Accepted Manuscript

Research papers

A nationwide regional flood frequency analysis at ungauged sites using ROI/GLS with copulas and super regions

Martin Durocher, Donald H. Burn, Shabnam Mostofi Zadeh

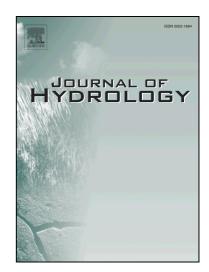
PII: S0022-1694(18)30772-8

DOI: https://doi.org/10.1016/j.jhydrol.2018.10.011

Reference: HYDROL 23175

To appear in: Journal of Hydrology

Received Date: 13 March 2018
Revised Date: 4 October 2018
Accepted Date: 6 October 2018



Please cite this article as: Durocher, M., Burn, D.H., Zadeh, S.M., A nationwide regional flood frequency analysis at ungauged sites using ROI/GLS with copulas and super regions, *Journal of Hydrology* (2018), doi: https://doi.org/10.1016/j.jhydrol.2018.10.011

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

A nationwide regional flood frequency analysis at ungauged sites using ROI/GLS with copulas and super regions

Martin Durocher^{1,*}, Donald H. Burn¹ and Shabnam Mostofi Zadeh¹

1 - Department of Civil and Environmental Engineering, University of Waterloo, 200

University Ave W, Waterloo (ON), Canada, N2L 3G1.

*Corresponding author: mduroche@uwaterloo.ca

Abstract

Region of influence is a common approach to estimate runoff information at ungauged locations. To estimate flood quantiles from annual maximum discharges, the Generalized Least Squares (GLS) framework has been recommended to account for unequal sampling variance and intersite correlation, which requires a proper evaluation of the sampling covariance structure. Since some jurisdictions do not have clear guidelines to perform this evaluation, a general procedure using copulas and a nonparametric intersite correlation model is investigated to estimate sampling covariance structure in situations where no common at-site distribution is imposed or when some paired sites do not have common periods of record. The investigated methodology is applied on 771 sites in Canada. The Normal copula is verified to be an adequate model that better fit paired observations than other types of extreme copulas. A sensitivity analysis is carried out to evaluate the impact of either ignoring, or considering a simpler form of, intersite correlation. Additionally, super regions are defined based on drainage area and mean annual precipitation to improve the calibration of pooling groups across large territories and a wide range of climate conditions. Performance criteria based on cross-validation revealed that using super regions and a combination of geographic distance and similarity between catchment descriptors improves the calibration of the pooling groups by providing more accurate estimates.

keywords: Canada, Floods, Regional frequency analysis, Generalized least squares, Region of influence, Ungauged.

1. Introduction

To estimate flood quantiles at a site of interest where little or no streamflow information is available, hydrologists and practitioners have relied on statistical models to predict runoff properties according to available catchment descriptors. Countries such as the United States (IACWD, 1982), the United Kingdom (Robson and Reed, 1999) and Australia (Rahman et al. 2016) have adopted clear guidelines to standardize how such analysis should be conducted. However, nationwide recommendations are not available in all countries, including Canada. Recently, the project FloodNet (2015) was created as an initiative to coordinate the efforts of experts in various fields, for improving the understanding of floods in Canada. In this line, the present study investigates the problem of performing Regional Frequency Analysis (RFA) to obtain flood quantiles at ungauged sites.

For Quebec and Ontario a previous study compared several combinations of delineation and prediction methods (GREHYS, 1996a, 1996b). One general conclusion was that approaches based on the notion of regions of influence (ROI), where each site is the center of its own pooling group, leads to better results than the delineation of fixed regions. This finding is corroborated by other studies reported in Canada (Burn, 1990; Ribeiro-Corréa et al., 1995; Zrinji and Burn, 1994) and outside Canada (Eng et al., 2007; Merz and Blöschl, 2005; Ouarda et al., 2008). Another conclusion was that the index-flood model (Dalrymple, 1960) performed similarly to the direct estimation of the flood quantiles by regression (Thomas and Benson, 1975). Similar comparisons were repeated in a separate context with similar results (Durocher et al., 2016a; Haddad and Rahman, 2012). Overall, the decision between these two approaches appears to be mostly conceptual. One may argue that the index-flood model provides a more coherent framework by determining the complete regional distribution, while the direct regression approach is more flexible and allows one to easily mix sites with different types of distributions. The present study focuses on the direct regression approach.

Flood frequency analysis is often performed over restricted geographic areas where boundaries are determined by practical considerations. For instance, Canadian studies are generally performed at the provincial level (El-Jabi et al., 2016; Gado and Nguyen, 2016; Sandrock et al., 1992), because water policies fall within these jurisdictions. However, political boundaries are arbitrary from a hydrological perspective and considering a larger database increases the amount of available information. On the other hand, a large country such as Canada can have diverse climatic and flood regimes (Buttle et al., 2016). Therefore, a nationwide analysis will also present additional challenges in the formation of the pooling groups and the calibration of regression models. Many studies have reported a relation between the sample moments and the drainage area or the mean annual precipitation (Basu and Srinivas, 2015; Blöschl and Sivapalan, 1997; Meigh et al., 1997). In a study including sites in Italy, Austria and Slovakia, Salinas et al. (2014a) showed that these two descriptors were proper surrogates for scale control and climate, which contribute to shape the flood generating process. A classification of sites into super regions based on these surrogates could provide a more meaningful solution for characterizing the outcomes of flood frequency analysis at a national level than can be obtained using political regions.

An important decision when using ROI is the choice of a similarity measure that is necessary for creating pooling groups centered around a target site. A distance between catchment descriptors is generally preferred over geographical distance as the topography of nearby catchments can change quickly and lead to distinct hydrological properties. Some studies have compared the usefulness of these two notions of distance and showed that better predictions are generally obtained when considering both of them simultaneously. In the United States, Eng et al.(2007) have successfully used a hierarchical approach where the closest sites in terms of distance between descriptors are selected inside a bounded geographical area. In contrast, Merz and Blöschl (2005) have found in an Austrian case study that superior predictive power was found when using a spatial interpolation technique inside regions where a minimum of similarity among the catchments was imposed. One possible explanation for these findings is that the set of available catchment descriptors is not sufficient to fully characterize the flood generating

process and hence, geographical location represents a surrogate for missing descriptors that evolves smoothly in space. In particular, for a study covering a large territory, such as Canada, the notion of geographical distance is likely to be related to climate characteristics.

When using a regression model to estimate flood quantiles at ungauged sites the response variable is not directly observed, but rather is estimated with different levels of uncertainty. Possible factors that contribute in creating variations in the sampling variance include record lengths and observation scales (Tasker, 1980). Additionally, large atmospheric systems that generate intense rainfall or build vast snowpacks can simultaneously affect many sites. Omitting the impact of intersite correlation in RFA does not create bias but does underestimate the model uncertainty (Bayazit and Önöz, 2004; Hosking and Wallis, 1988) and reduces the power of homogeneity tests (Castellarin et al., 2008). Consequently, ignoring the spatial component of the sampling can have important consequences on the decisions taken on the basis of a selected model (Douglas et al., 2000; Madsen and Rosbjerg, 1997). Generalized Least Squares (GLS) represents a natural approach to estimate parameters of a regression model that accounts for intersite correlation and unequal variances. The model described in Tasker and Stedinger (1989) and considered in several subsequent works (Haddad and Rahman, 2012; Kjeldsen and Jones, 2007; Madsen et al., 2002; Robson and Reed, 1999), separates the total error into sampling and model components. In addition to a better characterization of the source of variability, the approach was also shown to increase predictive accuracies (Reis et al., 2005; Stedinger and Tasker, 1985; Vogel and Kroll, 1990).

Although GLS is considered as a good practice based on several studies, it has not been largely employed in Canada. One reason to explain this low utilization may be the extra step required to evaluate the sampling covariance structure. Across Canada, different regions of dominant nival, pluvial and mixed flood regimes can be found (Burn et al., 2016). Therefore, it is reasonable to expect the existence of complex patterns of spatial dependencies. However, Kroll and Stedinger (1998) indicated that using a smoothed version of the intersite correlation structure has a relatively small impact on the predicted

variability. Consequently, several studies dealing with intersite correlations have preferred to accept some degree of approximation in the sampling covariance matrix by using either Taylor approximation or by assuming simpler spatial correlation models (Kjeldsen and Jones, 2004; Tasker and Stedinger, 1989). Moreover the evaluation of the sampling may depend on the type of at-site distributions selected (Griffis and Stedinger, 2007), which complicates the evaluation of the sampling covariance when no unique distribution is imposed.

