

Analysis of Sabine river flow data using semiparametric spline modeling

Soutir Bandyopadhyay^{*}, Arnab Maity[†]

November 4, 2010

Running Head: Analysis of Sabine river flow

ABSTRACT

In this article, a modeling approach for the mean annual flow in different segments of Sabine river, as released in the NHDPlus data in 2007, as a function of five predictor variables is described. Modeling flow is extremely complex and the deterministic flow models are widely used for that purpose. The justification for using these deterministic models comes from the fact that the flow is governed by some explicitly stated physical laws. In contrast, in this article, this complex issue is addressed from a completely statistical point of view. A semiparametric model is proposed to analyze the spatial distribution of the mean annual flow of Sabine river. Semiparametric additive models allow explicit consideration of the linear and nonlinear relations with relevant explanatory variables. We use a conditionally specified Gaussian model for the estimation of the univariate conditional distributions of flow to incorporate auxiliary information and this formulation does not require the target variable to be independent.

Keywords and phrases: Sabine River, semiparametric model, spline.

^{*}Department of Mathematics, Lehigh University, Bethlehem, PA 18015

[†]Department of Statistics, North Carolina State University, Raleigh, NC 27695

1 Introduction

One of the primary challenges for the professionals in water sectors is to meet multiple water demands within the constraint of limited freshwater supply. The necessity to integrate the ecosystem needs is also pronounced in water management. Proper ecosystem management is paramount to protect the ecological processes and biodiversity. It has been noted in some literature that demands for surface water are not expressed freely but rather controlled by water rights specifying the location and type of each allowed usage, the amount to be used and the priority date when the right is established (see for example, <http://www.oregonexplorer.info/willamette/>). Therefore, a good understanding of available water resources is needed for water professionals to achieve a sustainable water system that enriches both this generation and future, while considering the expected future climate and other relevant geographical and hydrological parameters.

As Mylevaganam and Srinivasan (2008) note, contemporary efforts in planning, designing and implementing resource management efforts are now at the catchment scale. The reason to exploit at the catchment scale is to allow management actions to be carried out unhindered until the magnitude of effect reaches to a point where regulation becomes necessary. It has also been mentioned in Ziemer (1994) that generalized regulations are usually not efficient and usually a higher level of regulation results in more streams being overprotected. The closer that the regulations can be tailored to the variables associated with the risk, the less likely that proposed management actions are curtailed needlessly, or, conversely, the less likely that the regulations are inadequate to protect a desired resource. Added to this, the effect of water resources allocation in the upstream of a river basin plays a crucial role in determining the state of the downstream water availability. The spatial connectivity of stream networks often plays a big role to avoid upstream-downstream conflict. Reliability of a catchment is also indirectly linked to the mean annual flow it conveys.

Further, the availability of hydrological data is also critical for water resources planning. Most drainage basins in this world do not have these data because of poorly developed hy-

47 hydrological networks (Oyebande (2001) and Rodda (2001)). It is also not feasible to establish
48 a flow measuring station on every drainage basin (Chiang et al. (2002)) and in addition the
49 sheer sizes of some countries make it impossible to develop adequate hydrological networks
50 and therefore most drainage basins are ungauged (Tucci et al. (1995)). Therefore, the need
51 for hydrological data has greatly increased as water resources which are in some cases scarce
52 have to be shared among competing uses.

53 Therefore, potential to predict water availability, in other words, mean annual flow at
54 a catchment scale considering all the influencing hydrological and geographical parameters
55 is paramount. This also greatly enhances the knowledge on hydrological characteristics of
56 ungauged basins for water resources planning purposes given the prevailing climate and other
57 conditions are of similar nature.

58 **1.1 Dataset**

59 In this section we give a brief description of the NHDPlus data. A more detailed description
60 of the NHDPlus can be found in the website of Center for Research in Water Resources
61 (<http://www.crwr.utexas.edu/gis/gishydro08/ArchHydro/NHDPlus.htm>).

62 According to NHDPlus Users Guide, NHDPlus (Horizon Systems, 2007) is an integrated
63 suite of application-ready geospatial data products, incorporating many of the best features
64 of the National Hydrography Dataset (NHD), the National Elevation Dataset (NED) and
65 the National Watershed Boundary Dataset (WBD) (Holtzschlag, 2009). NHDPlus dataset is
66 distributed for each region as shown in Figure 1. NHDPlus includes a stream network based
67 on the medium resolution NHD (1:100,000 scale), improved networking, feature naming and
68 “value-added attributes” (VAA). NHDPlus also includes elevation-derived catchments which
69 are produced using a drainage enforcement technique. The VAAs include greatly enhanced ca-
70 pabilities for upstream and downstream analysis and modeling. VAA-based routing techniques
71 are used to produce the NHDPlus cumulative drainage areas and land cover, temperature and
72 precipitation distributions. These cumulative attributes are used to estimate mean annual

73 flow and velocity. The objective of the study is to investigate and propose a lattice based
74 mean annual flow predictor for the NHDPlus dataset as released in 2007.

75 **1.2 Study Area and Sabine Basin Hydrology**

76 Detailed description of origin and flow of Sabine river and its hydrology is provided in Compre-
77 hensive Sabine Watershed Management Plan Report (1999), available at the official website of
78 Sabine River Authority of Texas. Below we briefly summarize some of the key points. We refer
79 the interested readers to the original report (located at [http://www.sra.dst.tx.us/srwmp/
80 comprehensive_plan/default.asp](http://www.sra.dst.tx.us/srwmp/comprehensive_plan/default.asp)) for more detailed description of origin, background and
81 hydrology of Sabine river.

82 Sabine River, a river in the southwestern United States, rises in northeastern Texas, flows
83 southeast and south, broadening near its mouth to form Sabine Lake and continues from Port
84 Arthur through Sabine Pass, a dredged navigable channel, to the Gulf of Mexico after a course
85 of 578 mi (930 km). It drains 10,400 sq mi (26,950 sq km) entirely in Texas and the Louisiana
86 Coastal Plain. The Sabine is a flat-water river that pumps about 6.8 million acre-feet into
87 the Gulf and is the single largest volume river in Texas in terms of its discharge. The water
88 has the tannin acid brown color that is common in East Texas rivers and streams.

89 The Sabine River Authority of Texas was created by the Legislature in 1949 as an official
90 agency of the State of Texas. The main purpose of this agency was to act as conservation
91 and reclamation district with responsibilities to control, store, preserve and distribute the
92 waters of the Sabine River and its tributary streams for useful purposes. The boundaries
93 were established by the Act of the Legislature and it comprise all of the area lying within
94 the watershed of the Sabine River and its tributary streams within the State of Texas. The
95 watershed area includes all parts of twenty-one counties. Figure 2(a) shows the total number
96 of catchments available in Texas. We consider the data set of catchments only in the Sabine
97 river basin (Figure 2(b)) containing 5,654 catchments.

98 The hydrology of Sabine river basin is characterized by diverse climatological, topographi-

99 cal and geological features as well as several climatological factors such as temperature, rainfall
100 and humidity. It is known that topography and geologic factors can affect runoff, evapora-
101 tion, sedimentation rates, reservoir storage capacity and water quality and define the river
102 system within the basin. As mentioned in Comprehensive Sabine Watershed Management
103 Plan (1999), the hydrology of the northern region of the basin is significantly different from
104 the southern region. These distinct regions are commonly referred to as the “Upper basin”
105 in the north and the “Lower basin” in the south, the division between the two areas being
106 the headwaters of Toledo Bend Reservoir. The Upper basin is characterized by cool winters,
107 hot summers and seasonal rainfall patterns. The Lower basin has a coastal climate with mild
108 winters, high annual rainfall and moderate to high humidity. However, in this paper, for the
109 modeling purpose we have not considered these two regions separately. We assume that even
110 if the two regions are distinct from the hydrological point of view but the flow at any catch-
111 ment can be modeled in the same way for both the regions. The flow at any catchment only
112 depends on a small number of the neighboring catchments and we assume that the hydrologi-
113 cal properties of a particular catchment is not significantly different from the properties of its
114 neighboring catchments. The geological factors affect the neighboring catchments similarly in
115 each region, and hence affecting the covariates (precipitation, temperature etc) similarly, but
116 not necessarily changing the dependence structure among the neighboring catchments.

