

Conference or Workshop Item (Paper)

Kjeldsen, Thomas R.; Jones, David A.. 2007 Recursive estimation of a hydrological regression model. In: *ASCE-EWRI, World Water & Environmental Resources Congress 2007, Tampa, FLORIDA, USA, 15-19 May 2007*. USA, ASCE-EWRI.

©2007 ASCE

This version available at <http://nora.nerc.ac.uk/2590/>

NERC has developed NORA to enable users to access research outputs wholly or partially funded by NERC. Copyright and other rights for material on this site are retained by the authors and/or other rights owners. Users should read the terms and conditions of use of this material at <http://nora.nerc.ac.uk/policies.html#access>

Contact CEH NORA team at
nora@ceh.ac.uk

Recursive estimation of a hydrological regression model

T. R. Kjeldsen¹, and D. A. Jones¹

¹*Centre for Ecology & Hydrology, Maclean Building, Crowmarsh Gifford, Wallingford, OX10 8BB, UK; email trkj@ceh.ac.uk, daj@ceh.ac.uk*

Abstract

The use of the generalised least square (GLS) technique for estimation of hydrological regression models has become good practice in hydrology. Through a regression model, a simple link between a particular hydrological variable and a set of catchment descriptors can be established. The regression residuals can be treated as the sum of sampling errors in the hydrological variable and errors in the regression model. This paper presents a recursive method for estimating a parameterised form of the cross correlation between the regression model errors, the variance of these errors and the regression model parameters. A re-weighted set of regression residuals can be defined such that the covariance of these residuals is essentially similar to that of the model error. The cross products of the re-weighted regression residuals, pooled within bins, can be used to identify a structure and to fit a parameterised form for the cross-correlations of the regression errors. The procedure has been tested successfully on annual maximum flow data from 602 catchments located throughout the UK.

Introduction

The use of linear regression models figures prominently among methods for deriving simple relationships between hydrological variables and a set of lumped catchment descriptors such as catchment area, annual average rainfall and soil type. This is partly due to regression models being computationally easy to use and being much less demanding with regards to data requirements than more detailed hydrological models. A well known hydrological variable is the index flood, as required by the index flood method for deriving flood frequency relationships (Stedinger *et al.*, 1993).

The objective of this study is to develop and implement an extended form of Generalised Least Squares (GLS) regression in which the regression residuals are treated as being the sum of two types of error, a sampling error and a modelling error, and where both types of error are spatially correlated. The procedure outlined here provides a direct non-parametric estimate of the relation to distance of the cross correlation between the regression modelling errors which can be used to identify and estimate a parametric form of this function. Overall, the procedure is a recursive one which provides estimates of the regression parameters and of the variance and correlation of the modelling errors, given that an initial separate analysis provides estimates of the variance and correlation of the sampling errors. The method has been tested on a dataset consisting of annual maximum instantaneous peak flow series from 602 rural catchments located throughout the UK.

Model Description

To relate the index flood variable from n different catchments to a set of catchment descriptors, consider a vector of sample (log transformed) median annual maximum floods, \mathbf{y} , where individual sites are denoted with a subscript i . Each sample value is described in terms of a population regression model and two individual error components representing the sampling, ε_i , and modelling, η_i , errors, respectively so that

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \varepsilon_i + \eta_i = \mathbf{x}_i^T \boldsymbol{\beta} + \omega_i \quad (1)$$

where $\boldsymbol{\beta}$ is a vector of regression model parameters and \mathbf{x}_i is a vector of catchment descriptors with a value of one in the first location. The covariance of the sampling errors is denoted $\boldsymbol{\Sigma}_\varepsilon$, the corresponding covariance of the modelling errors denoted $\boldsymbol{\Sigma}_\eta$, and the two errors are assumed mutually independent. Further, it is assumed that the elements along the diagonal of the modelling error covariance are identical and equal to σ_η^2 . In pioneering the use of the GLS procedure in hydrology, Stedinger and Tasker (1989) assumed the modelling covariance matrix to be of the form $\boldsymbol{\Sigma}_\eta = \sigma_\eta^2 \mathbf{I}$, i.e. there is an assumption of no cross correlation between the modelling errors. In contrast, the model formulated here assumes the cross correlation to be represented by the associated modelling error correlation matrix \mathbf{R}_η .

