

1 **Depth and homogeneity in regional flood frequency analysis**

2
3 F. Chebana* and T.B.M.J. Ouarda

4
5
6 *Industrial Chair in Statistical Hydrology/Canada Research Chair on the Estimation of*
7 *Hydrometeorological Variables,*
8 *INRS-ETE, 490 rue de la Couronne, Quebec (QC),*
9 *Canada G1K 9A9*
10

11
12
13
14
15
16
17 * **Corresponding author:** Tel: (418) 654-2542
18 Fax: (418) 654-2600
19 Email: fateh_chebana@ete.inrs.ca

20 July 11th 2008

1 **Abstract:**

2 Regional frequency analysis (RFA) consists generally in two steps: 1) delineation of
3 hydrological homogeneous regions, and 2) regional estimation. Existing regionalization methods
4 which adopt this two-step approach suffer from two principal drawbacks. First, the restriction of
5 the regional estimation to a particular region by excluding some sites can correspond to a loss of
6 some information. Second, the definition of a region generates a border effect problem. In order
7 to overcome these problems, a new method is proposed in the present paper. The proposed
8 method is based on three elements: (i) a weight function to treat the border effect problem, (ii) a
9 function to evaluate how “similar” each site is to the target one, and (iii) an iterative procedure to
10 improve estimation results. Element (ii) is treated using the statistical notion of depth functions
11 which is introduced to provide a ranking of stations in a multivariate context. Furthermore, the
12 properties of depth functions meet the characteristics sought in RFA. It is shown that the
13 proposed method is flexible and general and that traditional RFA methods represent special cases
14 of the depth-based approach corresponding to particular weight functions. A comparison is
15 carried out with the canonical correlation analysis (CCA) approach. Results indicate that the
16 depth-based approach performs better than CCA both in terms of relative bias and relative root
17 mean squares error.

18

1 **1. Introduction**

2 One of the problems encountered in hydrology is the lack of data, since the extreme
3 events we want to estimate are rare and record lengths are short. Consequently, statistical
4 inference is difficult in such sites. To overcome this problem, hydrologists have recourse to data
5 from other sites that are hydrologically similar to the target one. The estimation of extreme
6 hydrological events, such as floods, at sites where little or no data is available is the main aim of
7 regional frequency analysis. Delineation of homogeneous hydrological regions and regional
8 estimation are the two main steps in a regional flood frequency procedure. Several studies have
9 focused on the delineation of homogeneous regions (e.g., Burn, 1990; Hosking and Wallis, 1993;
10 Ouarda et al., 2006; and Chebana and Ouarda, 2007) and on regional estimation (e.g., Dalrymple,
11 1960; Stedinger and Tasker, 1986; Ouarda and Ashkar, 1994; Durrans and Tomic, 1996; Nguyen
12 and Pandey, 1996; Madsen and Rosbjerg, 1997; Alila, 1999, 2000 and Chokmani and Ouarda,
13 2004). An intercomparison of various regional flood estimation procedures was presented by
14 GREHYS (1996a,b) by coupling four methods for delineating homogenous regions and seven
15 regional estimation methods.

16

17 In a regional estimation procedure, one is interested in maximizing the amount of
18 transferred information. The delineation step corresponds usually to the exclusion of a number of
19 sites which may lead to a loss of some relevant information. Furthermore, the definition of a
20 region leads to the problem of the so-called “border effect”. This means that for two sites that are
21 very close but which are located on each side of the region limits, one is excluded while the other
22 one is included even though both sites offer similar information. This problem is not present
23 when the limits correspond to natural borders.

1 These elements motivate the development of a new method that overcomes some of these
2 drawbacks. The method developed in the present paper is based on the notion of depth function.
3 The depth function is a statistical notion developed in the seventies and which receives increasing
4 interest (e.g. Tukey, 1975; Liu, 1990; Liu and Singh, 1993; Rousseeuw and Hubert, 1999; Zuo
5 and Serfling, 2000; Mizera, 2002; Mizera and Müller, 2004; Zuo and Cui, 2005 and Lin and
6 Chen, 2006). The purpose for introducing the depth function is to provide an outward ordering of
7 points in a multivariate context. The most important properties of a depth function are (i) affine
8 invariance, (ii) maximality at center, (iii) monotonicity related to the deepest point, and (iv)
9 vanishing at infinity. These properties fit the constraints of regional flood frequency analysis.
10 Indeed, the affine invariance is useful to remove the scale effect when treating several variables.
11 The center, where the depth function is maximal, represents the target site. The monotonicity
12 related to the center point means that sites far from the target site are less important. Finally, the
13 very far sites have no importance or no contribution; this is the vanishing at infinity property.
14 Hence, the contribution of each site is related to its similarity to the target site.

15
16 The proposed procedure focuses on the estimation as a goal and avoids the delineation
17 step. Note that the delineation step is only an intermediate technical tool to estimate quantiles.
18 Different quantile estimation methods are proposed in the literature, such as the index flood
19 method and regressive models (see GREHYS, 1996a,b). In the present paper, the estimation is
20 based on the regressive model. The estimation of the regression parameters is obtained using a
21 weighted least squares method. A key element is related to the choice of the weights. These
22 weights are selected as functions of the site depth. Hence, the proposed method overcomes the
23 problems related to the border effect by using smooth weight functions. It also reduces the
24 problems related to the lack of data by using all available data in a more efficient manner.

1 Furthermore, traditional approaches represent special cases of depth-based approaches
2 corresponding to particular weight functions. Finally, the non requirement of data normality, the
3 availability of several kinds of depth functions and the smoothness of the weight functions
4 provide this method with a high level of flexibility.

5
6 Based on hydrological variables, on one hand, and on physio-meteorological
7 characteristics, on the other hand, regional regression is frequently integrated with the CCA
8 approach (Ouarda et al., 2000, 2001). The CCA-regression provides flood quantile estimates at
9 ungauged sites by using site physiographic characteristics. In order to study the performance of
10 the proposed depth-based approach, it is compared to the CCA-regression approach with optimal
11 neighborhoods (Ouarda et al., 2001). This comparison is based on a data set from 151 gauging
12 sites in the southern part of the province of Quebec, Canada. The specific quantiles
13 corresponding to 10- and 100-year return periods are estimated and a jackknife resampling
14 procedure is used to evaluate the estimation errors.

15
16 The paper is organized as follows. Brief presentations of the depth functions, the CCA
17 approach and the weighted least squares method are given in Section 2. Section 3 deals with the
18 proposed methodology in its general form. A case study is presented in Section 4 and the
19 developed approach is applied and compared to the CCA method in Section 5. Results and
20 discussions are reported in Section 6. Conclusions and future promising work are presented in the
21 last section.

22
23

1 **2. Background**

2 In this section, a brief description of the background material related to depth functions,
3 the CCA method and weighted regression analysis is presented.

5 **2.1 Depth function:**

6 Tukey (1975) presented the pioneering work in which the depth notion was introduced.
7 The author proposed a halfspace depth in order to define a multivariate analogous to the
8 univariate rank and order statistics. Later, several depth functions were formulated in an ad-hoc
9 manner. Zuo and Serfling (2000) standardized these definitions and classified existing examples
10 in the literature.

11
12 For a given cumulative distribution function F on R^d ($d \geq 1$) a corresponding depth
13 function is any bounded, nonnegative function $D(x;F)$ which provides a F -based center-outward
14 ordering of points x in \mathbb{R}^d that satisfies the following properties:

- 15 i. *Affine invariance*: the depth of a point $x \in R^d$ should not depend on the underlying
16 coordinate system or, in particular, on the scales of the underlying measurements.
- 17 ii. *Maximality at center*: for a distribution having a uniquely defined center (e.g., the point of
18 symmetry with respect to some notion of symmetry), the depth function should attain its
19 maximum value at this center.
- 20 iii. *Monotonicity relative to deepest point*: as a point $x \in R^d$ moves away from the deepest
21 point (the point at which the depth function attains its maximum value; in particular, for a
22 symmetric distribution, the center) along any fixed ray through the center, the depth at x
23 should decrease monotonically.

1 iv. *Vanishing at infinity*: the depth of a point x should approach zero as its norm $\|x\|$
2 approaches infinity.

3
4 A formal definition of depth functions, based on these properties, is given in Zuo and
5 Serfling (2000). In the following, we denote by Δ the class of cumulative distribution functions
6 on R^d and by F_Z the cumulative distribution function of a given random vector Z . Let the
7 mapping $D(\cdot, \cdot) : R^d \times \Delta \rightarrow R$ be bounded, non-negative, and satisfy the following conditions:

8 i. $D(Ax + b; F_{AX+b}) = D(x; F_X)$ holds for any random vector X in R^d , any
9 $d \times d$ nonsingular matrix A , and any d -vector b ;

10 ii. $D(\theta; F) = \sup_{x \in R^d} D(x; F)$ holds for any $F \in \Delta$ having center θ ;

11 iii. For any $F \in \Delta$ having deepest point θ , $D(x; F) \leq D(\theta + \alpha(x - \theta); F)$ holds for
12 $\alpha \in [0, 1]$; and

13 iv. $D(x; F)$ converges to 0 as the norm $\|x\|$ goes to infinity, for each $F \in \Delta$.

14 Then $D(\cdot; F)$ is called a statistical depth function.

15

16 Several kinds of depth functions are introduced in the literature. Here we present some of
17 the key ones:

18 1. *Mahalanobis depth*: It is defined on the basis of the Mahalanobis distance

19 $d_A^2(x, y) = (x - y)' A^{-1} (x - y)$ between two points $x, y \in R^d$ with respect to a positive
20 definite matrix A (Mahalanobis, 1936). The Mahalanobis depth is then given by:

21
$$MHD(x; F) = \frac{1}{1 + d_A^2(x, \mu)} \quad (1)$$

1 where F is a given distribution and μ and A are any corresponding location and
 2 covariance measures, respectively. Note that the Mahalanobis distance is used in the
 3 development of the CCA approach for regional flood frequency analysis (Ouarda et al.,
 4 2001). It is also important to note that the Mahalanobis depth function has values in the
 5 interval $[0,1]$. Hence, its values are more interpretable than those of the corresponding
 6 Mahalanobis distance.

7 2. L^2 depth: It is defined for a distribution F and $x \in R^d$ as:

$$8 \quad L^2D(x; F) = \frac{1}{1 + E \|x - X\|_{\Sigma^{-1}}} \quad (2)$$

9 where Σ is the covariance matrix of F and $\|x\|_M = \sqrt{x'Mx}$.