117 In this paper we are going to model the mean annual flow of Sabine river based on the NHD-
118 Plus data set released on 2007 in its different catchments based on several relevant variables
119 such as length, stream order, temperature, precipitation and slope. Detailed distributions of
120 these variables based on our data set are shown in tables 1 and 2.

121 The article is organized as follows. In Section 2, we discuss our methodology and implement
122 it to model the data. In Section 3 we discuss the implications of the fitted model.

123 2 Data analysis

124 The goal of the analysis performed here and the features of the data at hand give precise
125 indications about the model to be used. First, it is clear that if an event occurs in a region, it
126 is likely to affect the neighboring regions as well, *i.e.*, the events are spatially dependent. The
127 second aim consists in estimating the flow distribution as a function of explanatory variables
128 because flow can be related to a number of factors, for example, precipitation at the specified
129 catchment, temperature, slope of the region etc. For modeling purpose, the logarithmic trans-
130 formation of the flow values are considered as a function of relevant explanatory variables.
131 By doing this we implement the constraint that the response variable, flow of the river in a
132 catchment, is always a non-negative quantity.

133 Complex functional relations characterizing the flow of the river and their spatially de-
134 pendent structure lead to the adoption of a semiparametric lattice model. In this data we
135 have five covariates, namely, precipitation, temperature, slope, length and stream order of the
136 catchment. For this study, instead of taking the original values of the first four covariates, we
137 take their logarithmic values. The stream order is a variable only taking the values 1 - 11.

138 Considering that the logarithm of river flow is a continuous variable and the model should
139 include the auxiliary variables, it seems natural to resort to Gaussian models. However,
140 the classical Gaussian regression may not be completely adequate since the classical models
141 require the target variables to be independent. Thus some modifications are required to
142 incorporate the spatial dependence of the flow data. More specifically, we can use the well-
143 known *conditionally specified Gaussian models* (see e.g., Cressie (1993)) so that the spatial
144 dependence of the response variable can be taken into account by means of a “conditional
145 specification” model of spatial correlation. In such models one incorporate the fact that an
146 event observed in a certain geographic region depends on what happens in the neighboring
147 regions.

148 A model is said to be of a conditional specification type if the joint distribution function
149 of the units is built on the basis of the univariate conditional distributions. The conditionally

150 specified Gaussian model approach was first proposed by Besag (1974, 1977). However, the
 151 conditionally specified Gaussian model defines a structure of spatial dependence but does not
 152 allow to incorporate auxiliary variables. We have used a semiparametric additive model for
 153 the mean part in the conditionally specified Gaussian model. Linear effects of the length of
 154 the river in the catchment and the stream order of the catchment as well as the nonlinear
 155 effects of precipitation, temperature and slope at a given catchment are included in the model.

156 2.1 The semiparametric lattice model

157 The spatial models on lattices are analogues of time-series autoregressive models. In time
 158 domain the dependence relies upon the unidirectional flow of time where the spatial conditional
 159 approach expresses the dependence of a variable on its nearest neighbor regions. Let \mathbf{Y} be
 160 the $n \times 1$ vector of the dependent variable. The model can now be formalized by explicitly
 161 writing down the conditional distribution of the dependent variable at i -th catchment:

$$162 \quad f(Y_i|\{Y_j : j \neq i\}) = (\sqrt{2\pi}\sigma)^{-1} \exp \left[-\{Y_i - \mu_i - \gamma \sum_{j \in N_i} (Y_j - \mu_j)\}^2 / 2\sigma^2 \right],$$

163 where, $E(Y_i) = \mu_i$ for all $i = 1, \dots, n$, γ is the spatial dependence parameter, σ^2 denotes the
 164 conditional variance of Y_i given $\{Y_j : j \neq i\}$. In the above equation, N_i is the set of neighbors
 165 for the i -th catchment. For more detailed discussion about the conditionally specified Gaussian
 166 models, we refer the readers to Cressie (1993). To find the neighborhood structures we look
 167 at the “ToNode” and “FromNode” of each catchment. “ToNode” is a nationally unique ID
 168 for the to node (with correct coordinate direction, this is the downstream node) endpoint of
 169 the flow line. “FromNode” is the same with the upstream node. A number of catchments
 170 are said to be neighbors if the “ToNode” of the catchments are same as the “FromNode” of
 171 a particular catchment. The dependence parameter γ is estimated from the neighborhood

172 structure and the mean effect is modeled as

$$\begin{aligned}
 173 \quad \mu_i &= \beta_0 + \beta_1 \log(\text{length})_i + \sum_{k=2}^{11} \beta_{2k} I(\text{stream}_i = k) \\
 174 &\quad + f_1(\log(\text{precip})_i) + f_2(\log(\text{temp})_i) + f_3(\log(\text{slope})_i), \\
 175 &= \beta_0 + \beta_1 X_{1i} + \sum_{k=2}^{11} \beta_{2k} I(X_{2i} = k) + f_1(X_{3i}) + f_2(X_{4i}) + f_3(X_{5i}),
 \end{aligned}$$

176 where $f_1(\cdot)$, $f_2(\cdot)$ and $f_3(\cdot)$ are unknown functions describing the effects of precipitation, tem-
 177 perature and slope, respectively. We will use penalized splines (see Wand, 2003) to model
 178 these functions. The penalized regression splines representation of the smooth functions is
 179 given by:

$$180 \quad f_\ell(t_i) = \alpha_{1\ell} t_i + \alpha_{2\ell} t_i^2 + \cdots + \alpha_{p\ell} t_i^p + \sum_{j=1}^K \delta_{j+p,\ell} |t_i - \kappa_{j,\ell}|_+^p$$

181 for $\ell = 1, 2$ and 3 , where each knot $\kappa_{j,\ell}$ is associated to a coefficient $\delta_{j+p,\ell}$ and $x_+ =$
 182 $\max(0, x)$, $x \in \mathbb{R}$, where the coefficients $\delta_{j+p,\ell}$, $j = 1, \dots, K$ are to be penalized (Wand,
 183 2003). The number of knots and their positions can be obtained in an adaptive way although
 184 the sensitivity to this choice is quite low (Ruppert, 2001).

185 **Remark 1.** It is worth mentioning that there may be a situation where one encounters a dry
 186 season with significant occurrences of no precipitation leading to a number of zero values for
 187 the flow. In such a situation, one can use the two-stage model for non-negative variables with
 188 a mass point at zero, as described in Velarde et al.(2004). In this approach, a binary model is
 189 introduced to describe the presence or not of a zero level and then, conditional on observing a
 190 level different of zero, the quantity of the variable will be modeled. The probabilistic descrip-
 191 tion will be a mixture of a discrete and a continuous distribution, generically represented as
 192 $(1 - p) + pf(y|y \neq 0)$, where $p = Pr(Y > 0)$ denotes the probability of Y being greater than
 193 zero. For more detailed description of the zero-inflated model, see Lambert (1992), Ainsworth
 194 (2007).

195 2.2 Fitting the model

196 The penalized pseudo log likelihood is given by

$$197 \quad \mathcal{L} = - \sum_{i=1}^n \log(\sigma^2)/2 - \sum_{i=1}^n \left[- \{Y_i - (X_i\beta + Z_i\delta) - \gamma \sum_{j \in N_i} (Y_j - (X_j\beta + Z_j\delta))\}^2 / (2\sigma^2) \right]$$

$$198 \quad - \delta^T D \delta / 2,$$

199 where X and β denote the unpenalized part of the covariates and corresponding param-
 200 eters in the model; Z and δ denote the penalized part of the covariates associated with
 201 the penalized spline model (Wand, 2003) and corresponding parameters; the matrix $D =$
 202 $diag(\lambda_1 D_1, \lambda_2 D_2, \lambda_3 D_3)$ denotes the penalty matrix associated with δ with D_1, D_2 and D_3
 203 being the penalty matrices corresponding to individual functions $f_1(\cdot), f_2(\cdot)$ and $f_3(\cdot)$, and
 204 λ_1, λ_2 and λ_3 are respective smoothing parameters.