While estimates of the sampling error covariance can be obtained directly from the dataset, the covariance of the modelling errors has to be estimated as part of a recursive procedure. From an initial guess of the modelling error covariance, a set of regression residuals can be estimated. By re-weighting these residuals, it is possible to obtain a set of GLS residuals from which the modelling error variance can be estimated. By further re-weighting the GLS residuals, an estimate of the modelling error correlation matrix can be obtained. These recursive estimates can then be used to estimate a new regression model and a new set of regression residuals. This procedure is continued until the modelling error variance σ_η^2 has converged.

The first step in the recursive procedure is to define the covariance matrix of the vector $\boldsymbol{\omega}$ of total errors as

$$E\{\boldsymbol{\omega}\boldsymbol{\omega}^T\} = \boldsymbol{\Sigma}_\omega = \boldsymbol{\Sigma}_\eta + \boldsymbol{\Sigma}_\varepsilon = \sigma_\eta^2 (\mathbf{R}_\eta + \boldsymbol{\Sigma}_\varepsilon / \sigma_\eta^2) = \sigma_\eta^2 \mathbf{G}. \quad (2)$$

To implement the procedure, the expression in Eq. (2) is interpreted as representing the covariance of the total error in terms of σ_η^2 , being the value to be estimated from the present step of the recursive procedure, and of \mathbf{G} , a known matrix derived from values of σ_η^2 and \mathbf{R}_η , which are either initial guesses or the estimates obtained in the previous step. In the expressions developed below, Eq. (2) is taken temporarily to be valid even though an estimated value of \mathbf{G} is used.

It can be shown that the individual estimates of the overall residuals, $\hat{\omega}_i$, can be expressed in terms of the true underlying residuals as

$$\hat{\boldsymbol{\omega}} = \mathbf{V}\boldsymbol{\omega}, \quad \mathbf{V} = \mathbf{I} - \mathbf{X}(\mathbf{X}^T \mathbf{G}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{G}^{-1} \quad (3)$$

which enables the covariance matrix of the estimated regression residuals to be represented as

$$\boldsymbol{\Sigma}_{\hat{\boldsymbol{\omega}}} = E\{\hat{\boldsymbol{\omega}}\hat{\boldsymbol{\omega}}^T\} = \sigma_{\eta}^2 [\mathbf{G} - \mathbf{X}(\mathbf{X}^T \mathbf{G}^{-1} \mathbf{X})\mathbf{X}^T]. \quad (4)$$

GLS residuals

For Generalised Least Squares analysis, it is common to work with an alternative set of sample residuals, the GLS residuals. These residuals, $\tilde{\boldsymbol{\omega}}$, can be related to the “raw” sample residuals, $\hat{\boldsymbol{\omega}}$, in the following way. A matrix-square-root of the scaled covariance matrix \mathbf{G} is first required, and it convenient to work with the Cholesky decomposition:

$$\mathbf{G} = \mathbf{U}_{\mathbf{G}}^T \mathbf{U}_{\mathbf{G}} \quad (5)$$

where $\mathbf{U}_{\mathbf{G}}$ is an upper triangular matrix. The sample GLS residuals are defined as

$$\tilde{\boldsymbol{\omega}} = \mathbf{U}_{\mathbf{G}}^{-T} \hat{\boldsymbol{\omega}} = \mathbf{U}_{\mathbf{G}}^{-T} (\mathbf{y} - \hat{\mathbf{y}}) \quad (6)$$

Given the assumption that the value of \mathbf{G} being temporarily used is correct, an unbiased estimate of σ_{η}^2 is provided by

$$\hat{\sigma}_{\eta}^2 = (N - p)^{-1} \sum_{i=1}^N \tilde{\omega}_i^2 \quad (7)$$

and, given the assumption, this is the minimum variance unbiased estimate for σ_{η}^2 . The estimated value of σ_{η}^2 can then be carried forward to the next step of the recursion.

