



HAL
open science

A REVIEW ON QUANTILE REGRESSION FOR STOCHASTIC COMPUTER EXPERIMENTS

Léonard Torossian, Victor Picheny, Robert Faivre, Aurélien Garivier

► **To cite this version:**

Léonard Torossian, Victor Picheny, Robert Faivre, Aurélien Garivier. A REVIEW ON QUANTILE REGRESSION FOR STOCHASTIC COMPUTER EXPERIMENTS. 2019. hal-02006032v1

HAL Id: hal-02006032

<https://hal.science/hal-02006032v1>

Preprint submitted on 4 Feb 2019 (v1), last revised 20 Jan 2020 (v3)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A REVIEW ON QUANTILE REGRESSION FOR STOCHASTIC COMPUTER EXPERIMENTS

Léonard Torossian

MIAT, Université de Toulouse, INRA
and Institut de Mathématiques de Toulouse
leonard.torossian@inra.fr

Victor Picheny

PROWLER.io, 72 Hills Road, Cambridge
victor@prowler.io

Robert Faivre

MIAT, Université de Toulouse, INRA
robert.faivre@inra.fr

Aurélien Garivier

Univ. Lyon, ENS de Lyon
aurelien.garivier@ens-lyon.fr

January 29, 2019

ABSTRACT

We report on an empirical study of the main strategies for conditional quantile estimation in the context of stochastic computer experiments. To ensure adequate diversity, six metamodels are presented, divided into three categories based on order statistics, functional approaches, and those of Bayesian inspiration. The metamodels are tested on several problems characterized by the size of the training set, the input dimension, the quantile order and the value of the probability density function in the neighborhood of the quantile. The metamodels studied reveal good contrasts in our set of 480 experiments, enabling several patterns to be extracted. Based on our results, guidelines are proposed to allow users to select the best method for a given problem.

1 Introduction

Computer simulation models are now an essential tool for performance evaluation, decision making, quality control and uncertainty quantification of complex systems.

These models generally depend on multiple input variables that can be divided into two groups, one controllable the other uncontrollable. The end-user will decide based on the controllable variable, but the decisions should account for the effect of uncontrollable variables. For example, in pharmacology, the optimal drug posology depends on the drug formulation but also on the individual targeted (genetic characters, age, sex) and environmental interactions. The shelf life and performance of a manufacturing device depend on its design, but also on its environment and on uncertainties during to the manufacturing process. The performance of a deep neural network depends on its (hyper)parameters and on the quality of the training set.

In such systems, the links between the inputs and outputs may be too complex to be fully understood or to be formulated in a closed form. In this case, the system can be considered as a black box and formalized by an unknown function: $f : \mathcal{X} \times \Omega \rightarrow \mathbb{R}$, where $\mathcal{X} \subset \mathbb{R}^d$ denotes a compact space representing the controllable variables, and Ω denotes the probability space representing the uncontrollable variables. Note that some black boxes have their outputs in a multidimensional space but this aspect is beyond the scope of the present paper. Because the space Ω may be large and/or complex, we do not assume that it has any structure. In contrast to deterministic systems, for any fixed x , $f(x, \cdot)$ is a random variable of law $\mathbb{P}_{\mathcal{X}}$; hence, such systems are often referred to as stochastic black boxes.

Based on the standard constraints encountered in computer experiments, throughout this paper, we assume that the function is only accessible through pointwise evaluations $f(x, \omega)$; no structural information is available regarding f ; in addition, we take into account the fact that evaluations may be expensive, which drastically limits the number of possible calls to f .

In order to understand the behavior of the system of interest or to take optimal decisions, information is needed about \mathbb{P}_X . An intuitive approach is to use a simple Monte-Carlo technique and evaluate $f(x, \omega_1), \dots, f(x, \omega_n)$ to extract statistical moments, the empirical cumulative distribution function, etc. Unfortunately a stratified approach is not practical in the case of an f that is expensive to evaluate.

Instead, we focus on *surrogate models* (also referred to as *metamodels* or *statistical emulators*), which are appropriate approaches in a small data setting associated with a regularity hypothesis (with respect to \mathcal{X}) concerning targeted statistics. Among the vast choice of surrogate models [73, 81], the most popular ones include regression trees, Gaussian processes, support vector machines and neural networks. In the framework of stochastic black boxes, the standard approach consists in estimating the conditional expectation of f . This case has been extensively treated in the literature and many applications, including Bayesian optimization [64], have been developed. However, the conditional expectation is risk-neutral, whereas pharmacologists, manufacturers, asset managers, data scientists and agronomists need to evaluate the worst case scenarios associated with their decisions, for example.

Risk information can be introduced by using a surrogate expectation-variance model [32, 65] or by estimating the whole distribution via heteroscedastic Gaussian processes [36, 40]. However, such approaches usually imply that the shape of the distribution (e.g. normal, uniform, etc.) is the same for all $x \in \mathcal{X}$. Another possible approach would be to learn the whole distribution \mathbb{P}_X with no strong structural hypotheses [52], but this requires a large number of evaluations of f . Here, we focus on the conditional quantile estimation of order τ , a flexible tool to tackle cases in which the distribution of $f(x, \cdot)$ varies markedly in spread and shape with respect to $x \in \mathcal{X}$, and a risk-awareness tool in decision theory [3].

1.1 Paper Overview

Many metamodels originally designed to estimate conditional expectations have been adapted to estimate the conditional quantile. However, despite extensive literature on estimating the quantile in the presence of spatial structure, few studies have reported on the constraints associated with stochastic black boxes. The performance of a metamodel with high dimension input is treated in insufficient details, performance based on the number of points has rarely been tackled and, to our knowledge, dependence on specific aspects of the quantile functions has never been studied. The aim of the present paper is to review quantile regression methods under standard constraints related to the stochastic black box framework, to provide information on the performance of the selected methods, and to recommend which metamodel to use depending on the characteristics of the computer simulation model and the data.

A comprehensive review of quantile regression is of course beyond the scope of the present work. We limit our review to the approaches that are best suited for our framework, while insuring the necessary diversity of metamodels. In particular, we have chosen six metamodels that are representative of three main categories: approaches based on statistical order as K-nearest neighbors (KN) regression [6] and random forest (RF) regression [49], functional or frequentist approaches as neural networks (NN) regression [15] and regression in reproducing kernel Hilbert space (RK) [75], and Bayesian approaches based on Gaussian processes as Quantile Kriging (QK) [55] and the variational Bayesian (VB) regression [1]. Each category has some specificities in terms of theoretical basis, implementation and complexity. The methods are described in full in sections 3 to 6.1.

In order to identify the relevant areas of expertise of the different metamodels, we designed an original benchmark system based on three toy functions and an agronomical model [16]. The dimension of the problems ranged from 1 to 9 and the number of observations from 40 to 2000. Particular attention was paid to the performance of each metamodel according to the size of the learning set, the value of the probability density function at the targeted quantile $\tilde{f}(\cdot, q_\tau)$ and the dimension of the problem. Sections 6 and 7 describe the benchmark system and detail its implementation, with particular focus on the tuning of the hyperparameters of each method. Full results and discussion are to be found in Sections 8 and 9, respectively.

Figure 1 summarizes our findings. In a nutshell, quantile regression requires large budgets: while the rule-of-thumb for computer experiments is a budget (i.e. number of experiments) 10 times the dimension, in our problems, we found that no method was able to provide a relevant quantile estimate with a number of observations less than 50 times the dimension. For larger budgets, no method works uniformly better than any other. NN and VB are best when the budget is large. When the budget is smaller, RF, RK, KN are best when the pdf is small in the neighborhood of the quantile, in other words, when little information is available. However, VB outperforms all the other methods when more information is available, that is, when the pdf is large in the neighborhood of the quantile.

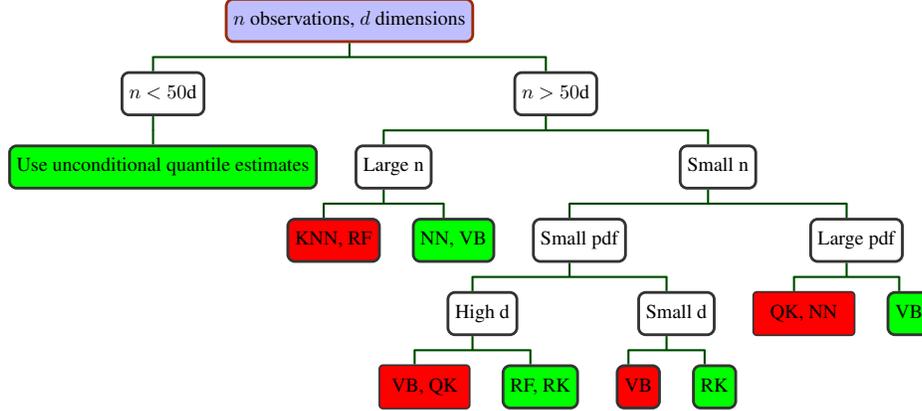


Figure 1: Method recommendation depending on the problem at hand (green: recommended methods, red: methods to avoid). KN: nearest-neighbors, RF: random forests, NN: neural networks, RK: RKHS regression, QK: quantile kriging, VB: variational Bayesian.

2 Quantile emulators and design of experiments

We first provide the necessary definitions, objects and properties related to the quantile. The quantile of order $\tau \in (0, 1)$ of a random variable Y can be defined either as the (generalized) inverse of a cumulative distribution function (CDF), or as the solution to an optimization problem:

$$q_\tau = \inf \{q : F(q) \geq \tau\} = \arg \inf_{q \in \mathbb{R}} \mathbb{E}_{\mathbb{P}_Y} [l_\tau(Y - q)], \quad (1)$$

where $F(\cdot)$ is the CDF of \mathbb{P}_Y and

$$l_\tau(\xi) = (\tau - \mathbb{1}_{(\xi < 0)})\xi, \quad \xi \in \mathbb{R} \quad (2)$$

is the so-called pinball loss [38] (Figure 2). In the following, we only consider situations in which F is continuous.

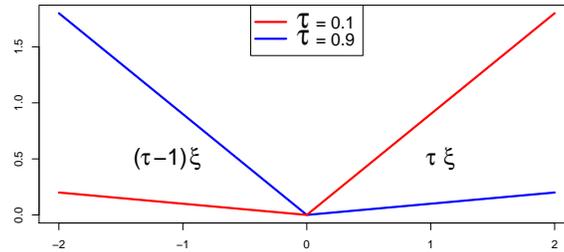


Figure 2: Pinball loss function with $\tau = 0.1$ and $\tau = 0.9$.

Given a finite observation set $\mathcal{Y}_n = (y_1, \dots, y_n)$ composed of i.i.d samples of Y , the empirical estimator of q_τ can thus be introduced in two different ways:

$$\hat{q}_\tau = \inf \{y_i \in \mathcal{Y}_n : \hat{F}(y_i) \geq \tau\} \quad (3)$$

$$= \arg \inf_{q \in \mathbb{R}} \frac{1}{n} \sum_{i=1}^n l_\tau(y_i - q), \quad (4)$$

where \hat{F} denotes the empirical CDF function. \hat{q}_τ coincides with the order statistic:

$$\hat{q}_\tau = y_{([n\tau])},$$

where $[n\tau]$ represents the smallest integer greater than or equal to $n\tau$ and $y_{(k)} = \mathcal{Y}_n^{(k)}$ is the k -th smallest value in the sample $\{y_1, \dots, y_n\}$. Similarly to (1), the conditional quantile of order $\tau \in (0, 1)$ can be defined in two equivalent ways:

$$q_\tau(x) = \{ \inf q : F(q|X = x) \geq \tau \} = \arg \inf_{q \in \mathcal{F}} \mathbb{E}_{\mathbb{P}_{X,Y}} [l_\tau(Y - q(X))], \quad (5)$$

where $F(\cdot|X = x)$ is the conditional CDF of $\mathbb{P}_{X,Y}$ and \mathcal{F} is a functional space containing q_τ .

In a quantile regression context, one only has access to a finite observation set $\mathcal{D}_n = \{(x_1, y_1), \dots, (x_n, y_n)\} = (\mathcal{X}_n, \mathcal{Y}_n)$ with \mathcal{X}_n a $n \times d$ matrix. Estimators for (5) are either based on the order statistic as in (3) (section 3), or on the pinball loss as in (4) (sections 4 and 5). Throughout this work, the observation set \mathcal{D}_n is fixed (we do not consider a dynamic or sequential framework). Following the standard approach used in computer experiments, the training points x_i are chosen according to a space-filling design [17] over a hyperrectangle. In particular, we assume that there are no repeated experiments: $x_i \neq x_j, \forall i \neq j$; most of the methods chosen in this survey (KN, RF, RK, NN, VB) work under that setting.

However, as a baseline approach, one may decide to use a stratified experimental design $\mathcal{D}_{n',r}$ with r i.i.d samples for a given $x_i, i = 1, \dots, n'$, extract pointwise quantile estimates using (3) and fit a standard metamodel to these estimates. The immediate drawback is that for the same budget ($n' \times r = n$) such experimental designs cover much less of the design space than a design with no repetition. The QK method is based on this approach.

3 Methods based on order statistics

A simple way to compute a quantile estimate is to take an order statistic of an i.i.d. sample. A possible approach is to emulate such a sample by selecting all the data points in the neighborhood of the query point x , and then by taking the order statistic of this subsample as an estimator for the conditional quantile. One may simply choose a subsample of \mathcal{D}_n based on a distance defined on \mathcal{X} : this is what the K -nearest neighbors approach does. It is common to use KN based on the Euclidean distance but of course any other distance can be used, such as Mahalanobis [80] or weighted Euclidean distance [28]. Alternatively, one may define a notion of neighborhood using some space partitioning of \mathcal{X} . That includes all the decision tree methods [14], in particular regression trees, bagging or random forest [49].

3.1 K -nearest neighbors

The K -nearest neighbors method was first proposed for the estimation of conditional expectations [71, 72]. Its extension to the conditional quantile estimation can be found in [6].

3.1.1 Quantile regression implementation

KN works as follows: define $\mathcal{X}_n^K(x^*)$ the subset of \mathcal{X}_n containing the K points that are the closest to the query point x^* . Define $\mathcal{Y}_n^K(x^*)$ the associated outputs, and define $\hat{F}(y|X = x^*)$ as the associated empirical CDF. According to (3), the conditional quantile of order τ is define as

$$\hat{q}_\tau(x^*) = \mathcal{Y}_n^K([K\tau]). \quad (6)$$

Algorithm 1 details the implementation of the KN method.

