

Université du Québec  
Institut National de la Recherche Scientifique  
Centre Eau Terre Environnement

**Application du modèle <<Random Forest Regression>> dans l'analyse  
fréquentielle régionale des évènements hydrologiques extrêmes**

Par  
Shitanshu Desai

Mémoire présenté pour l'obtention du grade de  
Maître ès sciences (M.Sc.) en sciences de l'eau

**Jury d'évaluation**

Président du jury et  
examineur interne

Fateh Chebana  
INRS-ETE

Examineur externe

Ali Assani  
Université du Québec à Trois-Rivières

Directeur de recherche

Taha BMJ Ouarda  
INRS-ETE

## AVANT-PROPOS

Ce mémoire de maîtrise par article est composé de deux chapitres. Le premier chapitre intitulé «Synthèse», fait état de la problématique et de la pertinence de mon sujet de recherche ainsi que de ma contribution à ce travail de recherche. Les détails sur la méthode théorique utilisée ainsi que les résultats obtenus sont présentés dans l'article au chapitre II de ce document. L'article a été soumis au journal «Journal of Hydrology» en Septembre 2018 et est en cours de révision. Le choix de cette revue vient du fait qu'elle est très prestigieuse.

Le titre et les auteurs de l'article sont les suivants :

**Titre :** Regional Hydrological Frequency Analysis at Ungauged Sites with Random Forest Regression

**Auteurs :** Shitanshu Desai, Taha B.M.J Ouarda

Les contributions des auteurs sont divisées comme suit :

**Shitanshu Desai :** Contribution à l'élaboration de la méthodologie, à la revue de la littérature, à la programmation du modèle, à la production des résultats, à l'analyse des résultats et à la rédaction de l'article.

**Taha B.M.J Ouarda :** Contribution au développement de la méthodologie, à l'interprétation des résultats et à la rédaction de l'article.



## REMERCIEMENTS

Je voudrais remercier mon directeur, le professeur Taha Ouarda, pour ses conseils et son encadrement. Ma gratitude va également à tous les membres du jury pour avoir accepté de juger ma thèse. Je voudrais remercier mes parents, Rakesh et Dipali et mon frère Akanshu. Je voudrais remercier ma petite amie Isha qui était toujours avec moi. J'aimerais remercier mes amis, Adnen, Amina, Habiba, Ramina et tous ceux que j'aurais pu aider dans ma maîtrise.





## RÉSUMÉ

Les inondations provoquent des dégâts importants sur les plans environnemental, économique et social. Il faudrait alors estimer adéquatement la fréquence des événements hydrologiques extrêmes pour une meilleure planification et gestion des ressources en eau ce qui assure la sécurité publique. L'analyse de fréquence des variables hydrologiques est une approche couramment utilisée lorsque l'information hydrologique est disponible au site d'intérêt. Cependant, il est souvent nécessaire d'estimer les événements extrêmes sur des sites dits non jaugés où aucune observation hydrologique n'est disponible. Au niveau de ces sites, l'analyse Fréquentielle Régionale (AFR) peut être utilisée pour estimer leurs caractéristiques des inondations. AFR comprend deux étapes principales. La première consiste à délimiter les régions hydrologiquement homogènes. Dans cette étude, nous utilisons l'analyse canonique des corrélations (ACC) pour la délimitation de ces régions. La deuxième étape consiste à appliquer un modèle d'estimation régional pour chaque région délimitée.

Traditionnellement, les modèles linéaires ont été utilisés dans l'AFR mais ils ne peuvent pas établir les relations complexes entre les variables. Ainsi, un certain nombre de techniques non linéaires ont été proposées dans la littérature comme les réseaux des neurones, les modèles additifs généralisés, etc. Cependant, « Random Forest (RF) », où «Forêts Aléatoires» en français, présente une technique générale et puissante, qui n'est pas utilisée souvent. La « Random Forest Regression (RFR) » », où «Forêts Aléatoires de régression» en français, est une technique d'apprentissage d'ensemble non linéaire, non paramétrique et capable d'établir des relations complexes ainsi que donner des estimations plus fiables.

L'objectif de la présente étude est d'introduire la technique RF dans l'estimation régionale des quantiles de crue. La RFR sert à établir des relations non linéaires entre les caractéristiques physio-météorologiques du bassin versant et les caractéristiques du débit, et à estimer les caractéristiques des crues dans les sites non jaugés. La RFR est également appliquée aux régions hydrologiquement homogènes, obtenues à l'aide de l'ACC (ACC-RFR), pour l'estimation de leurs quantiles de crue. Suite à une étude de cas à la province de Québec (Canada), une analyse comparative entre plusieurs approches d'estimation des quantiles de crue est effectuée. Les résultats indiquent que l'ACC-RFR donne les meilleurs résultats d'estimation parmi les modèles testés en termes d'erreur quadratique moyenne. L'utilisation de l'ACC, pour délimiter les zones hydrologiquement homogènes, améliore considérablement la performance de RFR.

# TABLE DES MATIÈRES

AVANT-PROPOS.....	III
REMERCIEMENTS.....	V
RÉSUMÉ.....	VII
LISTE DES TABLEAUX.....	IX
LISTE DES FIGURES.....	X
CHAPITRE I: SYNTHÈSE.....	1
1. INTRODUCTION.....	2
1.1 PROBLÈME.....	3
1.2 HYPOTHÈSE.....	3
1.3 OBJECTIF.....	4
1.4 ORIGINILATÈ.....	4
1.5 ORGANISATION DE LA SYNTHÈSE.....	4
2. METHODES.....	4
2.1 ANALYSE CANONIQUE DES CORRELATIONS (ACC).....	4
2.2 RANDOM FOREST (RF).....	5
2.3 CRITÈRE D'ÉVALUATION.....	5
3. ZONE D'ÉTUDE.....	6
4. RÉSULTATS.....	6
5. CONCLUSIONS ET RECOMMANDATIONS.....	7
6. RÉFÉRENCES.....	9
CHAPITRE II: ARTICLE.....	11
REGIONAL HYDROLOGICAL FREQUENCY ANALYSIS AT UNGUAGED SITES WITH RANDOM FOREST REGRESSION.....	12
1. INTRODUCTION.....	15
2. THEORETICAL BACKGROUND.....	17
2.1. Random Forest Regression.....	17
2.2. CCA approach in RFA.....	18
2.3. Selection of Methods for Comparison.....	19
2.4. Evaluation Metrics.....	20
2.5. Evaluation Procedure.....	20
3. CASE STUDY.....	20
4. RESULTS.....	21
5. CONCUSIONS.....	23
6. ACKNOWLEDGEMENTS.....	24
7. BIBLIOGRAPHY.....	25



# LISTE DES TABLEAUX

Table 1 : Descriptive Statistics of physio-meterological and Hydrological Variables.....28

Table 2: NASH, RMSE, RMSEr, BIAS and BIASr values for all the models. Best values for each quantile for the corresponding metrics are marked in bold.....29

Table 3: Feature Importance of Five Input Variables used for Specific Flood Quantile Estimation.....30

## LISTE DES FIGURES

Figure 1: Number of trees ( $n_{\text{estimators}}$ ) vs OOB error rate for 10, 50 and 100-year flood quantiles.....	31
Figure 2: Relative errors associated with quantiles $q_{50}$ calculated using RFR and CCA-RFR (the sites are ordered according to their area).....	32
Figure 3: A) $q_{10}$ , B) $q_{50}$ and C) $q_{100}$ estimation using RFR approach.....	33
Figure 4: A) $q_{10}$ , B) $q_{50}$ and C) $q_{100}$ estimation using CCA-RFR approach.....	34

## **CHAPITRE 1 : SYNTHÈSE**

# 1. INTRODUCTION

Les inondations représentent l'une des catastrophes naturelles les plus fréquentes ([Stefanidis and Stathis, 2013](#)). De 1995 à 2015, les inondations ont représenté plus de 43 % de toutes les catastrophes naturelles, touchant 2,3 milliards de personnes et en tuant 157 000 personnes ([Wallemacq, 2015](#)). Au cours de la même période, plus de 662 milliards de dollars américains de dommages économiques ont été subis dans le monde entier. Il est donc de la plus haute importance de prévoir adéquatement les caractéristiques de ces événements sur tous les sites.

En raison de leur coût élevé et de leur impact énorme, les installations hydrauliques doivent être conçues avec la plus grande précision et en utilisant toutes les informations disponibles. La compréhension adéquate de la nature probabiliste des débits de crue et l'estimation fiable des quantiles de conception sont donc cruciales sur un site où, par exemple, un barrage doit être construit. Toutefois, dans la pratique, les données hydrologiques ne sont disponibles qu'à un nombre limité d'endroits et nous sommes souvent tenus d'estimer les quantiles de conception à des endroits où il n'y a pas ou peu d'information disponible. Dans ce cas, l'analyse Fréquentielle Régionale (AFR) est couramment utilisée pour l'estimation des caractéristiques de débit.

L'AFR permet le transfert de l'information à partir des sites jaugés vers le site d'intérêt non jaugé. Cette approche comprend habituellement deux étapes principales. La première étape consiste à délimiter des régions hydrologiquement homogènes. Dans cette étape, les sites qui sont similaires selon certains critères d'homogénéité sont regroupés. Comme les sites d'une région homogène donnée sont similaires, l'information peut raisonnablement être transférée d'un site jaugé à un site non jaugé. La deuxième étape est l'application d'un modèle d'estimation régional dans chaque région délimitée ([Ouarda, 2013](#)). Les modèles d'estimation régionaux sont ensuite formés pour établir des relations fonctionnelles entre les caractéristiques physio-météorologiques des bassins et les caractéristiques d'écoulement des bassins non jaugés.

