

1 Exploratory functional flood frequency analysis and
2 outlier detection

3 Fateh Chebana[§], Sophie Dabo-Niang[‡] and Taha B.M.J Ouarda[§]

4 [§] Canada Research Chair on the Estimation of Hydrometeorological Variables, INRS-ETE,
5 490 rue de la Couronne, Quebec (QC), Canada G1K 9A9

6 [‡] Laboratoire EQUIPPE, Université Charles De Gaulle, Lille 3, maison de la recherche,
7 domaine du pont de bois, BP 60149, 59653 Villeneuve d'Ascq Cedex, France.

8 January, 20th 2012

9 [§] Corresponding author:

10 Tel: +1 (418) 654-2542

11 Fax: +1 (418) 654-2600

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

13 **Abstract:** The prevention of flood risks and the effective planning and management of water
14 resources require river flows to be continuously measured and analyzed at a number of stations.
15 For a given station, a hydrograph can be obtained as a graphical representation of the temporal
16 variation of flow over a period of time. The information provided by the hydrograph is essen-
17 tial to determine the severity of extreme events and their frequencies. A flood hydrograph is
18 commonly characterized by its peak, volume and duration. Traditional hydrological frequency
19 analysis (FA) approaches focused separately on each of these features in a univariate context.
20 Recent multivariate approaches considered these features jointly in order to take into account
21 their dependence structure. However, all these approaches are based on the analysis of a num-
22 ber of characteristics, and do not make use of the full information content of the hydrograph.
23 The objective of the present work is to propose a new framework for frequency analysis using
24 the hydrographs as curves: functional data. In this context, the whole hydrograph is considered
25 as one infinite dimensional observation. This context allows to provide more effective and effi-
26 cient estimates of the risk associated with extreme events. The proposed approach contributes
27 to addressing the problem of lack of data commonly encountered in hydrology by fully em-
28 ploying all the information contained in the hydrographs. A number of functional data analysis
29 tools are introduced and adapted to flood FA with a focus on exploratory analysis as a first stage
30 towards a complete functional flood FA. These methods, including data visualization, location
31 and scale measures, principal component analysis as well as outlier detection, are illustrated in
32 a real-world flood analysis case study from the province of Quebec, Canada.

33 **Key Words:** Functional data, frequency analysis, hydrology, flood, outliers, exploratory analy-
34 sis, principal component analysis.

1 Introduction

Extreme hydrological events such as floods, droughts and rain storms may have significant economic and social consequences. Hydrological frequency analysis (FA) procedures are essential and commonly used for the analysis and prediction of such extreme events, which have a direct impact on reservoir management and dam design. Flood FA is based on the estimation of the probability $P(X > x_T)$ of exceedence of the event x_T corresponding to a quantile of a given return period T e.g. $T = 10, 50$ or 100 years. The random variable X is commonly taken to be the peak of the flood which is the maximum of the daily streamflow series during a hydrological year or season. 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).

The accurate estimation of the risk associated with the design and operation of water infrastructures requires a good knowledge of flood characteristics. Indeed, an overestimation of the design flood leads to an over-sizing of hydraulic structures and, would therefore involve additional costs, while underestimation of design floods leads to material damages and loss of human lives. Flood FA is commonly employed to study this risk. It has been traditionally carried out for the analysis of flood peaks in a univariate context. The reader is referred, e.g. to Cunnane (1987) and Rao and Hamed (2000).

In general, a flood is described through a number of correlated characteristics, e.g. peak, volume and duration. The univariate treatment of each flood characteristic ignores their dependence structure. Consequently, the univariate framework is less representative of the phenomenon and reduces the risk estimation accuracy. Thereafter, several authors focused on the joint treatment of flood characteristics through the use of a number of multivariate techniques such as multi-

57 variate distributions and copulas (e.g. Yue et al., 1999; Shiau, 2003; Zhang and Singh, 2006;
58 Chebana and Ouarda, 2011a). Multivariate studies contributed to the improvement of the esti-
59 mation accuracy and provide information concerning the dependence structure between flood
60 characteristics. The multivariate framework is applied in several hydrological events, such as
61 floods, droughts and storms. For instance in floods, it is used for hydraulic structure design and
62 extreme event prediction purposes (see Chebana and Ouarda, 2011a for recent references).

63 Despite their usefulness, univariate and multivariate FA approaches have a number of limita-
64 tions and drawbacks. The separate or joint use of hydrograph characteristics constitutes a major
65 simplification of the real phenomenon. Furthermore, the way these characteristics can be deter-
66 mined is neither unique nor objective (in particular, flood starting and ending dates). In addition,
67 each flood characteristic can be seen as a real-valued transformation of the hydrograph, e.g. the
68 peak is the maximum. For hydrological applications, the bivariate setting is largely employed to
69 treat two hydrological variables. A limited number of studies deals with the trivariate one, e.g.
70 Serinaldi and Grimaldi (2007) and Zhang and Singh (2007). The trivariate models generally
71 suffer from less representativity and formulation complexity. Note that, in general, the number
72 of associated parameters grows up rapidly with the dimension of the model and therefore the
73 generated uncertainty increases. In addition, higher dimensions are not considered in hydrolog-
74 ical practice. Finally, given the lack of data in hydrology, working with a limited number of
75 extracted characteristics represents a loss of information in comparison to the overall available
76 series.

77 The main data source in FA is daily streamflow series, which during a year constitutes a hydro-
78 graph, from which the univariate and multivariate variables are extracted. The total information

79 that is available in a hydrograph is necessary for the effective planning of water resources and
80 for the design and management of hydraulic structures. The entire hydrograph, as a curve with
81 respect to time, can be considered as a single observation within the functional context. In the
82 univariate and the multivariate settings an observation is respectively a real value and a vector.
83 Therefore, the functional framework which treats the whole hydrograph as a functional obser-
84 vation (function or curve) is more representative of the real phenomena and makes better use of
85 available data. Figure 1 illustrates and summarizes the three frameworks.

86 In the hydrological literature, there were some efforts towards a representation of the hydro-
87 graph as a function, such as in the study of the design flood hydrograph, e.g. Yue et al. (2002),
88 and in the flow duration curve study e.g. by Castellarin et al. (2004) where the mean, median
89 and variation are presented as curves. These studies underlined the importance to consider the
90 shape of the hydrograph which is necessary, for instance, for water resources planning, design
91 and management. The shape of flood hydrographs for a given river may change according to
92 the observed storm or snowmelt events. More practical issues and examples related to the hy-
93 drograph can be found for instance in Yue et al. (2002) or Chow et al. (1988). Note that the
94 main flood characteristics, i.e. peak, volume and duration, can not completely capture the shape
95 of the hydrograph. The study of the hydrographs in Yue et al. (2002), and similar studies, are
96 simplistic and limited, as they approximated the flood hydrograph using a two-parameter beta
97 density and considered only single-peak hydrographs. On the other hand, the flow duration
98 curve approach (Castellarin et al., 2004) is in the univariate setting and the presented functional
99 elements (e.g. mean and median curves) are important but remain limited. The previous studies
100 show the need to introduce a statistical framework to study the whole hydrograph and to per-

101 form further statistical analysis. The functional framework is more general and more flexible
102 and can represent a large variety of hydrographs.

103 Functional data are becoming increasingly common in a variety of fields. This has sparked a
104 growing attention in the development of adapted statistical tools that allow to analyze such kind
105 of data. For instance, Ramsay and Silverman (2005), Ferraty and Vieu (2006) and Dabo-Niang
106 and Ferraty (2008) provided detailed surveys of a number of parametric and nonparametric
107 techniques for the analysis of functional data. In practice, the use of functional data analysis
108 (FDA) has benefited from the availability of the appropriate statistical tools and high perfor-
109 mance computers. Furthermore, the use of FDA allows to make the most of the information
110 contained in the functional data. The aims of FDA are mainly the same as in the classical
111 statistical analysis, e.g. representing and visualizing the data, studying variability and trends,
112 comparing different data sets, as well as modeling and predicting. The majority of classical
113 statistical techniques, such as principal component, linear models, confidence interval estima-
114 tion and outlier detection, were extended to the functional context (e.g. Ramsay and Silverman,
115 2005). The application of FDA has been successfully carried out, for instance, in the case of the
116 El Niño climatic phenomenon (Ferraty et al., 2005) and radar wave curve classification (Dabo-
117 Niang et al., 2007). Dabo-Niang et al. (2010) proposed a spatial heterogeneity index to compare
118 the effects of bioturbation on oxygen distribution. Delicado et al. (2008) and Monestiez and
119 Nerini (2008) considered spatial functional kriging methods to model different temperature se-
120 ries. Sea ice data are treated in the FDA context by Koulis et al. (2009).

121 The functional methodology constitutes a natural extension of univariate and multivariate hy-
122 drological FA approaches (see Figure 1). This new approach uses all available data by em-

123 ploying the whole hydrograph as a functional observation. In other words, FDA permits to
124 exhaustively analyze hydrological data by conducting one analysis on the whole data instead
125 of several univariate or multivariate analysis. In addition, the approach proposed by Yue et al.
126 (2002) can be generalized in the FDA context where it becomes more flexible and includes hy-
127 drographs with different shapes such as multi-peak ones.

128 Given the above arguments, for hydrological applications, the functional context could be seen
129 as an alternative framework to the univariate and multivariate ones, or it can also be employed
130 as a parallel complement to bring additional insight to those obtained by the two other frame-
131 works. The main objective of the present paper is to attract attention to the functional nature
132 of data that can be used in all statistical techniques for hydrological applications through the
133 FDA framework. A second objective is to introduce some of the FDA techniques, point out
134 their advantages and illustrate their applicability in the hydrological framework. In the present
135 paper, we focus on hydrological FA.

136 Four main steps are required in order to carry out a comprehensive hydrological FA: i) de-
137 scriptive and exploratory analysis and outlier detection, ii) verification of FA assumptions, i.e.
138 stationarity, homogeneity and independence, iii) modeling and estimation and iv) evaluation and
139 analysis of the risk. The first step (i) is commonly carried out in univariate hydrological FA as
140 pointed out, e.g. by Rao and Hamed (2000), Kite (1988) and Stedinger et al. (1993) whereas in
141 the multivariate framework it was investigated recently by Chebana and Ouarda (2011b). Con-
142 trary to the univariate setting, exploratory analysis in the multivariate and functional settings is
143 not straightforward and requires more efforts. Table 1 summarizes the four FA steps and their
144 status in each one of the three frameworks. It is indicated that the specific aim of the present

145 paper is to treat step (i) which deals with data visualization, location and scale measures as well
146 as outlier detection. A new non-graphical method to detect functional outliers is also proposed
147 in the present paper. The presented techniques are applied to floods based on daily streamflow
148 series from a basin in the province of Quebec, Canada.

149 Exploratory data analysis as a preliminary step of FA is useful for the comparison of hydrologi-
150 cal samples and for the selection of the appropriate model for hydrological variables. It consists
151 in a close inspection of the data to quantify and summarize the properties of the samples, for
152 instance, through location and scale measures. Outliers can have negative impacts on the se-
153 lection of the appropriate model as well as on the estimation of the associated parameters. In
154 order to base the inference on the right data set, detection and treatment of outliers are also
155 important elements of FA (Barnett and Lewis, 1998). Therefore, it is essential to start with the
156 basic analysis (step i) in order to perform a complete functional FA.

157 This paper is organized as follows. The theoretical background of functional statistical methods
158 is presented in Section 2 in its general form. In Section 3, the functional framework is adapted
159 to floods. The functional FA methods are applied, in Section 4, to a real-wold case study rep-
160 resenting daily streamflows from the province of Quebec, Canada. A discussion as well as a
161 comparison with multivariate FA are also reported in Section 4. Conclusions and perspectives
162 are presented in the last section.

163 **2 Functional data analysis background**

164 This section presents the general functional techniques. It is composed of four parts represent-
165 ing FDA phases: first, data smoothing is discussed, second location and scale parameters are

166 introduced, then functional principal component analysis (FPCA) is described and finally data
167 visualization and outlier detection methods are presented.

168 Data are generally measured in discrete time steps such as hours or days. Therefore, the first
169 phase in FDA consists in the conversion of observed discrete data to functional data. Once
170 the discrete data are transformed to curves, they can be analyzed in the functional framework.
171 In a descriptive statistical study, it is of interest to obtain estimates of the location and scale
172 parameters within FDA. The next phase in the considered FDA is to extract information from
173 functional data using FPCA where the corresponding scores to these components are the basis
174 for visualization and outlier detection.

175 **2.1 Data smoothing**

176 The objective of this step is to prepare data to be used in the FDA context. As a preparation
177 step of the data to be employed, it is analogous to the step of extracting peaks in univariate
178 FA or peak and volume series in the multivariate FA. Note that the statistical object of FDA is
179 a function (curve) as shown in Figure 1. However, the curves are not observed, instead, only
180 discrete measurements of the curves are available. In the case where data series are of good
181 quality and long enough records, one can simply interpolate the measurements to obtain the
182 curves, e.g. for rainfall series. Otherwise, smoothing can be required. This is typically the
183 case for diffusive processes like in the present study of floods. However, even in the first case,
184 smoothing could be necessary depending on the goal of the study (e.g. Ramsay and Silverman,
185 2005).

