | --- | 0.9931 |
|  |  | Outliers (years) | 1986;1988; 1990-1995 | --- | None |
|  | Empirical | Threshold | 0.8921 | 0.9695 | 0.9394 |
|  |  | Outliers (years) | None | None | None |

**Table 5:** Tukey depth and outlyingness values for the flood peak-volume series (Magpie)

| Year | $Q$ | $V$ | Tukey Depth | $O_{MD}$ | $O_S$ | $O_{HD}$ |
|---|---|---|---|---|---|---|
| 1979 | 886.67 | 2088.92 | 0.2692 | 0.0571 | 0.1361 | 0.4615 |
| 1980 | 849.67 | 2357.02 | 0.3846 | 0.1971 | 0.1567 | 0.2308 |
| 1981 | 1456.67 | 3909.14 | 0.0385 | **0.8851** | **0.9563** | **0.9231** |
| 1982 | 1270.00 | 2443.15 | 0.0385 | 0.8032 | 0.6246 | **0.9231** |
| 1983 | 974.67 | 3012.18 | 0.0769 | 0.6700 | 0.8500 | 0.8462 |
| 1984 | 1056.67 | 2751.69 | 0.1154 | 0.4713 | 0.6857 | 0.7692 |
| 1985 | 787.00 | 1574.21 | 0.1538 | 0.4623 | 0.4815 | 0.6923 |
| 1986 | 610.33 | 1536.34 | 0.1154 | 0.5306 | 0.6026 | 0.7692 |
| 1987 | 344.33 | 1069.86 | 0.0385 | 0.8225 | 0.9204 | **0.9231** |
| 1988 | 843.33 | 2374.49 | 0.3077 | 0.2390 | 0.2455 | 0.3846 |
| 1989 | 678.67 | 1534.53 | 0.1923 | 0.4534 | 0.5395 | 0.6154 |
| 1990 | 506.33 | 1752.06 | 0.0769 | 0.7223 | 0.5603 | 0.8462 |
| 1991 | 740.00 | 2260.57 | 0.1538 | 0.4461 | 0.3003 | 0.6923 |
| 1992 | 710.80 | 1128.71 | 0.0385 | 0.7223 | 0.8923 | **0.9231** |
| 1993 | 666.80 | 1407.32 | 0.1538 | 0.5400 | 0.6964 | 0.6923 |
| 1994 | 932.90 | 2722.55 | 0.1538 | 0.4802 | 0.6113 | 0.6923 |
| 1995 | 868.77 | 2192.44 | 0.3462 | 0.0068 | 0.0324 | 0.3077 |
| 1996 | 886.90 | 2476.36 | 0.3077 | 0.2644 | 0.3562 | 0.3846 |
| 1997 | 697.30 | 2665.87 | 0.0385 | 0.7817 | 0.6607 | **0.9231** |
| 1998 | 825.00 | 1843.60 | 0.3077 | 0.1963 | 0.2717 | 0.3846 |
| 1999 | 1306.67 | 2652.26 | 0.0385 | 0.8042 | 0.7450 | **0.9231** |
| 2000 | 858.90 | 2492.65 | 0.2308 | 0.3526 | 0.4095 | 0.5385 |
| 2001 | 732.50 | 1188.92 | 0.0769 | 0.7053 | 0.8076 | 0.8462 |
| 2002 | 999.60 | 1485.36 | 0.0385 | 0.8045 | 0.6758 | **0.9231** |
| 2003 | 1004.93 | 1883.80 | 0.1538 | 0.6236 | 0.4102 | 0.6923 |
| 2004 | 842.57 | 2802.32 | 0.0769 | 0.6783 | 0.7252 | 0.8462 |

Bold character indicates outlyingness of the detected outliers

Figure 1a : Bagplots using Tukey depth : (*Q*, *V*) left, (*D*, *V*) middle and (*Q*, *D*) right (Ashuapmushuan)



Figure 1b : Contour plots using Tukey depth : (*Q*, *V*) left, (*D*, *V*) middle and (*Q*, *D*) right (Ashuapmushuan)