The copula framework has gained popularity for describing non-traditional forms of spatial dependence (Bárdossy, 2006; Gräler and Pebesma, 2011). Common assumptions in RFA are to consider the spatial structure of a multivariate Normal distribution, which in terms of copula is equivalent to considering a Normal copula (Durocher et al., 2016b; Renard, 2011). The Normal copula has also been considered in RFA of extreme rainfall, but some studies have considered models based on the generalization of the extreme value theory to spatial extremes, called max-stable processes, to provide more realistic representation of the spatial dependence (Neves and Gomes, 2011; Shang et al., 2011; Westra and Sisson, 2011). In the copula framework, max-stable processes correspond to multivariate Husler-Reiss copula, which has non-negligible probabilities that two extreme events occur jointly, which is not the case of the Normal copula (Joe, 2015). Choosing a copula in RFA poses a similar dilemma to adopting, or not, the GEV distribution in at-site frequency analysis as it is motivated by asymptotic arguments that assumes that the maximums are taken over an infinite number of events. However, this assumption is not realistic in cold regions, because the annual maximum discharge is generally the result of one event, the spring snowmelt. The study of Wang et al. (2014) compared the performance of models based on max-stable processes and regional L-moment algorithm (Hosking and Wallis, 1997) for extreme precipitation in Switzerland. When correctly specified, the max-stable model improved the model fitting and the predicting capability, but when misspecified it was shown to lead to non-negligible bias, which underlines the importance of correctly choosing the copula when modeling spatial extremes.

The present study investigates the ROI/GLS framework when applied to a nationwide database that covers vast territories and includes a large spectrum of climate conditions. Different models for estimating the sampling covariance matrix are examined in light of the copula framework. Among them a nonparametric model is proposed, which does not assume specific at-site distributions or estimation methods and remains valid when few paired observations do not share common periods of record. In Canada, as far as the authors know, there is no study that validates the choice of a proper copula for intersite correlation between floods. One objective is to find such copula and to measure its impact on the estimation of flood quantiles. The notion of super regions is also introduced in the context of ungauged analysis to help with calibrating and understanding the outcomes of ROI/GLS regression models in terms of scale control and climate. Additionally, the combination of geographical distance and distance between descriptors is explored to find the right balance between them in the formation of pooling groups.

The present document is organized as follows. First, section 2 will describe the proposed ROI/GLS methodology and its components. In section 3, the methods are applied on a large dataset of gauged sites across Canada where different calibration of the regression models are examined in terms of quality of the fitting and predictive performance. Finally, further discussions and conclusions are provided in section 4.

2. Methodology

The present methodology has three main components. First, an at-site frequency analysis of the gauged sites is conducted to provide at-site estimates of flood quantiles. Second, an uncertainty analysis of the at-site flood quantiles with return periods 10 and 100 years (denoted Q10 and Q100) is carried out including the choice of a copula, the modeling of the intersite correlation and the estimation of a sampling covariance matrix by Monte-Carlo simulation. Third, relationships between the flood quantiles and catchment descriptors are characterized by the ROI/GLS approach and examined within super regions. The techniques included in these three components are described in more detail below.

2.1 At-site frequency analysis from annual maximums

At-site frequency analysis based on annual maximums has a long tradition in hydrology and remains one of the most common approaches for quantile estimation (Bezak et al., 2014). Theoretical arguments suggest the utilization of Generalized Extreme Values (GEV) distribution, which arises as the limiting distribution of blocks of maximums. However practical considerations often lead to the consideration of other types of distributions. Identification of a best distribution remains an active debate (Salas et al., 2013). Guidelines in the United States are to adopt the log-Pearson III (IACWD, 1982), while in the United Kingdom the Generalized Logistic distribution is recommended (Robson and Reed, 1999). In Europe, Salinas et al. (2014b) showed that even though the GEV often represents a good fit, it cannot accurately describe the complete spectrum of hydrological diversities. Therefore, the present study will prioritize the GEV distribution but will also consider alternative distributions when warranted; alternatives considered are the Gumbel, Gamma, Pearson III, Normal, Generalized Normal and Generalized Logistic.

A classical approach for estimating the parameters of a statistical distribution is to maximize the data likelihood (Coles, 2001). In the presence of a small sample size and heavy tails the maximum likelihood estimator (ML) can sometimes have erratic behavior (Smith, 1985). Therefore, an estimator based on the probability weighted moments, or equivalently the L-moments, is preferred (Hosking, 1990). In the present study, the selection of the at-site distribution is established from a procedure that is guided by the Akaike Information Criterion AIC = 2 k - 2l, where k represents the number of parameters and l the log likelihood (see, for instance, Di Baldassarre et al. (2009)). After identification of the best distribution in terms of AIC, the AIC for the best distribution is compared to the AIC for the GEV distribution. If the difference of AIC is less than one, the two distributions are assumed to fit the data equally well and the GEV is selected. A comparison study of several statistical distributions was performed by Zhang et al. (2018) and showed that the GEV is generally the best choice for Canadian Rivers. Therefore, the criterion of a AIC difference lesser than one for judging the equivalence with the GEV is based on practical

considerations and aims at selecting alternative distributions only when this choice is not supported by the data. Note that a threshold of one in this context is not very restrictive. In comparison, the addition of an extra parameter (for example passing from Gumbel to GEV) increases the AIC by two and some authors even suggest a difference greater than 4 to be a meaningful difference (Burnham and Anderson, 2002).

The uncertainty analysis of the flood quantile estimates based on different types of at-site distributions is not straightforward. In some analyses, the type of distribution is imposed, which allows the development of simple approximate formulas (Bayazit and Önöz, 2004; Griffis and Stedinger, 2007; Kjeldsen and Jones, 2004). This is, however, not the case of the present study. Instead, Monte-Carlo simulations are considered to approximate the sampling covariance matrix. The at-site frequency analysis is repeated 1000 times using samples generated by parametric bootstraps. The covariance matrix is then computed empirically from these samples. Although the method is computationally intensive, it is relatively simple to implement. The simulations require the specification of a statistical model that accounts for the at-site distributions and the intersite correlation.

2.2 Intersite correlation in the copula framework

A copula $C:[0,1]^d \to [0,1]$ is a multivariate distribution with uniform marginal distribution that respects some basic properties (Nelsen, 2006). The fundamental advantage of the copula approach is the separation of the dependence structure from marginal distributions. Commonly used multivariate distributions have the following copula representation

1)
$$G(x_1, K, x_d) = C[F_1(x_1), K, F_d(x_d)],$$

where the F_i are the marginal distributions evaluated at a vector (x_1, K, x_d) (Salvadori et al., 2007). For simulating a vector from G, one can obtain first a uniform vector from copula C and then transform the output to the desired marginal distribution using F_i^{-1} . This strategy is used to simulate one year of annual

maximums where the marginal distribution are the at-site distributions and the intersite correlation is described by a multivariate copula.

The Normal copula and t-copula characterize respectively the dependence of a multivariate Normal and Student distribution. As in classical multivariate theory, the Normal copula is the limit case of a t-copula when the degree of freedom is converging to infinity (Demarta and McNeil, 2005). Formally, for a uniform vector (u_1, K, u_d) the t-copula is defined by

2)
$$C_{\Sigma,\upsilon}(u_1, \mathbf{K}, u_d) = t_{\upsilon,d} \left[t_{\upsilon,1}^{-1}(u_1), \mathbf{K}, t_{\upsilon,1}^{-1}(u_d); \Sigma \right],$$

where $t_{\upsilon,d}$ is the distribution function of a standard Student distribution of dimension d having degrees of freedom υ and correlation matrix Σ . An important difference between the two copulas is a property called upper tail dependence that is defined between two random variables $X_i \sim F_i$ with i=1,2 as

3)
$$\lambda_{up} = \lim_{q \to 0^+} P \left[X_2 > F_2^{-1}(q) \mid X_1 > F_1^{-1}(q) \right].$$

For a t-copula, this property is controlled by the degrees of freedom where low values imply higher probabilities that two extreme events occur jointly. At the opposite extreme, for the Normal copula $\lambda_{up} = 0$, which means that two extreme values never occur together.

