

¹ **An approach to estimate non-parametric flow**
² **duration curves in ungauged basins**

D. Ganora,¹ P. Claps,¹ F. Laio,¹ and A. Viglione²

Daniele Ganora (daniele.ganora@polito.it)

Pierluigi Claps (pierluigi.claps@polito.it)

Francesco Laio (francesco.laio@polito.it)

Alberto Viglione (viglione@hydro.tuwien.ac.at)

¹Dipartimento di Idraulica, Trasporti ed
Infrastrutture Civili, Politecnico di Torino,
Torino, ITALY.

²Institut für Wasserbau und
Ingenieurhydrologie, Technische Universität,
Wien, AUSTRIA.

Abstract. A distance-based regional model is developed for the estimation of dimensionless flow duration curves in sites with no or limited available data. The curves are dimensionless because they are preliminarily normalised by an index value (e.g., the mean annual runoff). The model aims to represent this curve as a non-parametric object, rather than providing a parametric representation and trying to relate the parameter values to basin descriptors. This approach requires considering the (dis)similarity between all possible pairs of curves, and searching for characteristic distance values that can be related to basin descriptors, taken among geographic, geomorphologic and climatic parameters. The (dis)similarity between curves is computed using a predefined metric based on a linear norm and producing a distance matrix. This matrix is then related to analogous matrices of descriptors differences by means of linear regression models. The significant descriptors are identified and a cluster analysis is applied so that the sites can be grouped together. Each region is supposed to be characterized by a single dimensionless flow duration curve. The procedure is applied to 95 basins located in northwestern Italy and Switzerland. The performance in the regional estimation is assessed by means of a cross-validation procedure through comparison with “standard” parametric regional approaches based on two and three-parameter models. In most of the cases, the distance-based model produces better estimates of the flow duration curves, using only few catchment descriptors.

1. Introduction

The problem of estimating hydrological variables in ungauged basins has been the object of intense research activity in recent years [see e.g., *Sivapalan et al.*, 2003]. Regardless of the method used to perform such estimation, the underlying idea is to transfer the hydrological information from gauged to ungauged sites. When observations of the same variable at different measuring sites are available and are used for the estimation in ungauged sites, the related methods are called *regional methods*. *Regional frequency analysis* [e.g., *Hosking and Wallis*, 1997], where the interest is in the assessment of the frequency of occurrence of hydrological events, belongs to this class of methods. A frequently used regional approach is the *index-value method* [*Dalrymple*, 1960] in which it is implicitly assumed that the frequency distribution for different sites belonging to a homogeneous region is the same except for a site-specific scale factor, the index-value [see e.g., *Hosking and Wallis*, 1997, for details]. Hence, the estimation of the distribution for an ungauged site is obtained by separately estimating the index value and assigning the site to an homogeneous region, which entails assuming a dimensionless frequency distribution (a growth curve) to the site under analysis.

The frequency distribution is only one of the possible information that can be transferred using a regional approach. In this paper we deal with a specific descriptor of the runoff distribution in a basin: the *Flow Duration Curve* (FDC). A flow duration curve represents the flow in a stream rearranged to show the percentage of time during which a discharge value is equalled or exceeded. Strictly speaking this is not a probability curve, because discharge is correlated between successive time intervals and discharge characteristics

are dependent on the season; hence the probability that discharge on a particular day exceeds a specified value depends on the discharge on preceding days and on the time of the year [Mosley and McKerchar, 1993, p. 8.27]. However, a flow duration curve is often interpreted as the complement of the cumulative distribution function of the daily streamflow values at a site. The FDC also provides a graphical summary of streamflow variability and is often used in hydrologic studies for hydropower, water supply, irrigation planning and design, and water quality management (a review on many applications is provided by Smakhtin [2001]).

The empirical FDC is constructed from observed streamflow time series. These observations can have different time-scale resolution, although mean daily streamflow values are commonly used. The data are ranked in descending order and each ordered value is associated with an exceedance probability F , for example through a plotting position formula. If the FDC is constructed on the basis of the whole available data set, merging together all available years of data, it represents the variability of flow over the entire observation period. This representation is valid when the dataset is sufficiently long. A different approach, introduced by Vogel and Fennessey [1994], is to consider annual FDCs separately, i.e., to consider a different FDC for each year when data are available [e.g., Claps and Fiorentino, 1997; Iacobellis, 2008]. A parametric model able to represent both the total and the annual FDCs for gauged and ungauged sites has been proposed, for instance, by Castellarin et al. [2004b, 2007].

In the present work only total FDCs will be considered, adopting a non-parametric approach for their representation. The FDCs are modelled following the index-value approach, in which the flow duration curve $Q(F)$ is the product of two terms $Q(F) =$

69 $\mu \cdot q(F)$, where the *index flow* μ is the scale factor and the *dimensionless total flow*
 70 *duration curve* $q(F)$ represents the shape of the FDC. The present work focuses on the
 71 regionalization of the dimensionless curve, while the estimation of the index flow will not
 72 be treated. In section 2 we describe the distance-based method. This method is applied
 73 to a case study in section 3, where a set of basins located in North-Western Italy and
 74 Switzerland is investigated. The method's performances against alternative parametric
 75 methods are finally checked in section 4.

2. Distance-based Method

76 Leaving aside the index-flow estimation, the regional FDC analysis can be divided into
 77 two parts: the formation of cluster regions and the association of an ungauged site to one
 78 of them. Concerning the first point, the curves are grouped according to their similarity
 79 in terms of the basin descriptors that better “explain” the shape of the FDC. In standard
 80 approaches [e.g., *Fennessey and Vogel*, 1990; *Singh et al.*, 2001; *Holmes et al.*, 2002],
 81 this shape is represented in a parametric way. For instance, the coefficient of variation
 82 (CV) or the L -CV [*Hosking and Wallis*, 1997] of the curve can be used for this purpose.
 83 In this case, the selected parameter is related to basin descriptors through a linear or
 84 a more complex model. A regression analysis is performed with different combinations
 85 of descriptors, and those that are strongly related with the parameter are used for its
 86 estimation in ungauged sites.