Re-weighted GLS

To obtain an estimate of the modelling error correlation matrix \mathbf{R}_{η} , a re-weighted version of the GLS residuals is constructed: these can also be considered as a re-weighting of the raw residuals. In parallel with Eq. (5), a Cholesky decomposition of the correlation matrix is constructed, so that

$$\mathbf{R}_{\eta} = \mathbf{U}_{\eta}^T \mathbf{U}_{\eta} \quad (8)$$

where, again, \mathbf{U}_η is an upper triangular matrix. In implementing this scheme, the matrix \mathbf{R}_η used is the estimate available at the start of the particular step of the recursion. Then a set of re-weighted GLS residuals, $\tilde{\tilde{\boldsymbol{\omega}}}$, can be calculated as

$$\tilde{\tilde{\boldsymbol{\omega}}} = \mathbf{U}_\eta^T \tilde{\boldsymbol{\omega}} = \mathbf{U}_\eta^T \mathbf{U}_G^{-T} \hat{\boldsymbol{\omega}} = \mathbf{U}_\eta^T \mathbf{U}_G^{-T} (\mathbf{y} - \hat{\mathbf{y}}). \quad (9)$$

The covariance matrix for the re-weighted GLS residuals is given as

$$\begin{aligned} E\{\tilde{\tilde{\boldsymbol{\omega}}}\tilde{\tilde{\boldsymbol{\omega}}}^T\} &= E\{\mathbf{U}_\eta^T \mathbf{U}_G^{-T} \hat{\boldsymbol{\omega}} \hat{\boldsymbol{\omega}}^T \mathbf{U}_G^{-1} \mathbf{U}_\eta\} = \mathbf{U}_\eta^T \mathbf{U}_G^{-T} \boldsymbol{\Sigma}_\omega \mathbf{U}_G^{-1} \mathbf{U}_\eta \\ &= \sigma_\eta^2 \left[\mathbf{R}_\eta - \mathbf{U}_\eta^T \mathbf{U}_G^{-T} \mathbf{X} (\mathbf{X}^T \mathbf{G}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U}_G^{-1} \mathbf{U}_\eta \right] \end{aligned} \quad (10)$$

Thus, the raw residual vector, $\hat{\boldsymbol{\omega}}$, has been rescaled to form a revised residual vector, $\tilde{\tilde{\boldsymbol{\omega}}}$, which, apart from the use of estimated values to form the re-weighting matrix (\mathbf{G}), have a correlation matrix close to \mathbf{R}_η .

Case study

To test the recursive GLS procedure, a case study was undertaken which involved annual maximum instantaneous peak flow data from 602 rural catchments located throughout the UK. Each catchment is associated with 5 different catchment descriptors found in a previous study by the Institute of Hydrology (1999) to be useful for estimating the median of the annual maximum peak flow through regression modelling. A summary of the data is shown in Table 1 where AREA is the catchment area in km^2 , SAAR is the standard average annual rainfall (in mm) for the period 1961-1990, FARL is an index of flood attenuation due to reservoirs and lakes and both SPRHOST and RESHOST describe the hydrological properties of catchment soils. The FARL descriptor can take on values between zero and one and SPRHOST values are in the range between 0% and 60%.

Table 1: Data from 602 rural catchment located throughout the UK.

	Min	Mean	Max
Median, $\text{m}^3 \text{s}^{-1}$	0.2	92.7	981.4
Record length, years	4	33	117
AREA, km^2	1.6	335.1	4587.0
SAAR, mm	558	1162	2848
FARL, -	0.645	0.970	1.000
SPRHOST, %	5.1	37.4	59.9
RESHOST, -	-0.15	0.00	0.19

A further description of the catchment descriptors is provided by the Institute of Hydrology (1999) and Kjeldsen and Jones (2006). The actual regression model investigated in this study is based on log-transformed values of AREA, SAAR/1000, FARL and SPRHOST/100. In addition, a quadratic term, $\ln[\text{AREA}]^2$, and non-transformed values of RESHOST are included. For further background to the

variable selection, please refer to the comprehensive study reported by Institute of Hydrology (1999).