The validity of the choice of a copula C can be assessed by a goodness-of-fit test. For bivariate copulas, extensive Monte-Carlo simulations showed that the test based on Cramer Von Mises statistics generally leads to superior or competitive power in comparison to other alternatives (Berg, 2009; Genest et al., 2009). However, such tests require many observations to discriminate between similar copulas. The idea of assessing the quality of spatial models by examining paired observations inside a group of similar lag distances using copula was first suggested by Bárdossy (2006). Although this approach does not provide a formal test, rejection of the null hypothesis for some lag distance provides evidence of model misspecification. That strategy was later formalized by Durocher and Quessy (2017), who showed from simulation studies that reasonable power can be expected in realistic settings.

For t-copulas, the coefficients of correlation for Σ can be estimated by a moment-based estimator and once known, the degrees of freedom are estimated by maximum likelihood (Lindskog et al., 2003). Spearman rank correlation coefficient, or simply Spearman's rho ρ is defined as the correlation between the ranks of two variables. For the t-copula, there is a one-to-one relation between the Spearman's rho and one coefficient θ of Σ :

4)
$$\theta = 2\sin\left(\frac{\pi}{6}\rho\right)$$
.

In an ideal situation, all pairs of sites will have enough years of common record to ensure a reliable estimate of ρ and Σ . However, for different practical reasons, paired observations are recorded over different periods of time and so may prevent or lead to unreliable estimates of ρ . In that case, a spatial correlation model is necessary to have estimates at every pair of sites (Schabenberger and Gotway, 2004). A common choice of spatial correlation model is the power exponential model (POW) where the correlation function s in respect of distance h is

5)
$$s(h) = \begin{cases} (1-\tau)\exp\left[-3\left(\frac{h}{\alpha}\right)^{\gamma}\right] & h \neq 0\\ 1 & h = 0 \end{cases}$$

where $\alpha > 0$ controls the strength of the correlation, $0 \le \tau \le 1$ is a nugget effect and $0 < \gamma \le 2$ is a smoothing parameter. This correlation model is attractive for its simplicity, but may not adequately fit all complex situations. For this reason, a nonparametric model that is more flexible is also considered. Let h_{ij} be the distance (km) between a pair of sites i and j and define the average drainage area A_{ij} (km²), longitude x_{ij} and latitude y_{ij} for the pair (i, j). The nonparametric model characterizing the intersite correlation is

6)
$$g(\rho_{ii}) = f_h(h_{ii}) + f_A(A_{ii}) + f_{xy}(x_{ii}, y_{ii}) + e_{ii}$$

where g is the Fisher z-transformation

7)
$$g(\rho_{ij}) = \frac{1}{2} \log \left(\frac{1 + \rho_{ij}}{1 - \rho_{ij}} \right);$$

 f_{lb} f_A and f_{xy} are smooth continuous functions and e_{ij} is an error term. A model of this form falls under the umbrella of generalized additive model where an in-depth description is provided for instance in Wood (2006). The Fisher transformation is used to transform the empirical Spearman's rho ρ_{ij} between -1 and 1 to a near normal distribution. The smooth functions f_h , f_A and f_{xy} are thin plate regression splines, which are well suited for modeling spatial covariates (Wood, 2003). To avoid overfitting the estimation process is regularized by penalized least squares (Green and Silverman, 1993). The evolution of the Spearman's rho in respect of the distance is described by f_h , which plays a similar role to the POW model. The other components, f_A and f_{xy} , characterize other components of the intersite correlation. In particular, $f_{xy}(x_{ij}, y_{ij})$ adjusts the intersite correlation on the basis of the paired locations.

Estimation of Σ is deduced using equation (4) and fitted ρ_{ij} . However, the matrix Σ derived directly from the nonparametric model will not in general be positive definite, which leads to numerical problems and so the algorithm of Higham (2002) is used to find the nearest matrix that respects this condition.

2.3 Regression models using GLS

Flood quantiles are modeled by a multiple regression model at the logarithm scale for sites found inside a pooling group with regression equation

8)
$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\omega}$$

where $\mathbf{X} \in \mathbb{R}^p$ is a design matrix of relevant catchment descriptors and $\boldsymbol{\omega}$ is an error term. Due to the uncertainty resulting from at-site estimation, it is assumed that the response variable $\mathbf{y} = (y_1, \mathbf{K}, y_n)$ has a sampling error $\boldsymbol{\epsilon} = \epsilon_1, \dots, \epsilon_n$, with covariance matrix $E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}') = \Sigma$ as described above. Stedinger and

Tasker (1985) proposed considering a second term of error η of variance $\sigma_{\eta}^2 > 0$ that is independent and identically distributed. Overall, the total error $\omega = \eta + \epsilon$ has covariance matrix

9)
$$\Lambda(\sigma_{\eta}^2) = \sigma_{\eta}^2 I + \Sigma$$
,

which provides a better characterization of the multiple sources of variability. Notice that the total covariance matrix is dependent on an unknown parameter σ_{η}^2 that describes the model variance, *i.e.* the part of the total variance that is not due to sampling.

For a known model variance σ_{η}^2 , the total covariance matrix can be rewritten $\Lambda(\sigma_{\eta}^2) = \sigma_{\eta}^2 \mathbf{G}$. Computing $\mathbf{G} = \mathbf{U}^T \mathbf{U}$ by Cholesky decomposition allows reformulating the GLS problem as an ordinary least squares (OLS) problem with transformed response variable $\mathbf{y}^* = \mathbf{U}^{-T} \mathbf{y}$ and design matrix $\mathbf{X}^* = \mathbf{U}^{-T} \mathbf{X}$. Therefore, the GLS estimator and its covariance matrix are derived directly from classical OLS theory:

10)
$$\hat{\boldsymbol{\beta}} = \left[\mathbf{X}' \mathbf{G}^{-1} \mathbf{X} \right]^{-1} \mathbf{X} \mathbf{G}^{-1} \mathbf{y}$$

11)
$$\Sigma_{\hat{\boldsymbol{\beta}}} = \sigma_{\boldsymbol{\eta}}^2 \left[\mathbf{X}^T \mathbf{G}^{-1} \mathbf{X} \right]^{-1}$$

Moreover, the residuals can be linked to the GLS residuals

12)
$$\mathscr{O} = \mathbf{y}^* - \mathbf{X}^* \hat{\boldsymbol{\beta}} = \mathbf{U}^{-T} \boldsymbol{\omega}$$
,

which has residual variance corresponding to the model variance σ_{η}^2 . In general, the model variance σ_{η}^2 is unknown and a proper estimation can be obtained by iterative least squares (Kjeldsen and Jones, 2009). The procedure consists in estimating the model parameters $\boldsymbol{\beta}$ from an initial guess obtained by OLS and to update σ_{η}^2 using the empirical variance of GLS residuals in equation (12). These two steps are repeated until convergence.

2.4 Pooling groups and super regions

Pooling groups are formed of the M sites that are the closest to the target site. To that end, three types of distance are considered: The geographical distance (GEO), the Mahalanobis distance between catchment descriptors (PHY) and the canonical distance. The Mahalanobis distance is selected because it considers not only the scales of the catchment descriptors, but also their covariance structure, which accounts for information redundancy (Cunderlik and Burn, 2006). The study of Oudin et al. (2010) showed that regions based on the similarity among catchment descriptors do not always easily translate into similarity in terms of hydrological properties. Consequently, regions derived from these two notions may lead to very different groups of sites. Canonical correlation analysis (CCA) can be used to create new canonical coordinates in the hydrological space that are mutually independent and sequentially maximizes correlation with catchment descriptors. The canonical distance defined as the distance between the canonical coordinates was used in RFA to delineate homogenous regions (Ouarda et al., 2001; Spence et al., 2007) or to perform spatial interpolation of the outcomes of hydrological models (Hundecha et al., 2008). For ungauged analysis, the hydrological information is unknown, but a substitute target position is estimated using the correlation with the catchment descriptors. Therefore the quality of the canonical distance for pooling groups depends on the relevance of the canonical distance and the quality of the estimated target in the canonical space (Durocher et al., 2016a).