Dans cette étude, nous utilisons l'analyse canonique des Corrélations (ACC) au niveau l'étape de la délimitation. L'ACC a été utilisée dans un certain nombre d'études (Chebana et al. 2014 ; Shu et Ouarda, 2007). De plus, selon une étude publiée par [Ouarda et al. \(2007\)](#), la délimitation du voisinage avec l'ACC donne de meilleurs résultats. Donc le but de cette présente étude porte principalement sur l'amélioration de la précision d'estimation de l'approche AFR.

La plupart des techniques d'estimation au niveau de l'AFR sont des modèles linéaires. Cependant, les systèmes hydrologiques sont caractérisés par des processus complexes. Les méthodes linéaires ne parviennent pas à expliquer ces relations complexes entre les variables dépendantes et indépendantes. Ces modèles sont remplacés par plusieurs autres techniques non linéaires comme le « Artificial Neural Network (ANN) », les modèles additifs généralisés (GAM), l'ACC non linéaire, etc. Random Forest (RF) est l'une de ces techniques non linéaires. Lorsque la RF est utilisée à des fins d'approximation ou de régression fonctionnelle, elle est appelée « Random Forest Regression (RFR) » ou « Regression Forests » ([Breiman, 2001](#)). Dans RFR, à partir d'un ensemble donné de données, plusieurs échantillons sont tirés au hasard et des arbres de classification et de régression (CART) sont construits. Finalement, les résultats de tous ces arbres sont combinés et une estimation des variables cibles est obtenue en faisant la moyenne des résultats de chaque arbre.

Quelques études ont été réalisées en hydrologie par RFR. [Monira et al. \(2010a\)](#) et [Taksande and Mohod \(2015\)](#) ont respectivement utilisé les RF pour la prévision quotidienne et mensuelle des précipitations. [Chen et al. \(2012\)](#) ont utilisé les RF pour construire un modèle de prévision des sécheresses. [Nguyen et al. \(2015a\)](#) ont utilisé les RF pour prévoir les niveaux d'eau quotidiens. [Wang et al. \(2015\)](#) ont élaboré un modèle d'évaluation des risques d'inondation fondé sur les RF. Cependant, aucune étude n'a jamais été faite sur l'utilisation des RF dans l'approche de AFR. De plus, à la connaissance de l'auteur, aucune étude n'a été entreprise sur l'utilisation de l'ACC en conjonction avec les RF.

## 1.1 PROBLÈME

Les systèmes hydrologiques sont très complexes. Les modèles linéaires ne peuvent pas décrire ces relations entre les variables dépendantes et indépendantes. De plus, plusieurs modèles non linéaires complexes ont été utilisés dans l'AFR mais ils n'ont pas donné une bonne précision.

## 1.2 HYPOTHÈSE

« Random Forest » est une technique non paramétrique et non linéaire. Le couplage d'une telle technique avec l'ACC devrait accroître la précision et l'efficacité.

### **1.3 OBJECTIFS**

L'objectif de la présente étude est d'introduire la technique RFR pour l'estimation régionale des quantiles de crue. Bien que la RFR soit une technique célèbre et très précise, elle n'a pas été beaucoup utilisée en hydrologie.

### **1.4 ORIGINALITÉ**

À notre connaissance, il n'existe encore aucune étude combinant l'ACC et les RRF.

### **1.5 ORGANISATION DE LA SYNTHÈSE**

Cette synthèse est divisée en cinq parties principales. La section 1 présente l'introduction. La section 2 explique les méthodes. La section 3 fournit des renseignements sur l'étude de cas et les données utilisées. La section 4 présente les résultats. La section 5 présente les conclusions et les recommandations.

## **2. MÉTHODES**

L'objectif de la présente étude est d'introduire la technique RFR pour l'estimation régionale des quantiles de crues. La RFR sert à établir des relations non linéaires entre les caractéristiques physio-météorologiques du bassin et les caractéristiques des débits, et à estimer les caractéristiques des crues dans les sites non jaugés. La RFR est également appliquée aux voisinages hydrologiques identifiés à l'aide de l'ACC (ACC-RFR).

Les principes et les formules de base des différentes approches considérées dans cette étude sont présentés dans l'article. Dans les sous-sections suivantes, les différents outils statistiques utilisés sont résumés.

### **2.1 ANALYSE CANONIQUE DES CORRÉLATIONS (ACC)**

Soit  $X$  et  $Y$  un ensemble des variables aléatoires contenant, respectivement, des variables physio-météorologiques des bassins versants (e.g. la superficie du bassin, la pente, les caractéristiques des précipitations, etc.) et des variables

hydrologiques telles que les quantiles de crue. L'ACC permet d'identifier des paires de combinaisons linéaires de chaque ensemble de variables appelées variables canoniques. Il est fait de telle sorte que la corrélation entre la paire de variables soit maximisée et que la corrélation entre les variables canoniques hydrologiques/physiographiques distinctes soit nulle. Pour un site non jaugé (où l'information hydrologique peut ne pas être disponible), l'ACC permet de créer un voisinage formé des sites jaugés à partir desquels s'effectue le transfert de l'information hydrologique vers le site cible.

## 2.2 RANDOM FOREST (RF)

RFR est l'un des algorithmes d'apprentissage général les plus précis. RF performe mieux par rapport aux algorithmes d'arbre simple tel que CART. Il est rapide et présente des taux d'erreur comparables à ceux des algorithmes plus traditionnels. Dans la RFR, les prédicteurs « Regression Trees »  $h(x, \theta_k)$ ,  $k = 1 \dots M$  prennent des valeurs numériques en fonction des vecteurs aléatoires  $\{\theta_k\}$  ([Breiman, 2001](#)). Les données de formation sont tirées au hasard et indépendamment d'une distribution conjointe de  $(X, Y)$ , où le vecteur aléatoire  $X$  est l'entrée observée et le vecteur aléatoire  $Y$  est la sortie numérique attendue. Les arbres individuels sont cultivés à l'aide de l'algorithme CART (Classification and Regression Trees).

RFR possède deux caractéristiques importantes : le taux d'erreur « out-of-bag (OOB) » et l'importance variable. En général, nous utilisons environ les deux tiers des données d'un échantillon bootstrap et le tiers restant est exclu. C'est ce qu'on appelle les échantillons (OOB). L'erreur estimée sur ces échantillons omis est connue sous le nom de taux d'erreur OOB. Le taux d'erreur OOB peut être utilisé pour la validation ainsi que pour le calcul du nombre optimal d'arbres requis. Le « Variable Importance » variable permet d'identifier les variables prédicteurs les plus utiles pour la prédiction. L'importance variable peut être calculée à l'aide des RF en enregistrant les améliorations, à chaque nœud de chaque arbre de la forêt.

## 2.3 CRITÈRES D'ÉVALUATION

Dans la présente étude, le NASH (Nash–Sutcliffe model efficiency), la racine carrée de l'erreur quadratique moyenne (RMSE), la racine carrée de l'erreur quadratique moyenne relative (rRMSE), le biais moyen (BIAS) et le biais moyen relatif (rBIAS) sont utilisés pour évaluer les performances des modèles.

### 3. ZONE D'ÉTUDE

L'ensemble de données utilisées dans la présente étude comprend 151 stations hydrométriques situées dans le sud de la province de Québec-Canada. Les stations sont exploitées par Ministère du développement durable, de l'environnement et de la lutte contre les changements climatiques de la province de Québec. L'ensemble de données adoptées a été utilisé dans un certain nombre d'études d'AFR antérieures ([Chebana and Ouarda, 2008](#); [Chokmani and Ouarda, 2004](#); [Shu and Ouarda, 2007](#)), ce qui permet de faire une comparaison des résultats obtenus avec d'autres méthodes. Sur la base des travaux de [Chokmani and Ouarda \(2004\)](#) avec la même base de données, cinq variables physiométrologiques sont sélectionnées, dont trois sont des variables physiographiques et deux sont des variables météorologiques. Ces variables sont la superficie du bassin (AREA), la pente moyenne du bassin (MBS), la fraction de la superficie du bassin occupée par les lacs (FAL), les précipitations annuelles moyennes totales (AMP) et les degrés-jours annuels moyens au-dessus de 0° (AMD). Pour enlever les effets d'échelle, des quantiles spécifiques (quantiles divisés par la superficie du bassin) sont utilisés. Les quantiles de période des retours 100 ans, 50 ans et 10 ans (q100, q50 et q10 respectivement) sont les trois quantiles spécifiques utilisés dans la présente étude.

### 4. RÉSULTATS

Dans RF, la taille de l'ensemble de données, le nombre d'arbres et le nombre de variables utilisées dans chaque sous-ensemble ont un impact énorme sur le taux d'erreur. Comme la taille de l'ensemble de données n'est pas un paramètre ajustable, nous ajustons le nombre d'arbres et le nombre de variables à chaque sous-ensemble. Le nombre de variables à chaque sous-ensemble a été fixé à 3.