87 In the distance-based approach proposed here we consider the dimensionless FDC as a
 88 whole, without resorting to statistical descriptors of its shape. This means that a curve
 89 is not fitted by an analytical function, which would imply a parametric representation
 90 of the FDC. The multiregression approach can still be used to study the (dis)similarity

between pairs of basins. The procedure is synthetically described below as a sequence of logical steps, while details are provided in the following subsections:

1. for each couple of stations, a dissimilarity index between dimensionless curves is calculated using a predefined metric (section 2.1);

2. for each considered basin descriptor (e.g., area, mean elevation, mean slope, drainage path length, etc), the absolute value of the difference between its measure in two basins is used as the descriptor distance;

3. the distances between couples of FDCs (and between basin descriptors) are organized in distance matrices (section 2.2);

4. a multiregression approach is applied using the FDC distance matrix as the dependent variable, and the descriptor distance matrices as the independent variables; this serve to select the relevant basin descriptors (those associated to the best regression model) (section 2.2);

5. in the resulting descriptors' space, stations with similar descriptor values (small distances between descriptors) are grouped together into regions through a cluster analysis (section 2.3);

6. the regional dimensionless flow duration curve is estimated by taking the average of all the curves belonging to the cluster, as in the “graphical approaches” reviewed by *Castellarin et al.* [2004a] and references therein.

Critical points of this procedure, discussed more in detail in the following, are the choice of a suitable distance measure for the dimensionless flow duration curves, the identification of the best regression model between distance matrices, and the choice of the method of cluster analysis for the formation of the regions.

2.1. (Dis)similarity Between Curves

Let Q_s^* be the sequence of N_s daily discharges in the gauged station s , containing all the recorded values. Based on these data the scale factor μ_s is first computed as the average of the whole sequence. Then, the dimensionless sequence $q_s^* = Q_s^*/\mu_s$ is rearranged in descending order and each value $q_{i,s}$, with $i = 1, 2, \dots, N_s$, is associated to its exceedance probability (i.e., through the Weibull plotting position)

$$\left\{ \frac{1}{N_s + 1}, \frac{2}{N_s + 1}, \dots, \frac{N_s}{N_s + 1} \right\}. \quad (1)$$

The distance-based procedure proposed here is based on the comparison between couples of curves: for this purpose it is convenient the two curves have the same number of elements. Since total FDCs have generally different lengths, depending on the number of years they cover, we resample them to make the curves comparable. For this purpose, we resample the FDCs at the frequency values

$$\left\{ \frac{1}{365 + 1}, \frac{2}{365 + 1}, \dots, \frac{365}{365 + 1} \right\}, \quad (2)$$

obtaining a new representation of the FDC in the station s , $\{q_{1,s}, q_{2,s}, \dots, q_{365,s}\}$. Other sampling rates can be used to better sample particular parts of the curves. In this work we have also considered an alternative sampling method that produces 365 equally spaced values in the z -space, where z is the normal reduced variate (with zero mean and unit variance). Back-transforming these values to the frequency space, the 365 values are no more equally spaced but more concentrated around higher and lower frequencies. Figure

1 sketches two curves with different number of elements resampled with a constant and a
 2 z spacing in the frequency axis.

3 In our approach a measure of similarity between curves (hereafter termed *distance*) is
 4 required. Given two FDCs, relative to two gauging stations s_1 and s_2 , constituted by
 5 365 elements each: $\{q_{1,s_1}, q_{2,s_1}, \dots, q_{365,s_1}\}$ and $\{q_{1,s_2}, q_{2,s_2}, \dots, q_{365,s_2}\}$, a simple measure
 6 of their dissimilarity can be defined as the “distance” calculated by the norm of order one,

$$\delta_{s_1,s_2} = \sum_{i=1}^{365} |q_{i,s_1} - q_{i,s_2}|. \quad (3)$$

7 The value δ_{s_1,s_2} can be interpreted also as an approximation of the area between the
 8 curves. The computation of the distance according to equation (3) is exemplified in figure
 9 2 for two generic FDCs.

10 If n is the number of sites where data are available, the distance measures for each FDC
 11 pair are organized in a $n \times n$ distance matrix like:

$$\Delta = \begin{pmatrix} 0 & \delta_{1,2} & \dots & \delta_{1,n} \\ \delta_{2,1} & 0 & & \vdots \\ \vdots & & \ddots & \\ \delta_{n,1} & \dots & & 0 \end{pmatrix} \quad (4)$$

12 where the elements δ_{s_1,s_2} are distances between curves (calculated with equation (3)).
 13 Analogously, matrices like (4) can contain distances between catchment descriptors (if d_1
 14 is the value of the descriptor for basin 1 and d_2 for basin 2, then $\delta_{1,2} = |d_1 - d_2|$). Since
 15 the matrices are symmetric and with null diagonal values, after removing the redundant
 16 values, only $n(n-1)/2$ values per matrix are informative.

The distance measure of equation (3) not only depends on the resampling method but also on the “measurement space” considered for the representation of flows. For example, if the flows are transformed to provide a more convenient representation of the FDC, the distances δ_{s_1, s_2} are affected by the transformation. Three main representations of the FDC are considered in this work: (a) flow data plotted versus their corresponding plotting position, (b) log-transformed flows versus their corresponding plotting position and (c) log-normal probability plot (log-transformed flows versus normal reduced variate). There are no particular reasons to prefer a priori one of these representations, therefore all of them are considered in the case study and will be respectively referred as “linear representation”, “logarithmic representation” and “log-normal representation” (see figure 2). Three parametric functions will be used in a traditional regional FDC estimation exercise in section 4, for comparison to the distance-based procedure developed here.

2.2. Distance Matrices, Linear Regression and Mantel Test

In this section we show how to identify the catchment descriptors that, thanks to their relations with the FDCs, should be used for the formation of cluster regions. A different distance matrix, hereafter termed Δ_{X_i} , is determined for each descriptor, while the distance matrix for the dimensionless FDCs is called Δ_Y . The relation between the distance matrix Δ_Y and the various Δ_{X_i} is assessed using a multiregressive approach. Note that the multi-regressive approach based on distance matrices is not used to estimate FDC coefficients, but to identify the descriptors to be used in the following step for region creation. We start considering a simple linear model:

$$\Delta_Y = \beta_0 + \beta_1 \Delta_{X_1} + \dots + \beta_p \Delta_{X_p} + \epsilon \quad (5)$$

with p as the number of descriptors involved, β_i as the regression coefficients and ε the residual matrix. The best possible regression is selected through the *adjusted coefficient of determination*

$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1} \quad (6)$$

where R^2 is the standard coefficient of determination [e.g., *Kottegoda and Rosso, 1997*], p the number of descriptors and n the number of basins considered. The regression coefficients and R^2 can be computed in a standard way [*Legendre et al., 1994*], that is to say that it does not matter if the elements are organized in a distance matrix. However, in the formulation of the adjusted coefficient of determination it is better to use the value n (the number of basins) instead of $n(n - 1)/2$ that is the number of points involved in the regression (namely the number of distance values). This is due to the fact that the values inside the matrices are not mutually independent. Dependancy has another significant impact on the method. In particular, the validity of the tests used to assess the significance of the independent variables (e.g., the Student t test) is affected. A different significance test, as the *Mantel* test [*Mantel and Valand, 1970*], is then needed, which accounts for the non-independence of the elements in the distance matrices.