186 Let $\mathbf{Y}_i = (y_i(t_1), \dots, y_i(t_T))$, $i = 1, \dots, n$ be a set of n discrete observations where each $t_j \in$

187 $\mathcal{C} \subset \mathbb{R}^+$, $j = 1, \dots, T$ is the j th record time point from a given time subset \mathcal{C} . For a fixed
 188 observation i , each set of measurements $(y_i(t_1), \dots, y_i(t_T))$ is converted to be a functional data
 189 (curve) denoted $y_i(t)$ by using a smoothing technique where the index t covers the continuous
 190 subset \mathcal{C} . To this end, we suppose that the discrete observation $(y_i(t_j))_{j=1, \dots, T}$ is fitted using the
 191 regression model:

$$y_i(t_j) = x_i(t_j) + \epsilon_{ij} \quad i = 1, \dots, n \quad \text{and} \quad j = 1, \dots, T \quad (1)$$

192 where ϵ_{ij} are the errors and the functions $x_i(\cdot)$ are linear combinations of basis functions $\phi_k(\cdot)$,
 193 that permit to explain most of the variation contained in the functional observations:

$$x_i(t) = \sum_{k=0}^p \hat{c}_k^i \phi_k(t) \quad \text{for } t \in \mathcal{C}. \quad (2)$$

194 The functional data set $(y_i(t))_{i=1, \dots, n}$ are then given by:

$$y_i(t) = \hat{x}_i(t) = \sum_{k=0}^p \hat{c}_k^i \phi_k(t), \quad t \in \mathcal{C} \quad (3)$$

195 where the estimated coefficients \hat{c}_k^i are obtained by minimizing the following sum:

$$SSE(i) = \sum_{j=1}^T (y_i(t_j) - x_i(t_j))^2, \quad \text{for } i = 1, \dots, n \quad (4)$$

196 For more details, the reader is referred, for instance, to Ramsay and Silverman (2005). A num-
 197 ber of possible types of basis $\phi_k(\cdot)$ have been presented in the literature. Most of the practical
 198 situations are treated with the well-known basis, such as, polynomial, wavelet, Fourier and the
 199 various Spline versions. Fourier and B -Spline basis are widely employed in the FDA context.
 200 The functional representation uses Fourier series for periodic or near periodic data. When the
 201 data are far from being periodic, spline approximations are commonly used in FDA for most

202 problems involving non-periodic data (Ramsay and Silverman, 2005). Splines are more flexible
 203 but more complicated than Fourier series. The latter allows capturing the seasonal variability
 204 while the Spline series captures high and low values of the data (Ramsay and Silverman, 2005,
 205 Koulis et al., 2009). In general, the basis functions or the smoothing method to use should be
 206 based on objective considerations depending mainly on the nature of the data to be studied.
 207 Fourier basis functions $(\phi_k(\cdot))_{k=0,\dots,p}$ are defined by:

$$\phi_0(t) = 1, \phi_{2j-1}(t) = \sin(jwt), \phi_{2j}(t) = \cos(jwt), w = 2\pi/T. \quad (5)$$

208 Splines are piecewise polynomials defined on subintervals of the range of the observations.
 209 In each subinterval, the Spline is a polynomial function with a fixed degree but could be with
 210 different shapes. For instance, when the polynomial degree is three, we talk about cubic splines.
 211 For a comprehensive review about splines, the reader is referred to De Boor (2001).
 212 Note that the aim of using the above expansion (3) is to obtain smooth functions to be employed
 213 as observations in FDA. In this case, the expansion series need not to be interpreted since the
 214 interest is not to extract a signal from the whole series. However, the number of basis functions
 215 to be selected is important where a large number leads to over-fitting of the data while a small
 216 number leads to under-fitting. Hence, the smoothing degree of the obtained functions to be
 217 employed as observations depends on the aim of the analysis, e.g. in principal component
 218 analysis, the aim is to capture a large variability rather than to reach the peaks. For more
 219 flexibility, a penalty term can be added to (4) to ensure the regularity of the smoothed functions.
 220 More details can be found for instance in Ramsay and Silvermann (2005) and Wahba (1990).

2.2 Location and scale parameters for functional variables

In a descriptive statistical study, we generally begin by looking for centrality and dispersion properties of a given sample. A location parameter summarizes the data and indicates where most of the data are located. Scale parameters are useful to measure the dispersion of a sample and also to compare different samples. These notions are useful in hydrology since they appear in almost all commonly employed probability distributions and models. In hydrology, location curves can also be used to characterize a given basin and to proceed to comparison or grouping of a set of basins. The scale measures can be used in a similar way but at a second level. In the setting of real or multivariate random variables, this is usually done through the mean, median, mode, variance, covariance and correlation. To avoid the possibility of missing important information, it is generally recommended to employ more than one measure for each feature. For instance, by looking only at the mean of the sample one might miss a possible heterogeneity in the population which would be captured by the mode. Obviously, these same problems will also appear when one studies a sample composed of curves $\{y_i(t), t \in \mathcal{C}\}, i = 1, \dots, n$. In this setting, it is straightforward to define the mean curve $\bar{y}(\cdot)$ of the sample as:

$$\bar{y}(t) = \frac{1}{n} \sum_{i=1}^n y_i(t), \quad t \in \mathcal{C}. \quad (6)$$

One has to use this mean curve carefully according to the shape of the data. For instance, if the data exhibit a high roughness degree the mean curve could be less informative.

Robust and efficient alternatives to the sample mean are the median and the trimmed mean (e.g., Ouarda and Ashkar, 1998). In the functional context, both measures are based on the statistical notion of depth function which is initially introduced in the multivariate context. The aim of

241 depth functions is to extend the notion of ranking for a multivariate sample. These functions
 242 are introduced by Tukey (1975) and are introduced and applied to water sciences by Chebana
 243 and Ouarda (2008). Recently, the notion of depth has been extended to functional data (e.g.
 244 Fraiman and Muniz, 2001 and Febrero et al., 2008). The reader is referred to Chebana and
 245 Ouarda (2011b) for hydrological applications and a brief review and to Zuo and Serfling (2000)
 246 for a general and detailed description.

247 Fraiman and Muniz (2001) presented the definition of trimmed means in the functional con-
 248 text which are based on the empirical α -trimmed functional region. It is defined by $TR_\alpha :=$
 249 $\{x, D_n(x) \geq \alpha\}$ for $0 < \alpha < 1$ where $D_n(\cdot)$ is an empirical functional depth function, as the
 250 various ones defined, e.g. in Fraiman and Muniz (2001) and Febrero et al. (2008) where the
 251 corresponding formulations are explicitly given. A depth-based functional trimmed mean can
 252 be defined as the average over the $y_i(t)$, $i = 1, \dots, n$ that belong to the empirical trimmed region

253 $TR_\alpha:$
$$\bar{y}_\alpha(t) = \frac{1}{|TR_\alpha|} \sum_{y_i \in TR_\alpha} y_i(t), \quad t \in \mathcal{C} \quad (7)$$

254 where $|A|$ is the cardinal of the set A . For functional observations, the median curve is the
 255 deepest function in the sample $\{y_1, \dots, y_n\}$. It maximizes the depth function $D_n(\cdot)$:

$$Y_{median} = \operatorname{argmax}_{x \in \{y_1, \dots, y_n\}} D_n(x) \quad (8)$$

256 where $\operatorname{argmax}_{z \in A} g(z)$ stands for the element in the set A that maximizes the function g .

257 From a theoretical point of view, the mode as a location measure, when it exists, is the value
 258 that locally maximizes the probability density f of the underlying variable. Developments and
 259 applications related to nonparametric density estimation in this context can be found in Dabo-
 260 Niang et al. (2007). An estimator of the modal curve can be obtained by:

$$Y_{mode} = \operatorname{argmax}_{x \in \{y_1, \dots, y_n\}} f_n(x) \quad (9)$$

261 where f_n is an estimate of the density f .

262 The median and mode given respectively in (8) and (9) are natural extensions of their multi-
 263 variate counterparts. However, they are rarely used in practice because of their complex com-
 264 putations. Alternatively, they are commonly defined on the basis of the bivariate scores of a
 265 functional principal component analysis of the curves observations as described in Section 2.4
 266 below.

267 Variability is one of the important quantities to be evaluated and analyzed in statistics. For mul-
 268 tivariate data, the reader is referred to Liu et al. (1999) and Chebana and Ouarda (2011b). The
 269 simplest way to define a variance function in the functional context is by:

$$\operatorname{var}_y(t) = \frac{1}{n-1} \sum_{i=1}^n (y_i(t) - \bar{y}(t))^2, \quad t \in \mathcal{C} \quad (10)$$

270 The covariance function summarizes the dependence structure across different argument values:

271

$$\operatorname{cov}_y(s, t) = \frac{1}{n-1} \sum_{i=1}^n (y_i(s) - \bar{y}(s))(y_i(t) - \bar{y}(t)), \quad s, t \in \mathcal{C} \quad (11)$$

272 The variability of the functional sample is analyzed by plotting the surface $\operatorname{cov}_y(s, t)$ as a func-
 273 tion of s and t as well as the corresponding contour map.

274 Note that, for functional observations, several types of variability can occur such as the
 275 variability within the same observation or between the different observations. In addition, func-
 276 tional principal component analysis is also employed to explore the variability between obser-
 277 vations. The reader is referred to Ramsay and Silverman (2005) and the following sections for

278 a presentation of the functional principal component analysis.

279 **2.3 Functional principal component analysis (FPCA)**

280 Principal component analysis (PCA), as a multivariate procedure, is usually employed to re-
281 duce the dimensionality by defining new variables as linear combinations of the original ones
282 and which capture the maximum of the data variance. After converting the data into functions,
283 functional PCA (FPCA) allows to find new functions that reveal the most important type of
284 variation in the curve data. Note that these new functions cannot be in the Fourier or Spline
285 basis since their aim is not to smooth but rather to produce a reasonable summary of the data by
286 maximizing the capture of the variability. The FPCA method maximizes the sample variance
287 of the scores (defined below) subject to orthonormal constraints. It decomposes the centered
288 functional data in terms of an orthogonal basis as described in the following.

289 Let $y_i(t), i = 1, \dots, n$ be the functional observations obtained by smoothing the observed dis-
290 crete observations $(y_i(t_1), \dots, y_i(t_T)), i = 1, \dots, n$.

291 By definition, the mean curve is a way of variation common to most curves that can be
292 isolated by centering. Let $(y_i^*(t) = y_i(t) - \bar{y}(t))_{i=1, \dots, n}$ be the centered functional observations
293 where $\bar{y}(t)$ is the mean function of $(y_1(t), \dots, y_n(t))$ given by (6). A FPCA is then applied to
294 $(y_i^*(t))_{i=1, \dots, n}$ to create a small set of functions, called also harmonics, that reveals the most
295 important type of variation in the data.

296 The first principal component of $(y_i^*(t))_{i=1, \dots, n}$ denoted by $w_1(t)$ is a function such that the
297 variance of the corresponding real-valued scores $z_{i,1}$ written as:

$$z_{i,1} = \int_{\mathcal{C}} w_1(s) y_i^*(s) ds, \quad i = 1, \dots, n \quad (12)$$

298 is maximized under the constraint $\int_{\mathcal{C}} w_1(s)^2 ds = 1$. The next principal components $w_k(t)$ are
 299 obtained by maximizing the variance of the corresponding scores $z_{i,k}$:

$$z_{i,k} = \int_{\mathcal{C}} w_k(s) y_i^*(s) ds, \quad i = 1, \dots, n \quad (13)$$

300 under the constraints $\int_{\mathcal{C}} w_k(s) w_j(s) ds = 0$, $k \geq 2$, $k \neq j$. As in the multivariate setting, the
 301 interpretation of the principal component function w_k is slightly difficult as it depends on the
 302 type of data being used and may require nonstatistical considerations. A useful way consists
 303 in examining the plots of the overall mean function and perturbations around the mean based
 304 on w_k 's. The perturbation functions are obtained as suitable multiples of the considered w_k ,
 305 namely:

$$\bar{y} \pm 2\sigma_{\omega_k} * \omega_k, \quad k = 1, \dots, K \quad (14)$$

306 where σ_{ω_k} is the square root of the variance (eigenvalue) of the corresponding k th principal
 307 component. This presentation allows to isolate the perturbations about the mean across time
 308 and then assess the variability of the observations. Note that the principal components w_k are
 309 optimal, according to the maximization in (12) or (13), but are not unique. Therefore, any rota-
 310 tion with an orthogonal matrix of the w_k is also optimal and orthonormal. A well-known choice
 311 of such matrices is the VARIMAX. These rotated components can be useful for the interpreta-
 312 tion. More technical details can be found, for instance, in Ramsay and Silverman (2005). On
 313 the other hand, the regularity of the harmonics $w_k(\cdot)$ can be controlled. Rice and Silverman
 314 (1991) and Silverman (1996) extended this traditional functional PCA to the regularized FPCA
 315 (RFPCA) that maximizes the sample variance of the scores subject to penalized constraints.

316 **2.4 Functional data visualization and outlier detection methods**

317 In general, outliers represent gross errors, inconsistencies or unusual observations and should
318 be detected and treated (Barnett and Lewis, 1998). Univariate outliers are well defined and
319 their detection is straightforward (e.g. Hosking and Wallis, 1997; Rao and Hamed, 2000). This
320 topic is also relatively well developed in the multivariate setting (e.g. Dang and Serfling, 2010).
321 The identification and the treatment of outliers constitute an important component of the data
322 analysis before modeling. For hydrologic data, outlier detection is a common problem which
323 has received considerable attention in the univariate framework. In the multivariate setting, the
324 problem is well established in statistics. However, in the hydrologic field the concepts are much
325 less established. A pioneering work in this direction was recently presented by Chebana and
326 Ouarda (2011b). As it is the case in the univariate and multivariate settings, outliers may have
327 a serious effect on the modeling of functional data.

328 In this section, we focus on visualization methods that help to explore and examine certain fea-
329 tures, such as outliers, that might not have been apparent with summary statistics. Different
330 outlier detection methods exist in the functional context literature(e.g. Hardin and Rocke, 2005;
331 Febrero et al., 2007; Filzmoser et al., 2008). However, Hyndman and Shang (2010) showed, on
332 the basis of real data, that their methods are more able to detect outliers and computationally
333 faster. The methods proposed by Hyndman and Shang (2010) are graphical and consist first in
334 visualizing functional data through the rainbow plot, and then in identifying functional outliers
335 using the functional Bagplot and the functional highest density region (HDR) boxplot. The lat-
336 ter two methods can detect outlier curves that may lie outside the range of the majority of the
337 data, or may be within the range of the data but have a very different shape. These methods can

338 also exhibit curves having a combination of these features. In practice, depending on the nature
 339 of the data, the two outlier detection methods can give different results.

340 As pointed out by Jones and Rice (1992) and Sood et al. (2009), the considerable amount
 341 of information contained in the original functional data is captured by the first few principal
 342 components and scores. The outlier identification methods of Hyndman and Shang (2010) con-
 343 sidered here are based on these first two score vectors. Let $y_i(t)$, $w_k(t)$ and $z_{i,k}$ be respectively
 344 the smoothed observed curves, the principal component curves and the corresponding scores
 345 obtained from the FPCA decomposition (Section 2.3). Let $(z_{1,1}, \dots, z_{n,1})$ and $(z_{1,2}, \dots, z_{n,2})$ be
 346 the first two vector scores and $z_i = (z_{i,1}, z_{i,2})$ the bivariate score points. At the end of this
 347 section, a non-graphical outlier detection method is proposed on the basis of $z_i = (z_{i,1}, z_{i,2})$.

348 **2.4.1 Rainbow plot**

349 The rainbow plot, proposed by Hyndman and Shang (2010), is a simple presentation of all the
 350 data, with the only added feature being a color palette based on an ordering. In the functional
 351 context, this ordering is either based on functional depth or data density indices. These indices
 352 are evaluated from the bivariate score depths and kernel density. The bivariate score depth is
 353 given by:

$$OT_i = d(z_i, Z), \quad Z = \{z_j \in \mathbb{R}^2; j = 1, \dots, n\} \quad (15)$$

354 where $d(\cdot, \cdot)$ is the halfspace depth function introduced by Tukey (1975). Tukey's depth function
 355 at z_i is defined as the smallest number of data points included in a closed half-space containing z_i
 356 on its boundary. The observations are decreasingly ordered according to their depth values OT_i .
 357 The first ordered curve represents the *median curve*, while the last curve can be considered as the
 358 *outermost curve* in a sample of curves. As indicated in Section 2.2, this median curve based on

359 Tukey depth function of the bivariate principal scores z_i will be used in the following adaptation
360 to floods. Let θ be this Tukey bivariate depth median defined as $\theta = \arg \max_z d(z, Z)$. If there
361 are several maximizers, the Tukey bivariate depth median can be taken as their center of gravity.
362 The second way of ordering functional observations is based on the kernel density estimate (e.g.
363 Scott, 1992) at the bivariate principal component scores :

$$OD_i = \hat{f}(z_i) = \frac{1}{n} \sum_{j=1}^n \frac{1}{h_j} K\left(\frac{z_i - z_j}{h_j}\right), \quad i \neq j, i = 1, \dots, n \quad (16)$$

364 where $K(\cdot)$ is the kernel function and h_j is the bandwidth for the j^{th} bivariate score points $\{z_j\}$.
365 The functional data $\{y_i(t)\}$ are then ordered in a decreasing order with respect to OD_i . Hence,
366 the first curve with the highest OD is considered as the *modal curve* while the last curve with
367 the lowest OD can be considered as the most *unusual curve*. This modal curve will also be used
368 in the following application.

