



Non-crossing nonlinear regression quantiles by monotone composite quantile regression neural network, with application to rainfall extremes

Alex J. Cannon¹

Published online: 26 June 2018
© The Author(s) 2018

Abstract

The goal of quantile regression is to estimate conditional quantiles for specified values of quantile probability using linear or nonlinear regression equations. These estimates are prone to “quantile crossing”, where regression predictions for different quantile probabilities do not increase as probability increases. In the context of the environmental sciences, this could, for example, lead to estimates of the magnitude of a 10-year return period rainstorm that exceed the 20-year storm, or similar nonphysical results. This problem, as well as the potential for overfitting, is exacerbated for small to moderate sample sizes and for nonlinear quantile regression models. As a remedy, this study introduces a novel nonlinear quantile regression model, the monotone composite quantile regression neural network (MCQRNN), that (1) simultaneously estimates multiple non-crossing, nonlinear conditional quantile functions; (2) allows for optional monotonicity, positivity/non-negativity, and generalized additive model constraints; and (3) can be adapted to estimate standard least-squares regression and non-crossing expectile regression functions. First, the MCQRNN model is evaluated on synthetic data from multiple functions and error distributions using Monte Carlo simulations. MCQRNN outperforms the benchmark models, especially for non-normal error distributions. Next, the MCQRNN model is applied to real-world climate data by estimating rainfall Intensity–Duration–Frequency (IDF) curves at locations in Canada. IDF curves summarize the relationship between the intensity and occurrence frequency of extreme rainfall over storm durations ranging from minutes to a day. Because annual maximum rainfall intensity is a non-negative quantity that should increase monotonically as the occurrence frequency and storm duration decrease, monotonicity and non-negativity constraints are key constraints in IDF curve estimation. In comparison to standard QRNN models, the ability of the MCQRNN model to incorporate these constraints, in addition to non-crossing, leads to more robust and realistic estimates of extreme rainfall.

1 Introduction

Estimating regression quantiles—conditional quantiles of a response variable that depend on covariates in some form of regression equation—is a fundamental task in data-driven science. Focusing on the environmental sciences, quantile regression methods have been used to provide estimates of predictive uncertainty in forecast applications (Cawley et al. 2007); construct growth curves for organisms (Muggeo et al. 2013); relate soil moisture deficit with summer hot extremes (Hirschi et al. 2010); provide flood

frequency estimates (Ouali et al. 2016); estimate rainfall Intensity–Duration–Frequency (IDF) curves (Ouali and Cannon 2018); determine the relation between rainfall intensity and duration and landslide occurrence (Saito et al. 2010); estimate trends in climate, streamflow, and sea level data (Koenker and Schorfheide 1994; Barbosa 2008; Allamano et al. 2009; Roth et al. 2015); downscale atmospheric model outputs (Friederichs and Hense 2007; Cannon 2011; Ben Alaya et al. 2016); and determine scaling relationships between temperature and extreme precipitation (Wasko and Sharma 2014), among other applications.

Quantile regression equations can be linear or nonlinear. In most variants, including the original linear model (Koenker and Bassett 1978), conditional quantiles for specified quantile probabilities are estimated separately by different regression equations; together, these different equations can be used to build up a piecewise estimate of

✉ Alex J. Cannon
alex.cannon@canada.ca

¹ Climate Research Division, Environment and Climate Change Canada, Victoria, BC, Canada

the conditional response distribution. However, given finite samples, this flexibility can lead to “quantile crossing” where, for some values of the covariates, quantile regression predictions do not increase with the specified quantile probability τ . For instance, the $\tau_1 = 0.1$ -quantile (10th-percentile) estimate may be greater in magnitude than the $\tau_2 = 0.2$ -quantile (20th-percentile) estimate, which violates the property that the conditional quantile function be strictly monotonic. As Ouali et al. (2016) state, “crossing quantile regression is a serious modeling problem that may lead to an invalid response distribution”.

Three main approaches have been used to solve the quantile crossing problem: post-processing, stepwise estimation, and simultaneous estimation. In post-processing, non-crossing quantiles are enforced following model estimation by rearranging predictions so that they increase with increasing τ (Chernozhukov et al. 2010). In stepwise estimation, regression equations are constructed iteratively, with constraints added so that each subsequent quantile regression function does not cross the one estimated previously (Liu and Wu 2009; Muggeo et al. 2013). Finally, in simultaneous estimation, quantile regression equations for all desired values of τ are estimated at the same time, with additional constraints added to parameter optimization to ensure non-crossing (Takeuchi et al. 2006; Bondell et al. 2010; Liu and Wu 2011; Bang et al. 2016). Unlike sequential estimation, simultaneous estimation is attractive because it does not depend on the order in which quantiles are estimated. Furthermore, fitting for multiple values of τ simultaneously allows one to “borrow strength” across regression quantiles and improve overall model performance (Bang et al. 2016). This property is especially useful for nonlinear quantile regression models, which are more prone to overfitting and quantile crossing in the face of small to moderate sample sizes (Muggeo et al. 2013).

Baldwin (2006), paraphrasing Persson (2001), states “...while there is only one way to be linear, there are an uncountable infinity of ways to be nonlinear. One cannot check them all”. For a flexible nonlinear model like a neural network, imposing extra constraints, for example as informed by process knowledge, can be useful for narrowing the overall search space of potential nonlinearities. As a simple example, growth curves should increase monotonically with the age of the organism, which led Muggeo et al. (2013) to introduce a monotonicity constraint in addition to the non-crossing constraint. Similarly, Roth et al. (2015) applied nonlinear monotone quantile regression to describe non-decreasing trends in rainfall extremes. Takeuchi et al. (2006) developed a nonparametric, kernelized version of quantile regression with similarities to support vector machines; both non-crossing and monotonicity constraints are considered, with directions on the incorporation of other constraints, such as positivity and additivity constraints, also

provided. However, standard implementations of the kernel quantile regression model (e.g., Karatzoglou et al. 2004; Hofmeister 2017) are computationally costly, with complexity that is cubic in the number of samples, and do not explicitly implement the proposed constraints.

As an alternative, this study introduces an efficient, flexible nonlinear quantile regression model, the monotone composite quantile regression neural network (MCQRNN), that: (1) simultaneously estimates multiple non-crossing quantile functions; (2) allows for optional monotonicity, positivity/non-negativity, and additivity constraints, as well as fine-grained control on the degree of non-additivity; and (3) can be modified to estimate standard least-squares regression and non-crossing expectile regression functions. These features, which are combined into a single, unified framework, are made possible through a novel combination of elements drawn from the standard QRNN model (White 1992; Taylor 2000; Cannon 2011), the monotone multi-layer perceptron (MMLP) (Zhang and Zhang 1999; Lang 2005; Minin et al. 2010), the composite QRNN (CQRNN) (Xu et al. 2017), the expectile regression neural network (Jiang et al. 2017), and the generalized additive neural network (Potts 1999). To the best of the author’s knowledge, the MCQRNN model is the first neural network-based implementation of quantile regression that guarantees non-crossing of regression quantiles.

The MCQRNN model is developed in Sect. 2, starting from the MMLP model, leading to the MQRNN model, and then finally to the full MCQRNN. Approaches to enforce monotonicity, positivity/non-negativity, and generalized additive model constraints, as well as to estimate uncertainty in the conditional τ -quantile functions, are also provided. In Sect. 3, the MCQRNN model is compared via Monte Carlo simulation to standard MLP, QRNN, and CQRNN models using combinations of three functions and error distributions from Xu et al. (2017). In Sect. 4, the MCQRNN model is applied to real-world climate data by estimating IDF curves at ungauged locations in Canada based on annual maximum rainfall series at neighbouring gauging stations. IDF curves, which are used in the design of civil infrastructure such as culverts, storm sewers, dams, and bridges, summarize the relationship between the intensity and occurrence frequency of extreme rainfall over averaging durations ranging from minutes to a day (Canadian Standards Association 2012). The intensity of extreme rainfall, a non-negative quantity, should increase monotonically as the annual probability of occurrence decreases (e.g., from $1 - \tau = 0.5$ to 0.01 or, equivalently, a 2–100-year return period) and as the storm duration decreases (e.g., from 24-h to 5-min). Monotonicity and positivity/non-negativity constraints are thus key features of an IDF curve. MCQRNN IDF curve estimates are compared with those obtained by fitting separate QRNN models for each return period and duration, as done previously by Ouali and Cannon (2018). Finally, Sect. 5 provides closing remarks and suggestions for future research.

2 Modelling framework

2.1 Monotone multi-layer perceptron (MMLP)

The monotone composite quantile regression neural network (MCQRNN) model starts with the multi-layer perceptron (MLP) neural network with partial monotonicity constraints (Zhang and Zhang 1999) as its basis. For a data point with index t , the prediction $\hat{y}(t)$ from a monotone MLP (MMLP) is obtained as follows. First, the V covariates, each assumed to be standardized to zero mean and unit standard deviation, are separated into two groups: $x_{m \in M}(t)$ and $x_{i \in I}(t)$ with combined indices $\{M \cup I | 1, \dots, V, V = (\#M + \#I)\}$, where M is the set of indices for covariates with a monotone increasing relationship with the prediction, I is the corresponding set of indices for covariates without monotonicity constraints, and $\#$ denotes the number of set elements. Covariates are transformed into $j = 1, \dots, J$ hidden layer outputs

$$h_j(t) = f \left(\sum_{m \in M} x_m(t) \exp(W_{mj}^{(h)}) + \sum_{i \in I} x_i(t) W_{ij}^{(h)} + b_j^{(h)} \right) \tag{1}$$

where $\mathbf{W}^{(h)}$ is a $V \times J$ parameter matrix, $\mathbf{b}^{(h)}$ is a vector of J intercept parameters, and f is a smooth non-decreasing function, usually taken to be the hyperbolic tangent function. Finally, the model prediction is given as a weighted combination of the J hidden layer outputs

$$\hat{y}(t) = g \left(\sum_{j=1}^J h_j(t) \exp(w_j) + b \right) \tag{2}$$

where \mathbf{w} is a vector of J parameters, b is an intercept term, and g is a smooth non-decreasing inverse-link function.

Because both f and g are non-decreasing, partial monotonicity constraints (i.e., $\frac{\partial \hat{y}}{\partial x_m} \geq 0$ everywhere) can be imposed by ensuring that all parameters leading from each monotone-constrained covariate x_m are positive (Zhang and Zhang 1999), in this case by applying the exponential function to the corresponding elements of $\mathbf{W}^{(h)}$ and all elements of \mathbf{w} . Decreasing relationships can be imposed by multiplying covariates by -1 . Also, extra hidden layers of positive parameters can be added to the model. As pointed out by Lang (2005) and Minin et al. (2010), an additional hidden layer is required for the MMLP to maintain its universal function approximation capabilities. While multiple hidden layers are included in the software implementation by Cannon (2017), for sake of simplicity, this study only considers the single hidden layer architecture of Zhang and Zhang (1999). In practice, simple functional relationships can still be represented by a single hidden layer model.

If M is the empty set and the positivity constraint on the \mathbf{w} parameters is removed, this leads to the standard MLP model. If f and g are the identity function, the MMLP reduces to a linear model. If f is nonlinear, then the model can represent nonlinear relationships, including those involving interactions between covariates; the number of hidden layer outputs J further controls the potential complexity of the MLP mapping. All models in this study set f to be the hyperbolic tangent function.

Adjustable parameters ($\mathbf{W}^{(h)}$, $\mathbf{b}^{(h)}$, \mathbf{w} , b) in the MMLP are set by minimizing the least squares (LS) error function

$$E_{LS} = \frac{1}{N} \sum_{t=1}^N (y(t) - \hat{y}(t))^2 \tag{3}$$

over a training dataset with N data points $\{(\mathbf{x}(t), y(t)) | t = 1, \dots, N\}$, where $y(t)$ is the target value of the response variable. While LS regression is most common, different error functions are appropriate for different prediction tasks. Minimizing the LS error function is equivalent to maximum likelihood estimation for the conditional mean assuming a Gaussian error distribution with constant variance (i.e., a traditional regression task), while minimizing the least absolute error (LAE) function

$$E_{LAE} = \frac{1}{N} \sum_{t=1}^N |y(t) - \hat{y}(t)| \tag{4}$$

leads to a regression estimate for the conditional median (i.e., the $\tau = 0.5$ -quantile) (Koenker and Bassett 1978).