The present study considers the formation of pooling groups directly using one of the three distances mentioned, but also considering combinations of two distance measures using a hierarchical approach. Specifically, the procedure involves identifying for each site a subset of M_1 sites located the closest to the target based on the first distance and then forming a pooling group of size $M_2 \leq M_1$ based on the second distance measure. For instance, a subset of the 100 nearest sites is extracted using the GEO distance and then a pooling group of size 25 sites is selected among them based on the PHY distance. Such approach will be denoted as GEO-PHY distance. This strategy is similar to the approach of Eng et al. (2007) that used pooling groups based on the PHY distance, but where a fixed distance was used instead of a constant number of sites.

If a pooling group contains too few sites, it will lead to a large predictive variance, but including many sites that are not relevant to the site of interest may create bias. In addition to the size of the pooling group, the choice of the catchment descriptors can also influence the quality of the fitting. To guide the calibration of a pooling group, it is recommended to find settings that optimize the quality of the prediction (Reis et al., 2005). Let y_0 be the flood quantiles at an ungauged location that has catchment descriptors \mathbf{x}_0 . Using the GLS estimator in equation (10) the predictive variance is given by

13)
$$\sigma_{\eta}^2 \left[1 + \mathbf{x}_0^T \left(\mathbf{X}^T \mathbf{G}^{-1} \mathbf{X} \right)^{-1} \mathbf{x}_0 \right],$$

which can be optimized by comparing various combinations of pooling group size and catchment descriptors.

Another strategy for the calibration of the ROI/GLS model consists in selecting the same size and catchment descriptors for a group G of sites. Notice that these groups will only be used for calibration and do not affect the formation of the pooling groups. Although such groups could take different forms, the present study defines super regions according to site drainage area and mean annual precipitation. Similar super regions were considered by Salinas et al. (2014a), but instead of using 3 straightforward divisions (i.e., small, medium and large), the present study uses a hierarchical agglomerative clustering

method (Murtagh and Legendre, 2014; Ward, 1963) to provide more objective boundaries, while maintaining meaningful interpretation.

Inside a super region, the prediction performance associated with specific settings can be evaluated by cross-validation. Formally, let $y_{(i)}$ be predicted flood quantiles (log) obtained at the *i*-th site when it is considered ungauged. Optimal settings may be identified by minimizing a criterion based on the predicted residuals $y_i - y_{(i)}$. In this line, two common performance criteria are the Nash-Sutcliffe and the mean absolute deviation

14) NSH =
$$100 \times \left\{ 1 - \frac{\sum_{i \in G} w_i (y_i - y_{(i)})^2}{\sum_{i \in G} w_i (y_i - \overline{y})^2} \right\}$$
 and MAD = $\frac{\sum_{i \in G} w_i |y_i - y_{(i)}|}{\sum_{i \in G} w_i}$,

where w_i are weights such that $\sum_i w_i = 1$ and \overline{y} is the weighted average of the y_i . For the calibration of the ROI/GLS model, the weights are taken as the record length of a site, which gives more importance to sites with more data.

3. Results

3.1 Data

The annual maximums of river discharge are extracted from daily records provided by Water Survey of Canada (2017) and catchment descriptors are provided by Environment and Climate Change Canada. Figure 1 presents the locations of 771 selected sites that possess at least 20 years of records and do not exhibit significant trends according to the Mann Kendall test (Önöz and Bayazit, 2012). Note that from the initial dataset of 918 sites where the desired catchment descriptors were available, 147 sites were removed due to the presence of trends.

The concept of super regions is to improve the interpretability of the results by dividing sites with similar scale control and climate. After some experimentation, the sites were divided into 8 super regions. Figure 2 presents the correspondence between position in the descriptor space and their locations. Super region 1 uniquely includes sites from the Pacific coast having small to medium drainage area and the largest mean annual precipitation. Other wetter sites are found in super regions 2 to 4, which are located either in British Columbia or in the southeastern part of Canada. The drier ones correspond to super regions 5 to 8, which are mostly located in the prairie provinces (Alberta, Saskatchewan and Manitoba) and in the

northern part of Canada. In general, the largest watersheds are found in the more northerly locations and correspond to super regions 7 and 8. In particular, super region 7 is wetter than super region 8 and is located mostly in the north of Quebec, Ontario and British Columbia. One can see that British Columbia includes a large variety of rivers as it includes sites from the eight super regions. Table 1 presents descriptive statistics for the drainage area and mean annual precipitation as well as the other available catchment descriptors, which include: basin compactness, average slope, streamflow density, percentage of waterbody area and site elevation.

The fitting and the selection of the best distribution at each gauged site is performed as described in the methodology section. The GEV distribution (including Gumbel) was preferred in 62.5 % of the gauged sites, followed by the Pearson type III (including Gamma) with 29.7 %. The other distributions were selected in lower proportions. The flood quantiles associated with return levels are then computed for each site using the selected distribution.

3.2 Estimation of the sample covariance matrix

The correlation matrix of Spearman's rho coefficient is estimated using the nonparametric model in equation (6), where an adjusted coefficient of determination (R^2) of 52% is obtained, indicating that the model fits the data fairly well. The relative importance of each smooth term is assessed by examining the difference of R^2 when that component is removed. As expected, the most important smooth term is f_h associated with the distance and has a relative importance of 21%. It is followed respectively by f_{xy} and f_A with 13% and 2%. These values indicate that the intersite correlation is mostly influenced by the distance between sites, but that a non-negligible effect is also depending mostly on the location. In other words, the dependence among pairs of sites cannot be explained uniquely by the separating distance. The effect of f_{xy} is illustrated in Figure 3 using a map of the smoothed terms, even though discharge can only

be observed at unique points on rivers. One can see different zones where the Spearman's rho is lower or higher than what would be expected considering uniquely f_h , the effect of the distance between sites.

Figure 3 also shows the correlation coefficient in respect of the distance where the solid line represents the POW model fitted by weighted least squares using record lengths as weights; with parameters $\alpha = 905$, $\tau = 0.14$ and $\gamma = 0.81$. Notice that the coefficient of correlation in respect of the distance becomes (in average) stable after roughly 600 km, but never reaches zero. The POW model is biased after that point as it can be assumed that for separations beyond 600 km, the correlation becomes more and more negligible. From a practical point of view, after 600 km the correlation is relatively low (less than 0.2). It is then reasonable to assume that it decreases to values that are not statistically different from zero.

The choice of a copula is investigated to determine the nature of the dependence among near sites. As the power for rejecting the null hypothesis increases with the strength of the dependence, only paired observations separated by less than 50 km and having at least 40 years of common record are selected. In total, 109 paired observations are identified and are presented in Figure 1. Note that the minimum Spearman's rho for these paired observations is 0.7. For each pair, goodness-of-fit tests are performed using as null hypothesis one of the following copulas: Normal, t-copula, Husler-Reiss, Gumbel, Galambos and Clayton (see, for instance, Salvadori et al. (2007) for a description of these copulas). The results show that only the Clayton copula is rejected at a significance level 5% in a majority of the cases, while the other copulas are not rejected for all pairs. In the present study, the size of the sample remains relatively small and so the goodness-of-fit tests do not have high power of discrimination. Nevertheless, these results do not indicate any evidence that would suggest rejecting these copulas.

Figure 4 presents a comparison of the log-likelihood between the bivariate Normal copula in respect of two bivariate extreme copulas. One can see that for most paired observations, the Normal copula provides the best fit as evidenced by there being more points below than above the 45° line. In particular, it suggests a

preference for a multivariate Normal copula over a max-stable process to characterize the spatial dependence among all observations. The decision between the Normal copula or t-copula will affect the behavior in terms of tail dependence. Fitting a multivariate t-copula on all paired observations led to estimated degrees of freedom $\upsilon=15$, which is associated with a relatively low tail dependence. For instance, with correlation coefficient of 0.5 and 0.3 the tail dependence is respectively 0.03 and 0.01. This agrees with the better fit of the bivariate Normal copula in comparison to extreme copulas for which tail dependence is an important characteristic. Therefore, the Normal copula is adopted as a reasonable model for the rest of the analysis.