Nous avons réglé le nombre d'arbres en utilisant le taux d'erreur OOB. Le taux d'erreur OOB peut être utilisé pour le calcul du nombre optimal d'arbres requis. Dans cette étude, la plus faible taux d'erreur d'OOB correspond à 30 arbres. De plus, le critère de fractionnement considéré est l'erreur quadratique moyenne.

Les résultats sont décrits en détail dans l'article et sont brièvement discutés dans cette sous-section. Les résultats indiquent que le ACC-RFR est supérieur ou comparable aux autres modèles en se basant sur les différents critères d'évaluation des performances, sauf le critère NASH. De plus, l'ACC-RFR performe mieux que la RFR en se basant sur les différents critères autres que le NASH. Dans la ACC-RFR, seulement les stations des voisinages hydrologiques sont considérées pour la prévision, la taille de l'échantillon est considérablement



plus petite que l'ensemble de données originales. De plus, le critère de NASH est fortement influencé par le modèle utilisé ([Schaeffli and Gupta, 2007](#)). Ceci peut être à l'origine des mauvais résultats en terme de NASH.

Bien que nous ayons de faibles valeurs pour le critère NASH pour la RFR et l'ACC-RFR par rapport à d'autres modèles, nous pouvons observer que le ACC-RFR conduit aux meilleures estimations en termes de RMSE et RMSEr. La RMSE fournit une évaluation de la précision de la prédiction à une échelle absolue tandis que la RMSEr fait de même à une échelle relative. La RFR, combinée avec l'ACC, offre une meilleure prédiction que celle obtenue avec le modèle de base (RFR).

Le BIAS et le BIASr sont des critères d'évaluation utilisés pour déterminer si le modèle surestime ou sous-estime les diverses estimations. En général, c'est le modèle ACC-RFR qui présente le BIASr le plus faible de tous les modèles considérés et le BIASr est également comparable à ceux obtenus pour les modèles d' ACC-EANN et d'ACC-GAM. Il est également important de souligner qu'en termes de BIAS, le ACC-RFR surestime les quantiles de crue alors que le RFR les sous-estime. Cependant, lorsque BIASr est utilisé, tous les modèles sous-estiment les quantiles de crue.

Une autre expérience a été menée pour déterminer l'importance des variables prédictives individuelles pour l'estimation des quantiles de crue. La superficie du bassin est la variable physiographique la plus importante, suivie respectivement de la moyenne annuelle des précipitations totales (AMP) et de la moyenne annuelle des degrés-jours en-dessus de 0° C (AMD). La pente moyenne du bassin (MBS) est la quatrième. La fraction de la superficie couverte par les lacs (FAL) est la moins importante de toutes les variables physio-météorologiques.

## **5. CONCLUSIONS AND RECOMMENDATIONS**

La RF est couramment utilisé dans la classification des gènes, les banques, la médecine et le commerce électronique, etc. Cependant, jusqu'à présent, son application dans le domaine de l'hydrologie et en particulier dans l'AFR est encore limitée. La RFR est une approche non linéaire et non paramétrique qui a montré des meilleures performances dans d'autres domaines. Le but de cette étude est d'introduire les RFR dans l'AFR, puis d'appliquer les RFR aux voisinages délimités par l'ACC.

Le nombre d'arbres dans la RF pour cette étude a été fixé à 30. La comparaison avec d'autres modèles indique que, malgré que l'ACC-RFR présente des faibles valeurs de NASH, il est plus précis que les autres modèles en terme des RMSE et

RMSEr. Les résultats indiquent en outre que la RF donne des estimations plus précises lorsqu'elle est utilisée avec l'ACC.

Dans cette étude nous avons introduit l'approche de RF à l'AFR. L'utilisation de l'« Extremely Randomized Trees » et d'autres variantes de RF dans l'AFR devrait également être étudiée à l'avenir. Les activités de recherche futures devraient également porter sur l'utilisation des RF en combinaison avec d'autres techniques de délimitation telles que l'approche par région d'influence, les voisinages basées sur la notion statistique des fonctions de profondeur, etc. L'approche RF peut également être testée avec d'autres variables comme les débits d'étiages ou les sédiments en suspension.

## 6. RÉFÉRENCES

1. Breiman, L., 2001. Random Forests. *Machine Learning*, 45(1): 5-32. DOI:10.1023/a:1010933404324
2. Chebana, F., Ouarda, T.B., 2008. Depth and homogeneity in regional flood frequency analysis. *Water Resources Research*, 44(11). DOI: 10.1029/2007WR006771
3. Chen, J., Li, M., Wang, W., 2012. Statistical uncertainty estimation using random forests and its application to drought forecast. *Mathematical Problems in Engineering*, 2012. DOI: 10.1155/2012/915053
4. Chokmani, K., Ouarda, T.B.M.J., 2004. Physiographical space-based kriging for regional flood frequency estimation at ungauged sites. *Water Resources Research*, 40(12). DOI:10.1029/2003wr002983
5. Chebana, F., Charron, C., Ouarda, T.B.M.J., Martel, B., 2014. Regional Frequency Analysis at Ungauged Sites with the Generalized Additive Model. *Journal of Hydrometeorology*, 15(6): 2418-2428. DOI:10.1175/jhm-d-14-0060.1
6. Monira, S.S., Faisal, Z.M., Hirose, H., 2010. Comparison of artificially intelligent methods in short term rainfall forecast, 2010 13th International Conference on Computer and Information Technology (ICCIT), pp. 39-44. DOI:10.1109/iccitechn.2010.5723826
7. Nguyen, T.-T., Huu, Q.N., Li, M.J., 2015. Forecasting Time Series Water Levels on Mekong River Using Machine Learning Models. Seventh International Conference on Knowledge and Systems Engineering (KSE), pp. 292-297. DOI:10.1109/kse.2015.53.
8. Ouarda, T. et al., 2007. Regional flood frequency estimation at ungauged sites in the Balsas River Basin, Mexico, AGU Spring Meeting Abstracts.
9. Ouarda, T.B.M.J., 2013. Hydrological Frequency Analysis, Regional, *Encyclopedia of Environmetrics*. DOI:10.1002/9780470057339.vnn043
10. Schaefli, B., Gupta, H.V., 2007. Do Nash values have value? *Hydrological Processes*, 21(15): 2075-2080. DOI:10.1002/hyp.6825
11. Shu, C., Ouarda, T.B.M.J., 2007. Flood frequency analysis at ungauged sites using artificial neural networks in canonical correlation analysis physiographic space. *Water Resources Research*, 43(7). DOI:10.1029/2006wr005142
12. Stefanidis, S., Stathis, D., 2013. Assessment of flood hazard based on natural and anthropogenic factors using analytic hierarchy process (AHP). *Natural Hazards*, 68(2): 569-585. DOI:10.1007/s11069-013-0639-5
13. Taksande, A.A., Mohod, P., 2015. Applications of data mining in weather forecasting using frequent pattern growth algorithm. *IJSR*, 4(6): 3048-51.
14. Wallemacq, P.G.-S., Debarati & McClean, Denis & , CRED & , UNISDR, 2015. *The Human Cost of Weather Related Disasters - 1995 - 2015*. DOI:10.13140/RG.2.2.17677.33769
15. Wang, Z. et al., 2015. Flood hazard risk assessment model based on random forest. *Journal of Hydrology*, 527: 1130-1141. DOI:10.1016/j.jhydrol.2015.06.008



## **CHAPITRE 2 : ARTICLE**

# Regional Hydrological Frequency Analysis at Ungauged Sites with Random Forest Regression

Shitanshu Desai<sup>1,\*</sup>, Taha B. M. J. Ouarda<sup>1</sup>

<sup>1</sup> Canada Research Chair in Statistical Hydro-climatology,  
INRS-ETE,  
490 De la Couronne, Québec (QC),  
Canada. G1K 9A9.

\* Corresponding Author: [shitanshu.desai@ete.inrs.ca](mailto:shitanshu.desai@ete.inrs.ca); [situdesai@gmail.com](mailto:situdesai@gmail.com)

September 2018

## **ABSTRACT:**

Flood quantile estimation at sites with little or no data is important for the adequate planning and management of water resources. Regional Hydrological Frequency Analysis (RFA) deals with the estimation of hydrological variables at ungauged sites. Random Forest (RF) is an ensemble learning technique which uses multiple Classification and Regression Trees (CART) for classification, regression, and other tasks. The RF technique is gaining popularity in a number of fields because of its powerful non-linear and non-parametric nature. In the present study, we investigate the use of Random Forest Regression (RFR) in the estimation step of RFA based on a case study represented by data collected from 151 hydrometric stations from the province of Quebec, Canada. RFR is applied to the whole data set and to homogeneous regions of stations delineated by canonical correlation analysis (CCA). Using the Out-of-bag error rate feature of RF, the optimal number of trees for the dataset is calculated. The results of the application of the CCA based RFR model (CCA-RFR) are compared to results obtained with a number of other linear and non-linear RFA models. CCA-RFR leads to the best performance in terms of root mean squared error. The use of CCA to delineate neighborhoods improves considerably the performance of RFR. RFR is found to be simple to apply and more efficient than more complex models such as Artificial Neural Network-based models.

## **Keywords:**

Random Forest Regression, Canonical Correlation Analysis, Regional Flood Frequency Analysis, Ungauged basin, Machine Learning.