42

Figure 2: Location parameters: (*Q, V*) left, (*D, V*) middle and (*Q, D*) right. Top figures present the location parameters within the data and in the bottom figures a zoom is made to show the different location parameters (Ashuapmushuan)
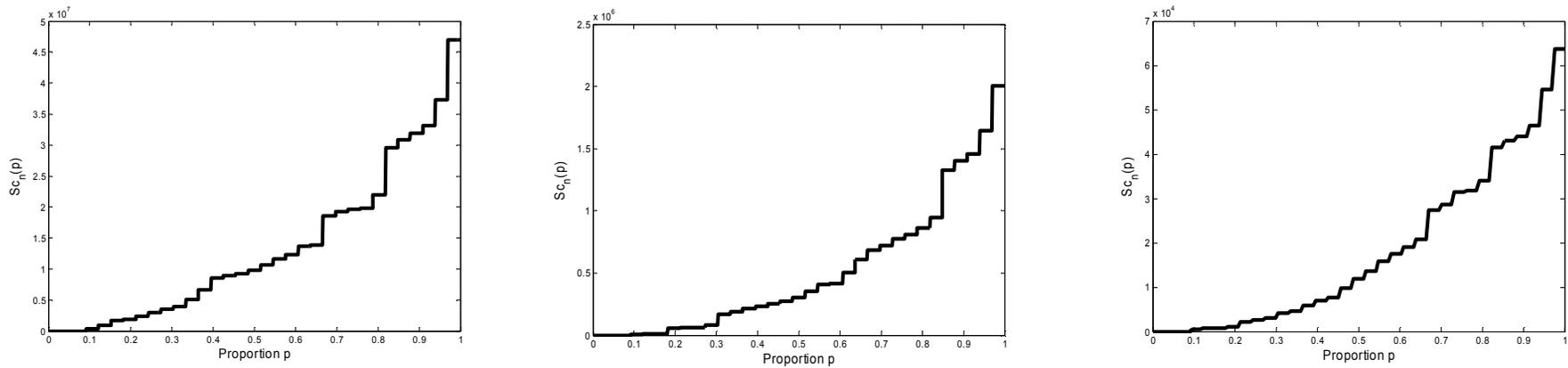
Figure 3 : Scalar scales using Tukey depth : $(Q, V)$ left, $(D, V)$ middle and $(Q, D)$ right (Ashuapmushuan)

Figure 4a : Spherical skewness using Tukey depth : (*Q*, *V*) left, (*D*, *V*) middle and (*Q*, *D*) right (Ashuapmushuan)



Figure 4b : Elliptical skewness using Tukey depth: (*Q*, *V*) left, (*D*, *V*) middle and (*Q*, *D*) right (Ashuapmushuan)

Figure 4c : Antipodal skewness using Tukey depth : (*Q, V*) left, (*D, V*) middle and (*Q, D*) right (Ashuapmushuan)



Figure 4d : Angular skewness using Tukey depth : (*Q, V*) left, (*D, V*) middle and (*Q, D*) right (Ashuapmushuan)

Figure 5a : Kurtosis measure with L(p) using Tukey depth : (*Q, V*) left, (*D, V*) middle and (*Q, D*) right (Ashuapmushuan)
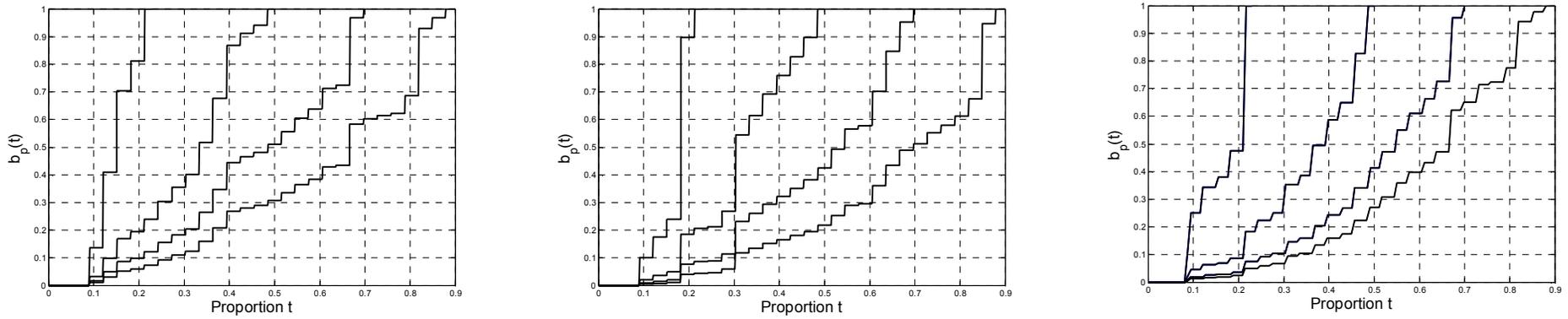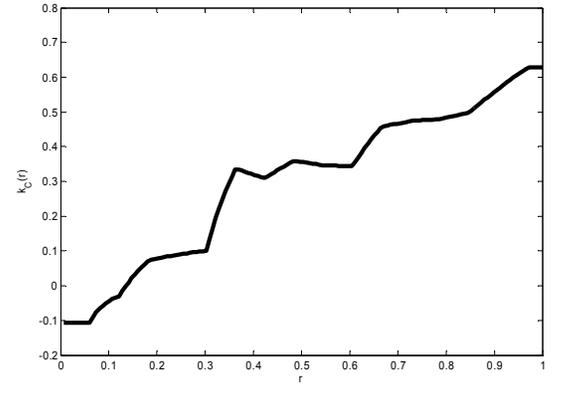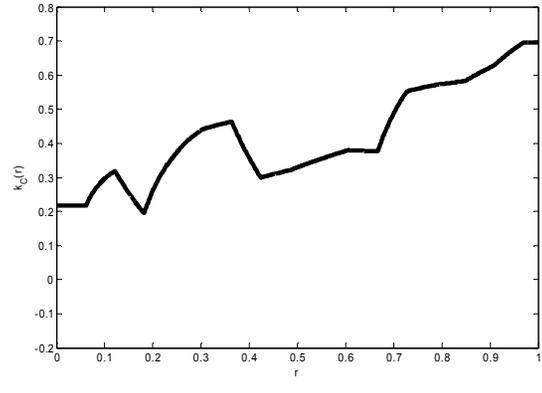


