# REGIONAL FREQUENCY ANALYSIS AT UNGAUGED SITES

# WITH THE GENERALIZED ADDITIVE MODEL

F. Chebana*[1], C. Charron[2], T.B.M.J. Ouarda[2,1] and B. Martel[1]

[1]INRS-ETE, University of Quebec, 490 de la Couronne, Québec (Qc), Canada, G1K 9A9

[2] Institute Center for Water and Environment (iWATER),
Masdar Institute of Science and Technology,
P.O.Box 54224, Abu Dhabi, UAE

*Corresponding author

Email: fateh.chebana@ete.inrs.ca

Tel: +1 418 654 2542

Revised version

July 2014

## Abstract

The log-linear regression model is one of the most commonly used models to estimate flood quantiles at ungauged sites within the regional frequency analysis (RFA) framework. However, hydrological processes are naturally complex in several aspects including nonlinearity. The aim of the present paper is to take into account this nonlinearity by introducing the generalized additive model (GAM) in the estimation step of RFA. A neighbourhood approach using canonical correlation analysis (CCA) is used to delineate homogenous regions. GAMs possess a number of advantages such as flexibility in shapes of the relationships as well as the distribution of the output variable. The regional model is applied on a dataset of 151 hydrometrical stations located in the province of Québec, Canada. A stepwise procedure is employed to select the appropriate physio-meteorological variables. A comparison is performed based on different elements (regional model, variable selection and delineation). Results indicate that models using GAM outperform models using the log-linear regression as well as other methods applied to this dataset. In addition, GAM is flexible and allows including and showing non linear effects of explanatory variables, in particular basin area effect (scale). Another finding is the reduced effect of CCA delineation when combined with GAM.

## Keywords

# 1. Introduction

Knowledge of flood characteristics is very important for resource management and design of hydraulic structures. Estimation of design flows is often needed at locations where little or no information is available. In this case, regional frequency analysis (RFA) is often used for the estimation of flow characteristics. Ouarda et al. (2008) presented a detailed review of the various available RFA methods (Blöschl et al. 2013). Generally, RFA is composed of two main steps: the identification of groups of hydrologically homogeneous basins and the application of a regional estimation method within each delineated region (GREHYS 1996a; Ouarda 2013). Since flow characteristics are highly dependent upon physiographical and meteorological basin characteristics, these can be used to estimate flood quantiles at un-gauged sites. The hydrological literature abounds with studies dealing with the development and evaluation of methods for the delineation of hydrological regions and for the study of their homogeneity. However, much less attention has been dedicated to the development of new regional estimation methods.

In the present study, canonical correlation analysis (CCA) is used to delineate homogenous regions. In GREHYS (1996b), it was shown that this method produced the best performances in comparison to other ones. Among RFA estimation methods, regression models and index-flood models are commonly used. GREHYS (1996b) showed that their performances are equivalent and are superior to other models. Generally, regression models such as linear regression models (LRM) or log-linear regression models (LLRM) are preferred for their simplicity and rapidity, as well as their performances. LLRM has been used in conjunction with CCA in many studies (Chokmani and Ouarda 2004; Ouarda et al. 2001). Linear models imply that the relations between the dependent variable (hydrologic) and the predictors (physio-meteorological) are linear. This is generally not realistic and can be problematic in some situations such as the effect

65  of the basin size on flood quantiles, where it is documented that small basins behave differently

66  than large ones. The basin hydrologic response is also not linearly related to the slope of the

67  basin, as larger basin slopes (which are often associated to smaller size basins) lead to much more

68  intense flood responses and very extreme specific peak values.

69  The generalized additive models, GAMs (Hastie and Tibshirani 1986) allow to take into account

70  possible nonlinearities which is not possible through linear models or by using simple variable

71  transformations such as log, power or square root. The use of a nonlinear model is justified by the

72  fact that hydrological processes are naturally nonlinear (Kundzewicz and Napiórkowski 1986;

73  Wittenberg 1999). Pandey and Nguyen (1999) compared a number of regional flood quantile

74  estimation methods for the power regression model (equivalently log-linear) and found that

75  nonlinear estimation methods (within the same power model) outperformed the log-linear one.

76  Shu and Ouarda (2007) used an artificial neural network approach, which represents a nonlinear

77  model, and obtained better results than with linear regression methods.

78  GAMs are an extension of the generalized linear models, GLMs (Nelder and Wedderburn 1972).

79  The latter brought flexibility to regression methods by allowing non-normal residuals as well as a

80  general link between predictors and the response variable. In addition, GAMs use non-parametric

81  smooth functions to link the dependant variable to the predictors. Therefore, they are more

82  flexible and can capture more realistically the relation between variables. GAMs have been

83  attracting high attention in statistical developments as well as in practical applications (Hastie and

84  Tibshirani 1986; Kauermann and Opsomer 2003; Marx and Eilers 1998; Morlini 2006;

85  Schindeler et al. 2009; Wood 2003). Recently, additional methodological developments and the

86  availability of implemented computer programs made GAMs increasingly popular in practical

87  research, mainly in the public health and epidemiology fields (Bayentin et al. 2010; Cans and

4

88  Lavergne 1995; Leitte et al. 2009; Rocklöv and Forsberg 2008; Vieira et al. 2009) and in

89  environmental studies (Borchers et al. 1997; Wen et al. 2011; Wood and Augustin 2002). In the

90  field of meteorology, GAMs were used to model the effect of traffic and meteorology on air

91  quality (Bertaccini et al. 2012), to predict air temperature from satellite surface temperature

92  (Kloog et al. 2012), as well as to model mean temperature in mountainous regions (Guan et al.

93  2009). In hydrological modeling, very few studies employed GAMs. For instance, Tisseuil et al.

94  (2010) used GLM and GAM for the statistical downscaling of general circulation model outputs

95  to local-scale river flows. GAMs were used to estimate nonlinear trends in water quality by

96  Morton and Henderson (2008) and in hydrological extreme series modeling by Ramesh and

97  Davison (2002). Recently, Asquith et al. (2013) employed GAMs to develop readily

98  implemented procedures for the estimation of discharge and velocity from selected predictors at

99  ungauged stream locations. However, to the author's best knowledge, GAMs have never been

100  used in the context of RFA of hydrological variables.