205 We first discuss the model fitting when σ^2 is known. To estimate the parameters, we will
 206 adopt a profiling approach as described in Cressie (1993). For a fixed value of γ , the score
 207 equations for β and δ are

$$208 \quad 0 = \sum_{i=1}^n X_i^{\#T}(\gamma) \{Y_i^{\#}(\gamma) - X_i^{\#}(\gamma)\beta - Z_i^{\#}(\gamma)\delta\},$$

$$209 \quad 0 = \sum_{i=1}^n Z_i^{\#T}(\gamma) \{Y_i^{\#}(\gamma) - X_i^{\#}(\gamma)\beta - Z_i^{\#}(\gamma)\delta\} / \sigma^2 - \lambda D \delta,$$

210 where we define $X_i^{\#}(\gamma) = X_i - \gamma \sum_{j \in N_i} X_j$ and similarly $Y_i^{\#}$ and $Z_i^{\#}$. We can rewrite the
 211 score equation in a matrix form

$$212 \quad 0 = \begin{bmatrix} X^{\#}(\gamma) \\ Z^{\#}(\gamma) \end{bmatrix} V^{-1} \{Y^{\#}(\gamma) - X^{\#}(\gamma)\beta - Z^{\#}(\gamma)\delta\} - D^{\#}(\beta^T, \delta^T)^T,$$

213 where $Y^{\#} = [Y_1^{\#}, \dots, Y_n^{\#}]^T$, $V = diag(\sigma^2, \dots, \sigma^2)$, $X^{\#} = [X_1^{\#T}, \dots, X_n^{\#T}]$ and similarly for
 214 $Z^{\#}$ and $D^{\#} = diag(0, D)$. Defining $W^{\#}(\gamma) = \begin{bmatrix} X^{\#}(\gamma) \\ Z^{\#}(\gamma) \end{bmatrix}$, we have

$$215 \quad \begin{bmatrix} \hat{\beta}(\gamma) \\ \hat{\delta}(\gamma) \end{bmatrix} = [W^{\#}(\gamma)V^{-1}W^{\#}(\gamma)^T + D^{\#}]^{-1}W^{\#}(\gamma)V^{-1}Y^{\#}(\gamma).$$

216 To estimate γ , we first construct the profile likelihood of γ :

$$217 \mathcal{L}_{\text{prof}}(\gamma) = -\{Y^\#(\gamma) - X^\#(\gamma)\hat{\beta}(\gamma) - Z^\#(\gamma)\hat{\delta}(\gamma)\}^T V^{-1} \{Y^\#(\gamma) - X^\#(\gamma)\hat{\beta}(\gamma) - Z^\#(\gamma)\hat{\delta}(\gamma)\}.$$

218 The estimate of γ is then constructed as

$$219 \hat{\gamma}_{\text{prof}} = \operatorname{argmax}_{\gamma} \mathcal{L}_{\text{prof}}(\gamma). \quad (2.1)$$

220 Since γ is an scalar parameter, this maximization problem is easy to solve in any standard
221 software. The final estimates are given by $\hat{\beta}_{\text{prof}} = \hat{\beta}(\hat{\gamma}_{\text{prof}})$ and $\hat{\delta}_{\text{prof}} = \hat{\delta}(\hat{\gamma}_{\text{prof}})$.

222 To estimate the variance, we first fit the model with $V = I$, that is, using working in-
223 dependence assumption. Let the resulting centered residuals be $\hat{\epsilon}_i$, $i = 1, \dots, n$. Then the
224 estimate $\hat{\sigma}^2$ can be obtained by taking the mean of the squares of the centered residuals, *i.e.*,
225 $\hat{\sigma}^2 = n^{-1} \sum_{j=1}^n \hat{\epsilon}_j^2$.

226 **Remark 2.** Instead of using the three-step approach above to estimate the model components,
227 one can use the maximum likelihood estimators where one maximizes the full likelihood with
228 respect to all parameters. This is reasonable from the theoretical point of view. However, we
229 encounter some computational problems and numerical instability issues while maximizing the
230 full likelihood with respect to all parameters possibly due to the fact that the maximization
231 needs to be jointly done on a large number of parameters. In contrast, for the profiling method
232 described above, estimation of β and δ is only a one-step procedure (with closed forms) for
233 each value of γ and is computationally much more efficient and fast. Thereafter, estimation
234 of γ is only a one-dimensional estimation problem. For further detail on this approach, see
235 Cressie (1993).

236 2.2.1 Smoothing parameter selection

237 Most smoothing parameter selection methods do not perform well in the presence of corre-
238 lated errors, as extensive research in the one dimensional case has shown; see Hart (1996)
239 and Opsomer, Wang and Yang (2001) for overviews. We adopt the approach as in Francisco-
240 Fernandez and Opsomer (2005). In that article the authors propose a smoothing parame-

241 ter selection method based on the generalized cross-validation (GCV) criterion (Craven and
 242 Wahba (1978)), suitably adjusted for the presence of spatial correlation. They consider select-
 243 ing the smoothing parameter $\lambda = (\lambda_1, \lambda_2, \lambda_3)$ that minimizes the following “bias-corrected”
 244 GCV criterion

$$245 \quad GCV(\lambda) = n^{-1} \sum_{i=1}^n \left(\frac{Y_i - \hat{Y}_i}{1 - n^{-1} \text{tr}(\mathbf{S}\Sigma)} \right)^2, \quad (2.2)$$

246 where \mathbf{S} is the $n \times n$ smoother matrix such that, $\hat{\mathbf{Y}} = \mathbf{S}\mathbf{Y}$, where $\hat{\mathbf{Y}} = [\hat{Y}_1, \dots, \hat{Y}_n]^T$ and
 247 \mathbf{Y} is defined similarly, and Σ the correlation matrix of the observations which is given by,
 248 $\Sigma = (I - C)^{-1}$, where C is a $n \times n$ matrix with (i, j) -th entry $C_{ij} = \gamma$, if the i -th and the j -th
 249 catchments are neighbors and $C_{ij} = 0$ otherwise. In general, Σ is unknown and we replace it
 250 with its estimate $\hat{\Sigma}$ in (2.2). To find $\hat{\Sigma}$, first note that the correlation matrix depends on the
 251 unknown parameter γ . We can first estimate γ using the profiling approach described earlier
 252 and then we can simply plug in the estimate in the expression for Σ to finally get an estimate
 253 of the correlation matrix. Thereafter, finding the minimizer of this function can be performed
 254 using numerical algorithms and can be easily implemented in standard statistical softwares.

255 **2.2.2 Knot selection**

256 The number of knots suitable to represent the nonlinear effect can be obtained as (Ngo and
 257 Wand (2004)). A reasonable default rule for the knot locations is: $\kappa_j = \{(j + 1)/(j + 2)\}$ th
 258 sample quantile of the unique \mathbf{x}_i 's, for $j = 1, \dots, K$. A simple default choice of K that usually
 259 works well is:

$$260 \quad K = \max \left\{ 5, \min \left(\frac{1}{4} \times \text{number of unique } \mathbf{x}_i\text{'s}, 35 \right) \right\}.$$

261 See Ruppert (2002) for further discussion on default knot specification. The number of knots
 262 and their positions can also be obtained in an adaptive way although the sensitivity to this
 263 choice is quite low (Ruppert (1997)).

