

1 Comparison of catchment grouping methods for flow 2 duration curve estimation at ungauged sites in France

3 E. Sauquet¹ and C. Catalogne¹

4 [1]{Cemagref, UR HHLY, 3 bis quai Chauveau - CP 220, F-69336 Lyon, France}

5 Correspondence to: E. Sauquet (eric.sauquet@cemagref.fr)

6 7 **Abstract**

8 The study aims at estimating flow duration curves (FDC) at ungauged sites in France and
9 quantifying the associated uncertainties using a large dataset of 1080 FDCs. The
10 interpolation procedure focuses here on 15 percentiles standardised by the mean annual
11 flow, which is supposed to be known at each site. In particular, this paper discusses the
12 relevance of different catchments grouping procedures on percentiles estimation by regional
13 regression models.

14 First, five parsimonious FDC parametric models were tested to approximate FDCs at gauged
15 sites. The results show that the model based on Empirical Orthogonal Functions (EOF)
16 expansion outperforms the other ones. In this model each FDC is interpreted as a linear
17 combination of regional amplitude functions with weights – the parameters of the model -
18 varying in space. Here, only one amplitude function was found sufficient to fit well most of the
19 observed curves. Thus the considered model requires only two parameters to be estimated
20 at ungauged locations.

21
22
23
24
25
26
27
28
29
30

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

Second, homogeneous regions were derived according to hydrological response on one hand, and geological, climatic and topographic characteristics on the other hand. Hydrological similarity was assessed through two simple indicators: the concavity index (*IC*) that represents the shape of the standardized dimensionless FDC and the seasonality ratio (*SR*) which is the ratio of summer and winter median flows. These variables were used as homogeneity criteria in three different methods for grouping catchments: (i) according to their membership in one of an *a priori* French classification into Hydro-Eco-Regions (HERs), (ii) by applying a regression tree clustering and (iii) by using ~~hydrological~~ neighbourhood obtained by canonical correlation analysis.

Finally, regression models between physiographic and/or climatic variables and the two parameters of the EOF model were derived considering all the data and thereafter for each group obtained through the tested grouping techniques. Results on percentiles estimation in cross validation show a significant benefit to form homogeneous regions before developing regressions, particularly when grouping methods use hydrogeological information.

Key words:

Flow duration curve, regional model, hydrological neighbourhood, France, Empirical Orthogonal Function

1 Introduction

A Flow Duration Curve (FDC) is the cumulative frequency distribution of observed flows during a period of interest (month, season, year, or entire period of record). It plots specified flows against their corresponding probability of exceedance that can be also interpreted as the percent of time these specified values are equalled or exceeded. FDC is a commonly used tool in water management applications, since it displays the full range of flows, including low flows and flood events (Vogel and Fennessey, 1995; Smakhtin, 2001). Here long-term flow duration curves were considered and derived from observed daily flows available at each site.

There have been numerous approaches for estimating FDC characteristics at ungauged locations, particularly low-flow percentiles, using regression equations under different climates (see Castellarin *et al.* (2007) for a recent review). Despite their interest for water management issues FDCs have until now received very little attention in France. The present study is to our knowledge the first attempt to develop regional flow duration models in this

1 country. Previous works have concentrated on mapping mean river flow statistics including
2 long-term mean annual and monthly flows (Sauquet, 2006; Sauquet *et al.*, 2008). ~~These~~
3 ~~results cannot be ignored.~~ A straightforward method for taking benefits from knowing the
4 mean annual flow qa is to consider percentiles expressed as proportions of the long-term
5 mean flow of the corresponding catchment as variables of interest. Regionalization can thus
6 focus on the shape of the FDC. The dimensionless FDC and the mean annual flow qa are
7 estimated separately and their combination provides the expected percentiles.

8 This approach, known as “index flow approach”, has been previously adopted by numerous
9 authors (e.g., Holmes *et al.*, 2002; Singh *et al.*, 2001; Castellarin *et al.*, 2004; Ganora *et al.*,
10 2009) leading to various procedures to estimate normalised percentiles. The simplest model
11 assumes that the shapes of the FDC at all sites within the study area ~~are~~ show a low
12 variability identical. In practice, dimensionless FDCs from monitored catchments within the
13 same region are pooled and averaged to create the representative shape. Since the
14 hypothesis of similarity may be too restrictive, the alternative way has been chosen here: a
15 reliable mathematical model with few parameters, which vary in space and are estimated at
16 gauging stations, approximates the dimensionless FDC. The main advantages of the
17 adopted approach are:

- 18 - ~~It~~ The choice of the index value ensures ~~consistency between river flow~~
19 ~~statistics~~ percentiles to be consistent with the mean annual runoff at ungauged sites, i.e.
20 estimates are expected to be in the range of qa (qa and percentiles) through the choice
21 of the index value.
- 22 - The only few parameters involved in the procedures can be easier to interpret and their
23 small number ~~It~~ reduces the computational effort in each step of the regionalisation
24 procedure ~~number of steps in the regionalisation procedure (only few parameters are~~
25 ~~involved in the procedures).~~
- 26 - It enables to distinguish the part related to the water balance (*i.e.* qa) from the
27 characteristic response of the catchment to climate to rainfall (*i.e.* the parameters of the
28 shape of the dimensionless FDC) and thus to better identify the most important sources
29 of spatial variability of FDC properties.

30 The last step of the procedure involved empirical relationships between the variables of
31 interest and basin descriptors. Indeed this approach is by far the most often employed in
32 regionalisation. In practice empirical formulas, usually established by ~~multivariate~~ multiple
33 regression, may perform poorly when applied at large scale due to high variability of
34 hydrological behaviours, providing estimates with large errors. A way to improve the
35 performance is to delineate homogeneous subregions assuming that pooled river

1 catchments with similar hydrological, physiographical and meteorological characteristics will
2 behave in a similar manner before developing separate regional regressions (Smakhtin,
3 2001).

4 The identification of homogeneous regions - both in theory and practice - has received much
5 attention in hydrology, but no general methodology has emerged. Hence different ways to
6 form homogeneous regions can be found in the literature, leading to fixed geographically
7 regions (either spatially contiguous or not) or ~~hydrological~~-neighbourhoods around each
8 target site. In the neighbourhood approach, each site is supposed to have its own
9 homogeneous region formed by gauging stations. Examples of contiguous regions defined
10 for estimating regional FDCs are provided by (Singh *et al.*, 2001) in the Himalayan region of
11 India based on a pre-existing partition into hydrometeorological subregions, and by (Laaha
12 and Blöschl, 2006a) in Austria where grouping according to seasonality indices was tested.
13 Geographically non-contiguous regions are usually identified using multivariate techniques
14 such as multiple regression, principal component analysis or classification procedures, all of
15 them incorporating catchment characteristics as well as flow statistics (e.g., Isik and Singh
16 (2008) in Turkey; Nathan and MacMahon (1990) in Australia; Laaha and Blöschl (2006a,
17 2006b) and Laaha *et al.* (2009) in Austria, Vezza *et al.* (2010) in Italy and Ganora *et al.*
18 (2009) in northwestern Italy and Switzerland). Two main neighbourhood methods are
19 commonly used. Both used auxiliary variables to define a hydrological catchment descriptors
20 space where distances are computed: the region of influence developed by Burn (1990a,b)
21 (e.g., Holmes *et al.* (2002) in the UK) and the canonical correlation analysis (CCA) promoted
22 by Ouarda *et al.* (2001).

23 Since the *a priori* efficiency of the grouping methods for regionalizing FDC characteristics is
24 unknown, we here assess the relative performance of three of them: (i) contiguous regions
25 obtained manually from expertise; (ii) regions obtained through Classification and Regression
26 Trees algorithm (CART) and (iii) neighbourhood based on canonical correlation analysis
27 (CCA). The choice of these methods was motivated (i) by a pre-existing partition established
28 in France to answer some basic questions related to the European Water Framework
29 directive, (ii) by published works demonstrating the potential of CART models in river flow
30 regime regionalisation in France (Snelder *et al.*, 2009) and (iii) by the wish to test a well-
31 established method formerly developed to address issues in flood estimation.

32 In this paper we successively investigate two main issues related to the choice of the most
33 adapted parametric model to fit observed dimensionless FDC at gauged sites and the way to
34 define homogeneous regions regardless of the interpolation procedure used to estimate FDC
35 characteristics. Regarding the last point, this study is in line with previous benchmark studies
36 on the performance of different grouping techniques for estimating low flow percentiles

1 (Laaha and Blöschl, 2006b; Vezza *et al.*, 2010). The paper is organised as follows. The study
2 area and data used are first presented in Sect. 2. Hereafter, Sect. 3 compares the various
3 mathematical models tested to approximate FDCs at gauged sites. Once the best performing
4 parametric model has been identified, the variable on which homogeneity is tested are
5 introduced in Sect. 4. Three approaches for delineating homogeneous regions are applied
6 and compared (Sect. 5). The results of the fitted regional regressions are discussed in Sect.
7 6 and some conclusions including future research directions are drawn in the final section.

9 **2 Study area and data**

10 Climate and geology are quite diverse in France (area approx. 550 000 km²): the northern
11 and western parts of France are under maritime temperate climate influences whereas
12 Mediterranean climate with hot and dry summer prevail in the south. In the latter areas,
13 rainfall and evaporation drive the seasonal variations of runoff, in contrast to mountainous
14 areas (high-altitude rivers in both the Pyrenees and the Alps) where snowmelt-fed regimes
15 are observed. From a geological standpoint, France is roughly composed of two major
16 geological formations: Hercynian crystalline impermeable substratum principally located in
17 the north-western part of France (Brittany) and in mountainous areas (Alps, Pyrenees and
18 Massif Central) and more or less permeable sedimentary rocks (limestone and clay) in flat
19 plain areas (e.g., in the northern part of France where large aquifers sustain flows).

20 The dataset (~~Fig. 1~~~~Fig. 1~~~~Fig. 1~~~~Fig. 1~~) consists in 1080 gauging stations among more than
21 3500 stations that are available in the French database HYDRO
22 (<http://www.hydro.eaufrance.fr/>). The following selection criteria were imposed to select these
23 gauging stations: (i) no significant human influence on flow, (ii) high quality of measurements,
24 (iii) record covering at least 18 years during the period 1970-2008 ~~(iii) high quality of~~
25 ~~measurements.~~ To help in the selection process qualitative metadata on the degree of
26 human influence on the river flow regime and on the uncertainty in discharge observations
27 provided by the monitoring authorities were gathered and interpreted. In addition we
28 investigated the presence of major reservoirs and water diversions upstream from the
29 gauging stations. Time-series were also examined to detect abnormal temporal patterns or
30 suspicious values in the data.

31 The final selection corresponds to an average density of about 2 gauging stations per 1000
32 km². The distribution of gauging stations across the country is however not uniform, with two
33 notable areas of low station density located in the northern part of France and south Brittany.
34 A total of 40% of the selected catchments have a record length varying between 35 and 45
35 years in most cases. Continuous observations during the period 1983-2000 are available for

1 90% of all selected stations, which ensures the temporal consistency of runoff statistics in
2 terms of climatic variability. The drainage areas vary in size between 1.4 and 109 930 km².
3 Most of the gauged catchments (44%) have areas from 100 to 500 km².