101  The objective of the present study is to introduce GAMs in a complete regional model to estimate

102  flood quantiles. A set of 151 basins in the province of Québec, Canada, is considered as case

103  study. It is used in combination with the neighborhood approach using CCA. A cross validation

104  is used to evaluate performances. In previous studies dealing with the estimation of flood

105  quantiles with the same dataset (Chokmani and Ouarda 2004; Kamali Nezhad et al. 2010; Shu

106  and Ouarda 2007), explanatory variables have been selected based on correlation with specific

107  quantiles. In the present study an attempt is made to select optimal variables with a stepwise

108  method. The regional model adopting GAM is compared with a model using LLRM, which is

109  commonly used in RFA. Comparisons are also carried out for models with and without the

110  delineation of homogenous regions with CCA, and also with and without the use of the stepwise

5

111    method for the selection of variables. The latter is important to separate the impacts of using the

112    GAM model and the stepwise variable selection procedure.

113    This paper is organized as follows. Section 2 presents the theoretical background on linear

114    regression models, GAMs and the CCA approach for the delineation of neighborhoods in RFA.

115    The considered dataset as well as the study design are presented in section 3. Section 4 includes

116    the obtained results, while the last section contains the conclusions of the study.

# 2. Theoretical Background

118    In this section, the required statistical tools are briefly presented and their use in RFA is

119    discussed.

## 2.1. Linear regression models

121    Regression analysis is used to find a relationship between a random variable $Y$, called the

122    response variable or dependant variable, and one or several random variables $X$, called the

123    explanatory or predictor variables (or independent variables). Let us define $\mathbf{X}$, a matrix whose

124    columns are $X_1, X_2, \ldots, X_m$, a set of $m$ explanatory variables. The linear regression model is

125    defined by:

126    $$Y = \beta_0 + \sum_{j=1}^{m} \beta_j X_j + \varepsilon \tag{1}$$

127    where $\beta_0$ and $\beta_j$ are unknown parameters and $\varepsilon$ is the error term which is assumed to be

128    normally distributed $N\left(0, \sigma^2\right)$. The model parameters are often estimated by the least squares

129    estimator $\hat{\beta} = \left(\mathbf{X'X}\right)^{-1} \mathbf{X'}Y$.

130 A power product model is generally used to express the relationship between flood quantiles and

131 explanatory variables (Ouarda et al. 2008; Pandey and Nguyen 1999). A log transformation

132 allows expressing this model as follows (log-linear model):

133
$$Y = \log(\beta_0) + \sum_{j=1}^{m} \beta_j \log(X_j) + \varepsilon \qquad (2)$$

134 Note that the log transformation introduces a bias in the prediction since the aim is the estimation

135 of the variable expectation rather than its logarithm (Girard et al. 2004).

## 2.2. Generalized additive models

137 The generalized linear models (GLMs) are a generalization of the well-known ordinary linear

138 model presented previously. They allows for a response distribution other than normal and for a

139 degree of nonlinearity in the model structure (Wood 2006). The GLM can be expressed as

140 follows:

141
$$g(Y) = \beta_0 + \sum_{j=1}^{m} \beta_j X_j + \varepsilon \qquad (3)$$

142 where $g$ is a monotonic link function, and $Y$ could have whatever distribution from the

143 exponential family which includes, for instance, Poisson, Binomial and Normal distributions.

144 For more flexibility, GLMs are themselves extended to GAMs by allowing non-parametric fits of

145 the $X_j$ where the linear forms are replaced by smooth functions $f_j$ (Hastie and Tibshirani 1986;

146 Wood 2006):

147
$$g(Y) = \alpha + \sum_{j=1}^{m} f_j(X_j) + \varepsilon \qquad (4)$$

7

148 GAM has several advantages over linear models. It is more flexible due to the smooth functions $f_j$

149 where there is no need for a transformation to achieve linearity. Hence, it is possible to identify

150 more realistically the effect of each explanatory variable $X_j$ on $Y$.

151 In order to estimate the smooth function $f_j$, a spline is used. A spline is a curve composed of

152 piecewise polynomial functions, joined together at points called knots. A number of spline types

153 have been proposed in the literature, such as cubic splines, P-splines and B-splines. The thin plate

154 regression splines have some advantages such as fast computation, lack of requirement for a

155 choice of knot locations, and optimality in approximation of the smoothing, for more details see

156 (Wood 2003, 2006). In the present study, the latter splines are considered.

157 In general, a smooth function $f_j$ can be defined by a set of $q$ spline basis functions $b_{ji}(x)$ such

158 that:

159
$$f_j(x) = \sum_{i=1}^{q} \beta_{ji} b_{ji}(x) \tag{5}$$

160 where $\beta_{ji}$ represents the smoothing coefficients related to the $j$th function. To avoid overfitting,

161 the estimator $\hat{\beta}$ of $\beta$ is obtained by maximizing the penalized log-likelihood:

162
$$l_p(\beta) = l(\beta) - \frac{1}{2} \sum_{j=1}^{m} \lambda_j \beta^T \mathbf{S}_j \beta \tag{6}$$

163 where $l_p(.)$ is the log-likelihood function, $\lambda_j$ is the smoothing parameter of the $j^{th}$ smooth

164 function $f_j$ and $\mathbf{S}_j$ is a matrix with known coefficients (Wood 2008). The parameter $\lambda_j$ controls

165 the smoothness degree of the curve $f_j$. Its value ranges from 0 to 1, with 0 corresponding to the

166 un-penalised case and 1 to the completely smoothed curve. The optimum value of $\lambda_j$ is a right

167  balance between best fitting and smoothing. The function $l_p(.)$ is maximized by the penalized

168  iteratively reweighted least squares, P-IRLS (Wood 2004). The smoothing parameter $\lambda$ can be

169  selected according to a criterion such as the generalized cross validation, GCV (Wahba 1985),

170  unbiased risk estimator, UBRE (Craven and Wahba 1978) or maximum likelihood (ML).

## 2.4.  CCA Approach in RFA

172  This section briefly presents the CCA approach and its connection to the delineation step of RFA.

173  This method is explained in more details in Ouarda et al. (2001) in the RFA context. Let us

174  define two sets of random variables $\mathbf{X} = \{X_1, X_2, ..., X_r\}$ and $\mathbf{Y} = \{Y_1, Y_2, ..., Y_s\}, s \geq r$. In the

175  present study, the set $\mathbf{X}$ contains basin physiographical and meteorological variables, e.g.

176  drainage area and mean annual precipitation, and $\mathbf{Y}$ contains basin hydrological variables such as

177  flood quantiles. In general, all variables should be standardized and transformed for normality.

178  Mainly, CCA aims to identify the dominant linear modes of covariability between the vectors $\mathbf{X}$

179  and $\mathbf{Y}$, and then make inference about $\mathbf{Y}$ given the vector $\mathbf{X}$.