## **LIST OF ABBREVIATIONS**

RFA : Regiona Frequency Analysis

CCA : Canonical Correlation Analysis

ANN : Artificial Neural Network

GAM : Generalized Additive Model

RF : Random Forest

RFR : Random Forest Regression

CART : Classification and Regression trees

CCA-RFR : Random Forest Regression with Canonical Correlation Analysis

OOB : out-of-bag

SANN : Single Artificial Neural Network

EANN : Ensemble Artificial Neural Network

CCA-SANN : Single Artificial Neural Network with Canonical Correlation Analysis

CCA-EANN : Ensemble Artificial Neural Network with Canonical Correlation Analysis

CCA-GAM : Generalized Additive Model with Canonical Correlation Analysis

RMSE : Root mean squared error

NASH : Nash Sutcliffe model efficiency criterion

RMSEr : Relative Root Mean Squared Error

BIAS : Mean Bias

BIASr : Relative Mean Bias

k-fold CV : K-fold Cross Validation

Area : Basin Area

MBS : Basin Mean Slope

FAL : Fraction of Basin Area Occupied by Lakes

AMP : Annual Mean Total Precipitation

AMD : Annual Mean Degree-days above 0°

q100, q50 and q10 : Specific Flood quantiles corresponding to 100, 50 and 10 year return periods

MDI : Mean Decrease in Impurity



## 1. Introduction

Floods represent one of the most commonly occurring natural disasters ([Stefanidis and Stathis, 2013](#)). Floods cause significant environmental, economic and social damages. In spite of all flood protection measures being taken, from 1990 to 2013, floods have caused damages of about 600 billion US dollars and close to 7 million deaths worldwide ([Wang et al., 2015](#)). Thus, it is of the utmost importance to adequately predict the characteristics of such events at all sites of interest.

However, hydrological information may not be available at certain sites of interest. At these “ ungauged sites”, Regional Frequency Analysis (RFA) can be used to develop estimates of flood characteristics. RFA allows transfer of information from gauged sites to the ungauged site of interest. RFA usually consists of two main steps. The first step is the delineation of homogeneous regions. In this step, sites that are similar according to some homogeneity criteria are grouped together. The rationale here is that as the sites within a given homogenous region are similar, information can reasonably be transferred from gauged to ungauged sites. The second step is the application of a regional estimation model within each delineated region ([Ouarda, 2013](#)). The regional estimation models are then trained to establish functional relationships between physio-meteorological basin characteristics and flow characteristics at ungauged basins.

Delineation can be done on the basis of geographical proximity, but that does not guarantee that such regions are homogenous in regards to their hydrologic response. In contrast, “Site focused” regionalization techniques (also called neighborhood-based techniques) have received much attention due to their effectiveness. In “Site focused” techniques, each site has a prospective set of catchments which form a homogenous region for that particular site. One such technique, Canonical Correlation Analysis (CCA) has been used for delineating homogenous regions in a number of studies ([Chebana et al., 2014](#); [Chokmani and Ouarda, 2004](#); [Ouarda et al., 2001](#); [Shu and Ouarda, 2007](#)). In the present study, CCA is used to delineate homogenous regions as [Ouarda et al. \(2007\)](#) indicated that it leads to superior performances.

Among the large number of RFA estimation methods proposed in the literature, linear models and their variants are commonly adopted because of their simplicity and the speed in which they can be trained as well as deployed. However, hydrological systems are characterized by complex processes and it is unrealistic to assume a linear relationship between physio-meteorological basin characteristics and flow characteristics. [Sivakumar and Singh \(2012\)](#) showed that the relationship between these variables is characterized by dominant non-linear relationships. [Pandey and Nguyen \(1999\)](#) and [Grover et al. \(2002\)](#) showed that non-linear regression models provide better performances for RFA.

Several non-linear techniques have been proposed in the literature. An Artificial Neural Network (ANN), a non-linear and a non-parametric approach modelled on the neurons present in the human brain, was used for solving several hydrological

problems such as regional flood frequency analysis, streamflow forecasting, rainfall-runoff modelling, flood forecasting, etc. ([Aziz et al., 2014](#); [Chokmani et al., 2008](#); [Huo et al., 2012](#); [Kumar et al., 2015](#); [Shu and Ouarda, 2007](#); [Tiwari and Chatterjee, 2018](#)). Generalized Additive Models (GAM) due to their considerable flexibility, are used in regional flood frequency analysis, water quality estimation, river discharge modeling, etc. ([Chebana et al., 2014](#); [Iddrisu et al., 2017](#); [Morton and Henderson, 2008](#); [Rahman et al., 2017](#)). Other non-linear approaches used RFA include Projection Pursuit Regression ([Durocher et al. \(2015\)](#)), Non-Linear CCA [Ouali et al. \(2015\)](#), and Adaptive Neuro-Fuzzy Inference Systems (ANFIS) (Shu and Ouarda, 2008).

Random Forest (RF), first proposed by [Breiman \(2001\)](#), is one such non-linear and non-parametric technique. It is a popular technique for classification, regression, variable selection, outlier detection and variable importance. When random forest is used for the purpose of function approximation or regression, it is called Random Forest Regression (RFR) or Regression Forests. In RFR, from a given set of data, multiple samples are randomly drawn and Classification and Regressions Trees (CART) are built. Eventually, the results of all such trees are combined and an estimate of target variables is obtained by averaging the outputs of individual trees.

A number of studies have been conducted in the field of hydrology using RFs. [Chen et al. \(2012\)](#) used RF to build a drought forecast model. [Nguyen et al. \(2015b\)](#) used RF to forecast daily water levels. [Monira et al. \(2010b\)](#) and [Taksande and Mohod \(2015\)](#) respectively used RF for daily and monthly rainfall forecasting. [Wang et al. \(2015\)](#) developed a flood hazard risk assessment model based on RF. RF represents a good alternative to Support Vector Machines ([Meyer et al., 2003](#); [Verikas et al., 2001](#)) and possesses a number of advantages including a reasonable amount of tolerance towards noise and outliers, high accuracy in forecasting and no overfitting problems.

The aim of the present study is to introduce the RF technique for regional flood quantile estimation. RFR is used to establish non-linear relationships between physio-meteorological basin characteristics and flow characteristics, and to estimate flood characteristics at ungauged sites. RFR is also applied to hydrological neighborhoods derived using CCA (CCA-RFR) for flood quantile estimation. A comparative analysis is carried out with several other approaches based on the application to a case study of data derived from the Province of Quebec, Canada.

The paper is organized as follows. In section 2, the theoretical background of RFR and CCA is presented along with the evaluation procedure and brief information about the models to be compared. The case study is presented in section 3 and the results are presented and discussed in section 4. Finally, the conclusions and recommendations for further research are presented in section 5.

## 2.1. Random Forest Regression

### 2.1.1. RFR Principle

Random Forest is an ensemble learning technique proposed by [Breiman \(2001\)](#). RFR is one of the most accurate general-purpose learning algorithms. Random Forest has been shown to give a very good performance while using few computational resources. RFR exhibits great performance improvement over single tree algorithms like CART. It is fast and has error rates comparable to more traditional and resource intensive algorithms.

In Random forest for regression, the tree predictors  $h(x, \theta_k)$ ,  $k = 1 \dots K$  take on numerical values depending on the random vectors  $\{\theta_k\}$  ([Breiman, 2001](#)). It is important to note that  $\{\theta_k\}$  are identically distributed and independent random vectors. The training data is randomly and independently drawn from a joint distribution of  $(X, Y)$ , where the random vector  $X$  is the observed input and the random vector  $Y$  is the expected numerical output. Individual trees are grown using the Classification and Regression Trees (CART) algorithm. Below is the algorithm for Random forest for regression as presented in [Trevor et al. \(2009\)](#).

- 
- (1) For  $b = 1$  to  $B$ :
    - (a) Draw a bootstrap sample  $Z^*$  of size  $N$  from training data.
    - (b) Grow a random-forest tree  $T_b$  to the bootstrapped data by recursively repeating the following steps for each terminal node of the tree, until the minimum node size  $n_{min}$  is reached.
      - (i) Select  $m$  variables at random from  $p$  variables.
      - (ii) Pick the best variable/split-point among the  $m$ .
      - (iii) Split the node into two daughter nodes.
  - (2) Output the ensemble of trees  $\{T_b\}_1^B$
- To make a prediction at a new point  $x$ :

$$\hat{f}_{rf} = \frac{1}{B} \sum_{b=1}^B T_b(x)$$

---

RFR possesses two important features, out-of-bag error rate, and variable importance. Generally, we use about two third of the data in a bootstrap sample and the rest one third are left out. These are known as out-of-bag (OOB) samples. The error estimated on these left out samples is known as OOB-error rate. OOB error rate can be used for validation purposes as well as for the calculation of the optimum number of trees required. Variable importance is a measure of which predictors are most useful for predicting the response variable. Variable importance can be computed using RFR by recording improvements, at each node in every tree in the forest.

Another advantage of using RFR is that it possesses an ‘acceptable’ tolerance to noise and outliers, as the input training sets are drawn by random bootstrap sampling, and as the nodes to be split are selected randomly. Also, as there is no correlation between individual trees and as each tree is allowed to grow to its maximum size, there is no overfitting of data. Consequently, the only parameter to be tuned is the number of trees or estimators.