2.2 Monotone quantile regression neural network (MQRNN)

The fundamental quantity of interest here is not just the median, but rather predictions $\hat{y}_\tau(t)$ of the conditional quantile associated with the quantile probability τ ($0 < \tau < 1$). In this context, combining the MMLP architecture from Sect. 2.1, as given by Eqs. (1) and (2),

$$\hat{y}_\tau(t) = g \left[\sum_{j=1}^J f \left(\sum_{m \in M} x_m(t) \exp(W_{mj}^{(h)}) + \sum_{i \in I} x_i(t) W_{ij}^{(h)} + b_j^{(h)} \right) \times \exp(w_j) + b \right], \tag{5}$$

with the quantile regression error function

$$E_\tau = \frac{1}{N} \sum_{t=1}^N \rho_\tau(y(t) - \hat{y}_\tau(t)) \tag{6}$$

where

$$\rho_\tau(\varepsilon) = \begin{cases} \tau \varepsilon & \varepsilon \geq 0 \\ (\tau - 1) \varepsilon & \varepsilon < 0 \end{cases} \tag{7}$$

leads to estimates \hat{y}_τ of the conditional τ -quantile function (Koenker and Bassett 1978). The resulting model is referred to as the MQRNN. When $\tau = 0.5$, Eq. (6) is, up to a constant scaling factor, the same as the LAE function (Eq. 4) that yields the conditional median; for $\tau \neq 0.5$, the asymmetric absolute value function gives different weight to positive/negative deviations. For example, fitting a model with $\tau = 0.95$ provides an estimate for the conditional 95th-percentile, i.e., a covariate-dependent probability of exceedance of 5%. Relaxing the monotonicity constraints gives the standard QRNN model as presented by Cannon (2011).

Parameters can be estimated by a gradient-based nonlinear optimization algorithm, with calculation of the gradient using backpropagation; given the simple relationship between Eqs. 4 and 6, the analytical expression for the gradient of the quantile regression error function follows from that of the LAE function (Hanson and Burr 1988). In this case, the derivative is undefined at the origin, which means that a smooth approximation is instead substituted for the exact quantile regression error function. Following Chen (2007) and Cannon (2011), a Huber-norm version of Eq. 7 replaces $\rho_\tau(\varepsilon)$ in the quantile regression error function. This approximation, denoted by (A), is given by

$$\rho_\tau^{(A)}(\varepsilon) = \begin{cases} \tau \varphi(\varepsilon) & \varepsilon \geq 0 \\ (\tau - 1) \varphi(\varepsilon) & \varepsilon < 0 \end{cases} \quad (8)$$

where the Huber function

$$\varphi(\varepsilon) = \begin{cases} \frac{\varepsilon^2}{2\alpha} & 0 \leq |\varepsilon| \leq \alpha \\ |\varepsilon| - \frac{\alpha}{2} & |\varepsilon| > \alpha \end{cases} \quad (9)$$

is a hybrid of the absolute value and squared error functions (Huber 1964).

The Huber function transitions smoothly from the squared error, which is applied around the origin ($\pm\alpha$) to ensure differentiability, and the absolute error. As $\alpha \rightarrow 0$, the approximate error function converges to the exact quantile regression error function. It should be noted that a slightly different approximation is used by Muggeo et al. (2012). Based on experimental results (not shown), both approximations ultimately provide models that are indistinguishable. However, the Huber function approximation is used here for its added ability to emulate the LS cost function. For sufficiently large α , all model deviations are squared and the approximate error function instead becomes an asymmetric version of the LS error function (Eq. 3). For $\tau = 0.5$ and large α , the error function is symmetric and is, up to a constant scaling factor, equal to the LS error function. For $\tau \neq 0.5$, the asymmetric LS error function results in an estimate of the conditional expectile function (Newey and

Powell 1987; Yao and Tong 1996; Waltrup et al. 2015). Hence, depending on values of α and τ , minimizing the approximate quantile regression error function can provide regression estimates for the conditional mean ($\alpha \gg 0$, $\tau = 0.5$), median ($\alpha \rightarrow 0$, $\tau = 0.5$), quantiles ($\alpha \rightarrow 0$, $0 < \tau < 1$), and expectiles ($\alpha \gg 0$, $0 < \tau < 1$) (Jiang et al. 2017). Unless noted otherwise, all subsequent references to $\rho_\tau^{(A)}$ and $E_\tau^{(A)}$ will refer to the conditional quantile form of the Huber function approximation.

Unlike linear regression, where the total number of model parameters is limited by the number of covariates V , the complexity of the MQRNN model also depends on the number of hidden layer outputs J . Model complexity, and hence J , should be set such that the model can generalize to new data, which, in practice, usually means avoiding overfitting to noise in the training dataset. Additionally, regularization terms that penalize the magnitude of the parameters, hence limiting the nonlinear modelling capability of the model, can be added to the error function

$$\tilde{E}_\tau^{(A)} = E_\tau^{(A)} + \lambda^{(h)} \frac{1}{VJ} \sum_{i=1}^V \sum_{j=1}^J (W_{ij}^{(h)})^2 + \lambda \frac{1}{J} \sum_{j=1}^J (w_j)^2 \quad (10)$$

where $\lambda^{(h)} \geq 0$ and $\lambda \geq 0$ are hyperparameters that control the size of the penalty applied to the elements of $\mathbf{W}^{(h)}$ and \mathbf{w} respectively. Values of J and, optionally, the $\lambda^{(h)}$ and λ hyperparameters are typically set by minimizing out-of-sample generalization error, for example as estimated via cross-validation or modified versions of an information criterion like the Akaike information criterion (QAIC) (Koenker and Schorfheide 1994; Doksum and Koo 2000)

$$\text{QAIC} = -2 \log(E_\tau) + 2p \quad (11)$$

where p is an estimate of the effective number of model parameters.

2.3 Monotone composite quantile regression neural network (MCQRNN)

The MQRNN model in Sect. 2.2 is specified for a single τ -quantile; no efforts are made to avoid quantile crossing for multiple estimates. To date, the simultaneous estimation of multiple τ -quantiles with guaranteed non-crossing has not been possible for QRNN models. However, simultaneous estimates for multiple values of τ are used in the composite QRNN (CQRNN) model proposed by Xu et al. (2017). CQRNN shares the same goal as the linear composite quantile regression (CQR) model (Zou and Yuan 2008), namely to borrow strength across multiple regression quantiles to improve the estimate of the true, unknown relationship between the

covariates and the response. This is especially valuable in situations where the error follows a heavy-tailed distribution. In CQR, the regression coefficients are shared across the different quantile regression models. Similarly, in CQRNN, the $\mathbf{W}^{(h)}$, $\mathbf{b}^{(h)}$, \mathbf{w} , b parameters are shared across the different QRNN models. Hence, the models are not explicitly trying to describe the full conditional response distribution, but rather a single τ -independent function that best describes the true covariate-response relationship. Structurally, the CQRNN model is the same as the QRNN model. The only difference is the quantile regression error function, which is now summed over K (usually equally spaced) values of τ

$$E_{C\tau}^{(A)} = \frac{1}{KN} \sum_{k=1}^K \sum_{t=1}^N \rho_{\tau_k}^{(A)}(y(t) - \hat{y}_{\tau_k}(t)) \tag{12}$$

where, for example, $\tau_k = \frac{k}{K+1}$ for $k = 1, 2, \dots, K$. Penalty terms can be added as in Eq. 10.

The MCQRNN model combines the MQRNN model architecture given by Eq. 5 with the composite quantile regression error function (Eq. 12) to simultaneously estimate non-crossing regression quantiles. To show how this is achieved, consider an $N \times \#I$ matrix of covariates \mathbf{X} , a corresponding response vector \mathbf{y} of length N , and the goal of estimating non-crossing quantile functions for $\tau_1 < \tau_2 < \dots < \tau_K$. First, create a new $\#M = 1$ monotone covariate vector $\mathbf{x}_m^{(S)}$ of length $S = KN$, where (S) denotes stacked data, by repeating each of the K specified τ values N times and stacking. Next, stack K copies of \mathbf{X} and concatenate with $\mathbf{x}_m^{(S)}$ to form a stacked covariate matrix $\mathbf{X}^{(S)}$ of dimension $S \times (1 + \#I)$. Finally stack K copies of \mathbf{y} to form $\mathbf{y}^{(S)}$. Taken together, this gives the stacked dataset

$$\mathbf{X}^{(S)} = \begin{bmatrix} \tau_1 & x_1(1) & \cdots & x_{\#I}(1) \\ \vdots & \vdots & \ddots & \vdots \\ \tau_1 & x_1(N) & \cdots & x_{\#I}(N) \\ \tau_2 & x_1(1) & \cdots & x_{\#I}(1) \\ \vdots & \vdots & \ddots & \vdots \\ \tau_2 & x_1(N) & \cdots & x_{\#I}(N) \\ \vdots & \vdots & \ddots & \vdots \\ \tau_K & x_1(1) & \cdots & x_{\#I}(1) \\ \vdots & \vdots & \ddots & \vdots \\ \tau_K & x_1(N) & \cdots & x_{\#I}(N) \end{bmatrix}, \mathbf{y}^{(S)} = \begin{bmatrix} y(1) \\ \vdots \\ y(N) \\ y(1) \\ \vdots \\ y(N) \\ \vdots \\ y(1) \\ \vdots \\ y(N) \end{bmatrix} \tag{13}$$

which is used to fit the MQRNN model. By treating the τ values as a monotone covariate, predictions $\hat{y}_\tau^{(S)}$ from Eq. 5 for fixed values of the non-monotone covariates are

guaranteed to increase with τ . Non-crossing is imposed by construction. Defining $\tau(s) = x_1^{(S)}(s)$, the composite quantile regression error function for the stacked data can be written as

$$E_{C\tau}^{(A,S)} = \sum_{s=1}^S \omega_{\tau(s)} \rho_{\tau(s)}^{(A)}(y^{(S)}(s) - \hat{y}_{\tau(s)}^{(S)}(s)) \tag{14}$$

where $\omega_{\tau(s)}$ are weights that can be used to allow regression quantiles for each τ_k to contribute different amounts to the total error (Jiang et al. 2012; Sun et al. 2013); constant weights $\omega_{\tau(s)} = 1/S$ lead to the standard composite quantile regression error function. Minimization of Eq. 14 results in the fitted MCQRNN model. (Note: non-crossing expectile regression models can be obtained by adjusting $\alpha \gg 0$ in $\rho_\tau^{(A)}$.) Following model estimation, conditional τ -quantile functions can be predicted for any value of $\tau_1 \leq \tau \leq \tau_K$ by entering the desired value of τ into the monotone covariate.