180  Consider the linear combinations $\mathbf{V}$ and $\mathbf{W}$ of the variables of $\mathbf{X}$ and $\mathbf{Y}$:

181  $$V = a_1 X_1 + a_2 X_2 + \cdots + a_r X_r = \mathbf{a'X} \text{ and } W = b_1 Y_1 + b_2 Y_2 + \cdots + b_s Y_s = \mathbf{b'Y} \qquad (7)$$

182  CCA allows to identify vectors $\mathbf{a}$ and $\mathbf{b}$ for which $\delta_{i,CCA} = corr(V_i, W_i) \quad i = 1, ..., p$ are maximized

183  as well as $corr(W_i, V_j) = 0, \quad i \neq j$ with unit variance.

184  For each basin $B_k$, $k = 1, ..., K$ within a given set of basins $B$, the corresponding values for $\mathbf{V}_i$

185  and $\mathbf{W}_i$ are denoted as $\mathbf{v}_{i,k}$ and $\mathbf{w}_{i,k}$. Let $\mathbf{v}_0$ denote the physio-meteorological canonical score

186  for a target site, associated to the obtained canonical variables. The vector $\mathbf{v}_0$ is known whereas

187  the interest is the estimation of the unknown hydrological canonical score $\mathbf{w}_0$. The

188 approximation can be obtained through $\Lambda \mathbf{v}_0$ such that $\Lambda = diag(\delta_{1,CCA},...,\delta_{p,CCA})$. This leads to the

189 definition of the 100(1-α)% confidence level neighbourhood for $\Lambda \mathbf{v}_0$ containing sites with

190 realizations $w$ of $W$ such that:

191 $$(\mathbf{w} - \Lambda \mathbf{v}_0)^T (I_p - \Lambda^2)^{-1} (\mathbf{w} - \Lambda \mathbf{v}_0) \leq \chi^2_{\alpha,p} \tag{8}$$

192 where $I_p$ is the $p \times p$ identity matrix and $\chi^2_{\alpha,p}$ is such that $P(\chi^2 \leq \chi^2_{\alpha,p}) = 1 - \alpha$. All the aspects

193 related to the CCA in the RFA context are developed in Ouarda et al. (2001).

## 3. Dataset and study design

195 The considered dataset has already been studied in the context of RFA in a number of previous

196 studies (Chebana and Ouarda 2008; Chokmani and Ouarda 2004; Kamali Nezhad et al. 2010; Shu

197 and Ouarda 2007), which provides an opportunity for comparative evaluation of the results. The

198 dataset consists of 151 hydrometric stations located in the southern half of the province of

199 Québec (between 45°N and 55°N), Canada. The hydrological variables are represented by

200 specific flood quantiles (quantiles divided by the basin area), denoted by $QS_{10}$, $QS_{50}$ and $QS_{100}$.

201 The physiographical and meteorological variables, available for each basin, are summarized in

202 Table 1. To avoid redundancy with the previously mentioned studies, details concerning the

203 dataset are not reported here. The reader is referred to the references listed above for information

204 concerning the geographic location of the stations and the scatter plots of the basins in the

205 canonical spaces.

206 The CCA in conjunction with LLRM has been proven to perform well (GREHYS 1996b).

207 However, it is suspected that the more general GAM approach can improve the estimations. In

208 this study, LLRM and GAM are compared as regional estimation models. The fitting of data for

209     GAM is performed with the *R* package *mgcv* (Wood 2004). Smooth parameters, $\lambda_j$ in (6), are

210     estimated with the P-IRLS procedure where the ML score is employed as criterion

211     Homogenous regions are delineated with the CCA method on the basis of the variables *BV,*

212     *PMBV, PLAC, PTMA* and *DJBZ.* These variables are selected on the basis of maximizing

213     correlations with the hydrological variables. Since CCA requires normality, these variables are

214     transformed for the regional analysis as in the previous studies for this region, i.e. a logarithmic

215     transformation for the hydrological variables, PMBV, PTMA and DJBZ, and a square root

216     transformation for PLAC. Figure 3 (not reported here to avoid repetition) in Shu and Ouarda

217     (2007) shows clear nonlinearities in different levels for some variables. This represents a

218     motivation for the use of the GAM model with the present dataset.

219     The design of the present study aims to check the performance of three elements: i) adoption of

220     the CCA delineation step or considering all stations, ii) consideration of the nonlinearity in the

221     regression model through either LLRM or GAM during the regional estimation step and iii) the

222     variable selection method (stepwise or correlation). This leads to 8 combinations denoted as

223     follows:

224     - LLRM|ALL|CORR: LLRM with all stations (no delineation) and with the 5 selected variables

225       (from correlation);

226     - LLRM|ALL|STPW: LLRM with  all stations (no delineation) and variables selected using the

227       stepwise method;

228     - LLRM|CCA|CORR: LLRM with homogeneous regions defined by CCA and with the 5

229       selected variables (from correlation);

230     - LLRM|CCA|STPW: LLRM with homogeneous regions defined by CCA and variables

231       selected using the stepwise method;

232   - GAM|ALL|CORR: GAM with all stations (no delineation) and with the 5 selected variables

233      (from correlation);

234   - GAM|ALL|STPW: GAM with all stations (no delineation) and variables selected using the

235      stepwise method;

236   - GAM|CCA|CORR: GAM with homogeneous regions defined by CCA and with the 5 selected

237      variables (from correlation);

238   - GAM|CCA|STPW: GAM with homogeneous regions defined by CCA and variables selected

239      using the stepwise method.

240   The selection method used in this study is the backward stepwise selection method. It starts with

241   an initial model including all available variables. The regression method is then applied with the

242   current model and the variable with the highest $p$-value is excluded, corresponding to the

243   hypothesis that $\beta_j = 0$ in (5) where $j$ is the $j$th variable. At each step, one variable is excluded.

244   The procedure ends when the $p$-values of all the remaining and significant variables are under a

245   given threshold (5%).

246   Once a model is established, its performance can be evaluated. A jackknife procedure is applied

247   to assess the performance of the models. In this procedure, gauged sites are in turn considered

248   ungauged in order to carry out regional estimation. This procedure allows assessing the following

249   performance criteria:

250   the coefficient of determination $\qquad R^2 = 1 - \dfrac{\sum\limits_{i=1}^{n}(z_i - \hat{z}_i)^2}{\sum\limits_{i=1}^{n}(z_i - \bar{z})^2}$ $\qquad\qquad$ (9)

251   the root mean square error $\qquad \mathrm{RMSE} = \sqrt{\dfrac{1}{n}\sum\limits_{i=1}^{n}(z_i - \hat{z}_i)^2}$ $\qquad\qquad$ (10)

252    the relative root mean square error    $rRMSE = 100\sqrt{\dfrac{1}{n}\sum_{i=1}^{n}\left[(z_i - \hat{z}_i)/z_i\right]^2}$    (11)

253    the mean bias    $BIAS = \dfrac{1}{n}\sum_{i=1}^{n}(z_i - \hat{z}_i)$    (12)

254    the relative mean bias    $rBIAS = 100\dfrac{1}{n}\sum_{i=1}^{n}(z_i - \hat{z}_i)\big/z_i$    (13)

255
256    where $z_i$ and $\hat{z}_i$ are respectively the local (at site) and regional quantile estimates at station $i$, $\overline{z}$

257    is the local mean of the hydrological variable and $n$ is the number of stations.