10 3. *Simplicial volume depth* (Oja, 1983): It is given through the expression :

$$11 \quad SVD^\alpha(x, F) = \left(1 + E \left[\left(\frac{\Delta(S[x, X_1, \dots, X_d])}{\sqrt{\det(\Sigma)}} \right)^\alpha \right] \right)^{-1} \quad \text{for } x \in R^d \quad (3)$$

12 where $\Delta(S[x, X_1, \dots, X_d])$ denotes the volume of the d -dimensional simplex
 13 $S[x, X_1, \dots, X_d]$, Σ is the covariance matrix of F and $\alpha > 0$. The quantity
 14 $\Delta(S[x, X_1, \dots, X_d])$ is a measure of the dispersion of the point cloud.

15 4. *Projection depth* (Liu, 1992): It is defined for $x \in R^d$ as :

$$16 \quad PD(x; F) = \left(1 + \sup_{\|u\|=1} \frac{|u'x - \text{Med}(u'X)|}{\text{MAD}(u'X)} \right)^{-1} \quad (4)$$

1 where X has distribution F , Med denotes the univariate median,
 2 $\text{MAD}(Y) = \text{Med}(|Y - \text{Med}(Y)|)$ represents the median absolute deviation for a random
 3 variable Y , and $\|\cdot\|$ is the Euclidian norm.

4 5. *Halfspace depth* (Tukey, 1975): It is defined for $x \in R^d$ with respect to a probability P on
 5 R^d as:

$$6 \quad HD(x; P) = \inf \{P(H) : H \text{ a closed halfspace that contains } x\} \quad (5)$$

7
 8 Note that, occasionally, some examples of depth functions do not meet some of the
 9 previous four properties in some special cases. For instance, the $L^2D(x; F)$ function only meets
 10 the above conditions under symmetric assumptions on the distribution F (see Zuo and Serfling,
 11 2000).

12
 13 A sample version of $D(x; F)$, denoted by $D_n(x) = D(x; \hat{F}_n)$, may be defined by replacing F
 14 with a suitable empirical function \hat{F}_n . The asymptotic properties of $D_n(x)$ are studied in several
 15 papers including Liu (1990), Arcones et al. (1994), Massé (2004) and Lin and Chen (2006). Liu
 16 and Singh (1993) established for the sample Mahalanobis depth function that
 17 $\sup_x |D_n(x) - D(x; F)|$ converges to zero almost surely as n goes to infinity, under suitable
 18 conditions on F . For convenience, the following notation is used for the Mahalanobis depth
 19 function in the next sections:

$$20 \quad MHD_A(x; \mu) = \frac{1}{1 + d_A^2(x, \mu)} \quad (6)$$

1 Depth functions are applied in several fields. For instance, Caplin and Nalebuff (1988,
2 1991a,b) employed depth notions in econometric and social studies. They were also applied in
3 industrial quality control by Liu and Singh (1993) and Liu (1995). Ghosh and Chaudhuri (2005)
4 investigated the use of depth functions in nonparametric discrimination analysis. Mizera and
5 Müller (2004) defined and studied the location-scale depth and gave some statistical applications.

6
7 The computation of some depth functions is complex and requires specific algorithms.
8 For instance, Miller et al. (2003) developed an algorithm for the computation of the halfspace
9 depth. However, to our knowledge, similar algorithms are not available for the projection depth.
10 The Mahalanobis depth is among the simplest ones to evaluate if the parameter μ and the
11 parameter matrix A are identified.

12
13 Liu et al. (1999) presented descriptive statistics, graphics and inference related to several
14 depth functions. A detailed description of the theoretical background of depth functions is
15 available in Zuo and Serfling (2000).

16
17 **2.2 Canonical Correlation Analysis (CCA)**

18 Canonical correlation analysis is concerned with the amount of linear relationship
19 between two sets of variables. Consider two normal random vectors X (physiographical and
20 meteorological variables) and Y (hydrological variables), CCA provides two sets of basis vectors
21 (called canonical variables), one for X and the other for Y . The main property of such vectors is
22 that the correlations between the projections of the variables onto these basis vectors are mutually

1 maximized (Muirhead, 1982). More precisely, let W and V be linear combinations of X and Y
 2 respectively, i.e.,

$$\begin{aligned} V &= a'X \\ W &= b'Y \end{aligned} \quad (7)$$

4 Let Σ be the covariance matrix of the variables X and Y , defined as:

$$\Sigma = \text{cov} \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{bmatrix} \Sigma_X & \Sigma_{XY} \\ \Sigma_{YX} & \Sigma_Y \end{bmatrix} \quad (8)$$

6 The correlation between W and V can then be calculated as:

$$\rho = \frac{a' \Sigma_{XY} b}{\sqrt{a' \Sigma_X a b' \Sigma_Y b}} \quad (9)$$

8 The goal of the CCA is to find the vectors a and b maximizing ρ subject to the constraint
 9 that W and V must have unit variances. Once the first pair of canonical variables is obtained,
 10 other pairs of canonical variables can be obtained in the uncorrelated directions to the previous
 11 ones by maximizing equation (9) subject to the constraint of unit variance. For more details
 12 concerning CCA application in regional flood frequency analysis, the reader is referred to Ouarda
 13 et al. (2001).

15 Based on the canonical hydrological variables W and physio-meteorological variables V ,
 16 the Mahalanobis distance for an ungauged target-site with given physiographical characteristics
 17 $V=v_0$ is given by:

$$D^2 = d_{I_p - \Lambda' \Lambda}^2(W, \Lambda v_0) = (W - \Lambda v_0)' (I_p - \Lambda' \Lambda)^{-1} (W - \Lambda v_0) \quad (10)$$

19 where I_p is the $p \times p$ identity matrix, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_p)$ with $\lambda_i = \text{corr}(V_i, W_i), i = 1, \dots, p$ and p is
 20 the rank of the covariance matrix Σ_{XY} .

1 To belong to the neighborhood of the target-site, at a $100(1 - \alpha)\%$ confidence level, a site
 2 should have a distance D^2 less than $\chi_{\alpha,p}^2$ where $P(\chi_p^2 \leq \chi_{\alpha,p}^2) = 1 - \alpha$ and χ_p^2 has a chi-squared
 3 distribution with p degrees of freedom. The Mahalanobis distance can also serve to indicate
 4 where, in the canonical space W that encompasses all possible realizations of the random variable
 5 W , would be found the realizations w of W for which V has realized v_0 .

6

7 **2.3 Weighted least squares estimation**

8 Most commonly, the power product model given by equation (10) is integrated with the
 9 CCA approach and is used to describe the relationship between flood quantiles QT of a return
 10 period T , and the physio-meteorological and basin characteristics A_1, \dots, A_r for a given region:

$$11 \quad QT = \alpha_0 A_1^{\alpha_1} A_2^{\alpha_2} \dots A_r^{\alpha_r} e \quad (11)$$

12 Taking s quantiles QT corresponding to s return periods, we construct a vector Y of the
 13 hydrological variables, that is $Y = (QT_1, QT_2, \dots, QT_s)$. Then, using the log-transformation and the
 14 matrix form, we obtain the multivariate log-linear model:

$$15 \quad \log Y = (\log X) \beta + \varepsilon \quad (12)$$

16 where $\log X = (1, \log A_1, \log A_2, \dots, \log A_r)$ is the $(r+1)$ vector of the physio-meteorological
 17 variables, β is the $(r+1) \times s$ matrix of parameters and ε represents the error vector with:

$$18 \quad E(\varepsilon) = 0 \text{ and } \text{Var}(\varepsilon) = \Gamma \quad (13)$$

19 If the number of sites in the region is denoted by N , the parameter β can be estimated, using the
 20 weighted least squares estimation method, by:

$$\begin{aligned} \hat{\beta}_w &= \arg \min_{\beta} \sum_{i=1}^N w_i (\log Y_i - \beta \log X_i)' (\log Y_i - \beta \log X_i) \\ &= ((\log X)' \Omega \log X)^{-1} (\log X)' \Omega \log Y \end{aligned} \quad (14)$$

and the matrix Γ is estimated by:

$$\hat{\Gamma} = \frac{(\log Y - \hat{\beta}_w \log X)(\log Y - \hat{\beta}_w \log X)'}{N - r - 1} \quad (15)$$

where $\Omega = \text{diag}(w_1, \dots, w_N)$ is the diagonal matrix composed by the weight elements w_1, \dots, w_N . A detailed description of multivariate regression analysis can be found in Rencher (2002).

3. Approach development

3.1 Description

Limited use of weighted least squares methods was made in the field of regional flood frequency analysis. Madsen and Rosbjerg (1997) used weighted least squares (WLS) and generalized least squares (GLS) methods in a regional flood estimation procedure that combines the index-flood concept with an empirical Bayes method. In the WLS and GLS methods, the weights are related to the variance and covariance of the errors in the regression model.

The approach proposed in the present work is focused directly on quantile estimation using the weighted least squares method to estimate regression parameters and does not use any delineation technique. The choice of the weights in equation (14) is very important for parameter estimation and hence for the predicted value \hat{Y} . In the un-weighted estimation, all weights are equal to one. However, if a region or a neighborhood has been defined, weights correspond to zero if the site is excluded from the region and one if it is included in the region. In the proposed

1 approach the weights are chosen differently. They are related to a weight function and a depth
2 function which will be developed in Section 3.2. This makes the methodology very flexible and
3 more general.

4
5 A special attention is given to the choice and the evaluation of the depth function. A
6 convenient depth function, which is related to the neighborhood approach, is the Mahalanobis
7 depth function (6). The value of μ is generally unknown for the ungauged site, and must be
8 estimated. The values of the depth function for the gauged sites are highly related to the quality
9 of μ estimates. In other words, the problem here is how to get the hydrological «reference value»
10 with respect to which the depths are computed. This reference value represents the deepest point.
11 Therefore, in order to get an accurate estimate of this «reference value», the proposed approach
12 utilizes an iterative estimation procedure based on the log-linear model. The iterative procedure
13 requires a start point, a criterion and a stopping condition. The approach is described below in its
14 general aspect and also with options for the iteration elements.

15
16 The iterative estimation procedure serves to improve the depth values and to make them
17 more accurate. This iterative technique has some similarities with the so-called One-step
18 estimator (see e.g., van der Vaart, 1998, pp. 71). Note that the iterative estimation procedure and
19 the way the weights are selected represent two elements that differentiate the proposed approach
20 from the WLS and GLS methods as applied in Madsen and Rosbjerg (1997). In the described
21 methodology, all available sites in the data set are used, without any restriction to a region or a
22 neighborhood. However, each site is associated to a weight related to its hydrological depth with

1 respect to the target-site. In that case the problem of the delineation of a region becomes rather a
 2 problem of a choice of weight and depth functions.

