ELSEVIER

Contents lists available at ScienceDirect

# Advances in Water Resources

journal homepage: www.elsevier.com/locate/advwatres





# On constructing limits-of-acceptability in watershed hydrology using decision trees

Abhinav Gupta <sup>a,\*</sup>, Rao S. Govindaraju <sup>b</sup>, Pin-Ching Li <sup>b</sup>, Venkatesh Merwade <sup>b</sup>

# ARTICLE INFO

# Keywords: Hydrological model Uncertainty Machine learning Runoff ratio Limits-of-acceptability Model validation

# ABSTRACT

A hydrological model incurs three types of uncertainties: measurement, structural and parametric uncertainty. For instance, in rainfall-runoff models, measurement uncertainty exists due to errors in measurements of rainfall and streamflow data. Structural uncertainty exists due to errors in mathematical representation of hydrological processes. Parametric uncertainty is a consequence of our inability to measure effective model parameters, limited data available to calibrate model parameters, and measurement and structural uncertainties. The existence of these predominantly epistemic uncertainties makes the model inference difficult. Limits-of-acceptability (LOA) framework has been proposed in the literature for model inference under a rejectionist framework. LOAs can be useful in model inference if they reflect the effect of errors in rainfall and streamflow measurements. In this study, the usefulness of quantile random forest (ORF) algorithm has been explored for constructing LOAs. LOAs obtained by QRF were compared to the uncertainty bounds obtained by rating-curve analysis and the LOAs obtained by runoff ratio method. Rating curve analysis yields uncertainty in streamflow measurements only and the runoff ratio method is expected to reflect uncertainty in rainfall and streamflow volume measurements. LOAs obtained by using QRF were found to envelop the uncertainty bounds due to streamflow measurement errors. LOAs obtained by QRF and runoff ratio methods were similar. Further, QRF LOAs were scrutinized in terms of their ability to reflect the effect of rainfall uncertainty, both qualitatively and quantitatively. Results indicate that ORF LOAs reflect the effect of rainfall uncertainty: increase in standard deviation with increase in mean streamflow values and decrease in coefficient of variation with increase in mean streamflow values. A mathematical analysis of the LOAs obtained by the QRF method is presented to provide a theoretical foundation.

# 1. Introduction

# 1.1. Background

In a generic hydrological model,

$$y = g(x, \theta) + \delta + \epsilon, \tag{1}$$

 $\delta$  and  $\varepsilon$  denote the *effect* of structural and measurement errors (Beven, 2005) in the estimation of time series of observed hydrologic variables (e.g., streamflow) y by the approximate model g. Here g denotes model inputs such as rainfall and temperature, and g denotes the set of model parameters. Measurement errors refer to errors in measurements of rainfall and streamflow, while structural errors refer to errors in the mathematical representation of hydrologic processes. Given a parameter set g, the structural and measurement errors are estimated based on the

residual time series  $\mathbf{y} - g(\mathbf{x}, \boldsymbol{\theta}^{\mathrm{s}})$ .

If an appropriate probability distribution over  $\delta$  and  $\epsilon$  may be assumed, the parameters of the distributions along with hydrologic model parameters can be obtained by using Bayes theorem (Kennedy and O'Hagan, 2001). However, the use of formal probability distributions has its own challenges (Beven and Smith, 2015). Often, a probability distribution over the *sum* of  $\delta$  and  $\epsilon$  is assumed, such as Gaussian or generalized Gaussian (Schoups and Vrugt, 2010; Ammann et al., 2019; Smith et al., 2015). But the residual time series can yield only an aggregate estimate of the effect of measurement and structural errors, that is, the quantities  $\delta$  and  $\epsilon$  are individually unidentifiable (Renard et al., 2010, 2011; Brynjarsdóttir and O'Hagan, 2014). Separate identification of structural and measurement errors is required to determine what part of modeling exercise needs to be addressed to reduce total uncertainty, the data or the model (e.g., Reichert and Mieleitner, 2009)

E-mail address: abhinav.gupta@dri.edu (A. Gupta).

https://doi.org/10.1016/j.advwatres.2023.104486

<sup>&</sup>lt;sup>a</sup> Division of Hydrologic Sciences, Desert Research Institute, Las Vegas, NV, USA

<sup>&</sup>lt;sup>b</sup> Lyles School of Civil Engineering, Purdue University, West Lafayette, IN, USA

<sup>\*</sup> Corresponding author.

and to facilitate rejection of bad models.

To identify structural uncertainty in a model, strong prior information about measurement uncertainties is required (Renard et al., 2010; McMillan et al., 2012; Brynjarsdóttir and O'Hagan, 2014; McMillan et al., 2018), and this information should be obtained before calibration and independent of the hydrologic model being used. Given information about measurement uncertainty and the residual time series corresponding to a model (or model parameters), a Bayesian characterization of structural uncertainty is possible in the sense that one can obtain a probabilistic estimate of the effect of structural uncertainty conditioned upon each possible realization of rainfall (and other inputs) and streamflow time series. Priors over measurement uncertainty are typically constructed by making aleatoric assumptions about the nature of these errors. For example, one can obtain information about random measurement uncertainty in streamflow by using rating curve analysis (Kiang et al., 2018; Petersen-Øverleir et al., 2009; Reitan and Petersen-Øverleir, 2009; Le Coz et al., 2014) or other probabilistic methods (de Oliveira and Vrugt, 2022). But epistemic uncertainties in streamflow, such as those introduced by extrapolation of rating curve to gauge heights well above the observations, may not be knowable. Reliable information about rainfall measurement uncertainty cannot be obtained in most situations. For instance, one may estimate the uncertainty in areal average rainfall by assuming that this uncertainty is dominated by spatial variability of rainfall and neglecting temporal errors and biases (Moulin et al., 2009; Renard et al., 2011). Spatial variability can be modeled using a statistical model such as Kriging, provided that enough data to estimate the parameters of the variogram are available. This is further complicated as the parameters of the variogram will change from event to event in unknown ways. Precipitation data also incur timing errors which can be significant if the precipitation gauges are sparse or are located outside the watershed.

If the observed event seem to violate the principle of mass balance (e. g., Beven and Westerberg, 2011), one may expect errors in the measurements of either rainfall data, or streamflow data, or both. Such time-periods in rainfall-runoff time series are referred to as disinformative (Beven and Westerberg, 2011) which should be discarded before model fitting. A disinformative event can introduce bias in the modeling effort because it violates mass balance, and also because it affects the antecedent conditions for subsequent events (Beven and Smith, 2015). Disinformative periods in a rainfall-runoff dataset may be identified as the ones with exceptionally high and low runoff ratios (Beven and Westerberg, 2011) where runoff ratio of an event is defined as the ratio of total event streamflow to total event rainfall. What is an exceptionally high or low value of runoff ratio may be determined using the knowledge about the rainfall-runoff response of the watershed. Several other attempts have been made to characterize the uncertainty in hydrologic data and hydrologic modeling (e.g., Kuczera and Parent, 1998; Kavetski et al., 2006a, 2006b; Gabellani et al., 2007; Gong et al., 2013; McMillan et al., 2018), but it still remains an unsolved problem because of dominantly epistemic nature of these errors. Recently, Gupta and Govindaraju (2022) noted that several methods have been proposed for uncertainty analysis in hydrology but there is no consensus on which method should be used.

Recently, the runoff ratio method has been proposed to construct limits-of-acceptability (LOA) bounds on streamflow that could then be used to identify behavioral models (Beven, 2019). A model (or a model parameter set) is considered behavioral if the streamflow simulated by it falls within the LOA at some predefined timesteps (Beven et al., 2022) depending on the purpose of the modeling exercise. It is clear that LOA should be such as to *encompass* the uncertainty due to measurement errors in rainfall and streamflow. Thus, a model that properly accounts for streamflow dynamics within the margin of measurement errors would not be rejected and will be considered behavioral.

LOAs have also been defined using flow duration curves (FDCs; Westerberg et al., 2011). In this method, measurement uncertainty over streamflow time series is obtained using rating-curve analysis.

Measurement uncertainty in streamflow is converted to an uncertainty bound over FDC. A model (or model parameter set) is considered behavioral if the FDC simulated by it falls into the FDC uncertainty bound. However, this method only compares the probability distribution of observed and simulated streamflows and removes the temporal information from the streamflow time series. Also, it does not account for rainfall measurement errors. In fact, most of the methods to derive LOAs are based on streamflow uncertainty only and neglect rainfall uncertainty (e.g., Krueger et al., 2010; Coxon et al., 2014). To the best of author's knowledge, the runoff ratio method is the only method that constructs LOAs while acknowledging uncertainty in both streamflow and rainfall measurements. The runoff ratio method also has some limitations as discussed below.

Fundamentally, the LOA method has been proposed in a rejectionist framework (Beven and Lane, 2019), which makes it different from Bayesian methods wherein no models are explicitly rejected. Frequentist statistics also provides a model rejection framework such as the likelihood ratio test (Neyman and Pearson, 1933), Fisherian hypothesis testing (Fisher, 1956) and, more recently, evidential testing (Royall, 2017; Lele, 2004). But these methods are based on aleatoric assumptions (as are Bayesian methods) about various uncertainties and, therefore, are difficult to justify in hydrologic applications. There have been a relatively few attempts in hydrology to use rigorous frequentist methods for model inference (but see Pande, 2013a, 2013b). The LOA framework provides an alternative to the formal statistical frameworks, as it combines the elements of Bayesian theory (parameter update as the models are tested against more data) and frequentist statistics (model rejection). LOA can also be applied in a purely Bayesian framework by defining an appropriate LOA-based likelihood function (e.g., Krueger et al., 2010). The aim of this study was to explore the potential of using machine learning algorithms called decision tree (DT) and, in particular, quantile random forest (QRF) in constructing LOAs in gauged and ungauged locations.

# 1.2. Runoff ratio method, and decision trees

In runoff ratio method, the rainfall and streamflow time series are divided into separate rainfall-runoff events. Then, the rainfall-runoff events with similar characteristics are pooled together. The main idea is that the two similar events should have similar runoff ratios. Of course, no two events are exactly similar, and there would be some differences in runoff ratios. But the large differences can be (at least partly) attributed to either rainfall and/or streamflow measurement errors. The differences between runoff-ratio values of two similar events may also result from imperfections in the methodology to compute runoff ratios. Multiplying a zero-loss streamflow event with runoff ratios of all the similar events would result in an ensemble of corresponding streamflow hydrographs. Zero-loss streamflow can be obtained by dividing the observed hydrograph by the corresponding runoff ratio. Beven (2019) suggested that the upper and lower bounds of these hydrographs be used as LOA over the rainfall-runoff event in question. The different hydrographs in the ensemble can be assigned a weight based on the similarity of the corresponding event with the event for which LOA is being constructed. This method is described in more detail

The advantage of the runoff ratio method is that it allows to define a distribution of streamflow hydrographs for a given rainfall event and antecedent conditions based on available data. A limitation of this method is that it is applicable to flashy watersheds only (Beven, 2019). Also, this method cannot account for potential timing errors in precipitation – it only accounts for errors in precipitation and streamflow volume and further can be applied only at an event timescale. Further, this method cannot be used to construct LOAs at ungauged locations where streamflow data are unavailable for computing runoff ratios.

These limitations can be addressed by using a Machine Learning (ML) method, while retaining the advantage of the runoff ratio method.

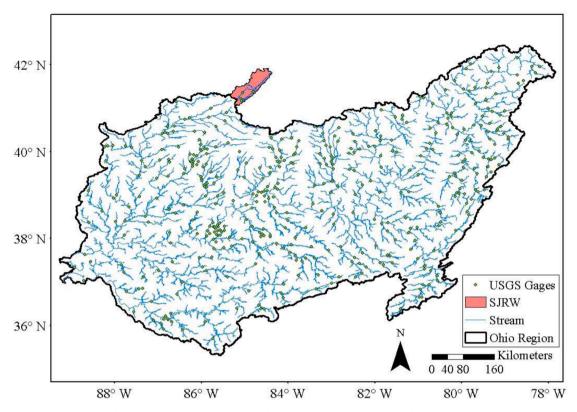


Fig. 1. Ohio river basin (ORB) and USGS streamflow stations (green dots). The watershed with red background is St. Joseph River Watershed (SJRW).

A direct mapping between relevant watershed attributes, meteorological data, and streamflow can be created by using an ML algorithm (e.g., Govindaraju, 2000; Zhang and Govindaraju, 2000, 2003; Iorgulescu and Beven, 2004; Shortridge et al., 2016; Kratzert et al., 2019). ML can be particularly useful in constructing LOAs for baseflow dominated watersheds where runoff ratio method is not applicable and to construct LOAs at ungauged locations. Further, the ML approach allows defining LOAs at the scale of available data. As discussed below, ML algorithms called decision trees (DTs) are particularly well-suited in this regard.

Another advantage of the ML approach is that data from several watersheds may be used to train the model and define LOAs. Data from different watersheds, however, may introduce disinformation because of watershed-specific epistemic uncertainties (Beven, 2020). But the hydrologically relevant information available from other watersheds may still be useful, especially when LOAs are to be constructed for an ungauged watershed. An ML algorithm such as DT will be able to identify hydrologically similar watersheds based on available watershed characteristics, albeit that watersheds characteristics are typically represented by spatially averaged indices neglecting their spatial variation. Thus, DTs are natural candidates to consider for constructing LOAs as discussed below.

The uncertainties in hydrological data are predominantly epistemic, which may change from event to event in unknown ways, and the true statistical behavior of uncertainties will not be generally represented by the available data. Therefore, DT would either overpredict or underpredict the effect of measurement errors. While overprediction is acceptable, underprediction may be problematic in many applications. Therefore, one needs to allow for outliers while validating the models using the LOA method (as in Beven et al., 2022). Further, the DT model would compensate for systematic biases. These systematic errors cannot be detected by a statistical approach. A bias term can be introduced in statistical models, but these models would not be able to differentiate between the bias in the data and the bias in the model simulations.

The classical method of finding uncertainty in the measurement of a phenomenon is to repeat the measurement process several times under

identical conditions. The repeated sampling method, however, is impossible for the measurements of environmental phenomena such as rainfall and streamflow (McMillan et al., 2012). But an approximate repeated sampling method may be implemented for environmental measurements. The main idea is to estimate the effect of measurement uncertainty using observations of rainfall-streamflow events under *similar* conditions across several different events and/or several different watersheds. The runoff ratio and DT methods can be thought of as approximate repeated sampling techniques.

Once the LOAs are obtained, either formal or informal Bayesian (Liu et al., 2009; Krueger et al., 2010; Beven and Lane, 2022) methodologies may be used for subsequent uncertainty analysis. In informal methods, one may define behavioral models (and model parameters) as ones that yield streamflow time series within the LOA. Thus, all the models with an inferior structure will eventually be rejected as more and more data are used (at least that is the expectation). One can also use the apparatus of formal Bayesian theory for model (or parameter) inference using the LOAs in Approximate Bayesian Computation framework (Nott et al., 2012; Sadegh and Vrugt, 2013; Vrugt and Sadegh, 2013; Vrugt and Beven, 2018).