4 The catchment characteristics selected for use in the delineation of hydrological regions and
5 in the development of regression equations were GIS-derived combining the SAFRAN high-
6 resolution atmospheric reanalysis (Quintana-Seguí *et al.*, 2008; Vidal *et al.*, 2010), a 1-km
7 grid digital elevation model and the associated drainage pattern (Sauquet, 2006). 18
8 catchment characteristics were selected for their possible influence on the shape of the
9 standardised flow duration curve. The variables considered in this study include the drainage
10 area (A), the coordinates of the centre of gravity (XG, YG), the mean catchment slope (Slp),
11 the three quartiles of the hypsometric curve (Z25, Z50 and Z75), the mean annual catchment
12 air temperature (TA), the mean summer catchment potential evapotranspiration (ETsummer)
13 using the formulation suggested by Oudin *et al.* (2005), the mean annual catchment actual
14 evapotranspiration (AETA) according to Turc formulation (1954), the mean annual catchment
15 precipitation (PA), the variance of the twelve mean monthly catchment precipitations
16 (VarPA), the mean seasonal precipitations (Pwinter, Pspring, Psummer and Pautumn), the
17 catchment yield (CY) defined by the ratio (PA-AETA)/qa and the fraction of the drainage
18 catchment with impermeable substratum (%Imp).

19 In addition, we used the Hydro-EcoRegion classification (HER) developed by Wasson *et al.*
20 (2002). The HERs delineation was performed by experts incorporating different aspects of
21 the geology, climate, physiography, drainage density, vegetation and topography of France.
22 In particular, HER is the result of the interpretation in terms of erosion resistance,
23 permeability, and hydrochemistry of a original geological map provided by the Bureau de
24 Recherches Géologiques et Minières (BRGM, 1996). The HER was not specifically
25 developed to discriminate river flow regimes. In the absence of quantitative information on
26 hydrogeology, HERs were considered the most reliable surrogate. This classification divides
27 France into 22 main regions (HER1) that are subdivided into 112 subregions (HER2). The
28 dominant class in terms of fraction of the drainage catchment underlain by each HER was
29 also computed.

30

31 **3 A parametric model for flow duration curve**

32 As suggested in Sect. 1, the identification of parsimonious models for summarizing FDCs is
33 advantageous to reduce the ~~computational effortsnumber of steps~~ in the regionalisation
34 procedure (only few parameters are required at ungauged sites to estimate dimensionless
35 FDC).

1 Numerous formulas have been suggested to approximate FDCs (e.g., Quimpo *et al.*, 1983;
 2 Franchini and Suppo, 1996; Yu *et al.*, 2002; Castellarin *et al.*, 2004; Li *et al.*, 2010). Four
 3 parametric functions including the exponential model (Eq. (1)), the logarithm model (Eq. (2)),
 4 the power law model (Eq. (3)) and the model suggested by Franchini and Suppo (Eq. (4))
 5 were tested on the dataset in this study. They approximate FDC at each site $i, i= 1, \dots, N$:

$$6 \quad Q_p(i) = b(i) e^{a(i)p} \quad (1)$$

$$7 \quad Q_p(i) = b(i) + a(i) \ln(p) \quad (2)$$

$$8 \quad Q_p(i) = b(i) p^{a(i)} \quad (3)$$

$$9 \quad Q_p(i) = b(i) + a(i) (1 - p)^{c(i)} \quad (4)$$

10 where Q_p is the p^{th} standardized dimensionless flow percentiles and $a(i)$, $b(i)$ and $c(i)$ are the
 11 parameters at location i .

12 In addition to these four analytical functions, we tested a different approach based on the
 13 discrete decomposition into Empirical Orthogonal Functions expansion (Holmström, 1963).
 14 This mathematical technique, also known as the Karhunen-Loeve transform, aims at
 15 extracting common patterns that represent a large fraction of the variability contained in a
 16 sample of N time series. EOF analysis has been already used for several purposes in
 17 hydrology (e.g., Hisdal and Tveito, 1991; Braud and Obled, 1991; Krasovskaia *et al.*, 1999).
 18 In this application, EOF analysis expresses logarithmically transformed FDC as a linear
 19 combination of M -shape functions β :

$$20 \quad \ln(Q_p(i)) = \gamma(i) + \sum_{m=1}^M \alpha_m(i) \beta_m(p), i = 1, \dots, N \quad (5)$$

21 where M is the number of flow percentiles describing the FDCs, N is the number of gauging
 22 stations, $\alpha_j(i), i= 1, \dots, N$ β_m is the m -th shape function and α_m is the weight associated with
 23 each m -th shape function. By definition $\beta_m, m= 1, \dots, M$ are M orthogonal functions with zero
 24 mean. This constraint leads to introduce the additional term:

$$25 \quad \gamma(i) = \sum \ln(Q_p(i)) / M \quad (6)$$

26 $\alpha_m, m= 1, \dots, M$ and $\gamma(i)$ are the parameters of the EOF model depending on the location of
 27 the site i and have to be estimated at ungauged sites. ~~are weights which vary with location~~
 28 ~~and β , are orthogonal functions with zero mean. Transforming~~ Note that the raw data has
 29 been logarithmically transformed ~~adopted~~ to avoid negative unrealistic estimates. The

1 interest in applying this method is to keep the most part of the dataset variance in a limited
2 number of shape functions. It is thus possible to truncate the series expansion to a subset of
3 $L < M$ functions to limit the number of model parameters without significant loss of information.

4 In this study, all models were calibrated using 15 standardized dimensionless percentiles Q_p ,
5 with respective exceedance probabilities $p = 1, 2, 5, 10, 20, 30, 40, 50, 60, 70, 80, 90, 95, 98$
6 and 99% of the observed FDCs. Analytical models parameters were optimised on
7 observations by applying ordinary least square procedures on logarithmically transformed
8 data to reduce the influence of the largest observations. Prior to optimization,
9 standardized dimensionless percentiles equalled to zero were replaced by 0.001 to apply the
10 logarithmic transformation.

11 The EOF decomposition applied on the dataset provides fourteen shape functions
12 characterized by different patterns. The first shape functions, with a contribution of 97.2% to
13 the total variance, represent the most common pattern of French FDCs. The other shape
14 functions stand for a negligible part of the total variance and allow readjustment for very
15 particular FDCs patterns. Considering these results, it was decided to keep only the first
16 shape function. Thus the number of the parameters for the EOF model is limited to two: the
17 mean of the log-transformed standardized dimensionless percentiles $\gamma_{\overline{\ln(Q)}}$ and the weight
18 associated with the first shape function α_1 .

19 The performance/uncertainty of each model was measured by the deviations from the 15
20 standardized dimensionless percentiles Q_p on which the five models are fitted. Unrealistic
21 values (negative) were also replaced by 0.001. Boxplots in Fig. 2~~Fig. 2~~Fig. 2~~Fig. 2~~ give a
22 graphical overview of the performance of each model. The median and the whiskers of the
23 boxplots measure the bias and the accuracy of the model, respectively. In addition, the fitted
24 curves are displayed on Fig. 3~~Fig. 3~~Fig. 3~~Fig. 3~~ for four gauged catchments representative of
25 the diversity of FDC patterns within the reference dataset. Results show that:

- 26 - None of the models are perfect; in particular, all the models fails to reproduce
27 correctly low-flow percentiles (relative errors may exceed 150% for some
28 catchments). One should note that this criterion is very selective for low values
29 (relative errors may reach large values when estimates are divided by a reference
30 value close to zero).
- 31 - The biases appear most pronounced for the power law model (Eq. (3)); low-flow
32 percentiles as well as high-flow percentiles tend to be largely overestimated.
- 33 - Comparable biases are found for the exponential model (Eq. (1)) and the Franchini
34 and Suppo model (Eq. (4)): standardized dimensionless percentiles Q_p are

1 underestimated for $p \leq 0.02$ and for $0.7 \leq p \leq 0.9$ whereas Q_p are overestimated for p
2 ≥ 0.98 and for $0.1 \leq p \leq 0.4$.

- 3 - The relative error range is smaller for the exponential model (Eq. (1)) and the
4 Franchini and Suppo model (Eq. (4)) for the two standardized dimensionless
5 percentiles ($p = 0.01, 0.02$). However, there is a systematic negative bias in estimated
6 high-flow standardized dimensionless percentiles.
- 7 - Results for the logarithm model (Eq. (2)) follow a very similar pattern to those for the
8 EOF model (Eq. (5)): on average, they both overestimate standardized dimensionless
9 percentiles with $0.4 \leq p \leq 0.8$ while high-flow and low-flow percentiles are
10 underestimated.

11 The degree of bias differs substantially depending on the fitted model. The power law model
12 (Eq. 3) provides the worst estimates in terms of relative error (bias and spread are the largest
13 among the models). Comparable biases are found for the exponential model (Eq. 1) and the
14 Franchini and Suppo model (Eq. (4)). The EOF model (Eq. (5)) appears to outperform among
15 the other models tested despite poor performance for high-flow percentiles. It performs
16 nearly as well as the logarithm model (Eq. (2)) but it also produces globally less biased
17 estimates (median relative errors are the closest to zero and most of the interquartile ranges
18 include zero for all the exceedance probabilities). The advantage of the EOF model is
19 probably a better flexibility (the other models are not enough flexible to reproduce possible
20 inflexion points in the observations) as it results from an empirical modelling of the shapes of
21 the FDC. Considering these results we finally kept the EOF model ~~is the only one to be kept~~
22 in the following steps. As an illustration Fig. 4 Fig. 4 Fig. 4 displays the spatial pattern of the
23 weight coefficient α_1 . The right panel shows how the shape of the FDC approximated by the
24 EOF model evolves as α_1 changes with γ fixed to zero. High values for α_1 correspond to
25 steep slopes of the FDC observed mainly along the Mediterranean and North-Atlantic coasts
26 whereas small values correspond to flat slopes of the FDC observed in the north part of
27 France where the river flow regime is governed by groundwater dynamics.

28

29 **4 Variables for testing hydrological homogeneity**

30 The application of grouping methods is conditioned by the prior definition of variables to
31 measure the degree of similarity between catchment behaviour and the level of homogeneity
32 within the region. The most obvious option would have been to derive groups based on the
33 two variables γ and α_1 to be interpolated. Nevertheless this choice is not optimal since these
34 values result from an approximation of FDCs. In addition working on empirical variables

1 independent of any analytical model was preferred for latter applications of the obtained
2 clusters. Several possible characteristics directly derived from river flow time series were
3 tested and two variables were finally chosen for their correlation to the shape of the FDC and
4 for their interpretation in terms of underlying hydrological processes.