3

4 **3.2 Computation algorithm**

5 The computation algorithm is based on the following estimate of the parameter β . It is the
 6 estimator given in equation (14) with particular weights w_i . To this end, let $\varphi(\cdot)$ be a positive
 7 increasing weight function, and D_N be a sample depth function. Then, the estimators (14) and (15)
 8 are given respectively by:

$$\begin{aligned}
 \hat{\beta}_{D_N, \varphi} &= \arg \min_{\beta} \sum_{i=1}^N \varphi(D_N(Y_i)) (\log Y_i - \beta \log X_i)' (\log Y_i - \beta \log X_i) \\
 &= \left((\log X)' \Omega_{D_N, \varphi} \log X \right)^{-1} (\log X)' \Omega_{D_N, \varphi} \log Y
 \end{aligned}
 \tag{16}$$

10 and

$$\hat{\Gamma}_{D_N, \varphi} = \frac{(\log Y - \hat{\beta}_{D_N, \varphi} \log X) (\log Y - \hat{\beta}_{D_N, \varphi} \log X)'}{N - r - 1}
 \tag{17}$$

12 where $\Omega_{D, \varphi} = \text{diag}(\varphi(D(Y_1)), \dots, \varphi(D(Y_N)))$.

13

14 The weight function is assumed to be increasing, to ensure that the deeper is the site, the more
 15 important it is and hence it receives a higher weight. It is important to indicate that the matrix in
 16 (17) contains the inter-quantile correlation rather than the inter-site correlation. The latter is taken
 17 into account in the GLS approach (see e.g. Madsen and Rosbjerg, 1997). It would be useful, in
 18 future efforts, to focus on the integration of spatial correlation in the depth-based approach.

19

1 Suppose i_0 is the index of the target-site. The algorithm is composed of the following
2 steps:

3 1. *Model*: Consider the model (12) relating $\log X$ and $\log Y$.

4 2. *Initial step*: In order to get a starting estimator \hat{Y}_{1,i_0} .

5 - Use a preliminary approach to estimate the model parameters. One possible option is the
6 uniform approach which allocates equal weights to all sites.

7 - Give the predicted initial value \hat{Y}_{1,i_0} .

8 3. k^{th} step ($k = 2, 3, \dots$): For each site i different from the target one i_0 ,

9 a. Consider the predicted value from the previous step ($k-1$), i.e. \hat{Y}_{k-1,i_0} using the

10 corresponding parameter estimators $\hat{\beta}_{k-1,i_0}$ and $\hat{\Gamma}_{k-1,i_0}$ given in (16) and (17) respectively,

11 where the notation is adapted to the iteration context.

12 b. Compute the depth of $\log Y_i$ with respect to $\log \hat{Y}_{k-1,i_0}$ between sites i and i_0 .

13 It is convenient to use the Mahalanobis depth $MHD_{\hat{\Gamma}_{k-1,i_0}}(\log Y_i; \log \hat{Y}_{k-1,i_0})$ given by (6).

14 c. Compute the weight w_i corresponding to site i . It is given for a weight function φ and

15 the Mahalanobis depth by $w_i = \varphi\left(MHD_{\hat{\Gamma}_{k-1,i_0}}(\log Y_i; \log \hat{Y}_{k-1,i_0})\right)$. Some examples of

16 weight functions are presented in the next section.

17 d. Estimate the model parameters using the weighted method.

18 In the log-linear model, the estimators $\hat{\beta}_{k,i_0}$ and $\hat{\Gamma}_{k,i_0}$ are expressed through (16) and (17)

19 respectively.

20 e. Give the predicted value \hat{Y}_{k,i_0} using X_{i_0} and $\hat{\beta}_{k,i_0}$ in model (12).

- 1 f. Assess the criterion to evaluate the model estimation quality. Several known criteria in
 2 regression analysis can be used, for instance, the relative bias (RB) and the relative root
 3 mean square error (RRMSE).
- 4 4. *Iteration*: Redo step 3 while the criterion is improving.
- 5 5. *End condition*: Stop the iteration whenever the criterion deteriorates or converges.
- 6 6. *Final result*: Get the estimators of the regression model parameters and hence the final
 7 predicted value.

8

9 **3.3 Particular cases of the proposed approach**

10 The proposed method is general and includes some known methods as special cases
 11 representing particular weight functions. Indeed:

- 12 1. The uniform approach which uses all sites corresponds to the weight function $\varphi_U \equiv 1$. In
 13 this case the estimator $\hat{\beta}$ is given by:

$$14 \quad \hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N w_i (\log Y_i - \beta \log X_i)^2 = \arg \min_{\beta} \sum_{i=1}^N (\log Y_i - \beta \log X_i)^2 \quad (18)$$

- 15 2. The traditional CCA approach with a given value of α corresponds to:

$$16 \quad \varphi_{CCA}(x) = 1_{[C_{\alpha,p}, 1]}(x) = \begin{cases} 1 & \text{if } x \in [C_{\alpha,p}, 1] \\ 0 & \text{else} \end{cases} \quad (19)$$

17 where $C_{\alpha,p} = \frac{1}{1 + \chi_{\alpha,p}^2}$ and $\chi_{\alpha,p}^2$ is the χ_p^2 quantile of order α . The function (19) can be

18 written in an informal way as:

$$19 \quad \begin{aligned} \varphi_{CCA}(\text{site depth}) &= 1_{[C_{\alpha,p}, 1]}(\text{site depth}) = 1_{[0, \chi_{\alpha,p}^2]}(\text{site distance}) \\ &= \begin{cases} 1 & \text{if the site is inside the neighborhood} \\ 0 & \text{if the site is outside the neighborhood} \end{cases} \end{aligned}$$

1 Therefore, the estimator $\hat{\beta}$ is the following:

$$\begin{aligned}
 \hat{\beta} &= \arg \min_{\beta} \sum_{i=1}^N w_i (\log Y_i - \beta \log X_i)^2 \\
 &= \arg \min_{\beta} \sum_{\text{neighborhood sites } i} 1 (\log Y_i - \beta \log X_i)^2 \\
 &\quad + \arg \min_{\beta} \sum_{\text{non neighborhood sites } i} 0 (\log Y_i - \beta \log X_i)^2
 \end{aligned} \tag{20}$$

3 3. The following weight function is used by Zuo et al. (2004):

$$\varphi_Z(x) = \begin{cases} \frac{\exp\left(-K\left(1 - (x/C)\right)^2\right) - \exp(-K)}{1 - \exp(-K)} & \text{if } x < C \\ 1 & \text{elsewhere} \end{cases} \tag{21}$$

5 where C and K are positive constant coefficients.

6 4. With an extra coefficient $s \geq 1$ in equation (21), Lin and Chen (2006) used the following
7 weight function:

$$\varphi_{LC}(x) = \begin{cases} \frac{\exp\left(-K\left(1 - (x/C)^s\right)^2\right) - \exp(-K)}{1 - \exp(-K)} & \text{if } x < C \\ 1 & \text{elsewhere} \end{cases} \tag{22}$$

9
10 In functions (21) and (22) the constant C defines the support of the weight function and K
11 represents the slope of the decay to zero. The CCA neighborhood approach may serve to
12 guide the choice of the constant C . Indeed, the choice of a constant C may depend on the

13 χ_p^2 quantile associated to the optimal α as $C_{\alpha,p} = \frac{1}{1 + \chi_{\alpha,p}^2}$.

14 5. A simple linear function can also be used as a weight function:

$$\varphi_{Linear}(x) = \begin{cases} 0 & \text{if } x \leq d_1 \\ \frac{x - d_1}{d_2 - d_1} & \text{if } d_1 \leq x \leq d_2 \\ 1 & \text{if } x \geq d_2 \end{cases} \quad (23)$$

with $d_2 > d_1 > 0$.

Special cases 1 and 2 along with the general depth-based approach are illustrated in Figure 1.

4. Case study

In this section, the approach proposed in Section 3 is applied on a real world data set and its performance is compared to that of the CCA approach. The case study on which the comparison is carried out concerns the hydrometric station network of the southern part of the province of Quebec, Canada. To be selected, each station in the data set must have a flood record of at least 15 years of data and its historical data must be homogenous, stationary and independent. The area of these catchments is larger than 200 km² and less than 100 000 km². Finally, a total of 151 stations located between the 45° N and the 55° N are selected. The geographical location of these stations is shown in Figure 2.

The variable selection is based on a previous study by Chokmani and Ouarda (2004). The selected variables are of three types: physiographical, meteorological, and hydrological. The physiographical variables are: basin area (AREA), mean basin slope (MBS) and the fraction of the basin area covered with lakes (FAL). The meteorological variables are annual mean total

1 precipitation (AMP) and annual mean degree days over 0° C (AMD). The hydrological variables
2 are represented by at-site flood quantiles QT corresponding to a return period T .

3
4 The data bases used in the present study and the at-site frequency analysis are presented in
5 Kouider et al. (2002). The at-site flood frequency procedure includes the use of statistical tests of
6 hypothesis to test for homogeneity, independence and stationarity, the fitting of several statistical
7 distributions to the station data, and the use of goodness-of-fit tests to identify the most
8 appropriate distribution for each station. It is important to mention that the approach proposed in
9 the present paper does not require the use of a unique regional distribution. It is hence possible to
10 use a different (most appropriate) distribution in each station. For more details concerning the at-
11 site frequency analysis, the reader is referred to Kouider et al. (2002).

12
13 Eaton et al. (2002) reported that, when modeling the physical mechanism of drainage
14 systems, scale effect should be eliminated from experiment data since it may have a negative
15 impact on the results. To reduce the scale effect, flood quantiles QT are standardized by the basin
16 area to obtain specific quantiles $QST = QT/AREA$. In this study, the 10-year (QS10) and the
17 100-year (QS100) specific flood quantiles are selected. The basic statistics of these variables are
18 summarized in Table 1.