The Mantel test was originally proposed by *Mantel and Valand* [1970] for analysis of correlation between distance matrices, and since then it has been widely improved and used with many different kinds of data. In fact, distance matrices have been frequently used in the biological and ecological sciences [e.g., *Legendre, 1993; Lichstein, 2007*]. The

simple Mantel test [Mantel and Valand, 1970] is used to evaluate the significance of the linear correlation between two distance matrices. This test is performed computing a statistic (usually the Pearson correlation coefficient) between all the pairwise elements of the two matrices. Its significance is tested by repeatedly permuting the objects in one of the matrices, and recomputing the correlation coefficient each time; permutations are performed simultaneously exchanging the rows and the columns of the matrices (e.g., if rows of indexes 2 and 10 are exchanged, also columns of indexes 2 and 10 have to be exchanged [see Legendre et al., 1994]). The significance of the statistic is assessed by comparing its original value to the distribution of values obtained from the permutations, which are considered as many realizations of the null hypothesis of no correlation.

The simple Mantel test can be extended to multiple predictor variables to be applied in multiple linear regression models as (5). The extension has been introduced by Smouse et al. [1986], discussed and improved by Legendre et al. [1994] and recently applied in the ecological field by Lichstein [2007]. Following the procedure of Lichstein [2007] each matrix, after removing redundant values, is unfolded into a vector of distances, and regression is performed in the classical way. Then, a null distribution is constructed permuting the elements only in the dependent variable distance matrix Δ_Y . Similarly to what described for the simple Mantel test, the rows and the columns of the matrix Δ_Y are permuted simultaneously and each regression coefficient is tested individually.

2.3. Cluster Analysis

The proposed procedure serves for the estimation of a FDC in an ungauged basin on the basis of curves relative to other basins. Given a large group of candidate “donor” basins, we want to extract a subset of basins that have geomorphologic and climatic

characteristics similar to those of the target site. The FDCs collected in these sites will be used for the estimation of the unknown curve. There are different regionalization techniques to choose the subset of basins, for example leading to the formation of fixed regions through cluster analysis [*Hosking and Wallis*, 1997; *Viglione et al.*, 2007b], or based on the method of the region of influence [ROI, *Burn*, 1990]. In this work we use the first approach, selecting fixed regions by splitting the descriptors space in non-overlapping areas by means of a cluster analysis. However, the generalization of the method to the ROI technique is straightforward. The definition of the descriptors space depends on the outcome of the multiregressive procedure described in section 2.2, that allows one to identify a group of significant geomorphoclimatic parameters.

The cluster analysis method used here is a mixed method in which the Ward hierarchical algorithm [*Ward*, 1963] is followed by a reallocation procedure that minimizes the dispersion within each cluster. The Ward algorithm is agglomerative; it starts with a configuration in which each element is a cluster itself, and progressively merges clusters in a way to produce the minimum information loss, measured as the sum of squared deviation of each element from its cluster centroid. We use the Ward algorithm because it is able to generate compact clusters with an evenly distributed number of elements. A disadvantage is that it does not allow elements reallocation, so that the final configuration could not be the optimal one. To avoid this inconvenience, a reallocation procedure is applied in concurrence with the agglomerative clustering. For instance, if the Ward clustering yields a final configuration with k clusters we compute the statistic

$$W = \sum_{i=1}^k \left(\sum_{j=1}^{n_i} D_{i,j}^2 \right), \quad (7)$$

where $D_{i,j}$ is the Euclidean distance between the j -th element of the i -th cluster and the cluster centroid, and n_i is the number of elements contained in the i -th cluster. An element is moved to another cluster if the new configuration provides a lower value of W . The procedure ends when W stops reducing after the reallocations, so that every element of a cluster is closer to its center of mass than to the centroid of the nearby cluster.

The reallocation procedure leads to an optimal configuration with k regions. A controversial point of the procedure is the choice of the optimal number of clusters. Usually, in regional analyses, the aim is to get the smallest possible number of homogeneous regions, so that each of them has a large enough number of elements. In this work, the selection of the ideal number of clusters is done investigating different k values and evaluating, for each configuration, a quality index. This index is computed by estimating (in cross-validation mode) the curves for all sites by using the regional model (with a given k), and computing the distance as in equation (3) where, in this case, s_1 is the measured curve and s_2 is the estimated one. This distance is adopted as an error measure and the overall mean error is used as a quality index to select the number of clusters. This method does not ensure that the clusters are homogeneous, because no homogeneity test is explicitly used.

After having subdivided the descriptors space in regions, one can proceed to the estimation of the flow duration curve in ungauged sites. For one such site one must first determine the values of the descriptors selected in the procedure of section 2.2. The descriptors at the ungauged site are entered as coordinates in the descriptors' space and the site is assigned to the cluster whose centroid is the closest to the basin descriptors. The

curves of all basins belonging to the selected cluster will be used to build the regional curve. This latter curve is simply estimated point by point as the average of the values of q relative to each duration for the curves belonging to the selected region, as in the graphical approach described in *Castellarin et al.* [2004a].

The descriptors used in the cluster analysis are preliminary standardized (i.e., converted into variables with zero mean and unit variance). Standardization of raw descriptors values avoids an unwanted weighting effect due to the different measurement units. If the descriptors are assumed to have different importance in the cluster creation, a procedure can be adopted to give different weights to each descriptor. Regression coefficients of equation (5) can be used to compare the relative effect of each descriptor distance matrix, if the distance matrices have been previously standardized: the greater the coefficient, the greater the relative effect of its descriptor distance matrix on the curve distance matrix, so that the coefficients can be used as weights. This weighted clustering procedure will be tested in the following sections by comparing it to the standard unweighted clustering.

After the regional curves have been determined, it is necessary to evaluate if they can be considered significantly different from each other, because otherwise the regions should be merged. To assess if two regional curves are significantly different, we use a procedure based on the distances between curves. First, a *reference distance* is computed as the median (or the mean) of the distances between each empirical curve and the regional one. Then, the distance matrix of the regional curves is computed and all its elements are compared against the reference distance: two regional curves are considered significantly different if their distance is greater than the reference distance, otherwise the two

clusters are merged together. This procedure is repeated until all the regional FDCs are significantly different.

Note that the reference distance and the distance matrix of the regional curves depend on the representation space on which the distances are calculated, hence different results are expected using different representation spaces.

3. Case Study: Distance-based Method Application

3.1. Hydrological and Geomorphologic Data

The application of the distance-based procedure for regional estimation of FDC has been carried out in the R statistical environment [*R Development Core Team*, 2007], integrated for Mantel test and cluster analysis with the *nsRFA* package [*Viglione*, 2007].