To illustrate, Fig. 1 shows results from a MCQRNN model ($J = 4$, $\lambda^{(h)} = 0.00001$, $\lambda = 0$, $K = 9$, $\tau = 0.1, 0.2, \dots, 0.9$) fit to 500 samples of synthetic data for the two functions from Bondell et al. (2010)

$$y_1 = 0.5 + 2x + \sin(2\pi x - 0.5) + \varepsilon \tag{15}$$

and

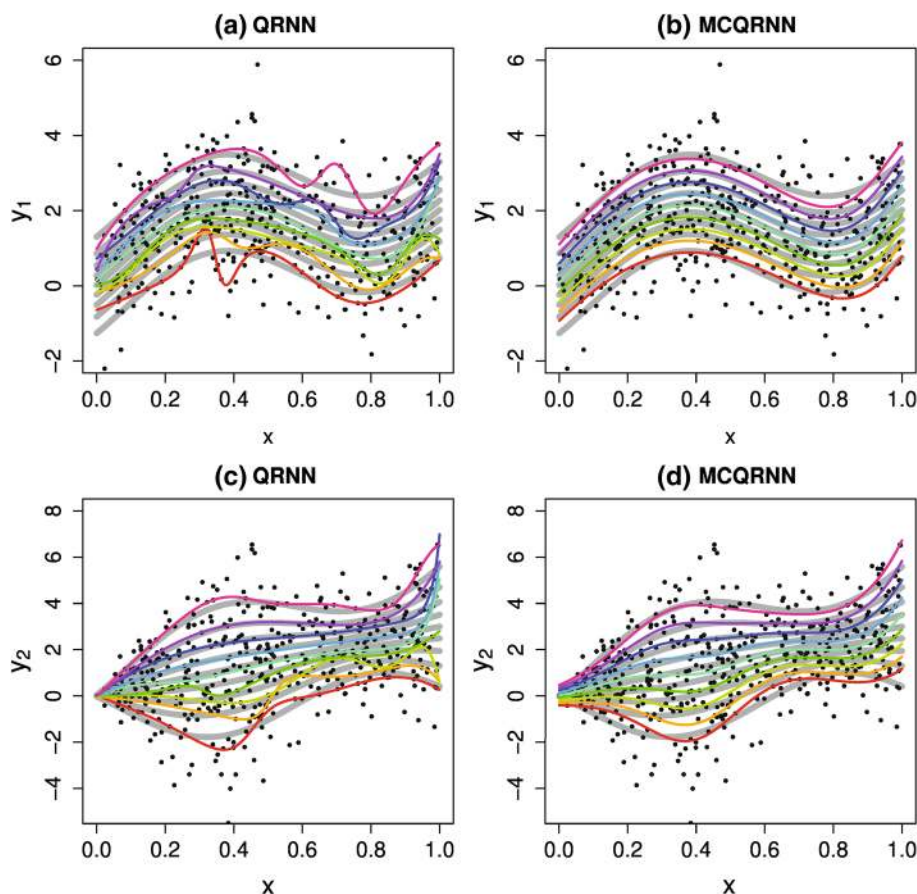
$$y_2 = 3x + [0.5 + 2x + \sin(2\pi x - 0.5)]\varepsilon \tag{16}$$

where x is drawn from the standard uniform distribution $x \sim U(0, 1)$ and ε from the standard normal distribution $\varepsilon \sim N(0, 1)$. All τ are weighted equally in Eq. 14 (i.e., values of $\omega_{\tau(s)}$ are constant). Results are compared with those from separate QRNN models ($J = 4$ and $\lambda^{(h)} = 0.00001$) for each τ -quantile. Quantile curves cross for QRNN, especially at the boundaries of the training data, whereas the MCQRNN model is able to simultaneously estimate multiple non-crossing quantile functions that correspond more closely to the true conditional quantile functions. While quantile crossing in QRNN models can be minimized by selecting and applying a suitable weight penalty (Cannon 2011), non-crossing cannot be guaranteed, whereas MCQRNN models impose this constraint by construction.

2.4 Additional constraints and uncertainty estimates

As mentioned above, constraints in addition to non-crossing of quantile functions may be useful for some MCQRNN modelling tasks. Partial monotonicity constraints for specified covariates can be imposed as described in Sect. 2.1; positivity or non-negativity constraints

Fig. 1 Predictions from QRNN (a, c) and MCQRNN (b, d) models fit to synthetic data (black points) generated by Eq. 15 (a, b) and Eq. 16 (c, d) are shown in rainbow colours. Plots of the true conditional quantile functions are shown by solid grey lines. The nine curves from bottom to top represent $\tau = 0.1, 0.2, \dots, 0.9$



can be added by setting g in Eq. 2 to the exponential or smooth ramp function (Cannon 2011), respectively; and covariate interactions can be restricted by the approach described in Appendix 1.

A form of the parametric bootstrap can be used to estimate uncertainty in the conditional τ -quantile functions. While the MCQRNN model is explicitly optimized for K specified values of τ , the use of the quantile probability as a monotone covariate means that conditional τ -quantile functions can be interpolated for any value of $\tau_1 \leq \tau \leq \tau_K$. Proper distribution, probability density, and quantile functions can then be constructed by assuming a parametric form for the tails of the distribution (Quiñonero Candela et al. 2006; Cannon 2011). The parametric bootstrap proceeds by drawing random samples from the resulting conditional distribution, refitting the MCQRNN model, making estimates of the conditional τ -quantiles, and repeating many times. Confidence intervals are estimated from the bootstrapped conditional τ -quantiles.

For illustration, examples of MCQRNN model outputs with positivity and monotonicity constraints, as well as confidence intervals obtained by the parametric bootstrap, are shown in Fig. 2 for the two Bondell et al. (2010) functions.

3 Monte Carlo simulation

Given the close relationship between the MCQRNN and CQRNN models, performance is first assessed via Monte Carlo simulation using the experimental setup adopted by Xu et al. (2017) for CQRNN. The MCQRNN model is compared with standard MLP, QRNN, and CQRNN models on datasets generated for three example functions:

$$\text{(example1)} \quad y = \sin(2x_1) + 2 \exp(-16x_2^2) + 0.5\varepsilon \quad (17)$$

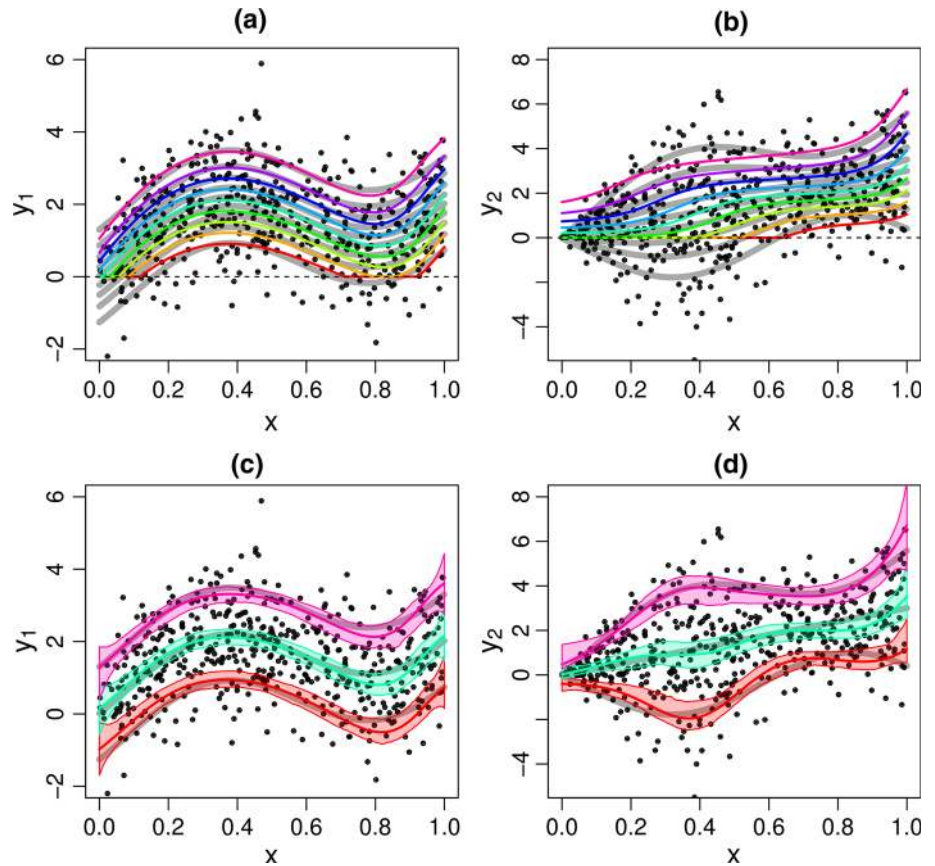
where $x_1 \sim N(0, 1)$ and $x_2 \sim N(0, 1)$;

$$\text{(example2)} \quad y = (1 - x + 2x^2) \exp(-0.5x^2) + \frac{(1 + 0.2x)}{5} \varepsilon \quad (18)$$

where $x \sim U(-4, 4)$; and

$$\begin{aligned} \text{(example3)} \quad y = & \frac{40 \exp\{8[(x_1 - 0.5)^2 + (x_2 - 0.5)^2]\}}{\exp\{8[(x_1 - 0.2)^2 + (x_2 - 0.7)^2]\}} \\ & + \exp\{8[(x_1 - 0.7)^2 + (x_2 - 0.7)^2]\}} + \varepsilon \quad (19) \end{aligned}$$

Fig. 2 As in Fig. 1b, d, but for MCQRNN models with additional **a** positivity constraints and **b** positivity and monotonicity constraints, respectively. **c, d** Estimates of 95% confidence intervals, based on 500 parametric bootstrap datasets, for the $\tau = 0.1, 0.5, 0.9$ -quantile regression curves shown in Fig. 1b, d



where $x_1 \sim U(0, 1)$ and $x_2 \sim U(0, 1)$. For each of the three functions, random errors are generated from three different distributions: the normal distribution $\varepsilon \sim N(0, 0.25)$, Student's t distribution with three degrees of freedom $\varepsilon \sim t(3)$, and the chi-squared distribution with three degrees of freedom $\varepsilon \sim \chi^2(3)$. Monte Carlo simulations are performed for the nine resulting datasets.

To evaluate the benefit of adding MCQRNN's non-crossing constraint to the simultaneous estimation of multiple regression quantiles, a second variant of CQRNN, referred to as CQRNN*, is included in the comparison. The CQRNN* model takes the same structure as MCQRNN, i.e., with τ values included as an extra input variable (Eq. 13). However, partial monotonicity constraints are removed from the τ -covariate; the exponential function is no longer applied to the relevant elements in $\mathbf{W}^{(h)}$ and all elements of \mathbf{w} . The resulting model provides estimates of multiple regression quantiles, but crossing can now occur. This differs from the CQRNN model of Xu et al. (2017), which estimates a single regression equation using the composite QR cost function, and MCQRNN, which additionally guarantees non-crossing of the multiple regression quantiles. Differences between the three models are

illustrated in Fig. 3 on the example 2 dataset with $\varepsilon \sim \chi^2(3)$ distributed noise.

For each example and error distribution in the Monte Carlo simulations, 400 samples are generated and split randomly into 200 training and 200 testing samples. Results for QRNN, MLP, CQRNN, CQRNN*, and MCQRNN models are compared by fitting to the training samples and evaluating on the testing samples. Simulations are repeated 1000 times. Following Xu et al. (2017), the number of hidden layer outputs in all models is set to $J = 4$ for example 1 and $J = 5$ for examples 2 and 3; for sake of simplicity, no weight penalty terms are added when fitting any of the models. (When comparing results with those reported by Xu et al. (2017), note that omitting weight penalty regularization here leads to smaller inter-model differences in performance within both the training and testing samples, which suggests potential instability in hyperparameter selection in the previous study.) The goal is to estimate the true functional relationship specified by Eqs. 17–19. The QRNN model is fit for $\tau = 0.5$, whereas CQRNN, CQRNN*, and MCQRNN models use $K = 19$ equally spaced values of τ . In the case of CQRNN* and MCQRNN, evaluations are based on an estimate of the conditional mean function obtained by taking the mean

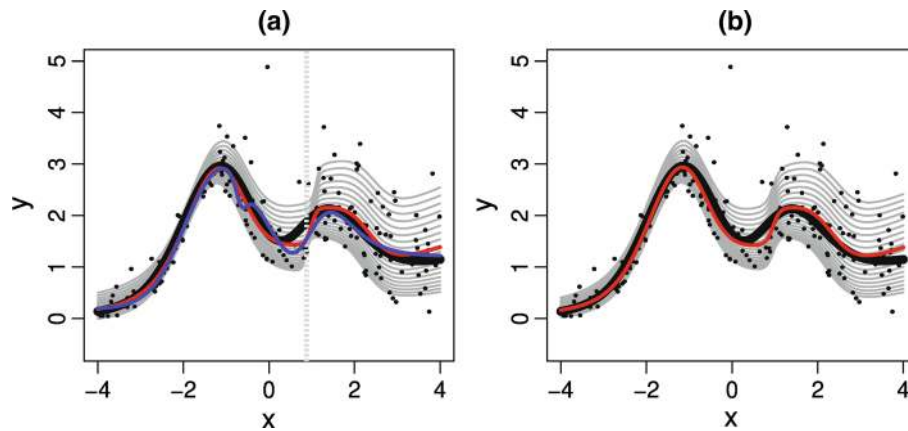


Fig. 3 Predictions from **a** CQRNN, CQRNN*, and **b** MCQRNN models on the example 2 dataset (Eq. 18) with $\varepsilon \sim \chi^2(3)$ distributed noise. Black dots show the synthetic training data and the thick black line indicates the true underlying function. Predictions of the conditional mean by CQRNN, CQRNN*, and MCQRNN are shown

by the blue line in (a), the red line in (a), and the red line in (b), respectively. For the CQRNN* and MCQRNN models, these values are obtained by taking the mean over predictions of the $K = 19\tau$ -quantiles shown in grey. Places where CQRNN* quantiles cross are indicated by vertical grey dashed lines

over predictions for the $K = 19$ τ -quantiles. Performance is measured by the root mean squared error (RMSE) between model predictions for the test samples and the actual values of y . For reference, training RMSE is also reported. Results are shown in Fig. 4.