## 4. Results and discussion

259    The CCA is applied on the dataset with the normalized variables *BV, PMBV, PLAC, PTMA* and

260    *DJBZ*. An optimal value of $\alpha = 0.05$ is obtained with the optimisation procedure of Ouarda et al.

261    (2001). This optimal value is used to delineate the neighborhood at each station. Each regional

262    model, when considering CCA delineation, uses the same neighbourhood for a given station.

263    When CCA is applied to the whole dataset, the two physiographical-meteorological canonical

264    variables are defined as:

265    $V_1 = 0.24\log(BV) - 0.07\log(PMBV) + 0.58\sqrt{PLAC} - 0.33\log(PTMA) - 0.03\log(DJBZ)$    (14)

266    $V_2 = 0.48\log(BV) - 0.25\log(PMBV) - 0.45\sqrt{PLAC} + 1.05\log(PTMA) + 1.10\log(DJBZ)$    (15)

267    and the two hydrological canonical variables are defined as:

268    $W_1 = 2.14\log(QS_{10}) - 13.14\log(QS_{50}) + 10.03\log(QS_{100})$    (16)

269    $W_2 = 6.27\log(QS_{10}) + 2.45\log(QS_{50}) - 8.84\log(QS_{100})$    (17)

270    The non-negligible values of the *BV* coefficient in $V_1$ and $V_2$ confirm the need to include *BV* in

271    the CCA despite the fact that specific hydrological quantiles are used.

272  The stepwise selection of variables is applied for each specific quantile separately and for each

273  regression model LLRM and GAM. Table 2 indicates that the selected variables are the same for

274  a given model and a given selection method, independently of whether CCA is used for

275  homogeneous region delineation. Therefore, the delineation step seems not to have an effect on

276  the selected variables.

277  The results of the application of the jackknife procedure for the performance evaluation of each

278  regional model are presented in Table 3. The best overall performances are obtained with

279  GAM|ALL|STPW and GAM|CCA|STPW with CCA leading to slightly better performances.

280  More precisely and in particular based on the rRMSE, GAM always performs better than LLRM

281  for combinations using the same variable selection approach and the same delineation approach

282  (CCA or ALL).

283  The use of CCA to delineate hydrologically homogeneous regions generally leads to

284  improvements in regional estimation in comparison to the ALL approach for the same selection

285  of variables and the same regression model (GAM or LLRM). However, when GAM is used, the

286  difference between CCA and ALL is not significant especially when using the stepwise

287  procedure for the selection of variables. These results show that the use of GAM makes the

288  procedure more robust and compensates for the advantages of using CCA. This is not the case for

289  LLRM where the use of CCA was shown to lead to significant improvements, see e.g. Chokmani

290  and Ouarda (2004). In other words, this indicates that the use of GAM reduces the importance of

291  delineating the appropriate hydrological neighborhood. A possible interpretation for this result is

292  that the consideration of non-linear formulations in the relation between the explanatory

293  physiographical and meteorological variables on one side and the hydrological variables on the

14

294   other side leads to a reduction of the weight of basins that are not hydrologically similar to the

295   target site.

296   The stepwise method for variable selection improves quantile estimations in comparison to those

297   obtained with the fixed 5 variables. This can be explained by the fact that the correlation-based

298   selection of physiographical and meteorological variables to be used in the model is mainly based

299   on a linear relationship between variables. It must also be noted that the variables are originally

300   selected for CCA purposes (delineation) rather than for regression modeling (estimation).

301   Figures 1 and 2 present the smooth functions $f_j$ of the response variable log(QS100) with the

302   explanatory variables of the fitted models GAM|ALL|CORR and GAM|ALL|STPW respectively.

303   It can be seen that the variables BV, PLAC, LAT and DJBZ show nonlinear relations.

304   Furthermore, the nonlinear relation is more complex for some variables. For instance, the

305   relationship between log(QS100) and DJBZ decreases for small values of DJBZ, increases for

306   midrange values and decreases again for high values of DJBZ. This result reflects the seasonality

307   effect of temperature, through DJBZ, on the flood regime. Another particular example of interest

308   concerns the BV variable. Indeed, it can be seen that small basins have a different effect than

309   moderate basins. This result is important since nonlinearity allows appropriately including the

310   variable BV in the model which eliminates the need to develop specific models for small,

311   moderate or large basins. Variables PMBV, LONG, PLMA and PTMA have approximately

312   linear relations.

313   In the present study, the proposed approach based on GAM is mainly compared with the basic

314   formulation of one of the most popular RFA approaches, which is the log-linear estimation model

315   combined with the CCA delineation approach. The comparison can be extended to other regional

316   flood frequency models, such as the ensemble artificial neural networks-CCA approach (EANN-

317 CCA) (Shu and Ouarda 2007; Shu and Ouarda 2008), the kriging-CCA approach (Chokmani and

318 Ouarda 2004), and the depth-based approach (Chebana and Ouarda 2008; Wazneh et al. 2013a,

319 2013b). In order to widen the comparison, results corresponding to the above approaches are

320 considered since they are already available for the data set considered in the present study. Table

321 4 summarizes the obtained results for all these methods. The results indicate that the GAM-based

322 approach outperforms significantly all the above listed approaches in terms of rRMSE. In terms

323 of rBIAS, the optimal depth-based approach seems to lead to slightly better results, although the

324 difference is not significant.

## 5. Conclusions

326 GAM is commonly used in health, epidemiological and environmental studies. However, it

327 remains unutilized in the field of hydrology, especially in RFA. The multiple linear regression

328 model is the most employed estimation model in RFA mainly because of its simplicity. However,

329 it assumes a log linear relationship between the response variable and the explanatory variables.

330 This assumption is not always true and does not reflect the complexity of the hydrological

331 processes involved. The purpose of the present study is first to introduce GAM in RFA and then

332 to compare its results with those obtained by LLRM. GAM is a flexible model that relaxes the

333 assumptions of the LLRM model (normality and linearity).