Algorithm 1 K -nearest neighbors

- 1: **Inputs:**
 $\mathcal{D}_n, \tau, K, \mathcal{X}_{test}$
 - 2: **for** each point in $x^* \in \mathcal{X}_{test}$ **do**
 - 3: Compute all the distances between x^* and \mathcal{X}_n
 - 4: Sort the computed distances
 - 5: Select the K -nearest points from x^*
 - 6: $\hat{q}_\tau(x^*) = \mathcal{Y}_n^K([K\tau])$
 - 7: **end for**
-

3.1.2 Computational complexity

For a naive implementation of such an estimator, one needs to compute $n \times N_{new}$ distances, where N_{new} is the number of query points, hence for a cost in $O(nN_{new}d)$. Moreover, sorting n distances in order to extract the K nearest points has a cost in $O(nN_{new} \log(n))$. Combining the two operations implies a complexity of order

$$O(nN_{new}d) + O(nN_{new} \log(n)).$$

Note that some algorithms have been proposed in order to reduce the computational time, for example by using GPUs [30] or by using tree search algorithms [4].

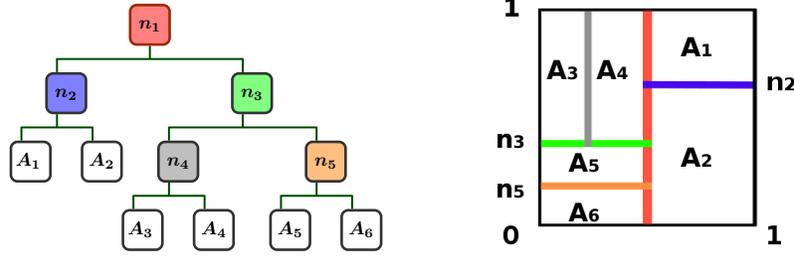


Figure 3: Left: a partitioning tree T . The nodes n_i ($1 \leq i \leq 5$) represent the splitting points, the A_i 's ($1 \leq i \leq 6$) represent the leaves. Right: $\mathcal{X} = [0, 1]^2$ as partitioned by T . The regression tree prediction is constant on each leaf A_i .

3.2 Random forests

Random forests were introduced by Breiman [13] for the estimation of conditional expectations. They have been used successfully for classification and regression, especially with problems where the number of variables is much larger than the number of observations [25].

3.2.1 Overview

The basic element of random forests is the *regression tree* T , a simple regressor built via a binary recursive partitioning process. Starting with all data in the same partition *i.e.* \mathcal{X} , the following sequential process is applied. At each step, the data is split in two, so that \mathcal{X} is partitioned in a way that can be represented by a tree as it is presented Figure 3.

Several splitting criteria can be chosen (see [34]). In [49], the splitting point x_S minimizes

$$C(x_s) = \sum_{x_i \leq x_s} (y_i - \bar{Y}_L)^2 + \sum_{x_j > x_s} (y_j - \bar{Y}_R)^2, \quad (7)$$

where \bar{Y}_L and \bar{Y}_R are the mean of the left and right sub-populations, respectively. Equation (7) applies when the x 's are real-valued. In the multidimensional case, the dimension d_S in which the split is performed has to be selected. The split then goes through x_S and perpendicularly to the direction d_S . There are several rules to stop the expansion of T . For instance, the process can be stopped when the population of each cell is inferior to a minimal size *nodesize*: then, each node becomes a terminal node or leaf. The result of the process is a partition of the input space into hyperrectangles $R(T)$. Like the KN method, the tree-based estimator is constant on each neighborhood. However the regression tree framework automatically builds neighborhoods from the data that should be adapted to each problem.

Despite their simplicity of construction and interpretation, regression trees are known to suffer from a high variance. To overcome this drawback, regression trees can be used with ensemble methods like bagging. Instead of using only one tree, bagging creates a set of tree $\mathcal{T}_N = \{T^1, \dots, T^N\}$ based on a bootstrap version $\mathcal{D}_{N,n} = \{((x_1, y_1), \dots, (x_n, y_n))\}_{t=1}^N$ of \mathcal{D}_n . Then the final model is created by averaging the results among all the trees.

Although bagging reduces the variance of the predictor, as the splitting criterion has to be optimized over all the input dimensions, computing (7) for each possible split is costly when the dimension is large. The random forest algorithm, a variant of bagging, constructs an ensemble of weak learners based on $\mathcal{D}_{N,n}$ and aggregates them. Unlike plain bagging, at each node evaluation, the algorithm uses only a subset of p covariables for the choice of the split dimensions. Because the p covariables are randomly chosen, the result of the process is a random partition $R(t)$ of \mathcal{X} constructed by the random tree T^t .

3.2.2 Quantile prediction

We present the extension proposed in [49] for conditional quantile regression. Let us define $l(x^*, t)$ the leaf obtained from the tree t containing a query point x^* and

$$\omega_i(x^*, t) = \frac{\mathbb{1}_{\{x_i \in l(x^*, t)\}}}{\#\{j : x_j \in l(x^*, t)\}}, \quad i = 1, \dots, n$$

$$\bar{\omega}_i(x^*) = \frac{1}{N} \sum_{t=1}^N \omega_i(x^*, t).$$

The $\bar{\omega}_i(x^*)$'s represent the weights illustrating the ‘‘proximity’’ between x^* and x_i . In the classical regression case, the estimator of the expectation is:

$$\hat{\mu}(x^*) = \sum_{i=1}^n \bar{\omega}_i(x^*) y_i. \quad (8)$$

In [49] the conditional quantile of order τ is defined as in (3) with the CDF estimator defined as

$$\hat{F}(y|X = x^*) = \sum_{i=1}^n \bar{\omega}_i(x^*) \mathbb{1}_{\{y_i \leq y\}}. \quad (9)$$

Algorithm 2 details the implementation of the RF method.

Algorithm 2 Random forest

```

1: Training
2: Inputs:
    $\mathcal{D}_n, N, p$ 
3: for each of the  $N$  trees do
4:   Uniformly sample with replacement  $n$  points in  $\mathcal{D}_n$  to create  $\mathcal{D}_{t,n}$ .
5:   Consider the cell  $R = \mathcal{X}$ .
6:   while any cell of the tree contains more than nodesize observations do
7:     for the cells containing more than nodesize observations do
8:       Uniformly sample without replacement  $p$  covariables in  $1, \dots, D$ .
9:       Compute the cell point among the  $p$  covariables that minimizes (7).
10:      Split the cell at this point perpendicularly to the selected covariable.
11:     end for
12:   end while
13: end for
14: Prediction
15: Inputs:
    $\mathcal{X}_{test}, \tau$ 
16: for each point in  $x^* \in \mathcal{X}_{test}$  do
17:   Compute  $\bar{\omega}_i(x^*)$ ,  $i = 1 \dots, n$ 
18:    $\hat{F}(y|X = x^*) = \sum_{i=1}^n \bar{\omega}_i(x^*) \mathbb{1}_{\{y_i \leq y\}}$ 
19:    $\hat{q}_\tau(x^*) = \inf \{y_i : \hat{F}(y_i|X = x^*) \geq \tau\}$ 
20: end for

```

3.2.3 Computational complexity

Assuming that the value of (7) can be computed sequentially for consecutive thresholds, the RF computation burden lies in the search of the splitting point that implies sorting the data. Sorting n variables has a complexity in $O(n \log(n))$. Thus, at each node the algorithm finds the best splitting points considering only $p \leq d$ covariables. This implies a complexity of $O(pn \log(n))$ per node. In addition, the depth of a tree is generally upper bounded by $\log(n)$. Then the computational cost of building a forest containing N trees under the criterion (7) is

$$O(Npn \log^2(n))$$

[44][82]. One may observe that RF are easy to parallelize and that contrary to KN the prediction time is very small once the forest is built.

4 Approaches based on functional analysis

Functional methods search directly for the function mapping the input to the output in a space fixed beforehand by the user.

Estimating any functional S of $\mathbb{P}_{X,Y}$ with this framework implies selecting a loss l (associated to S) and a function space \mathcal{H} . Thus, the estimator $\hat{S} \in \mathcal{H}$ is obtained as the minimizer of the empirical risk \mathcal{R}_e associated to l , *i.e.*

$$\hat{S} \in \arg \min_{s \in \mathcal{H}} \mathcal{R}_e[s] = \arg \min_{s \in \mathcal{H}} \frac{1}{n} \sum_{i=1}^n l(y_i - s(x_i)). \quad (10)$$

The functional space \mathcal{H} must be chosen large enough to extract some signal from the data. In addition, \mathcal{H} needs to have enough structure to make the optimization procedure feasible (at least numerically). In the literature, several formalisms such as linear regression [63], spline regression [46], support vector machine [79], neural networks [8] or deep neural networks [61] use structured functional spaces with different levels of flexibility.

However, using a too large \mathcal{H} can lead to overfitting, *i.e.* return predictors that are good only on the training set and generalize poorly. Overcoming overfitting requires some *regularization* [62] [85] [86], defining for example the regularized risk

$$\mathcal{R}_{r,e}[s] = \frac{1}{n} \sum_{i=1}^n l(y_i - s(x_i)) + \frac{\lambda}{\beta} \|s\|_{\mathcal{H}}^{\beta}, \quad (11)$$

where $\lambda \in \mathbb{R}^+$ is a penalization factor, $\beta \in \mathbb{R}^+$ and $\|\cdot\|_{\mathcal{H}}^2$ is either a norm for some methods (Section 4.2) or a measure of variability for others (Section 4.1). The parameter λ plays a major role, as it allows to tune the balance between bias and variance.

Classically, squared loss is used: it is perfectly suited to the estimation of the conditional expectation. Using the pinball loss (Eq. 2) instead allows to estimate quantiles. In this section we present two approaches based on Equation (11) with the pinball loss. The first one is regression using artificial neural networks (NN), a rich and versatile class of functions that has shown a high efficiency in several fields. The second approach is the generalized linear regression in reproducing kernel Hilbert spaces (RK). RK is a non-parametric regression method that has been much studied in the last decades (see [69]) since it appeared in the core of learning theory in the 1990's [62, 79].

4.1 Neural Networks

Artificial neural networks have been successfully used for a large variety of tasks such as classification, computer vision, music generation, and regression [8]. In the regression setting, feed-forward neural networks have shown outstanding achievements. Here we present quantile regression neural network [15] which is an adaptation of the traditional feed-forward neural network.

4.1.1 Overview

A feed-forward neural network is defined by its number of hidden layers H , its numbers of neurons per layer J_h , $1 \leq h \leq H$, and its activation functions g_h , $h = 1, \dots, H$. Given an input vector $x \in \mathbb{R}^d$ the information is fed to the hidden layer 1 composed of a fixed number of neurones J_1 . For each neurone $N_i^{(h_1)}$, $i = 1, \dots, J_1$, a scalar product (noted $\langle \cdot, \cdot \rangle$) is computed between the input vector $x = (x_1, \dots, x_d) \in \mathbb{R}^d$ and the weights $w_i^{(h_1)} = (w_{i1}^{(h_1)}, \dots, w_{id}^{(h_1)}) \in \mathbb{R}^d$ of the $N_i^{(h_1)}$ neurones. Then a bias term $b_i^{h_1} \in \mathbb{R}$ is added to the result of the scalar product. The result is composed with the activation function g_1 (linear or non-linear) which is typically the sigmoid or the ReLU function [61] and the result is given to the next layer where the same operation is processed until the information comes out from the output layer. For example, the output of a 3-layers NN at x^* is given by

$$s(x^*) = g_3 \left(\sum_{j=1}^{J_2} g_2 \left(\sum_{i=1}^{J_1} g_1 \left(\langle w_i^{(h_1)}, x^* \rangle + b_i^{(1)} \right) \omega_j^{(h_2)} + b_j^{(2)} \right) \omega_j^{(h_3)} + b^{(3)} \right). \quad (12)$$

The corresponding architecture can be found in Figure 4.

The architecture of the NN defines \mathcal{H} . Finding the right architecture is a very difficult problem which will not be treated in this paper. However a classical implementation procedure consists of creating a network large enough (able to overfit) and then using techniques such as early stopping, dropout, bootstrapping or risk regularization to avoid overfitting [67]. In [15], the following regularized risk is used:

$$\mathcal{R}_{r,e}[s] = \frac{1}{n} \sum_{i=1}^n l(y_i - s(x_i)) + \sum_{j=1}^H \frac{\lambda}{J_j} \sum_{z=1}^{J_j} \left\| w_z^{(h_j)} \right\|^2. \quad (13)$$

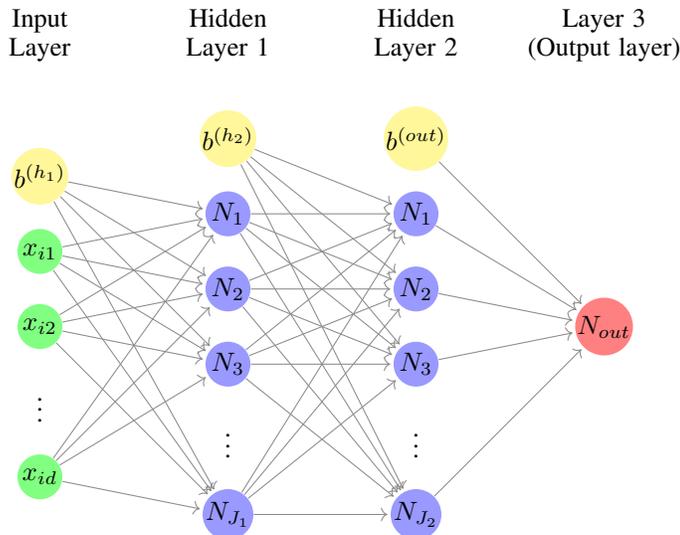


Figure 4: Architecture of 3-layer feedforward neural network.

4.1.2 Quantile regression

Minimizing Equation (13) (with respect to all the weights and biases) is in general challenging, as $\mathcal{R}_{r,e}$ is a highly multimodal function. It is mostly tackled using derivative-based algorithms and multi-starts (*i.e.* launching the optimization procedure m_s times with different starting points). In the case of quantile estimation case, the loss function is non-differentiable at the origin, which may cause problems to some numerical optimization procedures. To address this issue, [15] introduced a smooth version of the pinball loss function, defined as:

$$l_\tau^\eta(\xi) = \tau h^\eta(\xi) - \mathbf{1}_{\xi < 0},$$

where

$$h^\eta(u) = \begin{cases} \frac{\xi^2}{2\eta} & \text{if } 0 \leq |\xi| \leq \eta \\ |\xi| - \frac{\eta}{2} & \text{if } |\xi| \geq \eta. \end{cases} \quad (14)$$

Let us define \mathbf{w} the list containing the weights and bias of the network. To find \mathbf{w}^* , a minimizer of $\mathcal{R}_{r,e}$, the idea is to solve a series of problems using the smoothed loss instead of the pinball one with a sequence E_K corresponding to K decreasing values of η . The process begins with the optimization with the larger value η_1 . Once the optimization converges, the optimal weights are used as the initialization for the optimization with η_2 , and so on. The process stops when the weights based on $l_\tau^{\eta_K}$ are obtained. Finally, $\hat{q}_\tau(x^*)$ is given by the evaluation of the optimal network at x^* . Algorithm 3 details the implementation of the NN method.