Sampling Error. Both the diagonal as well as the off-diagonal elements of the sampling error covariance are estimated based on consideration of the asymptotic variance of the sampling median and are given as

$$\Sigma_{\varepsilon,ij} = \begin{cases} 4\beta_i^2 / n_i & i = j \\ 4\beta_i\beta_j \frac{n_{ij}}{n_i n_j} \rho_{y_i y_j} & i \neq j \end{cases} \quad (11)$$

where β_i is the scale parameter of the GLO distribution, standardised to have unit median, estimated using the method of L-moments as described by Institute of Hydrology (1999). Here n_{ij} denotes the number of years for which catchments i and j both have data, while n_i and n_j denote the total numbers of years of data for the two catchments separately. Note that there is a minor conflict between conventional notations used for the GLO distribution and for regression analysis in the use of “beta” with two distinct meanings. In addition, estimation of the off-diagonal elements requires estimates of the correlation coefficient between the log transformed median annual maximum flood for each pair of sites, $\rho_{y_i y_j}$.

A bootstrap experiment was carried out similar to the experiments used by Kjeldsen and Jones (2006) for investigating the cross-correlation between L-moment ratios. Bootstrapping is a technique where new samples are created from an original sample by randomly selecting (with replacement) observations from the original sample. Considering the annual maximum series of peak flow from the 602 rural catchments, a total of 11062 pairs of gauges with a minimum of 40 years of overlapping record were available. To investigate the cross-correlation between the log-median annual maximum peak flow and relate it to geographical distance between catchment centroids, each of these pairs were analysed in turn. For each station pair, a new bootstrap sample was created for each station by randomly (with replacement) selecting years in the overlapping record. From each selected year the joint pair of observations was transferred to the joint bootstrap sample, thereby preserving the cross-correlation between the annual maximum series of the two sites. The selection is continued until the new bootstrap sample has a record length equal to the length of the overlapping record in the original sample. From the joint bootstrap sample, the medians of the log transformed annual maximum peak flows are estimated for both sites and recorded. By creating 1000 new bootstrap samples for each station pair, the correlation between the medians can be estimated and linked to the distance between catchment centroids as

$$\rho_{y_i y_j} = \theta \exp(-\phi_1 d_{ij}) + (1 - \theta) \exp(-\phi_2 d_{ij}) \quad (12)$$

where d_{ij} is the distance (km^2) between centroids of catchments i and j . The three parameters θ , ϕ_1 and ϕ_2 are estimated using a least-squares technique. The outcome

of the bootstrapping experiment is shown in Figure 1, together with the correlation function that has been fitted.

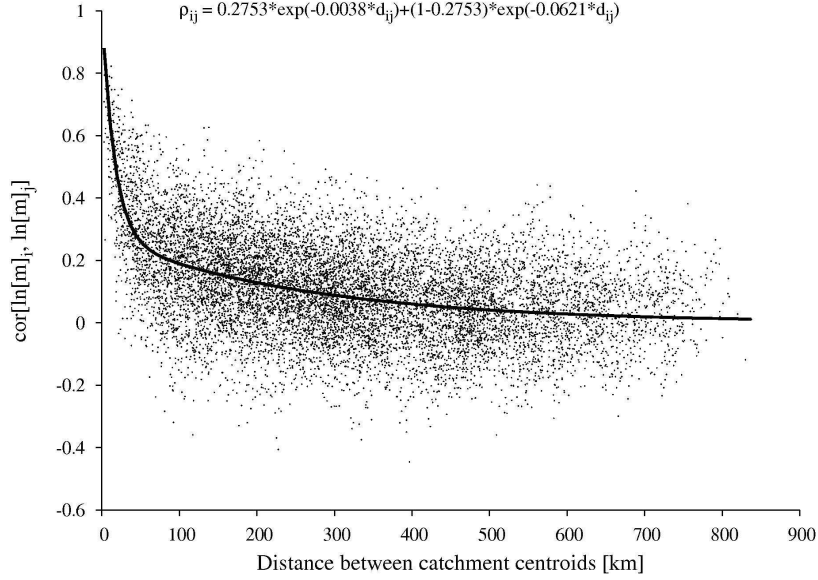


Figure 1: Correlation between sampling errors of log-transformed median annual maximum flood as a function of distance between catchment centroids.