264 2.3 Results

265 In order to analyze the spatial distribution of $\log(\text{flow})$ of Sabine river, we model the mean part
266 with a quadratic spline ($p = 2$). The entire analysis is based on the standardized variables.
267 The number of knots suitable to represent the nonlinear effects of $\log(\text{precip})$, $\log(\text{temp})$ and
268 $\log(\text{slope})$ are fixed to 35 and is obtained as Ngo and Wand (2004). The knots are considered to
269 be equidistant. The summary statistics for different variables used in this study are presented
270 in Table 1. Table 2 describes the mean and standard deviations of different variables for each
271 stream order. The penalty parameters for functions of precipitation, temperature and slope
272 are calculated as 0.53, 1.02, 0.06, respectively using the GCV criterion as described earlier.
273 For the fitting purposes, we select a grid of 51 equidistant points in the interval $[q_2(x), q_{98}(x)]$
274 for each $x = \text{precipitation, temperature and slope}$, where $q_2(x)$ and $q_{98}(x)$ denotes the 2nd and
275 98th quantiles of x . We estimate the functions on this grid and center the estimates so that
276 $\sum_{k=1}^{51} \hat{f}_\ell(g_{k,\ell}) = 0$, $\ell = 1, 2, 3$, where $g_{k,\ell}$, $k = 1, \dots, 51$ denotes the grid points corresponding to
277 that function. We plot each of the covariate's effect over the grid along with a 95% point-wise
278 confidence band. The estimated effects are presented in Figures 3 - 6. The slope parameter
279 associated with $\log(\text{length})$ is estimated to be 0.89 with a standard error is 0.01.

280 From the results, it is evident that mean annual flow increases as log length increases. This
281 justifies the spatial pattern of precipitation, draining capacity and their influence on stream
282 flow. From the effect of slope, we see that the significance of steepness of catchment slope on
283 river flow is also pronounced. From Figures 3 and 4, the logarithm values of precipitation and
284 temperature illuminate that these may alone not play a role in determining the stream flow.
285 It is evident that the estimated effects are showing nonlinear patterns within small limits in
286 the vertical axes. From Figure 3, we see a slight upward trend of effect of precipitation on
287 flow. However, the effect becomes flat at the right tail. The connotation is that the type of
288 land use pattern plays a remarkable role in abstracting the precipitation before it eventually
289 contributes to the stream flow. The stream order effects portray that the basin is of mixed
290 nature when it comes to its primary source. There is a possibility of streams being dried in

291 some sections of the basin (higher stream orders).

292 Regarding model diagnostics, Figures 7(a) and (b) represent the location-spread plot, that
293 is, $|\widehat{\Sigma}^{-1/2}(Y_i - \widehat{Y}_i)/\widehat{\sigma}|^{1/2}$ against \widehat{Y}_i , and residual versus predicted values plot for our fitted
294 model. It seems there might be heteroscedasticity present in the model. In view of this, we
295 refit our model using σ_i^2 in place of σ^2 , where the model is,

$$296 \quad \log(\sigma_i^2) = h_1(\log(\text{precipitation}_i)) + h_2(\log(\text{temperature}_i)) + h_3(\log(\text{slope}_i)).$$

297 We fit this additive model using squared centered residuals from an working independence fit
298 of data as response variable and specifying Gaussian likelihood and log-link function. The
299 estimates of σ_i^2 for each individual catchments $i = 1, \dots, n$ are presented in Figure 8 with
300 the horizontal dashed line denoting the estimated variance in the homoscedastic case. We
301 refit our model using this updated variance estimates. The results are very similar to those
302 in Figures 3 - 6 and hence we do not present them. It is interesting to mention that the
303 spread-location and residual-predicted values plots of the updated model (not shown here)
304 still show some signs of heteroscedasticity. We believe this is due to various other physical
305 factors and variables unaccounted in the data. For instance, there are various deterministic
306 relationships/physical models describing the relationship between precipitation, temperature
307 and slope to river flow. We only look at their relationship from a purely statistical point of
308 view and thus do not account for any such physical relationships. This is certainly an area of
309 interest and we hope to pursue this as a future direction of our research.

310 We also investigate several models apart from the above model. Table 3 describes these
311 models and their corresponding AIC values in the homoscedastic case. First column of
312 Table 3 describes the model, for example, the first entry of first column ‘stream.order +
313 length + f(precip)’ corresponds to the model where we include stream order and standard-
314 ized log(length) as linear covariates and standardized log(precipitation) as nonparametrically
315 modeled covariate. The second column of the table provides corresponding AIC values. It
316 is evident that among the models investigated, the model we fit above with all the variables
317 produces least AIC.

318 **3 Discussion**

319 This study defines a lattice based additive model relating catchment properties such as channel
320 slope, precipitation, temperature, length of the stream and stream order to mean annual flow.
321 Though the model is applied to analyze flow of Sabine rive, this type of model have general
322 applicability to other types of such flow network data. Other covariates can also be included
323 in the model if available.

324 There are several important implication of this model. As noted in Arnold et al. (2000),
325 base flow characteristics are essential for efficient development of groundwater resources, and
326 for minimizing pollution risks to connected surface water. Therefore the integrated approach
327 is necessary to enhance the sustainability of both surface water and ground water. It has
328 been also noted in Adane and Foerch (2006) that river systems are often augmented by their
329 base flows during lean seasons. The fitted values of stream order intercepts could be used
330 to form Base flow Index (BFI) providing a systematic way of assessing the proportion of
331 base flow in the total runoff of a catchment. It indicates the influence of soil and geology
332 on river flows and is important for low flow studies. In addition, extreme low flow events
333 are gradually earning more importance in the emerging field of ecohydrology and are more
334 diligently analyzed nowadays (Adane and Foerch, 2006). However, it is often difficult to get
335 recorded data on base flows of rivers because many of the catchments in developing countries
336 remain ungauged. Our work may provide an indication of the underlying baseflow given the
337 climatic and geographical conditions are similar.

338 In addition, our model can be used to estimate the rainfall elasticity. Typically, the
339 rainfall elasticity of stream flow is defined as the proportional change in mean annual stream
340 flow divided by the proportional change in mean annual rainfall (Chiew, 2006). However, this
341 definition assumes that the rate of change in flow relative to change in precipitation is the
342 same for any level of precipitation, that is, the relationship between flow and precipitation is
343 linear. One can use our model to estimate the relationship between flow and precipitation and
344 estimate the rainfall elasticity without being constrained by the linearity assumption and also

345 taking into account the change in flow due to other geographical and climatological factors.

346 The main limitations of this computation are that it does not consider changes in the
347 rainfall frequency and distribution, changes in vegetation characteristics under different cli-
348 matic conditions and potential feedbacks between the atmosphere and the land surface. We
349 also look at the problem from a purely statistical standpoint and do not take into account
350 the different deterministic models relating flow to other variables. One may take into account
351 these deterministic models into the statistical formulation to borrow strength and information
352 from them. This is one of the future directions of our research.

353 **Acknowledgments.** The authors thank S. Mylevaganam of the Spatial Sciences Labora-
354 tory in Texas A&M University for providing with the dataset and a number of constructive
355 suggestions regarding the hydrology of Sabine river.

356 References

357 Adane, A. and Foerch, G. (2006), “Catchment characteristics as predictors of base flow index
358 (BFI) in Wabi-shebele river basin, East Africa”, Proceedings of TROPENTAG 2006, October
359 11-13, 2006, Bonn, Germany.

360 Ainsworth, L. (2007), “Models and methods for spatial data: detecting outliers and handling
361 zero-inflated counts”, PhD Thesis, Simon Fraser University.

362 Arnold, J. G., Muttiah, R. S., Srinivasan, R., Allen, P. M. (2000), “Regional estimation of
363 base flow and groundwater recharge in the Upper Mississippi river basin”, Journal of Hydrol-
364 ogy, 227, 21 - 44.

365 Besag, J. (1974), “Spatial Interaction and the Statistical Analysis of Lattice Systems”, Jour-
366 nal of the Royal Statistical Society, Series B, 36(2), 192-236.

367 Besag, J. (1977), “Efficiency of pseudo likelihood estimation for simple Gaussian fields”,
368 Biometrika, 64, 616-618.

369 Comprehensive Sabine Watershed Management Plan Report (1999), available at http://www.sra.dst.tx.us/srwmp/comprehensive_plan/default.asp.