As expected, the MLP model, which is fit using the LS error function and hence is optimal for normally distributed errors with constant variance, tends to perform best for the three examples when $\varepsilon \sim N(0, 0.25)$. Differences are, however, small for both training and testing datasets. Median RMSE values for each of the models fall within 10% of MLP in all cases and the 90% inter-percentile ranges are typically comparable. For the two non-normal error distributions, $\varepsilon \sim t(3)$ and $\varepsilon \sim \chi^2(3)$, CQRNN* and MCQRNN models tend to outperform the other models on the testing datasets. Again, differences in median testing RMSE are small, especially among the QRNN-based models. In general, however, MLP performs worst, followed by QRNN and CQRNN, with CQRNN* and MCQRNN offering slight improvements. In terms of robustness, as measured by the 5th and 95th percentiles of testing RMSE, MLP is clearly least robust, while MCQRNN tends to perform best, especially for example 3. For this example and the two non-normal error distributions, MCQRNN also outperforms CQRNN*, which points to added value of the non-crossing constraint. Overall, the MCQRNN model performs well on the synthetic data from Xu et al. (2017). In the next section, the modelling framework is applied to real-world climate data. As a proof of concept, rainfall IDF curves are estimated by MCQRNN at ungauged locations in Canada and, following Ouali and Cannon (2018), results are compared against those obtained from QRNN models.

4 Rainfall IDF curves

4.1 Data

The design of some civil infrastructure—hydraulic, hydrological, and water resource structures—is based on the design flood, which is the flood hydrograph associated with a specified frequency of occurrence or return period. In the absence of gauged discharge data, rainfall data are instead used to generate a design storm, which can then be transformed into synthetic peak streamflows for the return period of interest. The design storm provides the temporal distribution of rainfall intensities associated with a specified return period and duration. The necessary information on the frequency of occurrence, duration, and intensity of rainstorms is compactly summarized in an IDF curve, and hence IDF curves are key sources of information for engineering design applications. IDF curves provided by Environment and Climate Change Canada (ECCC) summarize the relationship between annual maximum rainfall intensity for specified frequencies of occurrence (2-, 5-, 10-, 25-, 50- and 100-year return periods, i.e., $\tau = 0.5, 0.8, 0.9, 0.96, 0.98, 0.99$ -quantiles) and durations ($D = 5$ -, 10-, 15-, 30- and 60-min, 2-, 6-, 12- and 24-h) at locations in Canada with long records of short-duration rainfall rate observations. Annual maximum rainfall rate data for durations from 5-min to 24-h are archived by ECCC as part of the Engineering Climate Datasets (Environment and Climate Change Canada 2014). The rainfall rate dataset is based on tipping bucket rain gauge observations at 565 stations across Canada (Fig. 5). Record lengths range from 10 to 81-year, with a median length of 25-year. Information on the observing program, quality

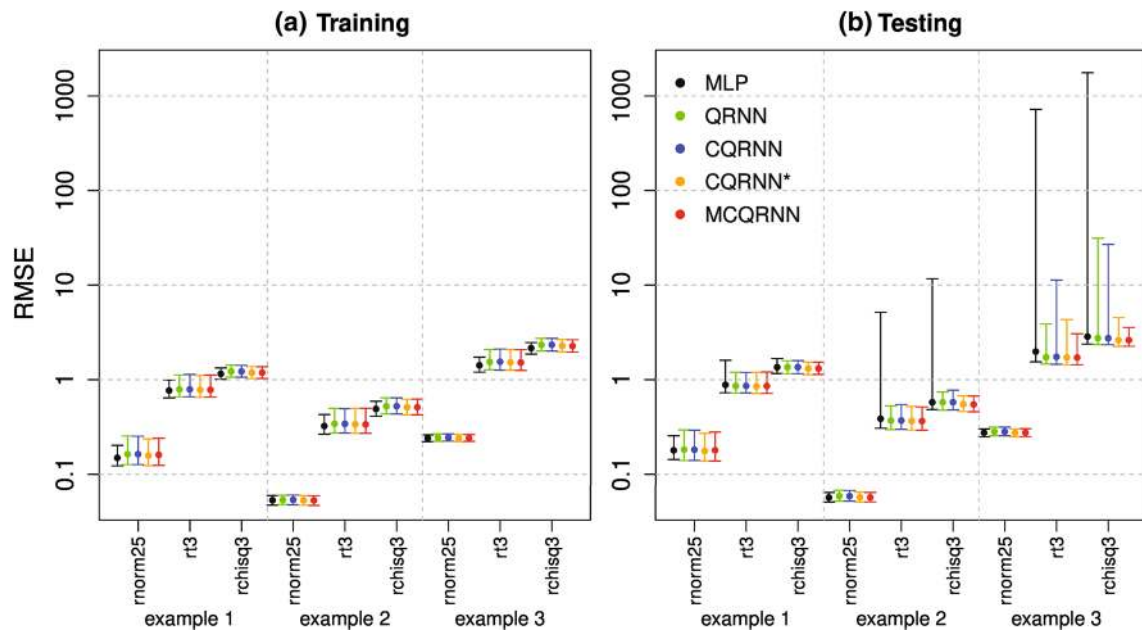


Fig. 4 Distribution of RMSE values over the 1000 Monte Carlo simulations for MLP (black), QRNN (green), CQRNN (blue), CQRNN* (orange) and MCQRNN (red) models in the **a** training and **b** testing datasets for examples 1, 2, and 3 from Xu et al. (2017)

with $N(0, 0.25)$ (rnorm25), $t(3)$ (rt3), and $\chi^2(3)$ (rchisq3) distributed noise. The central dot indicates the median RMSE and the lower and upper bars the 5th and 95th percentiles, respectively

control, and quality assurance methods is provided in detail by Shephard et al. (2014).

Official ECCC IDF curves are constructed by first fitting the parametric Gumbel distribution to annual maximum rainfall rate series at each site for each duration. At the majority of stations, the actual curves are then based on best fit linear interpolation equations between log-transformed duration and log-transformed Gumbel quantiles for each of the specified return periods. For reference, IDF curves for Victoria Intl A, a station on the southwest coast of British Columbia, Canada, are shown in Fig. 6. Points indicate return values of rainfall intensity obtained from the fitted Gumbel distribution for each combination of return period and duration; the IDF curves for each return period are based on log–log interpolating equations through these points, and hence plot as straight lines.

Naturally, the ECCC approach cannot provide quantile estimates for locations where short-duration rainfall observations are not recorded or available. Parametric extreme value distributions, fit in conjunction with regionalization or regional regression models, have been used to estimate IDF curves at ungauged locations in Canada by Alila (1999, 2000), Kuo et al. (2012) and Mailhot et al. (2013). As a non-parametric alternative to standard parametric approaches, Ouali and Cannon (2018) recently evaluated regional QRNN models for IDF curves at ungauged locations. While results suggest that the QRNN model can outperform standard parametric

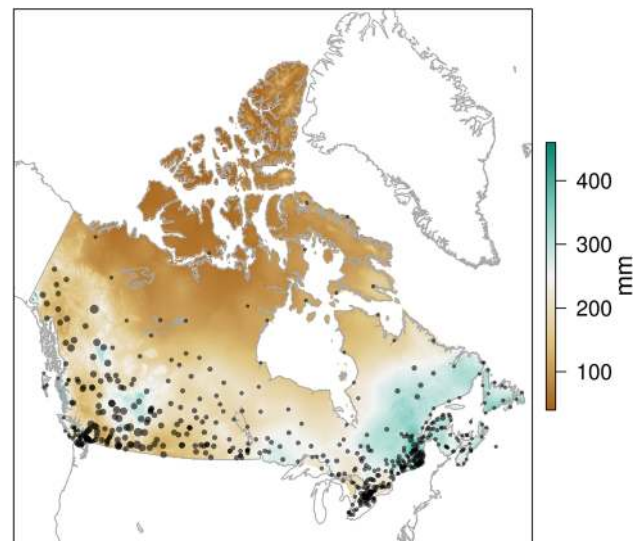


Fig. 5 Points (black circle) show locations of ECCC IDF curve stations; point size is proportional to station elevation. Shading indicates the climatological summer total precipitation (1971–2000)

methods, further improvements are still possible. In particular, Ouali and Cannon (2018) fit separate QRNN models for each τ -quantile and duration, which means that quantile crossing is possible; further, rainfall intensities may not increase as storm duration decreases. Instead, use of the MCQRNN is proposed to ensure non-crossing quantiles and a monotone decreasing relationship with

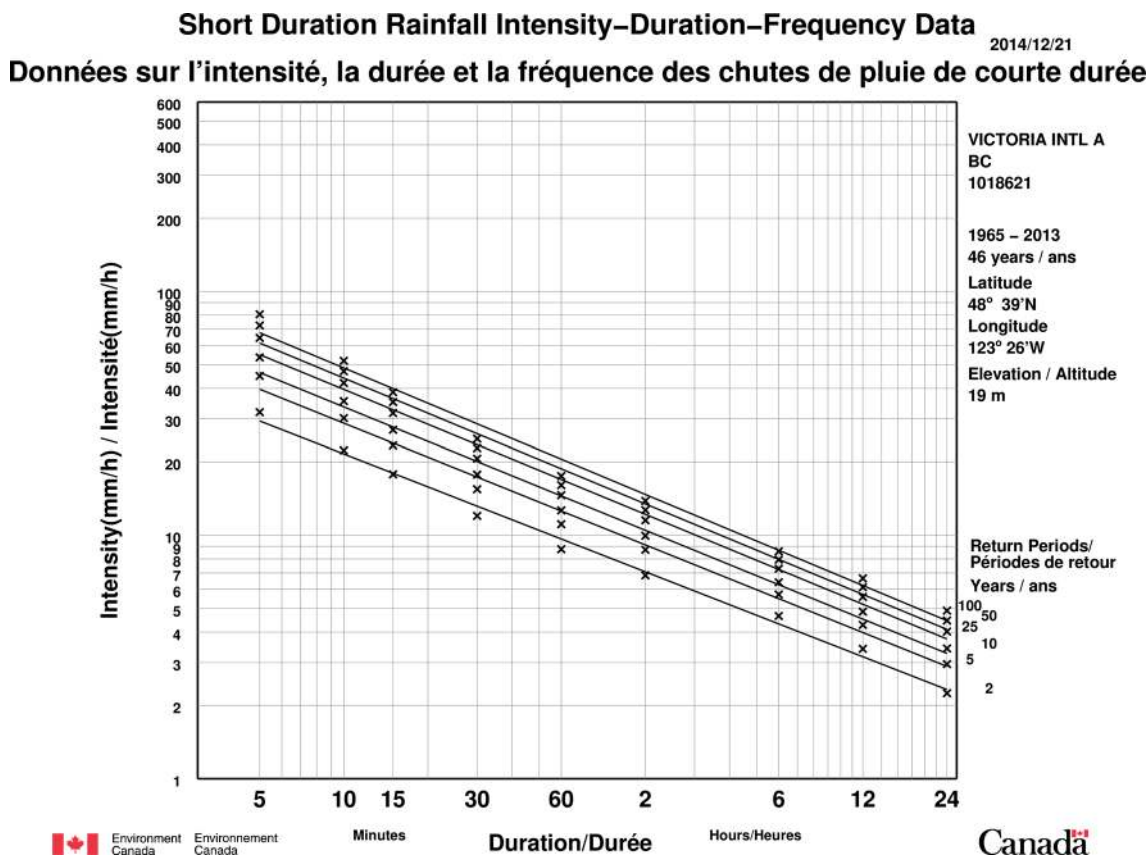


Fig. 6 Example ECCC IDF data for Victoria Intl A (station 1018621) in British Columbia, Canada. Points (×) show quantiles associated with 2-, 5-, 10-, 25-, 50- and 100-year (from bottom to top) return period intensities estimated by fitting the Gumbel distribution by the method of moments to annual maximum rainfall rate data for 5-, 10-,

15-, 30- and 60-min, 2-, 6-, 12- and 24-h durations (left to right). Lines are from best fit linear interpolation equations between log-transformed duration and log-transformed Gumbel quantiles for each return period

increasing storm duration. Estimation at ungauged sites typically relies on pooling gauged data from a homogeneous region around the site of interest, whether in geographic space or some derived hydroclimatological space (Ouarda et al. 2001), and then fitting a regression model linking spatial covariates with the short-duration rainfall rate response. As the focus of this study is on methods for conditional quantile estimation, and not the delineation of homogeneous regions, regionalizations here are based on a simple geographic region-of-influence (Burn 1990) in which data from the 80 nearest gauged sites are pooled together to form the training dataset for the site of interest. Following Aziz et al. (2014), this emphasizes the use of data from a large number of sites rather than the most homogeneous sites; it is then up to the regression model to infer relevant covariate-response relationships from within this larger pool of data. In areas with low station density, however, it is questionable whether any statistical regional

frequency analysis technique can be used to reliably estimate rainfall extremes. Performance in sparsely monitored regions will be explored as part of the subsequent model evaluation.