334 Results of this study indicate that significantly better estimations are obtained from regional

335 models with GAM. For some explanatory variables, the logarithmic relationship of the response

336 variable with the explanatory variables is not linear. Smooth curves allow for a more realistic

337 understanding of the true relationship between response and explanatory variables. The

338 performance gain is not significant using CCA in conjunction with GAM compared to LLMR.

339 This indicates that GAM is robust and is efficient in RFA even without use of a neighborhood

340    approach. Further efforts are required to generalize this conclusion and to test the benefits of

341    GAM modeling in other hydrological applications.

342    In summary, the use of GAM in RFA is valuable not only in terms of performance but also in

343    terms of other practical aspects (e.g. explicit formulation of the smooth functions, flexibility,

344    reduced number of assumptions, and less subjective choices).

## Acknowledgments

# References

Asquith, W. H., G. R. Herrmann, and T. G. Cleveland, 2013: Generalized Additive Regression Models of Discharge and Mean Velocity Associated with Direct-Runoff Conditions in Texas: Utility of the U.S. Geological Survey Discharge Measurement Database. *Journal of Hydrologic Engineering*, **18,** 1331-1348.

Bayentin, L., S. El Adlouni, T. B. M. J. Ouarda, P. Gosselin, B. Doyon, and F. Chebana, 2010: Spatial variability of climate effects on ischemic heart disease hospitalization rates for the period 1989-2006 in Quebec, Canada. *International Journal of Health Geographics*, **9**.

Bertaccini, P., V. Dukic, and R. Ignaccolo, 2012: Modeling the short-term effect of traffic and meteorology on air pollution in turin with generalized additive models. *Advances in Meteorology*, **2012**.

Blöschl, G., M. Sivapalan, T. Wagener, A. Viglione, and H. Savenije, 2013: *Runoff prediction in ungauged basins. Synthesis across processes, places and scales.* Cambridge University Press.

Borchers, D. L., S. T. Buckland, I. G. Priede, and S. Ahmadi, 1997: Improving the precision of the daily egg production method using generalized additive models. *Canadian Journal of Fisheries and Aquatic Sciences*, **54,** 2727-2742.

Cans, C., and C. Lavergne, 1995: De la régression logistique vers un modèle additif généralisé : un exemple d'application. *Revue de Statistique Appliquée*, **43,** 77-90. .

Chebana, F., and T. B. M. J. Ouarda, 2008: Depth and homogeneity in regional flood frequency analysis. *Water Resour. Res.*, **44,** W11422.

Chokmani, K., and T. B. J. M. Ouarda, 2004: Physiographical space-based kriging for regional flood frequency estimation at ungauged sites. *Water Resour. Res.*, **40,** W12514.

Craven, P., and G. Wahba, 1978: Smoothing noisy data with spline functions. *Numer. Math.*, **31,** 377-403.

Girard, C., T. B. M. J. Ouarda, and B. Bobée, 2004: Study of the bais in the log-linear regional estimation model. *Can. J. Civ. Eng.*, **31,** 361-368.

GREHYS, 1996a: Presentation and review of some methods for regional flood frequency analysis. *Journal of Hydrology*, **186,** 63-84.

——, 1996b: Inter-comparison of regional flood frequency procedures for Canadian rivers. *Journal of Hydrology*, **186,** 85-103.

Guan, B. T., H. W. Hsu, T. H. Wey, and L. S. Tsao, 2009: Modeling monthly mean temperatures for the mountain regions of Taiwan by generalized additive models. *Agricultural and Forest Meteorology*, **149,** 281-290.

Hastie, T., and R. Tibshirani, 1986: Generalized Additive Models. *Statistical Science*, **1,** 297-310.

Kamali Nezhad, M., K. Chokmani, T. Ouarda, M. Barbet, and P. Bruneau, 2010: Regional flood frequency analysis using residual kriging in physiographical space. *Hydrological Processes*, **24,** 2045-2055.

Kauermann, G., and J. D. Opsomer, 2003: Local Likelihood Estimation in Generalized Additive Models, **30,** 317-337.

Kloog, I., A. Chudnovsky, P. Koutrakis, and J. Schwartz, 2012: Temporal and spatial assessments of minimum air temperature using satellite surface temperature measurements in Massachusetts, USA. *Science of the Total Environment*, **432,** 85-92.

Kundzewicz, Z. W., and J. J. Napiórkowski, 1986: Non linear models of dynamic hydrology. *Hydrological Sciences Journal-Journal Des Sciences Hydrologiques*, **31,** 163-185.