4.1.3 Computational complexity

In [15] the optimization is based on a Newton method. Thus the procedure needs to inverse a Hessian matrix. Without sophistications, its cost is $O(s_{pb}^3)$ with s_{pb} the size of the problem *i.e.* the number of parameters to optimize. Note that using a high order method makes sense here because NN has few parameters (in contrast to deep learning methods). Moreover providing an upper bound on the number of iterations needed to reach an optimal point may be really hard in practice because of the non convexity of (13). In the non-convex case, there is no optimality guaranty and the optimization could be stuck in a local minima. However, it can be shown that the convergence near a local optimal point is at least super linear and may be quadratic (with some additional hypotheses). It implies, for each η , the number of iterations until $\mathcal{R}_{r,e}^\eta(\mathbf{w}) - \mathcal{R}_{r,e}^\eta(\mathbf{w}^*) \leq \epsilon$ is bounded above by

$$\frac{\mathcal{R}_{r,e}^\eta(\mathbf{w}_0) - \mathcal{R}_{r,e}^\eta(\mathbf{w}^*)}{\gamma} + \log_2 \log_2(\epsilon_0/\epsilon),$$

with γ the minimal decreasing rate, $\epsilon_0 = 2M_\eta^3/L_\eta^2$, M_η the strong convexity constant of $\mathcal{R}_{r,e}^\eta$ near \mathbf{w}^* and L_η the Hessian Lipschitz constant (see [11] page 489). As $\log_2 \log_2(\epsilon_0/\epsilon)$ increases very slowly with respect to ϵ , it is

Algorithm 3 Neural network

-
- 1: **Training**
 - 2: **Inputs:**
 $\mathcal{D}_n, \tau, \lambda, H, (J_1, \dots, J_H), (g_1, \dots, g_H), E_K$
 - 3: **Initialize:**
 Fix w_0 as the list containing the initial weights and biases.
 - 4: **for** $t = 1$ to K **do**
 - 5: $\epsilon \leftarrow E_K[t]$
 - 6: Starting the optimization procedure with w_0 and define

$$w^* = \arg \min_w \frac{1}{n} \sum_{t=1}^n l_\tau^\epsilon(y_i - q(x_i)) + \sum_{j=1}^H \frac{\lambda}{J_j} \sum_{i=1}^J \|w_i^{(h_j)}\|^2$$

with $q(x_i)$ the output of the network.

- 7: $w_0 \leftarrow w^*$
 - 8: **end for**
 - 9: **Prediction**
 - 10: **Inputs:**
 $\mathcal{X}_{test}, w^*, \lambda, H, (J_1, \dots, J_H), (g_1, \dots, g_H)$
 - 11: **for** each point in $x^* \in \mathcal{X}_{test}$ **do**
 - 12: Define $\hat{q}_\tau(x^*)$ as the output of the network evaluated at x^* .
 - 13: **end for**
-

possible to bound the number of iterations N typically by

$$\frac{\mathcal{R}_{r,e}^\eta(w_0) - \mathcal{R}_{r,e}^\eta(w^*)}{\gamma} + 6.$$

That means, near an optimal point, the complexity is $O(L_\eta n (Jd)^3)$, with J the total number of neurons. Then using a multistart procedure implies a complexity of

$$O(m_s L_{\eta^*} n (Jd)^3),$$

with $L_{\eta^*} = \max_{\eta_1, \dots, \eta_K} L_\eta$.

4.2 Generalized linear regression

Regression in RKHS was introduced for classification via Support Vector Machine by [22, 33], and has been naturally extended for the estimation of the conditional expectation [27, 58]. Since, many applications have been developed [69, 62], here we present the quantile regression in RKHS [75, 60].

4.2.1 RKHS introduction and formalism

Under the linear regression framework, S is assumed to be under the form $S(x) = x^T \alpha$, with α in \mathbb{R}^d . To stay in the same vein while creating non-linear responses, one can map the input space \mathcal{X} to a space of higher dimension \mathcal{H} (named the feature space), thanks to a feature map Φ . For example the feature space could be a polynomial space, in that case we are working with the spline framework [46]. For a large flexibility and few parameters, the feature space can even be chosen as an infinite dimensional space. In the following, $\Phi = (\phi_1, \phi_2, \phi_3, \dots)$ defines a feature map from \mathcal{X} to \mathcal{H} , where \mathcal{H} is the \mathbb{R} -Hilbert functional space defined as

$$\mathcal{H} = \left\{ s, s(x) = \sum_{i \in I} \alpha_i \phi_i(x), \text{ s.t. } \|s\|_{\mathcal{H}} < +\infty \right\},$$

$$\text{with } \|s\|_{\mathcal{H}} := \sqrt{\langle s, s \rangle_{\mathcal{H}}},$$

where I is the cardinal of \mathcal{H} . Under the hypothesis that S belongs to \mathcal{H} , S can be written as

$$S(x) = \sum_{i \in I} \alpha_i \phi_i(x). \tag{15}$$

Notice that without more hypothesis on \mathcal{H} , estimating S is difficult. In fact it is impossible to compute (15) directly because of the infinite sum. However, this issue may be tackled by using the RKHS formalism and the so-called *kernel trick* [62]. Let us first introduce the symmetric definite positive function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ such that:

$$k(x, x') = \langle \Phi(x'), \Phi(x) \rangle_{\mathcal{H}}. \quad (16)$$

Under this setting, \mathcal{H} is a RKHS with the reproducing kernel k , that means $\Phi(x) = k(\cdot, x) \in \mathcal{H}$ for all $x \in \mathcal{X}$ and the reproducing property

$$s(x) = \langle s, k(\cdot, x) \rangle_{\mathcal{H}} \quad (17)$$

holds for all $s \in \mathcal{H}$ and all $x \in \mathcal{X}$. It can be shown that working with a fixed kernel k is equivalent to working with its associated functional Hilbert space.

The kernel choice is based on kernel properties or assumptions made on the functional space. See for instance [69], chapter 4, for some kernel definitions and properties. In the following, \mathcal{H}_θ and k_θ denote respectively a RKHS and its kernel associated to the hyperparameters vector θ . $K_{x,x}^\theta \in \mathcal{R}^{n \times n}$ is the kernel matrix obtained via $K_{x,x}^\theta(i, j) = k_\theta(x_i, x_j)$.

From a theoretical point of view, the representer theorem implies that the minimizer \hat{S} of

$$\mathcal{R}_{r,e}[s] = \frac{1}{n} \sum_{i=1}^n l(y_i - s(x_i)) + \frac{\lambda}{2} \|s\|_{\mathcal{H}_\theta}^2$$

lives in $\mathcal{H}_{|X}^\theta = \text{span}\{\Phi(x_i) : i = 1, \dots, n\}$ with $\|s\|_{\mathcal{H}_{|X}^\theta}^2 = \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_\theta(x_j, x_i)$. Thanks to the reproducing property (17), and since $\hat{S} \in \mathcal{H}_{|X}^\theta$,

$$\hat{S}(x) = \left\langle \sum_{i=1}^n \alpha_i \Phi(x_i), \Phi(x) \right\rangle.$$

Using the definition (16), it is possible to rewrite \hat{S} as:

$$\hat{S}(x) = \sum_{i=1}^n \alpha_i k_\theta(x, x_i).$$

Hence, the original infinite dimensional problem (15) becomes an optimization problem over n coefficients $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_n) \in \mathbb{R}^n$. More precisely, finding \hat{S} is equivalent to minimizing in α the quantity

$$\frac{1}{n} \sum_{i=1}^n l(y_i - (\sum_{j=1}^n \alpha_j k_\theta(x_i, x_j))) + \frac{\lambda}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_\theta(x_j, x_i). \quad (18)$$

4.2.2 Quantile regression

Quantile regression in RKHS was introduced by [75], followed by several authors [42, 70, 20, 21, 60]. Quantile regression has two specificities compared to the general case. Firstly the loss l is defined as the pinball. Secondly, to ensure the quantile property, the intercept is not regularized. More precisely, we assume that

$$q_\tau(x) = g(x) + b \text{ with } g \in \mathcal{H}_\theta \text{ and } b \in \mathbb{R}.$$

and we consider the empirical regularized risk

$$\mathcal{R}_{r,e}[q] := \frac{1}{n} \sum_{i=1}^n l_\tau(y_i - q(x_i)) + \frac{\lambda}{2} \|g\|_{\mathcal{H}_\theta}^2. \quad (19)$$

Thus the representer theorem implies that \hat{q}_τ can be written under the form

$$\hat{q}_\tau(x^*) = \sum_{i=1}^n \alpha_i k_\theta(x^*, x_i) + b,$$

for a new query point x^* . Since (19) cannot be minimized analytically, a numerical minimization procedure is used. [22] followed by [75] introduced nonnegative variables $\xi^{(*)} \in \mathbb{R}^+$ to transform the original problem into

$$\mathcal{R}_{r,s}[q] := \frac{1}{n} \sum_{i=1}^n \tau \xi_i + (1 - \tau) \xi_i^* + \frac{\lambda}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j k_\theta(x_j, x_i),$$

subject to

$$y_i - \left(\sum_{j=1}^n \alpha_j k_\theta(x_i, x_j) + b \right) \leq \xi_i$$

and

$$\sum_{j=1}^n \alpha_j k_\theta(x_i, x_j) + b - y_i \leq \xi_i^*, \text{ where } \xi_i^*, \xi_i^* \geq 0.$$

Using a Lagrangian formulation, it can be shown that minimizing $\mathcal{R}_{r,e}$ is equivalent to the problem:

$$\begin{aligned} \min_{\alpha \in \mathbb{R}^n} \quad & \frac{1}{2} \alpha^T K_{x,x}^\theta \alpha - \alpha^T \mathbf{y} \\ \text{s.t.} \quad & \frac{1}{\lambda n} (\tau - 1) \leq \alpha_i \leq \frac{1}{\lambda n} \tau, \forall 1 \leq i \leq n \\ \text{and} \quad & \sum_{i=1}^n \alpha_i = 0. \end{aligned} \tag{20}$$

It is a quadratic optimization problem under linear constraint, for which many efficient solvers exist.

The value of b may be obtained from the Karush-Kuhn-Tucker slackness condition or fixed independently of the problem. A simple way to do so is to choose b as the τ -quantile of $(y_i - \sum_{j=1}^n \alpha_j k_\theta(x_i, x_j))_{1 \leq i \leq n}$. Algorithm 4 details the implementation of the RK method.

Algorithm 4 RKHS regression

1: **Training**

2: **Inputs:**

$$\mathcal{D}_n, \tau, \lambda, k_\theta$$

3: **Initialize:**

Compute the $n \times n$ matrix $K_{x,x}^\theta$

4: **Optimization:** Select α^* as

$$\begin{aligned} \alpha^* = \arg \min_{\alpha \in \mathbb{R}^n} \quad & \frac{1}{2} \alpha^T K_{x,x}^\theta \alpha - \alpha^T \mathbf{y} \\ \text{s.t.} \quad & \frac{1}{\lambda n} (\tau - 1) \leq \alpha_i \leq \frac{1}{\lambda n} \tau, \forall 1 \leq i \leq n \\ \text{and} \quad & \sum_{i=1}^n \alpha_i = 0 \end{aligned}$$

5: Define b as the τ -quantile of $(y_i - \sum_{j=1}^n \alpha_j k_\theta(x_i, x_j))_{1 \leq i \leq n}$

6: **Prediction**

7: **Inputs:**

$$\mathcal{X}_{test}, \alpha^*, k_\theta$$

8: **for** each point in $x^* \in \mathcal{X}_{test}$ **do**

9: compute $K_{x^*,x}^\theta$

10: $\hat{q}_\tau(x_{test}) = K_{x^*,x}^\theta \alpha^* + b$

11: **end for**

4.2.3 Computational complexity

Let us notice two things. Firstly, the minimal upper bound complexity for solving (20) is $O(n^3)$. Indeed solving (20) without the constraints is easier and it needs $O(n^3)$. Secondly the optimization problem (20) is convex, thus the optimum is global.

There are two main approaches for solving (20), the interior point method [11] and the iterative methods like libSVM [18]. The interior point method is based on the Newton algorithm, one method is the barrier method (see [11] p.590). It is shown that the number of iterations N for reaching a solution with precision ϵ is $O(\sqrt{n} \log(\frac{n}{\epsilon}))$. Moreover each

iteration of a Newton type algorithm costs $O(n^3)$ because it needs to inverse a Hessian. Thus, the complexity of an interior point method for finding a global solution with precision ϵ is

$$O\left(n^{7/2} \log\left(\frac{n}{\epsilon}\right)\right).$$

On another hand, iterative methods like libSVM transform the main problem into a smaller one. At each iteration the algorithm solves a subproblem in $O(n)$. Contrary to the interior point methods, the number of iterations depends explicitly on the matrix $K_{x,x}^\theta$. [43] shows that the number of iterations is

$$O\left(n^2 \kappa(K_{x,x}^\theta) \log(1/\epsilon)\right),$$

where $\kappa(K_{x,x}^\theta) = \lambda_{\max}(K_{x,x}^\theta) / \lambda_{\min}(K_{x,x}^\theta)$. Note that $\kappa(K_{x,x}^\theta)$ depends on the type of the kernel, it evolves in $O(n^s)$ with $s > 1$ an increasing value of the regularity of $K_{x,x}^\theta$ [19]. For more information about the eigenvalues of $K_{x,x}^\theta$ one can consult [12].

To summarize, it implies that the complexity of the libSVM method has an upper bound higher than the interior point algorithm. However, these algorithms are known to converge pretty fast. In practice, the upper bound is almost never reached, and thus the most important factor is the cost per iteration, rather than the number of iterations needed. That is the reason why libSVM is popular in this setting.

5 Bayesian approaches

Bayesian formalism has been used for a wide class of problems such as classification and regression [57], model averaging [10] and model selection [56].