Based on this experimental design, observed short-duration rainfall rate data i_D for multiple durations D are used as the response variable in the MCQRNN model and spatial variables available over the domain—including at the ungauged location—are used as covariates in the regression equations. In this study, five covariates ($\#I = 5$), including latitude (lat), longitude (lon), elevation (elev), and climatological total winter (DJF) and summer precipitation (JJA) (Fig. 5) (McKenney et al. 2011), are used alongside the two ($\#M = 2$) monotone covariates [τ and $-\log(D)$]. As an abbreviated example, stacked data matrices for a single site (s_1), two quantiles (τ_1 and τ_2), and two durations (D_1 and D_2), for N years of short-duration rainfall observations would take the form:

$$\begin{aligned}
 \mathbf{y}_{s_1}^{(S)} &= \begin{bmatrix} i_{D_1}(1) \\ \vdots \\ i_{D_1}(N) \\ i_{D_2}(1) \\ \vdots \\ i_{D_2}(N) \\ i_{D_1}(1) \\ \vdots \\ i_{D_1}(N) \\ i_{D_2}(1) \\ \vdots \\ i_{D_2}(N) \end{bmatrix}, \\
 \mathbf{X}_{s_1}^{(S)} &= \begin{bmatrix} \tau_1 & -\log(D_1) & \text{lat}(s_1) & \text{lon}(s_1) & \text{elev}(s_1) & \text{DJF}(s_1) & \text{JJA}(s_1) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & -\log(D_1) & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & -\log(D_2) & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \tau_1 & -\log(D_2) & \vdots & \vdots & \vdots & \vdots & \vdots \\ \tau_2 & -\log(D_1) & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & -\log(D_1) & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & -\log(D_2) & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \tau_2 & -\log(D_2) & \text{lat}(s_1) & \text{lon}(s_1) & \text{elev}(s_1) & \text{DJF}(s_1) & \text{JJA}(s_1) \end{bmatrix}. \tag{20}
 \end{aligned}$$

For a given site of interest, the full stacked training dataset is expanded to include data from the 80 nearest gauged sites, 6 values of τ (0.5, 0.8, 0.9, 0.96, 0.98, 0.99), and 9 durations (5-, 10-, 15-, 30- and 60-min, 2-, 6-, 12- and 24-h).

4.2 Cross-validation results

Regional MCQRNN and QRNN models for IDF curves are evaluated via leave-one-out cross-validation. Each of the 565 observing sites is treated, in turn, as being “un-gauged”, i.e., data from nearest 80 sites to each left-out site are used to fit the models, model predictions are made at

the left-out site, and model performance statistics are calculated based on the left-out data. Following Ouali and Cannon (2018), 54 separate QRNN models are fit for each site, one for each combination of the 9 durations ($D = 5\text{-min to 24-h}$) and 6 τ -quantiles ($\tau = 0.5\text{--}0.99$) reported in ECCC IDF curves. Each MCQRNN model combines data for all 9 values of D and fits non-crossing quantile curves for the 6 τ -quantiles simultaneously.

Non-negativity constraints are imposed in both QRNN and MCQRNN models by setting g to the smooth ramp function (Cannon 2011). Monotonicity constraints—increasing with τ and decreasing with D —are imposed in the MCQRNN model by adopting the MMLP architecture with

Table 1 Summary of cross-validated relative differences RD_τ (%) in quantile regression error stratified by duration D , for all stations, for MCQRNN models (a) without weighting and (b) with weighting proportional to $\log(D)$

Return period/duration	5-min	10-min	15-min	30-min	60-min	2-h	6-h	12-h	24-h
(a) Unweighted									
2	- 0.1	- 0.2	0	+ 0.1	- 0.1	+ 0.4	+ 1.5	+ 2.7	+ 4.8
5	- 0.1	+ 0.2	+ 0.3	- 0.6	- 0.4	- 0.3	+ 1.0	+ 0.5	+ 1.9
10	+ 0.2	+ 0.1	+ 0.2	- 0.8	- 0.6	- 0.8	+ 0.7	+ 1.8	+ 1.7
25	+ 0.2	- 1.0	- 1.4	- 1.1	- 1.6	- 1.4	+ 1.1	+ 0.3	+ 0.6
50	- 2.1	- 3.5	- 3.9	- 1.9	- 1.1	- 6.7	+ 0.9	+ 0.8	+ 2.9
100	- 4.0	- 2.4	- 4.6	- 4.7	+ 1.6	+ 0.9	+ 2.8	+ 4.3	+ 5.6
(b) $\log(D)$ weighting									
2	+ 0.3	- 0.3	- 0.1	0	- 0.3	- 0.3	+ 0.2	+ 1.3	+ 2.9
5	+ 0.2	+ 0.2	+ 0.3	- 0.7	- 0.6	- 0.7	+ 0.1	- 0.2	+ 1.1
10	0	- 0.1	+ 0.1	- 0.9	- 0.8	- 1.0	- 0.1	+ 1.0	+ 0.9
25	+ 0.1	- 1.0	- 1.6	- 1.3	- 1.5	- 1.6	+ 0.3	- 0.8	- 0.8
50	- 2.1	- 3.6	- 4.1	- 2.4	- 1.4	- 7.0	+ 0.1	- 0.8	+ 0.7
100	- 3.3	- 2.5	- 5.0	- 5.6	+ 0.6	+ 0.3	+ 1.6	+ 1.7	+ 1.9

In both cases, QRNN IDF curve predictions serve as the reference model. Bold values indicate combinations of return period and duration for which MCQRNN performs better (i.e., lower errors) than QRNN; combinations with worse performance are underlined

additional monotone covariates $[\tau$ and $-\log(D)]$. The optimum level of complexity for each kind of model is selected based on values of QAIC, here based on the composite QR error function (e.g., Xu et al. 2017), averaged over all sites, from candidates with $J = 1, 2, \dots, 5$ (Koenker and Schorfheide 1994; Doksum and Koo 2000; Xu et al. 2017). The number of hidden nodes J is fixed to the same value for all sites in the study domain. QAIC is minimized for QRNN models with $J = 1$ and MCQRNN models with $J = 3$.

Cross-validation results comparing the MCQRNN ($J = 3$) and QRNN ($J = 1$) models are reported in terms of relative differences in leave-one-out estimates of the quantile regression error function

$$RD_\tau = 100 \left(\frac{E_\tau^{(MCQRNN)} - E_\tau^{(QRNN)}}{E_\tau^{(QRNN)}} \right) \tag{21}$$

summed over all stations for each return period and duration. Values are shown in Table 1a. Because the underlying model architecture is, aside from different values of J and inclusion of monotonicity constraints, fundamentally the same for the QRNN and MCQRNN models, it is not surprising that the two perform similarly well. MCQRNN and QRNN errors fall within 5% of one another for nearly all combinations of return period and duration, although MCQRNN tends to perform slightly better for short durations ($D = 5\text{-min}$ to 2-h) and QRNN for longer durations ($D = 6\text{-}24\text{-h}$). Poorer performance of the MCQRNN model in these cases is partly attributable to the smaller rainfall intensities that are associated with long duration storms being weighted less in the CQR cost function

(Eq. 14) than the larger intensities that accompany short duration storms. This can be remedied by setting $\omega_{\tau(s)} \propto \log(D)$ in Eq. 14. Results for the MCQRNN model with weighting are shown in Table 1b. Weighting improves performance for longer durations, while having minimal impact on shorter durations. Further results will be reported for the weighted MCQRNN model.

Despite the similar levels of quantile error, the additional MCQRNN monotonicity constraints on τ and D leads to IDF curves that are guaranteed to increase as occurrence frequency and storm duration decrease, properties that need not be present for QRNN predictions. This is evident for Victoria Intl A (Fig. 7), where quantile crossing and non-monotone increasing behaviour with decreasing storm duration is noted for the 100-year QRNN model predictions (cf. Fig. 6).

Each of the QRNN ($J = 1$) models for the 54 combinations of τ and D contain $J(\#I + 1) + J + 1 = 1(5 + 1) + 1 + 1 = 8$ parameters or 432 parameters in total. Because it borrows strength over τ and D ($\#M = 2$), the MCQRNN ($J = 3$) model requires just $J(\#I + \#M + 1) + J + 1 = 3(5 + 2 + 1) + 3 + 1 = 28$ shared parameters for the same task. Given that the two models show similar levels of performance, parameters in the separate QRNN equations must be largely redundant. If model complexity is increased, for example to $J = 5$, the total number of estimated parameters is 1944 for QRNN (36 for each combination of τ and D) versus 46 for MCQRNN. By way of comparison, the at-site (rather than ungauged) ECCC IDF curves require estimation of 30 parameters (18 Gumbel distribution and 12 interpolation equation parameters).

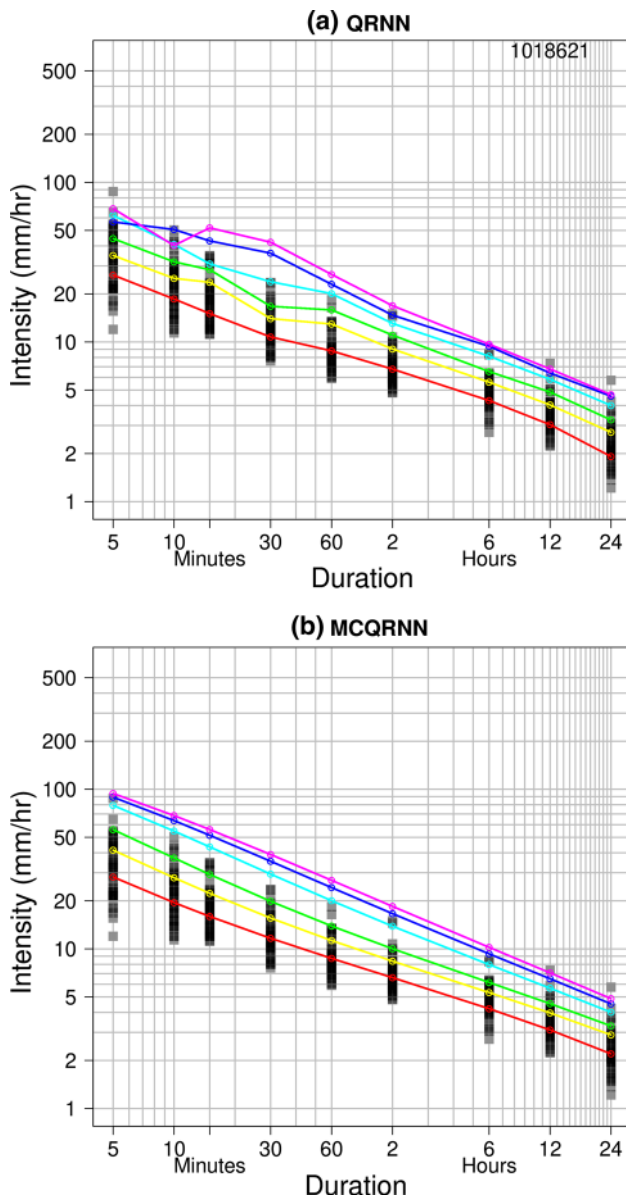


Fig. 7 Leave-one-out predictions of IDF curves for 2-, 5-, 10-, 25-, 50- and 100-year (in rainbow colours from bottom to top) return period intensities for Victoria Intl A (station 1018621) from **a** QRNN models and **b** MCQRNN model (cf. Fig. 6). Points (black square) show observed annual maximum rainfall rate data for 5-, 10-, 15-, 30- and 60-min, 2-, 6-, 12- and 24-h durations

Do the non-crossing/monotonicity constraints and ability to borrow strength provide a guard against overfitting if MCQRNN model complexity is misspecified? Fig. 8 shows relative differences RD_τ in cross-validated quantile regression error for MCQRNN and QRNN models with $J = 1, 2, \dots, 5$; in both cases, the optimal QRNN ($J = 1$) model serves as the reference. Consistent with results from QAIC model selection, cross-validated QRNN errors increase when $J > 1$. When using more than the recommended number of hidden nodes, the QRNN performs

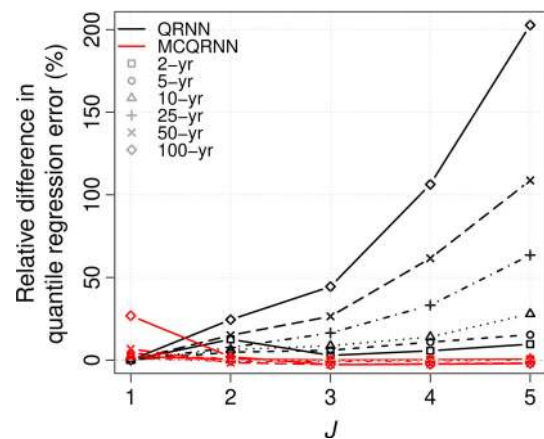


Fig. 8 Cross-validated relative differences RD_τ (%) in quantile regression error between MCQRNN and QRNN IDF curve predictions for $J = 1, 2, \dots, 5$ using QRNN ($J = 1$) as the reference model. Results are shown for 2-, 5-, 10-, 25-, 50- and 100-year return periods

poorly, especially for long return period estimates. However, for MCQRNN, in the absence of underfitting (i.e., $J = 1$), there is little penalty for specifying an overly complex model. Performance of the optimal MCQRNN ($J = 3$) model recommended by QAIC model selection is nearly identical to that of the misspecified $J = 5$ model. The non-crossing constraint provides strong regularization and resistance to overfitting.