396     Leitte, A. M., and Coauthors, 2009: Respiratory health, effects of ambient air pollution and its
397     modification by air humidity in Drobeta-Turnu Severin, Romania. *Science of The Total*
398     *Environment*, **407,** 4004-4011.
399     Marx, B. D., and P. H. C. Eilers, 1998: Direct generalized additive modeling with penalized
400     likelihood. *Computational Statistics & Data Analysis*, **28,** 193-209.
401     Morlini, I., 2006: On Multicollinearity and Concurvity in Some Nonlinear Multivariate Models.
402     *Statistical Methods &amp; Applications*, **15,** 3-26.
403     Morton, R., and B. L. Henderson, 2008: Estimation of nonlinear trends in water quality: An
404     improved approach using generalized additive models. *Water Resour. Res.*, **44**.
405     Nelder, J. A., and R. W. M. Wedderburn, 1972: Generalized Linear Models. *Journal of the Royal*
406     *Statistical Society. Series A (General)*, **135,** 370-384.
407     Ouarda, T. B. M. J., 2013: Regional Hydrological Frequency Analysis. *Encyclopedia of*
408     *Environmetrics*, John Wiley & Sons, Ltd.
409     Ouarda, T. B. M. J., A. St-Hilaire, and B. Bobée, 2008: Synthèse des développements récents en
410     analyse régionale des extrêmes hydrologiques / A review of recent developments in regional
411     frequency analysis of hydrological extremes. *Revue des sciences de l'eau / Journal of Water*
412     *science*, **21,** 219-232.
413     Ouarda, T. B. M. J., C. Girard, G. S. Cavadias, and B. Bobee, 2001: Regional flood frequency
414     estimation with canonical correlation analysis. *Journal of Hydrology*, **254,** 157-173.
415     Pandey, G. R., and V. T. V. Nguyen, 1999: A comparative study of regression based methods in
416     regional flood frequency analysis. *Journal of Hydrology*, **225,** 92-101.
417     Ramesh, N. I., and A. C. Davison, 2002: Local models for exploratory analysis of hydrological
418     extremes. *Journal of Hydrology*, **256,** 106-119.
419     Rocklöv, J., and B. Forsberg, 2008: The effect of temperature on mortality in Stockholm 1998-
420     2003: A study of lag structures and heatwave effects. *Scandinavian Journal of Public Health*, **36,**
421     516-523.
422     Schindeler, S., D. Muscatello, M. Ferson, K. Rogers, P. Grant, and T. Churches, 2009:
423     Evaluation of alternative respiratory syndromes for specific syndromic surveillance of influenza
424     and respiratory syncytial virus: a time series analysis. *BMC Infectious Diseases*, **9,** 190.
425     Shu, C., and T. B. J. M. Ouarda, 2007: Flood frequency analysis at ungauged sites using artificial
426     neural networks in canonical correlation analysis physiographic space. *Water Resour. Res.*, **43,**
427     W07438.
428     Shu, C., and T. B. M. J. Ouarda, 2008: Regional flood frequency analysis at ungauged sites using
429     the adaptive neuro-fuzzy inference system. *Journal of Hydrology*, **349,** 31-43.
430     Tisseuil, C., M. Vrac, S. Lek, and A. J. Wade, 2010: Statistical downscaling of river flows.
431     *Journal of Hydrology*, **385,** 279-291.
432     Vieira, V., T. Webster, J. Weinberg, and A. Aschengrau, 2009: Spatial analysis of bladder,
433     kidney, and pancreatic cancer on upper Cape Cod: an application of generalized additive models
434     to case-control data. *Environmental Health*, **8,** 3.
435     Wahba, G., 1985: A comparison of GCV and GML for choosing the smoothing parameter in the
436     generalized spline smoothing problem. *Ann. Stat.*, **13,** 1378-1402.
437     Wazneh, H., F. Chebana, and T. B. M. J. Ouarda, 2013a: Optimal depth-based regional frequency
438     analysis. *Hydrology and Earth System Sciences*, **17,** 2281-2296.
439     ——, 2013b: Depth-based regional index-flood model. *Water Resour. Res.*, **In presss**.
440     Wen, L., K. Rogers, N. Saintilan, and J. Ling, 2011: The influences of climate and hydrology on
441     population dynamics of waterbirds in the lower Murrumbidgee River floodplains in Southeast

442 Australia: Implications for environmental water management. *Ecological Modelling*, **222,** 154-
443 163.

444 Wittenberg, H., 1999: Baseflow recession and recharge as nonlinear storage processes.
445 *Hydrological Processes*, **13,** 715-726.

446 Wood, S. N., 2003: Thin plate regression splines. *J. R. Stat. Soc. Ser. B-Stat. Methodol.*, **65,** 95-
447 114.

448 ——, 2004: Stable and efficient multiple smoothing parameter estimation for generalized
449 additive models. *Journal of the American Statistical Association*, **99,** 673-686.

450 ——, 2006: *Generalized Additive Models: An Introduction with R.* Chapman and Hall/CRC
451 Press, 392 pp.

452 ——, 2008: Fast stable direct fitting and smoothness selection for generalized additive models. *J.*
453 *R. Stat. Soc. Ser. B-Stat. Methodol.*, **70,** 495-518.

454 Wood, S. N., and N. H. Augustin, 2002: GAMs with integrated model selection using penalized
455 regression splines and applications to environmental modelling. *Ecological Modelling*, **157,** 157-
456 177.

457    Table 1. Descriptive statistics of hydrological variables and physio-meteorological variables.

| Variable | Unit | Notation | Min | Moy | Max | SD |
|---|---|---|---|---|---|---|
| Specific flood of 10 year return period | m³/s.km² | $QS_{10}$ | 0.03 | 0.22 | 0.53 | 0.13 |
| Specific flood of 50 year return period | m³/s.km² | $QS_{50}$ | 0.03 | 0.28 | 0.77 | 0.18 |
| Specific flood of 100 year return period | m³/s.km² | $QS_{100}$ | 0.03 | 0.31 | 0.94 | 0.20 |
| Area of Watershed | km² | $BV$ | 208 | 6 265 | 96 600 | 11 713 |
| Length of main channel | km | $LCP$ | 17 | 157 | 855 | 142 |
| Slope of main channel | m/km | $PCP$ | 0.20 | 3.23 | 23.60 | 3.22 |
| Mean slope of watershed | ° | $PMBV$ | 0.96 | 2.43 | 6.81 | 0.99 |
| Percentage of the basin occupied by forest | % | $PFOR$ | 18.00 | 83.05 | 99.80 | 16.61 |
| Percentage of the basin occupied by lakes | % | $PLAC$ | 0.03 | 7.72 | 47.00 | 7.99 |
| Mean annual total precipitations | mm | $PTMA$ | 646 | 988 | 1 534 | 154 |
| Mean annual liquid precipitations | mm | $PLMA$ | 423 | 717 | 1625 | 176 |
| Mean annual solid precipitations | cm | $PSMA$ | 166 | 302 | 720 | 86 |
| Mean annual liquid precipitations during summer and fall | | $PLME$ | 306 | 455 | 664 | 72 |
| Mean annual degree-days over 0°C | dgr-day | $DJBZ$ | 8 589 | 16 346 | 29 631 | 5 385 |
| Latitude of the station | ° | $LAT$ | 45 | 48 | 54 | 2 |
| Longitude of the station | ° | $LONG$ | 58 | 72 | 79 | 4 |
| Altitude of the station | m | $ALT$ | 5 | 157 | 555 | 125 |

458

459    Table 2. Variables selected for each regional model.

| Regional Models | Quantile | Selected explanatory variables |
|---|---|---|
| [LLRM\|ALL\|STPW], [LLRM\|CCA\| STPW] | $QS_{10}$ | BV, PMBV, PFOR, PLAC, PLMA, DJBZ, LONG |
| | $QS_{50}$ | BV, PMBV, PFOR, PLAC, PLMA, LONG |
| | $QS_{100}$ | BV, PLAC, PLMA, LONG |
| [GAM\|ALL\|STPW], [GAM\|CCA\|STPW] | $QS_{10}$ | BV, PFOR, PLAC, PTMA, LAT, LONG |
| | $QS_{50}$ | BV, PLAC, PLMA, LAT, LONG |
| | $QS_{100}$ | BV, PLAC, PLMA, LAT, LONG |
| [LLRM\|ALL\|CORR], [LLRM\|CCA\|CORR], [GAM\|ALL\|CORR], [GAM\|ALL\|CORR] | $QS_{10}$ | BV, PMBV, PLAC, PTMA, DJBZ |
| | $QS_{50}$ | BV, PMBV, PLAC, PTMA, DJBZ |
| | $QS_{100}$ | BV, PMBV, PLAC, PTMA, DJBZ |

460

461

462  Table 3. Performances obtained with the eight combinations (model, delineation and variable

463  selection).