370

371 Chiang, S. M., Tsay, T. K. and Nix, S. J. (2002), "Hydrologic regionalization of watersheds.

372 I: Methodology", *Journal of Water Resources Planning and Management* 128(1), 3-11.

373 Chiew, H. S. F. (2006), "Estimation of rainfall elasticity of streamflow in Australia", *Hydro-*

374 *logical Sciences Journal*, 51(4), 613-625

375 Craven, P. and Wahba, G. (1978) , " Smoothing noisy data with spline functions: Estimating

376 the correct degree of smoothing by the method of generalized cross-validation ", *Numerische*

377 *Mathematik* 31(4), 377-403.

378 Cressie, N. (2003), "Statistics for spatial data", New York: Wiley.

379 Francisco-Fernandez, M. and Opsomer, J. (2005), "Smoothing Parameter Selection Methods

380 for Nonparametric Regression with Spatially Correlated Errors", *The Canadian Journal of*

381 *Statistics*, 33(2), 279-295.

382 Hart, J. D. (1996), "Some automated methods of smoothing time-dependent data", *Journal*

383 *of Nonparametric Statistics*, 6, 115-142.

384 Hastie, T. and Tibshirani, R. (1990), "Generalized Additive Models", Chapman and Hall.

385 Holschlag, D.J. (2009), "Application guide for AFINCH (analysis of flows in networks of

386 channels) described by NHDPlus," U.S. Geological Survey Scientific Investigations Report

387 2009-5188, 106 p.

388 Horizon Systems. (2007). "National Hydrography Dataset Plus." [http://www.horizon-systems.](http://www.horizon-systems.com/nhdplus/)

389 [com/nhdplus/](http://www.horizon-systems.com/nhdplus/)

390 Lambert, D. (1992), "Zero-inflated Poisson regression, with an application to defects in man-

391 ufacturing", *Technometrics*, 34(1), 1-14.

392 Mylevaganam, S. and Srinivasan, R. (2008), "Effect of grid sizes as subbasins on SWAT model

393 hydrologic and water quality predictions”, Research project (project ID: 2008TX306B). Avail-
394 able at <http://water.usgs.gov/wrri/08grants/progress/2008TX306B.pdf>

395 Ngo, L. and Wand, M. P., (2004), “Smoothing with Mixed model software”, Journal of Sta-
396 tistical software, 9:1.

397 Opsomer, J., Wang, Y. and Yang, Y. (2001), “Nonparametric Regression with Correlated
398 Errors”, Statistical Science, 16(2), 134-153.

399 Oyebande, L. (2001), “ Water problems in Africa-how can sciences help?”, Hydrological Sci-
400 ences Journal 46(6), 947-961.

401 Rodda, J. C. (2001), “Water under pressure”, Hydrological Sciences Journal 46(6), 841-853.

402 Ruppert, D. (1997), “Empirical-bias bandwidths for local polynomial nonparametric regres-
403 sion and density estimation”, Journal of American Statistical Association 92, 1049-1062.

404 Ruppert, D. (2002), “Selecting the Number of Knots for Penalized Splines”, Journal of Com-
405 putational and Graphical Statistics, 11(4), 735-757.

406 Tucci, C., Silveira A. and Sanchez, J. (1995), “Flow regionalization in the upper Paraguay
407 basin, Brazil”, Hydrological Sciences Journal, 40(4), 485-497.

408 Velarde, L. G. C., Migon, H. S., Pareira, B. D. B. (2004), “ Space-time modeling of rainfall
409 data”, Environmetrics, 15, 561–576.

410 Wand, M. P. (2003), “Smoothing and mixed models”, Computational Statistics 18, 223-249.

411 Wood, S. N. (2006), “Generalized Additive Models: An Introduction with R”, Chapman and
412 Hall/CRC Press.

413 Ziemer, R. (1994), “Cumulative effects assessment impact thresholds: myths and realities”,
414 Kennedy, Alan J., ed. Cumulative Effects Assessments in Canada: From Concept to Practice.
415 Alberta Association of Professional Biologists. Edmonton, Alberta, Canada.

variables	min.	Q_1	median	mean	Q_3	max.	st. dev.
log(flow)	-8.294	-1.791	-0.399	-0.820	0.560	3.304	2.009
log(precipitation)	6.912	6.985	7.111	7.111	7.230	7.307	0.119
log(temperature)	5.136	5.168	5.195	5.197	5.217	5.293	0.033
log(slope)	0	0.0002	0.002	0.004	0.005	0.211	0.007
log(length)	-4.510	-0.488	0.561	0.296	1.240	3.582	1.288

Table 1: *The summary statistics (minimum, first quartile (Q_1), median, mean, third quartile (Q_3) and maximum) for different variables.*

stream order	log(flow)	log(precip.)	log(temp.)	log(slope)	log(length)
1	-0.58(1.89)	7.10(0.11)	5.19(0.03)	0.010(0.006)	0.41(1.27)
2	-0.98(2.07)	7.11(0.12)	5.20(0.03)	0.002(0.005)	0.23(1.27)
3	-1.02(2.06)	7.12(0.12)	5.20(0.03)	0.001(0.009)	0.18(1.22)
4	-1.49(2.14)	7.14(0.13)	5.21(0.04)	0.001(0.002)	0.13(1.36)
5	-1.36(2.22)	7.08(0.13)	5.18(0.03)	0.001(0.005)	0.17(1.40)
6	-1.05(2.27)	7.13(0.10)	5.20(0.03)	0.001(0.006)	0.18(1.45)
7	-0.81(2.15)	7.24(0.05)	5.23(0.03)	0.002(0.010)	0.29(1.27)
8	-1.68(1.78)	7.30(0.01)	5.26(0.003)	0.000(0.001)	0.07(1.35)
9	-2.38(2.35)	7.30(0.01)	5.26(0.004)	0.0001(0.0002)	-0.28(1.42)
10	-1.92(1.91)	7.29(0.005)	5.27(0.005)	0.0002(0.0001)	-0.15(1.22)
11	-1.30(1.73)	7.11(0.12)	5.20(0.03)	0.002(0.005)	-0.04(1.02)

Table 2: *Means and standard deviations (in parentheses) of different variables for each stream order.*

Model	AIC
stream.order + length + f(precip)	6200.72
stream.order + length + f(temp)	6192.73
stream.order + length + f(slope)	6211.16
stream.order + length + f(precip) + f(temp)	6169.06
stream.order + length + f(precip) + f(slope)	6178.41
stream.order + length + f(temp) + f(slope)	6169.43
stream.order + length + f(precip) + f(temp) + f(slope)	6152.28