5 The first variable is directly related to empirical properties observed on FDCs. The analysis of
6 observed FDCs suggests that the 10th percentile is a breakpoint delineating two parts of the
7 curves: gradient tends to be higher in the upper branch (10% < p < 99%) than in the lower
8 branch (1% < p < 10%). On this basis, a concavity index is computed as follows:

$$9 \quad IC = \frac{Q_{10} - Q_{99}}{Q_1 - Q_{99}} \quad (6)$$

10 This descriptor is a measure of the contrast between low flow and high flow regime. A map of
11 the concavity index in France including the location of the selected stations is presented in
12 Fig. 5Fig. 5Fig. 5Fig. 4. The parameter takes values between 0 and 1. Values close to 1 are
13 observed where large aquifers (e.g., in the northern part of France) and storages in snow
14 pack (e.g., in the mountainous area) moderate the variability of daily flow. Values close to 0
15 are related to catchments exposed to contrasted climate (e.g., small catchments in the
16 Mediterranean area experiencing hot and dry summers and intense short rainy events in
17 autumn) and also to catchments with no storage capacity (e.g., on impermeable substratum)
18 resulting in severe low-flows and quick runoff responses to rainfall events. It is worth noting
19 that IC is well correlated with the parameters of the analytical FDC models (Fig. 4Fig. 4Fig.
20 4) and the average base flow index as well (not published here).

21 The second variable is a seasonality index. Laaha and Blöschl (2006a) demonstrated the
22 value of such a variable for regionalizing the low-flow percentile Q_{95} in Austria. Indeed,
23 grouping based on seasonality indices performed better than alternative groupings since
24 these indices enable to discriminate well low flow processes at the regional scale when
25 seasonal variability of runoff is high. Laaha and Blöschl (2006a) have used the ratio of the
26 95th percentile of the winter (December to March) FDC divided by the 95th percentile of the
27 summer (April to November) FDC. Since our objective encompasses low flows, a
28 Seasonality Ratio (SR) based on the medians was used here instead:

$$29 \quad SR = Q_{50}(\text{summer})/Q_{50}(\text{winter}) \quad (7)$$

30 $SR \approx 1$ relates to catchments with nearly uniform flows through the year, often when
31 significant groundwater contributions filter out seasonal climatic variability. Catchments
32 influenced by snowmelt-fed processes display $SR < 1$ whereas for typical rainfall-fed
33 catchments with low flow in summer and high flow in winter SR is above 1. SR is used here
34 as a complement to IC to better identifying the causes of low seasonal variability in runoff

1 (snow or groundwater storages). The variation in *SR* is governed by geology and air
2 temperature and consequently in France by topographic influences.

3 These two variables *IC* and *SR* are the flow characteristics used to delineate homogeneous
4 groups. Methods and results are presented in the subsequent section.

5

6 **5 Grouping methods**

7 **5.1 Methods**

8 **5.1.1 Visual grouping (VG)**

9 Non-overlapping regions of approximately homogeneous low-flow indices *SR* and *IC* have
10 first been identified visually. The starting point was the partition of France into 112 Hydro-
11 EcoRegions (HER2s) at the finest level ([Wasson et al., 2002](#)). These HER2s, [introduced in](#)
12 [Sect. 2](#), have been pooled based on hydrological expert knowledge.

13 The boundaries of HER2s have been first superimposed to the map displayed in [_Fig. 4Fig.](#)
14 [5Fig. 5Fig. 5](#). The most similar neighbouring HER2s have been progressively pooled by
15 respecting contiguity, minimizing the dispersion within each cluster and maximizing the
16 dissimilarity between the clusters based on visual inspection. The pooling process is far from
17 obvious. In particular, due to the uneven density of the reference network, some of the
18 HER2s contain too few stations to relate undoubtedly them to other neighbouring HER2s.
19 Hence we used additional information such as rough description of hydrogeology to merge
20 the ungauged HER2s with one of the adjacent clusters. Lastly, inspection of *SR* values led to
21 a partition of the preliminary groups into sub-groups of HER2s, homogenous in terms of
22 seasonality.

23 [Fig. 6Fig. 6Fig. 6Fig. 5](#) presents the division of France into 18 different regions so obtained.
24 Mixed regions may persist due to the heterogeneity at HER2 scale or due to the merging of
25 HER2s containing a small number of gauged sites to large clusters. The identified regions
26 include from 21 to 138 gauged sites and the average size is 57 (5% of the dataset).

27 **5.1.2 Regression Tree (RT)**

28 The aim of the analyses via tree-building algorithms is to predict dependent variables from a
29 set of factor effects. Classification and Regression Trees approaches perform successive
30 binary partitions of a given dataset according to decision variables. One advantage of this
31 method is its ability to handle qualitative data (e.g., membership to a specific class). In
32 general, RT leads to a set of if-then logical conditions as basis for classification. The
33 algorithm identifies the best possible predictors, starting from the most discriminating factors

1 and proceeding to the less important controls, to divide the clusters (*nodes*) into two
 2 successive parts. The optimal choices are determined recursively by increasing the
 3 homogeneity within the two resulting clusters. In this application the R software package
 4 *rpart* (Therneau and Atkinson, 2010) was used. The decision variables were selected
 5 automatically by the algorithm among the 19 catchment descriptors (*i.e.* including the
 6 dominant HER2) to ensure an optimal homogeneity of *IC* chosen as the dependent variable,
 7 in the successive clusters. The only constraint was to include at least 30 gauging stations in
 8 each region. At last 22 hydrological regions were identified with a mean number of 54
 9 gauging stations per region ([Fig. 7](#)[Fig. 7](#)[Fig. 7](#)[Fig. 6](#)).

10 5.1.3 Canonical Correlation Analysis (CCA)

11 Canonical Correlation Analysis (Hotelling, 1936) is a multivariate statistical method suited to
 12 study interrelations between two sets of variables. CCA has been previously suggested by
 13 Ouarda *et al.* (2001) as a neighbourhood definition method. CCA provides two sets of
 14 canonical variables $V_j, j = 1, \dots, k$ and $W_j, j = 1, \dots, k$ obtained as follows:

- 15 - $V_j, j = 1, \dots, k$ are linear combinations of k standardized hydrological variables
 16 $X_j, j = 1, \dots, k$.
- 17 - $W_j, j = 1, \dots, k$ are linear combinations of r standardized physiographic and climatic
 18 characteristics of the catchment $Y_j, j = 1, \dots, r$ ($k < r$).
- 19 - (V_j, W_j) have maximum correlation.
- 20 - $(V_i, V_j), (V_i, W_j)$ and (W_i, W_j) ($i \neq j$) are uncorrelated.

21 Theoretical developments show that the weight for V_i (resp. for W_i) is the i -th eigenvector
 22 $\Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{YY}^{-1} \Sigma'_{XY}$ (resp. $\Sigma_{YY}^{-1} \Sigma'_{XY} \Sigma_{XX}^{-1} \Sigma_{XY}$) where Σ_{XY} is the $k \times r$ covariance matrix and
 23 Σ'_{XY} the transpose of Σ_{XY} . Canonical variables $V_j, j = 1, \dots, k$ and $W_j, j = 1, \dots, k$ can be
 24 interpreted as coordinates in hydrological and catchment-related physical spaces,
 25 respectively. Knowing $Y_j, j = 1, \dots, r$ at ungauged location it is then possible to compute
 26 $W_j, j = 1, \dots, k$ and through the calculation of correlation coefficients between canonical
 27 variables (V_i, W_j) their possible proximity - according to Mahalanobis distance - to the
 28 gauged stations in the hydrological space, which delineates neighbourhood around each site.
 29 CCA has been formerly applied to regional flood frequency estimation (e.g., Ouarda *et al.*,
 30 2001; Chokmani and Ouarda, 2004; Shu and Ouarda, 2007). The present study is probably

1 one of the first published works on CCA application to predict FDCs at ungauged locations.
2 Here CCA was carried out between the two indicators *IC* and *SR* and all the catchment
3 descriptors (excepted dominant HER2, since traditional CCA cannot manage qualitative
4 variables). Geological description is thus reduced to the percentage of impervious areas. All
5 combinations of 2 to 18 variables among the 18 catchment characteristics detailed in
6 remaining basin descriptors (listed in Sect. 2) were tested and at last we retained a
7 combination of six characteristics which provides to the highest correlations between the first
8 two pairs of canonical variables, i.e. (V_1, W_1) and (V_2, W_2) ($p=2$). These catchment
9 characteristics relate to location (the coordinates of the centre of gravity), climate (the mean
10 annual catchment actual evapotranspiration and the variance of the twelve mean monthly
11 catchment precipitations), geology (the fraction of the drainage catchment with impermeable
12 substratum) and altitude (the third quartile of the hypsometric curve).

13 In addition to the variables involved in CCA, one should define the boundaries of the
14 neighbourhood to exclude gauging stations too far from the target site. Ouarda *et al.* (2001)
15 suggested a distance threshold depending on a given confidence level and on target site.
16 Preliminary tests showed the difficulty to define a satisfactory confidence level for our
17 dataset, in particular for very atypical sites for which too few similar sites are selected to
18 derive, thereafter, reliable regional regressions. Consequently we chose here to fix the
19 number of stations contributing to neighbourhood to 50, *i.e.* the 50 closest gauging stations
20 to the target site, to allow objective comparisons with the results of the two other grouping
21 methods.

22

23 5.2 Results

24

25 Fig. 6~~Fig. 6~~Fig. 6~~Fig. 5~~ and Fig. 7~~Fig. 7~~Fig. 7~~Fig. 6~~ present maps obtained by VG and RT,
26 respectively. One colour is assigned to each reach of the main river network (*i.e.* all locations
27 draining more than 50 km²). Displaying results from CCA on a map is not feasible since each
28 site has its own neighbourhood. The comparison between the two maps suggests that:

- 29 - The two procedures based on the same auxiliary variables lead to different divisions.
30 The spatial pattern provided by RT is patchier than the one obtained by VG: small
31 tributaries may belong to different classes than the main stem they flow into. —The
32 relative influence of the location is naturally moderate on class allocation since
33 mountainous basins in the Alps and Pyrenees are pooled together. This result is in
34 direct line with conclusions of previous studies dedicated to flood quantile estimation

1 (Merz and Blöschl, 2005; Ouarda et al., 2001) that concluded that geographical
2 proximity does not involve hydrological similarity.

- 3 - Common geographical groupings can be found e.g., in the north part of France (in
4 brown in
- 5 - ~~Fig. 6~~~~Fig. 6~~~~Fig. 6~~~~Fig. 5~~ and in cyan in ~~Fig. 7~~~~Fig. 7~~~~Fig. 7~~~~Fig. 6~~) and in the west part of
6 France (in orange in
- 7 - ~~Fig. 6~~~~Fig. 6~~~~Fig. 6~~~~Fig. 5~~ and in dark blue in ~~Fig. 7~~~~Fig. 7~~~~Fig. 7~~~~Fig. 6~~), supporting visually
8 the fact that the two partitions are not totally inconsistent.