Under the regression framework, a classical Bayesian model is defined as

$$y = S(x) + \varepsilon(x), \quad (21)$$

where S depends on parameters that have to be inferred, y is an observation and $\varepsilon(x)$ a noise term. Knowing \mathcal{D}_n , the goal is to estimate the posterior distribution $p(S|\mathcal{D}_n)$ at x^*

$$p(S(x^*)|\mathcal{D}_n) = \int p(S(x^*)|S, \zeta, \mathcal{D}_n) p(S, \zeta|\mathcal{D}_n) dS d\zeta,$$

with ζ a vector of parameters. The posterior $p(S, \zeta|\mathcal{D}_n)$ is a priori unknown and $p(S(x^*)|S, \zeta, \mathcal{D}_n)$ is the model hypothesis. According to the Bayes formula, the posterior can be written as

$$p(S, \zeta|\mathcal{D}_n) = \frac{p(\mathcal{Y}_n|S, \mathcal{X}_n, \zeta) p(S, \zeta)}{p(\mathcal{Y}_n|\mathcal{X}_n)}.$$

Because the normalizing constant is independent of S and ζ , considering only the likelihood and the prior is enough. We obtain

$$p(S, \zeta|\mathcal{D}_n) \propto p(\mathcal{Y}_n|S, \mathcal{X}_n, \zeta) p(S, \zeta).$$

with the likelihood

$$p(\mathcal{Y}_n|S, \mathcal{X}_n, \zeta) = \prod_{i=1}^n p(y_i|S, x_i, \zeta). \quad (22)$$

The Bayesian quantile regression framework can be summarized as follows. Starting from Equation (22), there are several estimation possibilities. The first idea was introduced in [83] where the authors worked under a linear framework, *i.e.* $q_\tau(x) = x^T \zeta$ but the linear hypothesis is too restrictive to treat the stochastic black box setting. [74] introduced a mixture modeling framework called Dirichlet process to perform nonlinear quantile regression. However the inference is performed with MCMC methods [31, 29], a procedure that is often costly. A possible alternative is the utilisation of Gaussian process (GP). GPs are powerful in a Bayesian context because of their flexibility and their tractability (GP are only characterized by their mean m and covariance k^θ). In this section we present QK and VB, two approaches using a Gaussian process (GP) as a prior for q_τ .

Note that while QK and VB use GP, both metamodels are intrinsically different. For QK, $q_\tau \sim \text{GP}$ is a consequence of the hypothesis made on ε . However VB uses the artificial assumption $q_\tau \sim \text{GP}$ in order to simplify the estimation procedure while ensuring flexibility.

5.1 Quantile kriging

Kriging takes its origins in geostatistics and spatial data interpolation [23, 68]. Since the 2000's, kriging drew attention of the machine learning community [57]. We present a very intuitive method that gives flexible quantile estimators based on data containing repetition and GPs [55].

5.1.1 Kriging introduction

In the kriging framework, ε in (21) is assumed to be Gaussian, *i.e.*

$$y_i = S(x_i) + \varepsilon_i \quad \text{with} \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_i^2), \quad (23)$$

and

$$S(x) = \Phi(x)^T \boldsymbol{\alpha},$$

with

$$\boldsymbol{\alpha} \sim \mathcal{N}(0, \Sigma_p), \quad \text{with} \quad \Sigma_p \in \mathbb{R}^{I \times I}. \quad (24)$$

Note that here $\zeta = \boldsymbol{\alpha}$. Under this assumption, the likelihood (22) is given by

$$\begin{aligned} p(\mathcal{Y}_n | \Phi, \mathcal{X}_n, \boldsymbol{\alpha}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(y_i - \Phi(x_i)^T \boldsymbol{\alpha})^2}{2\sigma_i^2}\right) \\ &= \mathcal{N}(\Phi(X)^T \boldsymbol{\alpha}, \boldsymbol{\sigma}^2 I), \end{aligned} \quad (25)$$

with $\boldsymbol{\sigma} = (\sigma_1, \dots, \sigma_n)$. It can be shown (see [57]) that under the hypotheses (23) and (24)

$$S(x^*) \sim \mathcal{N}(\bar{S}, \mathbb{V}_S(x^*)) \quad (26)$$

with

$$\begin{aligned} \bar{S}(x^*) &= K_{x^*, x}^\theta (K_{x, x}^\theta + B)^{-1} \mathcal{Y}_n, \\ \mathbb{V}_S(x^*) &= k_\theta(x^*, x^*) - K_{x^*, x}^\theta (K_{x, x}^\theta + B)^{-1} K_{x, x^*}, \end{aligned}$$

B being a diagonal matrix of size n with diagonal terms equals to σ_i^2 which represents the observation noise. As in Section 4.2, the covariance functions are usually chosen among a set of predefined ones that depend on a set of hyperparameters $\theta \in \mathbb{R}^{d+1}$.

The best hyperparameter θ^* can be selected as the maximizer of the marginal likelihood, which can be computed by marginalizing over S :

$$p(\mathcal{Y}_n | \mathcal{X}_n, \theta) = \int p(\mathcal{Y}_n | S, \mathcal{X}_n, \theta) p(S | \mathcal{X}_n, \theta) dS.$$

In addition, since $p(\mathcal{Y}_n | S) = \mathcal{N}(S, \sigma_n^2 I)$ it follows [57] that

$$\begin{aligned} p(\mathcal{Y}_n | \mathcal{X}_n, \theta) &= -\frac{1}{2} \mathcal{Y}_n^T (K_{x, x}^\theta + B)^{-1} \mathcal{Y}_n \\ &\quad -\frac{1}{2} \log |(K_{x, x}^\theta + B)| - \frac{n}{2} \log(2\pi), \end{aligned} \quad (27)$$

where $|K|$ is the determinant of the matrix K . Maximizing this likelihood with respect to θ is usually done using derivative-based algorithms, although the problem is non-convex and known to have several local maxima.

Different estimators of S may be extracted based on (26). Here \hat{S} is fixed as \bar{S}_{θ^*} . Note that this classical choice is made because the maximum a posteriori of a Gaussian distribution coincides with its mean.

5.1.2 Quantile kriging prediction

As q_τ is a latent quantity, the solution proposed in [55] is to consider the sample $\mathcal{D}_{n', r}$ which represents a design of experiments with n' different points that are repeated r times in order to obtain quantile observations. For each $x_i \in \mathcal{X}$, $1 \leq i \leq n'$, let us define:

$$y_{i, r} = (y_{i, 1}, \dots, y_{i, r})$$

and

$$\mathcal{D}_{n', \tau, r} = ((x_1, \hat{q}_\tau(x_1)), \dots, (x_n, \hat{q}_\tau(x_n))), \quad \text{with} \quad \hat{q}_\tau(x_i) = y_{i, \lfloor r\tau \rfloor}.$$

Following [55], let us assume that

$$\hat{q}_\tau(x_i) = q_\tau(x_i) + \varepsilon_i, \quad \text{with} \quad \varepsilon_i \sim \mathcal{N}(0, \sigma_i^2). \quad (28)$$

Algorithm 5 Quantile kriging

```

1: Training
2: Inputs:
    $\mathcal{D}_{n',r,\tau}, k_\theta$ 
3: Initialize:
   Compute the  $n' \times n'$  matrix  $K_{x,x}^\theta$ 
   Define the local estimators of the  $\tau$ -quantile
4: for  $i = 1$  to  $n'$  do
5:    $\hat{q}_\tau(x_i) \leftarrow y_{i,(\lceil r\tau \rceil)}$ 
6:   Estimate  $\sigma_i$  by bootstrap
7: end for
8: Define  $B = \text{Diag}(\sigma_1^2, \dots, \sigma_{n'}^2)$  and compute  $(K_{x,x}^\theta + B)^{-1}$ 
9: Define the kernel hyperparameters  $\theta^*$  as  $\arg \max_\theta$  of
   
$$-\frac{1}{2} \mathcal{Q}_{n'}^T (K_{x,x}^\theta + B)^{-1} \mathcal{Q}_{n'} - \frac{1}{2} \log |K_{x,x}^\theta + B| - \frac{n}{2} \log(2\pi)$$

10: Inputs:
    $\mathcal{X}_{test}, \mathcal{Q}_{n'}, \theta^*, B$ 
11: for each point in  $x^* \in \mathcal{X}_{test}$  do
12:   compute  $K_{x^*,x}^{\theta^*}$ 
13:    $\hat{q}_\tau^{\theta^*}(x^*) = K_{x^*,x}^{\theta^*} (K_{x,x}^{\theta^*} + B)^{-1} \hat{q}_\tau$ 
14: end for

```

Note that from a statistical point of view the assumption (28) is wrong however asymptotically consistent. The resulting estimator is

$$\hat{q}_\tau(x^*) = K_{x^*,x}^\theta (K_{x,x}^\theta + B)^{-1} \mathcal{Q}_{n'}, \quad (29)$$

with $\mathcal{Q}_{n'} = (\hat{q}_\tau(x_1), \dots, \hat{q}_\tau(x_{n'}))$ and $B = \text{diag}(\sigma_1^2, \dots, \sigma_{n'}^2)$.

There are several possibilities to evaluate the noise variances σ_i^2 . Here we choose to use a bootstrap technique (that is, generate bootstrapped samples of $y_{i,r}$, compute the corresponding $\hat{q}_\tau(x_i)$ values and take the variance over those values as the noise variance). The hyperparameters are selected based on (27) changing \mathcal{Y}_n by $\mathcal{Q}_{n'}$. Algorithm 5 details the implementation of the QK method.

5.1.3 Computational complexity

Optimizing (27) with a Newton type algorithm implies to inverse a $(d+1) \times (d+1)$ matrix. In addition, for each θ , obtaining the partial derivatives of (27) requires the computation of $(K_{x,x}^\theta + B)^{-1}$ [57]. Thus at each step of the algorithm, the complexity is $O(n'^3 + d^3)$. Assuming the starting point θ_{start} is close to an optimal θ^* , based on the same analysis as in section 4.1.3, the complexity to find θ^* is

$$O(L(d^3 + n'^3)).$$

Finally, obtaining $\hat{q}_\tau^{\theta^*}$ from (29) implies inverting the matrix $K_{x,x}^{\theta^*} + B$ that is in $O(n'^3)$. So the whole complexity is

$$O(L(d^3 + n'^3) + n'^3).$$

5.2 Bayesian variational regression

Quantile kriging requires repeated observations to obtain direct observations of the quantile and make the hypothesis of Gaussian errors acceptable. Variational approaches allow us to remove this critical constraint, while setting a more realistic statistical hypothesis on ε . Starting from the decomposition of Eq. 21, $\varepsilon(x)$ is now assumed to follow a Laplace asymmetric distribution [84, 45], implying:

$$p(y|q, \tau, \sigma, x) = \frac{\tau(1-\tau)}{\sigma} \exp\left(-\frac{l_\tau(y-q(x))}{\sigma}\right), \quad (30)$$

with the priors on q and σ that has to be fixed.

Such assumption may be justified by the fact that minimizing the empirical risk associated to the pinball loss is equivalent to maximizing the asymmetric Laplace likelihood, which is given by

$$p(\mathcal{Y}_n|q_\tau, \mathcal{X}_n, \theta) = \prod_{i=1}^n \frac{\tau(1-\tau)}{\sigma} \exp\left(-\frac{l_\tau(y_i - q_\tau(x_i))}{\sigma}\right). \quad (31)$$

Due to the non-linearity of the pinball loss, it is not possible to compute the likelihood (31) as in (25) and to use it for inference. To overcome this problem, [9] used a variational approach with an expectation-propagation (EP) algorithm [51], while [1] used a variational expectation-maximization (EM) algorithm which was found to perform slightly better.

5.2.1 Variational EM algorithm

The EM algorithm was introduced in [24] to compute maximum-likelihood estimates. Since then, it has been widely used in a large variety of fields [48]. Classically, the purpose of the EM algorithm is to find ζ that maximizes $p(\mathcal{Y}_n|\zeta)$ thanks to the introduction of the hidden variables z . Starting from the log-likelihood $\log(p(\mathcal{Y}_n|\zeta))$, thanks to Jensen's inequality, it is possible to show that:

$$\log(p(\mathcal{Y}_n|\zeta)) \geq \mathcal{L}(\tilde{p}, \zeta),$$

where

$$\mathcal{L}(\tilde{p}, \zeta) = \int \tilde{p}(z) \log\left(\frac{p(\mathcal{Y}_n, z|\zeta)}{\tilde{p}(z)}\right) dz.$$

Moreover, it can be shown that

$$\log p(\mathcal{Y}_n|\zeta) = \mathcal{L}(\tilde{p}, \zeta) + KL(\tilde{p}||p),$$

where KL is the Kullback-Leibler divergence:

$$KL(\tilde{p}||p) = - \int \tilde{p}(z) \log\left(\frac{p(z|\mathcal{Y}_n, \zeta)}{\tilde{p}(z)}\right) dz.$$

As presented Figure 5.2.1, the EM algorithm can be viewed as a two-step optimization technique. The lower bound \mathcal{L} is first maximized with respect to \tilde{p} (E-step), then with respect to the likelihood parameter ζ (M-step). Dealing

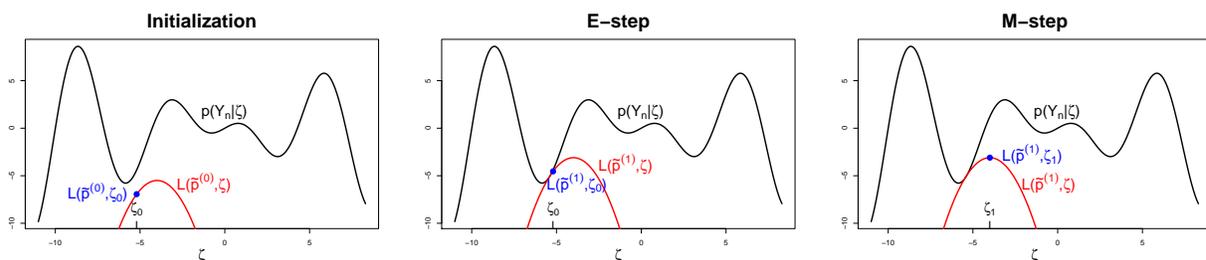


Figure 5: First steps of the variational EM algorithm.

with this traditional formalism implies to know exactly $p(z|\mathcal{Y}_n, \zeta)$, which is not always possible. The variational EM framework has been introduced to bypass this requirement [78]. In this formalism, $p(z|\mathcal{Y}_n, \zeta)$ is assumed to have a particular form. In the literature the factorization form

$$\tilde{p}(z) = \prod_{j=1}^M \tilde{p}_j(z_j)$$

has shown good results [35]. Under this assumption the $\tilde{p}_j^*(z_j)$ $j = 1, \dots, M$ that maximizes $\mathcal{L}(\tilde{p}, \zeta)$ can be written [78] under the form:

$$\tilde{p}_j^*(z_j) = \frac{\exp(\langle \log(p(z|\zeta)) \rangle_{i \neq j})}{\int \exp(\langle \log(p(z|\zeta)) \rangle_{i \neq j}) dz_j}, \quad (32)$$

with

$$\langle \log(p(z|\zeta)) \rangle_{i \neq j} = \int \log p(z|\zeta) \prod_{i \neq j} \tilde{p}_i dz_i.$$

5.2.2 Variational EM applied to quantile regression

Following [1], let us suppose that

$$\begin{aligned} q_\tau &\sim \mathcal{GP}(m(\cdot), K_\theta(\cdot, \cdot)) \\ \sigma &\sim IG(10^{-6}, 10^{-6}), \end{aligned}$$

with IG defining the inverse gamma distribution. Note that contrary to the formalism introduced with QK, here σ is taken as a random variable. Let us introduce an alternative definition of the Laplace distribution [45, 39]:

$$p(y_i|q_\tau, x_i, \sigma, \tau) = \int \mathcal{N}(y_i|\mu_i, \sigma_{y_i}) \exp(-w_i) dw, \quad (33)$$

where $\mu_{y_i} = q_\tau(x_i) + \frac{1-2\tau}{\tau(1-\tau)}\sigma w_i$, $\sigma_{y_i} = \frac{2}{\tau(1-\tau)}\sigma^2 w_i$ and w_i is distributed according to an exponential law of parameter 1.