# 1.3. Objectives

The objective of this study is to develop a method for constructing LOAs that can account for both precipitation and streamflow measurement errors and can be used for ungauged catchments. In this study, we ask if a variant of DT called quantile random forest (QRF) may be used to construct meaningful LOAs. A second question is if the LOAs obtained by QRF algorithm are comparable to those obtained by the runoff ratio method of Beven (2019).

The novelty of this study lies in using QRF model to construct LOAs that account for measurement uncertainty based on available data. To address the objective of this study, uncertainty bounds obtained by QRF model are scrutinized to check if they can be used as LOAs. The uncertainties in real world data are, however, unknown; therefore, it is

**Table 1**Predictor variables in machine learning models to estimate streamflow time series at a station in a river-network. Exploratory statistics in the third column represent (minimum, maximum, median, and mean).

Predictor variable	Description	Exploratory Statistics
Drainage area	Cumulative drainage area of	(7.74, 250260, 624,
(Km <sup>2</sup> )	streamflow station	4187)
Impervious Area* (%)	Percentage of impervious area	(1.92, 7.74, 6.36, 6.44)
Sand	Percentage of sand content	(6.34, 49.61, 20.97,
content**(%)		19.78)
Clay content (%)	Percentage of clay content	(15.88, 45.12,
		26.03, 27.58)
Conductivity	Average hydraulic conductivity of the	(0.01, 77.22, 0.19,
$(\mu m s^{-1})$	drainage area	3.51)
Permeability	Average permeability of the drainage	(1.02, 15.09, 3.87,
$(cmhr^{-1})$	area	4.82)
Rainfall***	Total daily rainfall during current and	_
	previous 1, 7, and 30 days	
Snowfall	Total Daily snowfall during current and	-
	previous 1 and 30 days	
Snow depth	Daily snow depth during current and	_
	previous 1 and 30 days	
Temperature	Average daily maximum and minimum	_
	temperature at current day	

<sup>\*</sup> Land-use data were collected from NLCD database.

Table 2
List of USGS stations used for testing the proposed method.
These stations are located St. Joseph River Watershed (SJRW).

USGS station	Drainage Area (km²)
04180500	2745.40
04180000	699.30
04179520	233.62
04178000	1579.90

impossible to check if the uncertainty bounds obtained by any method represent true uncertainties. Some characteristics of the uncertainties can be obtained by using statistical methods based on aleatoric assumptions; we test whether the QRF estimated LoAs reflect the effect of these uncertainties or not.

Further, this paper presents a mathematical analysis of the proposed hypothesis. The goal of the mathematical analysis is (1) to show how decision trees such as QRF can be used to encompass measurement uncertainties due to errors in rainfall and streamflow measurements, and (2) to clarify the logic and assumptions behind the proposed method.

In Section 2, the theory behind DTs and QRF algorithm are discussed along with the methodology to empirically test the proposed method. Section 3 discusses the results of the study. Section 4 presents a brief mathematical analysis of the QRF method in terms of defining LOAs. Section 5 concludes the paper.

# 2. Theory and methodology

# 2.1. Study area, data, and the models developed

In this study, data from Ohio river basin (ORB) were used to calibrate and validate the QRF model. This basin contains 431 USGS streamflow stations (Fig. 1). The streamflow data were downloaded from USGS website for all the 431 stations. Data for these watersheds are available from water year 2011 to 2020. Total drainage area of each USGS station was delineated on the 30m  $\times$  30m resolution digital elevation model

(Archuleta et al., 2017; U.S. Geological Survey, The National Map, 2017) by using the ArcHydro toolbox. For each of the drainage areas, predictor variables (listed in Table 1) were computed or collected. Climate data were collected over the study area from Historical Climate Network (HCN) stations available at National Centers for Environmental Information (NCEI) website.

To test the capability of the QRF model in capturing rainfall and streamflow measurement uncertainties, data from St. Joseph River Watershed (SJRW) were used as test cases. SJRW is located just above the ORB in Northwest as indicated in Fig. 1 (see also Figure B1 in Appendix B). The drainage areas of the SJRW watersheds are listed in Table 2. Specifically, QRF models were used to generate LOAs at four USGS streamflow stations located in SJRW.

Three kinds of QRF models were developed:

- (1) Gauged-single scenario: In this case, four individual QRF models were developed for each of the four SJRW watersheds using data from the watershed where the LOAs were to be constructed. For example, to construct LOAs at station 04180500, the data from only this station were used to train the QRF model. These models are referred to as "gauged-single models".
- (2) Gauged scenario: In this case, a QRF model was trained using data from both the ORB and the four SJRW watersheds. The model thus trained is referred to as "gauged model". Three kinds of models were developed in this scenario: (2a) QRF was trained using data from all the training watersheds (referred to 'gauged all'), (2b) QRF was trained using data from the 4 most similar watersheds to the watershed where LOAs are to be constructed (referred to 'gauged 4'), and (2c) QRF was trained using the data from the 20 most similar watersheds (referred to 'gauged 20').
- (3) Ungauged scenario: In this case, a QRF model was trained using data only from the ORB watersheds without using the SJRW data. The model thus trained will be referred to as "ungauged model".

Out of the 431 ORB stations, 80% of the stations were fixed for the calibration of QRF and the remaining stations were fixed for validation. Similar watersheds in the 'gauged scenario' were selected based on the watershed static attributes and the mean climate (mean precipitation and temperatures). The first two scenarios allow us to test the usefulness of QRF approach in constructing LOAs at a gauged location and the third scenario allows us to test the usefulness of the approach at ungauged locations. The comparison of the first two and the third scenario allows to test the usefulness of data across multiple watersheds in constructing LOAs.

# 2.2. Machine learning models to map predictor variables to streamflow

The main idea behind ML algorithms is to create a mapping between predictor and response variables (Friedman et al., 2001, chap. 2). For most watershed scale rainfall-runoff models, the set of predictor variables constitutes meteorological data, soil data, land-use data, etc. (Table 1), and the response variable typically is streamflow time series. Available data are divided into calibration and validation sets. The samples contained in calibration set are used to create a mapping such that a loss function, which is a function of the mapping, is minimized, and the samples contained in validation set are used to test the generalizability of the created mapping.

In this study, QRF was used to create a mapping between predictor and response variables (Breiman et al., 1984; Breiman, 2001). The basic building block of QRF is another ML algorithm called regression trees (Friedman et al., 2001, chap. 9; Iorgulescu and Beven, 2004). Regression trees create a non-linear mapping between predictor and response variables. In this method, the space of predictor variables is divided into *S* (contiguous) subregions, and in each subregion, the response variable is approximated by a unique function.

Let the set containing predictor and response variables be denoted by

<sup>\*\*</sup> Soil data were collected from STATSGO database.

 $<sup>^{\</sup>ast\ast\ast}$  Climate data were collected from Global Historical Climatology Network (GHCN) database.

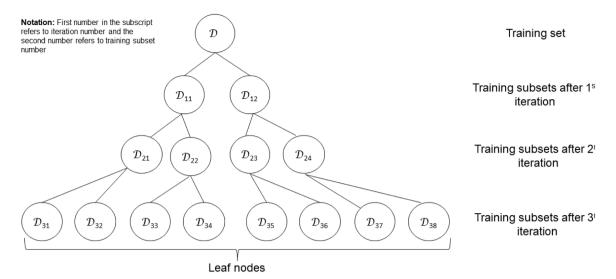


Fig. 2. Illustration of regression tree. In this hypothetical example, only three iterations were carried out to divide the training set into smaller subsets.

 $\mathscr{D}$ . Each element of  $\mathscr{D}$  represents a calibration/training sample. Let the  $i^{\text{th}}$  calibration sample be denoted by  $(x_i, y_i)$ , then  $\mathscr{D} = \{(x_1, y_1), (x_2, y_2), ..., (x_N, y_N)\}$  where N is the total number of calibration samples. The vector  $x_i$  is a p-vector where p denotes the number of predictor variables, that is,  $x_i = (x_{i1}, x_{i2}, ..., x_{ip})$ , and  $y_i$  is a scalar that denotes the response variable corresponding to the  $i^{\text{th}}$  sample. In this study, the  $i^{\text{th}}$  response variable is streamflow at the outlet of a watershed at a particular timestep. The  $i^{\text{th}}$  predictor vector includes static watershed attributes and meteorological data at multiple lags (Table 1). The regression tree is created using an iterative procedure. In the first iteration, the set  $\mathscr{D}$  is divided into two (or more) subsets based on a randomly selected  $j^{\text{th}}$  predictor variable. Let the two subsets be denoted by  $\mathscr{D}_{11}$  and  $\mathscr{D}_{12}$ , then

$$\mathscr{D}_{11} = \left\{ (\mathbf{x}_i, y_i) \middle| x_{ij} < x_{j,\text{thresh}} \right\},\tag{2}$$

$$\mathscr{D}_{12} = \mathscr{D} \backslash \mathscr{D}_{11},$$

where  $x_{i,\text{thresh}}$  denotes a randomly chosen threshold for  $j^{\text{th}}$  predictor variable. In the second iteration, the subsets  $\mathscr{D}_{11}$  and  $\mathscr{D}_{12}$  are further divided into smaller subsets, and so on for subsequent iterations. At the end of the iterative procedure, S smaller subsets of  $\mathcal{D}$  are obtained, and each subset occupies a distinct region of the predictor space. Thus, the regression tree algorithm divides the predictor space into S contiguous subregions. This method is referred to as regression trees because the process of division of training samples into S subsets can be visualized as creating a tree (Fig. 2, see also Friedman et al., 2001, pp. 268). The tree grows deeper with each iteration. Therefore, the number of iterations is also referred to as tree depth. Typically, a maximum value of tree depth, d, is assigned to avoid overfitting. The subsets obtained in the last iteration are also referred to as leaf nodes. It is clear that there is a relationship between the number of leaf nodes *S* and maximum tree depth *d*: an increase in d implies an increase in S. Note that once the tree is created, each subregion can be identified by a set of rules on predictor

After the tree is created, response of a sample with predictor vector  $\boldsymbol{x}$  is obtained as follows. The first step is to identify the subregion of the predictor space to which the vector  $\boldsymbol{x}$  belongs. Suppose that  $\boldsymbol{x}$ belongs to the  $i^{th}$  subregion corresponding to  $i^{th}$  training subset denoted by  $S_i$ . Then the response variable corresponding to  $\boldsymbol{x}$  is estimated as the average response of calibration samples contained in  $S_i$ 

$$\widehat{\mathbf{y}}(\mathbf{x}) = \frac{1}{L_i} \sum_{i=1}^{L_i} \mathbf{y}(\mathbf{x}_i), \tag{3}$$

where  $L_i$  denotes the number of samples in  $S_i$ . Regression trees are

developed so that the sum of square errors between observed and estimated responses is minimized (with some regularization to avoid overfitting). The averaging of data in the leaf node, however, neglects the variability in the data. Therefore, not just the average but the entire distribution  $y(x_j)$  for  $x_j \in S_i$  were used to construct LOAs as explained below.

The method of regression trees is particularly suitable for the purpose of creating LOAs because it mimics the function of an approximate repeated sampler by grouping similar calibration samples (similarity in predictor space) together based on several watershed attributes, thus enabling the accounting of measurement uncertainty due to errors in response and predictor variables. Regression trees have to be regularized to avoid overfitting; therefore, B regression trees are developed instead of a single one. Each of the B regression trees is created by randomly drawing K samples by bootstrapping from the calibration set  $\mathcal{D}$ . This, yields an ensemble  $Y(x) = \{\hat{y}_1(x), \hat{y}_2(x), ..., \hat{y}_B(x)\}$  of streamflow estimates corresponding to the predictor variable x where the  $b^{th}$ estimate  $y_b(x)$ , obtained by Eq. (3), corresponds to the  $b^{th}$  tree. The average of values in Y(x) is taken as the final estimate. This method is known as random forest (RF). In this study, the RF algorithm was used to create a mapping between predictor variables (listed in Table 1) and streamflow, and the streamflow in each subregion of the predictor space was estimated as the average streamflow of calibration samples in that subregion (Eq. (3)). But as mentioned above, taking averages of data in the leaf node neglects the variability in the leaf node which might contain important information about uncertainties. Therefore, quantile random forest (QRF) technique was used to construct LOAs, where quantiles instead of averages are computed. In this technique, the ensemble YORF is constructed by using the entire distribution of data in leaf nodes. If a given predictor, say x, falls into the  $i^{th}$  leaf node of the  $b^{th}$  tree, denoted by  $S_i^b$ , then the distribution of response variable in  $S_i^b$  can be represented

$$Y^{b}(\mathbf{x}) = \{ y_{i} | y_{i} \in S_{i}^{b}, \}. \tag{4}$$

Thus, we will have a distribution  $Y^b$  for each tree. Now, the data from each  $Y^{bs}$  can be combined to form an ensemble  $Y_{QRF}(x)$ 

$$Y_{QRF}(\mathbf{x}) = \{ y_j | y_j \in Y^b, b = 1, 2, \dots B \}.$$
 (5)

Note that the  $y_j$  values contained in  $Y_{QRF}$  are observed values not the estimates. QRF estimates different quantiles of the response for a given x by treating  $Y_{QRF}$  as the distribution of response. In this study, 2.5<sup>th</sup> and 97.5<sup>th</sup> percentiles obtained by QRF were used as lower and upper LOAs. We found that these percentiles were typically adequate for constructing

**Table 3**List of priors over rating curve parameters.

Parameter	Prior		
m	√{1,2,3} – discrete uniform		
$log a_k$	$\mathscr{U}(0,8)$		
$b_k$	$\mathscr{U}(0.5, 3.5)$		
$h_{0,1}$	$\mathscr{U}(-5,h_{min})$		
$h_{\mathrm{s},k}$	$\mathscr{U}(h_{s,k-1},h_{max})$		
$h_{0,k}$	$\mathscr{U}(-5,h_{s,k-1})$		
$h_{\mathrm{s},1}$	$\mathscr{U}(h_{0,1},h_{max})$		
ф	$\mathscr{I}\mathscr{G}(2,0.1)$		
β	$\mathscr{U}(-1,1)$		
γ	$\mathscr{G}(1/0.57, 0.57)$		
$\mathscr{U} = \text{Uniform}; \ \mathscr{G} = \text{Gamma}; \ \mathscr{I}\mathscr{G} = \text{Inverse Gamma}$			
Gamma distribution: $f(x) = \frac{1}{\Gamma(\alpha)\beta^{\alpha}}x^{\alpha}$	$\frac{x}{e} - \frac{x}{\beta}$		

LOAs in the sense that most of the observations were enveloped by the LOAs but a few flow values could not be enveloped. Therefore, in practical applications, more extreme percentiles might be appropriate for creating LOAs.