19
20 CCA application requires all variables to be transformed in order to be normalized and
21 standardized. The appropriate normalizing transformations for the present data set were obtained
22 by Chokmani and Ouarda (2004). A logarithmic transformation was used for the variables QS10,
23 QS100, AREA, MBS, AMP and AMD, and a square-root transformation for FAL. The traditional

1 CCA procedure is applied with values of the coefficient α ranging in the interval $[0, 1]$. An
2 optimal value of the coefficient α is selected according to minimum values of the relative bias
3 (RB) and the relative root mean square error (RRMSE) of the jackknife resampling procedure as
4 explained in Ouarda et al. (2001). The optimal value is found to be $\alpha = 0.25$ for the present case
5 study. The scatter plots of sites in the hydrological canonical space (W1,W2) and the physio-
6 meteorological canonical space (V1,V2) are illustrated in Figure 3.

7
8
9

10 **5. Study methodology**

11 The proposed depth-based approach described in Section 3 is applied to the above case
12 study, and is compared to the CCA approach. Other methods are also considered in the
13 comparison.

14

15 From equation (16), it can be seen that the depth computation method and the weight
16 function are the two main elements of the estimation in the proposed approach. The depth can be
17 computed in two ways: by the CCA Mahalanobis distance using directly equation (10) in
18 equation (6); or by the iterative algorithm described in Section 3. Equation (6) indicates that the
19 Mahalanobis depth and Mahalanobis distance are equivalent, that is, their values can be deduced
20 from each other. Various combinations of depth computation and weight selection methods are
21 considered in this study. The following methods are compared:

- 22 I. The uniform approach which uses all sites with the same importance.
- 23 II. The traditional CCA approach considered with the optimal value of α .

1 III. The depth-based approach considered with the following weight functions:

- 2 a. φ_Z with $K = 200$ and $C = 0.51$.
- 3 b. φ_{LC} with $s = 5$, $K = 200$ and $C = 0.51$.
- 4 c. φ_{LC} with $s = 2$, $K = 100$ and $C = 0.52$.
- 5 d. φ_{Linear} with $d_1 = 0.30$ and $d_2 = 0.80$.

6 IV. The CCA approach with iteration: It consists in the combination of the weight function
7 φ_{CCA} with the optimal value of $\alpha = 0.25$ along with the Mahalanobis depth
8 evaluated by the iterative algorithm. This combination can be seen as a special
9 case of (III) with the specific weight function φ_{CCA} .

10 V. The depth-based approach without iteration: The depth is evaluated from the CCA
11 Mahalanobis distance using equation (10). The weight function is φ_{LC} with the
12 following coefficients:

- 13 a. $s = 2$, $K = 100$ and $C = 0.30$, which is similar in shape to φ_{CCA} with $\alpha = 0.25$.
- 14 b. $s = 2$, $K = 100$ and $C = 0.52$, which is one of the weight functions
15 considered in (III).

16
17 The combinations (IV) and (V) are introduced to study the effect of the depth evaluation method
18 and the weight function selection. The considered combinations are summarized in the left-hand
19 part of Table 2 and the corresponding weight functions are illustrated in Figure 4.

20
21 In order to evaluate the performance of the various methods, a jackknife resampling
22 procedure is used. It consists in considering each site as an ungauged one by removing it

1 temporarily from the region. The criteria employed to evaluate the performances of the
2 approaches are the relative bias (RB) and the relative root mean square error (RRMSE) given
3 respectively by:

$$4 \quad RB = \frac{1}{N} \sum_{i=1}^N \left(\frac{y_i - \hat{y}_i}{y_i} \right) \quad (24)$$

$$5 \quad RRMSE = \sqrt{\frac{1}{N-1} \sum_{i=1}^N \left(\frac{y_i - \hat{y}_i}{y_i} \right)^2} \quad (25)$$

6 where y_i are the local realizations of the hydrological variable, \hat{y}_i are the regional estimates and
7 N is the number of sites in the data set.

8

9 **6. Results and discussion**

10 Results related to the various methods are summarized in Table 2. For all methods, results
11 indicate that the RB and RRMSE are smaller for QS10 than QS100. Generally in frequency
12 analysis, QST is more accurately estimated than QST' if $T < T'$ since for small return periods, the
13 corresponding quantile is close to the central body of the distribution. Hence, an important part of
14 the data contributes to its estimation. Table 2 shows also that the results of the uniform method (I)
15 are the worst. This confirms the need to use regional delineation techniques. The remaining
16 methods are classified according to the depth evaluation procedure (iteration or direct CCA). The
17 iterative depth evaluation leads to better results than the direct evaluation using CCA. Indeed, the
18 results of the depth-based approach (III) are the best, and are followed by those of method (IV).
19 In these two methods the depths are iteratively evaluated. Moreover, the differences in terms of
20 RB and RRMSE are not significant between the various combinations in (III).

21

1 In methods (II) and (V), depths are evaluated directly using the CCA Mahalanobis
2 distance. These methods lead to RB and RRMSE values that are larger than those obtained by
3 methods (III) and (IV). In particular, the RB and RRMSE of methods (II) and (V) are
4 significantly larger than those of methods (III). The RB and RRMSE of (IV) are slightly smaller
5 than those of (II). In these last two methods, the same weight function φ_{CCA} is used. Hence, the
6 CCA approach results can be slightly improved when depths are iteratively evaluated. However,
7 the results from (V.b) are not satisfactory compared to the other methods except the uniform one
8 (I). Note that combination (V.b) uses the same weight function than (III.c). Note also that the
9 results of combinations (V.a) and (II) are very similar. In both these methods the depth is
10 evaluated using the CCA Mahalanobis distance and the corresponding weight functions have
11 similar shapes (see Figure 4). The results of combination (III.d) and the shape of the
12 corresponding weight function φ_{Linear} suggest classifying the gauged sites into three classes: a
13 class of sites to be excluded from the regional estimation procedure; another class of
14 “intermediate” sites for which the contribution is gradually employed; and a last class of sites to
15 be fully included in the estimation.

16
17 In the following, selected detailed results are presented. Figure 5 presents the evolution of
18 the performance criteria as a function of the iteration number for methods (I), (II) and (III.c). For
19 the uniform (I) and the traditional CCA (II) approaches, the criteria are represented by straight
20 lines, since they are independent of the iteration number. For both QS10 and QS100, the criteria
21 values of method (III.c) improve with the iteration number. Moreover, they seem to converge,
22 after approximately 15 to 20 iterations, to the corresponding results given in Table 2. This
23 illustrates the superiority of the depth-based approach over the traditional CCA approach. Note

1 that the computer running time of all the 40 iterations of (III.c) is comparable to that of the
2 traditional CCA approach.

3
4 Figure 6 presents, for each site i , the relative error $(y_i - \hat{y}_i)/y_i$ related to the estimation
5 of the specific quantiles QS10 and QS100 with respect to the basin area. It concerns the results of
6 the CCA approach (II) and those of the 20th iteration of the depth-based approach (III.c). The 20th
7 iteration is selected, as an example, since the criteria converge after 15 to 20 iterations. Using
8 both estimation methods, large negative errors are observed for some sites such as number 46, 64,
9 66 and 148. In particular, sites 64 and 66 have small basin areas. It is generally observed that the
10 relative errors obtained from the depth-based approach (III.c) are smaller than those obtained
11 from the traditional CCA approach (II).

12
13 The following two elements can be used to explain some aspects related to the CCA
14 approach:

- 15 1. Under the normality condition imposed on X and Y , it is implicitly assumed that the
16 conditional canonical variable $(W | V = v_0)$ is $N(\Lambda v_0, I_p - \Lambda^2)$, see Ouarda et al. (2001).

17 Consequently:

$$18 \quad E(W | V = v_0) = \Lambda v_0 \quad \text{and} \quad \text{Var}(W | V = v_0) = I_p - \Lambda^2 \quad (26)$$

19 which suggests modeling the relationship between W and V by a linear regression model
20 as follows:

$$21 \quad W = \Lambda V + \varepsilon, \quad \text{with} \quad \varepsilon \sim N(0, I_p - \Lambda^2) \quad (27)$$

22 Usually, in the CCA approach, sites are presented in the hydrological canonical space
23 $(W1, W2)$ (see Figure 3). Following relation (27), it is of interest to present sites in the

1 canonical spaces (V1,W1) and (V2,W2). This is illustrated in Figure 7 for the considered
2 data set. It shows that the relationship between V1 and W1 can be acceptably considered
3 to be linear, and hence meets the model (27), whereas it is not the case for V2 and W2.
4 The illustration in the space (V1,W1) is useful to get prior information about the
5 estimation error for a given site.

- 6 2. The canonical hydrological value W_0 of the target-site is unknown. Hence, in the CCA
7 approach, the Mahalanobis distance (10) is computed with respect to its expectation Λv_0
8 (it represents also its estimator using (27)). This influences the computation of the
9 Mahalanobis distance, especially when the Mahalanobis distance
10 $(W_0 - \Lambda v_0)' (I_p - \Lambda^2)^{-1} (W_0 - \Lambda v_0)$ between W_0 and Λv_0 is large. This is the case, for
11 instance, for site number 66.

12
13 Figure 8 illustrates the Mahalanobis depths evaluated from (II) and from the 20th
14 iteration of (III.c) for a selection of sites. The sites considered in Figure 8 are selected from
15 Figures 6 and 7 to represent a wide variety of conditions. Indeed, from Figure 6, sites number 66,
16 122 and 148 are estimated with very high relative errors. These sites are located far from the
17 straight line relating V1 and W1 as shown in Figure 7.a. However, sites number 4, 92 and 141 are
18 accurately estimated as shown in Figure 6 and they are located near the regressive line in Figure
19 7.a. Figure 8 shows that the global shapes of the depth values by both methods are similar.
20 However, differences in depth magnitude are observed, especially for site 122. This site has a
21 very high value of basin area covered with lakes (FAL = 43%).

22

1 From the above results and analysis it appears that the first important element, for
2 quantile regression estimation, is the accurate computation of depth values. The second important
3 element is related to the selection of the weight function. Hence, the results of the classical
4 delineation can be improved with accurate depth values (combination IV). However, generally
5 less accurate depth values with an arbitrary weight function can not improve the classical results
6 (combinations V). Consequently, the best combination of these two elements is a smooth weight
7 function along with an iteratively evaluated Mahalanobis depth (combinations III).

8

9 **7. Conclusions and future research directions**

10 The present paper provides an adaptation of the statistical notion of “depth function” in
11 regional flood frequency analysis. The depth-based approach is introduced in order to overcome
12 some drawbacks of the classical methods and to improve their performances. Along with the
13 flexibility and generality of the proposed method, it is shown that its use leads to improvements
14 in traditional methods. Furthermore, the proposed methodology can be useful in other areas and
15 disciplines of water resources where the regression model is applicable.