The distribution of q_τ at a new point x^* is given by averaging the output of all Gaussian models with respect to the posterior $p(q_\tau, \sigma, w|x, Y)$:

$$p(q_\tau(x^*)|\mathcal{D}_n) = \int p(q_\tau(x^*)|q_\tau, \sigma, w, \mathcal{D}_n) p(q_\tau, \sigma, w|\mathcal{D}_n) dq d\sigma dw. \quad (34)$$

The main problem comes from the calculation of the posterior $p(q_\tau, \sigma, w|\mathcal{D}_n) \propto p(\mathcal{Y}_n|q_\tau, \sigma, w, \mathcal{X}_n) p(q_\tau, \sigma, w)$. The likelihood $p(\mathcal{Y}_n|q_\tau, \sigma, w, \mathcal{X}_n)$ is not analytically tractable. To overcome this difficulty, in [1] the authors use $z = (q_\tau = (m, K_\theta), \omega, \sigma)$ as hidden variables and $\zeta = \theta \in \mathbb{R}^{d+1}$ as parameters and the variational factorization approximation

$$p(q_\tau, \sigma, w|\mathcal{D}_n) \approx \tilde{p}(q_\tau, w, \sigma) = \tilde{p}(q_\tau) \tilde{p}(w) \tilde{p}(\sigma). \quad (35)$$

The EM algorithm provides a nice formalism here. Although the goal is to find θ such that $p(\mathcal{Y}_n|\mathcal{X}_n, \theta)$ is maximal, the algorithm estimates the underlying GP (i.e. $p(q_\tau, \sigma, w|\mathcal{D}_n)$) that is able to have a likelihood as large as possible. Then the estimated value $p(q_\tau, \sigma, w|\mathcal{D}_n)$ is plugged into (34) to provide the final quantity of interest.

E-step. We seek to optimize $\mathcal{L}(\tilde{p}(q_\tau, w, \sigma, \theta|\tau))$ with respect to q_τ, w, σ . Thanks to the property (32), it is possible to compute the optimal distribution that is given by

$$\tilde{p}(q_\tau) \sim \mathcal{N}(\mu_\theta, \Sigma_\theta),$$

where

$$\begin{aligned} \mu_\theta &= \frac{1}{2}\Sigma_\theta \left(\tau(1-\tau)\mathbb{E}\left(\frac{1}{\sigma^2}\right) \text{diag}\left(\mathbb{E}\left(\frac{1}{w_i}\right)\right)_{i=1,\dots,n} \mathcal{Y}_n - (1-2\tau)\mathbb{E}\left(\frac{1}{\sigma}\right) \right) \text{ and} \\ \Sigma_\theta &= \left(\frac{\tau(1-\tau)}{2}\mathbb{E}\left(\frac{1}{\sigma^2}\right) \text{diag}\left(\mathbb{E}\left(\frac{1}{w_i}\right)\right)_{i=1,\dots,n} + K_{x,x}^{\theta^{-1}} \right)^{-1}. \end{aligned}$$

The posterior on w_i is a Generalized Inverse Gaussian **GIG**($1/2, \alpha_i, \beta_i$) with :

$$\alpha_i = \left(\frac{(1-2\tau)^2}{2\tau(1-\tau)} + 2 \right)$$

and

$$\beta_i = \frac{\tau(1-\tau)}{2}\mathbb{E}\left(\frac{1}{\sigma^2}\right) \left(y_i^2 - 2y_i\mathbb{E}(q_\tau(x_i)) + \mathbb{E}(q_\tau(x_i)^2) \right).$$

Due to numerical problems, in [1] the authors use the restriction $\tilde{p}(\sigma) = IG(a, b)$. Finding the best a, b is done numerically. Then finding the best a, b is equivalent to maximizing:

$$\begin{aligned} J(a, b) &= (a - N - 10^{-6}) \log(b - \psi(a)) \\ &+ (b - \gamma) \frac{a}{b} - \delta \frac{a(a+1)}{b^2} - a \log(b) + \log \Gamma(a), \end{aligned}$$

with

$$\gamma = -\frac{1-2\tau}{2} \sum_{i=1}^n y_i - \mathbb{E}(q_\tau(x_i))$$

and

$$\delta = \frac{\alpha(1-\tau)}{4} \sum_{i=1}^n \mathbb{E} \left(\frac{1}{w_i} \right) \left(y_i^2 - 2y_i \mathbb{E}(q_\tau(x_i)) + \mathbb{E}(q_\tau(x_i)^2) \right).$$

M-step. The M step consists of maximizing $\mathcal{L}(\tilde{p}(q_\tau, w, \sigma, \theta|\tau)$ with respect to θ . Ignoring terms that do not depend on θ , we obtain the lower bound:

$$\begin{aligned} \tilde{\mathcal{L}}(\theta) &= \int \tilde{p}(q_\tau|\theta) \tilde{p}(w) \tilde{p}(\sigma) \log p(y|q_\tau, w, \sigma) p(q_\tau|\theta) d\sigma dw dq_\tau \\ &\quad - \int \tilde{p}(q_\tau|\theta) \log \tilde{p}(q_\tau|\theta) dq_\tau \\ &= \frac{1}{2} \left(\mu_\theta^T \Sigma_\theta^{-1} \mu_\theta - \log |\mathbb{E}(D^{-1}) + K_{x^*, x}^\theta| \right). \end{aligned} \quad (36)$$

The optimization of $\tilde{\mathcal{L}}$ with respect to θ is done using a numerical optimizer.

Recalling the goal is to compute (34), thanks to (35), we make the approximation:

$$p(q_\tau(x^*)|\mathcal{D}_n) \approx \int p(q_\tau(x^*)|q_\tau, \sigma, \mathcal{D}_n) \tilde{p}(q_\tau) \tilde{p}(\sigma) \tilde{p}(w) dq_\tau dw d\sigma.$$

Finally we obtain

$$q_\tau(x^*) \approx \mathcal{N}(\bar{q}_\tau(x^*), \mathbb{V}_q(x^*)),$$

where

$$\begin{aligned} \bar{q}_\tau(x^*) &= K_{x^*, x}^\theta K_{x, x}^{\theta -1} \mu_\theta \text{ and} \\ \mathbb{V}_q(x^*) &= K_{x^*, x}^\theta - K_{x^*, x}^\theta K_{x, x}^{\theta -1} K_{x^*, x}^{\theta T} + K_{x^*, x}^\theta K_{x, x}^{\theta -1} \Sigma_\theta K_{x, x}^{\theta -1} K_{x^*, x}^{\theta T}. \end{aligned}$$

Finally, as explained in section 5.1.2, the quantile estimator \hat{q}_τ is selected as \bar{q}_τ . Algorithm 6 details the implementation of the VB method.

5.2.3 Computational complexity

E-step. The complexity of this step is in

$$O(n^3).$$

In fact the algorithm computes Σ_θ that implies inverting a matrix of size $n \times n$.

M-step. Optimizing $\tilde{\mathcal{L}}$ with a Newton type algorithm costs $O(n^3 + d^3)$ at each iteration (for details refer to the optimization description of (27)). Assuming the starting point θ_{start} is close to an optimal θ^* , based on the same analysis as in section 4.1.3, the whole complexity is in

$$O(L(d^3 + n^3)).$$

Overall complexity. At each iteration of the EM algorithm, the computation cost is $O(L(d^3 + n^3) + n^3)$. The final complexity is obtained by multiplying by the number of iterations n_{it} of the EM algorithm. Thus, the overall complexity is in

$$O(n_{it}(L(d^3 + n^3) + n^3)).$$

6 Metamodel summary and implementation

In this section we detail our implementation procedure. After providing a summary of the six metamodels in Table 6.1, we present the packages we used and the hyperparameters we chose (which hyperparameters we set and which hyperparameters we optimized). We then describe the procedure we used to optimize the hyperparameters (optimization strategies and evaluation metrics).

6.1 Summary of the models

Table 6.1 lists the analytical expressions of the six metamodels, along with the associated underlying quantity.

Algorithm 6 Bayesian variational regression

```

1: Training
2: Inputs:
    $\mathcal{D}_n, \tau, k_{\theta_0}$ 
3: Initialize:
   Compute the  $n \times n$  matrix  $K_{x,x}^\theta$  and  $K_{x,x}^{\theta^{-1}}$ 
    $\theta = \theta_0$ 
4: for  $t = 1$  to  $n_{iter}$  do
5:   E-step
6:    $\Sigma_\theta = \left( \frac{\tau(1-\tau)}{2} \mathbb{E} \left( \frac{1}{\sigma^2} \right) \text{diag} \left( \mathbb{E} \left( \frac{1}{w_i} \right) \right)_{i=1,\dots,n} + K_{x,x}^{\theta^{-1}} \right)^{-1}$ 
7:    $\mu_\theta = \frac{1}{2} \Sigma_\theta \left( \tau(1-\tau) \mathbb{E} \left( \frac{1}{\sigma^2} \right) \text{diag} \left( \mathbb{E} \left( \frac{1}{w_i} \right) \right)_{i=1,\dots,n} \mathcal{Y}_n - (1-2\tau) \mathbb{E} \left( \frac{1}{\sigma} \right) \right)$ 
8:    $\alpha_i = \left( \frac{(1-2\tau)^2}{2\tau(1-\tau)} + 2 \right)$ 
9:    $\beta_i = \frac{\tau(1-\tau)}{2} \mathbb{E} \left( \frac{1}{\sigma^2} \right) \left( y_i^2 - 2y_i \mathbb{E}(q_\tau(x_i)) + \mathbb{E}(q_\tau(x_i)^2) \right)$ 
10:   $w_i \sim \mathbf{GIG}(0.5, \alpha_i, \beta_i) \quad i = 1, \dots, n$ 
11:   $(a, b) = \arg \max (a - N - 10^{-6}) \log(b - \psi(a)) + (b - \gamma) \frac{a}{b} - \delta \frac{a(a+1)}{b^2} - a \log(b) + \log \Gamma(a)$ 
12:   $\sigma \sim IG(a, b)$ 
13:  M-step
14:   $\theta = \arg \max \frac{1}{2} \left( \mu_\theta^T \Sigma_\theta^{-1} \mu_\theta - \log |\mathbb{E}(D^{-1}) + K_{x^*,x}^\theta| \right)$ 
15: end for
16: Prediction
17: Inputs:
    $\mathcal{X}_{test}, \theta, \mu_\theta$ 
18: for each point in  $x^* \in \mathcal{X}_{test}$  do
19:   compute  $K_{x^*,x}^\theta$ 
20:    $\hat{q}_\tau(x^*) = K_{x^*,x}^\theta K_{x,x}^{\theta^{-1}} \mu_\theta$ 
21: end for

```

6.2 Packages and hyperparameter choices

Each method depends on many parameters that can be tuned to improve performance, for example the choice of the kernel function and the value of its parameters for RK, QK and VB or the penalization factor for RK and NN. Here, to limit the computational burden, we chose to optimize only the most critical ones. When possible, for the other parameters, we applied the arbitrary choices and values made by the authors of the original papers. Most changes were made to improve robustness. Below, we describe our experimental settings, also summarized in Table 2.

Nearest Neighbors. We set $d(\cdot)$ as the Euclidean distance and optimized only the size K of the neighborhood.

Random forest. In this case, the only hyperparameter we optimized was the maximum size of the leaves. Regarding the number of trees, we noticed that the metamodel needs many more trees than are needed for the estimation of the expectation. In some problems, the metamodel needs up to 10,000 trees to stabilize the variance. Thus, in our experiments we set the number of trees at 10,000 in all cases. We set the number of dimensions considered for the split at the default choice $d/3$. The depth of the tree is not constrained and the splitting rule is based on Eq. 7. We used the R package *QuantregForest* [50].

Neural network. Based on [15], we set the number of hidden layers at one and the transfer function as the sigmoid. The optimization algorithm is a Newton method, we set E_K at $1/2^K$ with $K = 1, 2, 5, 10, 15, 20, 25, 30, 35$ and the number of multistarts to optimize the empirical risk at five. We optimized the number of neurons h_l in the hidden unit and the regularization parameter $\lambda \in \mathbb{R}_+$. The metamodel is generated using the R package *qrnn* [15].

Method	Definition of $\hat{q}_\tau(x^*)$	Related quantity	Complexity
KN	$y_{(\lfloor K\tau \rfloor)}(x^*)$	The K -nearest points from x^*	$O(nN_{new}(d + \log(n)))$
RF	$\inf\{y_i : \hat{F}(y_i X = x^*) \geq \tau\}$	$\hat{F}(y_i X = x^*) = \sum_{i=1}^n \bar{\omega}_i(x^*) \mathbf{1}_{\{y_i \leq y\}}$	$O(Npn \log^2(n))$
NN	For a 3 layer NN $g_3(\sum_{j=1}^{J_2} g_2(\sum_{i=1}^{J_1} g_1(\langle w_i^{(h_1)}, x^* \rangle + b_i^{(h_1)}))w_j^{(h_2)} + b_j^{(h_2)})w^{(h_3)} + b^{(h_3)}$	With $w_i^{(h_1)}, w_j^{(h_2)}, b_i^{(h_1)}, b_j^{(h_2)}, w^{(h_3)}, b^{(h_3)}$, $1 \leq i \leq J_1, 1 \leq j \leq J_2$, minimizing $\frac{1}{n} \sum_{t=1}^n l_\tau(y_i - \hat{q}_\tau(x_i)) + \sum_{j,i} \frac{\lambda}{J_j} \ w_i^{(h_j)}\ ^2$	$O(m_s L_{\eta^*} n (Jd)^3)$
RK	$\sum_{i=1}^n \alpha_i k_\theta(x^*, x_i) + b$	With $\alpha = (\alpha_1, \dots, \alpha_n)$ minimizing $\frac{1}{2} \alpha^T K_{x,x} \alpha - \alpha^T \mathbf{y}$ s.t $\frac{\tau - 1}{\lambda n} \leq \alpha_i \leq \frac{\tau}{\lambda n}, \forall 1 \leq i \leq n$ and $\sum_{i=1}^n \alpha_i = 0$ and b the τ -quantile of $(y_i - \sum_{j=1}^n \alpha_j k_\theta(x_i, x_j))_{1 \leq i \leq n}$	$O(n^{7/2} \log(\frac{n}{\epsilon}))$
QK	$K_{x^*,x}^\theta (K_{x,x}^\theta + B)^{-1} \hat{q}_\tau$	Maximizing the likelihood: $p(\mathcal{Q}_{n'} \mathcal{X}_{n'})$ $\hat{q}_\tau(x_i) = q_\tau(x_i) + \varepsilon_i \quad \varepsilon_i \sim \mathcal{N}(0, B_{ii})$	$O(L(d^3 + n^3) + n^3)$
VB	$K_{x^*,x}^\theta K_{x,x}^{\theta - 1} \mu$	Approached solution that maximize: $p(\mathcal{Y}_n \mathcal{X}_n)$ $y_i = q_\tau(x_i) + \varepsilon \quad \varepsilon \sim \text{ALP}(0, \sigma)$	$O(n_{it}(L(d^3 + n^3) + n^3))$

Table 1: Summary of the metamodels, included the definition of the estimators, the associated numerical quantity and the related computation complexity.