If the premise 'the ensemble of estimated streamflow represents only measurement uncertainty' were true, then in the absence of measurement errors the different streamflow estimates in the ensemble would be (approximately) identical. In practice, however, even in the absence of measurement errors, the streamflow estimates in the ensemble would be different because of several reasons:

- (1) Imperfections in creating the regression trees: These imperfections include selection of appropriate values of *B* (number of regression trees) and *S* (number of leaf nodes). A large value of *S* (or large value of maximum tree depth *d*) may result in an overestimation of measurement errors and conversely for a small value of *S* (or small value of *d*). In this study, optimal values of *B* and *d* along with minimum number of samples in a leaf node were estimated by computing the out-of-bag (OOB) error. The OOB error is the prediction error of calibrated RF from the left-out training set. An early stopping method searches for the optimal values of these parameters with the minimal OOB error.
- (2) Small calibration set which is inadequate to represent the population of measurement errors: Calibration sets should be large enough such that the variability in measurement errors (in rainfall and streamflow) is captured. In this study, data from a total of 431 ORB stations plus 4 SJRW stations were used, out of which data from a total of 344 stations were used for calibration.
- (3) The set of predictor variables used to train the ML algorithm is incomplete: If a relevant predictor variable is missed in the set of predictor variables, the uncertainty bound yielded by QRF would also contain structural errors. The predictor variables used in this study are listed in Table 1. Though these predictors variable are incomplete; they are still good enough to estimate the streamflow time series accurately in many watersheds, as evident by high NSE for some of the test stations shown in the results section.

Even after taking all the precautions, the LOAs created by QRF method would still contain structural errors. QRF would be able to construct better LOAs as the sample size increases. When the LOAs are to be constructed at a gauged location, the longer length of data at the location will be more important than the data from other watersheds. But data from other watersheds would be the only option when LOAs are to be constructed at an ungauged location.

The LOAs obtained by QRF were compared against the bounds obtained over streamflow measurements uncertainty which in turn were obtained by rating curve analysis. If the LOAs obtained by QRF indeed reflect the effects of measurement uncertainties in rainfall and streamflows, these should envelop the uncertainty bound obtained by rating

curve analysis. Also, we compared the bounds obtained by runoff ratio method to the bounds obtained by QRF method. Analysis of rating curve and runoff ratio results was carried out at the four USGS streamflow gauging stations within SJRW as indicated in Table 2. SJRW is located immediately Northwest of ORB as indicated in Fig. 1.

Moreover, the QRF LOA should also reflect the effects of measurement uncertainty in rainfall. In this study, the measurement uncertainty in areal average rainfall was obtained using an empirical approach. One challenge is that the rainfall uncertainty bounds cannot be directly compared to the LOAs since rainfall is processed through the watersheds in a highly non-linear fashion before it reaches the watershed outlet. There is no exact way of translating measurement uncertainty in rainfall to streamflow space. Therefore, in this study, the various realizations of rainfall were processed through the SCS curve-number (CN) formula for different values of CN to get an estimate of excess rainfall. Subsequently, coefficient of variation of streamflow ( $CV_Q$ ) were compared to the coefficient of variation of excess rainfall time series ( $CV_R$ ).

# 2.3. Rating curve analysis to quantify uncertainties in measured streamflow

The streamflow at a river cross-section is estimated using the observed relationship between measured gage heights at the cross-section and corresponding measured discharges; this relationship is referred to as rating curve (Herschy, 1993). Commonly, a rating curve is modeled as multiple power law segments (Le Coz et al., 2014):

$$\log(Q_{r}(h)) = \begin{cases} 0, & h \leq h_{0,1}, \\ \log a_{1} + b_{1}\log(h - h_{0,1}), & h_{0,1} \leq h \leq h_{s,1}, \\ \log a_{2} + b_{2}\log(h - h_{0,2}), & h_{s,1} \leq h \leq h_{s,2}, \\ \vdots & \vdots & \vdots \\ \log a_{m} + b_{m}\log(h - h_{0,m}), & h_{s, m-1} \leq h. \end{cases}$$
(6)

In Eq. (6),  $Q_r$  is the estimated streamflow, h is measured gage height,  $h_{0,1}$  is the cease-to-flow parameter of lowest power-law segment which corresponds to height of riverbed with respect to datum,  $h_{s,k}$  is the upper bound of  $k^{th}$ power-law segment on h axis,  $h_{0,k}$  is the cease-to-flow parameter of  $k^{th}$  segment,  $a_k$  and  $b_k$  are the multiplier and exponent parameters of the  $k^{th}$  segment, and m is the number of rating curve segments. Typically, several gage heights are measured during a day which are then converted to streamflow using the rating curve. Eq. (6) corresponds to Manning equation (Sturm, 2001) for flow in open channels (with the assumption that hydraulic radius is approximately equal to depth; Le Coz et al., 2014) and is a frequently used relationship in hydraulic modeling. Errors in gage height measurements may be assumed negligible (Reitan and Petersen-Øverleir, 2009). Thus, uncertainties in estimated streamflow are mainly due to errors in direct measurements of streamflow that are used to construct the rating curve. In this study, the following model was used to quantify the uncertainties in estimated streamflow:

$$Q(h) = Q_{\rm r}(h) + \epsilon_{\rm r},\tag{7}$$

where  $Q_r(h)$  is determined by Eq. (6),  $\varepsilon_r$  is the random measurement error in observed streamflow and Q(h) is the observed streamflow. Further, we assumed the  $\varepsilon_r$ 's at different time-steps to be distributed independently as skewed exponential power distribution (Fernández and Steel, 1998). Also, Q(h) was truncated at zero which makes the probability density of Q equal to

$$p_{Q}(Q) = \frac{\frac{2}{\gamma + \gamma^{-1}} \left\{ f_{\epsilon_{r}} \left( \frac{\epsilon_{r}}{\gamma} \right) \mathbf{I}_{[0,\infty)}(\epsilon_{r}) + f_{\epsilon_{r}}(\gamma \epsilon_{r}) \mathbf{I}_{(-\infty,0)}(\epsilon_{r}) \right\}}{1 - \Phi(0|Q_{r}, \phi, \beta, \gamma)} \mathbf{I}_{[0,\infty)}(Q), \tag{8}$$

where  $\gamma \in (0, \infty)$  is the skew parameter, I denotes the indicator function,  $\Phi(0|Q_r,\phi,\beta,\gamma)$  is the probability that the value of untruncated Q is less than zero, and  $f_{\epsilon_r}$  is the power exponential distribution with scale parameter  $\phi$  and shape parameter  $\beta \in (-1,1]$ ,

**Table 4**Number of rainfall-runoff events for each of the USGS stations in the SJRW.

USGS station	Number of rainfall-runoff events	
04180500	138	
04180000	148	
04179520	139	
04178000	146	

$$f_{\epsilon_{r}}(\epsilon_{r}) = \Gamma^{-1} \left( 1 + \frac{2}{1+\beta} \right) 2^{-\left(1 + \frac{2}{1+\beta}\right)} \phi^{-1} \exp\left( -\frac{1}{2} \left| \frac{\epsilon_{r}}{\phi} \right|^{\frac{2}{1+\beta}} \right). \tag{9}$$

The priors listed in Table 3 were used as weakly informative priors over parameters of the models  $Q_{\rm r}$  and  $\varepsilon_{\rm r}$ , following Reitan and Petersen-Øverleir (2009). Strictly, uniform priors over the parameters of  $Q_{\rm r}$  are not non-informative (Gupta et al., 2022). This difference, however, would have minimal effect on our analysis as we are concerned only with the width of uncertainty bounds over streamflow time series, not the probabilities assigned to different realizations of streamflow time series. Further, we have not imposed any upper limit on the distribution of streamflow. Very low (practically zero) probability will be assigned beyond a certain magnitude of Q (irrespective of the prior distribution used) – the results obtained for the four SJRW stations confirm that absence of upper limit does not have any effect on the obtained uncertainty bounds. Validity of the error model of Eq. (8) was assessed a-posteriori via QQ plots.

The aleatoric assumption made in the analysis may not be valid during the peak events. It has been shown using hydraulic modeling that uncertainty during peak events can be very high (Di Baldassarre and Montanari, 2009). These uncertainties are epistemic in nature rather than aleatoric, and, therefore, a formal statistical treatment of these uncertainties is difficult. To test how well the QRF LOAs envelop the streamflow uncertainty due to these epistemic sources, we computed the fraction of peaks enveloped by the QRF LOAs, if the true peaks were some multiple f of the observed peaks, with f varying from 1.1 to 2. We refer to this analysis as the multiplier analysis in this study. Only the peaks with flow values greater than 50-percentile were considered for this analysis.

The posterior distribution over parameters was computed using Delayed Rejection Adaptive Metropolis (DRAM) algorithm (Haario et al., 2006) in an approximate Bayes setting (Nott et al., 2012). The approximate Bayes computations facilitated faster convergence to a posterior distribution. This method of rating curve analysis is same as that of Reitan and Petersen-Øverleir (2009) except that they used a multiplicative error model instead of an additive error model. The multiplicative error model was considered unsuitable in this case because of the large range of streamflow values as opposed to that in Reitan and Petersen-Øverleir (2009) study: a multiplicative error model would result in unrealistically high uncertainties at larger values of observed streamflow. Additive error structure used in this study was found to be appropriate (by the way of QQ plot test) in the examples considered in this study. Convergence to posterior distribution was confirmed using R-diagnostic statistic (R<sub>d</sub>; Gelman and Rubin, 1992). Markov chains were assumed to converge to posterior distribution if  $R_d$ converged to a value below 1.1 and never increased on further simulations of the chains. The posterior distribution was further processed to remove the parameter sets that yielded large deviations between observed and estimated streamflow: the deviation between observed and estimated streamflow was measured using sum-of-square-errors. The computed posterior distribution over parameters (of both Q<sub>r</sub> and  $\varepsilon_r$ ) was used to simulate several streamflow time series that were assumed to represent random uncertainty in measurements of streamflow, as obtained by the rating-curve method.

# 2.4. Uncertainty bound in areal average rainfall

The uncertainty in areal average rainfall exists due to errors in rainfall measurements at a gauging station and due to spatial interpolation. Errors in rainfall measurements at a gauging station are difficult to obtain due to lack of a simple error model. The errors due to spatial interpolation are likely to dominate the total error in areal average rainfall (e.g., Renard et al., 2011). Therefore, the errors in rainfall measured at a gauging station are neglected in this study, and it is assumed that the errors in areal average rainfall exist solely due to spatial variability of rainfall. Several different models have been proposed to capture the spatial variation of rainfall such as cluster point Poisson processes (Waymire and Gupta, 1981a, b, c), random cascades (Gupta and Waymire, 1993), Kriging (Moulin et al., 2009), and conditional simulations (Renard et al., 2011). All these models treat rainfall as a random field in space-time domain. But most of these models are based on strict assumptions about the covariance of spatial rainfall or error structure which are not justifiable in practice. Even if the assumptions are approximately true, the rain gauge density is typically too small to reliably estimate the parameters of the covariance function. This issue is further complicated as the covariance structure may vary from event to event in unknown ways, depending upon the type of event. Therefore, in this study, an empirical approach was used to get an estimate of the uncertainty in areal average rainfall.

There were 6 rainfall gauging stations near the SJRW (locations on these stations are shown in Fig. B1) at which daily timescale data were available. Typically, data from the available rain gauges are used to compute a single areal average rainfall time series using the Thiessen polygon interpolation method. In this study, all the  $63=(2^6-1)$  different combination of the 6 rain gauges were used to produce 63 realizations of areal average rainfall using the Thiessen polygon method. These 63 realizations represent an estimate of uncertainty in areal average rainfall.

# 2.5. Uncertainty bounds using runoff ratio method

The QRF method does not allow one to incorporate a hydrologists' knowledge about a watershed to construct the measurement uncertainty bounds. One method that allows incorporation of such knowledge was proposed by Beven (2019) using runoff ratios of observed rainfall-runoff events. In this method, only the observed rainfall-runoff data (along with evaporation data) of the watershed in question are used to create LOAs. This method was used to derive LOA estimates that were then compared to the LOAs estimated by the ORF algorithm.

In the first step, the observed rainfall-runoff data were separated into different rainfall-runoff events. This kind of hydrograph separation requires estimation of the recession curve. To this end, the master recession curve (MRC) technique was used (Lamb and Beven, 1997) - MRC is a characteristic recession curve of the watershed (Tallaksen, 1995). Once an MRC is defined, the streamflow time series can be divided into different rainfall-runoff events. In this study, a rainfall value below 1mmday<sup>-1</sup> was considered negligible, and a new rainfall event was assumed to start if the rainfall was negligible for more than 7 consecutive days. For example, a new rainfall event started at time-step  $t_n$  if the rainfall values at the time-steps  $t_{n-1}$ , ..., and  $t_{n-7}$  were less than 1mmday<sup>-1</sup>. The streamflow hydrograph corresponding to each rainfall event was assumed to start at the beginning of the rainfall event and end just before the start of next rainfall period. Next, MRC was appropriately appended at the end of the streamflow hydrograph for each rainfall-runoff event. The number of rainfall-runoff events, thus obtained for four of the stations in SJRW, are listed in Table 4.

In the second step, the runoff ratio of each event was computed as the ratio of the total volume of event streamflow to the total volume of event rainfall, where 'event streamflow' refers to streamflow time series obtained after appending the MRC. This resulted in an ensemble of runoff ratios. In the third step, LOAs were computed over each of the rainfall-

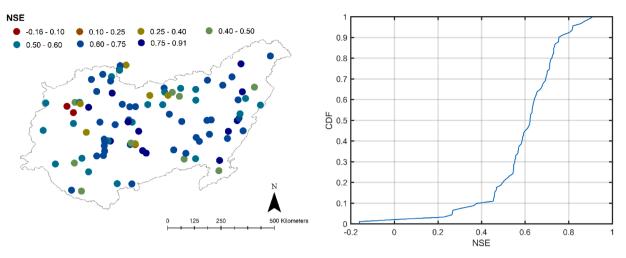


Fig. 3. (a) Spatial distribution of NSE values for the test set including ORB and SJRW station, and (b) cumulative distribution function (CDF) of the test NSE values. These NSE values were derived from ungauged model.