The matrix of sampling covariances for the flood quantiles (log) is obtained for several return periods using Monte-Carlo simulations as described in the methodology. Figure 5 (left panel) illustrates the correlation between the paired Q100 in respect of the distance and the solid line represents the fitted POW model. Similarly, the right part of Figure 5 shows fitted POW model for various return periods. One can see that the strength of the correlation decreases with the return period. On average, the correlation between Q100 becomes less than 0.2 after roughly 100 km, while intersite correlation between annual maximum reach that point more around 400 km. Note that the zero tail dependence property of the Normal copula implies that for very large return periods the correlation between flood quantiles will continue to decrease towards zero.

3.3 Calibration of the pooling groups

As described in the methodology, the calibration of the ROI/GLS models may be guided by the predictive variance, equation (13), or by cross-validation, equation (14). In the following, each site is treated in turn as ungauged and a regression model is fitted using several combinations of pooling group sizes and catchment descriptors. More precisely, the pooling group sizes $M_2 = 20,25,K$, 80 and all combinations of three or more catchment descriptors are tried. The initial subset inside which the pooling groups are formed is restricted to $M_1 = 100$ or 200.

Table 2 presents the prediction performance obtained using different distances and comparing the calibration obtained using the predictive variance (Individual) and the criterion MAD inside super regions. Note that only the best hierarchical distances are reported. For GEO-PHY and GEO-CCA this corresponds to M_1 = 100 sites, while for PHY-GEO, M_1 = 200 sites. When the predictive variance is used, the geographical distance led to better predictive performance than the other distance in terms of MAD for both Q10 and Q100. However, all the results are relatively similar in terms of NSH, except for the direct canonical distance that has poorer performance. Nevertheless, better predictive performances are always found when the models are calibrated using super regions. In these cases, the GEO-PHY and the PHY-GEO distance perform similarly with a slight advantage to GEO-PHY. This illustrates the clear advantage of combining these two distances. In particular, the hierarchical approach GEO-CCA performs substantially better than the direct use of the canonical distance.

One objective of the present study is to examine the advantage of a nationwide analysis in comparison to an analysis based on smaller geographical areas. To explore this impact, the calibration of the ROI/GLS model is split in four distinct administrative regions. The first region regroups the four Atlantic provinces (New Brunswick, Nova Scotia, Prince Edward Island and Newfoundland and Labrador) with 90 sites. The second region combines Ontario and Quebec with 180 sites, while the third region regroups British Columbia and Yukon Territory with 213 sites. The Prairies and the two remaining territories complete the fourth region with 288 sites. For each region, 3 super regions were delineated inside. Note that these settings do not represent all possibilities that involve administrative boundaries. Nevertheless, it should illustrate the potential impact this has on the calibration of ROI/GLS. When using a nationwide analysis, the performance criteria for Q100 using the GEO-PHY distance are NSH = 91.05 and MAD = 0.343. The same criteria are slightly inferior when using administrative boundaries with NSH = 90.77 and MAD = 0.351. Similar outcomes are observed with Q10 where performance criteria pass from NSH = 93.49 and MAD = 0.296 to NSH = 92.85 and MAD = 0.317, respectively.

The results obtained in the present study are coherent with previous studies for portions of Canada. In the province of Quebec (Durocher et al., 2016a) and Ontario (Grover et al., 2002) relative root mean square errors (RRMSE) of 0.435 and 0.347 were reported for Q100. When using similar metrics, the ROI/GLS method used in the present study finds RRMSE of 0.438 and 0.346. However, such comparison must be done with care as these studies don't consider exactly the same set of gauged sites and were conducted on different periods.

3.4 Uncertainty analysis

Overall the best ROI/GLS models are obtained using super regions and the GEO-PHY distance, which is now examined in more detail. Table 3 reports the sizes and catchment descriptors used in the pooling groups of each super region and Table 4 separates the predictive performance by super regions. The wetter super region 1 and the larger super region 8 are well fitted with NSH greater than 90 for both flood quantiles; these super regions also have the lowest MAD. Good predictive performance is also observed for the wetter super regions 2 to 4. The prediction of the drier sites mostly found in the Prairies and the north is, however, relatively less accurate, especially for the smaller watersheds in super regions 5 and 6, with respective NSH of 32.23 and 65.43. These lower predictive performances are largely due to the presence of problematic sites. Figure 6 presents the predicted residuals standardized by the predictive variance. For Q10, one potential outlier is found in super region 5 where the predictive residual is higher than 8 standard deviations. For Q100, the same site remains problematic and other potential outliers arise in super region 6. This finding suggests that smaller and drier sites are more likely to be problematic to estimate than those of the other super regions. When predictive residuals higher than 3 standard deviations are removed for Q100, the NSH of super regions 5 and 6 is much better and becomes 76.02 and 73.35, respectively. Note that the NSH criterion evaluates the predictive performance in comparison to the weighted average, but does not represent a measure of the absolute uncertainty. For instance, super region 3 includes medium sized catchments strongly concentrated in southern Ontario. The NSH of Q100 for this super region is 69.72 which is less than super region 7 with 84.85. However, the MAD of super

region 3 is 0.29 in comparison to 0.38 for super region 7, showing that estimations are overall less uncertain. Even after removing outliers, the estimation of the flood quantiles Q10 and Q100 for the drier super regions 5 to 7 remain less accurate with MAD greater than 0.33, while less than 0.29 for the other super regions.

A sensitivity analysis is conducted to measure the effect of the intersite correlation model used to obtain the sampling covariance matrix on variability of the estimated flood quantile. Let σ_{NN}^2 , σ_{POW}^2 and σ_{NP}^2 denote the predictive variance obtained assuming intersite independence, the POW model and the nonparametric model. The nonparametric model is more flexible and so its predictive variance should be closer to the true value. Consequently σ_{NP}^2 is considered as a benchmark and Figure 7 shows the ratio of variance $\sigma_{POW}^2/\sigma_{NP}^2$ (left) and $\sigma_{NND}^2/\sigma_{NP}^2$ (right). The ROI/GLS model was calibrated using the GEO distance (top) and the GEO-PHY distance (bottom). The ratios $\sigma_{POW}^2/\sigma_{NP}^2$ of both distances (left) show that a smoother version of the intersite correlation has a limited impact on the evaluation of the model uncertainties, while the assumption of independence tends to underestimate it. In particular, the underestimation of the model uncertainties is affecting more the drier and larger basins found in super regions 3, 6, 7 and 8. Figure 7 (bottom-right) indicates that even though the distance GEO-PHY mainly formed the pooling group based on the distance between catchment descriptors, there is a substantial impact on the estimation of predictive variance from the covariance model. However, this underestimation is as expected less important than when using the GEO distance (top-right) that forces stronger intersite correlations in the pooling groups.

4. Conclusions

The ROI/GLS method was investigated to estimate flood quantiles at ungauged locations using a large database of 771 sites across Canada. The calibration procedure provides a general guideline to apply ROI/GLS regression in situations where no direct formula is available for evaluating the sample

covariance matrix. A pairwise fitting of copulas was considered among close sites to show that the Normal copula is a reasonable model for the intersite correlation structure. In particular, it reveals that extreme type of spatial dependence is not generally the best option for characterizing the relationship among annual maximum streamflow in Canada. Similar findings should be expected in other cold regions, because floods are dominated by spring snowmelt. Such behavior goes in the opposite direction to the recent interest in the max-stable processes in the characterization of spatial dependence between extreme rainfall events.

A nonparametric model was used to estimate the sampling covariance matrix of various pairs of flood quantiles. This approach was selected to better estimate the associations between paired observations when there are few or no years of common record. The main objective of using GLS in RFA is to obtain accurate estimation of the uncertainties of flood quantiles. In this regard, it was found that a simpler intersite correlation structure characterized by a power exponential model in respect of the geographical distance does not largely affect the estimation of the predictive variance and led to similar evaluation of the predictive variance. However, ignoring intersite correlation was shown to underestimate the predictive variance substantially. This is especially true for the drier basins located mostly in the Prairies and in the northern part of the country.