Table 3: *AIC for different models investigated in the data analysis section.*

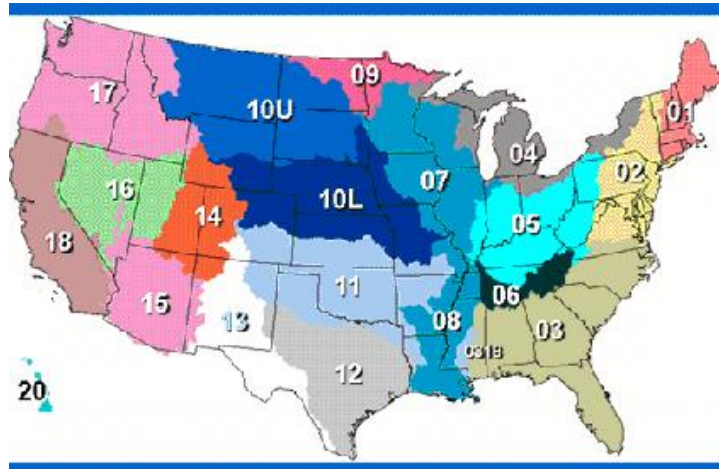


Figure 1: NHDPlus Region

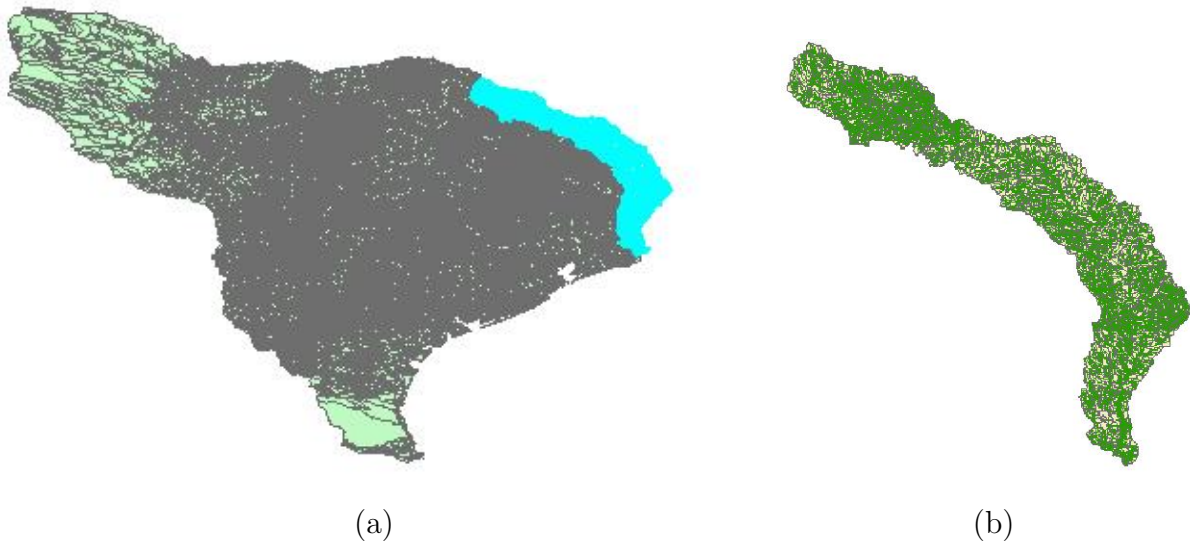


Figure 2: (a) River catchments in Texas, (b) Sabine river basin

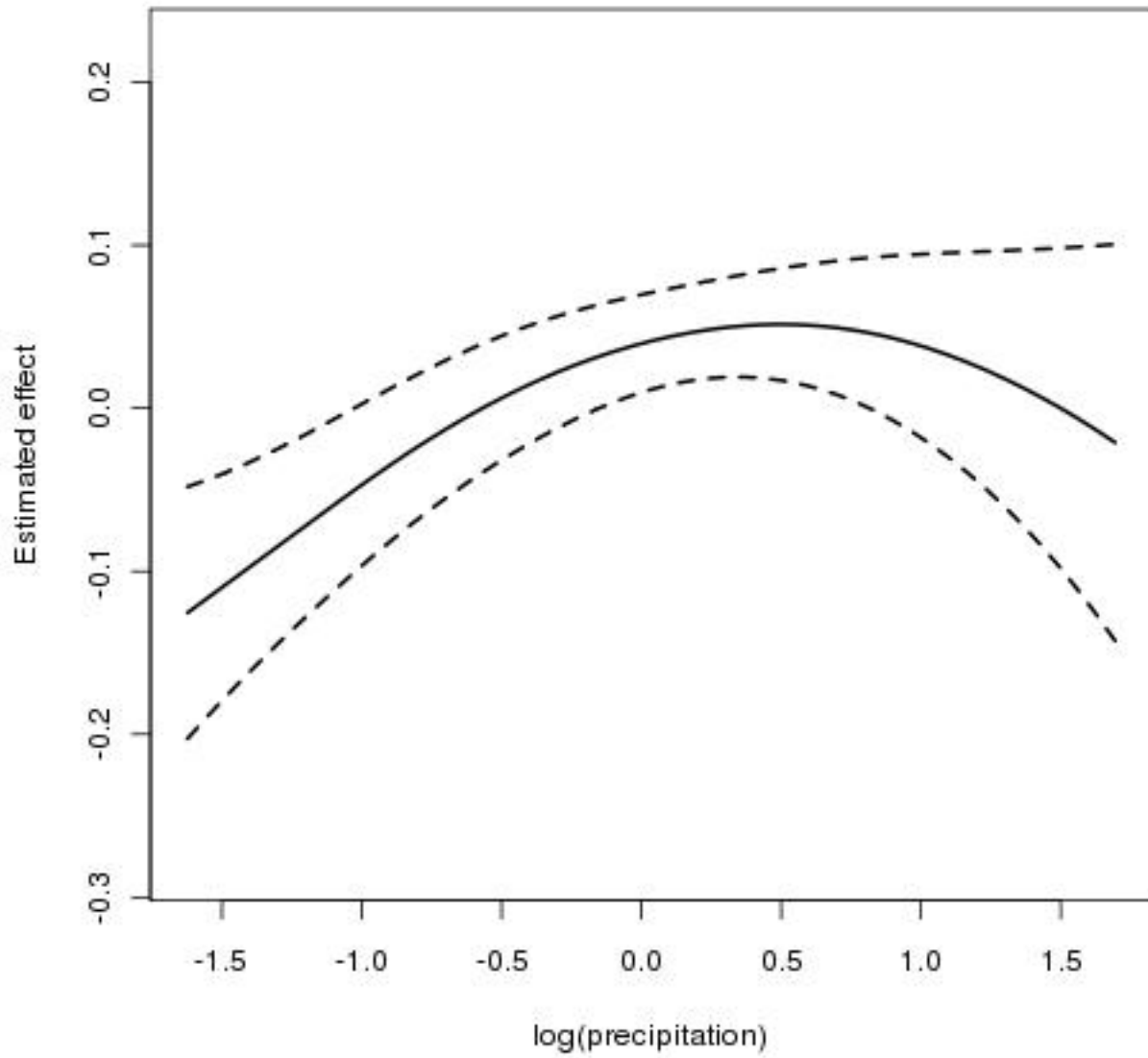


Figure 3: Estimated effect of the logarithm of the precipitation values

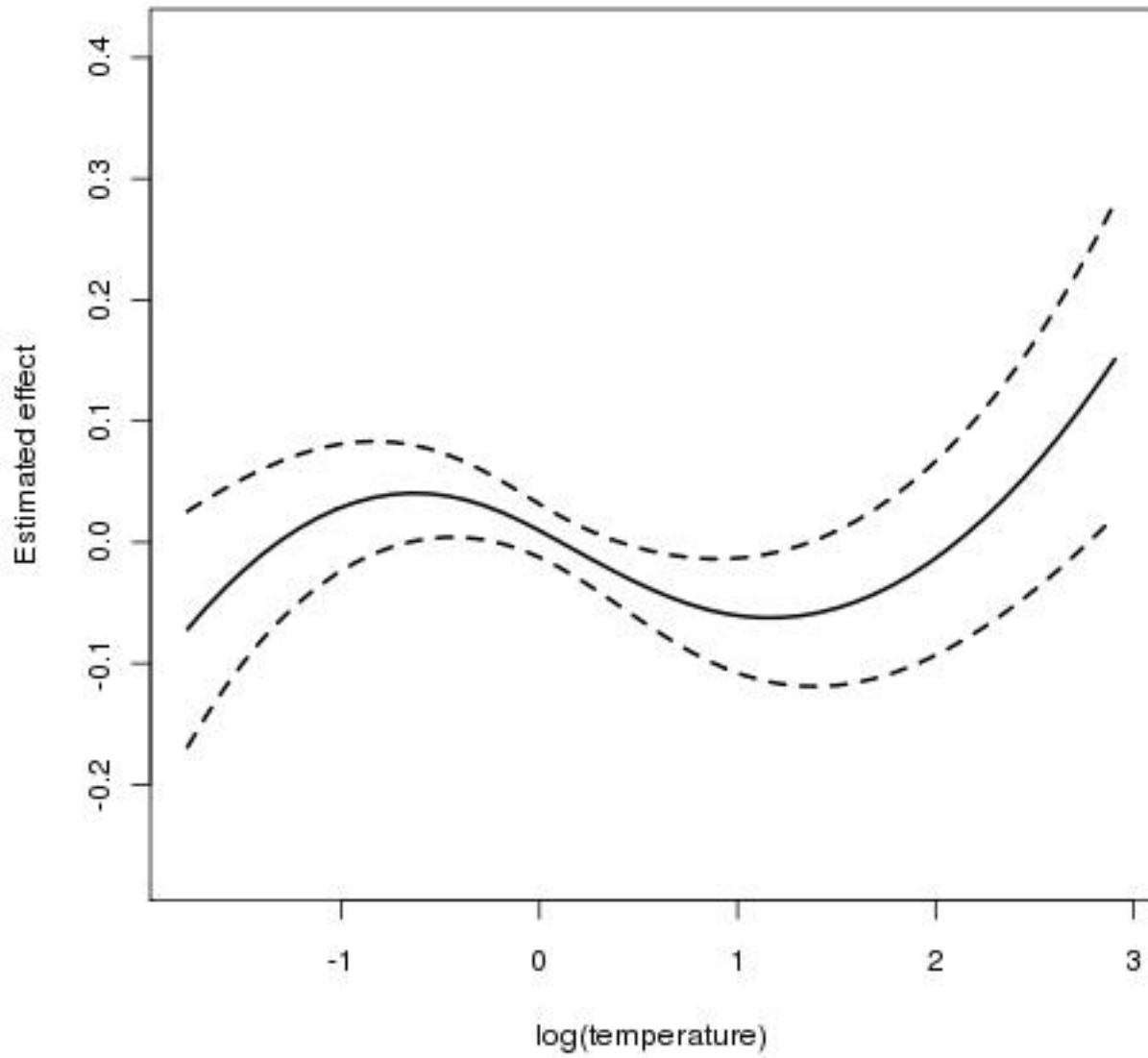