### 2.1.2. Classification and Regression Trees (CART)

CART decision tree is a binary recursion partitioning scheme which is capable of processing continuous and nominal attributes for regression and classification. In the present study, we use CART trees for regression. Regression trees are a nonparametric regression method that approximates real-valued functions. A regression tree is built using binary partitioning, where each node is iteratively split into two partitions or branches. Initially, all input variables are grouped into the same partition. Then mean squared error (MSE) is calculated and a split decision is taken. The split decision is taken based on Greedy minimization. The split which minimizes the MSE is selected and further that node is split into two off-springs. The splitting rule is then applied to each of the new offsprings. Each tree is grown to the largest possible extent which aids in better regression accuracy.

### 2.2. CCA approach in RFA

This section contains a brief discussion about CCA and its connection to the delineation step of RFA. Let  $X = \{X_1, X_2 \dots X_r\}$  be a random variable containing basin meteorological and physiographical variables, for eg. basin area, etc. and  $Y = \{Y_1, Y_2 \dots Y_r\}$  be a random variable containing basin hydrological variables like flood quantiles.

Consider linear combinations  $V$  and  $W$  of the variables  $X$  and  $Y$ :

$$V = a_1X_1 + a_2X_2 + \dots + a_rX_r = a'X \quad (1)$$

$$W = b_1Y_1 + b_2Y_2 + \dots + b_rY_r = b'Y \quad (2)$$

where  $a'$  and  $b'$  are transposes of vector  $a$  and  $b$  respectively. CCA enables identifying vectors  $a$  and  $b$  such that  $corr(V, W)$  is maximum with vectors  $V$  and  $W$  having unit variances. For each basin  $B_k$ , where  $k = 1, 2 \dots K$  from the set  $B$  of basins,  $v_{i,k}$  and  $w_{i,k}$  are corresponding values of  $V_i$  and  $W_i$ . We have the values of vector  $v_0$  and our aim is to estimate the unknown vector  $w_0$ , where  $v_0$  and  $w_0$  represent the canonical scores of physio-meteorological and hydrological variables respectively.

The approximation of the  $w_0$  vector can be obtained from a  $100(1 - \alpha)\%$  confidence interval about  $\lambda v_0$  by constituting all the realizations  $w$  of  $W$  where:

$$(w - \lambda v_0)'(I_p - \lambda^2)^{-1}(w - \lambda v_0) \leq \chi_{\alpha,p}^2, \quad (3)$$

is conditional on  $\chi_{\alpha,p}^2$  being  $P(\chi^2 \leq \chi_{\alpha,p}^2) = 1 - \alpha$ . For more detailed information concerning the algorithm, the reader is referred to ([Ouarda et al., 2001](#)).

### 2.3. Selection of Methods for Comparison

The RFR and CCA-RFR models are used to estimate the 100, 50 and 10-year flood quantiles. To evaluate the relative performances of these two approaches, they are compared to the following models:

- Canonical Correlation Analysis-Multiple linear regression model (CCA-MLR) ([Ouarda et al., 2001](#)). After selecting the optimal hydrological neighborhoods for each site using CCA analysis, multiple regression is used for regional flood estimation.
- Single Artificial Neural Network (SANN) ([Shu and Burn, 2004](#)). A single ANN is used to identify a functional relationship between physio-meteorological variables and flood quantiles.
- Ensemble ANN (EANN) ([Shu and Burn, 2004](#)). An ANN ensemble is created by bagging several single ANNs. This helps in improving the generalization ability of the SANN model. The final output is generated by taking the mean of the outputs of individual ANNs.
- Canonical Kriging Model (CCA-Kriging) ([Chokmani and Ouarda, 2004](#)). The physiographical space defined by CCA is used by the Kriging model to obtain regional flood estimates by interpolating data over that physiographic space. This method was shown to lead to comparable results to the traditional CCA model but is computationally less complicated.
- Single Artificial Neural Network in CCA physiographical space (CCA-SANN) ([Shu and Ouarda, 2007](#)). CCA is used to form the canonical physiographical space and then single ANN is applied to the data to form functional relationships between physiographical variables and flood quantiles.
- Ensemble ANN in CCA physiographical space (CCA-EANN) ([Shu and Ouarda, 2007](#)). In the CCA-EANN model, each component uses the same configuration as a Single ANN but the CCA-EANN is trained on bootstrapped sample data and the results are averaged out.
- Generalized Additive Model in conjunction with CCA (CCA-GAM) ([Chebana et al., 2014](#)). In the CCA-GAM approach, firstly backward stepwise selection is used to select the variables to be used in the model. Then GAM is applied to the neighborhoods delineated by CCA.

### 2.4. Evaluation Metrics

The following metrics are used to assess the quality of our regional flood analysis models. They are NASH (Nash Criterion), RMSE (Root mean squared error), RMSEr (Relative Root Mean Squared Error), BIAS (Mean Bias) and BIASr (Relative Mean Bias).

$$NASH = 1 - \frac{\sum_{i=1}^n (o_i - s_i)^2}{\sum_{i=1}^n (o_i - \bar{o})^2} \quad (4)$$

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (o_i - s_i)^2} \quad (5)$$

$$RMSEr = \sqrt{\frac{1}{n} \sum_{i=1}^n \left( \frac{o_i - s_i}{o_i} \right)^2} \quad (6)$$

$$BIAS = \frac{1}{n} \sum_{i=1}^n (o_i - s_i) \quad (7)$$

$$BIASr = \frac{1}{n} \sum_{i=1}^n \left( \frac{o_i - s_i}{o_i} \right) \quad (8)$$

where,  $o_i$  is the observed value at site  $i$ ,  $s_i$  is the simulated value using the model for site  $i$ ,  $\bar{o}$  is the mean of observed at-site values and  $n$  is the number of sites.

## 2.5. Evaluation Procedure

K-fold Cross Validation (k-fold CV) is used as the model validation technique in this work. In k-fold CV the data is split into  $k$  small and equal sets. A model is trained using  $k - 1$  folds as training data and then the model is validated using the remaining data. The performance thus reported by k-fold CV is the mean of the values computed in the loop.

The reason for using k-fold CV in the present study is that models trained with k-fold CV have lower variance than models trained with the jackknife validation procedure. In jackknife validation, there is more overlap between training folds as only one sample is omitted which means that almost the entire dataset is used for training. While in k-fold CV there is less overlap between training folds and thus it leads to smaller variability. Therefore, results obtained with jackknife might be better but the results obtained using k-fold CV are more robust.

## 3. Case Study

The dataset used in the present study consists of 151 hydrometric stations located in the southern part of the province of Quebec (between 45° and 55°N), Canada. The stations are operated by the Ministry of Environment of Quebec. The adopted dataset has been used in a number of previous RFA studies ([Chebana and Ouarda, 2008](#); [Chokmani and Ouarda, 2004](#); [Shu and Ouarda, 2007](#)) making it convenient for comparison of the results with those obtained with other methodologies.

On the basis of the work of [Chokmani and Ouarda \(2004\) with the same database](#), a total of five physio-meteorological variables are selected, of which three are physiographical and two are meteorological variables. These variables are the basin area (Area), the mean basin slope (MBS), the fraction of basin area occupied by lakes (FAL), the annual mean total precipitation (AMP) and the annual mean degree-days above 0° (AMD), respectively. A number of statistics of these data, like the minimum, mean, maximum and standard deviation are presented in table 1.

The database compiled by [Kouider et al. \(2002\)](#) is used to extract at-site flood estimates for all of the 151 gauging stations in the study area. The most appropriate statistical distribution is used to get flood quantile estimates for each site by fitting the distribution to observed flood data. To avoid negative scale effects, specific quantiles (quantiles divided by basin areas) are used. The 100-year, 50-year, and 10-year quantiles (q100, q50, and q10 respectively) are the three specific flood quantiles used in the present study.

The reader is directed to ([Shu and Ouarda, 2007](#)) for more details concerning the dataset, such as scatter plots of basins in canonical space and geographical location of stations, to avoid redundancy. According to the recommendations of [Shu and Ouarda \(2007\)](#), the logarithmic transformation is applied to the variables q10, q50, q100, Area, MBS, AMP and AMD and a square root transformation is applied to FAL.

#### 4. Results

In the present study, Scikit-learn module of Python is used to obtain the results ([Pedregosa et al., 2011](#)). In RF the size of the dataset, the number of trees (n\_estimators) and the number of variables at each split have a huge impact on the error rate. According to [Breiman \(2001\)](#), the number of variables at each split should be taken as the square root of the total number of variables, i.e. 2 in this study. As the size of the dataset is not a tunable parameter, only the number of trees is tuned in this study.

Figure 1 illustrates that the OOB error rate decreases as the number of trees increases. At around 30 trees the value levels off and there is almost no improvement after this point by increasing the number of trees. Therefore, the number of trees is fixed at 30 for the present study. It is also important to note that all the trees were allowed to grow to the maximum extent without pruning.

The results of the application of the two models RFR and CCA-RFR along with the models described in Section 2.3 to the dataset described in Section 3 are illustrated in Table 2. The bold font describes the best approach for that particular

flood quantile and the particular evaluation metric. Results indicate that CCA-RFR either outperforms or is comparable to other models in all the metrics except the NASH criterion. Also, CCA-RFR outperforms RFR in every metric other than NASH.