Cross-validation was used to evaluate and guide the calibration of the ROI/GLS models using three notions of distance. It was shown that the best choice was a hierarchical approach where first the 100 nearest sites to the target are identified according to the geographical distance and then a smaller pooling group is formed using the Mahalanobis distance between catchment descriptors. Additionally, super regions were delineated based on the drainage area and the mean annual precipitation to help the calibration of the ROI/GLS model. For each super region the same catchment descriptors and pooling group sizes were chosen. This strategy led to better predictive power than individually calibrating sites using the predictive variance. The effect of administrative boundaries on the calibration of the regression

models was also explored. The results showed that performing a nationwide RFA analysis resulted in ROI/GLS models with better predictive power. The concept of super regions was also found useful to calibrate the pooling groups and to better understand the quality of the flood estimations in respect of scale control and climate. The results indicated that the flood quantiles of the drier basins are estimated with greater uncertainty than the wetter ones. In particular, ROI/GLS resulted in rather poor predictions for some problematic sites in some of the smaller and drier basins.

Overall, the present study validates successful settings to carry out RFA using the ROI/GLS framework. In particular, it provides guidelines for estimating the sampling covariance matrix in a general context and using super regions to improve the calibration of the pooling groups. These outcomes are helpful to promote the adoption of GLS in RFA.

Acknowledgements

This work was supported by the Natural Science and Engineering Research Council (NSERC) Canadian FloodNet (# NETGP 451456 - 13). Computation were realized using the R environment (R Core Team, 2017). The authors would like to thank two anonymous reviewers for their valuable comments and suggestions to improve the quality of the paper.

References

- Rahman, A., Haddad, K., Kuczera, G., Weinman, E., 2016. Book 3: Peak flow estimation Ch. 3
 Regional flood methods, in: Australian Rainfall and Runoff A Guide to Flood Estimation.
 Commonwealth of Australia. http://book.arr.org.au.s3-website-ap-southeast-2.amazonaws.com/
- Basu, B., Srinivas, V.V., 2015. A recursive multi-scaling approach to regional flood frequency analysis. Journal of Hydrology 529, 373–383. https://doi.org/10.1016/j.jhydrol.2015.07.037
- Bárdossy, A., 2006. Copula-based geostatistical models for groundwater quality parameters. Water Resources Research 42. https://doi.org/10.1029/2005WR004754
- Bayazit, M., Önöz, B., 2004. Sampling variances of regional flood quantiles affected by intersite correlation. Journal of Hydrology 291, 42–51. https://doi.org/10.1016/j.jhydrol.2003.12.009
- Berg, D., 2009. Copula goodness-of-fit testing: an overview and power comparison. The European Journal of Finance 15, 675–701. https://doi.org/10.1080/13518470802697428

- Bezak, N., Brilly, M., Šraj, M., 2014. Comparison between the peaks-over-threshold method and the annual maximum method for flood frequency analysis. Hydrological Sciences Journal 59, 959–977. https://doi.org/10.1080/02626667.2013.831174
- Blöschl, G., Sivapalan, M., 1997. Process controls on regional flood frequency: Coefficient of variation and basin scale. Water Resour. Res. 33, 2967–2980. https://doi.org/10.1029/97WR00568
- Burn, D.H., 1990. An appraisal of the "region of influence" approach to flood frequency analysis. Hydrological Sciences Journal 35, 149–165. https://doi.org/10.1080/02626669009492415
- Burn, D.H., Whitfield, P.H., Sharif, M., 2016. Identification of changes in floods and flood regimes in Canada using a peaks over threshold approach. Hydrol. Process. 30, 3303–3314. https://doi.org/10.1002/hyp.10861
- Burnham, K.P., Anderson, D.R., 2002. Model Selection and Multimodel Inference: A Practical Information-Theoretic Approach, Springer Science & Business Media.
- Buttle, J.M., Allen, D.M., Caissie, D., Davison, B., Hayashi, M., Peters, D.L., Pomeroy, J.W., Simonovic, S., St-Hilaire, A., Whitfield, P.H., 2016. Flood processes in Canada: Regional and special aspects. Canadian Water Resources Journal / Revue canadienne des ressources hydriques 41, 7–30. https://doi.org/10.1080/07011784.2015.1131629
- Castellarin, A., Burn, D.H., Brath, A., 2008. Homogeneity testing: How homogeneous do heterogeneous cross-correlated regions seem? Journal of Hydrology 360, 67–76. https://doi.org/10.1016/j.jhydrol.2008.07.014
- Coles, S., 2001. An introduction to statistical modeling of extreme values. Springer Verlag.
- Cunderlik, J.M., Burn, D.H., 2006. Switching the pooling similarity distances: Mahalanobis for Euclidean. Water Resour. Res. 42, W03409. https://doi.org/10.1029/2005WR004245
- Dalrymple, T., 1960. Flood-frequency analysis. Survey Water-Supply Paper 1543.
- Demarta, S., McNeil, A.J., 2005. The t Copula and Related Copulas. International Statistical Review 73, 111–129. https://doi.org/10.1111/j.1751-5823.2005.tb00254.x
- Di Baldassarre, G., Laio, F., Montanari, A., 2009. Design flood estimation using model selection criteria. Physics and Chemistry of the Earth, Parts A/B/C, Recent developments of statistical tools for hydrological application 34, 606–611. https://doi.org/10.1016/j.pce.2008.10.066
- Douglas, E.M., Vogel, R.M., Kroll, C.N., 2000. Trends in floods and low flows in the United States: impact of spatial correlation. Journal of Hydrology 240, 90–105. https://doi.org/10.1016/S0022-1694(00)00336-X
- Durocher, M., Chebana, F., Ouarda, T.B.M.J., 2016a. Delineation of homogenous regions using hydrological variables predicted by projection pursuit regression. Hydrol. Earth Syst. Sci. 20, 4717–4729. https://doi.org/10.5194/hess-20-4717-2016
- Durocher, M., Chebana, F., Ouarda, T.B.M.J., 2016b. On the prediction of extreme flood quantiles at ungauged locations with spatial copula. Journal of Hydrology 533, 523–532. https://doi.org/10.1016/j.jhydrol.2015.12.029
- Durocher, M., Quessy, J.-F., 2017. Goodness-of-fit tests for copula-based spatial models. Environmetrics n/a-n/a. https://doi.org/10.1002/env.2445
- El-Jabi, N., Caissie, D., Turkkan, N., 2016. Flood analysis and flood projections under climate change in New Brunswick. Canadian Water Resources Journal 41, 319–330. https://doi.org/10.1080/07011784.2015.1071205