As the estimator of the at-site sampling variability of y in Eq. (11) involves an estimate of the median of the annual maximum peak flow itself, it was considered appropriate to replace the direct estimates of the GLO parameter β in Eq. (11) with corresponding estimates derived using an ordinary least squares (OLS) regression model linking $\ln[\beta_i]$ to a set of catchment descriptors as

$$\ln[\beta_i] = \theta_0 + \sum_{p=1}^P \theta_p \ln[x_{i,p}] + \gamma_i \quad (13)$$

where P is the total number of catchment descriptors used in the regression model, $x_{i,p}$ is the value of the p 'th catchment descriptor for the i 'th catchment and θ_p is the p 'th regression model parameter. Only a limited investigation has been made of the errors, γ_i : the results of the OLS regression are reported (Table 2) as if they can be assumed to be independent and normally distributed with mean zero and variance σ_γ^2 , whereas the errors are very likely to be correlated between the catchments. Thus the estimates of the standard errors of the regression parameters are likely to be too small. The use of OLS estimates at this stage rather than GLS estimates is not thought to be important.

Table 2: Summary statistics for regression model describing $\ln[\beta_i]$.

Coefficient	Parameter θ_p	Standard error	t-value	p-value
Intercept (θ_0)	-1.1221	0.0664	-16.906	$< 2^{-16}$
Ln[AREA]	-0.0816	0.0105	-7.783	$3.12 \cdot 10^{-14}$
Ln[SAAR/1000]	-0.4580	0.0401	-11.431	$< 2^{-16}$
Ln[BFIHOST]	0.1065	0.0520	2.049	0.0409
$\sigma_\gamma^2 = 0.107 \quad df = 598 \quad r^2 = 0.28$				

The regression model has an r^2 value of only 28% which indicates less predictive power than could have been hoped for, but this relates to the substantial sampling error in the estimates of the GLO scale parameters. To estimate the sampling covariance Σ_ε estimates of β obtained through Eq. (13) are substituted into Eq. (11).

Modelling error. The two components of the modelling error covariance (the modelling error variance, σ_η^2 , and the modelling error correlation matrix, \mathbf{R}_η) can now be estimated using the recursive procedure outlined in the previous section. While the procedure provides recursive estimates of both σ_η^2 and \mathbf{R}_η , at present only the former is used to determine if the procedure has converged. To start the iterations, it is necessary to make initial guesses of the values of both σ_η^2 and \mathbf{R}_η . In this study, a large value ($\sigma_\eta^2 = 100$) was chosen for the modelling error variance, and a unit matrix ($\mathbf{R}_\eta = \mathbf{I}$) for the modelling error correlation.

The first step recursion is to transform the estimated raw residuals, $\hat{\omega}$, into the corresponding GLS residuals, $\tilde{\omega}$, through Eq. (11). Using the residual sum of squares of the GLS residuals, an iterative estimate of the modelling error variance σ_η^2 is obtained from Eq. (7). Next, the off-diagonal elements of \mathbf{R}_η are estimated. By re-ordering Eq. (10) it can be ascertained that

$$\frac{E\{\tilde{\omega}\tilde{\omega}^T\}}{\sigma_\eta^2} + \mathbf{B} = \mathbf{R}_\eta \quad (14)$$

where the matrix \mathbf{B} is a bias correction given as $\mathbf{B} = \mathbf{U}_\eta^T \mathbf{U}_G^{-T} \mathbf{X} (\mathbf{X}^T \mathbf{G}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \mathbf{U}_G^{-1} \mathbf{U}_\eta$. If all catchment-pairs within a specified distance interval, between 1 km and 2 km say, are grouped together, and it is assumed that the correlation between the modelling errors depends only on distance, then an estimate of the average correlation for this distance interval, $r_{\eta,d}$ can be obtained as