Figure 2 illustrates the relative errors associated with quantiles q50 estimated using RFR and CCA-RFR. Figure 2 indicates that CCA-RFR performs better than RFR for large basins, while RFR outperforms CCA-RFR for very small basins. These smaller basins are associated with larger specific quantiles. Therefore we can attribute the low NASH scores associated to CCA-RFR to these smaller sites. Similarly, according to [McCuen et al. \(2006\)](#), the NASH criterion is sensitive to a number of factors including sample size and outliers. In CCA-RFR, as only the stations in the hydrological neighborhoods are considered for the prediction and training, the sample size is considerably smaller than the complete original dataset. Also, the NASH criterion is heavily influenced by the model used ([Schaefli and Gupta, 2007](#)). RFR provides a reasonable tolerance to outliers which can be seen in the RFR NASH values. However, as we use just the neighborhoods for CCA-RFR, the sample size is small and thus outliers have more effect than in the basic RFR model which leads to lower NASH values.

Although we have low values for the NASH criterion for both RFR and CCA-RFR in comparison to other models, we can observe that CCA-RFR leads to the best RMSE and RMSEr values among all the models studied in this work. RMSE provides an evaluation of prediction accuracy in the absolute scale while RMSEr does the same in relative terms. CCA based RFR provides better generalization ability than the basic RFR model. As RFRs are nonparametric data-driven approaches, they have limited scope for extrapolation beyond the observed data. Therefore, the combination of RFR along with CCA, a parametric model helps the performance of RFR. Consequently, even though the NASH value for CCA-RFR is lower than other models the prediction accuracy is not compromised and is rather improved.

The BIAS and BIASr are evaluation criteria used to determine whether the model overestimates or underestimates the various quantiles. In general, CCA-RFR has the lowest BIAS of all the models considered and BIASr is also comparable with CCA-EANN and CCA-GAM which have the best BIASr value. It is also important to point out that, in terms of BIAS, CCA-RFR overestimates flood quantiles while RFR underestimates them. However, when BIASr is used, all the models underestimate the flood quantiles.

Overall, it can be concluded that applying RFR to CCA delineated neighborhoods improves the results in comparison to RFR applied to the whole set of stations. This is consistent with the results of previous studies, such as [Chokmani and](#)



[Ouarda \(2004\)](#) and [Shu and Ouarda \(2007\)](#), which indicated that applying other estimation techniques to CCA delineated neighborhoods leads to better performances for the estimation of flood quantiles than their application to the whole set of stations in the database.

The scatter plots of regional estimates using RFR and CCA-RFR are shown in Figure 3 and Figure 4, respectively. As would be expected, we observe that the estimation error and bias are positively correlated with the return period. With the increase in return periods, bias and estimation error increase simultaneously. Also, the low NASH scores can be explained by high variation as seen in Figure 4. It is clear from the results that all models underestimate flood quantiles at sites with higher specific quantiles. These sites can be associated with smaller basins which have large specific quantiles ([Shu and Ouarda, 2007](#)).

An additional experiment is conducted to identify the importance of individual predictor variables for flood quantile estimation. In the python implementation of RFR, “Mean decrease in Impurity (MDI)” or “Gini importance” is used to calculate the importance of each variable on the accuracy of the model. MDI is defined as “total decrease in node impurity averaged over all the trees. Node impurity is weighted by the probability of reaching that node (which is approximated by the proportion of sample reaching that node)” ([Brieman et al., 1984](#)). The results are illustrated in Table 3. Basin Area (Area) is shown to be by far the most important physio-meteorological variable. Annual mean total precipitation (AMP) and Annual mean degree days over 0° C (AMD) are distant second and third, respectively. Mean Basin Slope (MBS) is fourth while the Fraction of Area covered by lakes (FAL) is the least important of all physio-meteorological variables.

## 5. Conclusions

RF has been commonly used in gene classification, banking, medicine, and E-commerce. However, so far it has not found much application in the field of hydrology and especially in RFA. Most common studies in RFA establish linear relationships between physio-meteorological variables and flood quantiles. However, these models do not generally explain the complex relationships between the response variable and the explanatory variables. Random forest, a non-linear and a non-parametric data-driven approach, is one such technique which has shown good performances in other fields in explaining such complex relationships. The purpose of this study is to first introduce RFR in RFA and then apply RFR to neighborhoods delineated by CCA.

The number of trees in the RF for this study was fixed at 30. Also, all the trees were allowed to grow to their maximum potential without pruning. The comparison

with other models indicates that, although CCA-RFR has a lower NASH score, it is more accurate than the other models. RFR is particularly more advantageous because of its low computational cost and high prediction quality. The results further indicate that the Random Forest, when used in conjunction with CCA, provides more robust and accurate results.

The research presented in this work is based on the introduction of the RF approach to RFA. The use of Extremely Randomized Trees and other variants of RF in RFA should also be attempted in the future. Future research activities should also focus on the use of RF in conjunction with other delineation techniques such as the Region of Influence approach, statistical depth functions, or projection pursuit regression. The effectiveness of the same techniques should also be investigated in the future using other data sets from different parts of the world and for other hydrological variables such as low flows or suspended sediments.

**6. Acknowledgments.** Financial support for this study was graciously provided by the Natural Sciences and Engineering Research Council (NSERC) of Canada. The authors would like to thank the Ministry of Sustainable Development, Environment, and Fight Against Climate Change of the Province of Quebec (MDDELCC) for the employed datasets.

## Bibliography

1. Aziz, K., Rahman, A., Fang, G., Shrestha, S., 2014. Application of artificial neural networks in regional flood frequency analysis: a case study for Australia. *Stochastic Environmental Research and Risk Assessment*, 28(3): 541-554. DOI:10.1007/s00477-013-0771-5
2. Breiman, L., 2001. Random Forests. *Machine Learning*, 45(1): 5-32. DOI:10.1023/a:1010933404324
3. Breiman, L., Friedman, J., Olshen, R., Stone, C., 1984. *Classification and Regression Trees*. Wadsworth. Inc, Pacific Grove, CA.
4. Chebana, F., Charron, C., Ouarda, T.B.M.J., Martel, B., 2014. Regional Frequency Analysis at Ungauged Sites with the Generalized Additive Model. *Journal of Hydrometeorology*, 15(6): 2418-2428. DOI:10.1175/jhm-d-14-0060.1
5. Chebana, F., Ouarda, T.B., 2008. Depth and homogeneity in regional flood frequency analysis. *Water resources research*, 44(11).
6. Chen, J., Li, M., Wang, W., 2012. Statistical uncertainty estimation using random forests and its application to drought forecast. *Mathematical Problems in Engineering*, 2012.
7. Chokmani, K., Ouarda, T.B., Hamilton, S., Ghedira, M.H., Gingras, H., 2008. Comparison of ice-affected streamflow estimates computed using artificial neural networks and multiple regression techniques. *Journal of Hydrology*, 349(3-4): 383-396. DOI:10.1016/j.jhydrol.2007.11.024
8. Chokmani, K., Ouarda, T.B.M.J., 2004. Physiographical space-based kriging for regional flood frequency estimation at ungauged sites. *Water Resources Research*, 40(12). DOI:10.1029/2003wr002983
9. Durocher, M., Chebana, F., Ouarda, T.B.M.J., 2015. A Nonlinear Approach to Regional Flood Frequency Analysis Using Projection Pursuit Regression. *Journal of Hydrometeorology*, 16(4): 1561-1574. DOI:10.1175/jhm-d-14-0227.1
10. Grover, P.L., Burn, D.H., Cunderlik, J.M., 2002. A comparison of index flood estimation procedures for ungauged catchments. *Canadian Journal of Civil Engineering*, 29(5): 734-741. DOI:10.1139/l02-065
11. Huo, Z. et al., 2012. Integrated neural networks for monthly river flow estimation in arid inland basin of Northwest China. *Journal of Hydrology*, 420-421: 159-170. DOI:10.1016/j.jhydrol.2011.11.054
12. Iddrisu, W.A., Nokoe, K.S., Luguterah, A., Antwi, E.O., 2017. Generalized Additive Mixed Modelling of River Discharge in the Black Volta River. *Open Journal of Statistics*, 07(04): 621-632. DOI:10.4236/ojs.2017.74043
13. Kouider, A., Gingras, H., Ouarda, T., Ristic-Rudolf, Z., Bobée, B., 2002. Analyse fréquentielle locale et régionale et cartographie des crues au Québec. Rep. R-627-el.
14. Kumar, R., Goel, N.K., Chatterjee, C., Nayak, P.C., 2015. Regional Flood Frequency Analysis using Soft Computing Techniques. *Water Resources Management*, 29(6): 1965-1978. DOI:10.1007/s11269-015-0922-1