9 To supplement this analysis, we examined the empirical distributions of both *SR* and *IC* per
10 regions (identified by a letter on the x-axis). Box plots are presented in ~~Fig. 8~~~~Fig. 8~~~~Fig. 8~~~~Fig. 7~~.
11 There is no obvious difference between the spread of *SR* and *IC*. The absence of
12 significant improvement in terms of homogeneity within each group (regarding the
13 interquartile provided by the empirical distribution of each variable) and in discrimination
14 between groups (regarding the differences between the medians of each groups for each
15 variable) is due to the valuable information contained in the Hydro-EcoRegions. Both
16 methods lead to two very distinct regions with high values for *IC*. As a proof the membership
17 to clusters of HER is chosen as the first splitting variables.

18 Regarding CCA we decided to compare results with published works in terms of correlation
19 structure. Fig. 10~~Fig. 10~~~~Fig. 10~~ indicates moderate correlations between the canonical
20 variables: $r_1 = 0.71$ between W_1 and V_1 and $r_2 = 0.57$ between W_2 and V_2 .

21 These values are lower than those obtained in regional flood quantile estimation by Ouarda
22 et al. (2001) in the Province of Ontario (Canada) (r_1 between 0.959 and 0.960 and r_2 between
23 0.279 and 0.422), by Haché et al. (2002) in the Saint-Maurice river region (Canada) ($r_1 =$
24 0.986 et $r_2 = 0.842$) and by Ouarda et al. (2008) in Mexico ($r_1 = 0.966$ and $r_2 = 0.247$).

25 In these studies the analysis of the weights associated with the hydrological variables X and
26 the catchment descriptors Y in the linear combinations shows that the high correlation rate r_1
27 principally depends on the strong link between one T -year flood quantile QT expressed in
28 m^3/s and the drainage area A . It reflects the dependence of the productivity of the basin in
29 terms of volume to the catchment size. On the contrary correlation coefficient r_2 is very weak
30 in most cases. It illustrates the difficulty in identifying relevant basin descriptors to explain the
31 residual spatial variability. As a result the identification of neighbouring catchments using the
32 Mahalanobis distance leads to cluster catchments of equivalent size (the weight of the
33 second pair of canonical variable (= r_2^2) is practically negligible in the calculation of the

1 distance) which is certainly the first (and obvious) factor of similarity between catchments
2 (but probably not sufficient to ensure homogeneity).

3 Here two dimensionless variables (*SR* and *IC*) mostly free from the scale effect have been
4 considered as the set of hydrological descriptors *X*. Even if can expect a slight influence of
5 the size of the catchment on the flatness of the FDC, *i.e.* on *IC*, results show that the
6 correlation between *IC* and the drainage area in the dataset is very weak and that the
7 introduction of *A* among the basin descriptors *Y* does not improve significantly the
8 correlations between the two first canonical variables. The highest coefficient of correlation
9 observed is just 0.34 between *SR* and the first quartile of the hypsometric curve. The context
10 for the definition of the first pair of canonical variable in our application is thus close to the
11 one met by Ouarda *et al.* (2001, 2008) and Haché *et al.* (2002) concerning the second pair of
12 the canonical variables. -No combination of catchment descriptors was found strongly
13 correlated with the two parameters *SR* and *IC*. As consequence the correspondence
14 between the hydrological space and the catchment-related physical space defined by CCA is
15 not guaranteed thereafter.

16 ~~Regarding CCA we decided to compare results with published works in terms of correlation~~
17 ~~structure. Fig. 9 indicates weak correlations between the canonical variables: $r_1 = 0.71$~~
18 ~~between W_1 and V_1 and $r_2 = 0.57$ between W_2 and V_2 . As comparison, for flood quantile~~
19 ~~estimation, Ouarda *et al.* (2001) obtained r_1 between 0.959 and 0.960 and r_2 between 0.279~~
20 ~~and 0.422 in an application in the Province of Ontario (Canada), Haché *et al.* (2002) obtained~~
21 ~~$r_1 = 0.986$ et $r_2 = 0.842$ in the Saint-Maurice river region (Canada) and Ouarda *et al.* (2008)~~
22 ~~obtained $r_1 = 0.966$ and $r_2 = 0.247$ in Mexico. These studies used at least one *T*-year flood~~
23 ~~quantile *QT* expressed in m^3/s as one of the hydrological variables and the drainage area *A*~~
24 ~~as one of the physiographical variables. Since catchment area is certainly the factor with the~~
25 ~~greatest influence on flood magnitude above climate, geology and land-use as one of the~~
26 ~~physiographical variables, CCA suggests automatically a first pair of canonical variables~~
27 ~~(V_1, W_1) highly correlated with *QT* and *A*, respectively. Roughly speaking, the presence of a~~
28 ~~strong link between one hydrological variable and one physiographical variable ensures at~~
29 ~~least one highly correlated pair of canonical variables. This is not the case here: the~~
30 ~~hydrological variables are two ratios free from scale effect and so *A* was excluded from the~~
31 ~~final variables involved in the definition of canonical variables and the highest coefficient of~~
32 ~~determination observed between *SR* and the first quartile of the hypsometric curve is just~~
33 ~~0.34. No statistical test (e.g. ANOVA, Laaha and Blöschl, 2006a) to check homogeneity in~~
34 ~~terms of FDC characteristics was performed on the clusters of gauged basins. Contrary to~~
35 ~~other applications (e.g. in Regional Flood Frequency Analysis for which a measure of~~
36 ~~regional heterogeneity is used to validate the derivation of a representative pooled growth~~

1 curve), we consider that statistical homogeneity (i.e. low variability around the mean values)
2 is not a necessary condition for ensuring accurate quantile estimates. Indeed an efficient
3 interpolation technique (e.g. an empirical formula) to predict the river flow characteristics of
4 interest could compensate the effect of heterogeneity present within the groups. Here
5 clustering is a way to remove the large scale variability due to dominant factors possibly
6 difficult to identify (e.g. hydrogeological properties) and interpolation procedures aim at
7 modelling thereafter the residual unexplained spatial variability at finer scale whatever the
8 homogeneity is. The next section presents the method considered to develop regional
9 regressions for each groupings approach. Their relative performances of each grouping
10 technique in terms of prediction of dimensionless FDC are compared.

11

12 **6 Regional regression**

13 **6.1 Method**

14 The homogeneous regions are now identified. ~~Multivariate~~ Multiple regression model
15 relations between the EOF model parameters and catchment descriptors can be developed.
16 Both linear and power form models dependences were investigated:

$$17 \alpha_1 = \lambda_0 + \sum_{j \in [1;18]} \lambda_j Y_j \quad (8)$$

$$18 \gamma = \lambda'_0 + \sum_{j \in [1;18]} \lambda'_j Y_j \quad (9)$$

$$19 \alpha_1 = \lambda_0 \prod_{j \in [1;18]} Y_j^{\lambda_j} \quad (10)$$

$$20 \gamma = \lambda'_0 \prod_{j \in [1;18]} Y_j^{\lambda'_j} \quad (11)$$

21 ~~Models-Parameters~~ $\lambda_j, j \in [0,18]$ and $\lambda'_j, j \in [0,18]$ were adjusted on observations to each
22 homogeneous group by the ordinary least squares method (using log transformed data to fit
23 power-form models).

24 In order to define the most appropriate model for each region, all combinations including one
25 to four variables among the 18 quantitative variables were tested and the 10 best regression
26 models in terms of adjusted coefficient of determination were retained. These models were
27 then refined/filtered through an interactive scheme: (i) outliers using Cook's distance were
28 removed first, (ii) the statistical properties of residuals (including normality and
29 homoscedasticity) were checked by visual inspection (only for the first two grouping
30 methods) and (iii) the robustness of each empirical formula was finally assessed by leave-

1 one-out cross-validation. The final models were selected regarding to the best value of the
2 coefficient of determination obtained ed by leave-one-out cross-validation.

3 **6.2 Results**

4 To measure the values of prior region delineation a global regression using the whole
5 available gauging stations dataset and the procedure described in Sect 4.56.1 was derived.

6 The descriptors involved are the elevation exceeded 25% of the catchment, the mean annual
7 catchment air temperature, the catchment yield and the fraction of the drainage catchment
8 with impermeable substratum. Note that two of them reflect the relevance of geological
9 properties to explain the variability of the parameters of the EOF model at large scale. The
10 analysis of the predictive models derived from VG and RT approaches demonstrates that:

11 - Linear and power form models are equally found.

12 - The performance of the regression as well as the set of relevant descriptors may vary
13 substantially from one region to another. R^2 ranges from 0 to 0.86, with the median equal
14 to 0.41. Most of the regressions involve four relevant basin descriptors.

15 - Regarding α_1 the four most important explanatory variables are the catchment yield CY ,
16 the drainage area A , the y -coordinate of the centre of gravity YG and the percentage of A
17 with impermeable substratum $\%Imp$. They are all involved in average in three empirical
18 formulas out of ten. Their presence is partly justified: YG may reflect the gradually
19 influence of the Mediterranean climate on flow variability from North to South; CY and
20 $\%Imp$ characterize more or less directly the effect of the geology (all things considered
21 the higher the fraction of impervious area, the sharper should be the FDC); lastly the
22 relevance of A can be justified if one assume that the flatness of the FDC probably
23 increases with the size of the basin due to larger storage capacities and due to
24 combinations of different river flow patterns originated from upstream tributaries.

25 The global predictive performance of each method in cross-validation (i.e. for all the sample)
26 was assessed using the root mean square error ($RMSE$) and the coefficient of determination
27 of the regression R^2 between observed and predicted values for the EOF model parameters,
28 $\gamma_{\ln(Q)}$ and α_1 . In addition to these statistics, scatter plots were drawn and inspected
29 visually to compare the spread of the predictions. These results are reported in the next four
30 figures (from Fig. 11~~Fig. 11~~Fig. 11~~Fig. 10~~ to Fig. 14~~Fig. 14~~Fig. 14~~Fig. 13~~). The two upper
31 panels plot estimated values against observed ones ($\gamma_{\ln(Q)}$ on the left and α_1 on the right).
32 Each point is related to one gauging station. A one-to-one line (in red) is added to each
33 graph. Absolute relative errors were also computed for each of the 15 selected

1 ~~standardizeddimensionless~~ percentiles Q_p and their empirical statistical distributions were
2 summarised by box plots displayed on the lower panel.

3 The cross validation results for the national regression are presented in ~~Fig. 11~~Fig. 11~~Fig.~~
4 ~~11~~Fig. 10. As expected the scores are unsatisfactory: dispersion is high around the one-to-
5 one line ($R^2 < 0.20$ for both EOF model parameters) and the low-flow percentiles were poorly
6 predicted. By comparison, the three next figures (~~Fig. 12~~Fig. 12~~Fig. 12~~Fig. 11 to ~~Fig. 14~~Fig.
7 ~~14~~Fig. 14~~Fig. 13~~) illustrate the performance of the three tested grouping methods and
8 suggest that:

- 9 - The regional regression based on the three grouping approaches is superior to global
10 regression like in Laaha and Blöschl (2006b) and in Vezza *et al.* (2010); results for all
11 models follow a similar pattern in terms of relative error on ~~standardizeddimensionless~~
12 percentiles: the highest errors are obtained for the lowest values.
- 13 - ~~RT is the best regionalisation method, and~~ VG performs nearly as well as RT with both
14 comparable R^2 and RMSE; however one should note that the estimations by VG
15 approach are probably heteroscedastic (the spread of errors increases along with α_1).
16 RT yields a little more accurate quantile estimates than VG when comparing the spread
17 of the relative absolute errors, i.e. the extent of the whiskers of the box plots in Figs. 11
18 and 12.
- 19 - CCA ~~outperforms~~ only slightly outperforms global regression. This finding is astonishing
20 since CCA is known as a very efficient regional estimation method.