369 The smoothed observations are presented with colors according to the values of OT and OD .
370 The curves close to the center are red while the most outlying curves are violet.

371 **2.4.2 Functional Bagplot**

372 The bivariate Bagplot is introduced by Rousseeuw et al. (1999) and is based on the halfspace
373 depth function. It is employed by Chebana and Ouarda (2011b) for multivariate hydrological
374 data. The functional Bagplot version is obtained from the bivariate Bagplot based on the first
375 two principal scores $z_i = (z_{i,1}, z_{i,2})$ given in Section 2.3. Each curve in the functional Bagplot is
376 associated with a point in the bivariate Bagplot. Similar to the bivariate Bagplot, the functional
377 Bagplot is composed by three elements: the Tukey median curve, the functional inner region
378 and the functional outer region. The inner region includes 50% of the observations whereas

379 the outer region covers either 95% or 99% of the observations. The outer region is obtained by
 380 inflating the inner region by a factor ϱ . Hyndman and Shang (2010) suggested that the factor ϱ
 381 could take the values 1.96 or 2.58 in order to include respectively 95% or 99% of the curves in
 382 the outer region. These values of ϱ correspond to the case where the bivariate scores follow the
 383 standard normal distribution. Finally, points outside the outer region are considered as outliers.

384 **2.4.3 Functional HDR boxplot**

385 The functional HDR boxplot corresponds to the bivariate HDR boxplot of Hyndman (1996)
 386 applied to the first two principal component scores $z_i \in \mathbb{R}^2$. The bivariate HDR boxplot is
 387 constructed using the bivariate kernel density estimate $\hat{f}(z)$. An HDR with order $\alpha \in (0, 1)$ is
 388 defined as:

$$R_{1-\alpha} = \{z \in \mathbb{R}^2 : \hat{f}(z) \geq f_{1-\alpha}\} \quad (17)$$

389 where $f_{1-\alpha}$ is such that $\int_{R_{1-\alpha}} \hat{f}(t) dt = 1 - \alpha$ and \hat{f} is defined by (16). An HDR can be seen as
 390 a density contour with expanding coverage decreasing with α . The associated bandwidth h_j in
 391 \hat{f} is selected by a smooth cross validation procedure (Duong and Hazelton, 2005).

392 The functional HDR boxplot is composed of the mode defined as $\arg \sup_z \hat{f}(z)$, the 50% inner
 393 region ($R_{50\%}$) and the 99% outer highest density region ($R_{1\%}$). For an HDR with 95% outer
 394 region, one can take $R_{5\%}$ instead of $R_{1\%}$. Curves excluded from the outer functional HDR are
 395 considered as outliers.

396 The difference between detecting outliers by the Bagplot and by the HDR boxplot lies mainly in
 397 the way the inner and outer regions are established. The Bagplot uses a depth function (Tukey)
 398 and the estimated median curve (based on the Tukey depth function of the first bivariate scores
 399 z_i) while the HDR uses the density estimate of the z_i and its mode. Hence, the most outlier

400 curves from HDR are unusual compared to the mode curve whereas those detected by the Bag-
401 plot are unusual with respect to the median curve.

402 In connection with the multivariate setting, as indicated in Chebana and Ouarda (2011b), the
403 points outside the fence of the Bagplot are considered as extremes rather than outliers. Chebana
404 and Ouarda (2011b) considered the approach proposed by Dang and Serfling (2010) to detect
405 outliers. This approach is based on the evaluation of the outlyingness of each observation, the
406 determination of a threshold and then the identification of the observations that exceed this
407 threshold are considered as outliers. The outlyingness values are simple decreasing transforma-
408 tions of depth functions. In the present study, we propose to consider this approach based on the
409 first two scores. A brief presentation of the approach is given in Chebana and Ouarda (2011b),
410 section 2.6.

411 The above graphical approaches should be considered as preliminary indications for suspected
412 observations. The latter could be seen as extremes rather than outliers (see e.g. Chebana and
413 Ouarda, 2011b). In addition, the approach by Dang and Serfling (2010) is based on the outlyin-
414 gness criteria and the corresponding threshold is empirical and not necessarily normally-based
415 (instead of the values of the inflating central region 1.96 or 2.58).

416 **3 Adaptation to floods**

417 The first and most important adaptation for floods lies in the nature of hydrological data. The
418 main data source in hydrology is daily flow from a given station. Flows can also sometimes be
419 available on an hourly, instantaneously or any other time scale. In the following, we focus on
420 daily data and we assume it is recorded during a number n of years of measurements, $\mathbf{Y}_i =$

421 $(y_i(t_1), \dots, y_i(t_T))'$, $i = 1, \dots, n$, $j = 1, \dots, T$, with $T = 365$ days and $y_i(t_j)$ is the flow measured
 422 at the day t_j of the i th year. The time subset index \mathcal{C} is then the interval $[1, 365]$. According
 423 to this kind of data, we have n discrete observations $\{y_i(t_j), j = 1, \dots, 365\}$, $i = 1, \dots, n$.
 424 The observation $\{y_i(t_j), j = 1, \dots, 365\}$ denotes the daily flow for the i th year. A functional
 425 observation constitutes a year starting from January 1st to December 31st. However, it can be
 426 cut out in different ways according to the seasonal characteristics of the geographical area of
 427 interest. For instance, for most parts of Canada, it is possible to define the March-June season
 428 for spring floods and the July-October season for fall floods.
 429 The discrete observed data $(y_i(t_j))_{j=1, \dots, T}$ are to be converted to smooth curves $y_i(t)$ as tempo-
 430 ral functions with a base period of $T = 365$ days and with $p = 52$ weeks non-constant basis
 431 functions as in (2). This smoothing can be obtained through the two well-known Fourier and
 432 B -Spline basis. Usually, the flow data of the whole series present some seasonal variability
 433 and periodicity over the annual cycle. Therefore, Fourier basis are preferred. Although the two
 434 smoothing methods do not give identical results, the differences between them in this adapta-
 435 tion are generally insignificant to affect interpretations. The choice of $p = 52$ can be justified
 436 to capture the flow variation within a week. Since in flood studies, the peak value is important,
 437 in order to ensure that the smooth curves reach the associated peaks, it may be recommended
 438 to consider values of p greater than 52. Nevertheless, this could lead to irregular curves which
 439 could not reasonably capture the entire flow variation.
 440 The nonparametric approach presented in Section 2.4.3, using the kernel density estimate of
 441 z_i 's, is employed for curve ordering and outlier detection and not for estimation purposes. Note
 442 that even though nonparametric approaches have been employed in hydrological FA in the uni-

443 variate context (see e.g. Adamowski and Feluch, 1990; Ouarda et al., 2001), they are still of
444 limited use for the hydraulic design of major structures (Singh and Strupczewski, 2002). In
445 addition, the mode as a location measure is useful to detect the presence of inhomogeneity in
446 the data. In hydrological FA, the mode is not commonly used since, generally, data should pass
447 a homogeneity test. Therefore, the fitted models should be unimodal.

448 Generally, in hydrology, there are two main sources of outliers. The data may be incorrect
449 and/or the circumstances around the measurement may have changed over time (Hosking and
450 Wallis, 1997). However, a detected outlier can also represent true but unusual observed data. In
451 the present functional context, outlier curves have different magnitudes and shapes compared
452 to the rest of the observed curves.

453 **4 Case study**

454 The methods described in Section 2 are applied to hydrological data series by using the adap-
455 tation presented in Section 3. In the following, the data are described, and functional as well
456 as analogous multivariate results are presented and discussed. More precisely, the conversion
457 of the discrete data to be employed as continuous functions is the first preliminary step. Then,
458 the different location functions are obtained and the variability of the sample is studied directly
459 as well as using the FPCA. The latter are also used for data visualization and as a preliminary
460 tool to identify outliers. These outliers are checked by the previously presented approaches
461 and interpreted on the basis of meteorological data. The last subsection provides results using
462 multivariate approaches for comparison purposes.

4.1 Data description and smoothing

The data series is a daily flow (m^3/s) from the Magpie station with reference number 073503. It is located at the outflow of the Magpie lake in the Côte-Nord region in the province of Quebec, Canada. The area of the drainage basin is 7 230 km² and the flow regime is natural. Data are available from 1979 to 2004. Figure 2 indicates the geographical location of the Magpie station.

According to the present dataset, we have $n = 26$ discrete observations $y_i(t_j)$, $t_j \in \mathcal{C} = [1, 365]$, $i = 1, \dots, n$. The i th discrete observation $\{y_i(t_j), j = 1, \dots, 365\}$ denotes the daily flow measurements for the i th year which is converted to a smooth curve $\{y_i(t), t \in \mathcal{C}\}$. This is done through the technique based on Fourier series expansion. This smooth representation of flow data is done with a 365-day base period ($T = 365$ days) and 52-week non-constant basis functions ($p = 52$). The obtained functional observations are given in Table 2 with the corresponding univariate and bivariate samples. This table allows to have an overall view of the data within the three frameworks.

Figure 3a illustrates the whole daily flow series. It shows that the data are nearly periodic. As indicated above, this periodicity can justify the use of Fourier basis. A number of observed hydrographs with the corresponding Fourier and B-splines smoothing curves are presented in Figure 3b. They show that the Fourier and B-Splines smoothing are similar and indicate also that the peaks are generally reached. Figure 4 displays the standard deviation of the residuals $\hat{\epsilon}_{ij} = y_i(t_j) - \hat{x}_i(t_j)$ over j after smoothing the flow data. It gives the residual variations across days, within each year. We observe that these errors are generally very small and do not exceed 32 m^3/s . The highest errors are associated with the years 1981, 1999 and 2002.

485 Note that, other values of p , both smaller and larger than 52, e.g. 4, 12, 90, 122, 182, 300,
486 were also tested. Even though, large values of p , e.g. close to the number of observations per
487 year (here 365), allow to reach almost all the daily flow points including the peaks, the obtained
488 curves are not smooth or regular enough and also do not allow to capture enough of the variance
489 by the first few principal components. Small values of p , e.g. 4,12 give a very bad quality of
490 smoothing, where a large amount of daily flow points are not reached, particularly the high and
491 low values. Therefore, it is reasonable to choose a number p which combines the quality of
492 smoothing (related to (4)) and a high percentage of explained variance by PCA analysis. In the
493 present application, the choice $p = 52$ fits reasonably the discrete data except for some extreme
494 points corresponding to a number of years (e.g. 1980, 1989 and 1993) where the resulting
495 differences between the real peaks and the smooth ones are less than $150 \text{ m}^3/\text{s}$, see Table 2.

496 **4.2 Functional results**

497 Figure 5 presents the smooth location curves (mean, median and mode). It shows that generally
498 the maximum flow occurs in late April and early May followed by a recession during May and
499 June. This phenomenon is common in Canada where floods are mainly caused by snow melt
500 during the Spring season. On the right tail of the curves, we observe a small flood which oc-
501 curs in the autumn and which is caused generally by liquid precipitations. This kind of flood is
502 exhibited by the mode. In both floods, spring or autumn, we observe that the mode is always
503 higher than the mean and the median. The mean seems to be more regular and can not reach
504 high flow values. Therefore, it is useful to consider all these location measures. These location
505 curves lead to different basin characterization through the whole event rather than just some of
506 its parts or summaries and therefore allow for comprehensive basin comparisons.

507 The bivariate (temporal) variance-covariance surface obtained from (11) as well as the corre-
508 sponding contour are presented in Figure 6. We observe that the main part of the variability
509 occurs in the middle of the year and it is negligible elsewhere. That is, the highest variability
510 occurs approximately between April and late June. This period corresponds approximately to
511 the highest flows. This measure has the advantage of providing information concerning the
512 variance structure and also when it occurs.

513 The principal components are obtained by FPCA on the centered observations $y^*(.)$. The vari-
514 ance rates accounted for by each one of the first four principal components are respectively
515 39.5%, 24.0%, 14.4% and 5.4%. These components account for 83.3% of the total variance of
516 the flow. The centered principal components are presented in Figure 7a. The perturbations of
517 these first four principal components about the mean, as given in (14), are presented in Figure
518 7b.

519 From Figure 7, where the first two principal components accumulate 63.5% of the total variance,
520 one can observe that the station flow is most variable between April and July. This variation
521 dominates the variation occurring between July and the end of the year and which is associated
522 with the third and fourth components, and represents 19.9% of the total variance. This finding
523 is, for all practical purposes, consistent with the one obtained from the variance-covariance sur-
524 face (Figure 6). More precisely, the first two principal components w_1 and w_2 are representative
525 of the spring floods whereas w_3 and w_4 are more likely to represent autumn floods.

526 The scores corresponding to the first four principal components are given in Table 3. Given
527 the high variation rate captured by the first principal components, the corresponding variation
528 indicates that the years for which the first or the second principal score is higher (resp. lower)

529 have higher (resp. lower) flow during April to July. Therefore, the year 1981 represents the
530 highest variability during this period followed by the year 1999. On the other hand, the small-
531 est variability of the flow during April-July is associated with the year 1987. Other years
532 could be considered also with low flow variability, such as 1986 and 2002. The flow variability
533 associated with the years 1981, 1986, 1987, 1999 and 2002 is unusual where some of the
534 corresponding curves (1981,1987, 1999) are displayed with the location curves in Figure 8.

535 In order to check the above unusual years, the outlier detection methods described in Section 2
536 are employed. Other functional methods are also tested, such as the functional depth method of
537 Febrero et al. (2007) and the Integrated squared error method of Hyndman and Ullah (2007).
538 However, these two methods gave either too many or no outliers. Hence, the corresponding
539 results are omitted.

540 Figure 9 presents the rainbow plots based on the bivariate depth ordering and the density order-
541 ing indices (15) and (16) respectively. The colors indicate the ordering of the curves where the
542 blue curves are the closest to the center. The red and black outlier curves correspond to 1981
543 and 1999 respectively. Results show that both methods lead to a similar ordering especially for
544 the years associated with high or low ordering.

545 The bivariate Bagplot associated with the first two principal scores as well as the corresponding
546 functional Bagplot for both 95% and 99% of probability coverage are presented in Figure 10.
547 We observe that the curve corresponding to the year 1981 is outside the outer bivariate Bagplot
548 region for both 95% and 99% cases. It corresponds to the red curve in the associated functional
549 Bagplot (Figure 10c,d). Hence, this year is considered as an outlier according to Tukey depth,
550 as described in Section 2. However, when considering the 95% Bagplot, the additional outlier

551 curve that is detected is the one corresponding to 1987 as shown in Figure 10b. Note that gener-
552 ally when outliers are relatively near the median, the functional Bagplot is not a good way to
553 detect them (Hyndman and Shang, 2010). Even though it is not the case here, it is also more
554 appropriate to use the functional HDR boxplot.