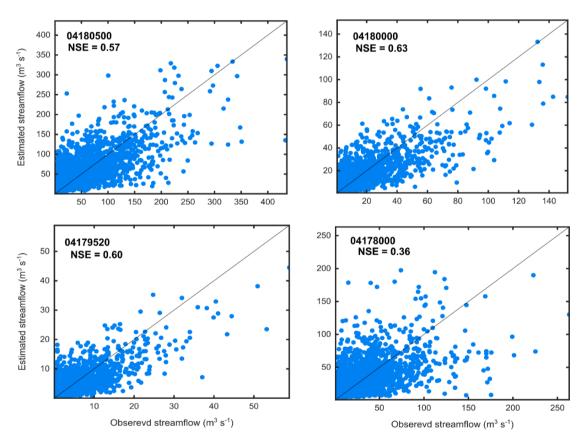


Fig. 4. Observed vs. estimated streamflows at four stations in St. Joseph River Watershed (SJRW). The estimated streamflow values were derived from ungauged model.

runoff events in an iterative manner. To construct the LOAs over the  $i^{th}$  event, the events in the ensemble similar to the  $i^{th}$  event were identified based on antecedent moisture condition and total volume of rainfall during the event. As an estimate of the antecedent moisture conditions, initial streamflow of the event was used. Thus, the events that were closest to the  $i^{th}$  event were identified by using the Mahalanobis distance between the events using these two variables (this is the k-nearest neighbor approach used by Beven, 2019). Appropriate value of the Mahalanobis distance to define the closeness of two events is a subjective decision. In this study, we first computed the Mahalanobis distance of the  $i^{th}$  event from rest of the events, and, then normalized the distance

values to lie between 0 and 1. Now, events similar to the  $i^{th}$  event may be defined as the events that are  $d_{\rm M,N}$  distance away from the  $i^{th}$  event, where  $d_{\rm M,N}$  denotes normalized Mahalanobis distance. Several values of  $d_{\rm M,N}$  were used to analyze the impact of this threshold on uncertainty bound. After the completion of the third step, one obtains runoff ratios of the  $i^{th}$  event and those of other  $N_i$  events that are similar to the  $i^{th}$  event. In addition to the k-nearest neighbor approach, we also used decision tree approach to group similar events again based on antecedent moisture condition and total rainfall volume. In what follows, the abbreviations RR-KNN and RR-QRF will be used to refer to runoff ratio method applied using k-nearest neighbor method and QRF method, respectively.

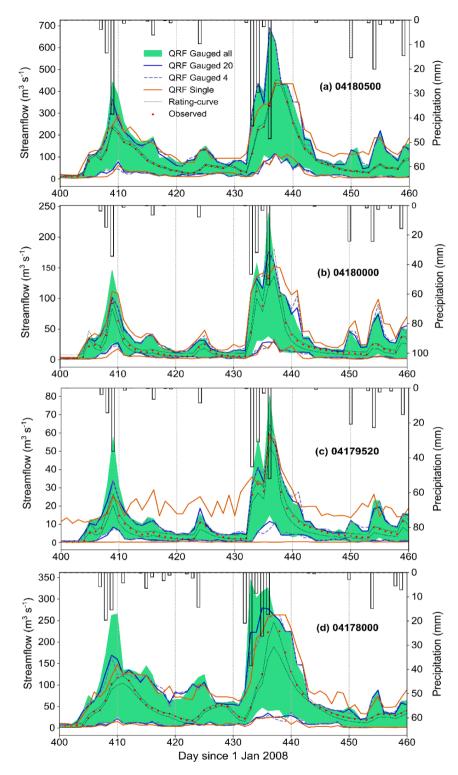


Fig. 5. LOAs obtained by quantile random forest (QRF) in different gauged scenarios: using all the training watersheds (green band), 20 most similar watershed including the four SJRW watersheds (blue lines), 4 SJRW watersheds (blue-dash lines), and gauged single model (orange-solid lines). Uncertainty bounds obtained by rating curve analysis (black-dash), and observed streamflow (red dots), along with precipitation are also shown.

In the fourth step, the streamflow time series of the  $i^{th}$  event was divided by its runoff ratio  $C_i$ , thus yielding a zero-loss streamflow time series of the  $i^{th}$  event that would have been observed if the runoff ratio of the  $i^{th}$  event was equal to 1. The zero-loss streamflow time series was then multiplied by the largest and smallest runoff ratios to obtain upper and lower bounds of LOA. In RR-KNN method, the largest and smallest runoff ratios were identified among the  $N_i$  runoff ratios of the events

similar to the  $i^{\rm th}$  event. In RR-QRF approach, the largest and smallest runoff ratios were the 100th and 0th percentiles in the leaf node to which the  $i^{\rm th}$  event belonged. RR-QRF approach is more objective than the RR-KNN approach since the value of  $d_{\rm M,N}$  needs to be specified subjectively in the latter. However, specification of appropriate percentiles in RR-QRF incurs some subjectivity.

**Table 5**Fraction of observations enveloped by the QRF LOAs.

	QRF ungauged	QRF gauged	QRF gauged-single
04180500	0.97	1.00	0.99
04180000	0.97	1.00	0.99
04179520	0.94	1.00	0.99
04178000	0.96	1.00	0.99

# 3. Experiments with rainfall-runoff data

# 3.1. An evaluation of decision tree (DTs) in terms of predicting streamflow

Fig. 3 shows the NSE values obtained by the RF ungauged model for the watersheds contained in the test set. NSE was greater than 0.60 for 55% of the watersheds and was greater than 0.5 for 80% of the test watersheds. There were some systematic patterns in the spatial distribution of NSE values. NSEs were typically higher in the eastern part of the basin than those in the western part. Most watersheds in the eastern ORB had NSEs greater than 0.5. For about 20% of all the test watersheds, the NSE was less than 0.5. It is likely that the RF algorithm could not identify the rainfall-runoff relationship in these watersheds, possibly because the hydrological behavior of these watersheds is not represented in the data. Overall, the performance of the RF model was deemed acceptable for majority of the watersheds for which NSE was greater than 0.50. It captured the rainfall-runoff dynamics in the sense that its response to input rainfall is hydrologically consistent. The term 'hydrologically consistent' is used to refer to an expected behavior of hydrological models: increasing streamflow with increasing rainfall under similar antecedent conditions. One question is if QRF model can be used to construct LOAs in a watershed where the NSE is low. We note that low NSE value can also be due to errors in streamflow or rainfall data. But still the LOAs obtained for these watersheds may not be reliably used for model inference. Fig. 4 shows the observed and predicted streamflow for the four stations located in SJRW. NSE was close to 0.6 for the three of the stations but was poor (=0.36) for station 04178000. These values seem adequate for constructing measurement uncertainty bounds except for station 04178000.

# 3.2. Limits-of-Acceptability (LOA) constructed by the QRF models

Fig. 5 shows the LOAs obtained by the QRF models trained under the first two scenarios (gauged-single and gauged) along with the uncertainty bounds obtained by the rating-curve analysis. Since rating curve analysis yields uncertainty due to errors in streamflow measurements only, LOAs obtained by QRF should envelop the uncertainty bound obtained by rating curve analysis as shown in Fig. 5. A similar observation was made for the majority of the cases (Table 5). Among the different QRF models (QRF-gauged-all, QRF-gauged-20, QRF-gauged-4, QRF-single), the LOAs obtained by the QRF-gauged models were widest and the LOAs obtained by the QRF-gauged-20 and QRF-gauged-4 models were typically close to each other. The QRF-single model yields very narrow LOAs at the two peaks shown (at time-steps 410 and 438). These two peaks are among the highest flow values observed in these watersheds implying that more data are required to construct reliable LOAs for these peaks. This illustrates the practical difficulty in constructing LOAs and highlights the need to allow for outliers when LOAs are used for model inference. There would not be enough data to estimate LOAs for events with return period greater than 2 to 10 years in many instances. The LOAs obtained by the three QRF-gauged models (ORF-gauged-all, ORF-gauged-20, ORF-gauged-4) were very similar except at a few time steps. As mentioned above, the 4 and 20 most similar watersheds to train the QRF model were identified using static watershed attributes. These static attributes are already used by the QRF method to partition the data into leaf nodes, which explains the similarity of LOAs obtained by the three QRF-gauged models.

The uncertainty bound obtained by rating curve analysis was significantly narrower at most of the time-steps indicating that errors in rainfall measurements contribute more to measurement uncertainty than do the errors in streamflow measurements. But the streamflow uncertainty bounds shown in Fig. 5 were obtained by making aleatoric assumptions. The peak streamflow values may contain larger uncertainties. Fig. 6 shows the fraction of peaks enveloped by upper bounds of LOAs if the observed peak magnitude were multiplied by a factor f. As the multiplier f increases, the fraction of peaks enveloped by the QRF uncertainty bound decreases. This decrease, however, occurs at different rates for the three models. Interestingly, the fractions of multiplied peaks enveloped by the LOAs were larger for the gauged-single model than the ones obtained by the gauged-all model. This is likely from timing errors in precipitation data as discussed below. The typical errors in peak streamflow have been reported to be 20–40% (Di

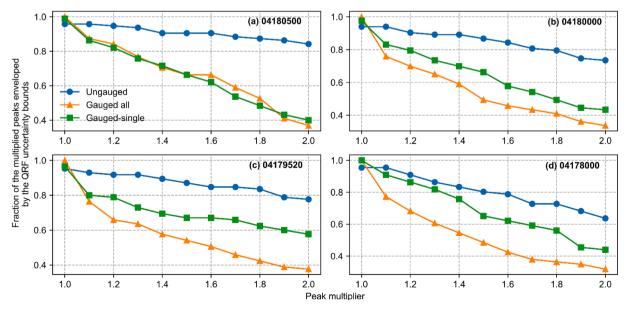


Fig. 6. Fraction of the peaks enveloped by the QRF uncertainty bounds if the observed peaks were 10-100% greater.

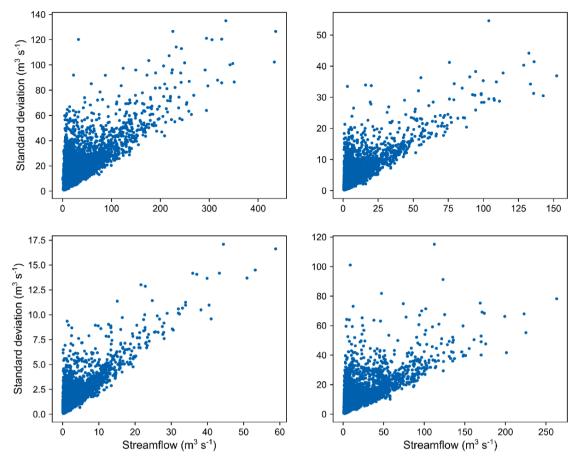


Fig. 7. Standard deviation of streamflow time series obtained by RF method plotted against observed streamflow data. The standard deviation increases with increase in streamflow value.

Baldassarre and Montanari, 2009); Fig. 6 shows that more than 55% of the peaks were enveloped in these ranges of errors by all the three models. Even for 100% errors, more than 30% of the peaks are enveloped by the QRF LOAs across the three models.

One of the characteristics of the LOA by the QRF method (Fig. 5) is that it is very wide at the time-steps corresponding to streamflow peaks and narrow at the time-steps where streamflow is small. Although not shown here, this pattern was visible throughout the study period. Fig. 7 shows the standard deviations of streamflow obtained by QRF method plotted against streamflow. The standard deviation increases as streamflow value increases in keeping with how rainfall uncertainty typically propagates to streamflow uncertainty (Moulin et al., 2009; Renard et al., 2011). These observations suggest that QRF is able to account for the effect of uncertainty due to rainfall and streamflow measurement errors.

One seeming discrepancy to the pattern discussed above is the wide LOA obtained by the QRF-gauged method between time-steps 410 and 420 even when the streamflow time series is in recession phase (Fig. 5) – this is especially the case for the stations 04180500 and 04178000. Data show that some rain did fall over the watershed at these time-steps (Fig. 5), and this rain event was similar in magnitude to the rain event that generated the streamflow peak at time-step 424. One possibility is that this rain event did not result in streamflow due to spatial location of the event (rain event might be far from the watershed outlet). The second possibility is that the rainfall measurement at the gauging station is erroneous. The third source of error is the unknown true intensity of rainfall: The observed rainfall data are at daily timescale and two events with same intensity at the daily timescale may have very different intensities at sub-daily timescales which will result in different hydrographs. These are examples of epistemic errors, and the exact reason for

these errors is difficult to know. In fact, we do not even know whether the measurement is actually erroneous. A good hydrological model forced with this rain event and uninformed by *true* spatial distribution and true intensity of rainfall will still generate a streamflow event (if the antecedent conditions allow). It would be unwise to reject this model if these errors indeed exist. This illustrates how QRF can account for epistemic errors. Similarly, at time-step 450, a wide LOA was obtained by the QRF method for three of the stations whereas streamflow time series is in recession phase. Again, a rainfall event was observed at this time-step which apparently did not result in a streamflow peak, and the same arguments apply.

In some of the events, timing errors between observed peak and QRF simulated peak were observed primarily in the LOAs created by the QRF-gauged model. An example of such timing errors may be seen at time-step 438 in Fig. 5. These timing errors occurred for less than 20 events per watershed (see also Fig. 11 where LOAs for a few other time-steps are also shown). For the five peak events shown in Fig. 5, timing error occurs only for one event for the three stations 04180500, 04180000, and 04179520. Out of the two major peaks at time-steps 410 and 438, timing errors at time-step 438 are present for stations 04180500 and 04180000.

There seem to be two possibilities behind these timing errors: (1) disinformation introduced by the data from other watersheds, or (2) timing errors in rainfall data. For the stations 04180500 and 04180000, there is zero lag between rainfall and QRF obtained streamflow peak at time-step 410. Meanwhile at time-step 438, a lag of 1–2 days between rainfall and streamflow peak is observed. Further, the rainfall event at time-step 438 is more intense (at daily timescale) and one would expect a smaller lag between rainfall and streamflow peaks for this event compared to the lag observed for the event at time-step 410. Therefore,

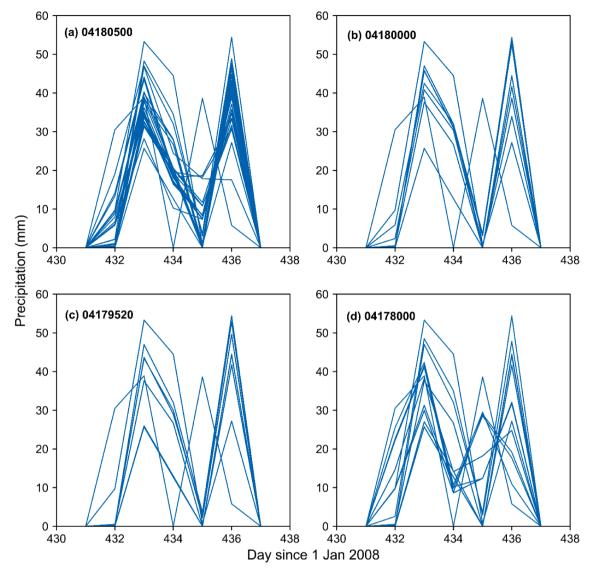


Fig. 8. Different realizations of areal average precipitation.