21 To understand the unexpected performance for CCA, we performed additional computations
22 and compared the neighbourhoods defined by CCA to the expected ones, ideally defined in
23 the hydrological space. We first verified that regional regressions obtained with the expected
24 neighbourhoods were suited to estimates the EOF model parameters. Results showed very
25 satisfactory performances (R^2 reaches 0.63 and 0.69 for γ and α_1 respectively). This high
26 difference between performances is certainly due to the fact that the selected neighbours by
27 CCA were almost never the expected ones: for the 50 closest gauged basins, the
28 concordance between the neighbourhoods predicted by CCA and the theoretical ones are
29 weak. It confirms that the correlation between canonical variables is not strong enough to
30 guarantee the correspondence between the physiographical and hydrological spaces and
31 thus to ensure the efficiency of CCA. As mentioned before it probably points out the lack of
32 efficient catchment characteristics to strengthen the link between the two spaces – certainly
33 characteristics explicitly linked to hydrogeology since the application of the two other
34 methods differs only by the introduction of such a variable (*i.e.* dominant HER2).

1

2 7 Conclusion

3 In this study, a regionalisation method is suggested to estimate flow duration characteristics.
4 The developed approach supposes that the mean annual flow qa is known before estimating
5 FDCs at ungauged sites. Efforts have been therefore concentrated on the estimation of the
6 shape of the normalised FDC using a large data set of FDC derived from 1080 gauging
7 stations.

8 First, a parametric and parcimonious model based on EOF decomposition has been
9 developed to fit the observed shapes of the FDC. A comparison to other models referenced
10 in the literature demonstrates that the EOF model leads to the best estimates at gauging
11 stations. A reason could be that, conversely to the empirical approach, analytical formulas
12 are not flexible enough to accommodate the full range of observed shapes. Thus it would be
13 unrealistic to support the idea of one parametric model adapted to all the hydrological
14 conditions.

15 In a second step, different grouping techniques for identifying homogeneous regions and
16 developing separate regression models have been compared. Two of the grouping
17 procedures, VG and RT, with comparable performance, demonstrate the significant gain to
18 develop regional regressions. One should note that the RT classification procedure has the
19 advantage to be automatic and objective whereas heterogeneity may persist in the VG
20 groups that could explain its ranking (2nd). Nevertheless a large portion of the variance
21 remains unexplained. Further effort could be devoted to the interpolation of the residuals.
22 One could apply techniques such as adapted kriging (Sauquet, 2006), Top-Kriging (Skøien *et*
23 *al.*, 2006) or physiographical space based interpolation (Castiglioni *et al.*, 2009) for this
24 purpose.

25 Despite ~~thea~~ greatest flexibility in neighbourhood selection, *i.e.* a neighbourhood is defined
26 individually for each target site, the third and last grouping method, CCA, performed poorly.
27 These bad and unexpected scores for CCA may result from the difficulty to obtain a sufficient
28 correlation link between hydrological and physiographical spaces in the absence of relevant
29 characteristics to describe the hydrogeological properties within the catchments. Indeed, for
30 the other two grouping techniques hydrogeology is summarized by one qualitative variable,
31 *i.e.* the class of the dominant HER2, which provides sufficient information to increase
32 homogeneity within regions and to ensure more efficient regional regressions. As a result the
33 application of CCA in predefined regions with homogeneous hydrogeological properties
34 should be investigated to compare equitably CCA to other methods on the same bases.

1 **Acknowledgments**

2 The authors wish to thank the French National Agency for Water and Aquatic Environments
3 (ONEMA) for partially funding the present work.

4

5 **References**

6 B.R.G.M.: Carte géologique au 1/1.000.000^{ème}, 1996.

7 Braud, I. and Obled, C.: On the use of Empirical Orthogonal Function (EOF) analysis in the
8 simulation of random fields, *Stochastic Hydrol. Hydraul.*, 5, 125-134, 1991.

9 Burn, D. H.: An appraisal of the “region of influence” approach to flood frequency analysis,
10 *Hydrolog. Sci. J.*, 35, 149–165, 1990a.

11 Burn, D. H.: Evaluation of regional flood frequency analysis with a region of influence
12 approach, *Water Resour. Res.*, 26, 2257–2265, 1990b.

13 Castellarin, A., Camorani, G., and Brath, A.: Predicting annual and long-term flow-duration
14 curves in ungauged basins, *Adv. Water Resour.*, 30, 937–953, 2007.

15 Castellarin, A., Galeati, G., Brandimarte, L., Montanari, A., and Brath, A.: Regional flow-
16 duration curves : reliability for ungauged basins, *Adv. Water Resour.*, 27, 953-965, 2004.

17 Castiglioni, S., Castellarin, A. and Montanari, A.: Prediction of low-flow indices in ungauged
18 basins through physiographical space-based interpolation. *J. Hydrol.* 378(3–4), 272–280,
19 2009.

20 Chokmani, K. and Ouarda, T. B. M. J.: Physiographical space-based kriging for regional
21 flood frequency estimation at ungauged sites, *Water Resour. Res.*, 40, W12514,
22 doi:10.1029/2003WR002983, 2004.

23 Franchini, M. and Suppo, M.: Regional analysis of flow duration curves for a limestone
24 region, *Water Resource Management*, 10, 199–218, 1996.

25 Ganora, D., Claps, P., Laio, F., and Viglione, A.: An approach to estimate nonparametric flow
26 duration curves in ungauged basins, *Water Resour. Res.*, 45, W10418,
27 doi:10.1029/2008WR007472, 2009.

28 Haché M., Ouarda T. B. M. J., Bruneau, P., and Bobee, B.: Regional estimation by canonical
29 correlation analysis: Comparison of types of hydrological variables, *Can. J. Civ. Eng.*, 29(6),
30 899-910, 2002.

- 1 Hisdal, H. and Tveito, O. E.: Generation of runoff series at ungauged locations using
2 Empirical Orthogonal Functions in combination with kriging, *Stochastic Hydrol. Hydraul.*, 6,
3 255-269, 1991.
- 4 Holmes, M. G. R., Young, A. R., Gustard, A., and Grew, R.: A region of influence approach to
5 predicting flow duration curves within ungauged catchments, *Hydrol. Earth Syst. Sci.*, 6(4),
6 721-731, 2002.
- 7 Holmström, I.: On a method for parametric representation of the state of the atmosphere,
8 *Tellus*, 15, 127–149, 1963.
- 9 Hotelling, H.: Relations between two sets of variates. *Biometrika*, 28, 321–377, 1936.
- 10 Isik, S. and Singh, V. P.: Hydrologic regionalization of watersheds in Turkey. *J. Hydrol. Eng.*,
11 13(9), 824-834, 2008.
- 12 Krasovskaia, I., Gottschalk, L. and Kundzewicz, Z. W.: Dimensionality of Scandinavian river
13 flow regimes, *Hydrol. Sci. J.*, 44(5), 705-723, 1999.
- 14 Laaha, G. and Blöschl, G.: Seasonality indices for regionalizing low flows, *Hydrol. Process.*,
15 20, 3851–3878, 2006a.
- 16 Laaha, G. and Blöschl, G.: A comparison of low flow regionalisation methods—catchment
17 grouping. *J. Hydrol.* 323(1–4), 193–214, 2006b.
- 18 Li, M., Shao, Q., Zhang, L., and Chiew, F. H. S.: A new regionalization approach and its
19 application to predict flow duration curve in ungauged basins, *J. Hydrol.*, 389(1-2), 137-145,
20 2010.
- 21 Merz, R. and Blöschl, G.: Flood frequency regionalisation—spatial proximity vs catchment
22 attributes. *J. Hydrol.*, 302(1–4), 283–306, 2005.
- 23 Nathan, R. J. and McMahon T. A.: Identification of homogeneous regions for the purpose of
24 regionalization. *J. Hydrol.*, 121, 217-238, 1990.
- 25 Ouarda, T. B. M. J., Bâ, K. M., Diaz-Delgado, C., Cârsteanu, A., Chokmani, K., Gingras, H.,
26 Quentin, E., Trujillo, E., and Bobée, B.: Intercomparison of regional flood frequency
27 estimation methods at ungauged sites for a Mexican case study, *J. Hydrol.*, 348(1-2), 40-58,
28 2008.
- 29 Ouarda, T. B. M. J., Girard, C., Cavadias, G. S., and Bobée, B.: Regional flood frequency
30 estimation with canonical correlation analysis, *J. Hydrol.*, 254(1-4), 157-173, 2001.
- 31 Oudin, L., Hervieu, F., Michel, C., Perrin, C., Andreassian, V., Anctil, F., and Loumagne, C.:
32 Which potential evapotranspiration input for a lumped rainfall-runoff model? Part 2 Towards a

1 simple and efficient potential evapotranspiration model for rainfallrunoff modelling. *J. Hydrol.*,
2 303, 290-306, 2005.

3 Quimpo, R. G., Alejandrino, A. A., and McNally, T. A.: Regionalised flow duration curves for
4 Philippines, *J. Water Resour. Pl.*, 109(4), 320–330, 1983.

5 Quintana-Seguí, P., Le Moigne, P., Durand, Y., Martin, E., Habets, F., Baillon, M., Canellas,
6 C., Franchisteguy, L., and Morel, S.: Analysis of near surface atmospheric variables:
7 Validation of the SAFRAN analysis over France, *J. Appl. Meteor. Climatol.*, 47, 92-107, 2008.

8 Sauquet, E., Gottschalk, L., and Krasovskaia, I.: Estimating mean monthly runoff at
9 ungauged locations: an application to France, *Hydrology Research*, 39(5-6), 403 – 423,
10 2008.

11 Sauquet, E.: Mapping mean annual river flows: geostatistical developments for incorporating
12 river network dependencies, *J. Hydrol.*, 331(1-2), 300 – 314, 2006.

13 Shu, C. and Ouarda, T. B. M. J.: Flood frequency analysis at ungauged sites using artificial
14 neural networks in canonical correlation analysis physiographic space, *Water Resour. Res.*,
15 43, W07438, 938, doi:10.1029/2006WR005142, 2007.

16 Singh, R. D., Mishra, S. K., and Chowdhary, H.: Regional flow–duration models for large
17 number of ungauged Himalayan catchments for planning microhydro projects, *J. Hydrol.*
18 *Eng.*, 6(4), 310–316, 2001.