Available data include 95 river basins located in northwestern Italy (36 basins of Piemonte and Valle d'Aosta regions) and in Switzerland (59 basins); the geographical location of the gauging stations is shown in figure 3. Italian flow data derive from the publications of the former Italian Hydrographic Service and include series lengths ranging between 7 years and 41 years. Hydrological and geomorphological variables relative to Italian basins are included in the widest CUBIST database [*CUBIST Team*, 2007] that contain such data for more than 500 basins in Italy. The catchment area of Northwestern Italy basins ranges between 22 and 7983 km², and their average elevation ranges from 494 to 2694 m a.s.l. Switzerland data are included in the Reference Hydrometric Network (SHRN) provided by the BAFU (Bundesamtes für UmweltSwiss) and include daily streamflow series with a minimum length of 18 years and a maximum length of 99 years. The catchment area of Switzerland basins ranges between 7 and 616 km², while their average elevation varies from 475 to 2847 m a.s.l. Geomorphological characteris-

tics of each basin has been obtained from a digital terrain model (about 90m cell grid) provided by NASA [2000] with automatic procedures originally developed by *Rigon and Zanotti* [2002] under a GRASS GIS environment. For the complete list of basins considered, whose codes are referred in figure 3, and their geomorphologic variables see auxiliary material at <http://www.idrologia.polito.it/~ganora>.

3.2. Procedure Setting

Several linear regression models between distance matrices have been investigated using relation (5). They are built using different combination of:

1. Curve distance matrices Δ_Y : the three representations described in section 2.1 and figure 2 (linear, logarithmic and log-normal plot) are considered;
2. Descriptors distance matrices Δ_X : all possible combination from one to five descriptors have been taken into account.

Regression models are ordered in terms of R_{adj}^2 values and tested for significance with the multiple Mantel test, with a significance level of 0.05. Furthermore, a test against multicollinearity has been performed in order to exclude variables with redundant information [Montgomery et al., 2001].

For the linear representation, best results are obtained with four and three descriptors. Lower R_{adj}^2 values arise from simpler models with only two descriptors. In the logarithmic space, the best model is again characterized by four descriptors, but in this case simpler models with two parameters have comparable R_{adj}^2 . In the log-normal space none of the solutions accepted after testing are based on more than two descriptors. We decided to adopt models with two parameters because of their higher robustness (see table 1). The R_{adj}^2 values obtained with regression models with distance matrices are very low, although

the descriptors result to be statistically significant. In this regard it is important to remind that regressions are only used for the selection of the suitable descriptors and not for direct estimation.

Table 1 shows the three best models for each representation with two descriptors, where all the models have been tested for significance of regression coefficients with the Mantel test with a level of significance of 0.05. It appears that, considering together the three representations of different curve distance matrices, the most significant descriptors are always the same: the minimum basin elevation (Hmin), the mean elevation (H), the mean hillslope length (MHL), the mean basin slope (Slo) and the modified basin slope (Pm, which is the ratio between the median elevation and the square root of the area). A summary of the range of these descriptors is reported in table 2. This suggests to adopt the same set of descriptors with all the three representation spaces; Hmin and MHL has been selected. The adoption of these two descriptors is coherent with the typology of investigated basins. In fact, since we are considering mainly mountain basins, the elevation descriptor is expected to be relevant because of its strong relation to snow-accumulation and snowmelt mechanisms. Similarly, the hillslope mean length provides a synthetic description of runoff routing mechanisms.

3.3. Regions Definition

The second step, after the choice of the suitable descriptors, is to pool the catchments together with the cluster analysis, as described in section 2.3. The procedure is applied to both the weighted and the unweighed cluster configurations. For all the three representation spaces, the unweighed procedure often demonstrates better performances, while the weighted procedure leads to marginal, if any, improvements that do not justify its use.

Following the criteria mentioned in section 2.3 and considering all the three representation spaces, the suggested number of clusters obtained for Italian and Switzerland data is four.

This configuration is then checked, to assess if the regional FDCs are significantly different, using the procedure described in section 2.3 for all the three representation spaces.

The FDCs of the original four clusters cannot be considered significantly different from each other, neither in the linear space, nor in the other two logarithmic spaces. Thus, for each representation space, the two most similar clusters are merged together. The new configurations with three clusters can be accepted in the linear space only. Applying again the procedure for the logarithmic and log-normal space we obtain two configurations consisting of two clusters each. To select one among these different configurations of clusters, we perform the following cross-evaluation: for each set of clusters (e.g., the one obtained in the linear space), we check if the difference between the regional FDCs is significant in the other representation spaces (i.e., also in the logarithmic and log-normal spaces). Based on this cross-evaluation, we choose the configuration with 2 clusters obtained in the logarithmic space, which is represented in Figures 4 and 5. Hence, this latter configuration will be used as the result of the distance-based model.

The final regions obtained are shown in figure 4. Curves belonging to each cluster are grouped together and the regional curves are derived as the average of all curves belonging to the region. Figure 5 shows the regional curves (black lines) obtained from curves belonging to the cluster (grey lines) in the log-normal space. Although every curve bundle appears to be quite wide, regional curves are able to represent two characteristic behaviors. In fact, we can observe an almost straight curve and a “S” shaped curve.

A quantitative representation of model quality and estimation errors is reported in the following section, where a comparison against some parametric methods is performed.

4. Comparison with Parametric Models

The distance-based regional procedure developed in this work is tested against some standard parametric regional models. In general, the choice of the reference model is not trivial and more than one function can be used to describe the FDCs. For this purpose, a useful tool is the L -moments ratio diagram of figure 6 [Hosking and Wallis, 1997] where one plots the L_{CA} (coefficient of L -skewness) of each dimensionless FDC versus its corresponding L_{kur} (coefficient of L -kurtosis). The lines represent the domain of the distributions over the $L_{CA} - L_{kur}$ space and can help one to identify the distribution to be used. This approach has been followed, for example, by *Castellarin et al.* [2007].

In this work, the analysis is performed over a database of 95 basins that have very different characteristics in terms of L_{CA} and L_{kur} , as figure 6 shows. The scattering of the points make the choice of the distribution rather difficult. For this reason, different parametric models are used for the comparison with the distance-based procedure.

Each parameter θ of a parametric model is related to the catchments' descriptors d by a linear model of the form

$$\theta = a_0 + a_1 \cdot d_1 + a_2 \cdot d_2 + \dots + a_n \cdot d_n + \varepsilon. \quad (8)$$

The first step is to identify a suitable regional model to estimate the generic parameter for an ungauged station, where θ is previously estimated at each station s using a suitable technique. The resulting parameters θ_s are then related to descriptor data (raw data, not

distances) for all the catchments (not classified in regions) to identify a regional model (regression) able to describe them. Many linear models of the form of equation (8) are considered and validated with a t-Student test followed by a multicollinearity (VIF) test and subsequently ordered by their values of R_{adj}^2 [e.g., *Montgomery et al.*, 2001].