Figure 5b : Kurtosis measure with L*(p) using Tukey depth : (*Q, V*) left, (*D, V*) middle and (*Q, D*) right (Ashuapmushuan)

Figure 5c : Kurtosis measure with shrinkage using Tukey depth : (*Q*, *V*) left, (*D*, *V*) middle and (*Q*, *D*) right (Ashuapmushuan)



Figure 5d : Kurtosis measure with fan plots using Tukey depth : (*Q*, *V*) left, (*D*, *V*) middle and (*Q*, *D*) right (Ashuapmushuan)
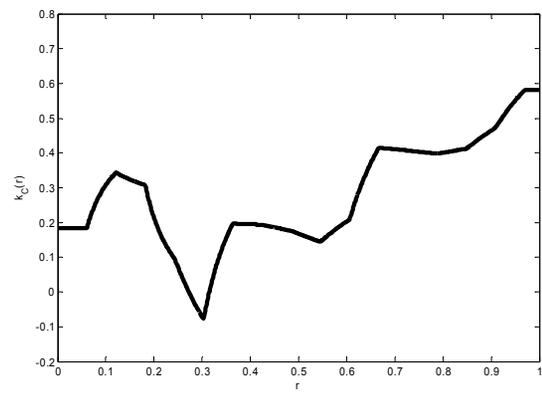
Figure 5e : Kurtosis measure with quantile using Tukey depth : ($Q$, $V$) left, ($D$, $V$) middle and ($Q$, $D$) right (Ashuapmushuan)
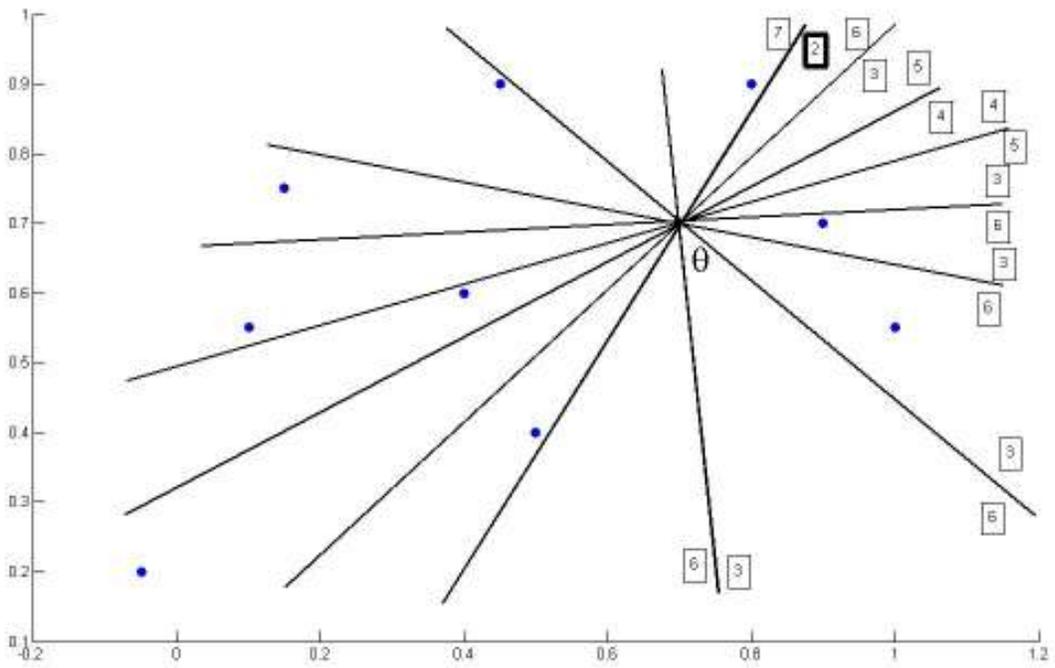
Figure 6: Half-space depth evaluation for the point $\theta$ in an arbitrary generated sample. The numbers in boxes represent the number of points in the associated half-space. The minimum value is 2 which gives the depth value of $\theta$ which is equal to as 2 divided by the sample size.