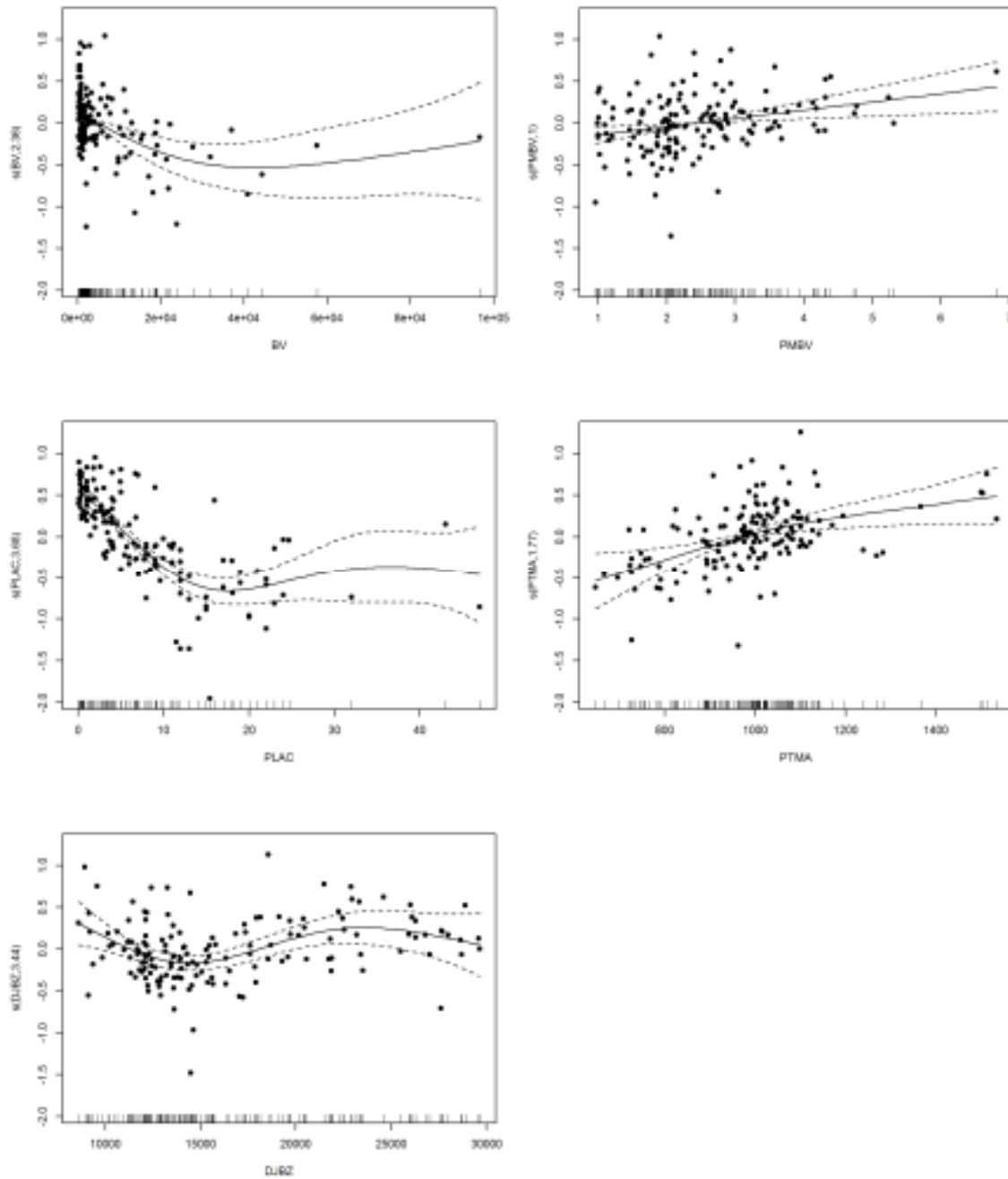| | | LLRM | | | | GAM | | | |
| | | ALL | | CCA | | ALL | | CCA | |
| | Quantiles | CORR | STPW | CORR | STPW | CORR | STPW | CORR | STPW |
|---|---|---|---|---|---|---|---|---|---|
| $R^2$ | $QS_{10}$ | 0.62 | 0.63 | 0.76 | 0.78 | 0.77 | **0.82** | 0.79 | **0.82** |
| | $QS_{50}$ | 0.56 | 0.63 | 0.68 | 0.72 | 0.68 | 0.75 | 0.73 | **0.76** |
| | $QS_{100}$ | 0.53 | 0.53 | 0.64 | 0.65 | 0.65 | **0.72** | 0.69 | 0.67 |
| RMSE | $QS_{10}$ | 0.078 | 0.077 | 0.062 | 0.060 | 0.061 | **0.054** | 0.059 | **0.054** |
| (m3/s.km2) | $QS_{50}$ | 0.117 | 0.108 | 0.100 | 0.094 | 0.099 | 0.088 | 0.092 | **0.087** |
| | $QS_{100}$ | 0.137 | 0.137 | 0.120 | 0.118 | 0.118 | **0.106** | 0.112 | 0.115 |
| rRMSE | $QS_{10}$ | 51.4 | 48.7 | 44.2 | 41.5 | 41.4 | 37.6 | 39.1 | **33.7** |
| (%) | $QS_{50}$ | 56.4 | 55.5 | 48.5 | 48.9 | 47.0 | **41.0** | 43.4 | 43.5 |
| | $QS_{100}$ | 58.9 | 60.0 | 50.7 | 50.9 | 49.3 | 42.1 | 45.6 | **37.0** |
| BIAS | $QS_{10}$ | -0.006 | -0.005 | -0.012 | -0.009 | 0.007 | **0.004** | 0.009 | 0.009 |
| (m3/s.km2) | $QS_{50}$ | -0.010 | -0.011 | -0.021 | -0.015 | 0.013 | 0.009 | 0.018 | **-0.003** |
| | $QS_{100}$ | -0.013 | -0.015 | -0.026 | -0.022 | 0.016 | **0.011** | 0.023 | 0.043 |
| rBIAS | $QS_{10}$ | 7.6 | 7.4 | 5.6 | 5.3 | -5.4 | -5.1 | -4.8 | **-3.5** |
| (%) | $QS_{50}$ | 8.9 | 8.8 | 6.0 | 7.5 | -6.8 | -6.1 | **-4.7** | -11.4 |
| | $QS_{100}$ | 9.6 | 10.0 | 6.3 | 7.7 | -7.6 | -6.5 | -4.9 | **3.4** |