The models considered here are the two-parameter log-normal distribution (LN2), the three-parameter Pearson type III (PE3) and the generalized Pareto (GPA) distributions. The log-normal model is represented by the relation

$$\log(q) = \theta_1 + \theta_2 \cdot z \quad (9)$$

where z is the quantile of a normal distribution with zero mean and unit variance corresponding to each flow's plotting position values. In the log-normal probability representation, equation (9) is a straight line whose coefficients θ_1 and θ_2 can be estimated with a least squares linear regression.

The GPA probability density function is defined as

$$f(q) = \theta_2^{-1} \exp[-(1 - \theta_3)y], \quad (10)$$

with $y = -\theta_3^{-1} \log[1 - \theta_3(q - \theta_1)/\theta_2]$ if $\theta_3 \neq 0$ and $y = (q - \theta_1)/\theta_2$ if $\theta_3 = 0$, where θ_1 , θ_2 and θ_3 are the location, scale and shape parameter, respectively; the PE3 probability density function is defined as

$$f(q) = \frac{(q - \theta_1)^{\theta_2 - 1} \exp[-(q - \theta_1)/\theta_3]}{\theta_3^{\theta_2} \Gamma(\theta_2)}, \quad (11)$$

414

415 where θ_1 , θ_2 and θ_3 are the location, scale and shape parameter, respectively, and $\Gamma(\cdot)$ is
 416 the gamma function. For details about these distributions and for parameters estimation
 417 we refer to *Hosking and Wallis* [1997] and *Viglione* [2007]. The regional estimation of the
 418 models' parameters use the descriptors listed in table 3, whose definitions [*Viglione et al.*,
 419 2007a] are provided in the auxiliary material.

420 Our model and the parametric ones are all tested using a cross-validation approach in
 421 which one station is considered ungauged and its data are removed from the database.
 422 The models are then recalibrated using only the remaining data, and the unknown curve
 423 is estimated. After this procedure is repeated for all basins, one can compute, for each
 424 basin, the error measure $\delta_{\text{MOD,EMP}}$ as the distance between the estimated FDC and its
 425 empirical counterpart.

426 The non-parametric FDC representation method performs better than the parametric
 427 models for most of the analyzed basins, independently of the representation space con-
 428 sidered. Figure 7 shows a comparison between the errors $\delta_{\text{MOD,EMP}}$ calculated with the
 429 parametric and the distance-based approaches. Each parametric model is able to well
 430 describe only a subset of the studied basins (see figure 6), which is probably the reason
 431 why they demonstrate similar and non excellent performances when applied to the whole
 432 dataset.

5. Conclusions

433 The procedure for dimensionless flow duration curves estimation in ungauged basins
 434 developed in this work hinges on the concept of distance, that quantitatively represents
 435 the dissimilarity between curves and catchment's descriptors. This approach, based on

distance matrices, allows one to account for a FDC as a whole object, avoiding the description of the curve by means of a parametric function. Moreover, no assumptions on the shape of the FDCs is made. This is an important feature when one has to manage at the same time curves described by a simple geometry (e.g., almost straight lines in the log-normal probability plot) and curves with more complex behavior (e.g., “S” shaped curves). In fact, complex shapes can be well described by a parametric model only using an high number of parameters, that sometimes can not guarantee a robust parameters estimation.

The results obtained by means of the distance-based model (non-parametric representation of the FDC) applied to our dataset are comparable, and many times better, than the estimation yielded by classical parametric models of the same or greater complexity. These results are obtained on the basis of only two descriptors, while the log-normal model requires six descriptors for the assessment of two parameters, and the PE3 and GPA models require 8 and 10 descriptors to estimate their three parameters.

The main advantage of the method based on distance matrices is its ability in dealing with curves. For instance, the regionalization method proposed here could be improved considering also “complex” catchment descriptors as the hypsographic curve, or climatic information like the precipitation regime curve.

Acknowledgments. The study has been supported by the Italian Ministry of Education through the grants no. 2006089189 and no. 2007HBTS85. Comments and suggestions from A. Castellarin, T. Torgersen, E. Martins and two anonymous reviewers are gratefully acknowledged.

References

- 458 Burn, D. H. (1990), Evaluation of regional flood frequency analysis with a region of
459 influence approach, *Water Resources Research*, 26(10), 2257–2265.
- 460 Castellarin, A., G. Galeati, L. Brandimarte, A. Montanari, and A. Brath (2004a), Regional
461 flow-duration curves: reliability for ungauged basins, *Advances in Water Resources*,
462 27(10), 953 – 965.
- 463 Castellarin, A., R. Vogel, and A. Brath (2004b), A stochastic index flow model of flow
464 duration curves, *Water Resources Research*, 40(3).
- 465 Castellarin, A., G. Camorani, and A. Brath (2007), Predicting annual and long-term
466 flow-duration curves in ungauged basins, *Advances in Water Resources*, 30(4), 937 –
467 953.
- 468 Claps, P., and M. Fiorentino (1997), *Integrated Approach to Environmental Data Manage-*
469 *ment Systems, NATO-ASI series*, vol. 2 (31), chap. Probabilistic Flow Duration Curvers
470 for use in Environmental Planning and Management, pp. 255–266, Harmancioglu et al.,
471 Kluwer, Dordrecht, The Netherlands.
- 472 CUBIST Team (2007), Cubist project: Characterisation of ungauged basins by integrated
473 use of hydrological techniques, Geophysical Research Abstracts, Vol. 10, EGU2008-A-
474 12048, 2008 SRef-ID: 1607-7962/gra/EGU2008-A-12048 EGU General Assembly 2008.
- 475 Dalrymple, T. (1960), *Flood frequency analyses, Water Supply Paper*, vol. 1543-A, U.S.
476 Geological Survey, Reston, Va.
- 477 Fennessey, N., and R. Vogel (1990), Regional flow-duration curves for ungauged sites in
478 massachusetts, *Journal of Water Resources Planning and Management-ASCE*, 116(4),
479 530 – 549.