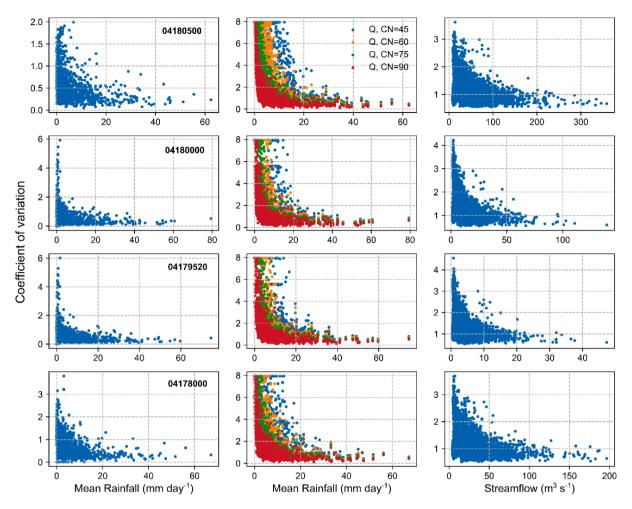
it seems more likely that the computed areal average rainfall has timing errors for this event. We note that it is also possible that the sub-daily timescale intensity of the event at time-step 438 was low which would justify the delay in peak. This is again an example of epistemic uncertainty.

The same arguments apply for the timing errors observed at the station 04178000, especially at time-step 424 where the lag between computed areal rainfall peak and observed streamflow is 3 days. The timing errors at the station 04178000 were more frequent which is partly the reason for poor validation NSE value at this station (Fig. 4). It is worth noting that the timing errors between LOAs constructed by QRF-single model and observed streamflow were typically absent. It is possible that the model has compensated for timing errors in precipitation.

Potential for the timing errors in rainfall around time-step 438 is also illustrated in Fig. 8 which shows all the different realizations of the areal average rainfall. For the majority of the realizations, the second precipitation peak occurs at time-step 436 while for a few realizations the second peak occurs at time-step 435. These realizations were constructed using six gauging stations which are located outside but near the SJRW watershed (Fig. B1). If data from more stations were available, some of the realization might have very well shown the second peak at the time-step 437.

It is worth noting that the information about consistent timing error may not be revealed by QRF-single model as it will learn this as a behavior of the watershed. Thus, this analysis illustrates the usefulness of data from different watersheds in constructing LOAs. Further, the analysis also illustrates how the LOAs constructed using decision trees may potentially capture the effect of timing errors. It is possible that the timing errors between the observed and QRF (gauged model) simulated peaks occur because of disinformation introduced by data from other watersheds. Therefore, it seems prudent to construct multiple LOAs using data from different sets of watersheds and use a combination of these LOAs for model inference.

Fig. 9 shows the  $CV_R$  (coefficient of variation) of areal average rainfall obtained by using the empirical approach described above. The  $CV_R$  values decrease as areal average rainfall increases, at all the stations. At first one may attribute this behavior to standard deviation of rainfall being constant irrespective of the mean rainfall value. However, it was observed that standard deviation of areal average rainfall increases with increasing mean rainfall values (now shown) similar to the standard deviation of streamflow. The CV values of excess rainfall, obtained by SCS-CN method, also follow the same pattern as areal average rainfall. But the  $CV_S$  corresponding to excess rainfall were typically higher than the  $CV_S$  corresponding to areal average rainfall. The difference between excess and areal average rainfall  $CV_S$  become smaller



**Fig. 9.** Coefficients of variation (CV) of areal average rainfall (left), excess rainfall for different values of *CN* (middle), and the CV of streamflow obtained by RF in ungauged scenario (right). In the legend, *Q* refers to excess rainfall obtained by using SCS-CN method for different value of the parameter *CN*. Each row refers to one basin.

for higher values of areal average rainfall. Many of the small non-zero areal average rainfall values produce no excess rainfall; increased number of zeros in excess rainfall increases the CV.

Fig. 9 shows that variation of CV<sub>O</sub> with streamflow follows the same pattern as that of variation of  $CV_R$  with mean areal average rainfall;  $CV_Q$ decreases as mean streamflow increases. For all the four stations, the magnitudes of CVs are of similar order for the areal average rainfall and streamflow time series. Another pattern in  $CV_R$  plots is that there is a larger (smaller) scatter in these values when mean rainfall is small (large). The same pattern can be seen in streamflow values also. The rainfall time series is transformed non-linearly through a watershed to yield streamflow. The same rainfall event can result in very different streamflow hydrograph depending upon the spatial distribution of rain within the watershed and antecedent moisture conditions. Thus, for a given rainfall magnitude, many different values of streamflow are possible which explains the larger scatter in CV<sub>O</sub>. Fig. 9 indicates that the statistical structure of RF uncertainty bound reflects the effect of rainfall uncertainty. Overall, these results combined with the results discussed above indicate that the DTs could account for the effect of uncertainty due to errors in rainfall and streamflow measurements.

Further, it can be argued that any model with heteroscedastic error structure would result in uncertainty bounds as shown in Fig. 5. The QRF method does not enforce heteroscedastic error structure, rather this error structure was identified by the algorithm from the data. The experiments with synthetic data showed (results not shown) that if the errors are homoscedastic, QRF produces homoscedastic error structure,

and if the errors are heteroscedastic, QRF produces a heteroscedastic error structure. LOAs shown in Fig. 5 do not represent measurement uncertainty only – it is likely that structural errors of QRF model are also contributing to these bounds.

# 3.3. How do QRF LOAs compare to the LOAs obtained by the runoff ratio method?

Fig. 10 shows the LOAs obtained by the runoff ratio method, along with the ensemble of runoff ratios at four of the gauging stations in SJRW. Ideally, the runoff ratios should lie between 0 and 1. The errors in rainfall and streamflow measurements, and inexactness of hydrograph separation method, however, may result in values of runoff ratios greater than one (Beven and Westerberg, 2011). Indeed, a few rainfall-runoff events had runoff ratio values greater than 2 which are likely to have occurred due to significant biases in rainfall measurements. These periods can be referred to as disinformative periods (Beven and Westerberg, 2011) which should not be used for parameter estimation and uncertainty analysis. In this study, however, these events were kept for further analysis as the final aim is to compare the bounds obtained by different methods. It may be noted that QRF will not recognize such disinformative periods but it will yield appropriate uncertainty bound for these events making it unlikely that a good model will be rejected by using the LOAs obtained by the QRF algorithm even if it includes disinformative periods. For example, if a rainfall event has large negative bias, QRF will identify this event as similar to other events

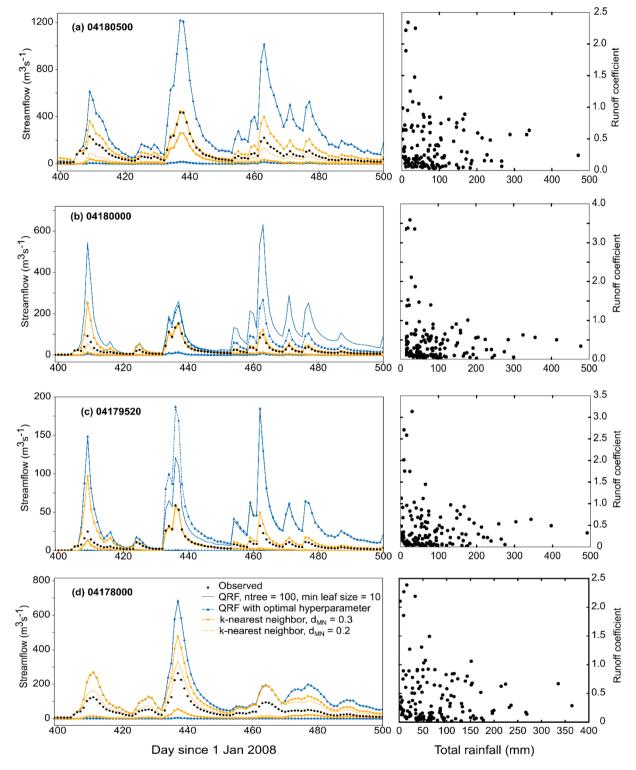


Fig. 10. LOAs obtained by runoff ratio method (left) and runoff ratios plotted against total rainfall of the each of the rainfall-runoff events (right).

with small rainfall and the LOAs for this event will span a large range of streamflow values.

Fig. 10 shows the LOAs obtained by using the runoff ratio method where similar events were selected using KNN method (with two different distance thresholds  $d_{\rm M,N}=0.2$ and 0.3) and by using QRF method. One expects the LOAs to envelop all the observations and the uncertainty bounds to become wider as the value of  $d_{\rm M,N}$  increases. This is indeed observed in Fig. 10 with the following special case: the observations coincide with the upper LOA at a few time-steps for small  $d_{\rm M,N}$ 

values. These cases occur because of the small number of rainfall-runoff events available at a station and even smaller number of similar rainfall-runoff events; this prohibits the construction of robust LOAs. LOAs obtained by RR-QRF method were typically wider than the those obtained by the RR-KNN method which is partly a consequence of using 0% and 100% percentile values of data in the leaf node for defining these bounds (see Section 2.5).

QRF-gauged algorithm yielded tigher LOAs compared to those obtained by runoff ratio method for a few time-steps (Fig. 11). But at other

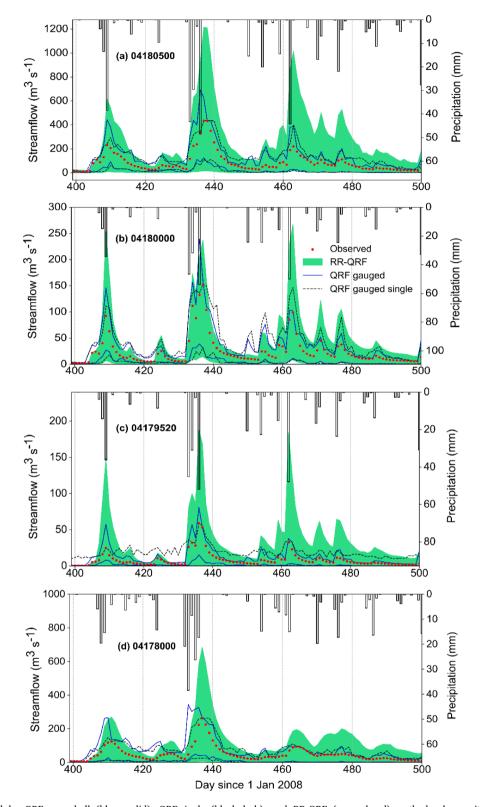


Fig. 11. LOAs obtained by QRF-gauged-all (blue -solid), QRF-single (black-dash), and RR-QRF (green band) methods along with observed precipitation and streamflow.

times-steps, e.g., between 400 and 420, the QRF LOAs were wider. There is one general similarity between the LOAs obtained by QRF and runoff ratio method: the width of both LOAs increase or decrease almost synchronously in time (except for a few timing errors, see above for a discussion of this issue). This gives us further confidence that the LOAs obtained by QRF are able to capture general patterns of measurement

uncertainty. If the patterns of LOAs obtained by QRF and runoff ratio method were significantly different, that would have disproved the usefulness of QRF in constructing LOAs.

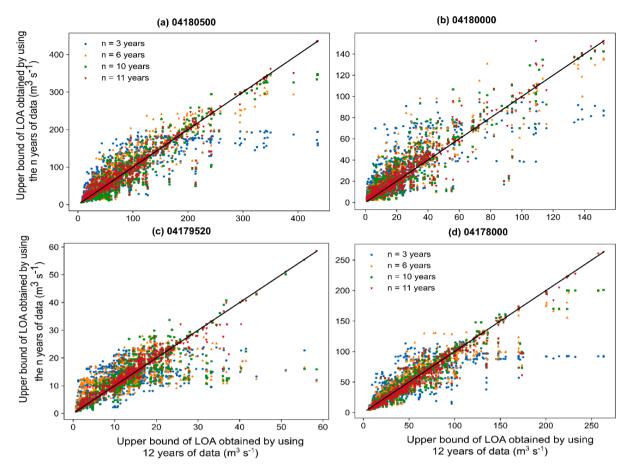


Fig. 12. Convergence properties of LOAs obtained by the QRF algorithm.

# 3.4. Convergence of LOAs obtained by QRF algorithm

To test the convergence properties of QRF estimated LOAs with increasing length of data, several QRF models were developed using different lengths of training data. In these experiments, data from only that watershed where LOAs are to be constructed were used, i.e., gauged-single models were developed. For each of the four test watersheds, 12 different gauged-single models were developed using 1, 2, ..., 12 years of data. Fig. 12 shows the 97.5th percentiles of LOAs that obtained using different amounts of data. For three stations (04180500, 04180000, and 04178000), LOA estimates at high flow time-steps started to converge when more than three years of data were used, but there were a few high flow time-steps where LOAs did not converge. At station 04179520, the convergence of LOAs seems to be much slower than the convergence at other stations. LOAs appear to be converging for low flows as well but more data are required to achieve the final bounds.

# 3.5. Limits-of-acceptability (LOA) created using the QRF ungauged model

One of the major advantages of the QRF algorithm is that it can be used to construct LOAs at ungauged locations. Fig. 13 shows the LOAs constructed by the QRF ungauged model, along with LOAs constructed by the other models for comparison. The LOAs obtained by the QRF-ungauged model were typically wider than the LOAs obtained by the other models. The timing errors between LOAs and observed streamflow can also be observed for the QRF-ungauged model.

At time step 406, there exists a widening of LOAs along with a very small peak in observed streamflow, but the observed precipitation is either zero or negligible. This is clearly because of an error in precipitation magnitude. It is likely that there was a small amount of precipitation in the watershed which was not recorded by the precipitation

gauges. There were a few other such events where very small observed precipitation corresponded to a significant observed streamflow resulting in very high runoff ratios (as discussed above). Therefore, depending upon the precipitation magnitudes during current and previous time steps, QRF predicts a peak in streamflow. Such peaks would not have any impact on model inference in the sense that a hydrological model would not produce streamflow peaks in the absence of rainfall and the simulated streamflows would always be enveloped by the LOAs at these time steps.

Fig. 6 shows that more than 60% of the multiplied peaks were enveloped by the QRF LOA even for 100% errors (f=2) for the ungauged model. The analysis suggests that LOAs obtained by the ungauged model are very conservative. This is desirable when the LOAs are to be constructed at an ungauged location so as to include a large number of rainfall-runoff behaviors. The results of this analysis are encouraging in terms of usefulness of QRF approach in creating LOAs at both gauged and ungauged locations.