- Eng, P. C. Milly, Gary D. Tasker, 2007. Flood Regionalization: A Hybrid Geographic and Predictor-Variable Region-of-Influence Regression Method. Journal of Hydrologic Engineering 12, 585–591. https://doi.org/10.1061/(ASCE)1084-0699(2007)12:6(585)
- Floodnet, 2015. Floodnet NSERC Network Enhanced flood forecasting and management capacity in Canada [WWW Document]. URL http://www.nsercfloodnet.ca/ (accessed 7.26.17).
- Gado, T.A., Nguyen, V.-T.-V., 2016. Comparison of Homogenous Region Delineation Approaches for Regional Flood Frequency Analysis at Ungauged Sites. https://doi.org/10.1061/(ASCE)HE.1943-5584.0001312
- Genest, C., Rémillard, B., Beaudoin, D., 2009. Goodness-of-fit tests for copulas: A review and a power study. Insurance: Mathematics and economics 44, 199–213. https://doi.org/10.1016/j.insmatheco.2007.10.005
- Gräler, B., Pebesma, E., 2011. The pair-copula construction for spatial data: a new approach to model spatial dependency. Procedia Environmental Sciences 7, 206–211. https://doi.org/10.1016/j.proenv.2011.07.036
- Green, P.J., Silverman, B.W., 1993. Nonparametric regression and generalized linear models: a roughness penalty approach. Chapman & Hall/CRC.
- GREHYS, 1996a. Presentation and review of some methods for regional flood frequency analysis. Journal of Hydrology 186, 63–84.
- GREHYS, 1996b. Inter-comparison of regional flood frequency procedures for canadian rivers. Journal of hydrology(Amsterdam) 186, 85–103.
- Griffis, V.W., Stedinger, J.R., 2007. The use of GLS regression in regional hydrologic analyses. Journal of Hydrology 344, 82–95. https://doi.org/10.1016/j.jhydrol.2007.06.023
- Grover, P.L., Burn, D.H., Cunderlik, J.M., 2002. A comparison of index flood estimation procedures for ungauged catchments. Can. J. Civ. Eng. 29, 734–741. https://doi.org/10.1139/l02-065
- Haddad, K., Rahman, A., 2012. Regional flood frequency analysis in eastern Australia: Bayesian GLS regression-based methods within fixed region and ROI framework Quantile Regression vs. Parameter Regression Technique. Journal of Hydrology 430–431, 142–161. https://doi.org/10.1016/j.jhydrol.2012.02.012
- Higham, N.J., 2002. Computing the nearest correlation matrix—a problem from finance. IMA J Numer Anal 22, 329–343. https://doi.org/10.1093/imanum/22.3.329
- Hosking, J.R.M., 1990. L-Moments: Analysis and Estimation of Distributions Using Linear Combinations of Order Statistics. Journal of the Royal Statistical Society. Series B (Methodological) 52, 105–124.
- Hosking, J.R.M., Wallis, J.R., 1997. Regional frequency analysis: an approach based on L-moments. Cambridge Univ Pr.
- Hosking, J.R.M., Wallis, J.R., 1988. The effect of intersite dependence on regional flood frequency analysis. Water Resour. Res. 24, 588–600. https://doi.org/10.1029/WR024i004p00588
- Hundecha, Y., Ouarda, T.B.M.J., Bárdossy, A., 2008. Regional estimation of parameters of a rainfall-runoff model at ungauged watersheds using the "spatial" structures of the parameters within a canonical physiographic-climatic space. Water Resources Research 44. https://doi.org/10.1029/2006WR005439
- Interagency Committee on Water Data (IACWD), 1982. Guidelines for determining flood flow frequency: Bulletin 17-B (revised and corrected). Hydrol. Subcomm.

- Joe, H., 2015. Dependence modeling with copulas, Monographs on Statistics and Applied Probability. CRC Press, Boca Raton, FL.
- Kjeldsen, T.R., Jones, D.A., 2009. An exploratory analysis of error components in hydrological regression modeling. Water Resour. Res. 45, W02407. https://doi.org/10.1029/2007WR006283
- Kjeldsen, T.R., Jones, D.A., 2007. Estimation of an index flood using data transfer in the UK. Hydrological Sciences Journal 52, 86–98. https://doi.org/10.1623/hysj.52.1.86
- Kjeldsen, T.R., Jones, D.A., 2004. Sampling variance of flood quantiles from the generalised logistic distribution estimated using the method of L-moments. Hydrol. Earth Syst. Sci. 8, 183–190. https://doi.org/10.5194/hess-8-183-2004
- Kroll, C.N., Stedinger, J.R., 1998. Regional hydrologic analysis: Ordinary and generalized least squares revisited. Water Resour. Res. 34, 121–128. https://doi.org/10.1029/97WR02685
- Lindskog, F., Mcneil, A., Schmock, U., 2003. Kendall's tau for elliptical distributions. Credit risk: Measurement, evaluation and management 149–156.
- Madsen, H., Mikkelsen, P.S., Rosbjerg, D., Harremoës, P., 2002. Regional estimation of rainfall intensity-duration-frequency curves using generalized least squares regression of partial duration series statistics. Water Resour. Res. 38, 1239. https://doi.org/10.1029/2001WR001125
- Madsen, H., Rosbjerg, D., 1997. Generalized least squares and empirical bayes estimation in regional partial duration series index-flood modeling. Water Resour. Res. 33, 771–781. https://doi.org/10.1029/96WR03850
- Meigh, J.R., Farquharson, F.A.K., Sutcliffe, J.V., 1997. A worldwide comparison of regional flood estimation methods and climate. Hydrological Sciences Journal 42, 225–244. https://doi.org/10.1080/02626669709492022
- Merz, R., Blöschl, G., 2005. Flood frequency regionalisation—spatial proximity vs. catchment attributes. Journal of Hydrology 302, 283–306. https://doi.org/10.1016/j.jhydrol.2004.07.018
- Murtagh, F., Legendre, P., 2014. Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion? J Classif 31, 274–295. https://doi.org/10.1007/s00357-014-9161-z
- Nelsen, R.B., 2006. An introduction to copulas. Springer.
- Neves, M., Gomes, D., 2011. Geostatistics for spatial extremes. A case study of maximum annual rainfall in Portugal**. Procedia Environmental Sciences 7, 246–251. https://doi.org/10.1016/j.proenv.2011.07.043
- Önöz, B., Bayazit, M., 2012. Block bootstrap for Mann–Kendall trend test of serially dependent data. Hydrol. Process. 26, 3552–3560. https://doi.org/10.1002/hyp.8438
- Ouarda, T.B.M.J., Bâ, K.M., Diaz-Delgado, C., Cârsteanu, A., Chokmani, K., Gingras, H., Quentin, E., Trujillo, E., Bobée, B., 2008. Intercomparison of regional flood frequency estimation methods at ungauged sites for a Mexican case study. Journal of Hydrology 348, 40–58. https://doi.org/10.1016/j.jhydrol.2007.09.031
- Ouarda, T.B.M.J., Girard, C., Cavadias, G.S., Bobée, B., 2001. Regional flood frequency estimation with canonical correlation analysis. Journal of Hydrology 254, 157–173. https://doi.org/10.1016/S0022-1694(01)00488-7
- Oudin, L., Kay, A., Andréassian, V., Perrin, C., 2010. Are seemingly physically similar catchments truly hydrologically similar? Water Resources Research 46, n/a–n/a. https://doi.org/10.1029/2009WR008887

- R Core Team, 2017. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria.
- Reis, D.S., Stedinger, J.R., Martins, E.S., 2005. Bayesian generalized least squares regression with application to log Pearson type 3 regional skew estimation. Water Resour. Res. 41, W10419. https://doi.org/10.1029/2004WR003445
- Renard, B., 2011. A Bayesian hierarchical approach to regional frequency analysis. Water Resour. Res 47, W11513. https://doi.org/10.1029/2010WR010089
- Ribeiro-Corréa, J., Cavadias, G.S., Clément, B., Rousselle, J., 1995. Identification of hydrological neighborhoods using canonical correlation analysis. Journal of Hydrology 173, 71–89. https://doi.org/10.1016/0022-1694(95)02719-6
- Robson, A., Reed, D., 1999. Flood estimation handbook. Institute of Hydrology, Wallingford.
- Salas, J.D., Heo, J.H., Lee, D.J., Burlando, P., 2013. Quantifying the Uncertainty of Return Period and Risk in Hydrologic Design. Journal of Hydrologic Engineering 18, 518–526. https://doi.org/10.1061/(ASCE)HE.1943-5584.0000613
- Salinas, J.L., Castellarin, A., Kohnová, S., Kjeldsen, T.R., 2014a. Regional parent flood frequency distributions in Europe Part 2: Climate and scale controls. Hydrol. Earth Syst. Sci. 18, 4391–4401. https://doi.org/10.5194/hess-18-4391-2014
- Salinas, J.L., Castellarin, A., Viglione, A., Kohnová, S., Kjeldsen, T.R., 2014b. Regional parent flood frequency distributions in Europe Part 1: Is the GEV model suitable as a pan-European parent? Hydrol. Earth Syst. Sci. 18, 4381–4389. https://doi.org/10.5194/hess-18-4381-2014
- Salvadori, G., De Michele, C., Kottegoda, N., Rosso, R., 2007. Extremes in nature: an approach using copulas. Springer Verlag.
- Sandrock, G., Viraraghavan, T., Fuller, G.A., 1992. Estimation of Peak Flows for Natural Ungauged Watersheds in Southern Saskatchewan. Canadian Water Resources Journal / Revue canadienne des ressources hydriques 17, 21–31. https://doi.org/10.4296/cwrj1701021
- Schabenberger, O., Gotway, C.A., 2004. Statistical methods for spatial data analysis. CRC Press.
- Shang, H., Yan, J., Zhang, X., 2011. El Niño–Southern Oscillation influence on winter maximum daily precipitation in California in a spatial model. Water Resour. Res. 47, W11507. https://doi.org/10.1029/2011WR010415
- Smith, R.L., 1985. Maximum likelihood estimation in a class of nonregular cases. Biometrika 72, 67–90. https://doi.org/10.1093/biomet/72.1.67
- Spence, C., Saso, P., Rausch, J., 2007. Quantifying the impact of hydrometric network reductions on regional streamflow prediction in Northern Canada. Canadian Water Resources Journal 32, 1+.
- Stedinger, J.R., Tasker, G.D., 1985. Regional Hydrologic Analysis: 1. Ordinary, Weighted, and Generalized Least Squares Compared. Water Resour. Res. 21, 1421–1432. https://doi.org/10.1029/WR021i009p01421
- Tasker, G., Stedinger, J., 1989. An operational GLS model for hydrologic regression. Journal of Hydrology 111, 361–375. https://doi.org/10.1016/0022-1694(89)90268-0
- Tasker, G.D., 1980. Hydrologic regression with weighted least squares. Water Resources Research 16, 1107–1113.
- Thomas, D., Benson, M., 1975. Generalization of streamflow characteristics from drainage-basin characteristics. US Geological Survey Water-Supply Paper.