- Holmes, M., A. Young, A. Gustard, and R. Grew (2002), A region of influence approach to predicting flow duration curves within ungauged catchments, *Hydrology and Earth System Sciences*, 6(4), 721 – 731.
- Hosking, J., and J. Wallis (1997), *Regional Frequency Analysis: An Approach Based on L-Moments*, Cambridge University Press.
- Iacobellis, V. (2008), Probabilistic model for the estimation of t year flow duration curves, *Water Resources Research*, 44(2), W02,413.
- Kottegoda, N. T., and R. Rosso (1997), *Statistics, Probability, and Reliability for Civil and Environmental Engineers*, international ed., McGraw-Hill Companies.
- Legendre, P. (1993), Spatial autocorrelation - trouble or new paradigm, *ECOLOGY*, 74(6), 1659 – 1673.
- Legendre, P., F. Lapointe, and P. Casgrain (1994), Modeling brain evolution from behavior - a permutational regression approach, *Evolution*, 48(5), 1487 – 1499.
- Lichstein, J. (2007), Multiple regression on distance matrices: a multivariate spatial analysis tool, *Plant Ecology*, 188(2), 117 – 131.
- Mantel, N., and R. Valand (1970), A technique of nonparametric multivariate analysis, *Biometrics*, 27, 209 – 220.
- Montgomery, D., E. Peck, and G. Vining (2001), *Introduction to linear regression analysis*, third ed., Wiley Series in Probability and Statistics.
- Mosley, M. P., and A. I. McKerchar (1993), *HandBook of Hydrology*, chap. 8, p. 39, McGraw-Hill Companies, international edition.
- NASA (2000), SRTM.

- 502 R Development Core Team (2007), *R: A Language and Environment for Statistical Com-*
503 *puting*, R Foundation for Statistical Computing, Vienna, Austria, ISBN 3-900051-07-0.
- 504 Rigon, R., and F. Zanotti (2002), *The Fluid Turtle Library, Users and Programmers*
505 *Guide*, University of Trento, Italy.
- 506 Singh, R., S. Mishra, and H. Chowdhary (2001), Regional flow-duration models for large
507 number of ungauged himalayan catchments for planning microhydro projects, *Journal*
508 *of Hydrologic Engineering*, 6(4), 310 – 316.
- 509 Sivapalan, M., et al. (2003), IAHS decade on predictions in ungauged basins (pub), 2003-
510 2012: Shaping an exciting future for the hydrological sciences, *Hydrological Sciences -*
511 *Journal - des Sciences Hydrologiques*, 48(6), 857–880.
- 512 Smakhtin, V. (2001), Low flow hydrology: a review, *Journal of Hydrology*, 240(3-4), 147
513 – 186.
- 514 Smouse, P., J. Long, and R. Sokal (1986), Multiple-regression and correlation extensions
515 of the mantel test of matrix correspondence, *Systematic Zoology*, 35(4), 627 – 632.
- 516 Viglione, A. (2007), *nsRFA: Non-supervised Regional Frequency Analysis*, r package ver-
517 sion 0.4-5.
- 518 Viglione, A., P. Claps, and F. Laio (2007a), *Water resources assessment and management*
519 *under water scarcity scenarios*, chap. Mean annual runoff estimation in North-Western
520 Italy, pp. 97–121, La Loggia et al., CSDU, Milano.
- 521 Viglione, A., F. Laio, and P. Claps (2007b), A comparison of homogeneity tests for regional
522 frequency analysis, *Water Resources Research*, 43(3), W03,428.
- 523 Vogel, R., and N. Fennessey (1994), Flow-duration curves .2. new interpretation and
524 confidence-intervals, *Journal of Water Resources Planning and Management-ASCE*,

525 $120(4)$, 485 – 504.

526 Ward, J. (1963), Hierarchical grouping to optimize an objective function, *Journal of the*
527 *American Statistical Association*, 58, 236–244.

Figures captions

Table 1 Regression models with two descriptors that well describe the relationship between curve distance matrix and descriptors distance matrices. All the models pass the Mantel test (significance of regression coefficients) with a level of significance of 0.05 and the VIF test (multicollinearity) with threshold equal to 5. The curve distance matrix is calculated in three different representation spaces: the linear, the logarithmic and the log-normal one.

Table 2 Brief description and range of variation of the descriptors used by the distance-based models (see table 1).

Table 3 Descriptors used to estimate the parametric model's parameters with level of significance (Student test) and Variance Inflation Factor test.

Figure 1 Comparison of dimensionless flow duration curves. Sampling points with constant spacing in frequency representation (a), and with a denser presence on the FDC tails due to normal transformation (b)

Figure 2 Distance between two FDCs calculated following equation (3). The three panels show a pair of FDCs in three different representation spaces: panel (a) is the linear representation (flow values versus exceedance frequency); panel (b) is the logarithmic representation in which discharges are log-transformed; panel (c) represents the log-normal probability plot in which the abscissa is the normal reduced variate z .

Figure 3 Geographical location of the gauging stations of the 95 catchments considered in the study. Basins 1 to 59 belong to Switzerland, while the remaining ones are located in the Northwestern part of Italy, in Piemonte and Valle d'Aosta regions. For additional information see auxiliary material at <http://www.idrologia.polito.it/~ganora>.

Figure 4 Disjoint regions in the space of catchment descriptors: Hmin is the minimum basin elevation and MHL is the mean hillslope length. The dashed lines represent the boundaries between the 4 clusters obtained before merging the clusters whose FDCs cannot be considered significantly different. The final 2 disjoint regions are separated by the solid line.

Figure 5 Flow duration curves grouped by cluster (in grey) and corresponding regional curves (in black).

Figure 6 L-moments ratio diagram for the dimensionless FDCs of the 95 basins (filled circles for Switzerland data and white circles for Italian data). The lines indicate different theoretical three-parameter distributions: generalized logistic (GLO), generalized extreme-value (GEV), generalized Pareto (GPA), lognormal (LN3), Pearson type III (PE3).

Figure 7 Quality of estimated dimensionless FDCs by the distance-based method compared with the log-normal model (a), the generalized Pareto model (b) and the Pearson type III model (c). The distance between the empirical curve and the estimated one $\delta_{\text{MOD,EMP}}$ is reported in the scatter plot for each considered basin. The solid line represents the ratio 1:1 between the errors, while dashed lines delimit the areas where errors for the distance-based model are twice the parametric ones, and viceversa. Points above the solid line represent curves better estimated by the distance-based method; points above the upper dashed line represent curves much better estimated by the distance-based method.