# 4. Logic behind the proposed method

In this section, a mathematical argument is presented for using DTs for constructing LOAs. We hypothesize that if *infinite* amount of hydrological data are available, DT-estimated LOA will reflect the effect of uncertainty due to errors in rainfall and streamflow measurements. This is a hypothetical scenario (as infinite data are never available) but it serves to illustrate the usefulness of DTs in constructing LOAs and provides a theoretical basis. In practical cases, the DTs would also reflect variability due to other sources. As the number of calibration samples approaches *infinity*, the error incurred by a DT approaches optimal Bayes error (Denil et al., 2014) which is the irreducible part of the error due to inherent variability in the process and due to measurement errors (both

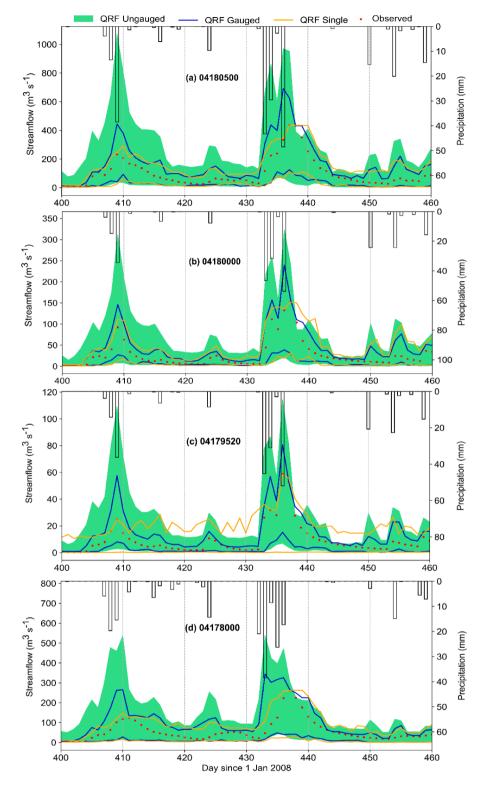


Fig. 13. LOAs obtained by quantile random forest (QRF) in ungauged scenario (green), by QRF in gauged-all scenario (blue), by QRF in gauged-single scenario (orange), along with observed streamflow and precipitation.

epistemic and aleatoric). Assuming, for the sake of discussion, that there is no inherent variability in the hydrological processes (more on this below), then errors incurred by a decision tree approach measurement error as the samples size increases. Thus, the results of Denil et al. (2014) suggest that decision tree can be used to account for measurement uncertainty, even if it holds only for the hypothetical case of infinite data. However, it may not be immediately clear how the uncertainty bounds

obtained by decision trees represent measurement uncertainty in case of infinite sample size. Here, we answer this question and elucidate the logic behind the proposed hypothesis. A formal analysis of the proposed hypothesis is provided in Appendix A.

First, consider the case where only the streamflow measurements are uncertain, and the rainfall measurements are free of errors. Further, assume that the errors in streamflow measurements are unbiased. As the

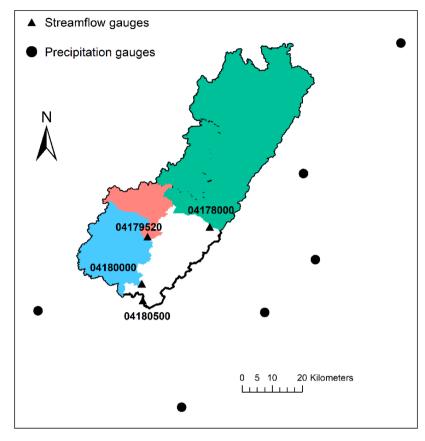


Fig. B1. Four sub watersheds located in St. Joseph River Watershed (SJRW) along with the precipitation gauges.

sample size increases, the diameter of each leaf node approaches zero, that is, predictor vectors contained in a leaf node are approximately equal (a formal proof of this statement is given in Appendix A). The true streamflow values corresponding to predictor vectors contained in a leaf node are approximately equal and any variations in the observed streamflow would be due to measurement errors. Thus, given an infinite sample, the minimum and maximum values contained in the leaf node represent lower and upper bounds over streamflow, and the difference between these bounds is due to measurement uncertainty. A formal analysis of this case is given in Section A.1.

Second, consider the case where only the rainfall measurements are uncertain, and the streamflow measurements are error free. In this case also, the diameter of a leaf node would approach zero (for the same reason as in the first case), and predictor vectors contained in a leaf node would be near identical, as the sample size approaches infinite. But, due to measurement errors, the underlying true values of predictor vectors contained in a leaf node would be different (more precisely, the projections of predictor vectors on rainfall subspace will be different). Since there exists a streamflow value corresponding to each true predictor vector, the set of streamflow values corresponding to true predictor vectors in a leaf node would represent the effect of measurement uncertainty in predictor vector on streamflow. A formal analysis of the second case is given in Section A.2.

Third, consider the case where both rainfall and streamflow measurements are corrupted by errors. The logic behind this case is similar to the logic discussed above for the first and second cases. A formal analysis of this case is given in Section A.3.

Finally, we elaborate on inherent variability in hydrological processes. The mathematical analyses provided above, and in the Appendix A, implicitly assume that the predictors variables used to train the decision tree are *complete* in the sense that predictor variables contain all the information that is required to predict streamflow. This, however, is

not possible since the physical structure of the watershed itself will be changing continuously, albeit only slowly with intermittent large disruptions, which will change the hydrological response of the watershed. This can be referred to as the inherent uncertainty in hydrological processes which is irreducible. Therefore, given an infinite sample, decision trees would also account for this inherent variability along with the measurement uncertainty.

Both measurement uncertainty and inherent variability are generally dominated by epistemic errors. Since, to construct LOAs, only the upper and lower bounds on errors are required for a given rainfall-runoff event, it is sufficient that the errors incurred in a given event fall in the range of the errors incurred from other similar events. Further, since the errors are epistemic and available data are finite in practice, it is possible that the errors of some events do not fall in the range of errors represented in the data; therefore, accommodation for such outliers needs to be made while using LOAs for model inference. Typically, 5% of the observations are allowed to fall outside the estimated uncertainty bound. In hydrological applications, these 5% outliers might well include the timesteps that one is most interested in (e.g., high flows for flood modeling). Therefore, a posteriori analysis of outliers should be carried out. A model can be declared unfit-for-purpose if all or most of the 5% outliers belong to the timesteps of interest. It is possible that all the models are rejected as unfit-for-purpose but nevertheless a model is required for some urgent practical application; in this case, some of the rejected models with least deviation from the LOAs might be used and the inverse of the magnitude of deviation can be used as the weight of that model in decision making. Alternatively, instead of defining fraction of outliers beforehand, one can report the accepted models for different fractions of outliers.

# 5. Summary and conclusions

Separation of structural and measurement uncertainty was recognized as one of the twenty-three unsolved problems in hydrology by Blöschl et al. (2019). The only way to address this problem is to estimate measurement uncertainty before model calibration. This is a difficult task given that statistical properties of rainfall and streamflow measurement uncertainty are poorly understood, especially those of rainfall measurements. There exist two dominant philosophies to address this problem: (1) to assume statistical distributions over measurement uncertainty due to both rainfall and streamflow errors, and (2) to construct limits-of-acceptability (LOA) that provide some bounds on measurement uncertainty before any modeling exercise. LOA has been used within the GLUE framework. However, both of these philosophies may also be combined together in Approximate Bayes Computation (ABC) framework. LOA can also be used in a purely Bayesian framework by defining a likelihood function that penalizes the simulations based on their deviations from the LOA defined through a suitable metric. The aim of this paper was to test the capability of decision tree algorithms in creating LOAs that provide meaningful bounds on measurement uncertainty.

In this study, quantile random forest (QRF) method was used to construct LOAs. The advantages of the QRF method are as follows: (1) it can reflect the effect of both precipitation and streamflow measurement uncertainty, (2) it can account for timing errors in precipitation, (3) it can be applied at the timescale of available data, and (4) it can be used to construct LOAs at ungauged catchments. The results show that the LOAs obtained by using QRF enveloped the uncertainty bounds over streamflow observations. Measurement uncertainty in streamflow due to aleatory variability was found to be very small. It was shown that the statistical structure of QRF uncertainty bound was similar to an uncertainty bound obtained by propagating rainfall uncertainty through a hydrological model. Some observations include:

- (1) Standard deviations of streamflow obtained by the QRF method increase with increasing values of observed streamflow.
- (2) CVs of simulated rainfall time series and QRF uncertainty bound follow the same pattern: they decrease with increasing value of rainfall and streamflow, respectively.
- (3) The general pattern of increase and decrease of width of uncertainty bound was similar for QRF and runoff ratio methods.

The QRF method does not contain any mechanism that induces the uncertainty bounds to follow any pre-determined patterns. Therefore, existence of these patterns suggests the QRF method is able to identify *some* of the characteristics of measurement uncertainty from data. We cannot conclude that all the characteristics of measurement uncertainty were identified because QRF is unable to extract all the hydrological information from available data for the four SJRW watersheds used as test cases in this study. Indeed, this is likely to be the case for most watersheds since data on all the factors determining the hydrological response of a watershed are not available.

A timing error between observed streamflow and the LOAs obtained by the QRF method was observed in all four test watersheds (Figs. 5 and 11) in gauged-all and ungauged cases. These timing errors are likely due to timing errors in precipitation data. Figs. 5 and 11 show that QRF can compensate for consistent precipitation timing errors in a watershed in gauged-single case. Thus, data from other similar watersheds can be useful in constructing LOAs that capture the effects of precipitation timing errors. In general, the shorter the length of data available to construct LOAs, the more the data from other similar watersheds will be required. The issue of choosing similar watersheds is discussed below. Another possible reason for timing errors in gauged-all case is that data from other watersheds may have introduced disinformation into the LOAs. Therefore, it appears that LOAs should be constructed using data from several sets of watersheds so that the effect of both the potential timing errors and disinformation can be accommodated. This will, in

general, mean a larger number of behavioral models and higher predictive uncertainty. Overall, the results of this paper indicate potential for the QRF approach for constructing LOAs at both gauged and ungauged locations.

In the *hypothetical* scenario, when infinite amount of hydrological data are available, the QRF algorithm can actually reflect the effects of measurement uncertainty as shown in the mathematical analysis in Appendix A. This analysis used the following main assumptions to prove the proposed hypothesis:

- The relationship between predictor and response variables is oneto-one.
- (2) The mapping between predictor and response variable is continuous.
- (3) The errors in predictor and response variables are unbiased but otherwise the errors could be either aleatoric or epistemic.
- (4) Error can be assumed independently and identically distributed within a leaf node.

We note that assumption 1 was made for mathematical convenience. A similar analysis can be carried out without this assumption. For a finite sample size, the uncertainty bounds obtained by a decision tree include contributions from structural uncertainty (of QRF method) along with measurement uncertainty.

A major advantage of QRF method (and indeed the LOA approach) is that it is a non-parametric approach for constructing LOAs and does not resort to strong assumptions on the statistical nature of streamflow and rainfall measurement errors. Overall, the QRF method offers promise as a powerful tool in hydrological model inference.

Rainfall-runoff data may also contain disinformative periods. To identify disinformation and biases, one requires physical understanding of the rainfall-runoff processes. Runoff ratio method is an example of using process-based knowledge to identify biases, but it is not applicable for baseflow dominated catchments and cannot be applied at ungauged locations. Moreover, runoff ratio method can identify the effect of errors in streamflow and precipitation volume – it cannot identify precipitation timing errors. QRF method addresses these limitations of the runoff ratio method. QRF will not explicitly identify disinformative periods, but it will likely define LOAs for the disinformative periods such that a good model would not be rejected because of these periods.

Further, as noted, it is possible that data from other watersheds introduce disinformation into the constructed LOAs. An interesting future problem in this respect would be to combine QRF method with catchment similarity analysis such that data from only the watersheds which are known to be hydrologically similar to the parent watershed (where LOAs are be constructed) are used. This would potentially reduce the disinformation introduced by the data from other catchments while yielding meaningful LOAs. This technique can be particularly useful for prediction in ungauged basins. In this paper, catchment characteristics (in the form of spatially averaged indices such as mean slope, mean soil properties etc.) were used in the QRF method to identify similar catchments. However, methods based on hydrological process understanding (e.g., Wagener et al., 2007) may prove to be better at identifying similar catchments.

One can also use other ML algorithms for creating LOAs in addition to the QRF method. Given a finite amount of data in practical applications, different algorithms would extract different information from available data and hence a different estimate of LOAs would be obtained. A combination of these different LOAs will be more desirable for model inference (a problem to be explored in future).

# Data availability

All the data used in this work are publicly available and can be downloaded from the DOI https://zenodo.org/record/7697209#. ZAJTxh\_MKUk

# CRediT authorship contribution statement

**Abhinav Gupta:** Conceptualization, Data curation, Formal analysis, Methodology, Software, Validation, Visualization, Writing – original draft, Writing – review & editing. **Rao S. Govindaraju:** Conceptualization, Project administration, Resources, Supervision, Visualization, Writing – review & editing. **Pin-Ching Li:** Methodology, Software. **Venkatesh Merwade:** Writing – review & editing.

# **Declaration of Competing Interest**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence

the work reported in this paper.

# Acknowledgements

An earlier version of the paper was substantially revised based on the review comments of Keith Beven. We are grateful to him for his insightful comments on the paper. Part of this manuscript was prepared when AG was a doctoral candidate at Purdue University. The rest of the manuscript was written when AG was a Maki Postdoctoral Associate at DRI where he was supported by DRI's Postdoctoral Support funds and Sulo and Aileen Maki postdoctoral fellowship. This support is gratefully acknowledged. PCL was supported by NSF award #1835822 and 2118329 during the period of this work.

# Appendix A: Mathematical analysis of the proposed hypothesis

In this section, a heuristic mathematical analysis in the support of the proposed hypothesis is provided. The aim of the analysis is to clarify the assumptions behind the hypothesis and limitations in practical implementation. Specifically, we show why the data in leaf nodes of a decision tree can be used to capture measurement uncertainty and under what condition structural uncertainty would be small. The analysis is divided into three parts for convenience: (1) when measurement errors occur in streamflow measurements only, (2) when measurement errors occur in rainfall measurements only, and (3) when both rainfall and streamflow measurements incur errors. We note that the analysis provided below is valid for both aleatoric and epistemic errors.

# A.1. Case 1: Only streamflow measurements are uncertain

First, we provide the analysis of the proposed hypothesis under the restriction that only the streamflow measurements contain errors and rainfall measurements are free of errors. Let  $\mathscr X$  denote the predictor space,  $x \in \mathscr X$  denotes a point in the predictor space, and  $\mathscr N_d(x)$  denote the d-neighborhood of x in  $\mathscr X$  where d is a suitable distance metric. Further, let us define by  $\mathscr Y$  the set containing error-corrupted value of a response variable as

$$\mathscr{Y} = \{y(x)|x \in \mathscr{X}\}\tag{A1}$$

Since y(x) is an error-corrupted value, it can be written as

$$y(x) = y_1(x) + \epsilon, \tag{A2}$$

where  $y_t(x)$  denotes the true but unobserved value of the response variable and  $\varepsilon$  denotes the measurement error in y. Here,  $\varepsilon$  represent a general error term which can be a function of x and/or y.

The data contained in a leaf node of a decision tree may be approximated as a neighborhood of the points close to its center. For example, if a leaf node constitutes the set  $\mathscr{X}_k = \{x_i | x_i \in \mathscr{X}\}_{i=1}^k$ , and the point  $x_m \in \mathscr{X}_k$  is close to its center; then  $\mathscr{X}_k$  can be treated as a neighborhood of  $x_m$ . To define a neighborhood, a distance metric is needed, and distance metric chosen defines the shape of neighborhood. In the analysis presented below, a different distance metric might be required for different leaf nodes of the decision tree. This does not pose any challenge to the generality of the analysis. The approximation of a leaf node by the d-neighborhood is made for the sake of mathematical convenience so that the analysis is manageable. Similar assumptions have been made by other authors (e.g., Denil et al., 2014).