19 Skøien, J., Merz, R., and Blöschl, G.: Top-kriging—geostatistics on stream networks. *Hydrol*
20 *Earth Syst. Sci.*, 10, 277–287, 2006.

21 Smakhtin, V. U.: Low flow hydrology: a review. *J. Hydrol.*, 240, 147–186, 2001.

22 Snelder, T. H., Lamouroux, N., Leathwick, J. R., Pella, H., Sauquet, E., Shankar, U.:
23 Predictive mapping of the natural flow regimes of France, *J. Hydrol.*, 373, 57–67, 2009

24 Therneau, T. M. and Atkinson, B.: rpart: Recursive Partitioning, R package. [http://cran.r-](http://cran.r-project.org/web/packages/rpart/)
25 [project.org/web/packages/rpart/](http://cran.r-project.org/web/packages/rpart/), last access March 2011, 2010.

26 Turc, L.: Water balance of soil: relation between precipitation, evapotranspiration and runoff,
27 *Annals Agronom.*, 5, 491–595, 1954.

28 Vezza, P., Comoglio, C., Rosso, M., and Viglione, A.: Low Flows Regionalization in North-
29 Western Italy. *Water Resour. Manage.*, 24, 4049–4074, 2010.

30 Vidal, J. P., Martin, E., Baillon, M., Franchistéguy, L., and Soubeyroux, J. M.: A 50-year high-
31 resolution atmospheric reanalysis over France with the Safran system, *International Journal*
32 *of Climatology*, 30(11), 1627-1644. doi: 10.1002/joc.2003, 2010.

1 Vogel, R. M. and Fennessey, N. M.: Flow duration curves II: a review of applications in water
2 resources planning. *Water Resour. Bull.*, 31(6), 1029–39, 1995.

3 Wasson, J. G., Chandesris, A., Pella, H., and Blanc, L.: Typology and reference conditions
4 for surface water bodies in France: the hydro-ecoregion approach, *TemaNord*, 566, 37–41,
5 2002.

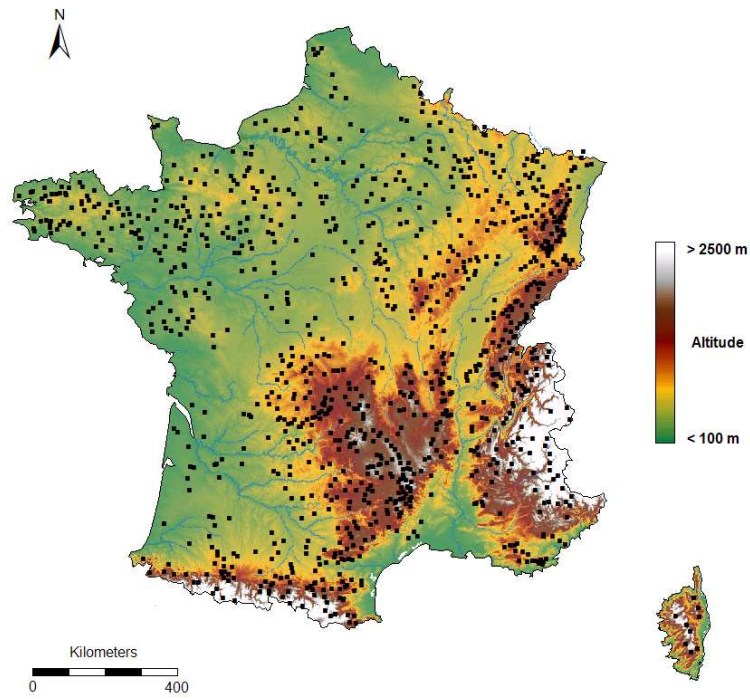
6 Yu, P. S., Yang, T. C., and Wang, Y. C.: Uncertainty analysis of regional flow duration
7 curves, *J. Water Resour. Plann. Manage*, 128(6), 424–30, 2002.

8

9

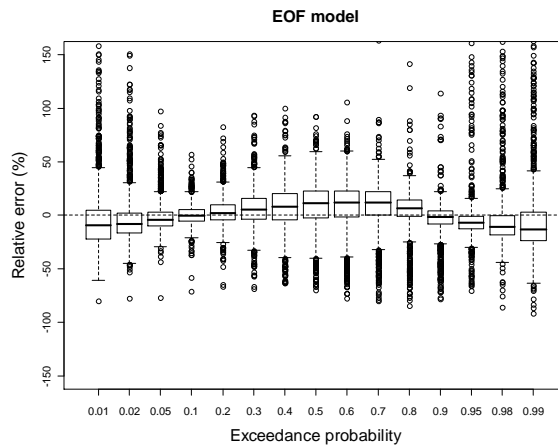
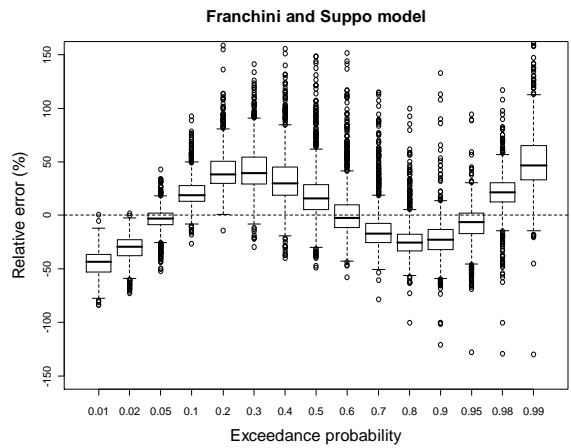
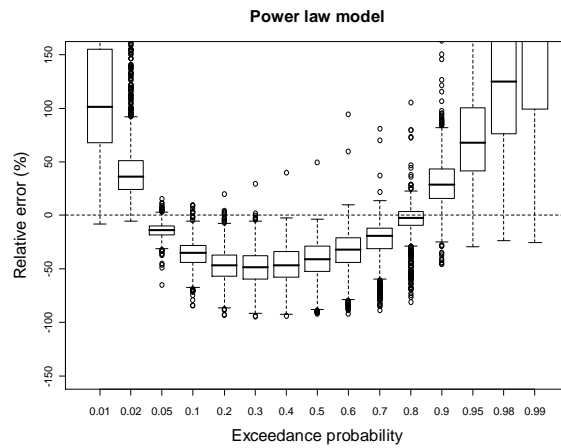
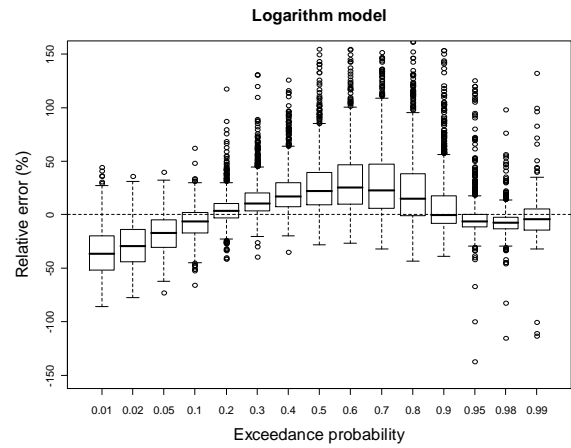
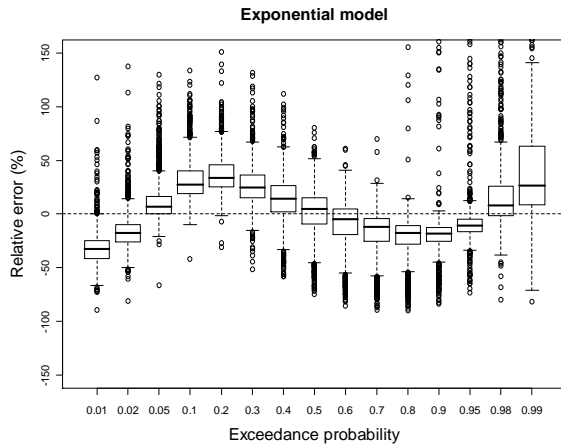
10

1 Figures



2
3 Fig. 1. Study area and gauging stations identified by their respective centre of gravity (black
4 square).

5
6
7
8
9
10
11
12



1

2

3

4

5 Fig. 2. Empirical distribution of the relative error for each percentile and each model. The
 6 boxplots are defined by the first quartile, the median and the third quartile. The whiskers
 7 extend to 1.5 the interquartile range; open circles indicate outliers.

8

9

10

11

12

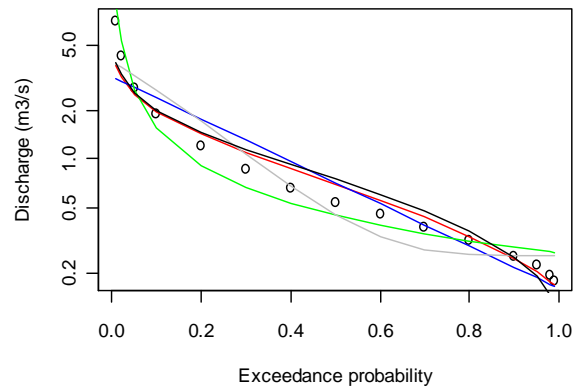
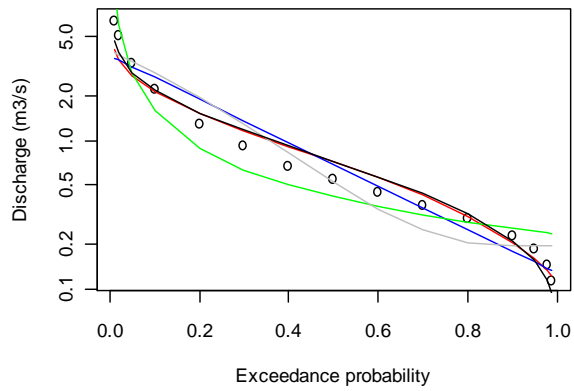
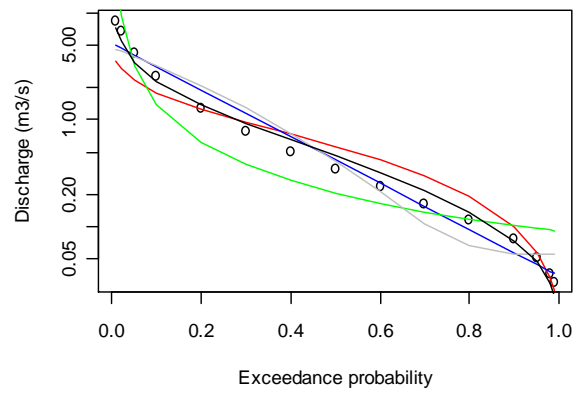
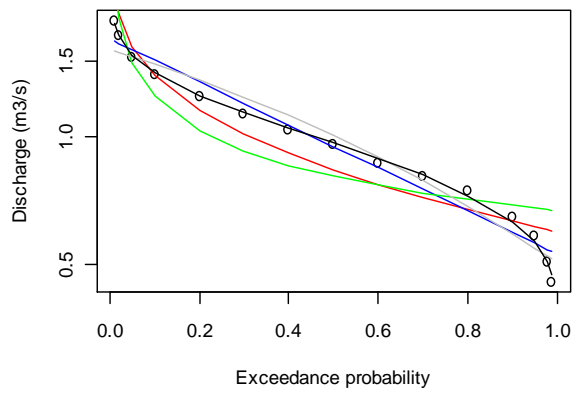
13