$$r_{\eta,d} = \frac{1}{n_d} \sum_{k=1}^{n_d} \left[\frac{(\tilde{\omega}_i \tilde{\omega}_j)_k}{\sigma_\eta^2} + (b_{ij})_k \right] \quad (15)$$

where b_{ij} are elements in the bias correction matrix \mathbf{B} and where k in the summation represents the k 'th (out of n_d) pair of catchments i and j whose inter-centroid distance is in the d 'th bin. The final step in the recursive procedure is to fit a distance-based

function to the estimates of $r_{\eta,d}$. In this study, distances up to 800 km were investigated with interval lengths of 4 km, i.e. a total of 200 bins. Based on initial trial runs it was found that the weighted sum of two exponential-type functions

$$r_{\eta,d} = \psi \exp[-\varphi_1 d] + (1 - \psi) \exp[-\varphi_2 d], \quad (16)$$

gives a reasonable fit to the average correlation values derived in Eq. (15) when fitted using a simple least square technique. Here φ_1 , φ_2 and ψ are model parameters and d is the distance. When finally constructing the recursive estimate of the modelling error correlation matrix \mathbf{R}_η , the value of each off-diagonal element is derived from Eq. (16) substituting d with the actual distance d_{ij} between the catchment-pair being considered.

Results. With a tolerance level of 10^{-4} , a total of 23 iterations were needed for the modelling error variance to have converged. The resulting regression model statistics are shown in Table 3 and the estimated average bias corrected residual cross-product, $r_{\eta,d}$, along with the fitted weighted exponential function are shown in Figure 2. The double exponential function in Figure 2 has been fitted to data from the 200 bins. To ensure convergence of the recursive procedure, it was necessary to replace the estimate of σ_η^2 in Eq. (15) with $(\sigma_{\eta,i-1}^2 + \sigma_{\eta,i}^2)/2$, where the extra subscript i indicate the iteration. Figure 2 also shows the average residual cross-product between the GLS residuals. For these residuals, the bias corrected cross-products are expected to be close to zero, which appears reasonable from a visual inspection of Figure 2.

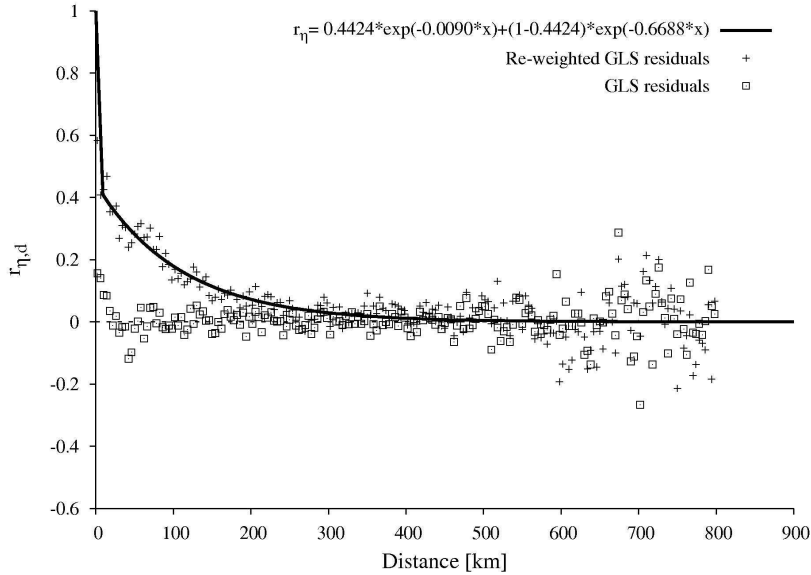


Figure 2: Bias corrected correlations for both GLS residuals and re-weighted GLS residuals.