464  Best performances are in bold character for each criterion and quantile

465

466  Table 4. Results of several RFA approaches applied to the same data set considered in this study

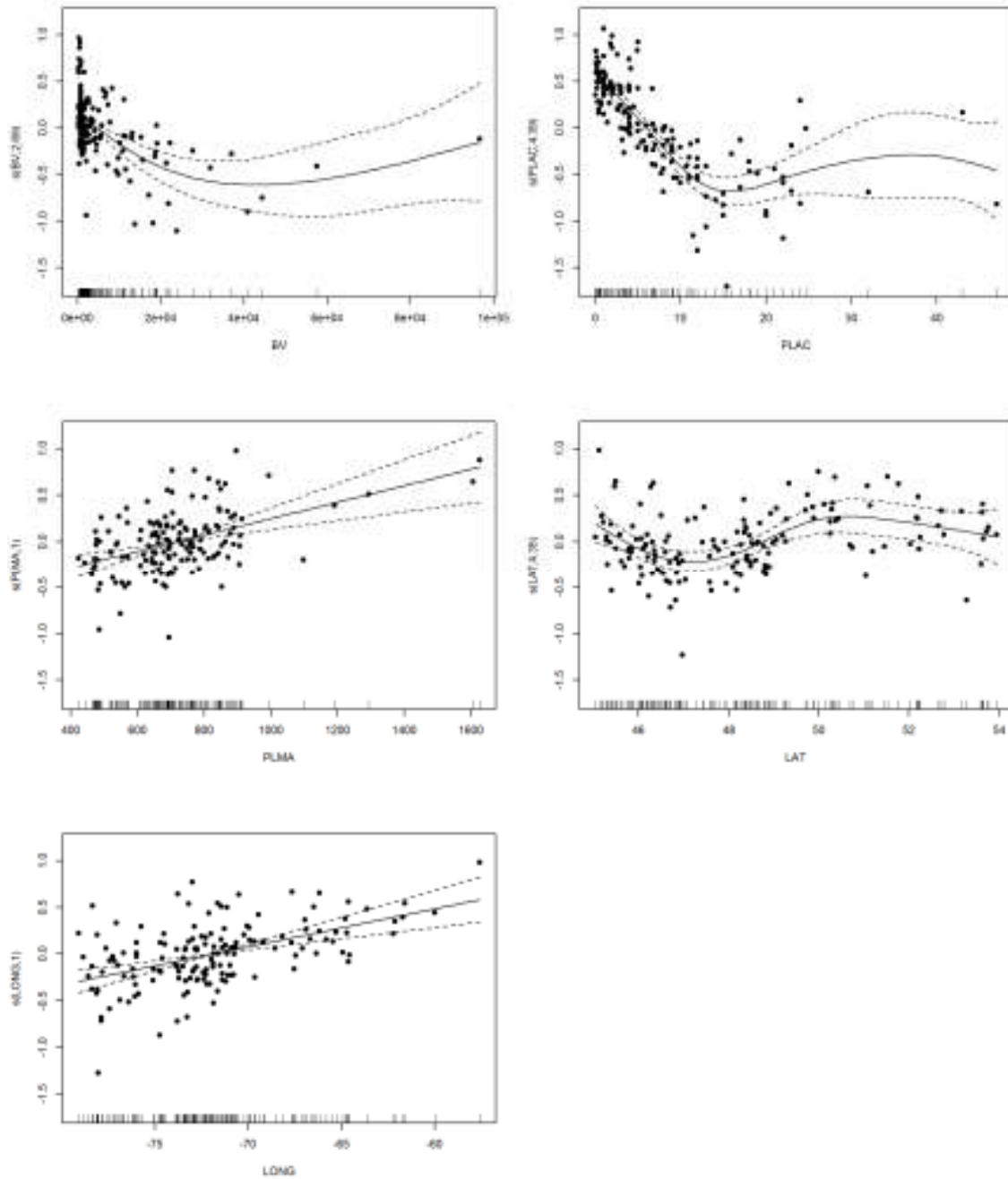| | | $QS_{10}$ | | $QS_{100}$ | |
| Method | References | rBIAS (%) | rRMSE (%) | rBIAS (%) | rRMSE (%) |
|---|---|---|---|---|---|
| Linear regression | Table 3 above | -9 | 55 | -11 | 64 |
| Nonlinear regression | Shu and Ouarda 2008 | -9 | 61 | -12 | 70 |
| Nonlinear regression with regionalization approach | Shu and Ouarda 2008 | -19 | 67 | -24 | 79 |
| Linear regression-CCA | Table 3 above | -7 | 44 | -8 | 52 |
| Kriging in the CCA Physiographical Space | Chokmani and Ouarda 2004 | -20 | 66 | -27 | 86 |
| Kriging in the PCA Physiographical Space | Chokmani and Ouarda 2004 | -16 | 51 | -23 | 70 |
| Adaptive Neuro-Fuzzy Inference Systems | Shu and Ouarda 2008 | -8 | 57 | -14 | 64 |
| Artificial Neural Networks | Shu and Ouarda 2008 | -8 | 53 | -10 | 60 |
| Single Artificial Neural Networks-CCA space | Shu and Ouarda 2007 | -5 | 38 | -4 | 46 |
| Ensemble Artificial Neural Networks | Shu and Ouarda 2007 | -7 | 44 | -10 | 60 |
| Ensemble Artificial Neural Networks -CCA space | Shu and Ouarda 2007 | -5 | 37 | -6 | 45 |
| Optimal depth-based approach | Wazneh et al. 2013a | **-3** | 38 | **-2** | 44 |
| GAM\|CCA\|STPW | Table 3 above | -3.5 | **33.7** | 3.4 | **37** |

Best results are in bold character

22

467

Figure 1. Smooth functions of $QS_{100}$ for the explanatory variables included in the regional model
GAM|ALL|CORR. The dotted lines represent the 95% confidence intervals. The y-axes are
named s(*var*,*edf*) where *var* is the name of the explanatory variable and *edf* is the estimated
degree of freedom of the smooth.

472

Figure 2. Smooth functions of $QS_{100}$ for the explanatory variables included in the regional model GAM|ALL|STPW. The dotted lines represent the 95% confidence intervals. The y-axes are named s(*var*,*edf*) where *var* is the name of the explanatory variable and *edf* is the estimated degree of freedom of the smooth.