555 The bivariate HDR and the associated functional HDR boxplots of the smooth flow curves are
556 presented in Figure 11 for both 95% and 99% of probability coverage. The only detected outlier
557 with 99% coverage probability is 1981 which is outside the bivariate HDR outer region. In the
558 present case, we can deduce that the flow corresponding to the year 1981 is the most outlier,
559 has a different magnitude and shape compared to the other curves and is not near the median.
560 Hence, we can conclude that 1981 is an effective outlier according to the HDR Boxplot. When
561 the probability coverage is 95%, another outlier is detected and corresponds to the year 1999 as
562 shown in Figure 11b. This curve is closer to the median than the curve corresponding to 1981
563 (Figure 8), that is why the functional HDR boxplot is more able to detect it as outlier than the
564 functional Bagplot.

565 As discussed in Section 2.4, the HDR boxplot and the Bagplot are graphical outlier detection
566 methods and their thresholds are based on normality. Therefore, the above detected years can be
567 considered as extreme curves and could be outliers. The approach developed by Dang and Ser-
568 fling (2010) is applied on the first two functional principal component scores Z of the dataset.
569 We evaluated Spatial, Mahalanobis and Tukey outlyingness functions for the bivariate score
570 series. The corresponding thresholds are obtained by selecting the ratio of false outliers to 15%
571 and the true number of outliers as 5 (the same choices as in Chebana and Ouarda, 2011b and
572 Section 4.3 below). Hence, the threshold corresponds to the 0.97-quantile of the outlyingness

573 values. Figure 12 presents the detected outliers. We observe that the Tukey outlyingness func-
574 tion detects several years as outliers (including 1981, 1987, 1999 and 2002) whereas the year
575 1981 is detected by the three outlyingness functions. In addition, the year 1987 corresponds
576 to the second largest Spatial and Mahalanobis outlyingness values and its value is very close
577 to 1999 with the Mahalanobis function. Note that Tukey outlyingness is not recommended by
578 Dang and Serfling (2010). Therefore, the year 1981 can be considered as an effective outlier to
579 be checked. The years 1987 and 1999 could be detected by Spatial and Mahalanobis outlying-
580 ness and considering a larger true number of outliers than 5 (with values of 5%, 10% and 20%
581 of the ratio of false outliers, the results remain the same). Note that the above suspected years
582 of 1986 and 2002 can be considered as extremes and not outliers.

583 Even though the curve of 1981 is the only effective outlier, in the following we examine also
584 the years 1987 and 1999 since they are close to the thresholds. We observe from Figure 8 that
585 the curves of 1981, 1987 and 1999 are clearly different from the location curves and from the
586 general shape of curves. Indeed, based on the corresponding hydrographs, the curve of 1981 is
587 characterized by very high peak and volume whereas 1987 seems to correspond to a dry year
588 since the flow was the lowest during the Spring season. The flood corresponding to the year
589 1999 has also a high peak, although lower than the one corresponding to 1981.

590 The detected outliers can be explained on the basis of meteorological data. The following in-
591 terpretations are drawn on the basis of the data available in Environment Canada's Web site
592 http://www.climat.meteo.gc.ca/climateData/canada_f.html. For 1981, which corresponds to the
593 most important flood for this basin, an important amount of snow was accumulated in early
594 Winter (October-November to January) followed by thaw and rain during February-March. For

595 the outlier corresponding to 1987, the comparison with the preceding and following years re-
596 veals that during the fall of 1987 there was much less rain and the temperatures were very cold,
597 whereas the end of Winter was hot. Hence, all the snow melted earlier compared to other years.
598 The curve of 1999 is relatively higher than the location curves and corresponds to an important
599 quantity of snow on the ground with high temperatures in March. In conclusion, the above
600 detected years seem to be actually observed and do not correspond to incorrect measurements
601 or circumstance changes over time. Hence, these observations should be kept and employed
602 for further analysis. However, it is recommended to use robust statistical methods to avoid
603 sensitivity of the obtained results (e.g. modeling and risk evaluation) to outliers.

604 **4.3 Multivariate results**

605 For comparison purposes, a multivariate study based on Chebana and Ouarda (2011b) is carried
606 out on the present dataset. We focus here on the flood peak Q and the flood volume V as they
607 are among the most important and studied flood characteristics (e.g. Yue et al., 1999 and Shiau,
608 2003). The bivariate series (Q, V) , given in the first three columns of Table 4, are obtained from
609 the daily flow series using an automated version of the algorithm of Pacher (2006). Note that the
610 multivariate approaches presented in Chebana and Ouarda (2011b) are mainly based on depth
611 functions. The Tukey depth function is considered in the present section. The corresponding
612 depth values of each bivariate observation are reported in the fourth column of Table 4. The
613 location and scale results are presented in Table 5. Results with other measures (such as the
614 trimmed mean and dispersion) are obtained but not presented due to space limitations and in
615 order to maintain the coherence with the FDA approach.

616 We observe that Q and V of the bivariate median correspond to those of the median curve

617 obtained in the previous section. Indeed, in both multivariate and functional frameworks, the
618 median corresponds to the year 1980. However, the Q and V of the bivariate mean vector are
619 quite different from those resulting from the mean curve. The mean vector is ($Q = 859.15$, V
620 $= 2138.70$) whereas, when using Pacher's (2006) algorithm, the Q and V of mean curve are
621 respectively 673.09 and 2230.46. We observe also that the difference is larger for the peak than
622 for the volume. This result could be explained by the effect of the detected outliers on the mean
623 which is not the case for the median. Note that the outliers do not necessarily have the same
624 impact in the multivariate and the functional frameworks.

625 Figure 13a presents the bivariate (Q, V) -Bagplot where the median, the central and the outer
626 regions are indicated as well as some particular observations (corresponding to years suspected
627 as outliers from Section 4.2). Note that the outer region is obtained by inflating the central
628 region by a factor of 3 instead of 1.96 or 2.58 as in the functional Bagplot (Figures 10a,b). We
629 observe that the shape of the bivariate (Q, V) -Bagplot is not similar to the functional Bagplot
630 and to the HDR boxplot based on the first two functional principal component scores $z_i =$
631 $(z_{i,1}, z_{i,2})$. The values in Tables 3 and 4 as well as the corresponding figures (Figures 10a,b,
632 13a) indicate that the first two functional principal component scores z_i capture the information
633 from the hydrograph in a different way than do (Q, V) . The former are based on an optimization
634 procedure whereas the latter have physical significance. Nevertheless, both ways are useful to
635 understand flood dynamics and should be used in a complementary manner. This finding should
636 be studied more thoroughly in future research by considering a number of case studies.

637 The bivariate (Q, V) -variability is evaluated both in a matrix form (Table 5) and by using scalar
638 curve (Figure 13b). Note that the variability is particularly useful when comparing at least two

639 data sets for the same kind of series (e.g. same variable or same vector). It is appropriate to
640 compare the univariate peak scale with the functional one since the flood peak has the same
641 unit and scale as the daily flow which is not the case for the volume. Hence, we observe that
642 the peak variance has the same magnitude as in the functional context as it can be seen from
643 Table 5 and Figure 6. One can also appropriately standardize the Q and V variables in order to
644 compare the variances of the vector (Q, V) and the functional context.

645 The procedure employed in Chebana and Ouarda (2011b) for outlier detection is based on depth
646 outlyingness measures and the corresponding thresholds. The reader is referred to Chebana and
647 Ouarda (2011b) or Dang and Serfling (2010) for more details about the outlyingness expressions
648 and threshold determination. In the present section, three outlyingness measures are evaluated
649 on the (Q, V) series, i.e. Tukey (TO), Mahalanobis (MO) and Spatial (SO). Their values are
650 presented in the last three columns of Table 4. To obtain the threshold that the outlyingness
651 of an outlier should exceed, we considered a ratio of false outliers of 15% among the allowed
652 ones and we also allowed 5 true outliers (the same choices as in Chebana and Ouarda, 2011b).
653 Hence, the threshold corresponds to the empirical 97%-quantile of the outlyingness values. The
654 obtained threshold values are 0.9231, 0.8676 and 0.9462 respectively for TO, MO and SO.
655 Consequently, 1981 is detected by all measures, 1987 is detected only by TO and it has also the
656 second highest outlying value by MO and SO but smaller than the corresponding thresholds.
657 The measure TO detects several other outliers, including 1999 and 2002, which all have the
658 same TO value (equal to the threshold). However, if a quantile of order higher than 97% is
659 considered, by modifying the parameters related to the threshold, the TO does not detect any
660 outliers. These results are consistent with those of the functional framework in the sense that the

661 most unusual observations are detected in both frameworks. However, the proposed approach
662 that consists in applying the Dang and Serfling (2010) technique on the first two score series z_i
663 seems to be justified and more reliable.

664 **5 Summary and concluding remarks**

665 The first aim of the present paper is to introduce the functional framework to hydrological
666 applications based on the curve nature of the data to be employed and analyzed. The FDA
667 framework can be seen as a natural extension of the multivariate FA where the latter is gaining
668 popularity and usefulness in meteorological and hydrological studies. In the present study we
669 introduced a number of FDA notions and techniques and adapted them to the hydrological
670 context, and more specifically to floods. The techniques within the first functional FA step deal
671 with visualization, location estimation, variability quantification, principal component analysis
672 and outlier detection. A new non-graphical (numerical) outlier detection method is proposed
673 which combines multivariate and functional techniques.

674 An application is carried out to demonstrate the potential of employing FDA techniques in
675 hydrology. The application deals with the natural streamflow series of the Magpie station in the
676 province of Quebec, Canada. Results regarding location measures such as mean, median and
677 modal curves, are obtained. The variability is studied as a simple bivariate function surface and
678 also by using principal component analysis. Outlier curves are identified by the most efficient
679 methods and interpretations are given based on meteorological data. For comparison purposes, a
680 brief bivariate study of flood peak and volume is carried out. Even though FDA is an extension
681 of multivariate analysis, it is recommended to perform both approaches to obtain a complete

682 understanding of floods and to make the appropriate decisions.

683 From the elements discussed in the introduction and the results obtained in the case study, the
684 following concluding remarks can be drawn and a number of limitations and perspectives are
685 given:

686 **I) Drawbacks of previous approaches:** The following drawbacks represent the motivation and
687 the need to introduce the functional framework in hydrological applications:

- 688 1. The separate or joint use of hydrograph characteristics constitutes a major simplification
689 of the real phenomenon;
- 690 2. Given the lack of data in hydrology, working with a limited number of extracted charac-
691 teristics represents a loss of a part of the available information;
- 692 3. The way these characteristics are determined is neither unique nor objective;
- 693 4. The multivariate analysis is a simplification of the hydrological phenomena since it is
694 based on flood characteristics which are simple transformations of the hydrograph;
- 695 5. In the multivariate setting, the complexity of the models, the fitting and estimation diffi-
696 culty, the number of parameters and the associated uncertainty increase with the dimen-
697 sion;
- 698 6. The importance of the hydrograph shape is shown in studies such as Yue et al. (2002)
699 where the approaches approximating the flood hydrograph using probability densities are
700 limited for instance to single-peak hydrographs;
- 701 7. The main flood characteristics, peak, volume and duration, can not completely capture
702 the shape of the hydrograph;

703 8. Even though, in the flow duration curve studies, e.g. Castellarin et al. (2004), a number
704 of functional elements, such as mean and median curves, are presented, they are limited
705 and do not have a functional statistical foundation;

706 **II) Conceptual advantages of the functional framework:** The functional framework presents
707 some general advantages which contribute to overcome the previous drawbacks at different
708 levels:

709 1. The functional framework treats the whole hydrograph as a functional observation (func-
710 tion or curve) which is more representative of the real phenomena;

711 2. It employs the maximum of the available information in the data where the impact of the
712 lack of data in hydrology can be reduced;

713 3. The functional framework is more general and more flexible and can represent a large
714 variety of hydrographs;

715 4. The functional methodology constitutes a natural extension of univariate and multivariate
716 hydrological FA approaches;

717 5. The location curves and functional scale measures can be used to characterize a given
718 basin and to proceed to comparison or grouping of a set of basins;

719 6. FDA allows to perform a single analysis on the whole data instead of several univariate
720 or multivariate analysis;

721 7. The approaches dealing with hydrograph shape, e.g. the one proposed by Yue et al.
722 (2002), can be generalized in the FDA context using smoothing techniques;

723 8. The functional setting avoids the definition and the evaluation of flood characteristics.
724 Therefore, it does not require specific algorithms and avoids subjective evaluations; and
725 the associated uncertainty can be reduced in the analysis;

726 **III) Concluding remarks from the application:** The following points are drawn as specific
727 results of the FDA application to the case study:

728 1. The location curves (mean, median and mode) give more information concerning the hy-
729 drological regime in the basin than the univariate and multivariate approaches by adding
730 temporal aspects. These curves allow to summarize the information contained in the data
731 for a given basin, and hence make comparisons between basins and group basins with
732 similar regime;

733 2. The bivariate (temporal) variance-covariance surface as well as the corresponding contour
734 give an additional insight to the hydrological regime variability than the real-value or
735 matrix in the univariate and multivariate contexts;

736 3. In addition to quantifying the variability, functional scale measures indicate when it oc-
737 curs;

738 4. The case study results show that the mode is useful to characterize high flood values,
739 the variability is very high during spring season and the principal components are shown
740 to describe the variability in spring floods and autumn floods. The detected outliers are
741 checked to be real observations and therefore it is suggested to use robust methods in any
742 further analysis;

743 5. The first two functional principal components capture the information from the hydro-

744 graph in a different way than do (Q, V) . Nevertheless, both ways are useful to understand
745 flood dynamics and should be used in a complementary manner;

746 6. The FPCAs represent a new way to distinguish the different flood events in a given year.
747 Indeed, the few first principal components can be used to identify where in the hydro-
748 graph the variation dominates and can be used to characterize flood events, e.g the first
749 two principal components are representative of the spring floods whereas the two others
750 represent autumn floods;

751 7. In the functional context, outlier curves have different magnitudes and shapes compared
752 to the rest of the observed curves. In the univariate and multivariate settings, the shape is
753 not considered and can not be captured even by using several variables;

754 8. The functional results obtained in this study are generally coherent with those of the
755 multivariate analysis but give more insight to the hydrological phenomena such as in
756 terms of location measures, variability and principal components;

757 **IV) Limitations and perspectives of the functional framework:** The present study presented
758 exploratory functional tools that are important on their own and it constitutes also a basis for
759 the next steps for a reliable FDA-based hydrological FA, especially in terms of model selec-
760 tion and risk evaluation. Several perspectives are promising and can be carried out in future
761 research:

762 1. Although the study focused on floods, the presented FDA methodology can be adapted
763 and applied to treat other hydro-meteorological variables such as droughts, precipitations,
764 storms and heat waves;

- 765 2. FDA relies on the smoothing step. Therefore, a careful inspection of the resulting curves
766 is recommended, for instance, to ensure the regularity of the smoothed functions, to reach
767 a majority or special points such as peaks or to capture enough of the variance by the
768 first few principal components. Even though a number of elements in this direction are
769 given in the present study, it could be of interest to develop general criteria and objective
770 choices depending on the objective of the analysis;
- 771 3. The classification of the curves of a given site as well as the clustering of sites in a region
772 on the basis of the full hydrograph are also topics of interest;
- 773 4. Inferential aspects, such as modeling for prediction purposes, represent also important
774 issues for future research efforts;
- 775 5. Future investigations should also deal with hypothesis testing as well as regression mod-
776 eling.