Results reported so far have compared leave-one-out cross-validation performance of the MCQRNN and QRNN models. This does not provide any indication of how well the ungauged predictions compare with those estimated by the at-site ECCC IDF curve procedure, i.e., by fitting the Gumbel distribution and log linear interpolating equations to observed annual maxima at each station. Following Ouali and Cannon (2018), the ability of the MCQRNN to replicate the at-site ECCC IDF curves is measured by the quantile regression error ratio

$$R_\tau = \frac{E_\tau^{(ECCC)}}{E_\tau^{(MCQRNN)}} \tag{22}$$

where $E_\tau^{(ECCC)}$ is the in-sample, at-site quantile regression error of the ECCC IDF curve interpolating equations. A value of 1 means that ungauged MCQRNN predictions reach the same level of error as the at-site ECCC IDF curves. Note: even though the ECCC IDF curves are calculated from observations at each station, it is possible for R_τ to exceed 1 as the annual maximum rainfall data may deviate from the assumed Gumbel distribution and log linear form of the interpolating equations. Results are summarized in Table 2. Values of R_τ greater than 0.9—based on the 10% relative error threshold recommended by Mishra et al. (2012) for acceptable model simulations of urban rainfall extremes—are found for 41 of the 54

Table 2 Summary of quantile regression error ratio R_τ stratified by duration D between at-site ECCC IDF curves and ungauged MCQRNN predictions for all stations

Return period/duration	5-min	10-min	15-min	30-min	60-min	2-h	6-h	12-h	24-h
2	1.05	0.97	0.98	0.99	0.99	0.98	0.95	0.94	0.97
5	1.06	0.96	0.97	0.99	0.99	0.98	0.94	0.93	0.95
10	1.05	0.94	0.95	0.99	0.99	0.97	0.92	0.90	0.93
25	1.03	0.91	0.91	0.99	0.98	0.97	0.89	0.85	0.88
50	1.02	0.90	0.89	0.95	0.97	0.95	0.86	0.79	0.84
100	0.99	0.87	0.85	0.89	0.94	0.91	0.78	0.74	0.78

Values ≥ 0.9 are shown in bold

combinations of of D and τ , including all return periods from 2 to 10-year. More broadly, values exceed 0.7 for all combinations of D and τ .

As shown in Fig. 5, stations are not evenly distributed across Canada; northern latitudes, in particular, are very sparsely gauged. Does MCQRNN performance depend on station density? Values of R_τ , stratified by the median distance of each ungauged station to its 80 neighbours, are shown in Fig. 9. As expected, errors are nearly equivalent ($R_\tau > 0.975$) to the at-site estimates in areas of high station density (median distances < 100 km). Modest performance declines are noted ($R_\tau > 0.875$) with increasing median distance up to 500 km, beyond which performance degrades more substantially, especially for the longest return periods ($R_{\tau=0.99} < 0.8$). The viability of ungauged estimation should be evaluated carefully in areas of low station density.

5 Conclusion

This study introduces a novel form of quantile regression that can be used to simultaneously estimate multiple non-crossing, nonlinear quantile regression functions. MCQRNN is the first neural network-based quantile regression model that guarantees non-crossing of regression quantiles. The model architecture, which is based on the standard MLP neural network, also allows optional monotonicity, positivity/non-negativity, and generalized additive model constraints to be imposed in a straightforward manner. As an extension, a simple way to control the strength of non-additive relationships is also provided. The Huber function approximation to the QR error function means that standard least-squares regression and non-crossing expectile regression functions can be estimated using the same model architecture.

Given its close relationship to composite QR models, MCQRNN is first evaluated using the Monte Carlo simulation experiments adopted by Xu et al. (2017) to demonstrate the CQRNN model. In comparison to MLP, QRNN, and CQRNN models, MCQRNN is more robust than the benchmark models, especially for non-normal error

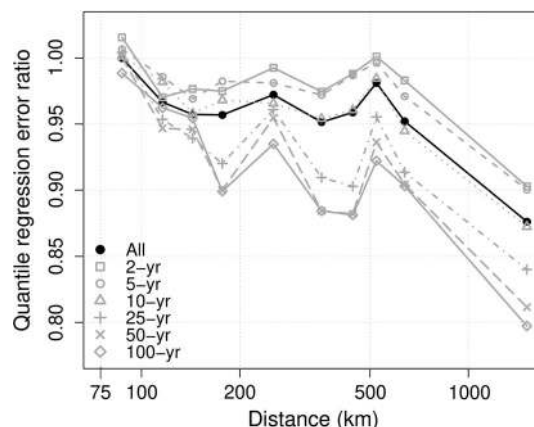
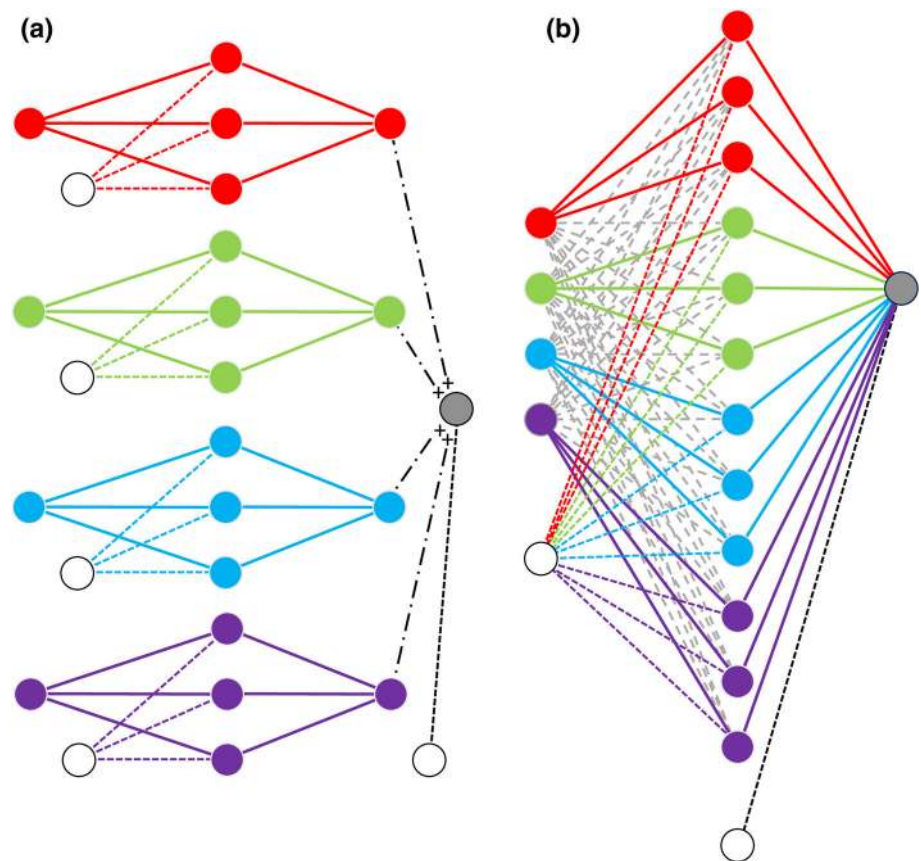


Fig. 9 Mean quantile regression error ratio R_τ between at-site ECCC IDF curves and leave-one-out cross-validated MCQRNN predictions; values of R_τ are stratified according to the median distance between the left-out station and its 80 neighbouring stations. Each of the 10 distance groupings contains an approximately equal numbers of stations (56 or 57)

distributions. Next, the MCQRNN model is evaluated on real-world climate data by estimating rainfall IDF curves in Canada. Cross-validation results suggest that the MCQRNN effectively borrows strength across different storm durations and return periods, which results in a model that is robust against overfitting. In comparison to standard QRNN, the ability of the MCQRNN model to incorporate monotonicity constraints—rainfall intensity should increase monotonically as the occurrence frequency and storm duration decrease—leads to more realistic estimates of extreme rainfall at ungauged sites. While promising, use of the MCQRNN for IDF curve estimation is presented here as a proof of concept. Other avenues of research include a more principled consideration of regionalization (Ouarda et al. 2001), other covariates (Madsen et al. 2017), and comparison against a wider range of nonlinear methods (Ouali et al. 2017). The MCQRNN model architecture is extremely flexible and many of its features are also not explored in this study. For example, the use of different weights for each τ in the composite QR error function (Jiang et al. 2012; Sun et al. 2013), multiple hidden layers, and the ability to estimate

Fig. 10 Schematic representations of **a** the generalized additive neural network architecture from Potts (1999) and **b** additivity constraints applied to a fully-connected MLP via a binary mask $\mathbf{A}^{(h)}$ applied to elements of $\mathbf{W}^{(h)}$. Parameters that have been set to zero by $\mathbf{A}^{(h)}$ are represented by dashed grey lines. Nonzero $\mathbf{W}^{(h)}$, \mathbf{w} parameters are represented by solid coloured lines, $\mathbf{b}^{(h)}$ parameters by dashed coloured lines, and b by dashed black lines



non-crossing, nonlinear expectile regression functions (Jiang et al. 2017) are left for future research.

Finally, code implementing the MCQRNN model is freely available from the Comprehensive R Archive Network as part of the qrnn package.

Acknowledgements The author would like to thank Dae Il Jeong, William Hsieh, Dhoucha Ouali, and the anonymous reviewers for their constructive feedback, and Cuixia Jiang for sharing their CQRNN computer code. The Comprehensive R Archive Network (CRAN) is acknowledged for hosting the qrnn package <https://CRAN.R-project.org/package=qrnn> for the R programming language and environment for statistical computing and graphics.

Open Access This article is distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made.

Appendix 1: Additive MLP models and control over non-additivity

As shown by Potts (1999), the MLP architecture used by the MCQRNN model can represent generalized additive relationships, i.e., where the model output depends on

linear combinations of unknown smooth functions applied to each covariate in turn. Each covariate is associated with its own MLP, separate from those for the other covariates (Fig. 10a), which means that interactions between covariates are neglected. The resulting model is easy to interpret, as contributions from covariates can be analyzed in isolation.