Method	Hyperparameters
KN	number of neighbors $K \in \mathbb{N}^*$
RF	node size $ns \in \mathbb{N}^*$
NN	regularization $\lambda \in \mathbb{R}_+$, $J_1 \in \mathbb{N}^*$
RK	regularization $\lambda \in \mathbb{R}_+$, lengthscales $\theta \in \mathbb{R}_+^d$
QK	length scale and variance $\theta \in \mathbb{R}_+^{d+1}$
VB	length scale and variance $\theta \in \mathbb{R}_+^{d+1}$

Table 2: Hyperparameters optimized on our benchmark.

Regression in RKHS. The kernel was set as a Matérn 5/2. We optimized the length scale parameters $\theta \in \mathbb{R}_+^d$ and the regularization hyperparameter $\lambda \in \mathbb{R}_+$. Optimization (20) is done with the quadratic optimizer quadprog [77].

Quantile Kriging. The kernel was set as a Matérn 5/2. The number of repetitions was set according to the total number of observations (see Table 4). We optimized the length scale hyperparameter $\theta \in \mathbb{R}_+^d$ and variance hyperparameter $\in \mathbb{R}_+$. QK is implemented in the R package *DiceKriging* [59].

Variational regression. The kernel was set as a Matérn 5/2, the number of EM iterations at 50. We optimized the length scale hyperparameter $\theta \in \mathbb{R}_+^d$ and variance hyperparameter $\in \mathbb{R}_+$. The implementation is based on the Matlab code provided in [1].

6.3 Tuning the hyperparameters

In the previous section, we defined the hyperparameters we wanted to optimize for each method. In fact, once the type of metamodel is chosen, the quantile estimator is given by a function $\hat{q}_\Theta : \mathcal{X} \rightarrow \mathbb{R}$ where $\Theta \in \mathbb{R}^v$ are called hyperparameters and v is metamodel dependent. Hyperparameter optimization (also known as model selection) is an essential procedure when dealing with non-parametric estimators. Although \hat{q}_Θ may be very efficient on \mathcal{D}_n , the prediction may perform very poorly on an independent dataset \mathcal{D}'_p . The goal is to find the Θ that provides the best possible generalized estimator. In the following, we present the validation metric used to optimize the hyperparameter values and detail the hyperparameters optimization procedure associated with each method.

6.3.1 Metrics

In the standard conditional expectation estimation, the validation and performance metrics are both based on $\|\hat{m}_\Theta - m\|_{L^2}$, where \hat{m}_Θ is the estimator and m the targeted value. With the quantile estimation procedure the two metrics are no longer the same. The goal is to find \hat{q}_Θ such that $\|\hat{q}_\Theta - q\|_{L^2}$ is as small as possible. However, q is unobserved so the validation metric cannot be based on the L^2 norm. Here we present two metrics able to measure the generalization capacity of a quantile metamodel. Bayesian metamodels (QK and VB) have their own validation metric, this is the likelihood function that can be maximized with respect to Θ . For quantile kriging, we use (27) while in the variational approach we use (36). The optimal hyperparameters are then:

$$\Theta_{mv}^* = \arg \max_{\Theta} p(\mathcal{Y}_n | \mathcal{X}_n, \Theta). \quad (37)$$

The second metric available for all metamodels is k -fold cross-validation associated with the pinball loss. The metric can be computed as follows. First, the data are split into k parts, then the model is trained on \mathcal{D}_{-j} , the training set without the j -th part. The performance is evaluated on the remaining part \mathcal{D}_j . As the quantile minimizes the pinball loss (on \mathcal{D}_j), the evaluation metric is

$$E_{cv}(\hat{q}_\tau^\Theta) = \frac{1}{k} \sum_{j=1}^k \frac{1}{n'_j} \sum_{i=1}^{n'_j} l_\tau((y_i - \hat{q}_\tau^\Theta(x_i))), \quad (38)$$

where n'_j is the number of observations in each fold. The optimal cross-validation hyperparameters are then:

$$\Theta_{cv}^* \in \arg \min_{\Theta} E_{cv}(\hat{q}_\tau^\Theta).$$

In our experiments, we chose $k = 5$ to limit the computational cost. However, we observed empirically that choosing a larger k did not substantially modify the performances of the metamodels. Although cross-validation is available for QK and VB, we chose to stay in the spirit of the methods and to only use likelihood to select hyperparameters.

6.3.2 Hyperparameters optimization procedure

Both likelihood functions come with analytical derivatives, enabling the use of gradient-based algorithms. However, since both functions are multi-modal, multi-start techniques are necessary to avoid being trapped in local optima. We ran $n_{start} = 20d$ optimization procedures from different starting points θ_{start} and chose the set of starting points based on a *maximin* Latin hypercube design.

For QK, the BFGS algorithm is used to optimize (27). For VB, two derivative-based optimizers are used alternately for the E and M-steps. Since each step may lead the algorithm toward a local minimum, we chose to apply the multi-start strategy in the entire EM procedure.

Optimization of the cross-validation metric (38) is done under the black box framework, since no structural, derivative or even regularity information is available. Hence, all optimizations are carried out using the branch-and-bound algorithm named Simultaneous Optimistic Optimization (SOO) [53]. SOO is a global optimizer, hence robust to local minima.

We used [76] to parallelize the computations.

6.3.3 Oracle metamodels

Each method presented in this paper is a trade-off between power and the difficulty of finding good hyperparameters. A good method should be powerful (i.e. provide flexible fits) but easy to tune. In order to assess the strengths and

weaknesses of the hyperparameter tuning methods in addition to standard metamodels, we provide what we call the *oracle metamodels* for each problem. Instead of using the cross-validation or likelihood metric, the oracle tunings are directly based on the evaluation metric (42). In a sense, they provide an upper bound on the performance of each method. This allows us to show which metamodels have the potential to tackle the problems and which are intrinsically too rigid or make poor use of information. In addition, this allows us to directly assess the quality of the validation procedure.

7 Benchmark design and experimental setting

Many factors can affect the efficiency of methods to estimate the right quantile. For our benchmark system, we considered four models or test cases to evaluate the performance of the six metamodels. We decided to focus primarily on the dimensionality of the problem, the number of training points available, the quantile order and the pdf value at the targeted quantile for test cases in which the distribution shape (say, from short-tailed to long-tailed) and the distribution spread (i.e. level of heteroscedasticity) can vary significantly. Our two objectives were to:

1. discover if there is a single best method for all factors variations considered or specific choices depending on the configuration at hand, and
2. assess the performance of the quantile regression, and in particular, the configurations for which the current state-of-the-art is insufficient.

A full 3D factorial experimental design was used to analyze the efficiency of the metamodels, the 3 factors being the test case (4 test cases), the number of training points (4 levels) and the quantile order (0.1, 0.5 and 0.9). We used part of this complete design to focus our analyses on the characteristics of the test cases (dimension and pdf distribution as shape or heteroscedasticity).

7.1 Test cases and numerical experiments

Test case 1 is a 1D toy problem on $[-1, 1]$ defined as

$$Y = 5 \sin(8x) + (0.2 + 3x^3)\xi, \quad (39)$$

with $\xi = \eta \mathbb{1}_{\eta < 0} + 7\eta \mathbb{1}_{\eta \geq 0}$ where $\eta \sim \mathcal{N}(0, 1)$.

Test case 2 is a 2D toy problem on $[-5, 5] \times [-3, 3]$ based on the Griewank function [26], defined as

$$Y = \left[\sum_{i=1}^2 \frac{x_i^2}{4000} - \prod_{i=1}^2 \cos\left(\frac{x_i}{\sqrt{i}}\right) + 1 \right] \xi, \quad (40)$$

with $\xi = \eta \mathbb{1}_{\eta < 0} + 5\eta \mathbb{1}_{\eta \geq 0}$ where $\eta \sim \mathcal{N}(0, 1)$.

Test case 3 is a 1D toy problem based on the Michalewicz function [26] on $[0, 4]$, defined as

$$Y = -2 \sin(x) \sin^{30}\left(\frac{x^2}{\pi}\right) - \frac{0.1 \cos(\pi x/10)^3}{|-2 \sin(x) \sin^{30}\left(\frac{x^2}{\pi}\right) + 2|} \xi^2, \quad (41)$$

with $\xi = 3\eta \mathbb{1}_{\eta < 0} + 6\eta \mathbb{1}_{\eta \geq 0}$ where $\eta \sim \mathcal{N}(0, 1)$.

Note that on those three toy problems, the random term ξ is defined such that the resulting distribution of Y would be strongly asymmetric. As ξ is also multiplied by a factor that depends on x , the distribution of Y is also heteroscedastic. The three toy problems are represented in Figure 6.

Test case 4 is based on the agronomical model SUNFLO, a process-based model which was developed to simulate sunflower grain yield (in tons per hectare) as a function of climatic time series, environment (soil and climate), management practices and genetic diversity. The full description of the model is available in [16]. In the regression model we consider \mathcal{X} corresponding to nine macroscopic traits that characterize the sunflower variety. Although SUNFLO is a deterministic model, for each simulation the climatic time series are randomly chosen within a database containing 190 years of weather records, which makes the output stochastic (see also [54] for more details).

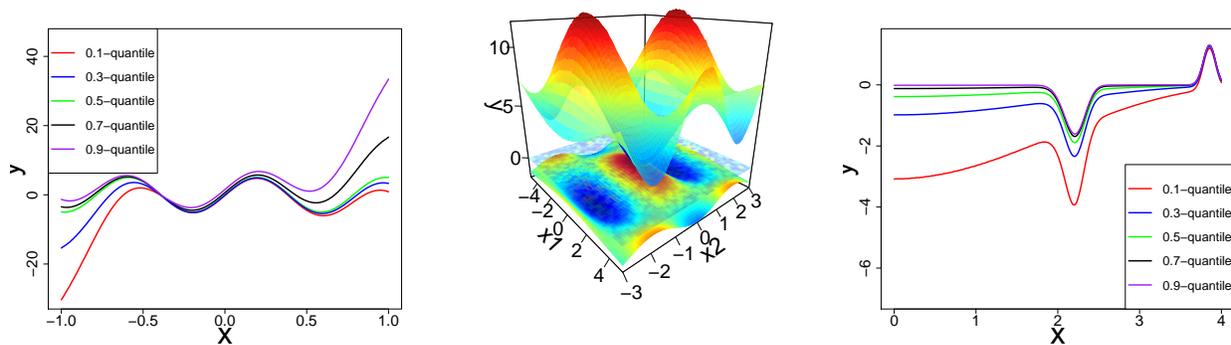


Figure 6: Illustration of the test cases (left: test case 1, center: test case 2, right: test case 3). For test case 1 and 3 the 0.1, 0.3, 0.5, 0.7, 0.9 -quantiles are represented. For test case 2, only the 0.1, 0.5, 0.9-quantiles are represented.

Method		Toy problem 1	Toy problem 2	Toy problem 3	Sunflo
Dimension		1	2	1	9
Heteroscedasticity		very strong	very strong	very strong	weak
Shape variation		very strong	weak	weak	strong
pdf value near the τ -quantile	$\tau = 0.1$	fluctuant	small	globally very small	large
	$\tau = 0.5$	fluctuant	large	small	large
	$\tau = 0.9$	fluctuant	very small	very large	small

Table 3: Summary of the characteristic of the problems.

Numerical experiments. On all problems, we consider four sample sizes. Those sizes depend on the dimension and are empirically chosen so that the smallest size corresponds to the minimal information required by the metamodels to work and the largest size is chosen keeping in mind the potentially high cost of simulators. Besides, our focus is on computer experiments, where data sizes rarely exceed thousands of points. For the 1D problems, the points are generated on a uniform grid. For the 2D and 9D problems, the observations are taken on a latin hypercube design optimized for a *maximin* criterion to ensure space-filling [47]. The same samples are used by all methods except QK, as it requires repetitions. For QK, the number of distinct points and number of repetitions depends on the budget. The different sample sizes are reported in Table 4. Finally, for each sample size and problem, 10 samples are drawn in order to assess robustness with respect to the data.

7.2 Structuration between the questions and the numerical setting

Factors. Three factors are explicit in our benchmark system: the number of training points, the problem dimension and the quantile level. The other factors depend on the characteristics of the problem concerned: shape variation, pdf value at the quantile, level of heteroscedasticity. For all four test cases, we consider three quantile levels: 0.1, 0.5 and 0.9. Note that due to the asymmetry of the problems, learning for the 0.1 and 0.9 quantiles is not equivalent in terms of difficulty. Indeed, when the response is heteroscedastic (a variance/spread depending on \mathcal{X}) and/or when the shape of $\mathbb{P}_X(Y)$ varies in \mathcal{X} , the pdf $\tilde{f}(x, q_\tau)$ may vary in \mathcal{X} . Intuitively, quantiles with large pdf values are easier to learn, as the data points may be closer to them. Figure 7 illustrates this effect. Table 3 summarizes the characteristics of our design concerning the number of training points with respect to the dimension of the test case. To make our results easier to analyze, we divided the problems into groups that allow us to focus on subsets of factors.

Dimension	Data size (no repetitions)				Data size (with repetitions)			
	40	80	160	320	5 (8)	10 (8)	10 (16)	16 (20)
1	40	80	160	320	5 (8)	10 (8)	10 (16)	16 (20)
2	100	200	400	800	10 (10)	20 (10)	25 (16)	40 (20)
9	250	500	1000	2000	25 (10)	50 (10)	100 (10)	100 (20)

Table 4: Data sizes for the different problems. The number in parentheses are the number of repetitions for QK.