Table 1. *Ganora et al.* [2009]

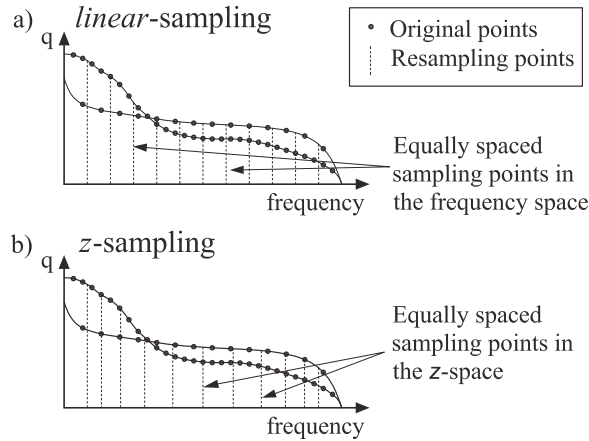
Best relation	Representation space		
	Linear	Logarithmic	Log-normal
1st	H + MHL	Hmin + MHL	Hmin + MHL
2nd	Hmin + MHL	Hmin + Pm	H + MHL
3rd	H + Slo	Pm + MHL	Pm + MHL

Table 2. *Ganora et al.* [2009]

Descriptor	Definition	Min	Mean	Max
H	mean elevation of the drainage basin above sea level (m)	475	1665	2847
Hmin	minimum elevation of the drainage basin above sea level (m)	82	839	1974
MHL	mean hillslope length (m)	584.1	759.5	973.6
Slo	average of the slope values associated to each pixel in the DEM of the drainage basin (%)	4	39.9	61.6
Pm	mean large-scale slope (%)	0.8	15.7	50.1

Table 3. *Ganora et al.* [2009]

Model	Parameter	Descriptors	Student	VIF	R_{adj}^2
Lognormal	θ_1	asp, Cc	< 0.05	< 5	0.12
	θ_2	Xb, PLDP, slo, MHL	< 0.05	< 5	0.17
GPA	θ_1	Xmax, PLDP, slo, MHL	< 0.02	< 5	0.28
	θ_2	Ymin, IPS25, cos(or)	< 0.05	< 5	0.54
	θ_3	Xc, Yc, IPS50	< 0.05	< 5	0.39
PE3	θ_1	Xmax, PLDP, slo, Cc, MHL	< 0.02	< 5	0.39
	θ_2	Xc, Ymin, IPS100, Cc	< 0.05	< 5	0.31
	θ_3	Ymin, PS50	< 0.02	< 5	0.28

**Figure 1.** *Ganora et al.* [2009]

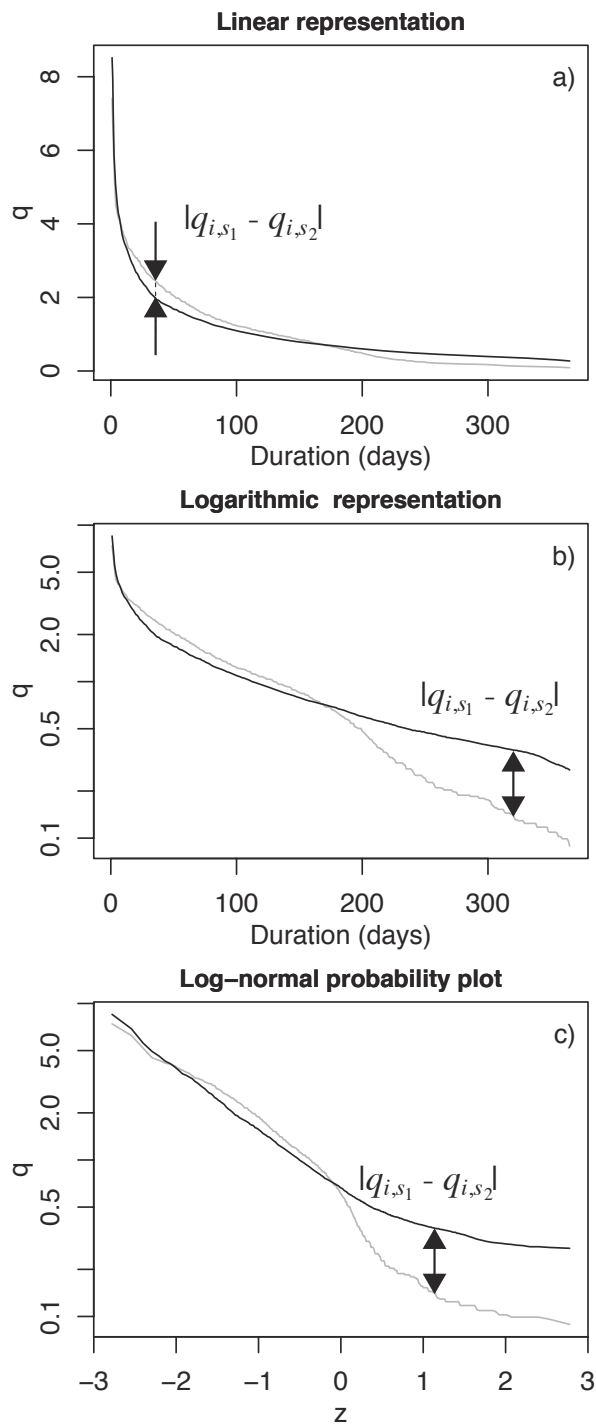


Figure 2. *Ganora et al.* [2009]

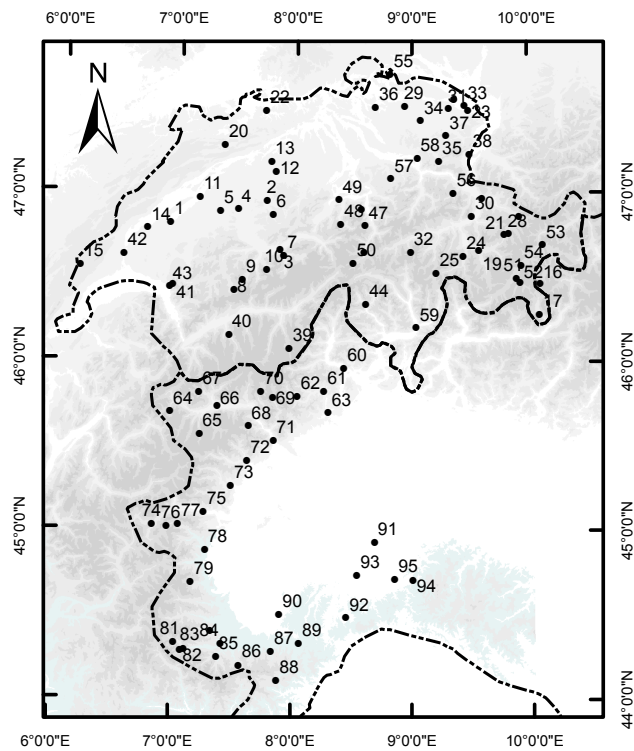


Figure 3. *Ganora et al.* [2009]