16

17 Introducing the statistical notion of depth function in regional flood frequency analysis
18 leads to several directions for future work including the following:

- 19 – The determination of an optimal weight function remains an element to be developed. One
20 option is to restrict the optimization problem to a specific parametric class of weight
21 functions, such as φ_{LC} or φ_{Linear} . Then, according to the optimal value of the considered
22 criterion, the coefficients of the weight function can be obtained. More flexible S-shaped

- 1 weight functions (such as the Gompertz function) represent a promising class to be
2 considered for optimization.
- 3 – The estimation in the index flood model (Dalrymple, 1960) should be developed following a
4 similar approach to the regression model. A comparison to the region of influence approach
5 (Burn, 1990) is to be considered. Note that the region of influence approach uses special
6 weight functions.
- 7 – In the present paper, the Mahalanobis depth function is considered for its simplicity and its
8 link to the CCA approach. It is of interest to consider other depth functions and to compare
9 the corresponding results. Note that Lin and Chen (2006) presented an estimator that is
10 similar to the one given by equation (16) using the projection depth function (equation 4).
11 They studied its asymptotic properties, including consistency, normality and robustness.
- 12 – The estimation of the uncertainty associated to regional quantiles represents an important
13 topic that is not often treated in the literature. It would be of value, in future efforts, to study
14 the uncertainty associated to the approach proposed in the present paper. It is important to
15 mention that some of the methods in the literature (such as Generalized and Weighted least
16 squares approaches) allow to carry out an analysis of the uncertainty (see e.g. Madsen and
17 Rosbjerg, 1997). However, with these approaches, the estimation of the uncertainty of the T-
18 year event represents only the uncertainty on the regressive part without integrating the
19 uncertainty associated to the definition of the homogeneous regions. The depth-based method
20 does not suffer from this problem and hence it requires only one global uncertainty evaluation.
- 21 – Information concerning the spatial correlation between the various sites should be integrated
22 in the depth-based flood frequency procedure.

23

1 **Acknowledgments**

2 Financial support for this study was graciously provided by the Natural Sciences and
3 Engineering Research Council (NSERC) of Canada, and the Canada Research Chair Program.

4 The authors wish to thank the Editor in Chief, the Associate Editor and the two anonymous
5 reviewers for their useful comments which led to the improvement of the paper.

6

1 **Bibliography**

- 2 Arcones, M. A.; Chen, Z. and Giné, E. (1994) Estimators Related to *U*-Processes with
3 Applications to Multivariate Medians: Asymptotic Normality. *Ann. Statist.*, **22**, 1460-1477.
- 4 Alila, Y. (1999) A hierarchical approach for the regionalization of precipitation annual maxima in
5 Canada. *J. Geophys. Res.*, **104**, 31,645-31,655.
- 6 Alila, Y. (2000) Regional rainfall depth-duration-frequency equations for Canada. *Water Resour.*
7 *Res.*, **36**, 1767-1778.
- 8 Burn, D. H. (1990) Evaluation of regional flood frequency analysis with a region of influence
9 approach. *Water Resour. Res.*, **26**, 2257-2265.
- 10 Caplin, A. and Nalebuff, B. (1988) On 64%-majority rule. *Econometrica*, **56**, 787-814.
- 11 Caplin, A. and Nalebuff, B. (1991a) Aggregation and social choice: A mean voter theorem.
12 *Econometrica*, **59**, 1-23.
- 13 Caplin, A. and Nalebuff, B. (1991b) Aggregation and imperfect competition: On the existence of
14 equilibrium. *Econometrica*, **59**, 25-59.
- 15 Chebana, F. and Ouarda, T.B.M.J. (2007) Multivariate *L*-moment homogeneity test. *Water*
16 *Resour. Res.*, **43**, W08406, doi:10.1029/2006WR005639.
- 17 Chokmani, K. and Ouarda T.B.M.J. (2004) Physiographical space-based kriging for regional
18 flood frequency estimation at ungauged sites. *Water Resour. Res.*, **40**, W12514, doi:
19 10.1029/2003WR002983.
- 20 Dalrymple, T. (1960) Flood frequency methods. *United States Geological Survey Water-Supply*
21 *Paper*, **1543** A, 11-51.
- 22 Durrans, S. R. and Tomic, S. (1996) Regionalization of low-flow frequency estimates: an
23 Alabama case study. *Water Resour. Bull.*, **32**, 23-37.

- 1 Eaton, B.; Church, M. and Ham, D. (2002) Scaling and regionalization of flood flows in British
2 Columbia, Canada. *Hydrol. Processes*, **16**, 3245– 3263.
- 3 Ghosh, A. K. and Chaudhuri, P. (2005) On Maximum Depth and Related Classifiers. *Scand. J.*
4 *Statist.*, **32**, 327–350.
- 5 GREHYS (1996a) Presentation and review of some methods for regional flood frequency
6 analysis. *J. Hydrology*, **186**, 63-84.
- 7 GREHYS (1996b) Inter-comparison of regional flood frequency procedures for Canadian rivers.
8 *J. Hydrology*, **186**, 85-103.
- 9 Hosking, J. R. M. and Wallis, J. R. (1993) Some statistics useful in regional frequency analysis.
10 *Water Resour. Res.*, **29**, 271-282.
- 11 Kouider, A., Gingras, H.; Ouarda, T. B. M. J.; Ristic-Rudolf, Z. and Bobée, B. (2002) Analyse
12 fréquentielle locale et régionale et cartographie des crues au Québec [In French], Rep. R-627-
13 el, Eau, Terre, et Environ., INRS, Ste-Foy, Que., Canada.
- 14 Lin, L. and Chen, M. H. (2006) Robust estimating equation based on statistical depth. *Statist.*
15 *Papers*, **47**, 263-278.
- 16 Liu, R. Y. (1990) On a notion of data depth based on random simplices. *Ann. Statist.*, **18**, 405-
17 414.
- 18 Liu, R. Y. (1992) Data depth and multivariate rank tests. In *L-1 Statistics and Related Methods* (Y.
19 Dodge, ed.) 279 294. North-Holland, Amsterdam.
- 20 Liu, R. Y. and Singh, K. (1993) A quality index based on data depth and multivariate rank tests.
21 *J. Amer. Statist. Assoc.*, **88**, 257-260.
- 22 Liu, R. Y. (1995) Control Charts for Multivariate Processes. *J. Amer. Statist. Assoc.*, **90**, 1380-
23 1388.

- 1 Liu, R. Y.; Parelius, J. M. and Singh, K. (1999) Multivariate analysis by data depth: descriptive
2 statistics, graphics and inference, (with discussion and a rejoinder by Liu and Singh). *Ann.*
3 *Statist.*, **27**, 783-858.
- 4 Madsen, H. and Rosbjerg, D. (1997) Generalized least squares and empirical Bayes estimation in
5 regional partial duration series index-flood modeling. *Water Resources Research*, **33**, 771-
6 781.
- 7 Mahalanobis, P. C. (1936) On the generalized distance in statistics. *Proc. Nat. Acad. Sci. India*,
8 **12**, 49-55.
- 9 Massé, J.-C. (2004) Asymptotics for the Tukey depth process, with an application to a
10 multivariate trimmed mean. *Bernoulli*, **10**, 1-23.
- 11 Miller, K. ; Ramaswami, S. ; Rousseeuw, P. ; Sellares, T. ; Souvaine, D. ; Streinu, I. and Struyf,
12 A. (2003) Efficient Computation of Location Depth Contours by Methods of Combinatorial
13 Geometry. *Statistics and Computing*, **13**, 153-162.
- 14 Mizera, I. (2002) On depth and deep points: a calculus. *Ann. Statist.*, **30**, 1681–1736.
- 15 Mizera, I. and Müller, C. H. (2004) Location-Scale Depth (with discussion). *J. Amer. Statist.*
16 *Assoc.*, **99**, 949-966.
- 17 Muirhead, R. J. (1982) *Aspect of Multivariate Statistical Theory*. John Wiley, Hoboken, N. J.
- 18 Nguyen, V.-T.-V. and Pandey, G. (1996) A new approach to regional estimation of floods in
19 Quebec. In: Delisle, C.E., Bouchard, M.A. (Eds.), Proceedings of the 49th Annual
20 Conference of the CWRA, June 26–28, Quebec City. Collection Environnement de l'U. de
21 M., 587–596.
- 22 Oja, H. (1983) Descriptive statistics for multivariate distributions. *Statist. Probab. Lett.*, **1**, 327-
23 332.

- 1 Ouarda, T. B. M. J. and Ashkar, F. (1994) Regional multiple regression flood frequency
2 estimation by the Peaks-Over-Threshold method. Internal Report, Department of
3 Mathematics, University of Moncton, Moncton, New Brunswick, Canada, Research report
4 for the Strategic Grant No. STR0118482 of NSERC, 20 pp.
- 5 Ouarda, T. B. M. J.; Haché, M.; Bruneau, P. and Bobée, B. (2000) Regional Flood Peak and
6 Volume Estimation in a Northern Canadian Basin. *ASCE J. Cold. Reg. Engrg.*, **14**, 176-191.
- 7 Ouarda, T. B. M. J.; Girard, C.; Cavadias, G. S. and Bobée, B. (2001) Regional flood frequency
8 estimation with canonical correlation analysis. *J. Hydrology*, **254**, 157-173.
- 9 Ouarda, T. B. M. J.; Cunderlik, J. M.; St-Hilaire, A.; Barbet, M.; Bruneau, P.; Bobée, B. (2006)
10 Data-based comparison of seasonality-based regional flood frequency methods. *J.*
11 *Hydrology*, **330**, 329-339.
- 12 Rencher, A. C. (2002) *Methods of multivariate analysis*. Second edition. Wiley Series in
13 Probability and Statistics. John Wiley & Sons, New York.
- 14 Rousseeuw, P. J. and Hubert, M. (1999) Regression depth. *J. Amer. Statist. Assoc.*, **94**, 389-433.
- 15 Stedinger, J.R. and Tasker, G. (1986) Regional hydrologic analysis, 2, Model-error estimators,
16 estimation of sigma and log Pearson type 3 distributions. *Water Resour. Res.*, **22**, 1487-
17 1499.
- 18 Tukey, J. (1975) Mathematics and picturing data. In *Proceedings of the 1975 International*
19 *Congress of Mathematics*, **2**, 523-531.
- 20 van der Vaart, A.W. (1998) *Asymptotic statistics*. Cambridge Series in Statistical and
21 Probabilistic Mathematics. Cambridge University Press, Cambridge.
- 22 Zuo, Y. and Serfling, R. (2000) General notions of statistical depth function. *Ann. Statist.*, **28**
23 461-482.