From Sect. 2.1—removing partial monotonicity constraints for sake of simplicity—this is equivalent to representing the hidden layer outputs in the form

$$h_j(t) = f\left(\sum_{i \in I} x_i(t) A_{ij}^{(h)} W_{ij}^{(h)} + b_j^{(h)}\right) \tag{23}$$

where $\mathbf{A}^{(h)}$ is an appropriate binary mask. For example, for a model with $\#I = 4$ covariates and $J = 3 (\#I) = 12$ hidden layer outputs, as shown in Fig. 10, the mask that enforces additive relationships is given by

$$\mathbf{A}^{(h)} = \begin{bmatrix} 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{bmatrix} \tag{24}$$

Each of the covariates x_i is passed through a smooth function defined, in this example, by a linear combination

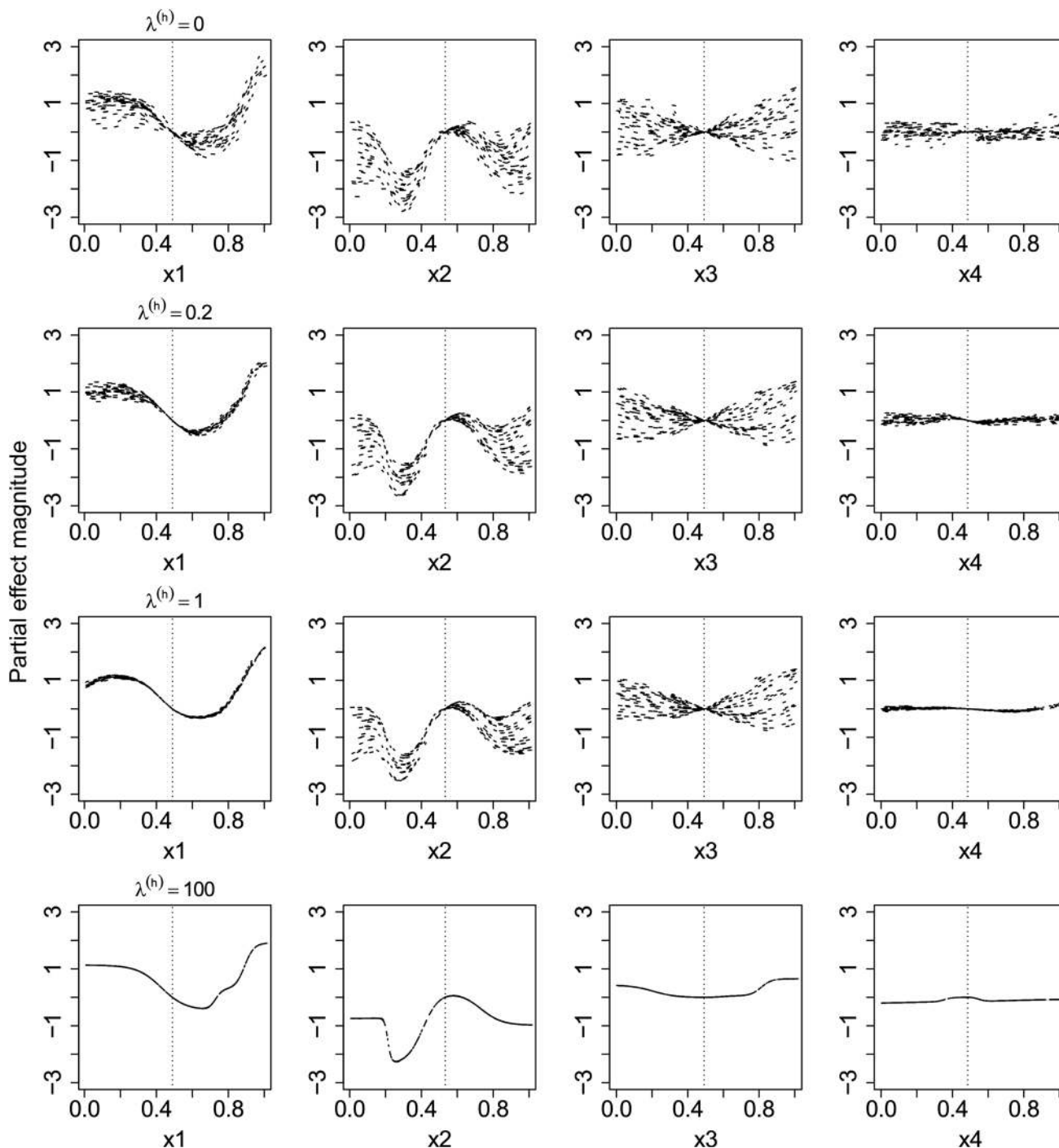


Fig. 11 Modified generalized additive model plots (Plate 1999) shows partial effects for covariates $x_1, x_2, x_3,$ and x_4 from MLP models ($\lambda^{(h)} = 0, 0.2, 1, 100$) fit to synthetic data generated by Eq. 27

of 3 hidden layer outputs. For a given covariate, the other hidden layer outputs, and hence covariates, do not contribute to the output because the additive mask multiplies the corresponding elements of $\mathbf{W}^{(h)}$ by zero (Fig. 10b).

A means of controlling non-additivity in a Gaussian process model was presented by Plate (1999). It was shown

that control over interactions in a flexible nonlinear model—allowing for models that range from being fully additive to those that do not constrain covariate interactions—can be beneficial for modelling tasks where interpretability and prediction performance are both important. Similar fine-grained control can be added to models based

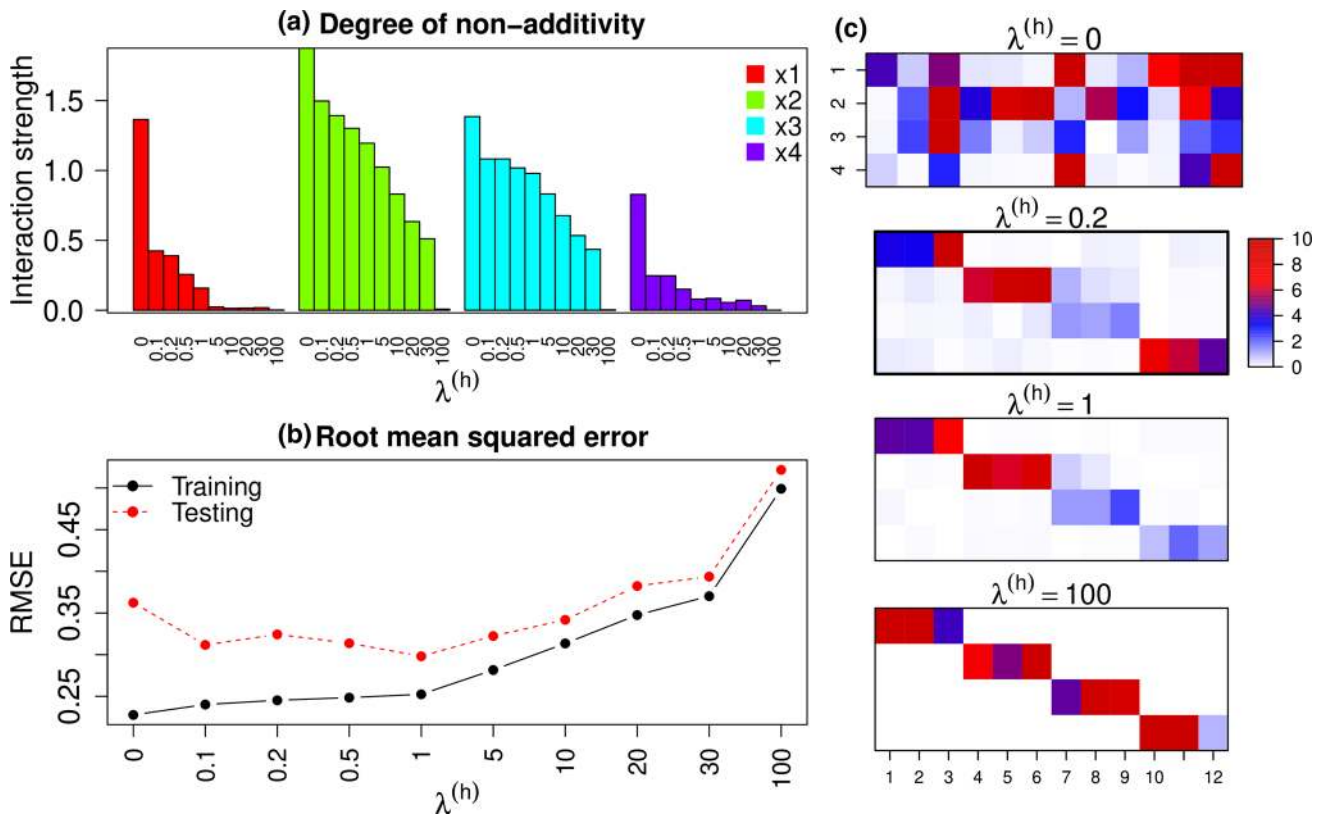


Fig. 12 **a** Interaction strength for covariates $x_1, x_2, x_3,$ and x_4 (Plate 1999), **b** training and testing RMSE, and **c** absolute magnitudes of $\mathbf{W}^{(h)}$ elements (cf. Eq. 26) associated with different values of $\lambda^{(h)}$

on the MLP architecture by removing $\mathbf{A}^{(h)}$ from Eq. 23 and instead modifying the error function

$$\tilde{E}_\tau^{(A)} = E_\tau^{(A)} + \lambda^{(h)} \frac{1}{VJ} \sum_{i=1}^V \sum_{j=1}^J L_{ij}^{(h)} \left(W_{ij}^{(h)} \right)^2 + \lambda \frac{1}{J} \sum_{j=1}^J (w_j)^2 \tag{25}$$

where

$$\mathbf{L}^{(h)} = \begin{bmatrix} 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 \end{bmatrix} \tag{26}$$

contains the logical negation of elements in the $\mathbf{A}^{(h)}$ matrix that would be applied in a fully-additive model. In effect, the first penalty term now applies only to elements of $\mathbf{W}^{(h)}$ responsible for controlling interactions between covariates; larger values of $\lambda^{(h)}$ will therefore suppress non-additive relationships.

To demonstrate, consider MLP models fit using the modified cost function (Eq. 25) to synthetic data generated by the function from Plate (1999)

$$y = 0.925\phi(x_1, x_2) + 2.248(x_2 + x_3 - 1)^3 + \varepsilon \tag{27}$$

where

$$\phi(x_1, x_2) = 1.3356 \left\{ 1.5(1 - x_1) + \exp(2x_1 - 1) \sin \left[3\pi(x_1 - 0.6)^2 \right] + \exp[3(x_2 - 0.5)] \sin \left[4\pi(x_2 - 0.9)^2 \right] \right\} \tag{28}$$

Covariate x_1 has a purely additive and nonlinear relationship with the response, while covariates x_2 and x_3 have an interactive, nonlinear relationship. A fourth covariate x_4 , which is irrelevant and does not contribute to the response, is also included. Two datasets are created: training data with 300 samples and testing data with 100,000 samples. Each of the four covariates is drawn from a uniform distribution $U(0, 1)$ and $\varepsilon \sim N(0, 0.5)$.

Figure 11 shows generalized additive model plots—modified following Plate (1999) so that non-additive relationships are indicated by vertical spread in points—for MLP models with $\lambda^{(h)} = 0, 0.2, 1, 100$. Values of $\lambda^{(h)} = 0, 0.2$ lead to spurious interactions for x_1 and x_4 , whereas $\lambda^{(h)} = 100$ suppresses the true interactions between x_2 and x_3 . $\lambda^{(h)} = 1$ appears to strike the appropriate balance, leading to a MLP model with a nonlinear additive relationship for x_1 , interactions for x_2 and x_3 , and no relationship between x_4 and the response. These results are

reflected in the measure of interaction strength, training and testing RMSE, and magnitudes of $\mathbf{W}^{(h)}$ elements shown in Fig. 12. The MLP with $\lambda^{(h)} = 1$ gives the lowest testing RMSE. This model has strong measured interactions for covariates x_2 and x_3 , which are associated with nonzero elements of $\mathbf{W}^{(h)}$.