Figure 4: Estimated effect of the logarithm of the temperature values

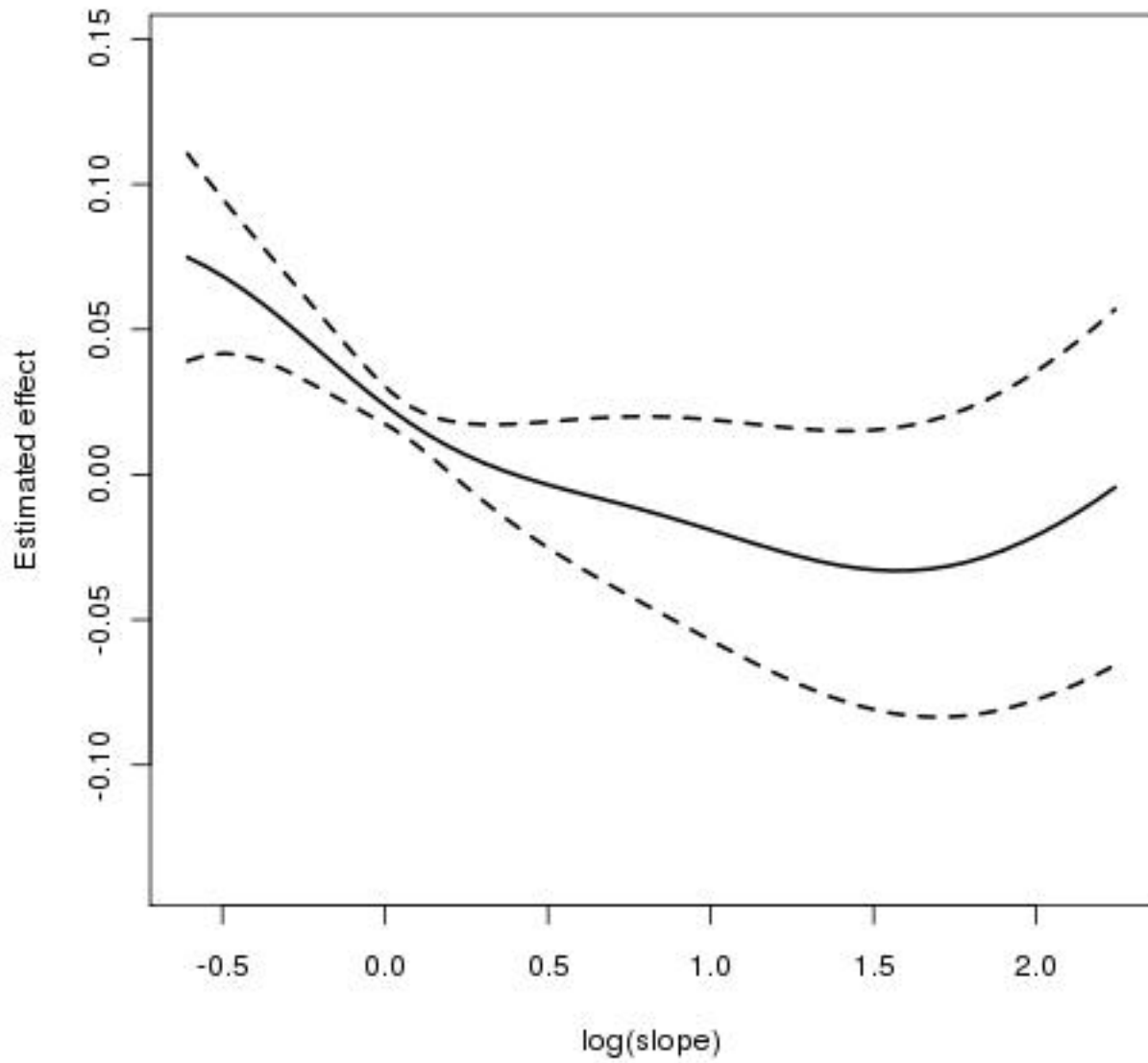


Figure 5: Estimated effect of the logarithm of the slope values

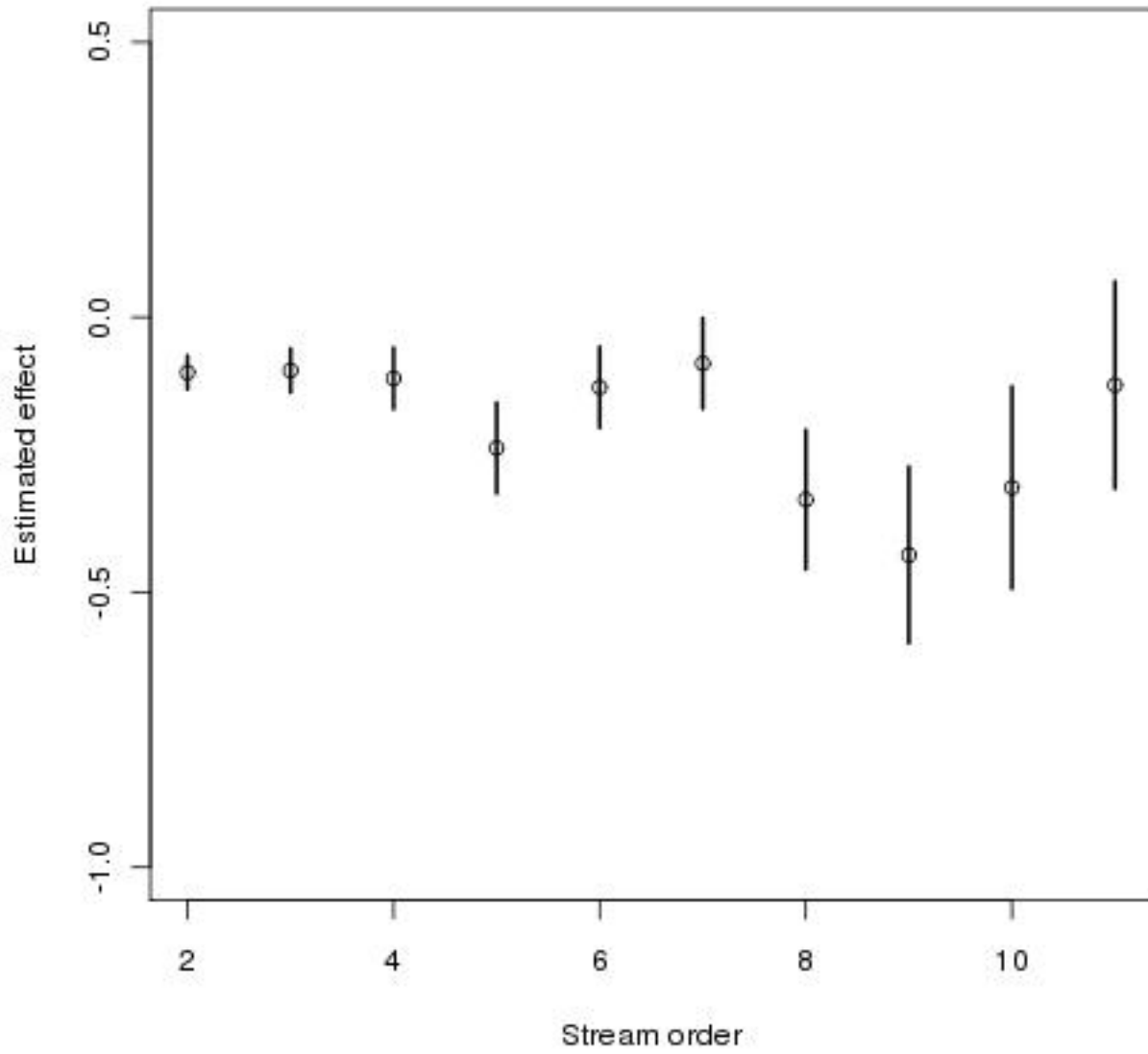
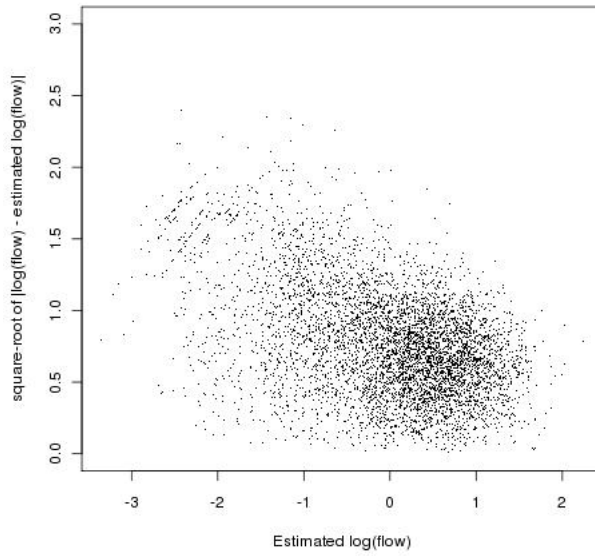
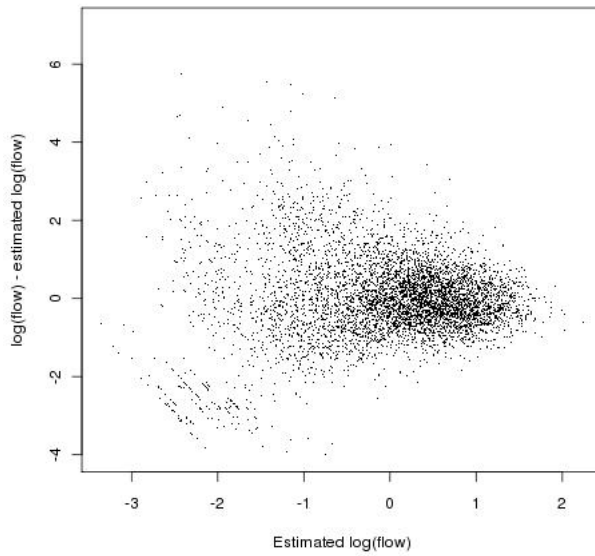


Figure 6: Fitted values of the flow values with the stream order



(a)



(b)