14

15

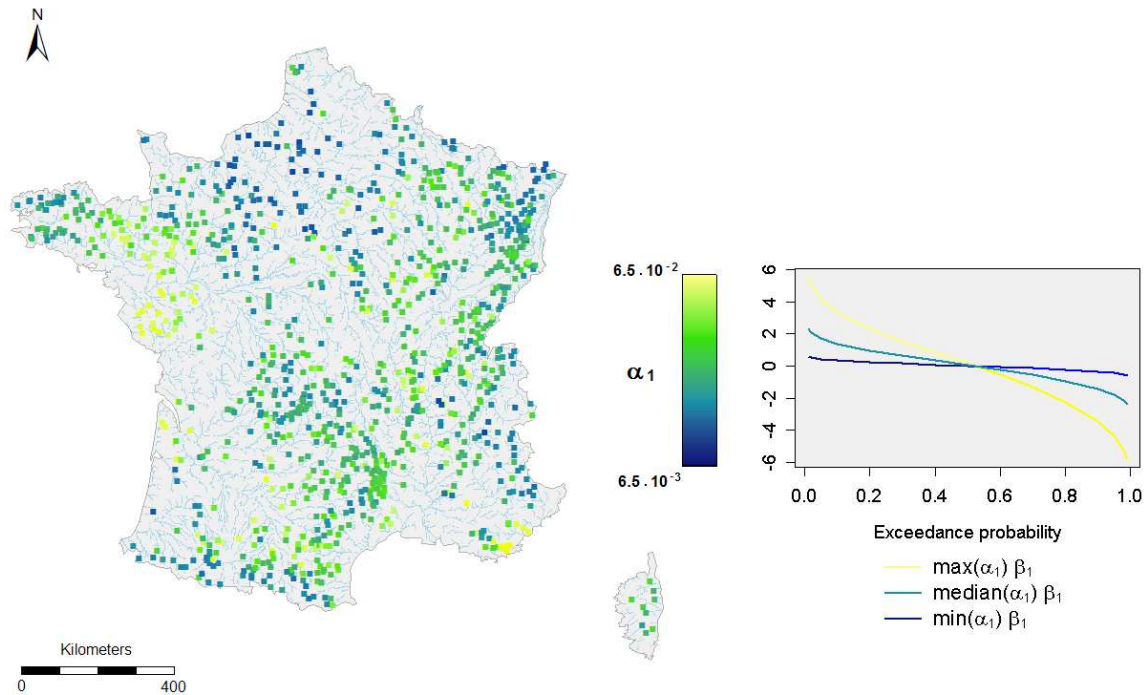
16

17



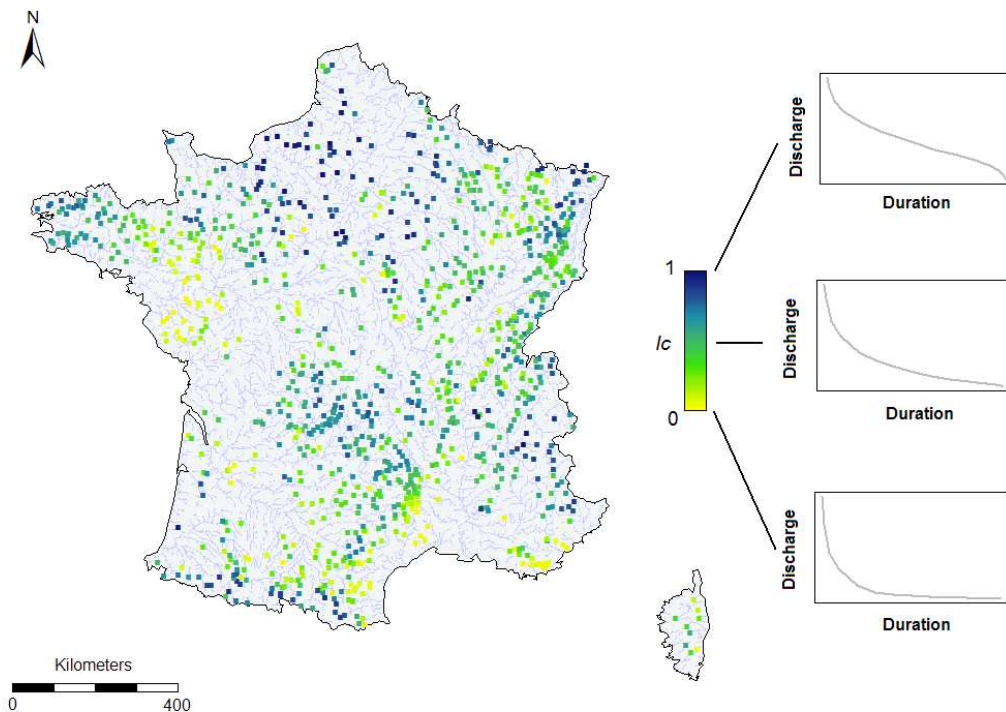
1
 2 Fig. 3. Comparison of observed (open circle) and modeled flow duration curves (logarithm
 3 (red), exponential (blue), power law (green), Franchini and Suppo (grey), EOF (black)).

4



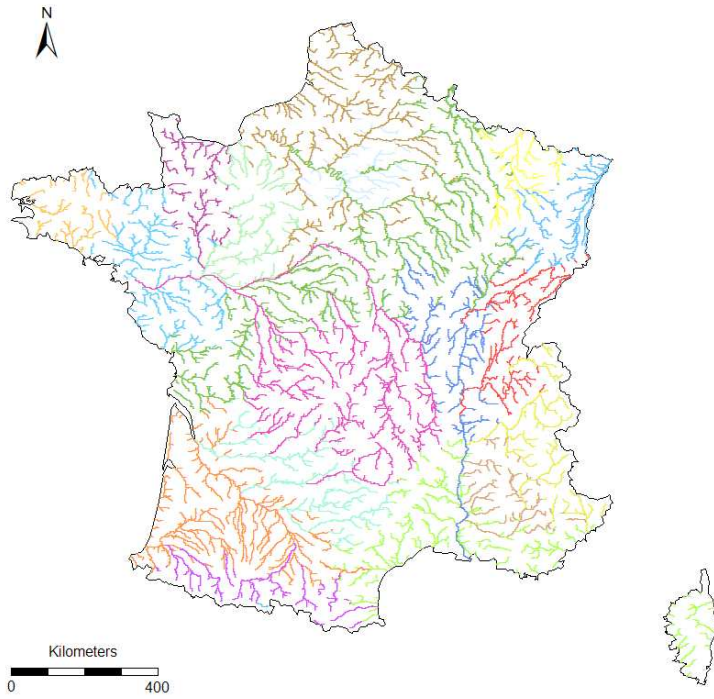
1
 2 Fig. 4. Spatial distribution of the weight α_1 observed at gauged catchments identified by the
 3 location of their centre of gravity.

4



5
 6 Fig. 5. Spatial distribution of the concavity index IC observed at gauged catchments identified
 7 by the location of their centre of gravity.