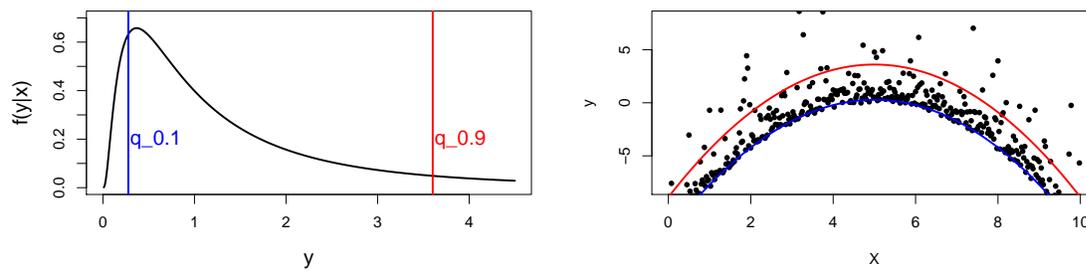


Figure 7: Left: log-normal density function with $\mu = 0$ and $\sigma = 1$. Right: a sample generated by the function $f(x) = \xi - (x - 5)^2/2$, with ξ following the density represented on the left. The 0.9 (resp. the 0.1) quantile is represented in red (resp. in blue). One can notice that more information is available in areas with large pdf (i.e. for the 0.1 quantile) than areas with small pdf.

Table 3 summarizes the characteristics of our design concerning the number of training points with respect to the dimension of the test case. To make our results easier to analyze, we divided the problems into groups that allow us to focus on subsets of factors.

Focus 1: is there a universal winner? To provide a universal ranking of the methods, we use all test cases, training points and quantile levels. As highlighted in Table 3, we created a set of different problems representative of a large number of characteristics that could be met dealing with any quantile regression problem. Note that our benchmark system is slightly biased towards small-dimensional problems, since only a fourth of the cases have a dimension higher than two.

Focus 2: what behavior with respect to the dimension, the number of training points, quantile order and pdf value? To analyze the effects of these factors on the performance of the methods, we combine toy problems 2 and 3 and the SUNFLO model. Toy problem 1 is excluded from the group because the pdf value near all the studied quantiles cannot be classified as large or small.

7.3 Performance evaluation and comparison metrics

Assessing the performance of quantile regression is not an easy task when only limited data are available. Here, since we are considering toy problems, the true quantile values can be computed with precision, so we can evaluate the L^2 error for each emulator:

$$E_{L^2}(\hat{q}) = \sum_{i=1}^{n_{\text{test}}} (\hat{q}(x_i) - q(x_i))^2, \quad (42)$$

where n_{test} is the size of the test set. We chose $n_{\text{test}} = 250$ for the 1D problems and $n_{\text{test}} = 4000$ for the two other. Now, since the problems vary in difficulty and in their response range (Figure 6), $E_{L^2}(\hat{q})$ cannot be aggregated directly over several problems or configurations. To do so, we normalize this error by the error obtained by a constant model (the constant being taken as the quantile of the sample):

$$E_{cq}(\hat{q}) = \sqrt{\frac{E_{L^2}(\hat{q})}{E_{L^2}(CQ)}} \times 100, \quad (43)$$

where CQ stands for constant quantile. Note that this metric closely resembles the $Q2$ criterion.

As an alternative criterion, we consider the ranks of the metamodels based on their $L2$ error. Although ranks do not provide information regarding the range of errors, they are insensitive to scaling issues, which makes aggregation between configurations more sensible. They allow us to assess whereas any method consistently outperforms others, regardless of overall performance.

8 Results

8.1 Focus 1: overall performance and ranks

First, we consider the overall performance and ranks, integrated over all runs. They are shown as boxplots in Figure 8. Based on these ranks (Figure 8, left), VB appears to be the best solution since it is ranked either 1st or 2nd in 50% of the problems. QK is the worst since its median rank is equal to five, which suggests that repeating observations is too much of a disadvantage in the framework considered here. However, all boxplots range between 1 and 6, indicating

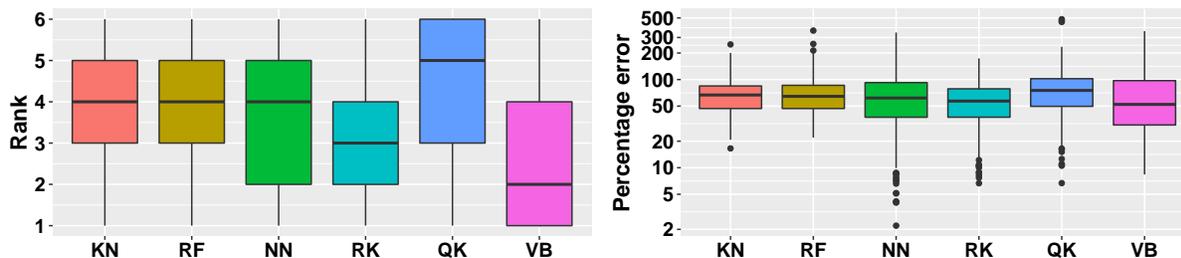


Figure 8: Boxplots of ranks and E_{CQ} error over the entire benchmark. Note that for clarity, the right boxplots do not contain the error of the median for the toy problem 2.

that no method is outperformed by another on all problems. This finding is reinforced by the performance boxplots (Figure 8, right), where all median performances are similar (VB being the best and QK the worst), and the variance is very large. Indeed, the errors range from 2% (of the error achieved by a constant metamodel) to 500%, all methods experiencing cases with more than 100% error (i.e. situations where they are worse than the constant metamodel).

8.2 Focus 2: dimension, number of training points, quantile order and pdf value

8.2.1 Performance according to the constant quantile

In this section, we analyze the performance of the methods with respect to the pdf value and the number of points. *Sample size:* Figure 9 shows the performances of the methods grouped according to the size of the sample. As expected, the performances increase with the size of the sample. For size 1 ($n \approx 50d$), the distribution of E_{CQ} of all the metamodels is almost centered around 100%, implying that these correspond to limit cases for quantile regression since the metamodels do not outperform the constant metamodel (although in some cases the error is as small as 40%). For size 4 ($n \approx 300d$), the median performance is roughly 50% (twice as accurate as the constant metamodel). BV, RK and especially NN experience situations with very accurate models. However, all the methods also experience bad performances (error greater than 100%) in the large sample regime. Unfortunately, from Figure 9 we can conclude that no method is sufficiently robust in all cases.

Pdf value. Figure 10 groups performance with respect to sample size (either small, i.e. level 1 and 2 or large, i.e. level 3 and 4) and pdf value (according to Table 3). According to the Figure 10, the performance depends to a great extent on the pdf value in the neighborhood of the targeted quantile. With a small n , the pdf value has no significant impact on the median of the performance (except for VB) but it does have an impact on the lower bound of the error. With large samples, both the median and the lower bound of the error depend on the pdf value. Metamodels may be very good when the pdf is large, for example 20 times better than the constant metamodel for NN whereas the error appears to have a lower bound when the pdf is small even with large n . In addition, for a problem with a small n and a large pdf, the performance is similar to the performance for problems with a large n and a small pdf (Figure 10, see the two columns in the center).

8.2.2 Rank

Pdf value. Figure 11 shows clearly that when the pdf is large, VB is the best model while when the pdf is small, VB is less good than RN, RK and KNN. This observation is supported by Figure 10 which reveals a strong contrast between the performance of the VB method. QK is poor in both cases, whereas RK performs comparatively better with small pdf.

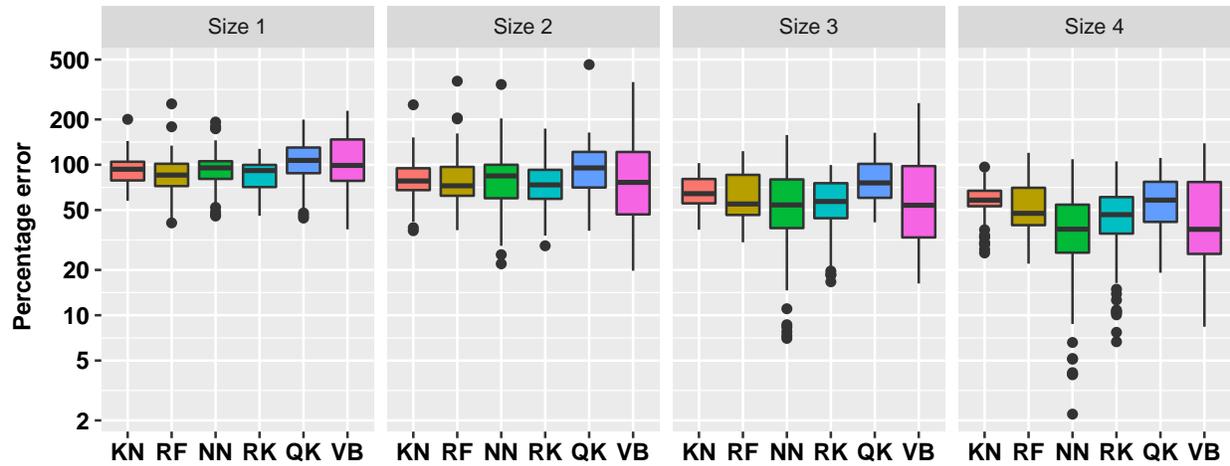


Figure 9: Error according to the sample size

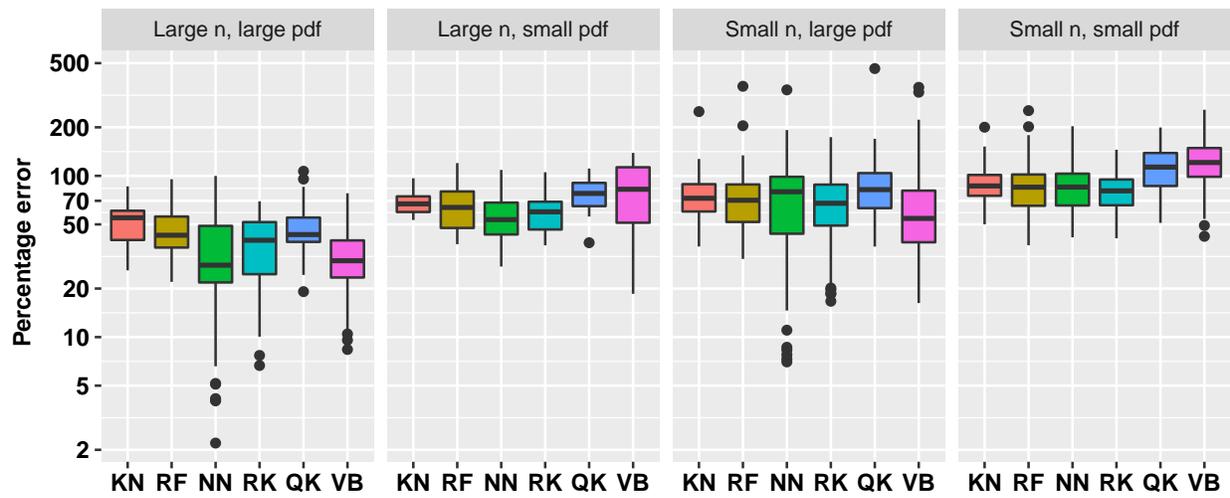


Figure 10: Error according to the size of the training set and the pdf value.

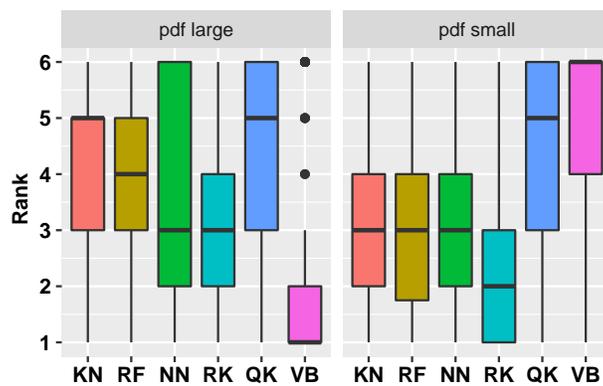


Figure 11: Rank according to the pdf value

Sample size. Figure 12 shows that the number of points has a major impact on the ranking of some methods. The ranking of QK and VB is relatively insensitive to the size of the sample. The other methods are less discriminating

when the sample size is small than when it is large. With a small sample KN, RF, NN and RK are comparable, whereas when sample size increases, NN and RK clearly outperform KN and RF. For the largest size, NN is slightly better than VB.

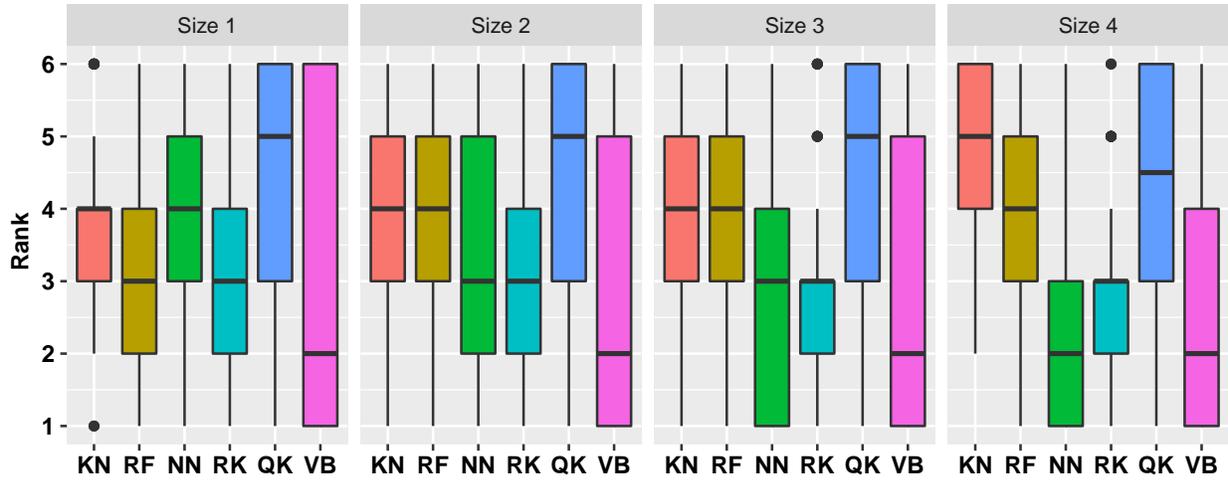


Figure 12: Rank according to the size of the sample

Dimension. Figure 13 groups performance based on dimension. The first contrast is the permutation between RK and NN. With a small dimension, RK is better than NN but the performance of NN improves with an increase in dimension. With small dimensions, RF and KN are comparable, but with high dimensions, RF outperforms KN.

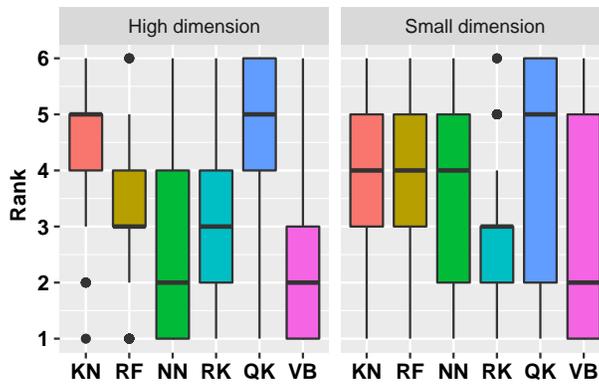


Figure 13: Rank according to the dimension

High dimension, small pdf. Figure 14 shows an extreme case in which the pdf is small but the dimension is high. With a small n , the best method is clearly RF followed by RK and KN. VB and QK are not well ranked. With a larger n , as mentioned above, NN and VB are better but with large variance, while RF, RK and KNN rank less well.