- 1 Zuo, Y.; Cui, H. J. and Young, D. (2004) Influence function and maximum bias of projection
- 2 depth based estimators. *Ann. Statist.*, **32**, 189-218.
- 3 Zuo, Y. and Cui, H. (2005) Depth weighted scatter estimators. *Ann. Statist.*, **33**, 381–413.
- 4

1 **Table 1. Descriptive statistics of hydrological, physiographical and meteorological variables**

Variable	Unit	Min	Mean	Max	Standard deviation
MBS	%	0.96	2.43	6.81	0.99
FAL	%	0.00	7.72	47.00	7.99
AMP	mm	646	988	1534	154
AMD	degree-day	8589	16346	29631	5382
AREA	km ²	208	6255	96600	11716
QS10	m ³ /s.km ²	0.03	0.22	0.53	0.13
QS100	m ³ /s.km ²	0.03	0.31	0.94	0.20

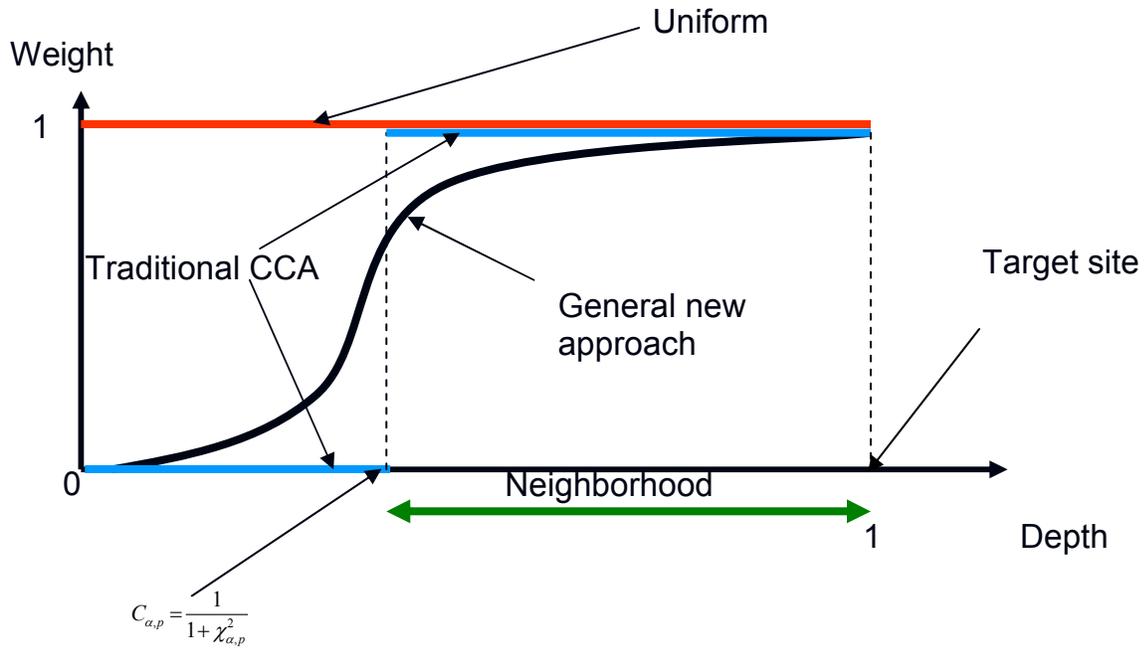
2

1 **Table 2. Quantile estimation results with the various methods**

2

			<i>QS10</i>		<i>QS100</i>		
Combination approach	Depth computation	Weight function	RB (%)	RRMSE (%)	RB (%)	RRMSE (%)	
(I)	Uniform	-	φ_U	-8.60	54.94	-10.93	64.05
(II)	Traditional CCA	CCA	φ_{CCA} optimal α	-7.54	44.62	-8.14	51.84
(III)	Depth-based	Iteration					
	a)		φ_Z K = 200, C = 0.51	-4.37	39.17	-3.51	44.92
	b)		φ_{LC} s=5, K = 200, C = 0.51	-4.47	38.78	-3.75	44.53
	c)		φ_{LC} s=2, K = 100, C = 0.52	-5.22	38.94	-4.68	44.72
	d)		φ_{Linear} d ₁ = 0.30, d ₂ = 0.80	-3.83	38.78	-2.52	44.57
(IV)	CCA with iteration	Iteration	φ_{CCA} optimal α	-5.52	42.27	-4.98	49.02
(V)	Depth-based without iteration	CCA					
	a)		φ_{LC} s = 2, K = 100, C = 0.30	-7.75	44.81	-8.38	51.82
	b)		φ_{LC} s = 2, K = 100, C = 0.52	-8.26	47.63	-9.88	57.02