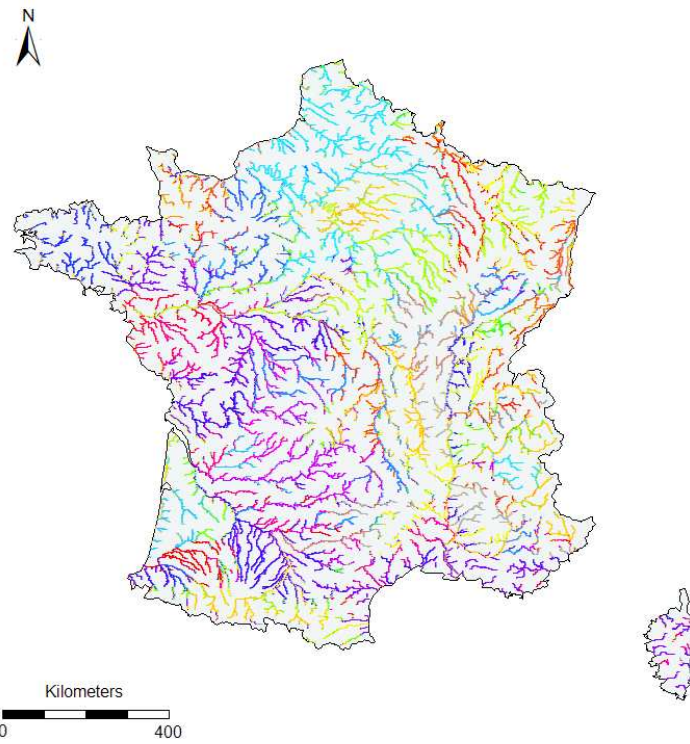
1



2

3 Fig. 6. Results of classification based on visual grouping (VG).

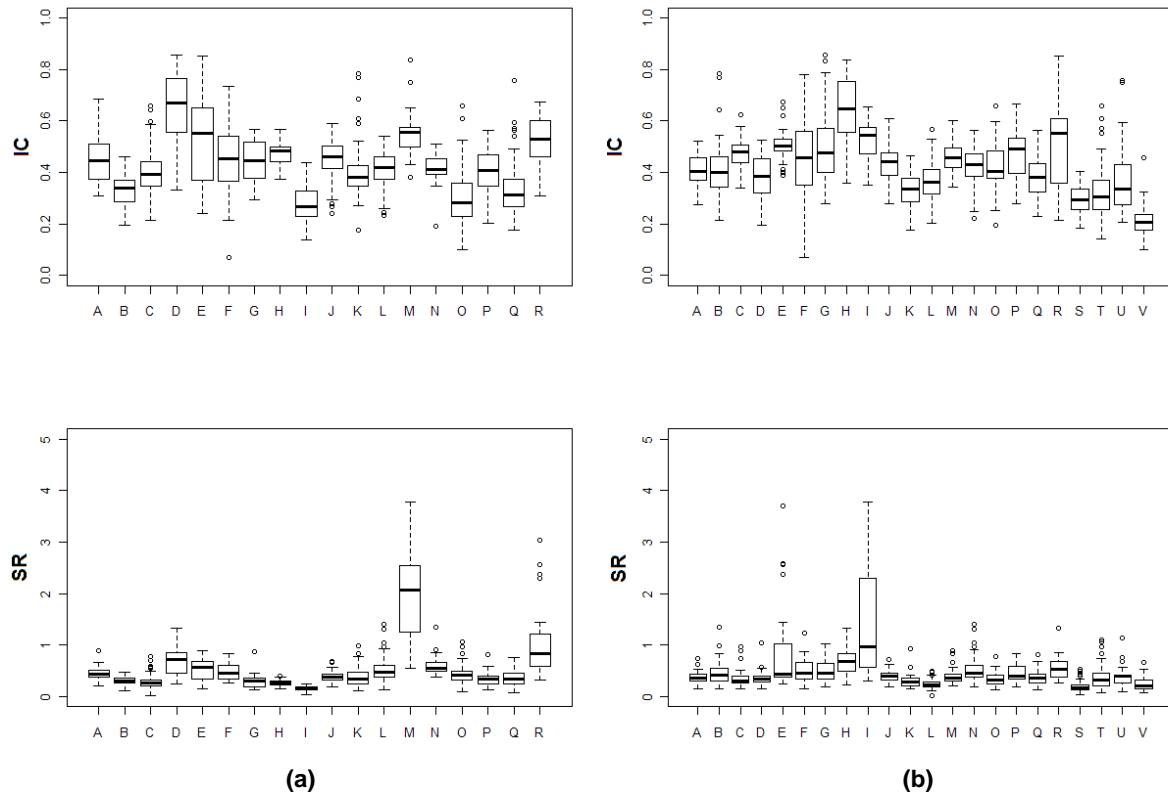
4



5

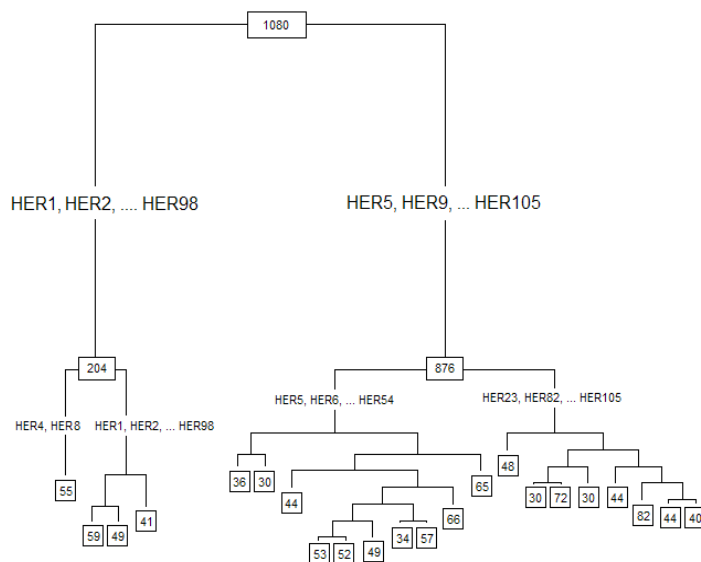
6 Fig. 7. Results of classification based on regression trees (RT).

7



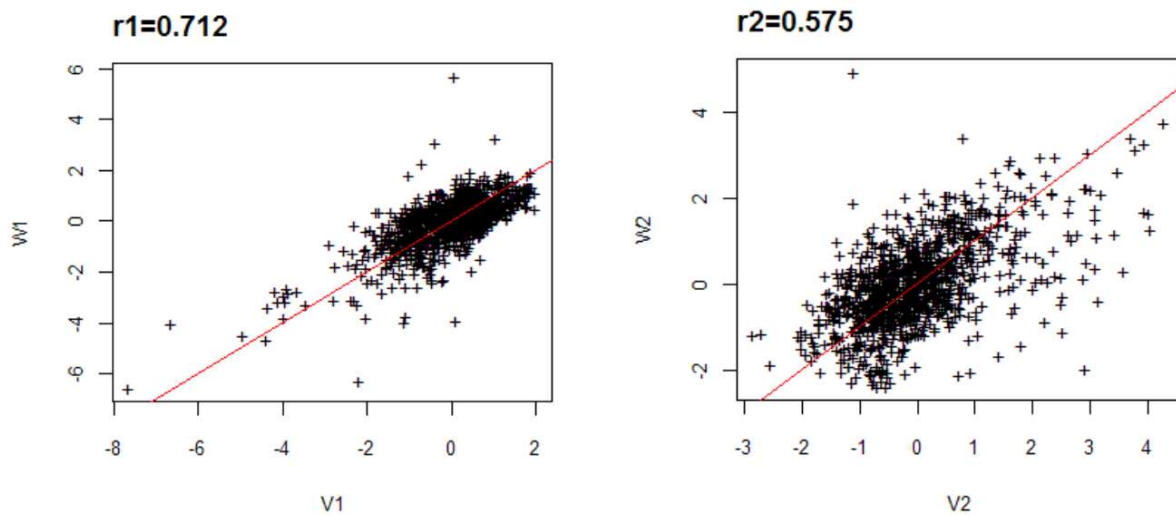
1
2
3
4

Fig. 8. Empirical distributions of the two hydrological indicators for each cluster according to VG (a) and RT (b).



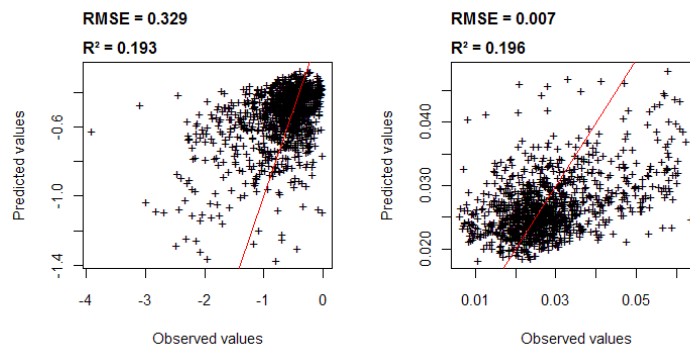
5
6
7
8

Fig. 9. Regression tree model (the numbers at each node of the tree and the name of the first splitting variables are reported in the boxes and in the middle of the branches, respectively).

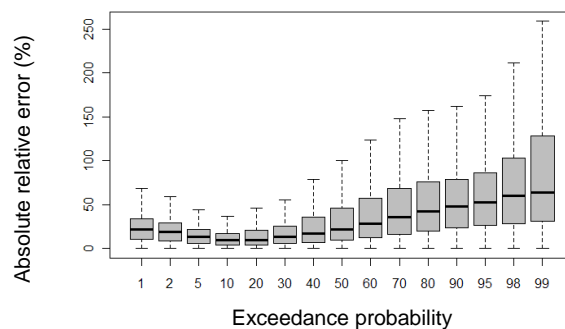


1
2

3 Fig. 1040109. Correspondence between the position of the gauged sites in the hydrological
4 space and the catchment descriptors space - Correlation between canonical variables. V_1 and
5 V_2 (resp. W_1 and W_2) are the two first canonical variables of hydrological space (resp. of
6 basin descriptors space)



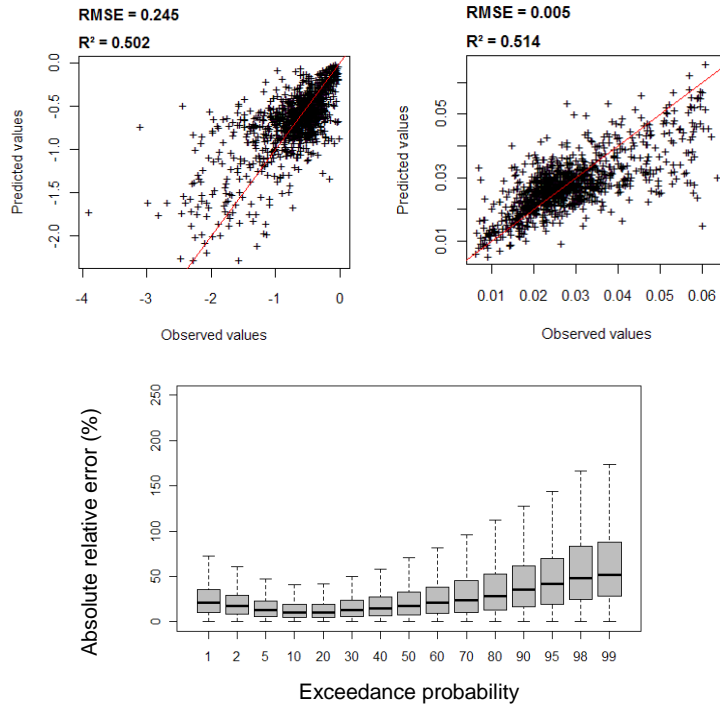
7



8

9 Fig. 11. Results for the global regression model.

10

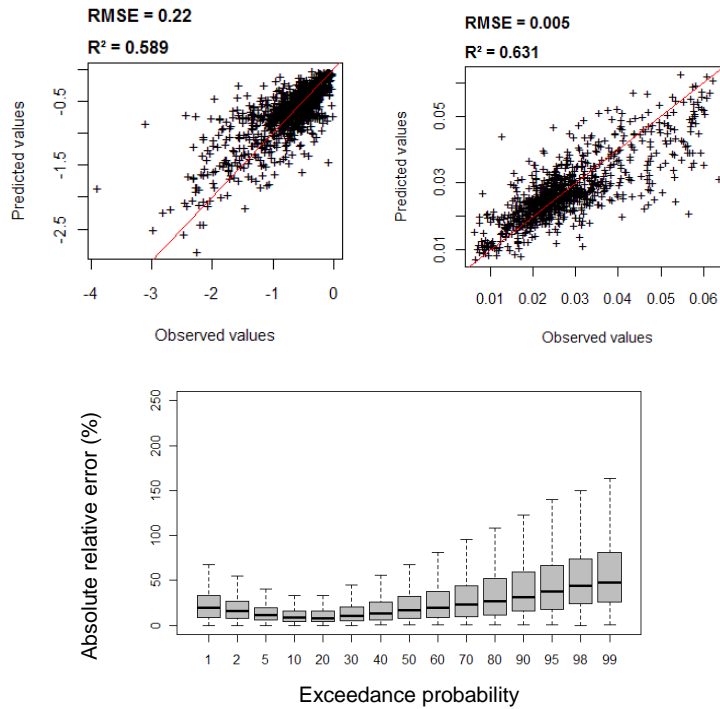


1

2

3 Fig. 12. Results for the regional regression model applied to visual grouping.

4



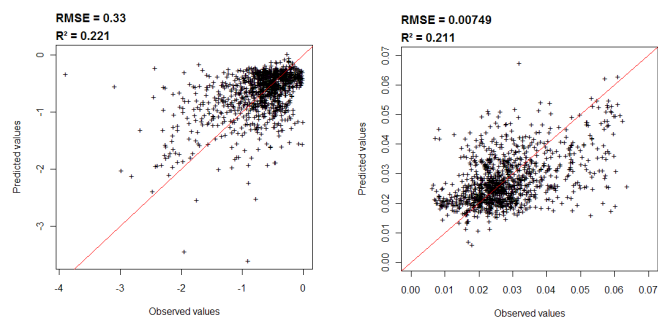
5

6

7 Fig. 13. Results for the regional regression model applied to groups derived from RT.

8

1



2

3 Fig. 14. Results for the regional regression model applied to neighbourhoods derived from
4 CCA.

5