Figure 7: Results from data analysis. Presented are (a) the location-spread plot and (b) plot of the scaled residual versus predicted values.

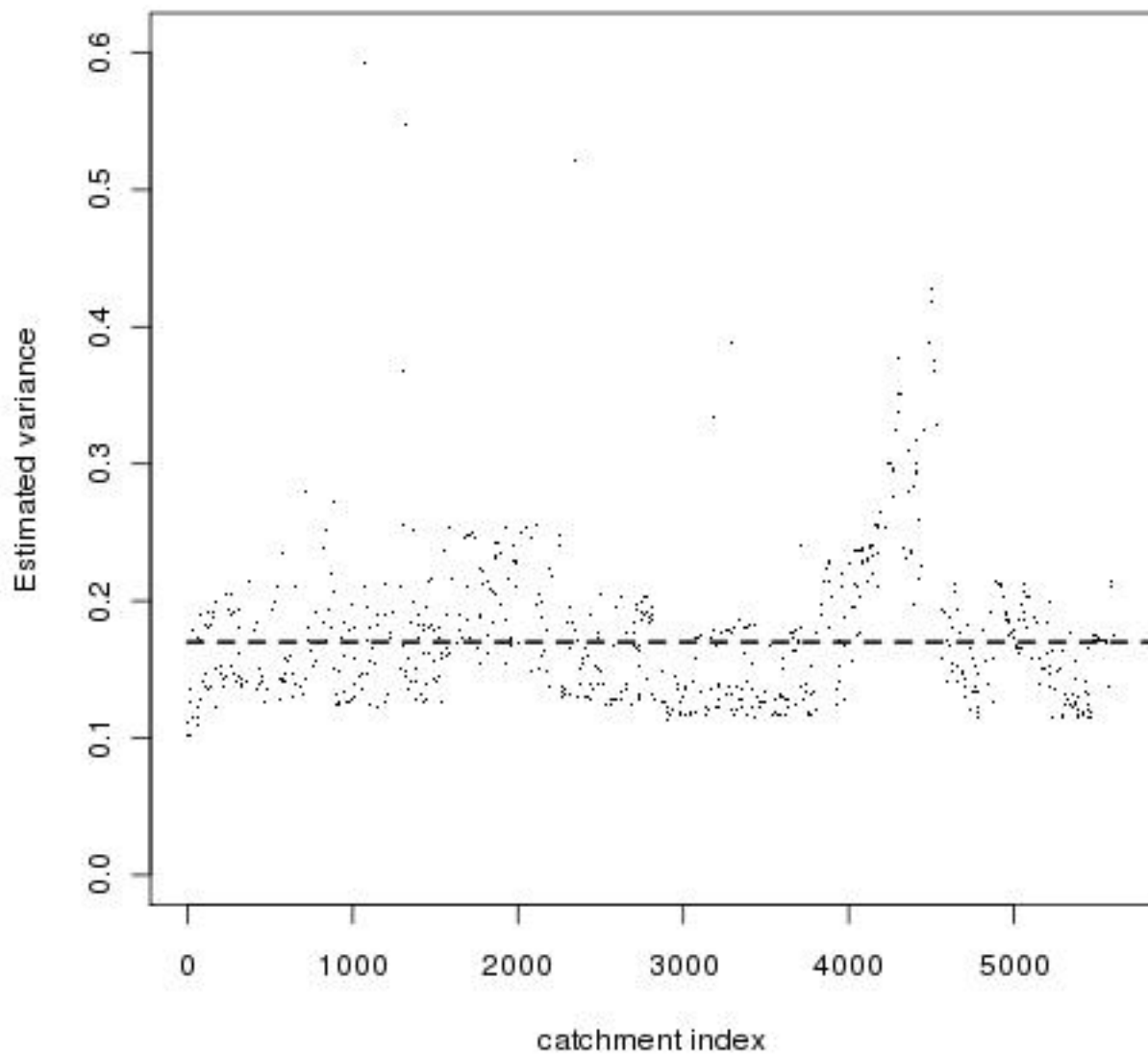


Figure 8: Results from data analysis. Plotted are the estimated variances for each catchment (points) as estimated from the heteroscedastic model and estimated common variance as derived from fitting the homoscedastic model (dashed line).