777 **Acknowledgement** Financial support for this study was graciously provided by the Natural Sci-
778 ences and Engineering Research Council (NSERC) of Canada, and the Canada Research Chair
779 Program. The authors would also like to thank the authors of the *FDA R* and *Rainbow R* pack-
780 ages. The authors wish to thank the Editor, Associate Editor and three anonymous reviewers
781 for their useful comments which led to considerable improvements in the paper.

782 **References**

- 783 Adamowski, K. and Feluch, W. (1990), Nonparametric flood frequency analysis with historical
784 information. *ASCE Journal of Hydraulic Engineering*, 116, 1035-1047.
- 785 Barnett, V., and T. Lewis (1998), *Outliers in statistical data*, 3rd ed., Wiley, Chichester.

- 786 Barnett, V. (2004), *Environmental statistics : methods and applications*, Wiley, Chichester,
787 England ; Hoboken, N.J.
- 788 Castellarin, A., Vogel, R. and Brath, A. (2004), A stochastic index flow model of flow duration
789 curves. *Water Resources Research*, 40, W03104.
- 790 Chebana, F. and Ouarda, T.B.M.J. (2007), Multivariate L-moment homogeneity test. *Water*
791 *Resources Research*, 43, W08406. doi:10.1029/2006WR005639.
- 792 Chebana, F. and Ouarda, T.B.M.J. (2008), Depth and homogeneity in regional flood frequency
793 analysis. *Water Resources Research*, 44, W11422. doi : 10.1029/2007WR006771
- 794 Chebana, F., Ouarda, T.B.M.J. and Duong, T.C. (2010), Testing for stationarity, homogeneity
795 and trend in multivariate hydrologic time series: a review, *INRS-ETE, report I-271*: Quebec.
796 p. 73.
- 797 Chebana, F. and Ouarda T.B.M.J. (2011a), Multivariate quantiles in hydrological frequency
798 analysis. *Environmetrics*, 22, 63-78, DOI: 10.1002/env.1027
- 799 Chebana, F. and Ouarda, T.B.M.J. (2011b), Depth-based multivariate descriptive statistics with
800 hydrological applications. *J. Geophys. Res.*, 116, D10120, doi:10.1029/2010JD015338.
- 801 Chow, V. T., Maidment, D. R. and L. R. Mays (1988), *Applied Hydrology*, McGraw-Hill, New
802 York.
- 803 Cuevas, A. and Febrero, M. and Fraiman, R. (2006), On the use of bootstrap for estimating
804 functions with functional data, *Computational Statistics and Data Analysis*, 51, 1063-1074.
- 805 Cunnane C. (1987), Review of statistical models for flood frequency estimation. In *Hydrologic*
806 *Frequency Modeling*, Singh V.P. (ed.). Reidel, Dordrecht: The Netherlands; 49-95.
- 807 Dabo-Niang, S., F. Ferraty and P. Vieu (2007), On the using of modal curves for radar wave-
808 forms classification. *Computational statistics and data analysis*, 51, 4878-4890.
- 809 Dabo-Niang, S. and F. Ferraty (2008), *Functional and Operatorial Statistics*, Springer Verlag,
810 New-York.
- 811 Dabo-Niang, S.; Yao, A.F., Pischedda, L.; Cuny, P. and Gilbert, F. (2010), Spatial mode esti-
812 mation for functional random fields with application to bioturbation problem. *Stochastic En-*
813 *vironmental Research and Risk Assessment*, 24 (4), 487-497.
- 814 Dang, X. and Serfling, R. (2010), Nonparametric depth-based multivariate outlier identifiers,
815 and masking robustness properties. *Journal of Statistical Planning and Inference*, 140, 198-
816 213.
- 817 de Boor, C. (2001). *A Practical Guide to Splines*, Revised Edition. New York, Springer.
- 818 Delaigle, A. and Hall, P. (2010), Defining probability density for a distribution of random
819 functions. *Annals of Statistics*, 38, 1171-1193.

- 820 Delicado, P., Giraldo, R. and Mateu, J. (2008), Point-wise Kriging for spatial prediction of
821 functional data. In *Functional and Operatorial Statistics*. Editors S. Dabo-Niang and F. Ferraty,
822 Physica-Verlag, 135-141.
- 823 Duong, T. and Hazelton, M. L. (2005), Cross-validation bandwidth matrices for multivariate
824 kernel density estimation. *Scandinavian Journal of Statistics*, 32, 485-506.
- 825 Febrero, M., Galeano, P. and Gonzalez-Manteiga, W. (2007), A functional analysis of NOx
826 levels: location and scale estimation and outlier detection, *Computational Statistics*, 22, 411-
827 427.
- 828 Febrero, M., Galeano, P. and Gonzalez-Manteiga, W. (2008), Outlier detection in functional
829 data by depth measures, with application to identify abnormal NOx levels, *Environmetrics*, 19,
830 331- 345.
- 831 Ferraty, F. and Vieu, P. (2006), *Nonparametric functional data analysis*. Springer-Verlag, New
832 York.
- 833 Ferraty, F.; Rabhi, A. and Vieu, P. (2005), Conditional quantiles for functional dependent
834 data with application to the climatic El-NINO phenomenon. *Sankhya: The Indian Journal of*
835 *Statistics*, 67, 378-398.
- 836 Filzmoser, P., Maronna, R. and Werner, M. (2008), Outlier identification in high dimensions.
837 *Computational Statistics and Data Analysis*, 52, 1694-1711.
- 838 Fraiman, R. and Muniz, G. (2001), Trimmed means for functional data. *Test*, 10, 419-440.
- 839 Hardin, J. and Rocke, D. M. (2005), The distribution of robust distances. *Journal of Computa-*
840 *tional and Graphical Statistics*, 14(4), 928-946.
- 841 Hosking, J. R. M. and Wallis, J. R. (1997), *Regional Frequency Analysis: An Approach Based*
842 *on L- Moments*. Cambridge University Press.
- 843 Hyndman, R. J. (1996), Computing and graphing highest density regions. *The American Statis-*
844 *tician*, 50, 120-126.
- 845 Hyndman, R. J. and Ullah, M. S. (2007), Robust forecasting of mortality and fertility rates: A
846 functional data approach. *Computational Statistics and Data Analysis*, 51, 4942-4956.
- 847 Hyndman, R. J. and Shang, H.L. (2010), Rainbow plots, bagplots and boxplots for functional
848 data. *Journal of Computational and Graphical Statistics*, 19, 29-45..
- 849 Jones, M. C. and Rice, J. A. (1992), Displaying the important features of large collections of
850 similar curves. *The American Statistician*, 46, 140-145.
- 851 Khaliq, M.N., Ouarda, T.B.M.J., Gachon, P., Sushama, L. and St-Hilaire, A. (2009), Identi-
852 fication of hydrological trends in the presence of serial and cross correlations: A review of
853 selected methods and their application to annual flow regimes of Canadian rivers. *Journal of*
854 *Hydrology*, 368, 117-130.

- 855 Kite, G.W. (1988), *Frequency and Risk Analysis in Hydrology*. Water Resources Publications,
856 Littleton, Colorado
- 857 Koulis, T.; Thompson, M.E. and LeDrew, E. (2009), A spatio-temporal model for Antarctic
858 sea ice formation. *Environmetrics.*, 20, 68-85.
- 859 Kundzewicz, Z.W., Graczyk, D., Maurer, T., Pinskiwar, I., Radziejewsky, M., Svensson, C.,
860 and Szwed, M. (2005), Trend detection in river flow series: 1. Annual maximum flow. *Hydrol.*
861 *Sci. J.* 50, 797-810.
- 862 Liu, R. Y., Parelius, J. M. and Singh, K. (1999), Multivariate analysis by data depth: Descrip-
863 tive statistics, graphics and inference. *The Annals of Statistic*, 27, 783-858.
- 864 Mizuta, M. (2006), Discrete functional data analysis. In: *Proceedings in Computational Statis-*
865 *tics*. Physica-Verlag, Springer, Berlin, 361-369.
- 866 Monestiez, P. and Nerini (2008), A Cokriging method for spatial functional data with appli-
867 cations in oceanology. In *Functional and Operatorial Statistics*. Editors S. Dabo-Niang and F.
868 Ferraty, Physica-Verlag, 237-242.
- 869 Ouarda, T. B. M. J. and Ashkar, F. (1998), Effect of trimming on LP III flood quantile esti-
870 mates. *Journal of Hydrologic Engineering*, 3, 33-42.
- 871 Ouarda, T. B. M. J.; Girard, C.; Cavadias, G. S. and Bobée, B. (2001), Regional flood fre-
872 quency estimation with canonical correlation analysis. *Journal of Hydrology*, 254, 157-173.
- 873 Pacher, G. W. (2006), Détermination objective des paramètres des hydrogrammes [In French].
874 Personal note.
- 875 Ramsay, J. O. and Silverman, B. W. (2002), *Applied Functional Data Analysis: methods and*
876 *case studies*. Springer, New York, London.
- 877 Ramsay, J. O. and Silverman, B. W. (2005), *Functional Data Analysis*, 2nd edn, Springer, New
878 York.
- 879 Rao, A. R. and Hamed, K.H. (2000), *Flood Frequency Analysis*. CRC Press: Boca Raton.
- 880 Rice, J.A. and Silvermann, B.W. (1991), Estimating the mean and covariance structure non-
881 parametrically when data are curves. *Journal of the Royal Statistical Society B*, 53, 233-243.
- 882 Rousseeuw, P. J.; Ruts, I. and Tukey, J. W. (1999), The bagplot: A bivariate boxplot. *The*
883 *American Statistician*, 53, 382-387.
- 884 Salvadori, G., De Michele, C., Kottegoda, N.T. and Rosso, R. (2007), *Extremes in Nature: An*
885 *Approach Using Copulas*. Springer.
- 886 Scott, D.W. (1992), *Multivariate density estimation: theory, practice, and visualization*, Wiley,
887 New York.

- 888 Serinaldi, F. and Grimaldi, S. (2007), Fully nested 3-copula: Procedure and application on
889 hydrological data. *Journal of Hydrologic Engineering*, 12, 420-430.
- 890 Shiau, J. T. (2003), Return period of bivariate distributed extreme hydrological events. *Stochas-
891 tic Environmental Research and Risk Assessment*, 17, 42-57.
- 892 Silvermann, B.W. (1996), Smoothed functional principal components analysis by choice of
893 norm. *The Annals of Statistics*, 24, 1-24.
- 894 Singh, V.P. and Strupczewski, W.G. (2002), On the status of flood frequency analysis. *Hydro-
895 logical Processes*, 16, 3737-3740.
- 896 Sood, A., James, G. M. and Tellis, G. J. (2009), Functional regression: a new model for
897 predicting market penetration of new products, *Marketing Science*, 28, 36-51.
- 898 Stedinger, J.R., Vogel, R.M. and Foufoula-Georgiou, E. (1993), *Frequency Analysis of Extreme
899 Events*. In: Handbook of Hydrology, D.R. Maidment (Editor). McGraw-Hill Inc., New York,
900 New York.
- 901 Tukey, J. W. (1975), Mathematics and the picturing of data, in Proceedings of the International
902 Congress of Mathematicians (Vancouver, B. C., 1974), 2, 523-531.
- 903 Wahba, G. (1990), *Splines Models for Observational data*, Philadelphia, SIAM.
- 904 Warner, R. M. (2008), *Applied statistics : from bivariate through multivariate techniques*,
905 SAGE Publications, Thousand Oaks, Calif.
- 906 Yue, S.; Ouarda, T. B. M. J.; Bobée, B.; Legendre, P. and Bruneau, P. (1999), The Gumbel
907 mixed model for flood frequency analysis. *Journal of Hydrology*, 226, 88-100.
- 908 Yue, S., Ouarda, T. B. M. J., Bobée, B., Legendre, P. and Bruneau, P. (2002), Approach for
909 describing statistical properties of flood hydrograph, *Journal of Hydrologic Engineering*, 7,
910 147.
- 911 Zhang, L., and Singh V.P. (2006), Bivariate flood frequency analysis using the copula method.
912 *Journal of Hydrologic Engineering*, 11, 150-164.
- 913 Zhang, L. and Singh, V.P. (2007), Trivariate flood frequency analysis using the Gumbel-
914 Hougaard copula. *Journal of Hydrologic Engineering*, 12, 431-439.
- 915 Zuo, Y. and Serfling, R. (2000), General notions of statistical depth function. *The Annals of
916 Statistics*, 28, 461-482.

FA Steps	Framework		
	Univariate	Multivariate	Functional
i) Exploratory analysis & outlier detection	Large literature : Cunnane 1987 Kite 1988 Stedinger et al 1993 Rao & Hamed 2000	Very sparse literature : Chebana & Ouarda 2011b	The specific aim of the present paper
ii) Checking the FA assumptions: stationarity homogeneity independence	Large literature: Yue et al 2002 Kundzewicz et al 2005 Khaliq et al 2009	Very sparse literature : Chebana et al 2010	To be developed
iii) Modeling & estimation	Large literature : Cunnane 1987 Bobée & Ashkar 1991	Large recent literature: Shiau 2003 Zhang & Singh 2006 Salvadori et al 2007	To be developed
iv) Risk evaluation & analysis	Large literature : Chow et al. 1988	Little but growing literature : Shiau 2003 Chebana & Ouarda 2011a	To be developed

Table 1: FA steps in the three frameworks.

Note: in the univariate framework, step (i) is straitforward and is generally not treated separately;

The references are given only as examples from the literature for space limitation.