References

- Alila Y (1999) A hierarchical approach for the regionalization of precipitation annual maxima in Canada. *J Geophys Res Atmos* 104(D24):31645–31655. <https://doi.org/10.1029/1999JD900764>
- Alila Y (2000) Regional rainfall depth–duration–frequency equations for Canada. *Water Resour Res* 36(7):1767–1778. <https://doi.org/10.1029/2000WR900046>
- Allamano P, Claps P, Laio F (2009) Global warming increases flood risk in mountainous areas. *Geophys Res Lett* 36(24):L24404. <https://doi.org/10.1029/2009GL041395>
- Aziz K, Rahman A, Fang G, Shrestha S (2014) Application of artificial neural networks in regional flood frequency analysis: a case study for Australia. *Stoch Env Res Risk Assess* 28(3):541–554. <https://doi.org/10.1007/s00477-013-0771-5>
- Baldwin RE (2006) In or Out: Does it matter? An evidence-based analysis of the Euro’s trade effects, chap. 2. Centre for Economic Policy Research (CEPR), London, p 110
- Bang S, Cho H, Jhun M (2016) Simultaneous estimation for non-crossing multiple quantile regression with right censored data. *Stat Comput* 26(1–2):131–147. <https://doi.org/10.1007/s11222-014-9482-0>
- Barbosa SM (2008) Quantile trends in Baltic sea level. *Geophys Res Lett* 35(22):L22704. <https://doi.org/10.1029/2008GL035182>
- Ben Alaya M, Chebana F, Ouarda T (2016) Multisite and multivariable statistical downscaling using a Gaussian copula quantile regression model. *Clim Dyn* 47(5–6):1383–1397. <https://doi.org/10.1007/s00382-015-2908-3>
- Bondell HD, Reich BJ, Wang H (2010) Noncrossing quantile regression curve estimation. *Biometrika* 97(4):825–838. <https://doi.org/10.1093/biomet/asq048>
- Burn DH (1990) Evaluation of regional flood frequency analysis with a region of influence approach. *Water Resour Res* 26(10):2257–2265. <https://doi.org/10.1029/WR026i10p02257>
- Canadian Standards Association (2012) PLUS 4013 (2nd ed.)—Technical guide: development, interpretation and use of rainfall intensity–duration–frequency (IDF) information: guideline for Canadian water resources practitioners. Canadian Standards Association, Mississauga
- Cannon AJ (2017) QRNN: Quantile regression neural network. R Package Version 2.0.2
- Cannon AJ (2011) Quantile regression neural networks: implementation in R and application to precipitation downscaling. *Comput Geosci* 37(9):1277–1284. <https://doi.org/10.1016/j.cageo.2010.07.005>
- Cawley GC, Janacek GJ, Haylock MR, Dorling SR (2007) Predictive uncertainty in environmental modelling. *Neural Netw* 20(4):537–549. <https://doi.org/10.1016/j.neunet.2007.04.024>
- Chen C (2007) A finite smoothing algorithm for quantile regression. *J Comput Graph Stat* 16(1):136–164. <https://doi.org/10.1198/106186007X180336>
- Chernozhukov V, Fernández-Val I, Galichon A (2010) Quantile and probability curves without crossing. *Econometrica* 78(3):1093–1125. <https://doi.org/10.3982/ECTA7880>
- Doksum K, Koo J-Y (2000) On spline estimators and prediction intervals in nonparametric regression. *Comput Stat Data Anal* 35(1):67–82. [https://doi.org/10.1016/S0167-9473\(99\)00116-4](https://doi.org/10.1016/S0167-9473(99)00116-4)
- Environment and Climate Change Canada (2014) Intensity–duration–frequency (IDF) files v2.30
- Friederichs P, Hense A (2007) Statistical downscaling of extreme precipitation events using censored quantile regression. *Mon Weather Rev* 135(6):2365–2378. <https://doi.org/10.1175/MWR3403.1>
- Hanson SJ, Burr DJ (1988) Minkowski-r back-propagation: learning in connectionist models with non-Euclidian error signals. In: *Neural information processing systems*, pp 348–357
- Hirschi M, Seneviratne SI, Alexandrov V, Boberg F, Boroneant C, Christensen OB, Formayer H, Orlowsky B, Stepanek P (2010) Observational evidence for soil-moisture impact on hot extremes in southeastern Europe. *Nature Geosci* 4(1):17–21. <https://doi.org/10.1038/ngeo1032>
- Hofmeister T (2017) qrsvm: SVM quantile regression with the pinball loss. R Package Version 0.2.1
- Huber PJ (1964) Robust estimation of a location parameter. *Ann Math Stat* 35(1):73–101
- Jiang X, Jiang J, Song X (2012) Oracle model selection for nonlinear models based on weighted composite quantile regression. *Stat Sin* 22:1479–1506. <https://doi.org/10.5705/ss.2010.203>
- Jiang C, Jiang M, Xu Q, Huang X (2017) Expectile regression neural network model with applications. *Neurocomputing* 247:73–86. <https://doi.org/10.1016/j.neucom.2017.03.040>
- Karatzoglou A, Smola A, Hornik K, Zeileis A (2004) kernlab: An S4 package for kernel methods in R. *J Stat Softw* 11(9):1–20
- Koenker R, Bassett G Jr (1978) Regression quantiles. *Econometrica* 46:33–50
- Koenker R, Schorfheide F (1994) Quantile spline models for global temperature change. *Clim Change* 28(4):395–404. <https://doi.org/10.1007/BF01104081>
- Kuo C-C, Gan TY, Chan S (2012) Regional intensity–duration–frequency curves derived from ensemble empirical mode decomposition and scaling property. *J Hydrol Eng* 18(1):66–74. [https://doi.org/10.1061/\(ASCE\)HE.1943-5584.0000612](https://doi.org/10.1061/(ASCE)HE.1943-5584.0000612)
- Lang B (2005) Monotonic multi-layer perceptron networks as universal approximators. In: *Artificial neural networks: formal models and their applications-ICANN*, vol 2005. pp 31–37. https://doi.org/10.1007/11550907_6
- Liu Y, Wu Y (2009) Stepwise multiple quantile regression estimation using non-crossing constraints. *Stat Interface* 2(3):299–310. <https://doi.org/10.4310/SII.2009.v2.n3.a4>
- Liu Y, Wu Y (2011) Simultaneous multiple non-crossing quantile regression estimation using kernel constraints. *J Nonparametr Stat* 23(2):415–437. <https://doi.org/10.1080/10485252.2010.537336>
- Madsen H, Gregersen IB, Rosbjerg D, Arnbjerg-Nielsen K (2017) Regional frequency analysis of short duration rainfall extremes using gridded daily rainfall data as co-variate. *Water Sci Technol* 75(8):1971–1981. <https://doi.org/10.2166/wst.2017.089>
- Mailhot A, Lachance-Cloutier S, Talbot G, Favre A-C (2013) Regional estimates of intense rainfall based on the Peak–Over–Threshold (POT) approach. *J Hydrol* 476:188–199. <https://doi.org/10.1016/j.jhydrol.2012.10.036>
- McKenney DW, Hutchinson MF, Papadopol P, Lawrence K, Pedlar J, Campbell K, Milewska E, Hopkinson RF, Price D, Owen T (2011) Customized spatial climate models for North America. *Bull Am Meteorol Soc* 92(12):1611–1622. <https://doi.org/10.1175/2011BAMS3132.1>
- Minin A, Velikova M, Lang B, Daniels H (2010) Comparison of universal approximators incorporating partial monotonicity by

- structure. *Neural Netw* 23(4):471–475. <https://doi.org/10.1016/j.neunet.2009.09.002>
- Mishra V, Dominguez F, Lettenmaier DP (2012) Urban precipitation extremes: How reliable are regional climate models? *Geophys Res Lett* 39:L03407. <https://doi.org/10.1029/2011GL050658>
- Muggeo VM, Sciandra M, Augugliaro L (2012) Quantile regression via iterative least squares computations. *J Stat Comput Simul* 82(11):1557–1569. <https://doi.org/10.1080/00949655.2011.583650>
- Muggeo VM, Sciandra M, Tomasello A, Calvo S (2013) Estimating growth charts via nonparametric quantile regression: a practical framework with application in ecology. *Environ Ecol Stat* 20(4):519–531. <https://doi.org/10.1007/s10651-012-0232-1>
- Newey WK, Powell JL (1987) Asymmetric least squares estimation and testing. *Econometrica* 55:819–847
- Ouali D, Chebana F, Ouarda T (2016) Quantile regression in regional frequency analysis: a better exploitation of the available information. *J Hydrometeorol* 17(6):1869–1883. <https://doi.org/10.1175/JHM-D-15-0187.1>
- Ouali D, Chebana F, Ouarda T (2017) Fully nonlinear statistical and machine-learning approaches for hydrological frequency estimation at ungauged sites. *J Adv Model Earth Syst* 9(2):1292–1306. <https://doi.org/10.1002/2016MS000830>
- Ouali D, Cannon AJ (2018) Estimation of rainfall intensity–duration–frequency curves at ungauged locations using quantile regression methods. *Stoch Environ Res Risk Assess*. <https://doi.org/10.1007/s00477-018-1564-7>
- Ouarda TB, Girard C, Cavadias GS, Bobée B (2001) Regional flood frequency estimation with canonical correlation analysis. *J Hydrol* 254(1):157–173. [https://doi.org/10.1016/S0022-1694\(01\)00488-7](https://doi.org/10.1016/S0022-1694(01)00488-7)
- Persson T (2001) Currency unions and trade: how large is the treatment effect? *Econ Policy* 33:435–448
- Plate TA (1999) Accuracy versus interpretability in flexible modeling: implementing a tradeoff using Gaussian process models. *Behaviormetrika* 26(1):29–50
- Potts WJ (1999) Generalized additive neural networks. In: *Proceedings of the fifth ACM SIGKDD international conference on knowledge discovery and data mining*. ACM, pp 194–200
- Quiñonero Candela J, Rasmussen CE, Sinz F, Bousquet O, Schölkopf B (2006) Evaluating predictive uncertainty challenge. *Lect Notes Comput Sci* 3944:1–27. https://doi.org/10.1007/11736790_1
- Roth M, Buishand T, Jongbloed G (2015) Trends in moderate rainfall extremes: a regional monotone regression approach. *J Clim* 28(22):8760–8769. <https://doi.org/10.1175/JCLI-D-14-00685.1>
- Saito H, Nakayama D, Matsuyama H (2010) Relationship between the initiation of a shallow landslide and rainfall intensity–duration thresholds in Japan. *Geomorphology* 118(1):167–175. <https://doi.org/10.1016/j.geomorph.2009.12.016>
- Shephard MW, Mekis E, Morris RJ, Feng Y, Zhang X, Kilcup K, Fleetwood R (2014) Trends in Canadian short-duration extreme rainfall: including an intensity–duration–frequency perspective. *Atmos Ocean* 52(5):398–417. <https://doi.org/10.1080/07055900.2014.969677>
- Sun J, Gai Y, Lin L (2013) Weighted local linear composite quantile estimation for the case of general error distributions. *J Stat Plan Inference* 143(6):1049–1063. <https://doi.org/10.1016/j.jspi.2013.01.002>
- Takeuchi I, Le QV, Sears TD, Smola AJ (2006) Nonparametric quantile estimation. *J Mach Learn Res* 7(Jul):1231–1264
- Taylor JW (2000) A quantile regression neural network approach to estimating the conditional density of multiperiod returns. *J Forecast* 19(4):299–311. [https://doi.org/10.1002/1099-131X\(200007\)19:4%3c299::AID-FOR775%3e3.0.CO;2-V](https://doi.org/10.1002/1099-131X(200007)19:4%3c299::AID-FOR775%3e3.0.CO;2-V)
- Waltrup LS, Sobotka F, Kneib T, Kauermann G (2015) Expectile and quantile regression—David and Goliath? *Stat Model* 15(5):433–456. <https://doi.org/10.1177/1471082X14561155>
- Wasko C, Sharma A (2014) Quantile regression for investigating scaling of extreme precipitation with temperature. *Water Resour Res* 50(4):3608–3614. <https://doi.org/10.1002/2013WR015194>
- White H (1992) Nonparametric estimation of conditional quantiles using neural networks. In: Page C, LePage R (eds) *Computing science and statistics*. Springer, pp 190–199. https://doi.org/10.1007/978-1-4612-2856-1_25
- Xu Q, Deng K, Jiang C, Sun F, Huang X (2017) Composite quantile regression neural network with applications. *Expert Syst Appl* 76:129–139. <https://doi.org/10.1016/j.eswa.2017.01.054>
- Yao Q, Tong H (1996) Asymmetric least squares regression estimation: a nonparametric approach. *J Nonparametr Stat* 6(2–3):273–292. <https://doi.org/10.1080/10485259608832675>
- Zhang H, Zhang Z (1999) Feedforward networks with monotone constraints. In: *IJCNN'99, International joint conference on neural networks*, vol 3. IEEE, pp 1820–1823. <https://doi.org/10.1109/IJCNN.1999.832655>
- Zou H, Yuan M (2008) Composite quantile regression and the oracle model selection theory. *Ann Stat* 36:1108–1126. <https://doi.org/10.1214/07-AOS507>