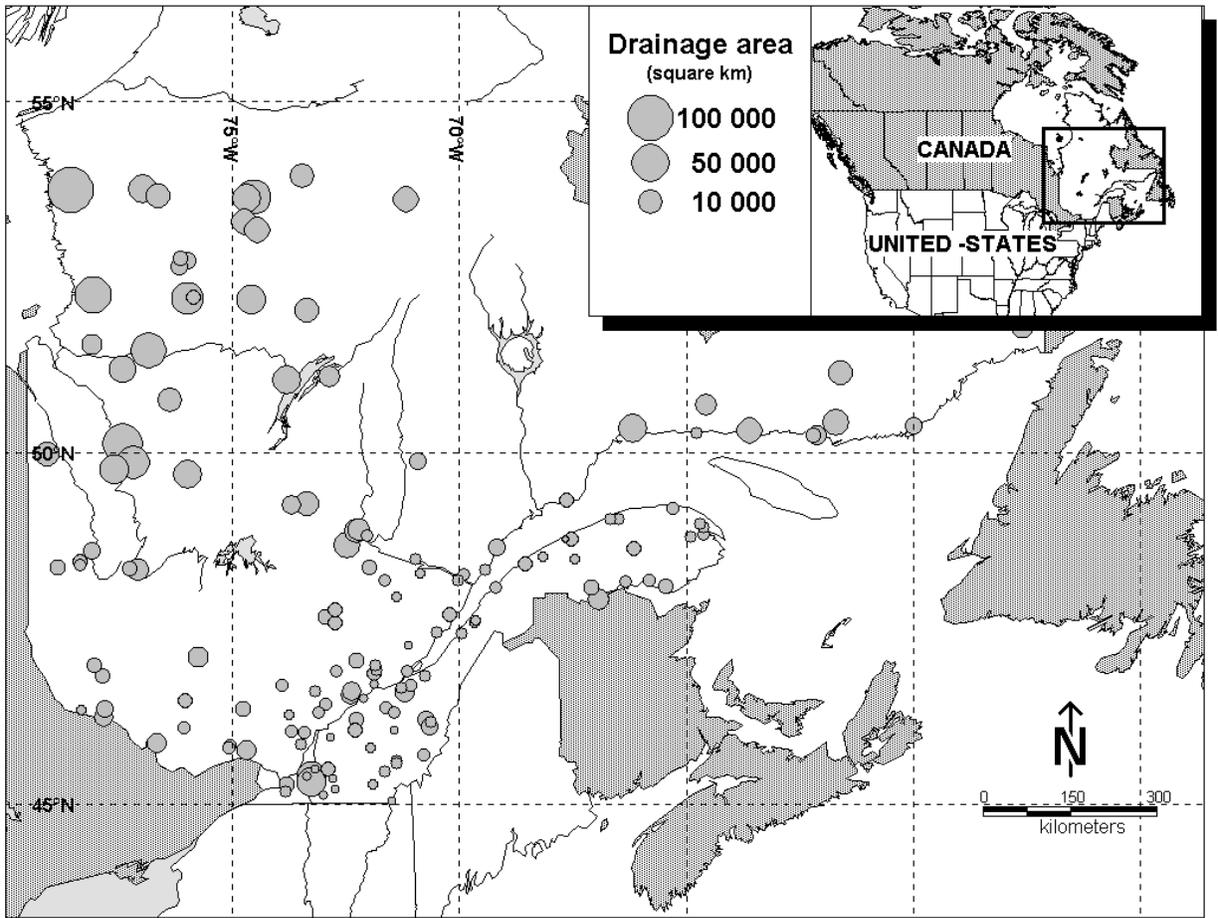
3



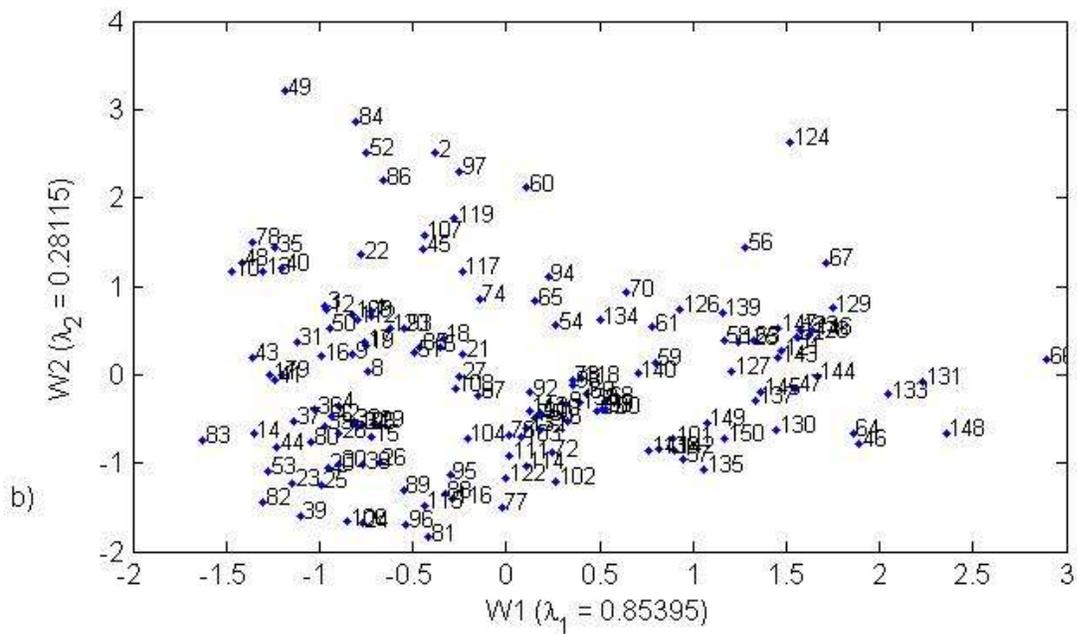
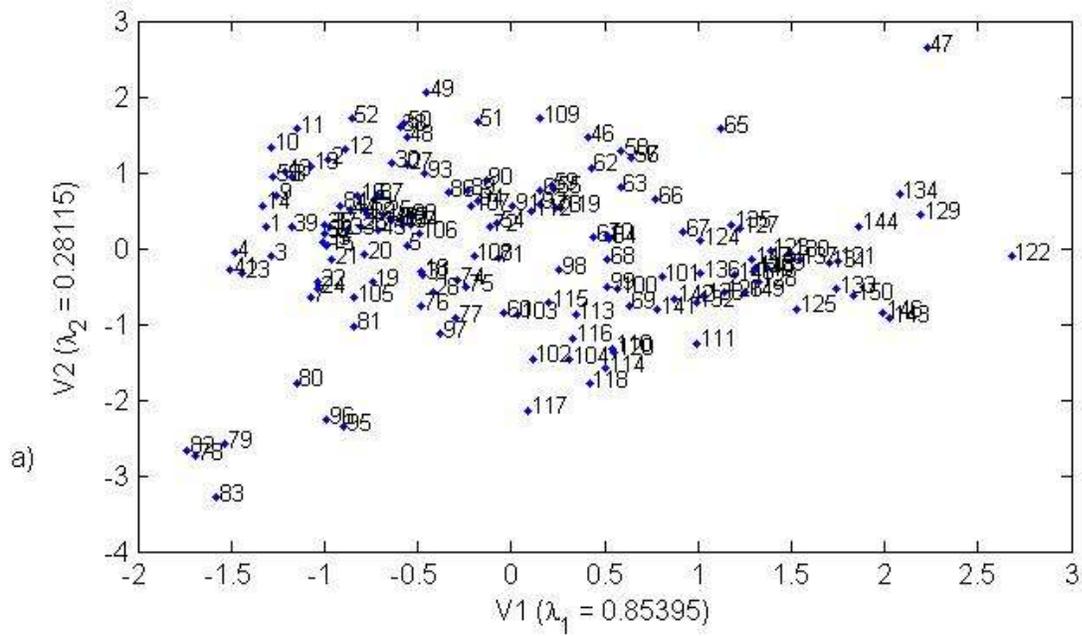
1

2 **Figure 1. Illustration of different regionalization approaches with the corresponding weight**

3 **functions**



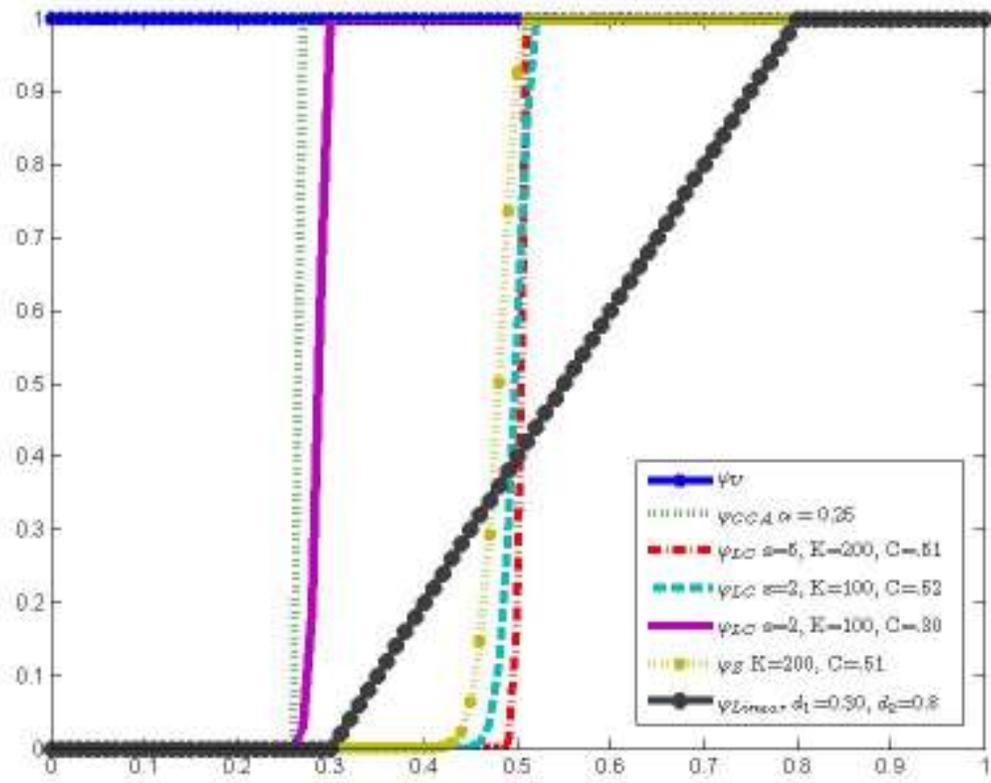
1
 2 **Figure 2. Geographical location of the studied sites in the province of Quebec, Canada**
 3



1

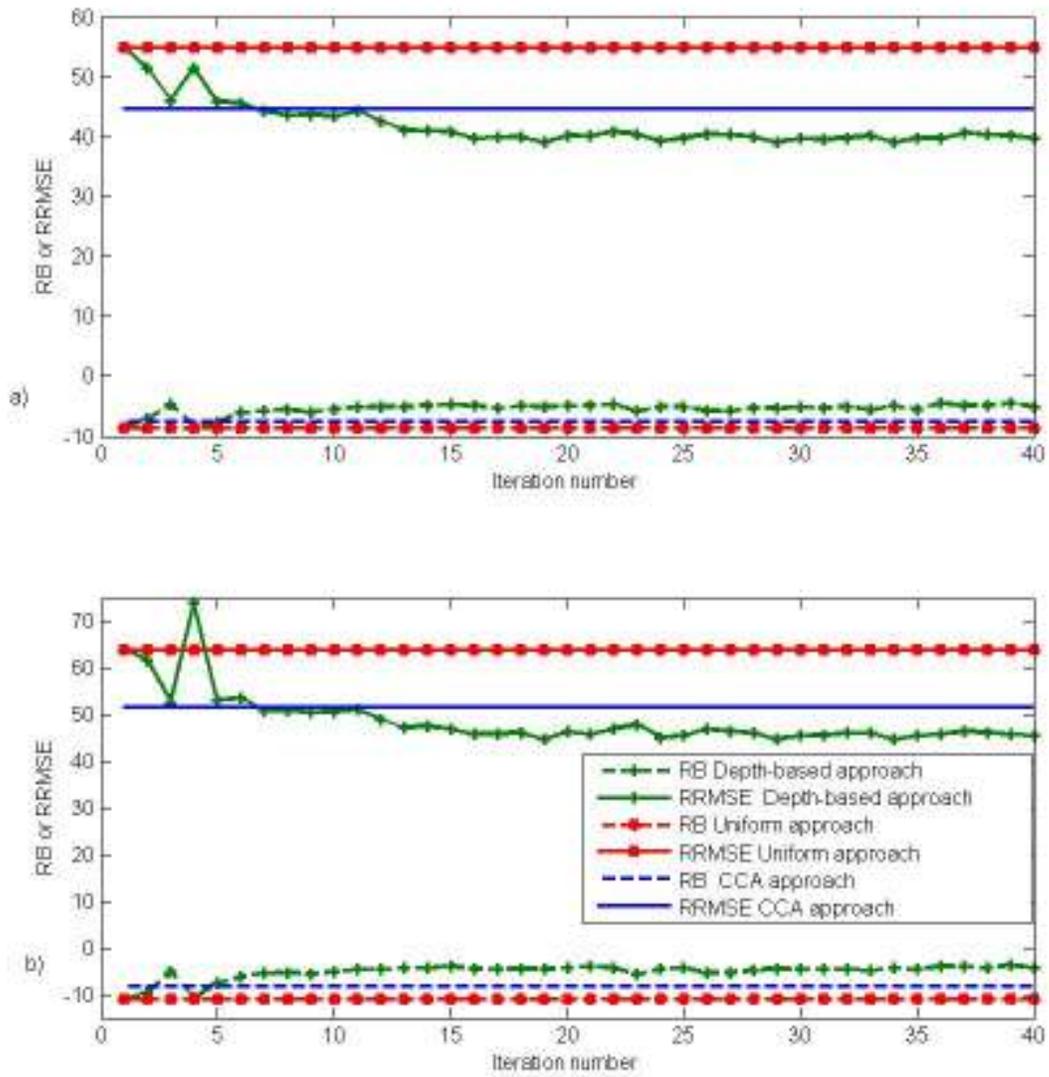
2 **Figure 3. Data set in the canonical spaces (a) physio-meteorological (b) hydrological**