- Vogel, R.M., Kroll, C.N., 1990. Generalized Low-Flow Frequency Relationships for Ungaged Sites in Massachusetts1. JAWRA Journal of the American Water Resources Association 26, 241–253. https://doi.org/10.1111/j.1752-1688.1990.tb01367.x
- Wang, Z., Yan, J., Zhang, X., 2014. Incorporating spatial dependence in regional frequency analysis. Water Resour. Res. 50, 9570–9585. https://doi.org/10.1002/2013WR014849
- Ward, J.H., 1963. Hierarchical Grouping to Optimize an Objective Function. Journal of the American Statistical Association 58, 236–244. https://doi.org/10.1080/01621459.1963.10500845
- Westra, S., Sisson, S.A., 2011. Detection of non-stationarity in precipitation extremes using a max-stable process model. Journal of Hydrology 406, 119–128. https://doi.org/10.1016/j.jhydrol.2011.06.014
- Wood, S., 2006. Generalized additive models: an introduction with R. Chapman & Hall/CRC.
- Wood, S.N., 2003. Thin plate regression splines. Journal of the Royal Statistical Society: Series B (Statistical Methodology) 65, 95–114. https://doi.org/10.1111/1467-9868.00374
- WSC, 2017. Water Survey of Canada [WWW Document]. URL http://www.wsc.ec.gc.ca/applications/H2O/index-eng.cfm
- Zhang, Z., Stadnyk, T.A., Burn, D.H., 2018. Identification of a broadly accepted statistical distribution for at-site flood frequency analysis in Canada. Presented at the 71rst CWRA National Conference, May 28 June 1, Victoria, BC, Canada.
- Zrinji, Z., Burn, D.H., 1994. Flood frequency analysis for ungauged sites using a region of influence approach. Journal of Hydrology 153, 1–21. https://doi.org/10.1016/0022-1694(94)90184-8

Tables

Table 1: Summary statistics for runoff and catchment descriptors.

Variables	Abrv.	Min	Q1	Med	Avg	Q3	Max	Table 2:
Record length (yr)		20	25	36	39	48	111	Performa
mean of annual maximum discharge (m³/s)		0.2	13.4	45.6	206.9	174.1	5068.3	nce criteria
Basin area (km²)	area	1	146	460	2829	1992	48867	for flood
Basin compactness	comp	0.4	1.7	2.5	2.6	3.4	6.3	quantiles
Basin mean slope (%)	slope	< 0.1	1.2	3.6	10.5	17.1	59.0	using the ROI/GL
Waterbody area (%)	wb	< 0.1	0.4	1.3	3.7	4.5	38.3	S model
Stream density (km ⁻¹)	dens	< 0.1	0.6	1.0	1.2	1.6	3.4	with
Elevation at site (m)	elev	1	181	382	474	731	1699	different
Mean annual precipitation (mm)	map	213	498	761	836	1052	3216	distance. Calibrati

on using predictive variance (individual) and super region.

		Indivi	dual	Super Region		
Variable	Distance	NSH	MAD	NSH	MAD	
Q10	PHY	91.50	0.351	93.21	0.309	
	GEO-PHY	92.05	0.337	93.49	0.296	
	GEO	92.27	0.322	93.34	0.302	
	PHY-GEO	91.53	0.343	93.48	0.299	
	CCA	88.29	0.422	89.28	0.400	

	GEO	-CCA	91.93	0.346	93.31	0.308	* Bold indicates best results in each column		
Q100	PHY		88.74	0.411	90.72	0.362	Table 3: Pooling group sizes and catchment		
	GEO	-PHY	89.15	0.392	91.05	0.343	descriptors for ROI/GLS model using		
	GEO		89.57	0.378	90.80	0.351	GEO-РНҮ.		
		-GEO	89.40	0.388	91.03	0.347	Table 4: Performance criteria for flood		
	CCA		84.42	0.488	85.96	0.457	quantiles using the ROI/GLS model with a		
		-CCA	88.69	0.406	90.61	0.360	GEO-PHY distance detailed by super		
	Super	-cca					region.		
	region	Size	Catchme	nt desci	riptors		108-011		
Q10	1		_		elev + map)			
	2		area + con	_	_		Figures		
	3 4		area + slop area + slop		elev + map)	Figure 1: Locations of 771		
	5		area + slop area + slop		_		gauged sites in Canada and		
	6		area + den		-		average locations of 109 pairs		
	7		area + dens + slope + wb + elev + map of sites separated by less than						
	8		area + dens + slope + wb + elev + map 50 km and having 40 years of						
Q100	1		area + slope + wb + elev + map area + comp + slope + wb common record.						
	2 3			-	-				
	4		area + slo area + slo				Figure 2: Super regions in		
	5		area + slo				geograpical (top) and		
	6		area + wb				descriptor space (bottom).		
	7				elev + ma		Figure 3: At left, correlation		
C	8	60 ;		mp + dei	ns + slope +	wb + elev	coefficient estimated from the		
Super	Q10		Q100		nonpara	metric i	nodel. The dashed line represents the fitted		
region	NSH	MAD	NSH	MAD			right, a represention of the component f_{xy} of		
1	96.40	0.188	93.13	0.236			ric model on a grid of locations.		
2	89.29	0.243	83.25	0.271	the nong	Jai ailieti	ic model on a grid of locations.		
3	78.72	0.240	69.72	0.290	Figure 4	l: Log-li	kelihood of fitted copulas for the 109 paired		
4	83.83	0.223	75.93	0.293	sites in I	igure 1.			
5	31.00	0.618	32.23	0.625	T2:	5. A4 1.	64		
6	75.41	0.406	65.43	0.476	estimate		eft, correlation between paired Q100 (log) Monte-Carlo simulations using the		
7	88.44	0.353	84.85	0.380		•	model. The solid line represents the fitted		
8	93.99	0.215	90.92	0.248			at right, the POW models are reported by		
Total	93.56	0.213	91.17	0.343			as a reference, the dashed line indicates the		
Total	75.50	0.270	71.17	0.573	_		intersite correlation found in Figure 3.		

Figure 6: Standardized predictive residuals using GEO-PHY and super regions.

Figure 7: Comparison of the ratio of predictive variance for Q10. The denominator is the predictive variance deduced from the nonparametric model, while the numerator is deduced from the POW model (left) and the assumption of independence (right). Pooling groups are formed using GEO distance (top) and GEO-PHY bottom.

- Characterization of the intersite correlation using copulas.
- Improved calibration of ROI/GLS models using super regions and hybrid distances.
- Uncertainty analysis of regional flood quantile estimates in Canada.