Year	z_1	z_2	z_3	z_4
1979	-1180.34	1174.70	1457.11	373.28
1980	42.34	329.59	-115.61	121.06
1981	2613.26	2046.54	-448.91	225.00
1982	1947.70	-704.01	500.88	-600.22
1983	-1673.67	1095.67	1936.68	-133.87
1984	843.83	822.51	-10.20	-238.71
1985	903.82	-1171.34	0.07	-419.37
1986	-1737.16	-185.44	466.40	36.39
1987	-1671.02	-1615.92	285.93	10.93
1988	-130.59	393.79	-732.23	18.73
1989	-633.73	-176.40	-663.42	87.96
1990	-669.66	-519.85	-375.08	-324.09
1991	529.43	-604.15	-281.41	-52.87
1992	-465.75	-725.79	-342.08	944.50
1993	-374.77	-753.53	-315.64	19.32
1994	1058.79	68.18	451.04	1281.09
1995	-268.43	463.72	-933.12	-268.28
1996	748.05	235.06	560.82	-575.86
1997	1085.56	-516.08	447.61	482.64
1998	-1557.15	428.41	-504.03	-165.38
1999	-1306.02	1809.15	-621.20	-641.45
2000	879.38	-20.22	145.23	-178.76
2001	-1173.90	-702.48	-837.66	133.55
2002	1134.07	-1692.79	796.87	-433.89
2003	-67.68	120.028	-1087.29	315.00
2004	1123.65	400.67	219.25	-16.70

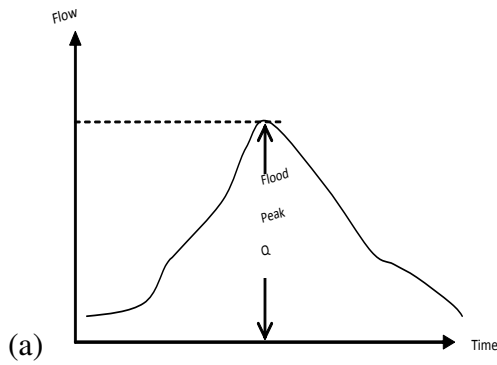
Table 2: First four principal component scores. The bold characters indicate the largest and the smallest values for the first and the second component.

Year	Peak	Volume	TD	MO	SO	TO
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

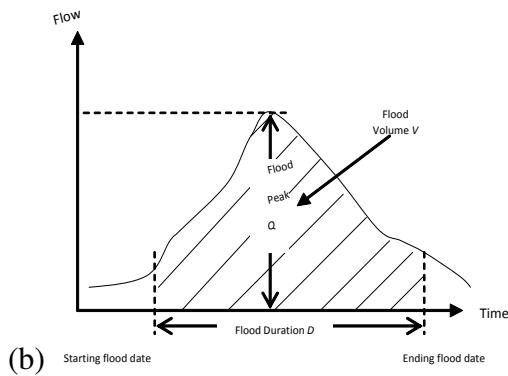
Table 3: Multivariate results for flood peak and volume. TD: Tukey Depth, MO: Mahalanobis Outlyingness, SO: Spatial Outlyingness, and TO: Tukey Outlyingness. Bold characters indicate the values of the outlying measure corresponding to the detected outlier.

	Peak	Volume
Mean (vector)	859.15	2138.70
Tukey median (vector)	847.72	2216.22
Dispersion (matrix)	57316.61	113915.10
	113915.10	457040.80

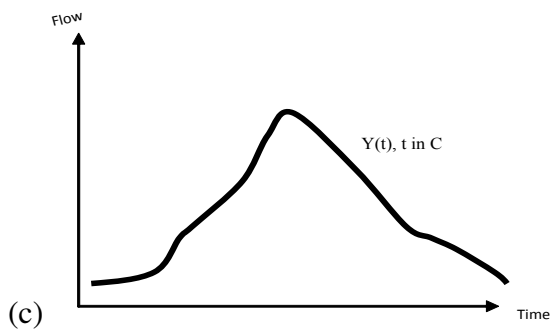
Table 4: Multivariate results for flood peak and volume: location and scale parameters.



Framework: univariate
 Variable : peak Q
 Nature : real value
 Series : q_1, \dots, q_n
 Duration: n years
 References: Rao and Hamad (2000)
 Cunnane (1987)



Framework: multivariate
 Variables: peak Q , volume V and duration D
 (the most important and most studied)
 Nature: vector
 Series : e.g. $(q_1, v_1, d_1), \dots, (q_n, v_n, d_n)$
 Duration: n years
 References: Yue et al. (1999)
 Shiau (2003)
 Zhang and Singh (2006)
 Chebana and Ouarda (2011a)



Framework: functional
 Variable: $Y(t), t \in \mathcal{C}$,
 the whole hydrograph
 Nature: function
 Series: $(Y_1(t), \dots, Y_n(t)), t \in \mathcal{C}$
 Duration: n years
 References: The object of the present paper

Figure 1: Illustration of the different approaches (a) univariate (b) multivariate and (c) functional with the corresponding types of variables, series and a number of references.

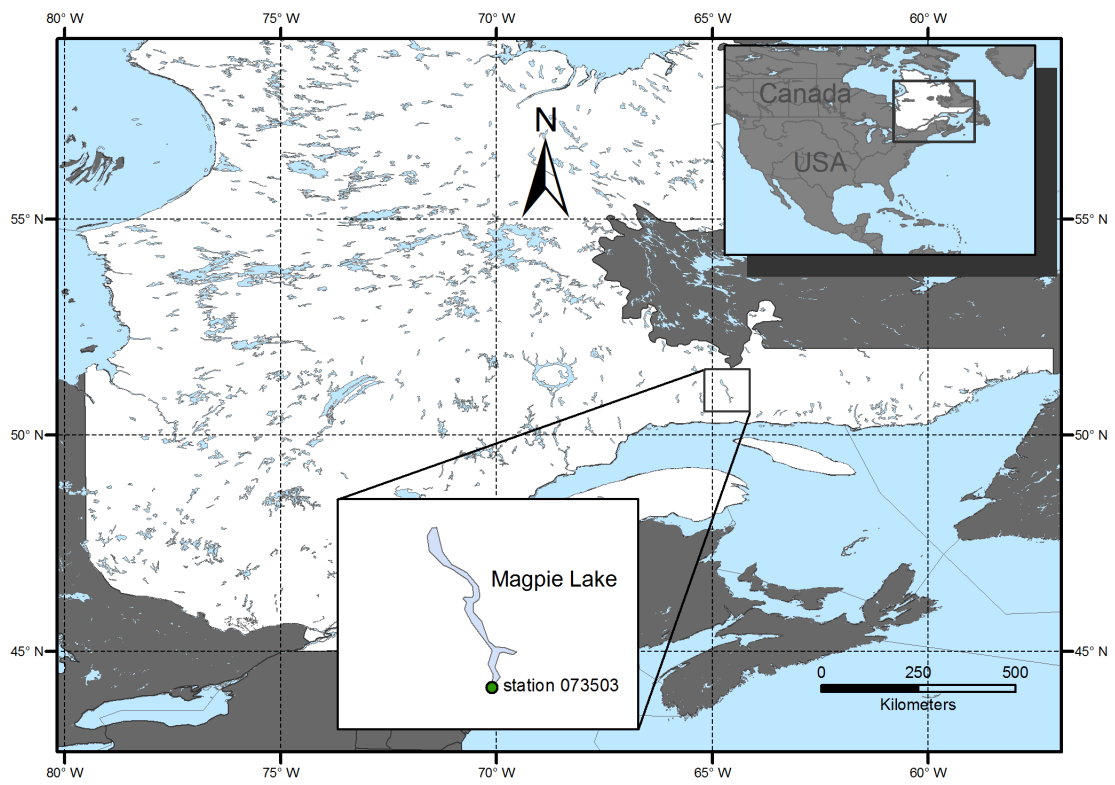


Figure 2: Geographical location of the Magpie station.

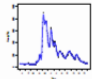
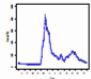
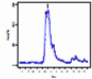
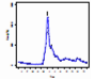
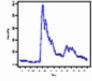
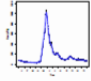
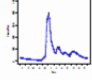
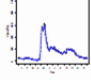
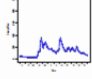
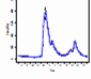
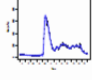
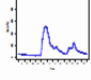
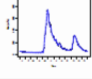
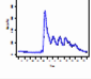
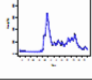
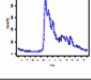
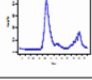
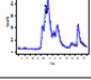
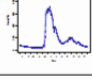
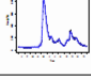
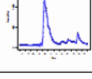
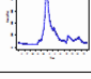
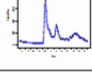
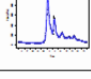
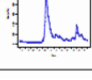
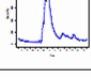
Year	Univ. Q	Biv. $\begin{pmatrix} Q \\ V \end{pmatrix}$	Func. Hydrograph	Year	Univ. Q	Biv. $\begin{pmatrix} Q \\ V \end{pmatrix}$	Func. Hydrograph
1979	886.7	$\begin{pmatrix} 886.7 \\ 2088.9 \end{pmatrix}$		1992	710.8	$\begin{pmatrix} 710.8 \\ 1128.7 \end{pmatrix}$	
1980	849.7	$\begin{pmatrix} 849.7 \\ 2357.0 \end{pmatrix}$		1993	666.8	$\begin{pmatrix} 666.8 \\ 1407.3 \end{pmatrix}$	
1981	1456.7	$\begin{pmatrix} 1456.7 \\ 3909.1 \end{pmatrix}$		1994	932.9	$\begin{pmatrix} 932.9 \\ 2722.5 \end{pmatrix}$	
1982	1270.0	$\begin{pmatrix} 1270.0 \\ 2443.1 \end{pmatrix}$		1995	868.8	$\begin{pmatrix} 868.8 \\ 2192.4 \end{pmatrix}$	
1983	974.7	$\begin{pmatrix} 974.7 \\ 3012.2 \end{pmatrix}$		1996	886.9	$\begin{pmatrix} 886.9 \\ 2476.4 \end{pmatrix}$	
1984	1056.7	$\begin{pmatrix} 1056.7 \\ 2751.7 \end{pmatrix}$		1997	697.3	$\begin{pmatrix} 697.3 \\ 2665.9 \end{pmatrix}$	
1985	787.0	$\begin{pmatrix} 787.0 \\ 1574.2 \end{pmatrix}$		1998	825.0	$\begin{pmatrix} 825.0 \\ 1843.6 \end{pmatrix}$	
1986	610.3	$\begin{pmatrix} 610.3 \\ 1536.3 \end{pmatrix}$		1999	1306.7	$\begin{pmatrix} 1306.7 \\ 2652.3 \end{pmatrix}$	
1987	344.3	$\begin{pmatrix} 344.3 \\ 1069.9 \end{pmatrix}$		2000	858.9	$\begin{pmatrix} 858.9 \\ 2492.6 \end{pmatrix}$	
1988	843.3	$\begin{pmatrix} 843.3 \\ 2374.5 \end{pmatrix}$		2001	732.5	$\begin{pmatrix} 732.5 \\ 1188.9 \end{pmatrix}$	
1989	678.7	$\begin{pmatrix} 678.7 \\ 1534.5 \end{pmatrix}$		2002	999.6	$\begin{pmatrix} 999.6 \\ 1485.4 \end{pmatrix}$	
1990	506.3	$\begin{pmatrix} 506.3 \\ 1752.1 \end{pmatrix}$		2003	1004.9	$\begin{pmatrix} 1004.9 \\ 1883.8 \end{pmatrix}$	
1991	740.0	$\begin{pmatrix} 740.0 \\ 2260.6 \end{pmatrix}$		2004	842.6	$\begin{pmatrix} 842.6 \\ 2802.3 \end{pmatrix}$	

Figure 3: Data for each one of the three frameworks: univariate, bivariate and functional.

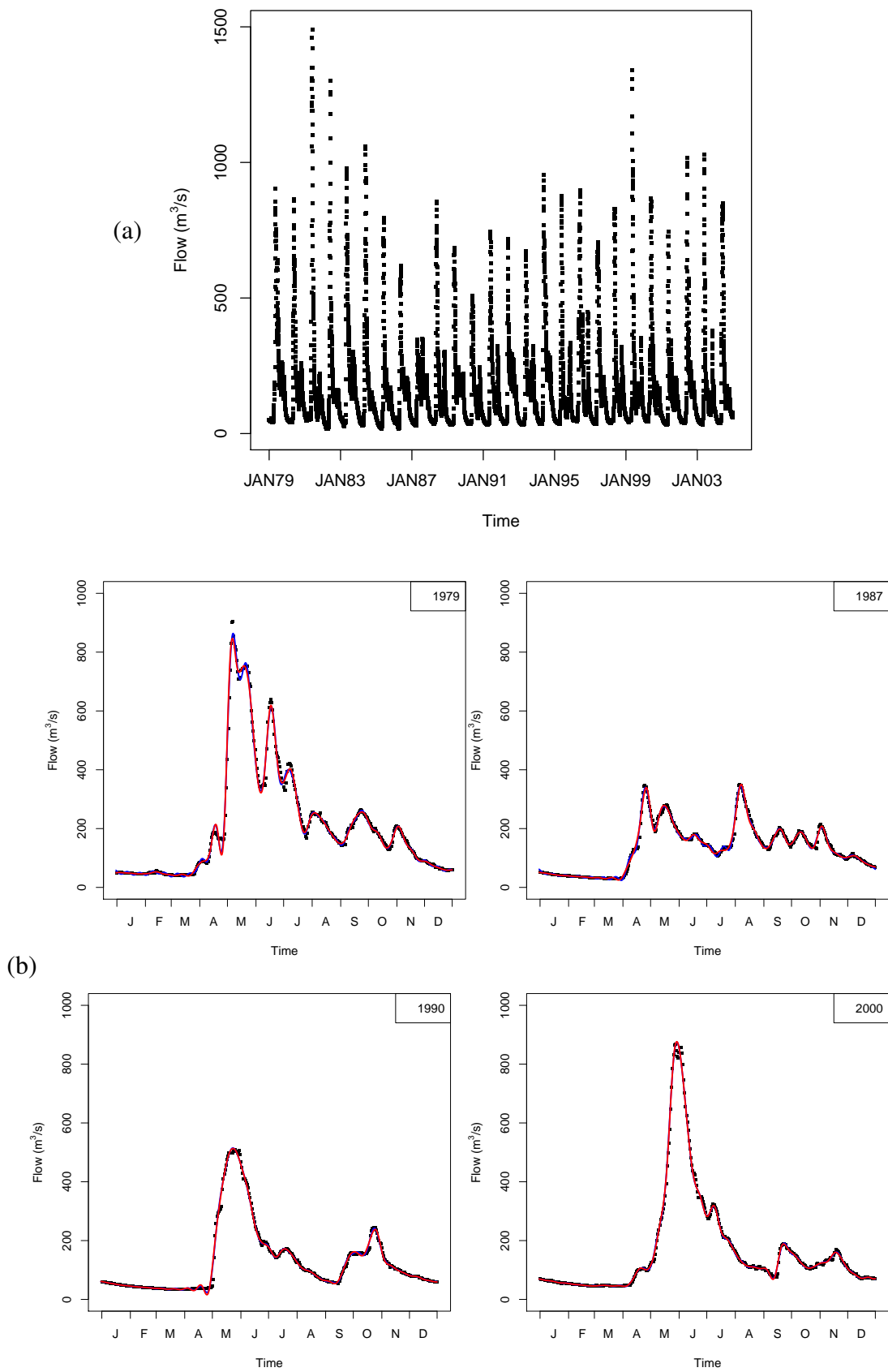


Figure 4: The representation of all the data (in (a)) and illustration of discrete hydrographs and the corresponding smoothing curves (Fourier in blue and B-Splines in red) for some selected years (in (b))