3

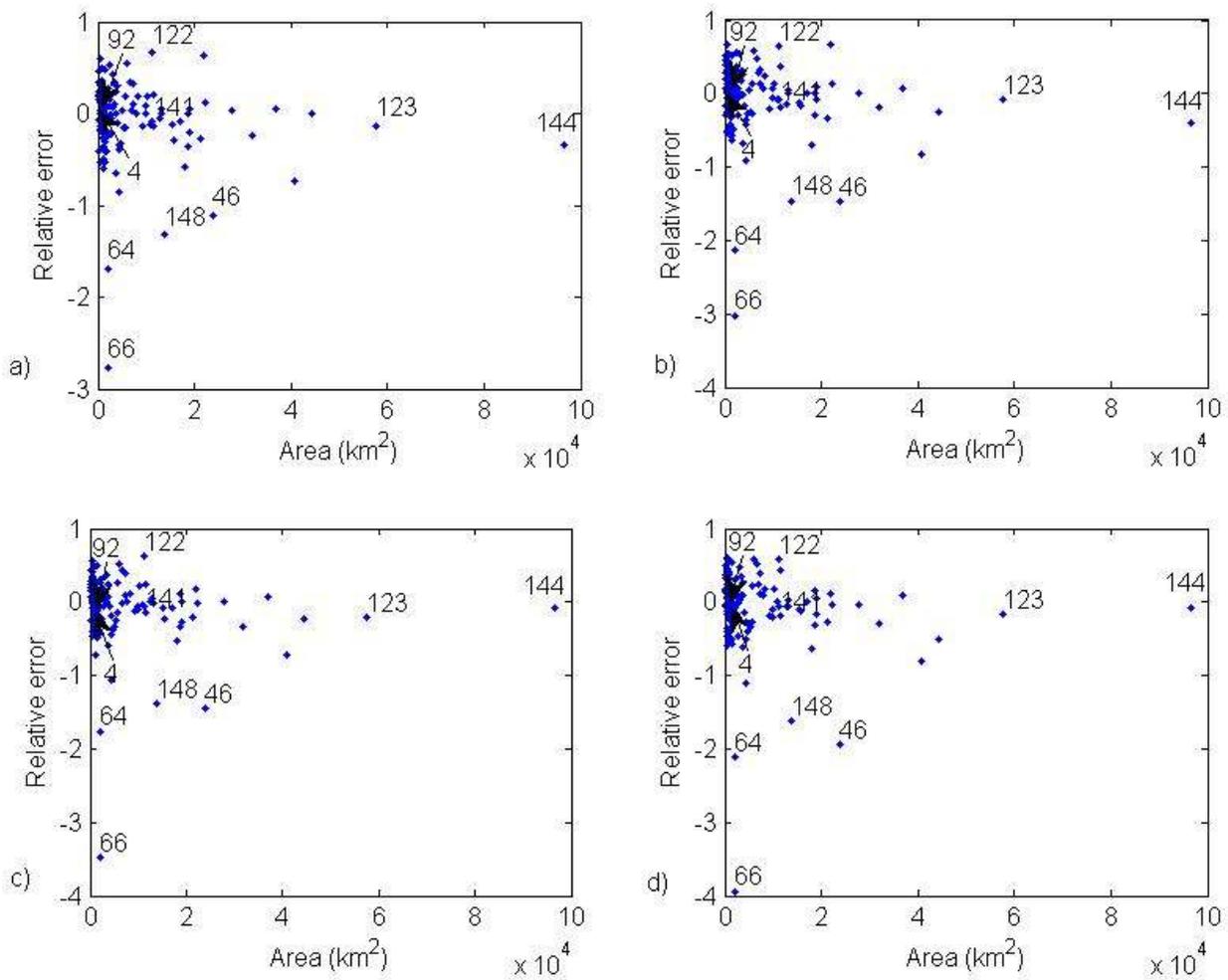


1
2
3

Figure 4. Weight functions used in the comparative study

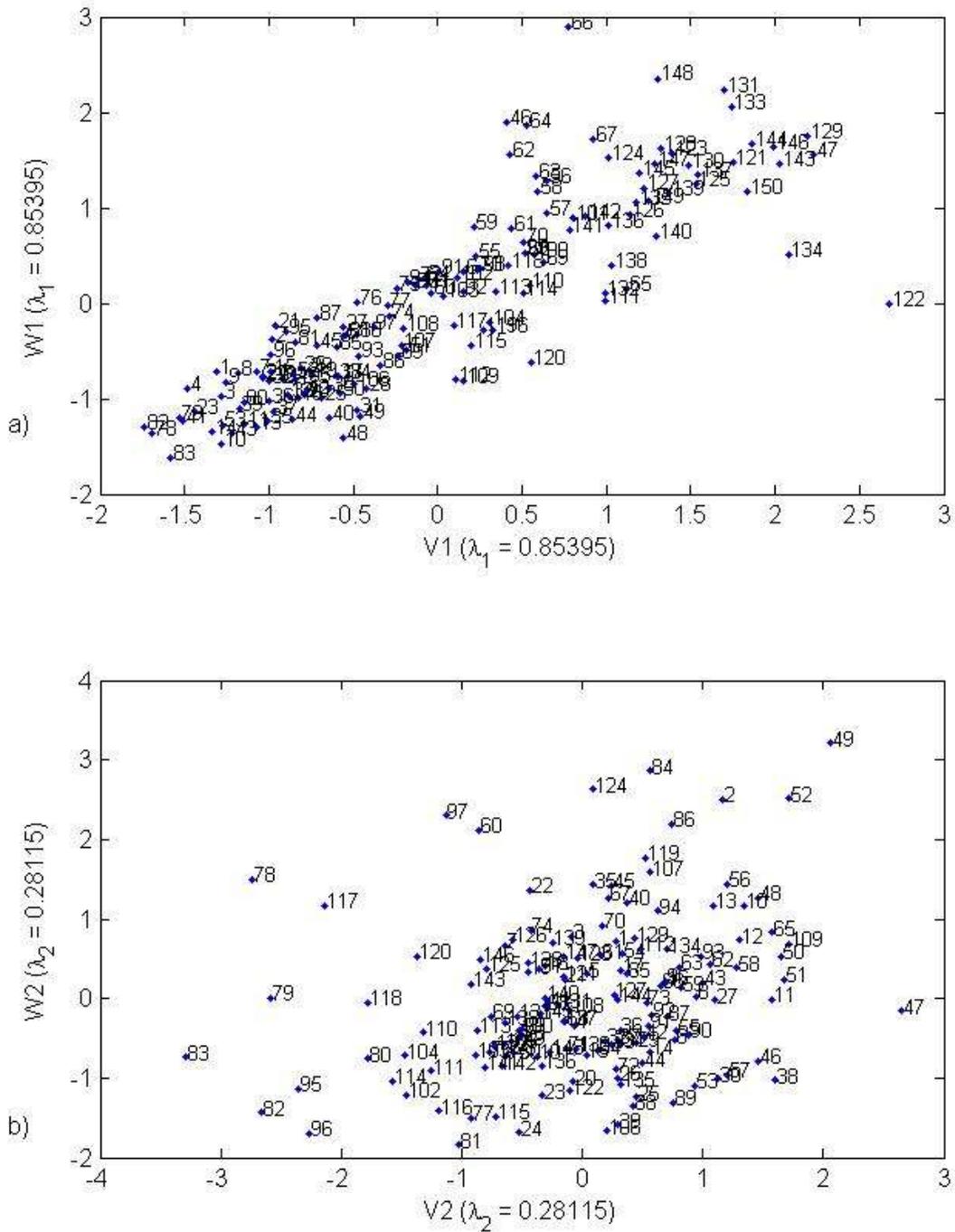


1
 2 **Figure 5. RB and RRMSE using the uniform (I), the CCA (II) and the depth-based (III.c)**
 3 **methods for the estimation of (a) QS10 and (b) QS100**

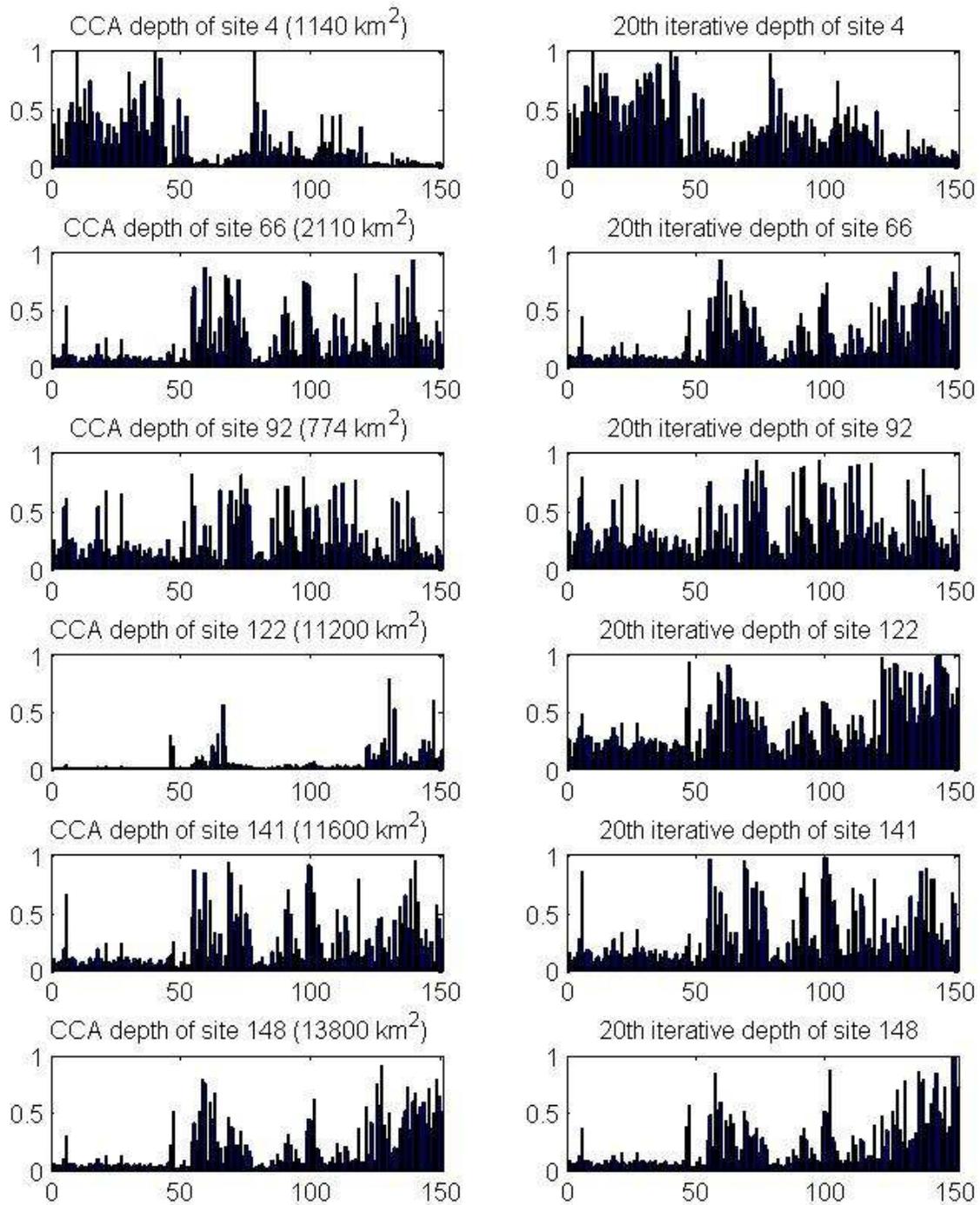


1
2 **Figure 6. Quantile estimation relative errors of: a) QS10 from the 20th iteration of the**
3 **method (III.c), b) QS100 from the 20th iteration of the method (III.c), c) QS10 from**
4 **the CCA method (II), d) QS100 from the CCA method (II). Selected sites are**
5 **indicated by their respective numbers**

6



1
 2 **Figure 7. Scatter plot of the data set in the canonical spaces (V1,W1) and (V2,W2)**



1
 2 **Figure 8. Depth values from the traditional CCA (II) and from the 20th iteration of the**
 3 **method (III.c) for a selection of sites. Values between parentheses represent the area of**
 4 **the corresponding site**