Assumption 1. The mapping between predictor and response variables is continuous.

**Assumption 2.** The relationship between probability distribution of  $\varepsilon$  with x and y does not change significantly in a c-ball,  $\mathcal{D}_{\varepsilon}(x)$ ,

$$\mathscr{B}_c(\mathbf{x}) = \{x_i | d(\mathbf{x}, \mathbf{x}_i) \le c\},\tag{A3}$$

where c is a sufficiently small number. In other words, the distribution of  $\varepsilon$  changes slowly over  $\mathscr{X}$ .

**Assumption 3**. Without loss of generality, we assume that the relationship between true values of predictor and true values of response variables is one-to-one. This assumption is also made for analytical convenience.

**Assumption 4**. The expected value of  $\varepsilon$  is zero.

**Assumption 5**. The response variable *y* varies smoothly with the predictor variable *x*. This is particularly true for rainfall runoff models where unit increase in rainfall can result in a maximum of unit increase in streamflow, all else being equal.

For every 
$$\mathbf{x}_d \in \mathcal{N}_d(\mathbf{x})$$
, there exists a  $\mathbf{y}_d \in \mathcal{Y}$  by definition of  $\mathcal{Y}$ . By virtue of Eq. (A2),  $\mathbf{y}_d(\mathbf{x}_d) = \mathbf{y}_t(\mathbf{x}_d) + \epsilon$ . Define  $\mathcal{Y}_d$  as  $\mathcal{Y}_d = \{y_d(\mathbf{x}_d) | \mathbf{x}_d \in \mathcal{N}_d(\mathbf{x})\},$  (A4)

and define  $\mathcal{Y}_{d,t}$  as

$$\mathcal{Y}_{d,t} = \{y_t(\mathbf{x}_d) | \mathbf{x}_d \in \mathcal{N}_d(\mathbf{x})\},\tag{A5}$$

Further, define the quantity

$$\overline{y}_d(\mathbf{x}) = \frac{1}{Vol\{N_d(\mathbf{x})\}} \int y_d(\mathbf{x}) d\mathbf{x}$$
(A6)

where Vol denotes volume.

**Assertion 1.** : The quantity  $\bar{y}_d(x)$ , defined in Eq. (A6), approaches the true value  $y_t(x)$  as the number of samples increases.

The proof of this assertion, along with technical conditions, can be found in Brieman et al., (1984) and Denil et al. (2014). These references do not directly consider errors in measurement, but the proofs provided in these references are still valid provided assumption 4 holds. If assumption 4 is not valid, then the prediction error obtained by a decision tree approaches the optimal Bayes error. Note that the discrete version of the Eq. (A6) is the response variable estimated by the RF algorithm. Therefore, the structural errors in RF estimate would decrease arbitrarily as the sample size increases.

**Assertion 2.** : The diameter of the  $\mathcal{Y}_{d,t}$  is small, if the sample size is large. In other words, the maximum difference between the  $y_t$  values contained in  $\mathcal{Y}_{d,t}$  would be small. Let this difference be denoted by  $dia(\mathcal{Y}_{d,t})$ .

We note the following

- a decision tree aims to create leaf nodes so as to minimize some measure of prediction error (such as mean-square error) on test set,
- the estimated response by the decision tree is the average of the response values contained in a leaf node given by Eq. (A6), and
- the leaf nodes create a partition of the predictor space  $\mathcal{X}$ , i.e., the subsets created by the leaf nodes are disjoint and cover the predictor space.

These requirements are met only if the quantity  $dia(\mathcal{Y}_{d,t}(x))$  is small for each x. (Here,  $\mathcal{Y}_{d,t}$  is denoted as a function of the argument x.) For, consider n points  $x_1, x_2, ..., x_n \in \mathcal{X}$  that constitute the training set with corresponding neighborhoods  $N_d(x_1), N_d(x_2), ..., N_d(x_n)$ . Denote the number of leaf nodes created by the decision tree by m. Clearly,  $m \le n$ . Further, consider the expression for mean-square error,

$$MSE_n = \frac{1}{n} \sum_{i=1}^{n} \{ y(\mathbf{x}_i) - \bar{y}(\mathbf{x}_i) \}^2,$$
 (A7)

where  $\overline{y}(x_i)$  is estimated response given by Equation (23). The expression (A7) is minimized when each term in the summation is minimized.

If  $m \ll n$ , there will be many out of n points that would fall into the same leaf node and, therefore, will have identical estimate of the response. Thus,  $MSE_n$  would not be minimized. This seems to imply that for  $MSE_n$  to be minimized we need m=n. Due to measurement errors, however, minimization of  $MSE_n$  on training set may not result in minimization of  $MSE_n$  on test set. And making m=n is likely to result in overfitting. Therefore, to satisfy the three conditions above), the value of m must be less than n but not much smaller than n. As n increases, m should also increase; otherwise, m would become much smaller than n. (Technically, this condition translates to the following:  $m \to \infty$  and  $m/n \to 0$ , as  $n \to \infty$ ). In decision tree language, as n increases, the predictor space would be split into smaller and smaller partitioning subregions, i.e., diameter of the leaf nodes would become smaller and smaller. Hence, it follows that  $dia(N_d) \to 0$ , as  $n \to \infty$ .

If diameter of  $N_d(x)$  is small, then by assumption 5 and the assumption that values in  $N_d(x)$  are error free, the  $dia(y_{d,i})$  is also small.

In summary, if the sample size is large, then the decision tree would be able to create small leaf nodes in order to minimize mean-square error. More technically, for  $n > N_a$ , and  $\delta > 0$ 

$$dia(\mathcal{Y}_{d,t}) < \delta, \tag{A8}$$

where  $N_a$  is some arbitrary large value.

**Theorem 1.** The set  $\mathcal{Y}_d$  approximately captures measurement uncertainty in response variable if the sample size is large.

**Proof.** The minimum value contained in  $\mathscr{Y}_d$  is greater than or equal to  $\min(\mathscr{Y}_{d,t}) + \varepsilon_l$  and the maximum value contained in  $\mathscr{Y}_d$  is less than or equal to  $\max(\mathscr{Y}_{d,t}) + \varepsilon_l$ . Here,  $\varepsilon_l$  denotes a value in the left tail of the distribution of  $\varepsilon$  such that probability of  $\varepsilon$  taking a value less than or equal to  $\varepsilon_l$  is  $\gamma_l$ . Similarly,  $\varepsilon_u$  denotes a value in the right tail of the distribution of  $\varepsilon$  such that probability of  $\varepsilon$  taking a value greater than or equal to  $\varepsilon_u$  is  $\gamma_u$ . Note that  $\varepsilon_l$  and  $\varepsilon_u$  are likely to be negative and positive quantities, respectively.

By assertion 2, the difference between  $\max(\mathcal{Y}_{d,t})$  and  $\min(\mathcal{Y}_{d,t})$  is small for large n, and, therefore,

$$\min(\mathcal{Y}_{d,t}) \approx \max(\mathcal{Y}_{d,t}) \approx y_t(\mathbf{x}). \tag{A9}$$

Using Eq. (A9), the minimum and maximum values contained in  $\mathcal{Y}_d$  may be approximated by  $y_t(x) + \varepsilon_l$  and  $y_t(x) + \varepsilon_l$ . These lower and upper bounds represent the bounds on measurement uncertainty due to errors in streamflow measurements. As sample size increases, the probabilities  $\gamma_l$  and  $\gamma_l$  would approach zero, the approximation (A9) would become more accurate, and, thus, the proposed hypothesis would become more accurate. This completed the analysis of the 1st case.

In the preceding paragraph, we argued mathematically that as the sample size increases and the neighborhood  $\mathcal{N}_d(x)$  becomes smaller, the set  $\mathcal{Y}_d$  represents measurement uncertainty in y more accurately. In reality,  $\mathcal{N}_d(x)$  cannot be arbitrarily small and the sample size is finite – thus  $\mathcal{Y}_d$  represents both measurement and structural uncertainty. However, the structural uncertainty would still be small if the sample size is large enough so as to create small leaf nodes (see Assertion 2 above and Eq. (A8)). Practically speaking, one can aim only for the modest goal of obtaining an uncertainty bound where majority of width is due to measurement uncertainty. Fortunately, this is useful in practice in the construction of LOAs as it helps avoid type-1 errors (rejecting models with good structures) at the cost of a few type-2 errors (accepting a few models with bad structures). This is a desirable property of the LOAs (Beven, 2019).

# A.2. Case 2: Only rainfall measurements are uncertain

Let  $\mathscr X$  denote the predictor space,  $x \in \mathscr X$  denote a point in the predictor space, and  $\mathscr N_d(x)$  denote the d-neighborhood of x in  $\mathscr X$  where d is a suitable distance metric. Here, x represents a vector containing rainfall and other relevant predictor variables. Let  $x_r$  denote the component of x containing error corrupted current and time-lagged rainfall values.  $x_r$  can be written as

$$x_{\rm r} = x_{\rm r,t} + \epsilon_{\rm x,r},\tag{A10}$$

where  $x_{r,t}$  is the true value and  $\varepsilon_x$  is the error in  $x_r$ . Denote by  $\mathcal{Y}$  the set containing y values as defined in Eq. (A1).

**Assumption 6.** The expected value of  $\varepsilon_x$  is zero.

**Assumption 7.** We assume that the probability distribution of  $\varepsilon_x$  varies slowly within  $\mathcal{N}_d(x)$ . The probability distribution of  $\varepsilon_x$  can be assumed independent and identically distributed within  $\mathcal{N}_d(x)$ .

For each  $x \in \mathcal{N}_d(x)$ , there exists a true value  $x_t$  and corresponding to each  $x_t$ , there exists a  $y_t$  value. Thus, we can define a set  $\mathcal{Y}_d$  similar to that defined in Eq. (A4), only difference being that the x values are error corrupted in this case.

**Assertion 3.** : The diameter of  $N_d(x)$  approaches zero as the sample size increases.

This assertion follows from the proof of assertion 2.

**Assertion 4.** : The true value of the values contained in  $\mathcal{N}_d(x)$  approximate the probability distribution of x, for large sample large.

Following assertion 3, it is reasonable to assume that values contained in  $N_d(x)$  are approximately equal, that is, any  $x_d \in N_d(x)$  is approximately equal to x. But the values contained in  $N_d(x)$  are error corrupted; therefore, the true value corresponding to any  $x_d \in N_d(x)$  can be written as

$$\mathbf{x}_{d,t} = \mathbf{x}_d - \epsilon_{\mathbf{x}} = \mathbf{x} - \epsilon_{\mathbf{x}}. \tag{A11}$$

From Eq. (A11), it is clear that  $x_{d,t}$  is a random variable with mean value x and larger moments defined by  $\varepsilon_x$ . Hence, assertion 4 follows.

**Corollary 1.** The minimum and maximum values contained in  $\mathcal{N}_d(x)$  can be approximated by  $\mathbf{x} + \varepsilon_{\mathbf{x},\mathbf{l}}$  and  $\mathbf{x} + \varepsilon_{\mathbf{x},\mathbf{u}}$  respectively. Here,  $\varepsilon_{\mathbf{x},\mathbf{l}}$  and  $\varepsilon_{\mathbf{x},\mathbf{u}}$  are defined similarly as  $\varepsilon_{\mathbf{l}}$  and  $\varepsilon_{\mathbf{u}}$  are defined in theorem 1. Again,  $\varepsilon_{\mathbf{x},\mathbf{l}}$  and  $\varepsilon_{\mathbf{x},\mathbf{u}}$  are likely to be negative and positive quantities, respectively.

**Assertion 5.** : There exists a one-to-one mapping between  $N_d(x)$  and  $\mathcal{Y}_d$ .

It can be seen from Eq. (A11) that there exists a *unique* true value corresponding to each  $x_d \in N_d(x)$ . For two values contained in  $N_d(x)$  to be identical, the value of  $\varepsilon_x$  will have to be identical; but the probability of such an event is practically zero (less than some arbitrarily small  $\delta > 0$  to be more precise).

By assumption 3, there exists a one-to-one relationship between true value of predictor and response variables; therefore, there must exist a one-one mapping between  $N_d(x)$  and  $\mathcal{Y}_d$ .

**Theorem 2.** The set  $\mathcal{Y}_d$  provides the effect of measurement uncertainty in rainfall on streamflow  $y_t(x)$ .

The truth in this assertion stems from one-to-one mapping between the elements of  $\mathcal{N}_d(x)$  and  $\mathcal{Y}_d$  (Assertion 5). And since by assertion 4,  $\mathcal{N}_d(x)$  provides measurement uncertainty in x,  $\mathcal{Y}_d$  yields the effect of measurement uncertainty in x on y(x).

The set  $N_d(x)$  contains several elements with approximately the same value x. But these values are error corrupted; the underlying true values will differ due to measurement uncertainty in x. For each unique true value in  $N_d(x)$ , there exists a unique value of y in  $\mathcal{Y}_d$ . When we observe an error corrupted value x, the corresponding response can be any value contained in  $\mathcal{Y}_d$  depending upon the error in x. Therefore, the LOA corresponding to x should be  $(\min(\mathcal{Y}_d), \max(\mathcal{Y}_d))$ .

This completes the analysis of 2nd case.

The above analysis is valid in the case of large number of samples. With finite samples,  $\mathcal{Y}_d$  would capture measurement uncertainty and structural uncertainty because the diameter of  $N_d(x)$  would not be small. But a sufficiently large number of samples would result in small structural uncertainty.

# A.3. Case 3: Both streamflow and rainfall measurements are uncertain

Here, we consider the case where both the rainfall and streamflow measurements are corrupted by errors. This case is a combination of case 1 and case 2. The notations and assumptions are same as in previous two cases. Consider  $x_d \in N_d(x)$  and the corresponding response variable  $y_d \in \mathcal{Y}_d$ . The error corrupted  $x_d$  and  $y_d$  can be represented by Eqs. (A2) and (A10), respectively.

**Theorem 3.** The set  $\mathcal{Y}_d$  provides lower and upper measurement bounds due to errors in response measurements and the effect of errors in predictor measurements, if the sample size is large.

From Theorem 2, clearly  $\mathcal{Y}_{d,t}$  would yield the effect of errors in predictor variable measurements. Here,  $\mathcal{Y}_{d,t}$  is defined as in Eq. (A5). Further, note that since response measurement is also error-corrupted, the values contained in  $\mathcal{Y}_d$  can be written as

$$y_d(\mathbf{x}_d) = y_t(\mathbf{x}_d - \epsilon_x) + \epsilon(y_t), \tag{A12}$$

where  $\mathbf{x}_d \in N_d(\mathbf{x})$  and  $\mathbf{y}_d \in \mathcal{Y}_d$  are error-corrupted values,  $\mathbf{y}_t$  and  $\mathbf{x}_d - \varepsilon_x$  are true values of predictor and response variables, respectively. The term  $\varepsilon$  represents measurement error in response variable which is a function of  $\mathbf{y}_t$ . Here,  $\varepsilon$  cannot be assumed independent of  $\mathbf{y}_t$  values since the variation of  $\mathbf{y}_t$  within  $\mathcal{Y}_{d,t}$  is large in this case as opposed to that in case 1.