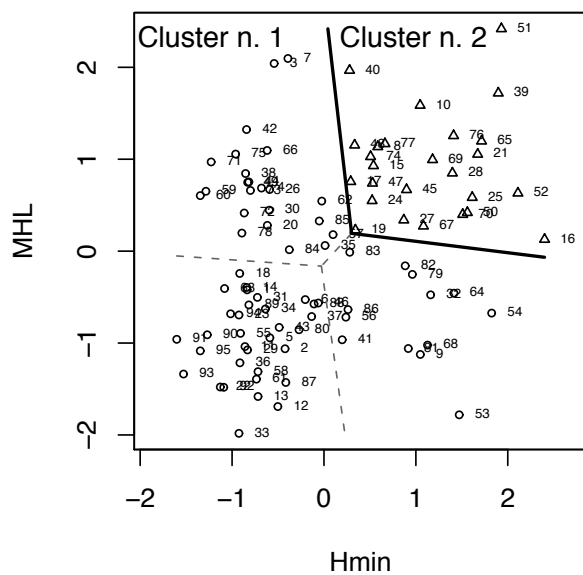


Figure 4. *Ganora et al.* [2009]

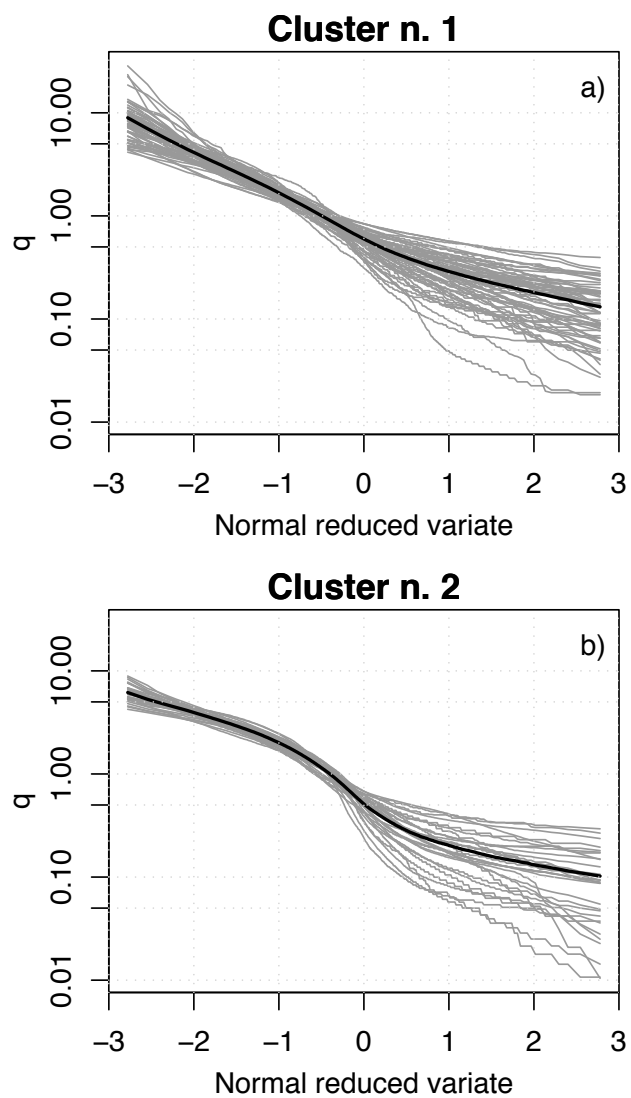


Figure 5. *Ganora et al.* [2009]

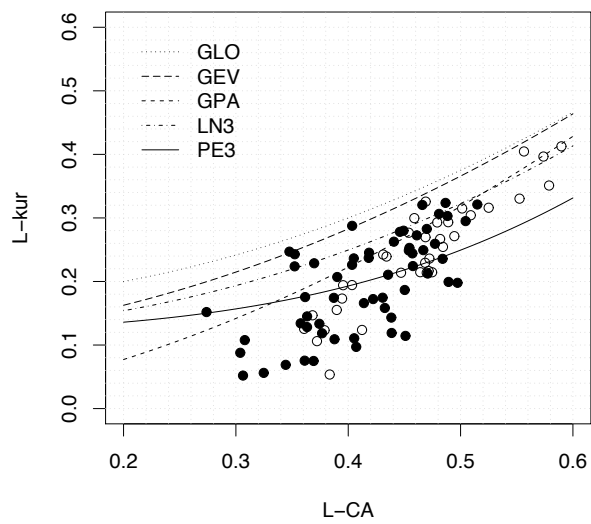


Figure 6. *Ganora et al.* [2009]

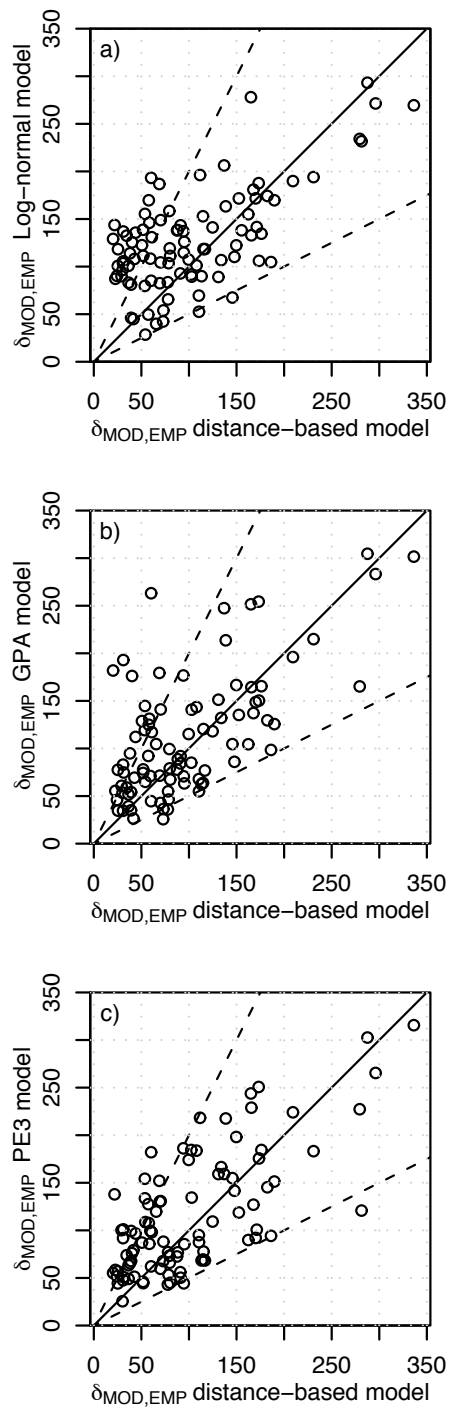


Figure 7. Ganora et al. [2009]