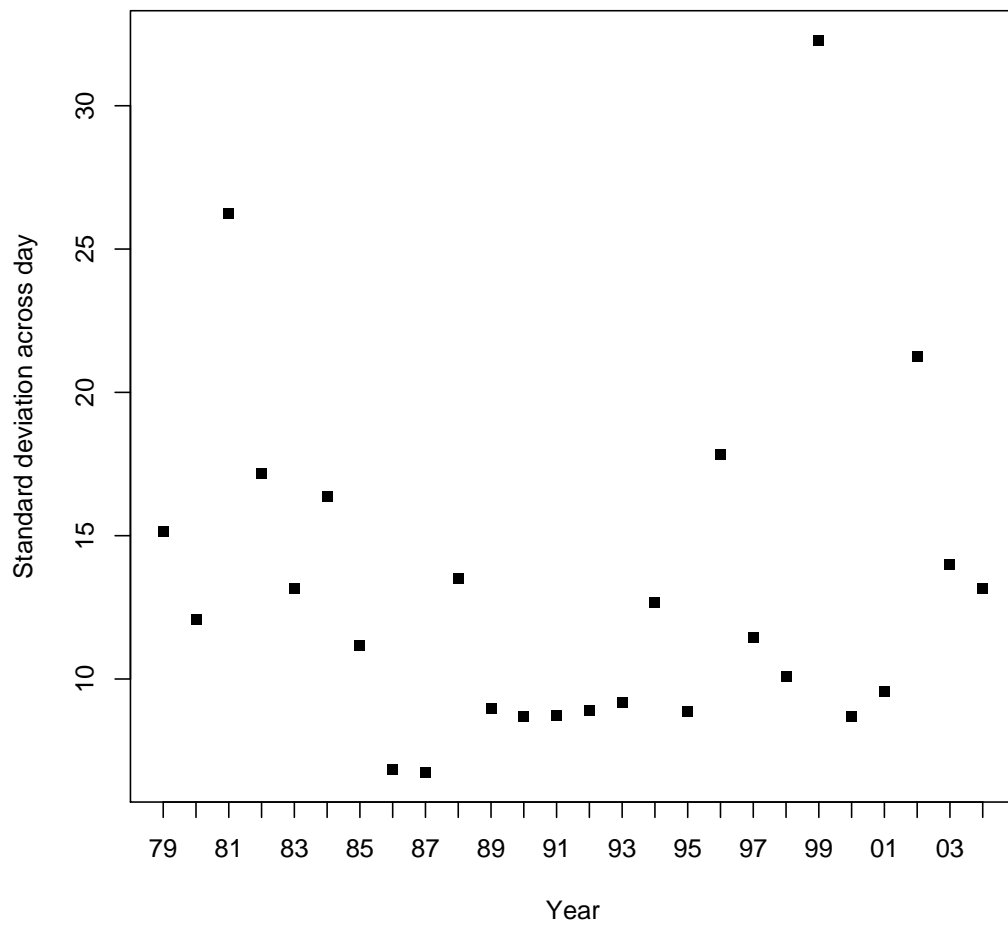


Figure 5: Standard deviations of the residuals from the smooth flows across day

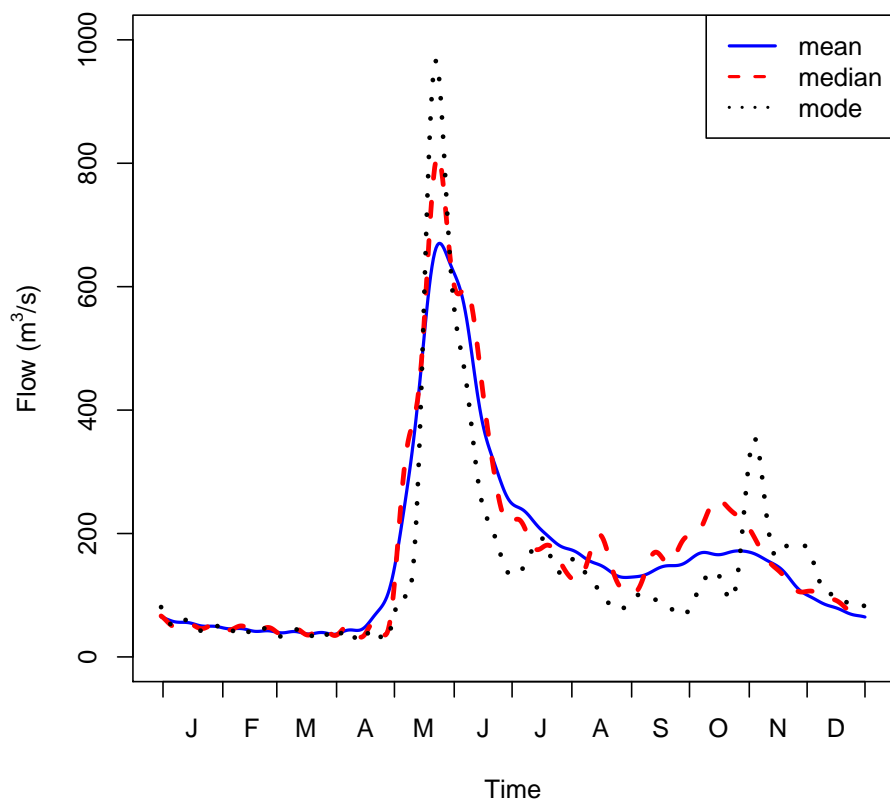
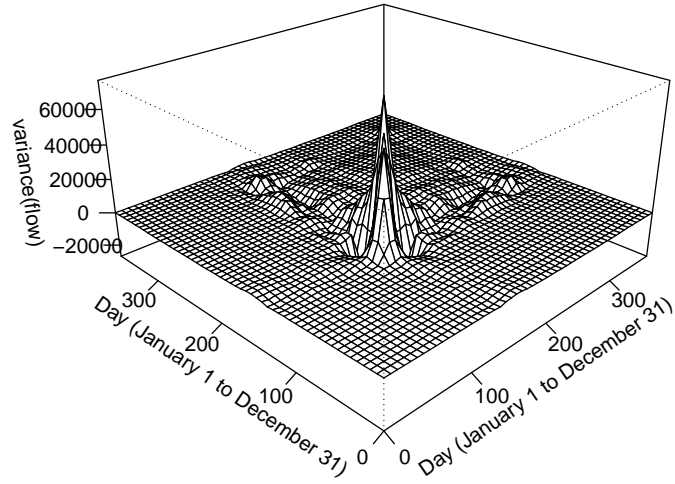
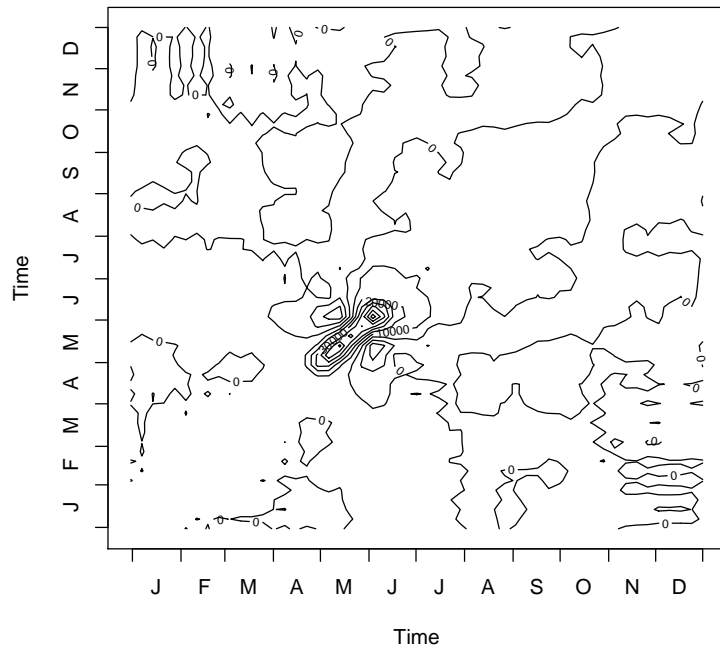


Figure 6: Fourier smoothed location curves : the mean, the median and the mode



(a)



(b)

Figure 7: Estimated variance-covariance (a) surface of the flow curves for years 1979 to 2004 and (b) the corresponding contour map

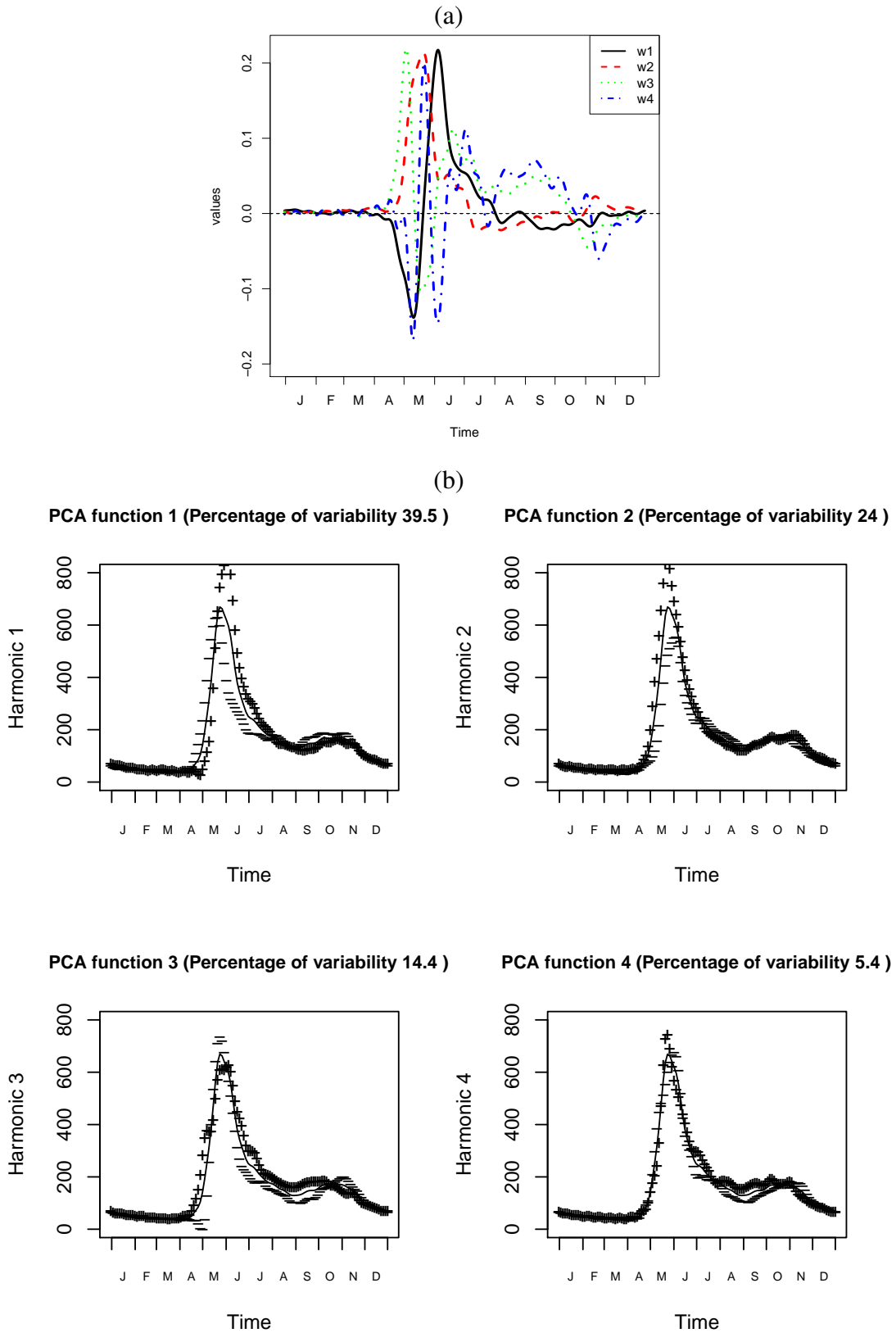


Figure 8: First four smoothed principal components: (a) centered components; (b) components with variation about the mean \bar{y} . Negative and positive perturbations are indicated respectively by the minus (-) and plus (+) symbols.

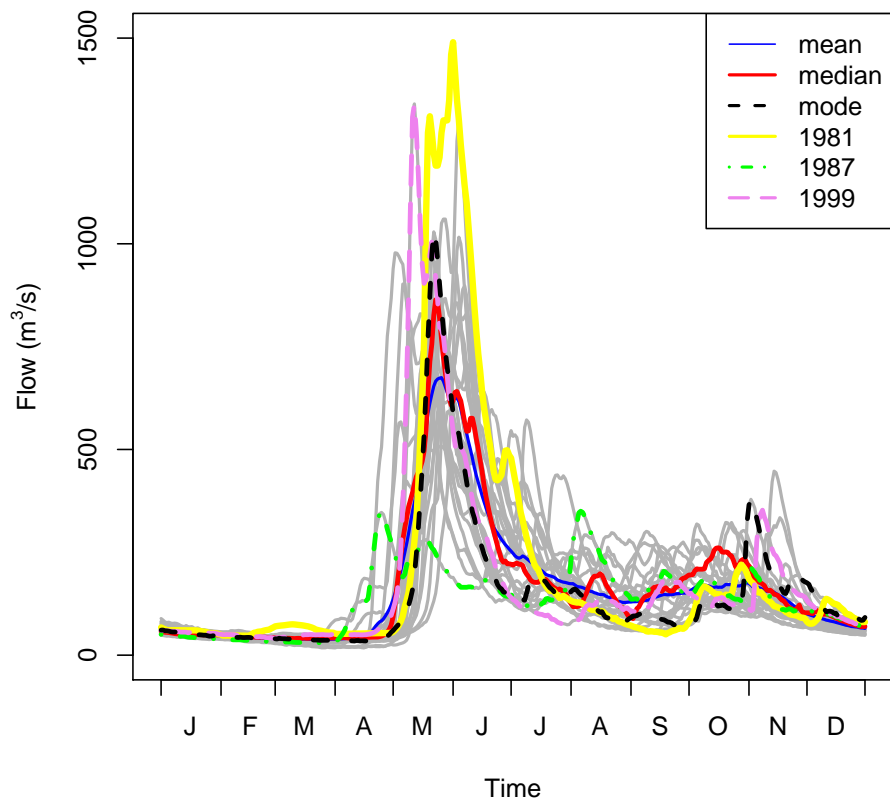
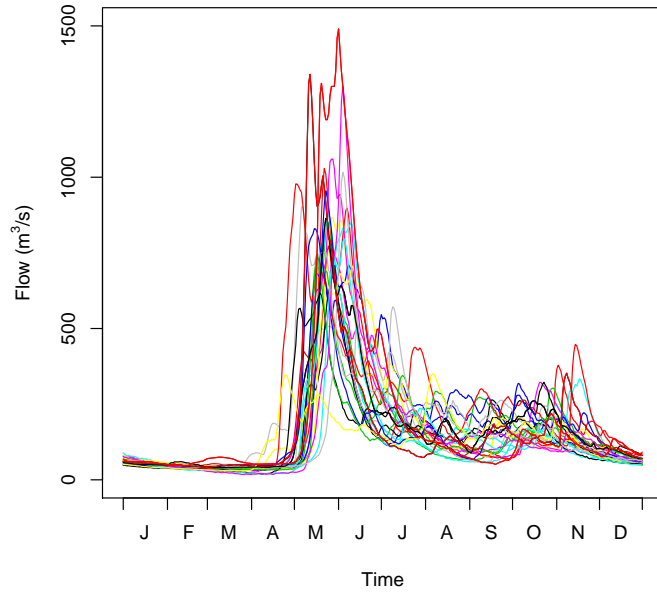
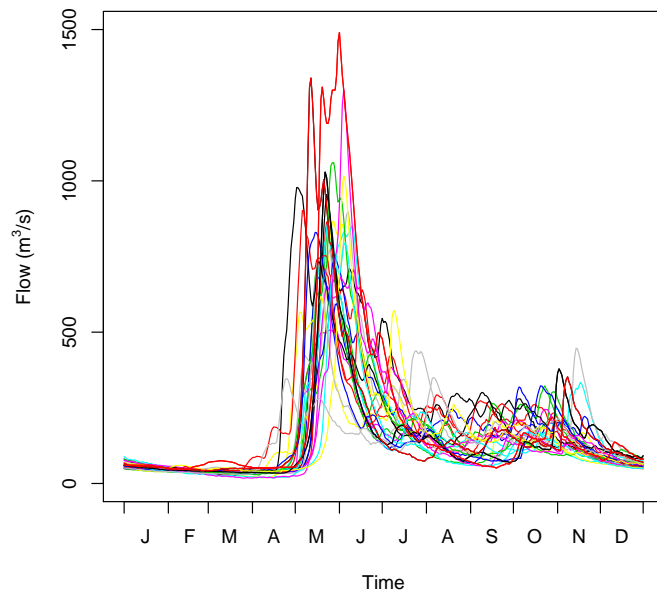


Figure 9: Curves corresponding to the suspected years (based on principal component scores) with the mean, median and mode curves



(a)



(b)

Figure 10: Rainbow plots of the flow curves for years 1979 to 2004 using (a) the bivariate score depth and (b) the kernel density estimate.