Denote the set containing true value  $y_t$  corresponding to each true value in  $N_d(x)$  by  $\mathcal{Y}_{d,t}$ , as in Eq. (A5). Then, the minimum and maximum values contained in  $\mathcal{Y}_d$  are  $\min(\mathcal{Y}_{d,t}) + \epsilon_1(\min(\mathcal{Y}_{d,t}))$  and  $\max(\mathcal{Y}_{d,t}) + \epsilon_u(\max(\mathcal{Y}_{d,t}))$ . Here,  $\epsilon_1(\min(\mathcal{Y}_{d,t}))$  is the value of  $\epsilon_1(\min(\mathcal{Y}_{d,t}))$  in the left tail of the

distribution such that probability of  $\epsilon(\min(\mathscr{Y}_{d,t}))$  taking a value less than  $\epsilon_l(\min(\mathscr{Y}_{d,t}))$  is  $\gamma_l$ . The term  $\epsilon_u(\max(\mathscr{Y}_{d,t}))$  is defined similarly. For large sample, the probability  $\gamma_l$  will approach 0. The quantities  $\min(\mathscr{Y}_{d,t}) + \epsilon_l(\min(\mathscr{Y}_{d,t}))$  and  $\max(\mathscr{Y}_{d,t}) + \epsilon_u(\max(\mathscr{Y}_{d,t}))$  are lower and upper bounds of total measurement uncertainty due to errors in predictor and response variables.

This completes the proof of case 3.

# Appendix B

### References

- Ammann, L., Fenicia, F., Reichert, P., 2019. A likelihood framework for deterministic hydrological models and the importance of non-stationary autocorrelation. Hydrol. Earth Syst. Sci. 23 (4), 2147–2172.
- Archuleta, C.M., Constance, E.W., Arundel, S.T., Lowe, A.J., Mantey, K.S., and Phillips, L. A., 2017, The National Map Seamless Digital Elevation Model specifications: U.S. Geological Survey Techniques and Methods, book 11, chap. B9, 39, 10.3133/tm11
- Beven, K., 2005. On the concept of model structural error. Water Sci. Technol. 52 (6), 167–175.
- Beven, K., 2019. Towards a methodology for testing models as hypotheses in the inexact sciences. Proc. R. Soc. A 475 (2224), 20180862.
- Beven, K., 2020. Deep learning, hydrological processes and the uniqueness of place. Hydrol. Process. 34 (16), 3608–3613.
- Beven, K., Lane, S., 2019. Invalidation of models and fitness-for-purpose: a rejectionist approach. Comput. Simul. Validation: Fundam. Concepts, Methodol. Frameworks, Philos. Perspect. 145–171.
- Beven, K., Lane, S., 2022. On (in) validating environmental models. 1. Principles for formulating a Turing-like Test for determining when a model is fit-for purpose. Hydrol. Process. 36 (10), e14704.
- Beven, K., Westerberg, I., 2011. On red herrings and real herrings: disinformation and information in hydrological inference. Hydrol. Process. 25 (10), 1676–1680.
- Beven, K., Lane, S., Page, T., Kretzschmar, A., Hankin, B., Smith, P., Chappell, N., 2022. On (in) validating environmental models. 2. Implementation of a Turing-like test to modelling hydrological processes. Hydrol. Process. 36 (10), e14703.
- Beven, K., Smith, P., 2015. Concepts of information content and likelihood in parameter calibration for hydrological simulation models. J. Hydrol. Eng. 20 (1), A4014010
- Blöschl, G., Bierkens, M.F., Chambel, A., Cudennec, C., Destouni, G., Fiori, A., Renner, M., 2019. Twenty-three unsolved problems in hydrology (UPH)–a community perspective. Hydrol. Sci. J. 64 (10), 1141–1158.
- Breiman, L., 2001. Random forests. Mach. Learn. 45, 5-32.
- Breiman, L., Friedman, J.H., Olshen, R.A., Stone, C.J., 1984. Classification and Regression Trees. Routledge.
- Brynjarsdóttir, J., O'Hagan, A., 2014. Learning about physical parameters: the importance of model discrepancy. Inverse Probl. 30 (11), 114007.
- Coxon, G., Freer, J., Wagener, T., Odoni, N.A., Clark, M., 2014. Diagnostic evaluation of multiple hypotheses of hydrological behaviour in a limits-of-acceptability framework for 24 UK catchments. Hydrol. Process. 28 (25), 6135–6150.
- de Oliveira, D.Y., Vrugt, J.A., 2022. The treatment of uncertainty in hydrometric observations: a probabilistic description of streamflow records. Water Resour. Res. 58 (11), e2022WR032263.
- Denil, M., Matheson, D., De Freitas, N., 2014. Narrowing the gap: random forests in theory and in practice. In: International Conference on Machine Learning. PMLR, pp. 665–673.
- Di Baldassarre, G., Montanari, A., 2009. Uncertainty in river discharge observations: a quantitative analysis. Hydrol. Earth Syst. Sci. 13 (6), 913–921.
- Fernández, C., Steel, M.F., 1998. On Bayesian modeling of fat tails and skewness. J. Am. Stat. Assoc. 93 (441), 359–371.
- Fisher, R.A., 1956. Statistical Methods and Scientific Inference, 3rd ed. McMillan, London.
- Friedman, J., Hastie, T., Tibshirani, R., 2001. The Elements of Statistical Learning. Springer Series in Statistics, New York. Vol. 1, No. 10.
- Gabellani, S., Boni, G., Ferraris, L., Von Hardenberg, J., Provenzale, A., 2007. Propagation of uncertainty from rainfall to runoff: a case study with a stochastic rainfall generator. Adv. Water Resour. 30 (10), 2061–2071.
- Gelman, A., Rubin, D.B., 1992. Inference from iterative simulation using multiple sequences. Stat. Sci. 7 (4), 457–472.
- Gong, W., Gupta, H.V., Yang, D., Sricharan, K., Hero III, A.O., 2013. Estimating epistemic and aleatory uncertainties during hydrologic modeling: an information theoretic approach. Water Resour. Res. 49 (4), 2253–2273.
- Govindaraju, R.S., 2000. Artificial neural networks in hydrology. II: hydrologic applications. J. Hydrol. Eng. 5 (2), 124–137.
- Gupta, A., Govindaraju, R.S., 2022. Uncertainty quantification in watershed hydrology: which method to use? J. Hydrol., 128749
- Gupta, A., Govindaraju, R.S., Morbidelli, R., Corradini, C., 2022. The Role of Prior Probabilities on Parameter Estimation in Hydrological Models. Water Resour. Res., e2021WR031291
- Gupta, V.K., Waymire, E.C., 1993. A statistical analysis of mesoscale rainfall as a random cascade. J. Appl. Meteorol. Climatol. 32 (2), 251–267.
- Haario, H., Laine, M., Mira, A., Saksman, E., 2006. DRAM: efficient adaptive MCMC. Stat. Comput. 16 (4), 339–354.

- Herschy, R., 1993. The stage-discharge relation. Flow Meas. Instrum. 4 (1), 11–15. lorgulescu, I., Beven, K.J., 2004. Nonparametric direct mapping of rainfall-runoff relationships: an alternative approach to data analysis and modeling? Water Resour.
- Kavetski, D., Kuczera, G., Franks, S.W., 2006a. Bayesian analysis of input uncertainty in hydrological modeling: 1. Theory. Water Resour. Res. 42 (3).
- Kavetski, D., Kuczera, G., Franks, S.W., 2006b. Bayesian analysis of input uncertainty in hydrological modeling: 2. Application. Water Resour. Res. 42 (3).
- Kennedy, M.C., O'Hagan, A, 2001. Bayesian calibration of computer models. J. R. Stat. Soc.: Ser. B (Stat. Methodol.) 63 (3), 425–464.
- Kiang, J.E., Gazoorian, C., McMillan, H., Coxon, G., Le Coz, J., Westerberg, I.K., Mason, R., 2018. A comparison of methods for streamflow uncertainty estimation. Water Resour. Res. 54 (10), 7149–7176.
- Kratzert, F., Klotz, D., Herrnegger, M., Sampson, A.K., Hochreiter, S., Nearing, G.S., 2019. Toward improved predictions in ungauged basins: exploiting the power of machine learning. Water Resour. Res.
- Krueger, T., Freer, J., Quinton, J.N., Macleod, C.J., Bilotta, G.S., Brazier, R.E., Haygarth, P.M., 2010. Ensemble evaluation of hydrological model hypotheses. Water Resour. Res. 46 (7).
- Kuczera, G., Parent, E., 1998. Monte Carlo assessment of parameter uncertainty in conceptual catchment models: the Metropolis algorithm. J. Hydrol. 211 (1–4), 69–85
- Lamb, R., Beven, K., 1997. Using interactive recession curve analysis to specify a general catchment storage model. Hydrol. Earth Syst. Sci. 1 (1), 101–113.
- Le Coz, J., Renard, B., Bonnifait, L., Branger, F., Le Boursicaud, R, 2014. Combining hydraulic knowledge and uncertain gaugings in the estimation of hydrometric rating curves: a Bayesian approach. J. Hydrol. 509, 573–587.
- Lele, S.R., 2004. Evidence functions and the optimality of the law of likelihood. Nat. Sci. Evid.: Stat., Philos., Empir. Considerations 191–216.
- Liu, Y., Freer, J., Beven, K., Matgen, P., 2009. Towards a limits of acceptability approach to the calibration of hydrological models: extending observation error. J. Hydrol. 367 (1–2), 93–103.
- McMillan, H.K., Westerberg, I.K., Krueger, T., 2018. Hydrological data uncertainty and its implications. Wiley Interdiscip. Rev.: Water 5 (6), e1319.
- McMillan, H., Krueger, T., Freer, J., 2012. Benchmarking observational uncertainties for hydrology: rainfall, river discharge and water quality. Hydrol. Process. 26 (26), 4078–4111.
- Moulin, L., Gaume, E., Obled, C., 2009. Uncertainties on mean areal precipitation: assessment and impact on streamflow simulations. Hydrol. Earth Syst. Sci. 13 (2), 99–114.
- Neyman, J., Pearson, E.S., 1933. IX. On the problem of the most efficient tests of statistical hypotheses. Philos. Trans. R. Soc. Lond. Ser. A, Containing Pap. Math. Phys. Character 231 (694–706), 289–337.
- Nott, D.J., Marshall, L., Brown, J., 2012. Generalized likelihood uncertainty estimation (GLUE) and approximate Bayesian computation: what's the connection? Water Resour. Res. 48 (12).
- Pande, S., 2013a. Quantile hydrologic model selection and model structure deficiency assessment: 1. Theory. Water Resour. Res. 49 (9), 5631–5657.
- Pande, S., 2013b. Quantile hydrologic model selection and model structure deficiency assessment: 2. Applications. Water Resour. Res. 49 (9), 5658–5673.
- Petersen-Øverleir, A., Soot, A., Reitan, T., 2009. Bayesian rating curve inference as a streamflow data quality assessment tool. Water Resour. Manage. 23 (9), 1835–1842.
- Reichert, P., Mieleitner, J., 2009. Analyzing input and structural uncertainty of nonlinear dynamic models with stochastic, time-dependent parameters. Water Resour. Res. 45 (10).
- Reitan, T., Petersen-Øverleir, A., 2009. Bayesian methods for estimating multi-segment discharge rating curves. Stochastic Environ. Res. Risk Assess. 23 (5), 627–642.
- Renard, B., Kavetski, D., Kuczera, G., Thyer, M., Franks, S.W., 2010. Understanding predictive uncertainty in hydrologic modeling: the challenge of identifying input and structural errors. Water Resour. Res. 46 (5).
- Renard, B., Kavetski, D., Leblois, E., Thyer, M., Kuczera, G., Franks, S.W., 2011. Toward a reliable decomposition of predictive uncertainty in hydrological modeling: characterizing rainfall errors using conditional simulation. Water Resour. Res. 47 (11)
- Royall, R., 2017. Statistical Evidence: A Likelihood Paradigm. Routledge.
- Sadegh, M., Vrugt, J.A., 2013. Bridging the gap between GLUE and formal statistical approaches: approximate Bayesian computation. Hydrol. Earth Syst. Sci. 17 (12), 4831–4850.
- Schoups, G., Vrugt, J.A., 2010. A formal likelihood function for parameter and predictive inference of hydrologic models with correlated, heteroscedastic, and non-Gaussian errors. Water Resour. Res. 46 (10).

- Shortridge, J.E., Guikema, S.D., Zaitchik, B.F., 2016. Machine learning methods for empirical streamflow simulation: a comparison of model accuracy, interpretability, and uncertainty in seasonal watersheds. Hydrol. Earth Syst. Sci. 20 (7).
- Smith, T., Marshall, L., Sharma, A., 2015. Modeling residual hydrologic errors with Bayesian inference. J. Hydrol. 528, 29–37.
- Sturm, T.W., 2001. Open Channel Hydraulics. McGraw-Hill, New York. Tallaksen, L.M., 1995. A review of baseflow recession analysis. J. Hydrol. 165 (1–4), 349–370.
- Vrugt, J.A., Beven, K.J., 2018. Embracing equifinality with efficiency: limits of Acceptability sampling using the DREAM (LOA) algorithm. J. Hydrol. 559, 954–971.
- Vrugt, J.A., Sadegh, M., 2013. Toward diagnostic model calibration and evaluation: approximate Bayesian computation. Water Resour. Res. 49 (7), 4335–4345.
- Wagener, T., Sivapalan, M., Troch, P., Woods, R., 2007. Catchment classification and hydrologic similarity. Geogr. Compass 1 (4), 901–931.
- Waymire, E.D., Gupta, V.K., 1981a. The mathematical structure of rainfall representations: 1. A review of the stochastic rainfall models. Water Resour. Res. 17 (5), 1261–1272.

- Waymire, E.D., Gupta, V.K., 1981c. The mathematical structure of rainfall representations: 3. Some applications of the point process theory to rainfall processes. Water Resour. Res. 17 (5), 1287–1294.
- Waymire, E.D., Gupta, V.K., 1981b. The mathematical structure of rainfall representations: 2. A review of the theory of point processes. Water Resour. Res. 17 (5), 1273–1285.
- Westerberg, I.K., Guerrero, J.L., Younger, P.M., Beven, K.J., Seibert, J., Halldin, S., Xu, C. Y., 2011. Calibration of hydrological models using flow-duration curves. Hydrol. Earth Syst. Sci. 15 (7), 2205–2227.
- Zhang, B., Govindaraju, R.S., 2000. Prediction of watershed runoff using Bayesian concepts and modular neural networks. Water Resour. Res. 36 (3), 753–762.
- Zhang, B., Govindaraju, R.S., 2003. Geomorphology-based artificial neural networks (GANNs) for estimation of direct runoff over watersheds. J. Hydrol. 273 (1–4), 18–34