Note that the procedure used for assigning pairs of catchments to the different distance bins is likely to result in a different number of pairs in each bin. In fact, the increase in the scatter of both types of residuals at short and large distances is probably due to the relatively smaller number of catchment-pairs allocated to these bins.

Table 3: Summary statistics for regression model describing $\ln[y_i]$.

Coefficient	Parameter β_p	Standard error	t-value	p-value
Intercept (β_0)	0.1010	0.4631	0.218	0.827
Ln[AREA]	0.9967	0.1509	6.605	$8.82 \cdot 10^{-11}$
Ln[AREA] ²	-0.0140	0.0148	-0.946	0.345
Ln[SAAR/1000]	1.7505	0.2300	7.611	$1.07 \cdot 10^{-13}$
Ln[FARL]	3.7763	0.6984	5.407	$9.29 \cdot 10^{-8}$
Ln[SPRHOST/100]	1.1228	0.1203	9.333	$< 2 \cdot 10^{-16}$
RESHOST	-3.7959	0.1008	-3.764	$1.84 \cdot 10^{-4}$
$\sigma_\eta^2 = 0.155 \quad df = 595 \quad r^2 = 0.938$				

The results in Table 3 indicate that most catchment descriptors included in the regression model have coefficients significantly different from zero. One possible exception is the Ln[AREA]² term. However, for the purpose of testing the recursive GLS procedure, no further attempts to adjust the regression model was undertaken.

Conclusion

This paper has outlined the development of a recursive procedure for estimating hydrological regression models. The procedure is considered an extension of the GLS model presented by Stedinger and Tasker (1989) by allowing the regression model errors to be cross-correlated. Initial testing of the procedure, on a dataset consisting of annual maximum series of instantaneous flow from 602 catchments located throughout the UK, has provided promising results in terms of estimating the regression model parameters. Some problems with regards to non-convergence in certain instances are still evident and require further attention. These problems are particularly evident for very simple regression models where the index flood is modelled using only very few catchment descriptors, for example using AREA only.

The procedure provides a method for verifying the existence or not of correlation between the modelling errors. Once a functional form for this cross correlation has been identified, it is likely that a more efficient procedure for estimating the regression model parameters (and, indeed the overall set of model parameters) can be developed using maximum-likelihood or Bayesian techniques. However, the more exploratory approach described here has two benefits. Firstly it allows consideration to be given to other ways of defining a distance to be used within the correlation function, for example taking into account river-network connectivity. Secondly, it allows some extra quality control of large datasets to be

made through the investigation of any anomalous correlations calculated for the distance-bins.

In the UK it is recommended practice that estimates of the index flood obtained at an ungauged site by using a regression model should, if possible, be adjusted through transfer of data from a nearby similar gauged catchment. This was on the basis that regression errors at nearby catchments were expected to be similar. It was shown by Kjeldsen and Jones (2007) that the best use of the transferred data depends on the correlation between the regression modelling errors at the two sites, as does the benefit obtained from the transfer in terms of improved prediction variance. Thus, the functional form of the model error correlation estimated in this study, i.e. the estimated form of Eq. (11) shown in Figure 2, can potentially become part of an improved procedure for data transfer in estimation of the index flood in the UK.

References

- Institute of Hydrology (1999) *Flood Estimation Handbook*, 5 Volumes, Institute of Hydrology, Wallingford, UK.
- Kjeldsen T. R. and D. A. Jones (2006) Prediction uncertainty in a median-based index flood method using L moments, *Water Resources Res.*, 42, W07414, doi:10.1029/2005WR004069
- Kjeldsen T. R. and D. A. Jones (2007) Estimation of the index flood using data transfer in the UK. *Hydrological Sciences Journal*, 52(1) (to appear).
- Stedinger, J. R. and G. D. Tasker (1989) An operational GLS model for hydrological regression. *Journal of Hydrology*, 111, 361-375.
- Stedinger, J. R., R. M. Vogel and E. Foufoula-Georgiou (1993) Frequency analysis of hydrological extreme events. In: *Handbook of Hydrology* (ed. D. Maidment), Chapter 18, McGraw-Hill, New York.