15. McCuen, R.H., Knight, Z., Cutter, A.G., 2006. Evaluation of the Nash–Sutcliffe efficiency index. *Journal of Hydrologic Engineering*, 11(6): 597-602. DOI:doi:10.1061/(ASCE)1084-0699(2006)11:6(597)
16. Meyer, D., Leisch, F., Hornik, K., 2003. The support vector machine under test. *Neurocomputing*, 55(1-2): 169-186. DOI:10.1016/s0925-2312(03)00431-4
17. Monira, S.S., Faisal, Z.M., Hirose, H., 2010a. Comparison of artificially intelligent methods in short term rainfall forecast, 2010 13th International Conference on Computer and Information Technology (ICCIT), pp. 39-44. DOI:10.1109/iccitechn.2010.5723826
18. Monira, S.S., Faisal, Z.M., Hirose, H., 2010b. Comparison of artificially intelligent methods in short term rainfall forecast, *Computer and Information Technology (ICCIT)*, 2010 13th International Conference on. IEEE, pp. 39-44. DOI:10.1109/ICCITECHN.2010.5723826
19. Morton, R., Henderson, B.L., 2008. Estimation of nonlinear trends in water quality: An improved approach using generalized additive models. *Water Resources Research*, 44(7). DOI:10.1029/2007wr006191
20. Nguyen, T.-T., Huu, Q.N., Li, M.J., 2015a. Forecasting Time Series Water Levels on Mekong River Using Machine Learning Models, 2015 Seventh International Conference on Knowledge and Systems Engineering (KSE), pp. 292-297. DOI:10.1109/kse.2015.53
21. Nguyen, T.-T., Huu, Q.N., Li, M.J., 2015b. Forecasting time series water levels on Mekong river using machine learning models, *Knowledge and Systems Engineering (KSE)*, 2015 Seventh International Conference on. IEEE, pp. 292-297. DOI:10.1109/KSE.2015.53
22. Ouali, D., Chebana, F., Ouarda, T.B.M.J., 2015. Non-linear canonical correlation analysis in regional frequency analysis. *Stochastic Environmental Research and Risk Assessment*, 30(2): 449-462. DOI:10.1007/s00477-015-1092-7
23. Ouarda, T. et al., 2007. Regional flood frequency estimation at ungauged sites in the Balsas River Basin, Mexico, AGU Spring Meeting Abstracts.
24. Ouarda, T.B.M.J., 2013. Hydrological Frequency Analysis, *Regional, Encyclopedia of Environmetrics*. DOI:10.1002/9780470057339.vnn043
25. 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(1-4): 157-173. DOI:10.1016/s0022-1694(01)00488-7
26. Pandey, G., Nguyen, V.-T.-V., 1999. A comparative study of regression based methods in regional flood frequency analysis. *Journal of Hydrology*, 225(1-2): 92-101. DOI:10.1016/S0022-1694(99)00135-3
27. Pedregosa, F. et al., 2011. Scikit-learn: Machine learning in Python. *Journal of machine learning research*, 12(Oct): 2825-2830.
28. Rahman, A., Charron, C., Ouarda, T.B.M.J., Chebana, F., 2017. Development of regional flood frequency analysis techniques using generalized additive models for Australia. *Stochastic Environmental Research and Risk Assessment*, 32(1): 123-139. DOI:10.1007/s00477-017-1384-1

29. Schaeffli, B., Gupta, H.V., 2007. Do Nash values have value? *Hydrological Processes*, 21(15): 2075-2080. DOI:10.1002/hyp.6825
30. Shu, C., Burn, D.H., 2004. Artificial neural network ensembles and their application in pooled flood frequency analysis. *Water Resources Research*, 40(9). DOI:10.1029/2003wr002816
31. Shu, C., Ouarda, T.B.M.J., 2007. Flood frequency analysis at ungauged sites using artificial neural networks in canonical correlation analysis physiographic space. *Water Resources Research*, 43(7). DOI:10.1029/2006wr005142
32. Sivakumar, B., Singh, V.P., 2012. Hydrologic system complexity and nonlinear dynamic concepts for a catchment classification framework. *Hydrology and Earth System Sciences*, 16(11): 4119-4131. DOI:10.5194/hess-16-4119-2012
33. Stefanidis, S., Stathis, D., 2013. Assessment of flood hazard based on natural and anthropogenic factors using analytic hierarchy process (AHP). *Natural Hazards*, 68(2): 569-585. DOI:10.1007/s11069-013-0639-5
34. Taksande, A.A., Mohod, P., 2015. Applications of data mining in weather forecasting using frequent pattern growth algorithm. *IJSR*, 4(6): 3048-51.
35. Tiwari, M.K., Chatterjee, C., 2018. Flood Forecasting and Uncertainty Assessment Using Wavelet- and Bootstrap-Based Neural Networks, *Handbook of Research on Predictive Modeling and Optimization Methods in Science and Engineering. Advances in Computational Intelligence and Robotics*, pp. 74-93. DOI:10.4018/978-1-5225-4766-2.ch004
36. Trevor, H., Robert, T., JH, F., 2009. *The elements of statistical learning: data mining, inference, and prediction*. New York, NY: Springer.
37. Verikas, A., Gelzinis, A., Malmqvist, K., 2001. Using unlabelled data to train a multilayer perceptron. *Neural Processing Letters*, 14(3): 179-201. DOI:10.1023/A:1012707515770
38. Wallemacq, P.G.-S., Debarati & McClean, Denis & , UNISDR, 2015. *The Human Cost of Weather Related Disasters - 1995 - 2015*. DOI:10.13140/RG.2.2.17677.33769
39. Wang, Z. et al., 2015. Flood hazard risk assessment model based on random forest. *Journal of Hydrology*, 527: 1130-1141. DOI:10.1016/j.jhydrol.2015.06.008

## LIST OF TABLES AND FIGURE CAPTIONS

Table 1 : Descriptive Statistics of physio-meterological and Hydrological Variables.

Table 2: NASH, RMSE, RMSEr, BIAS and BIASr values for all the models. Best values for each quantile for the corresponding metrics are marked in bold.

Table 3: Feature Importance of Five Input Variables used for Specific Flood Quantile Estimation.

Figure 1: Number of trees (n\_estimators) vs OOB error rate for 10, 50 and 100-year flood quantiles.

Figure 2: Relative errors associated with quantiles q50 calculated using RFR and CCA-RFR (the sites are ordered according to their area)

Figure 3: A) q10, B) q50 and C) q100 estimation using RFR approach.

Figure 4: A) q10, B) q50 and C) q100 estimation using CCA-RFR approach.

Table 1: Descriptive Statistics of physio-meterological and Hydrological Variables.

Variables	Minimum	Mean	Maximum	Standard deviation
q10 (m <sup>3</sup> /s.km <sup>2</sup> )	0.03	0.31	0.94	0.20
q50 (m <sup>3</sup> /s.km <sup>2</sup> )	0.03	0.28	0.77	0.18
q100 (m <sup>3</sup> /s.km <sup>2</sup> )	0.03	0.22	0.53	0.13
Area (km <sup>2</sup> )	208	6255	96600	11716
MBS (%)	0.96	2.43	6.81	0.99
FAL (%)	0.00	7.72	47.00	7.99
AMP (mm)	646	988	1534	154
AMD (degree day)	8589	16346	29631	5382

Table 2: NASH, RMSE, RMSEr, BIAS and BIASr values for all models. Best values for each quantile for the corresponding metrics are marked in bold.

	Hydrologic al Variables	CCA- SANN	CCA- EANN	CCA- Kriging	CCA- MLR	SANN	EANN	CCA- GAM	RFR	CCA- RFR
NASH	q10	0.82	<b>0.84</b>	0.78	0.78	0.75	0.78	0.82	0.721	0.577
	q50	0.78	<b>0.8</b>	0.72	0.72	0.69	0.72	0.76	0.657	0.532
	q100	0.77	<b>0.78</b>	0.7	0.68	0.66	0.69	0.67	0.644	0.507
RMSE	q10	0.053	0.05	0.05	0.059	0.06	0.058	0.054	0.063	<b>0.049</b>
	q50	0.082	0.079	0.093	0.094	0.098	0.093	0.087	0.089	<b>0.07</b>
	q100	0.095	0.093	0.11	0.112	0.115	0.109	0.115	0.099	<b>0.08</b>
RMSEr	q10	38	37	51	43	47	44	33.7	80.74	<b>29.44</b>
	q50	44	43	64	49	55	53	43.5	93.39	<b>33.27</b>
	q100	46	45	70	51	64	60	37.0	96.45	<b>35.02</b>
BIAS	q10	0.006	0.005	-0.004	<b>0.001</b>	0.006	0.004	0.009	-	0.002
	q50	0.009	0.009	-0.007	0.005	0.01	0.009	<b>-0.003</b>	-	<b>0.003</b>
	q100	0.013	0.012	-0.008	0.007	0.015	0.013	0.043	-0.019	<b>0.004</b>
BIASr	q10	-5	-5	-16	-9	-7	-7	<b>-3.5</b>	-21.12	-6.64
	q50	-7	<b>-5</b>	-21	-11	-8	-8	-11.4	-25.97	-8.14
	q100	-7	-6	-23	-11	-11	-10	<b>3.4</b>	-27.85	-8.89



Table 3: Feature Importance of Five Input Variables used for Specific Flood Quantile Estimation.