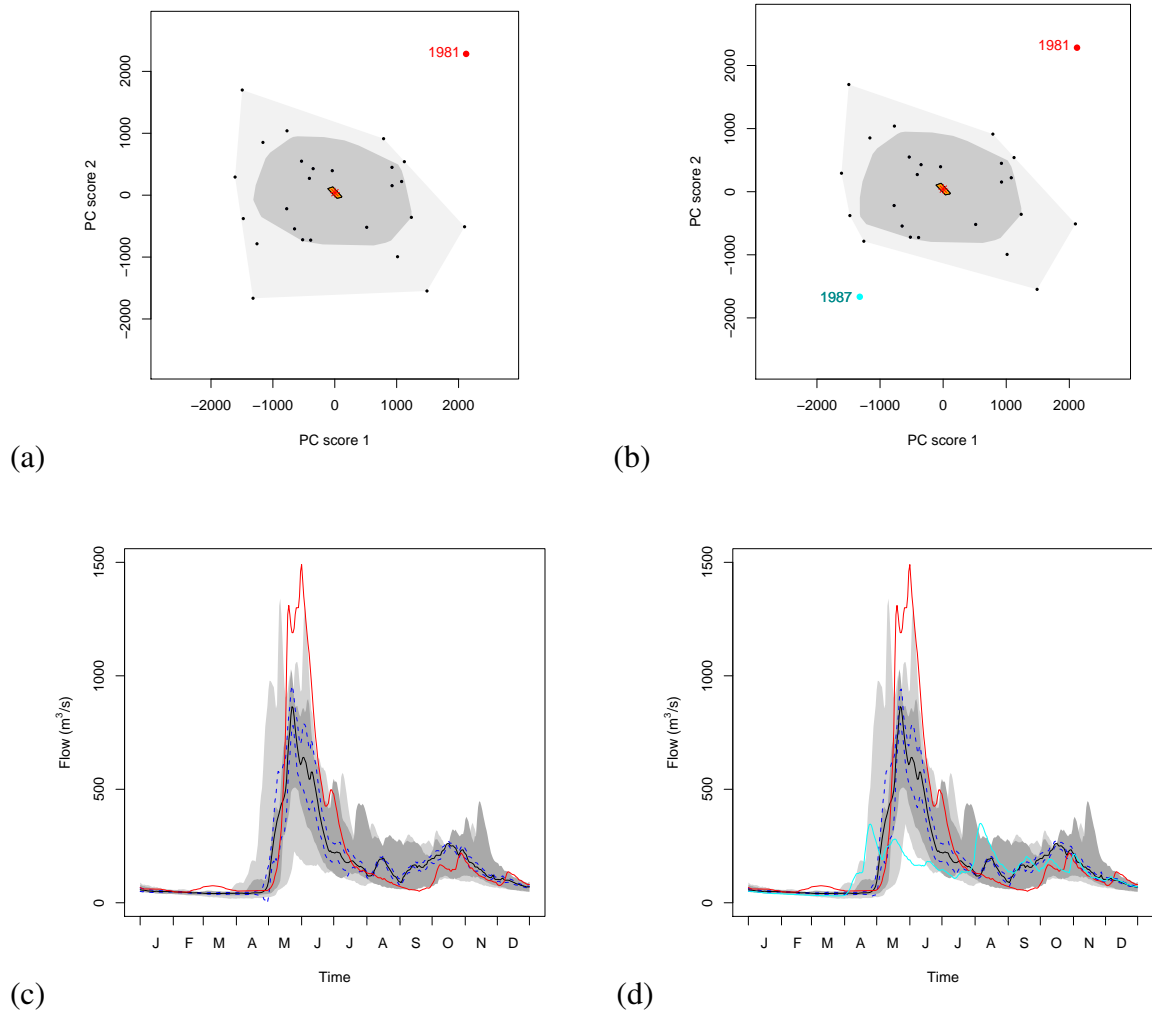
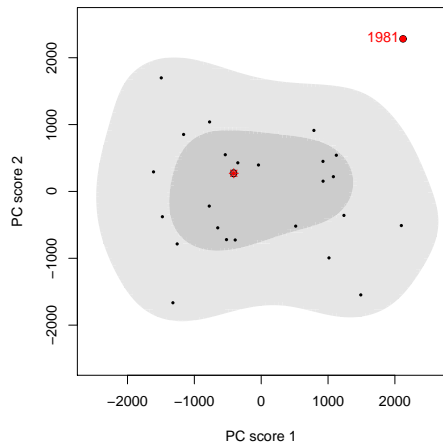
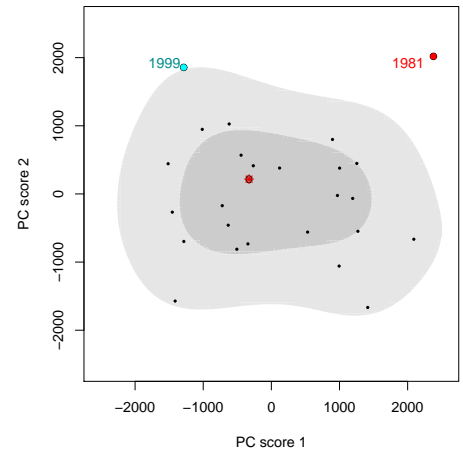


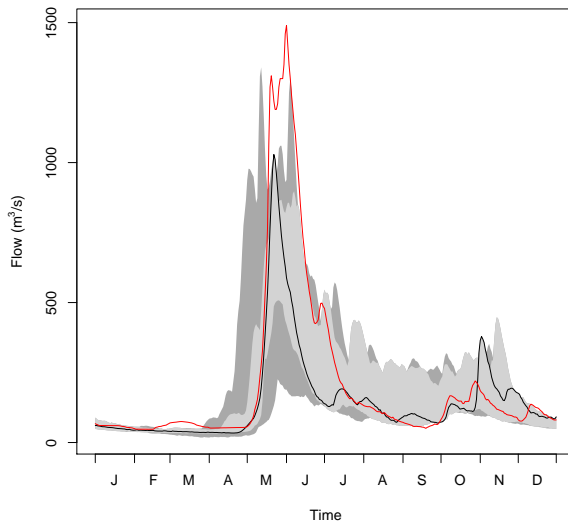
Figure 11: Bivariate score Bagplot with (a) 99% and (b) 95% of probability coverage and the corresponding functional Bagplot with (c) 99% and (d) 95% of probability coverage. The solid black curve shows the median curve and in blue are presented its 95% or 99% point-wise confidence intervals while in (a) and (b) the red asterisk is the Tukey median of the bivariate principal scores



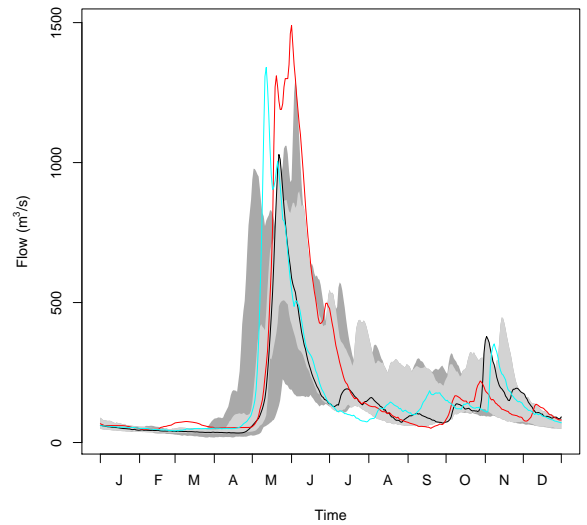
(a)



(b)



(c)



(d)

Figure 12: Bivariate score HDR boxplot with (a) 99% and (b) 95% of probability coverage and the corresponding functional HDR boxplot with (c) 99% and (d) 95% of probability coverage. The solid black curve shows the modal curve and in blue are presented its 95% or 99% point-wise confidence intervals while in (a) and (b) the red asterisk is the mode

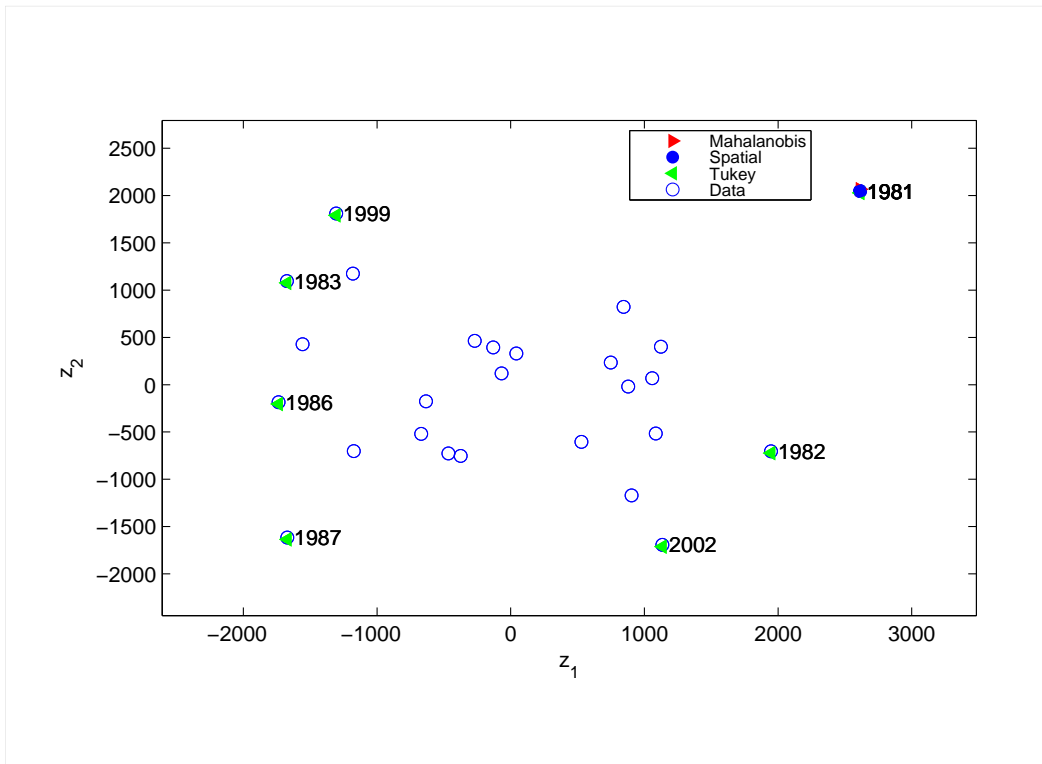
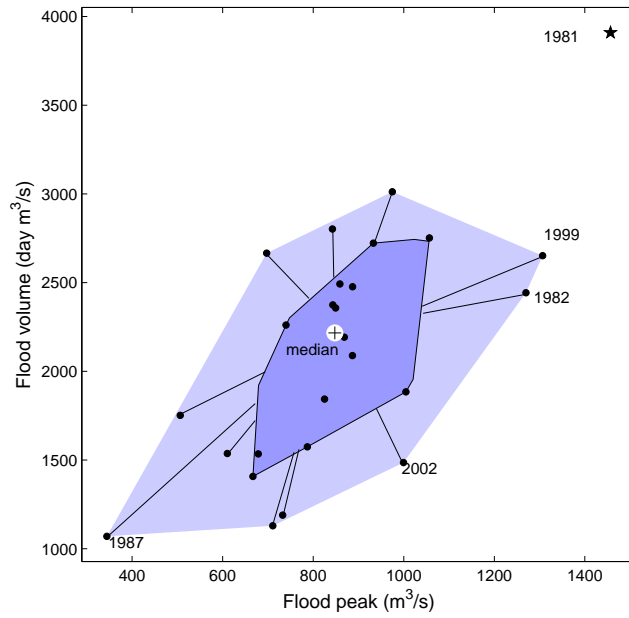
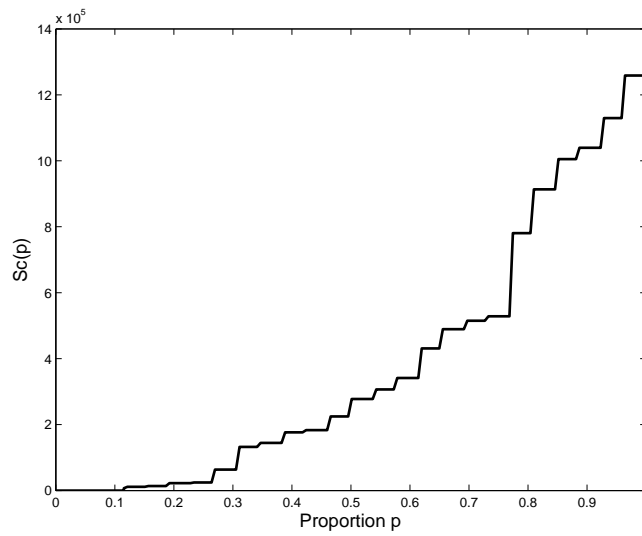


Figure 13: Outlier detection using the Outlyingness approach applied on the first two scores



(a)



(b)

Figure 14: Bivariate results : (a) Bagplot with the median and some particular years and (b) Scalar scale function