Input Variables	Relative Importance, %		
	q10	q50	q100
Area	87.175962	88.536998	78.255387
MBS	1.398922	0.656013	0.999385
FAL	1.102604	0.707512	0.575426
AMP	8.861476	7.71858	17.89112
AMD	1.461035	2.380897	2.278683



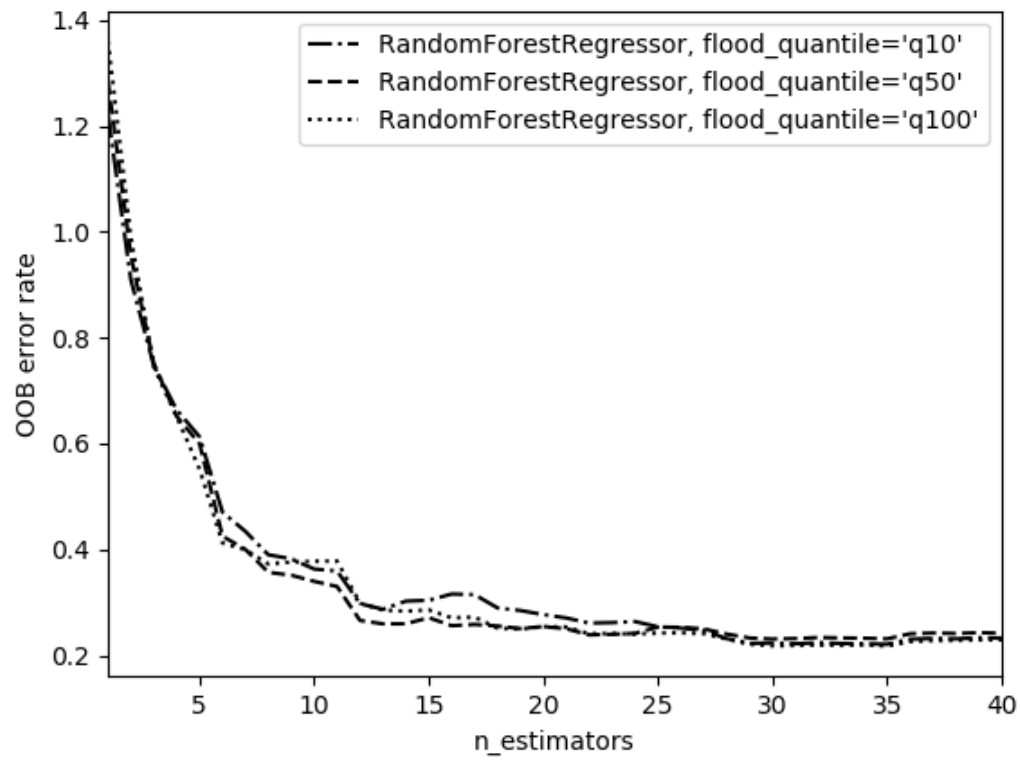


Figure 1: Number of trees ( $n_{\text{estimators}}$ ) vs OOB error rate for 10, 50 and 100-year flood quantiles.

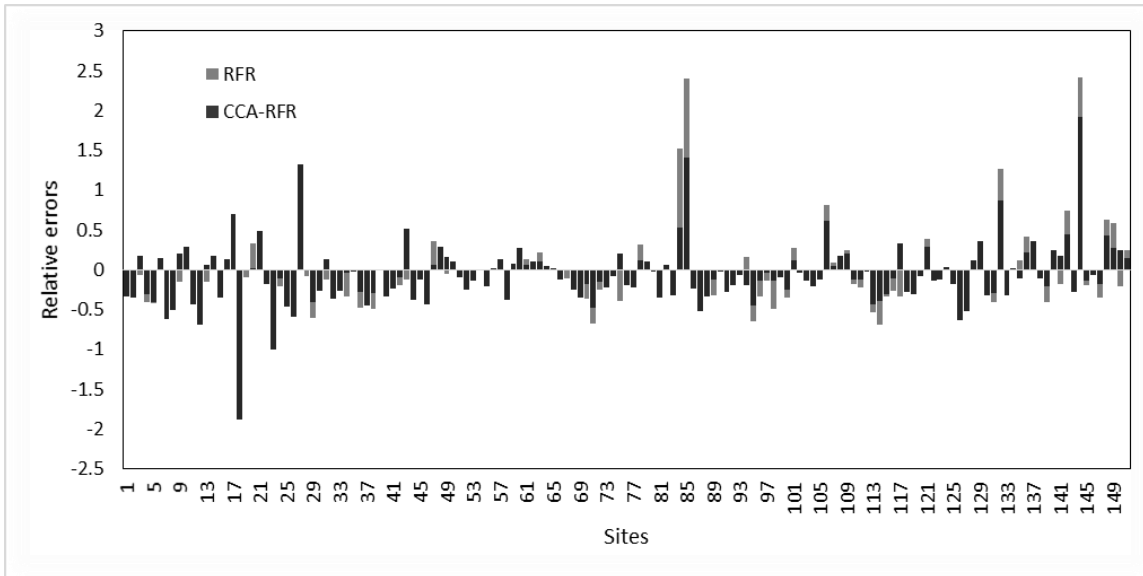
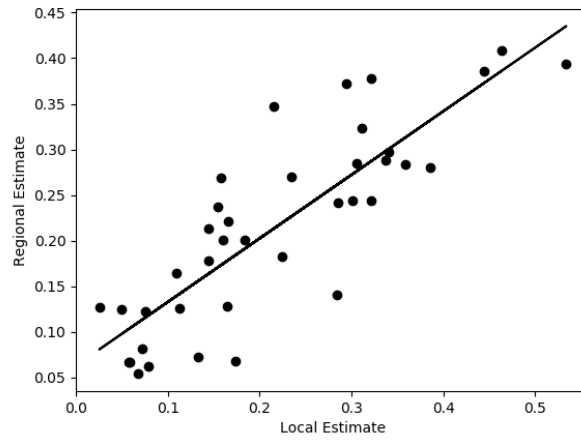
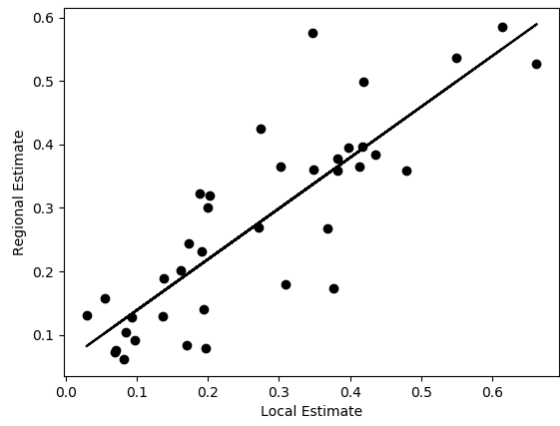


Figure 2: Relative errors associated with quantiles  $q_{50}$  calculated using RFR and CCA-RFR (the sites are ordered according to the increasing area)

A) q10 estimation



B) q50 estimation



C) q100 estimation

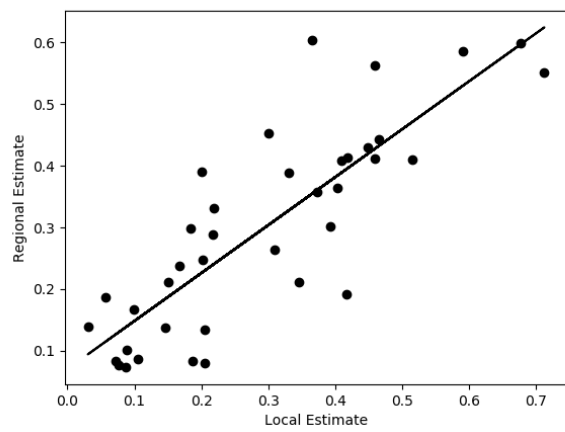
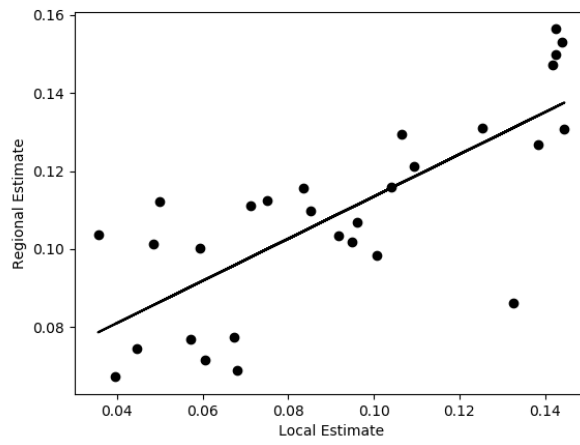
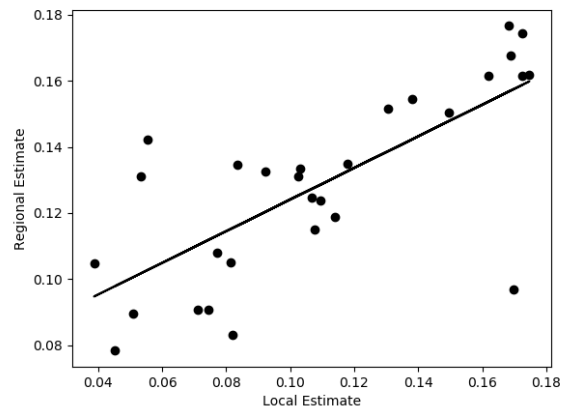


Figure 3: Estimation using the RFR approach

A) q10 estimation



B) q50 estimation



C) q100 estimation

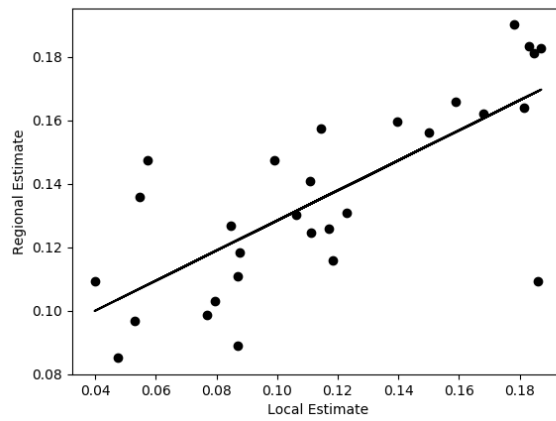


Figure 4: Estimation using the CCA-RFR approach