8.2.3 Summary

Figure 15 is a summary of our results and provide additional details to those shown in Figure 1. The main problem is divided into 12 leaves, defined by factor levels. For each leaf, the methods are ranked from one to six according to the mean of their rank, and the expected E_{cq} is provided in the center of each square.

Figure 15 clearly shows that VB is the best method for a subset of problems, but should be avoided in certain cases. RK and NN are reasonable choices overall, although both ranked 5th and 6th on two leaves. RF appears here to be less efficient, which is to be expected considering that its known area of best performance is high dimension / large sample. KN and QK again appear as the worst solutions, except QK when it estimates the median.

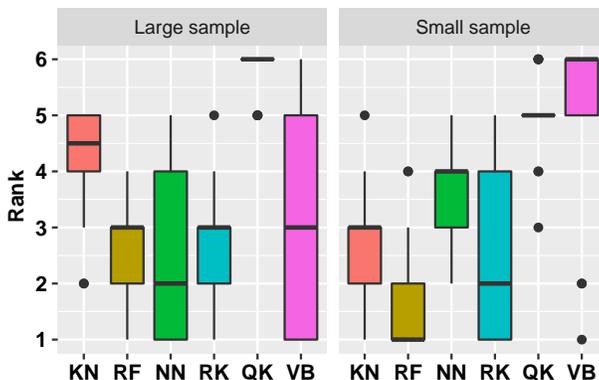


Figure 14: Rank associated to the case where little information is available, *i.e.* high dimension and small pdf

9 Discussion and perspective

9.1 Effect of hyperparameter tuning

In the following we define

$$\Delta E(\hat{q}_\tau^\ominus) = E_{cq}(\hat{q}_\tau^\ominus) - E_{cq}(\hat{q}_\tau^{\ominus*}),$$

the performance gap between the regular metamodel and its oracle performance (the loss in performance between actual hyperparameter tuning and oracle tuning). Figure 16 gives the average values of ΔE aggregated respectively over all problems and only aggregated over the problems with a large pdf, and considering the effect of dimension and sample size. In high dimension, the easiest methods to tune are KN, NN, and RF. In our study, KN and RF have a single hyperparameter to tune regardless of the dimension, and NN has two. This is clearly an advantage in terms of robustness in high dimension. The other methods are kernel-based and require the tuning of at least $d + 1$ hyperparameters. This consistently affects RK and QK, but affects VB only in the case of small pdf, while it is the most stable method in the case of large pdf.

With a small dimension, all the methods have roughly the same number of hyperparameters. The most noticeable change compared to the case of a high dimension is the good performance of RK, while NN becomes comparatively the most difficult method to train.

9.2 On the methods behavior

Statistical order methods. It is clear from Figure 13 that KN performs poorly in a high dimension. This may be due to the irrelevance of the Euclidean distance when there is a significant increase in dimension. RF clearly outperforms KN in this situation, as it is able to produce better neighborhoods than the Euclidean distance. Overall, (compared with the other methods), RF performance increases with dimension. This may be due to the fact that it has fewer hyperparameters to tune.

Functional methods. Figure 12 shows that NN works poorly in a small sample setting, but is one of the best methods when the sample is large. This result reflects the high flexibility of NN. Too much flexibility leads to overfitting when the sample is small. In contrast, when the number of points is large, NNs are able to fit the data very well (e.g. Figure 9, Size 4). According to Figures 9 and 12, RK is a robust method. Its robustness in both small and large data settings can be attributed in part to the selected kernel. If the selected kernel is sufficiently smooth (here continuous and derivable), the resulting metamodel cannot produce instable results. However it seems (Figure 9, Size 3 and 4) that this lack of flexibility may affect the performance with an increase in the size of the data set. In this case, more flexible methods like NN may outperform RK. The contrast between RK and NN shown in Figure 13 can be explained by the level of difficulty associated with each method involved in finding good hyperparameters (as explained above).

Bayesian models. Figure 15 shows that the QK method under-performs on extreme quantiles but is better on the median. One possible explanation is the erroneous assumption in Eq. (23) that the noise is centered, which is more critical for extreme quantiles. The generally poor performance of QK may be due to the fact that, in our case, using

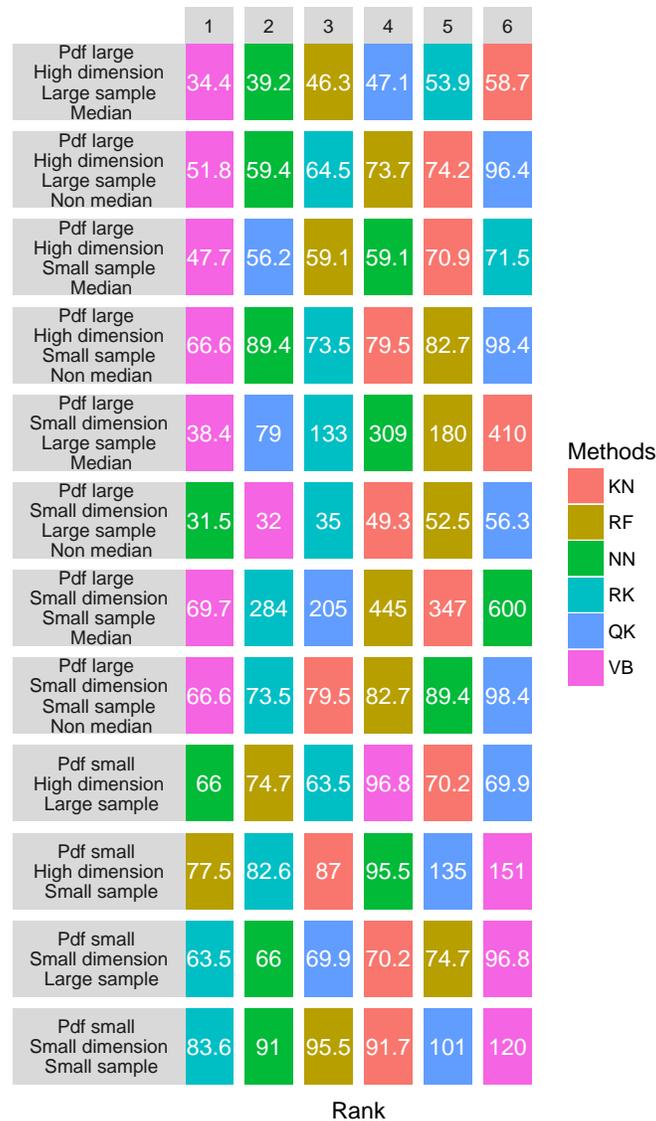


Figure 15: Summary of the rank and the error of the methods. The methods are ranked with respect to the expectation of their rank on the associated class of problems. The expectation of E_{CQ} is provided in white at the center of the cell.

repetitions in the experimental design is sub-optimal. In addition, the increasingly bad performance of QK with an increase in dimension (Figure 9) is likely a consequence of the fact that empty areas become larger in high dimensions.

VB is one of the best methods presented in our paper. Figures 11 and 16 show that VB is also the most dependent on the pdf value. When the pdf is large, it is the best method whereas when the pdf is small, it may be the worst. The explanation lies in the philosophy of the model. In the case of RN and RK, the model complexity (i.e. smoothness) is almost entirely related to the regularity of parameter λ that is selected by cross-validation. Hence, the model cannot excessively overfit and cannot perform very poorly. With Bayesian methods, the regularization is included in the model hypothesis: in our setting, the quantile is assumed to be a Gaussian process with covariance function $k_\theta(\cdot, \cdot)$, so θ performs the regularization. We observed that if the local quantity of information (roughly the product of the number of point times of the pdf value in the neighborhood of the quantile) is too small, the metamodel tends to interpolate the available data. When sufficient information is available, the optimization of the marginal likelihood provides θ values that allow a good trade-off between flexibility and smoothness. This is likely the reason why VB is easily beaten by RN and RK when the pdf is small.

	1	2	3	4	5	6		1	2	3	4	5	6	
High dimension Large sample	2.9	3.2	4.8	9.6	11	13		0.73	1.5	3.1	5.3	6.4	11	Methods ■ KN ■ RF ■ NN ■ RK ■ QK ■ VB
High dimension Small sample	7.2	9.6	12	19	23	31		4.6	7.1	7.9	12	18	23	
Small dimension Large sample	9.5	9.5	9.9	10	13	15		3.7	8	8.4	9.5	9.6	14	
Small dimension Small sample	13	15	16	22	29	30		14	14	16	18	30	32	

Figure 16: Average ΔE aggregated over all problems (left) and over the problems with a large pdf only (right), arranged in increasing order. For each method the rank is provide on the top of each figure.

9.3 Varying shape and heteroscedasticity.

In the previous section, the pdf was assumed to be roughly constant in \mathcal{X} in a neighborhood of the quantile. However, if the shape of $\mathbb{P}_X(Y)$ or the variance of Y (heteroscedasticity) vary w.r.t. \mathcal{X} , then $\hat{f}(x, q_\tau)$ may vary in \mathcal{X} .

Figure 17 illustrates the ability of RF, RK and VB to estimate quantiles of a distribution with a strongly varying shape (toy problem 1, with heavy left-tail for small x to heavy right-tail for large x). In this problem, as depicted in Figure 17 (top row), the quantiles are not perfectly estimated but the metamodels provide good indications about the shape of the true distribution. However, as can be seen in Figure 17 (bottom row), the methods can present strong instabilities. Here, for a sample virtually indistinguishable from the one leading to accurate estimates, the median estimates largely overestimate the true values for large x values. Such instabilities can be partly imputed to the difficulty of the task. However, this is also because no method is actually designed to deal with strongly varying pdf, as we explain below.

An "ideal" method would be almost interpolant for a very large pdf but only loosely fit the data when the pdf is small. However, most of the methods presented here rely on a single hyperparameter to tune the trade-off between data fitting and generalization: the number of neighbors for KN, minimum leaf size for RF and the penalization factor for NN and RK. As a result, the selected hyperparameters are the ones that are best on average. Theoretically, this is not the case for the Bayesian approaches: QK accounts for it *via* the error variance σ_i^2 computed by bootstrap, and the weights w_i (Eq. 33) allow VB to attribute different confidence levels" to the observations. However in practice, both methods fail to tune the values accurately, as we illustrate below. Figure 18 shows the three quantiles of toy problem 3 and their corresponding RF, RK and VB estimates. For $\tau = 0.1$ in particular, the pdf ranges from very small (x close to 0) to very large (x close to 4). Here, RF and RK use a trade-off that globally captures the trend of the quantile, but cannot capture the small hill in the case of large x . Inversely, VB perfectly fits this region but dramatic overfitting occurs on the rest of the domain.

Finally, Figure 19 illustrates that this is not an issue of hyperparameter tuning. For each method, we show the oracle estimate, a tuning that tends to underfit and another that tends to overfit. One can see that no tuning is entirely satisfactory, since capturing the region with high pdf leads to overfitting on the rest of the domain and vice-versa.

We believe that further research is necessary to obtain estimators that intrinsically account for strong heteroscedasticity and varying shape. One possible direction is the use of stacking, in the spirit of [66]. Under the stacking framework the final estimator could be

$$\hat{q}(x) = \sum_{i=1}^N g_i(x) \hat{q}_{\theta_i}(x),$$

where $\{\hat{q}_{\theta_i}\}_{1 \leq i \leq N}$ is a set of metamodel and $\{g_i(x)\}_{1 \leq i \leq N}$ is a set of weight functions. Choosing $\{\hat{q}_{\theta_i}\}_{1 \leq i \leq N}$ such that they correspond to different pdf values might provide more flexible estimates.

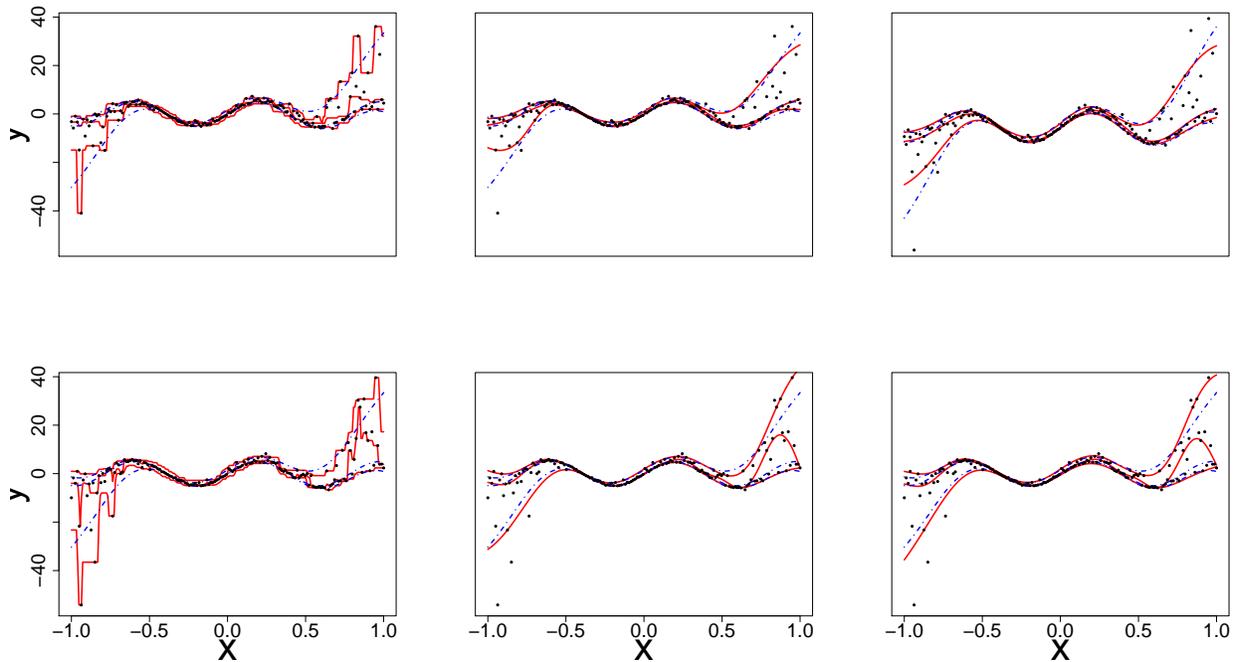


Figure 17: Quantiles estimates using RF (left), RK (middle), VB (right) for two 160-point samples (top and bottom rows, resp.) of the toy problem 1. Dots: observations; plain red lines: metamodels for the 0.1, 0.5, 0.9 quantile estimates; dotted blue lines: actual quantiles.

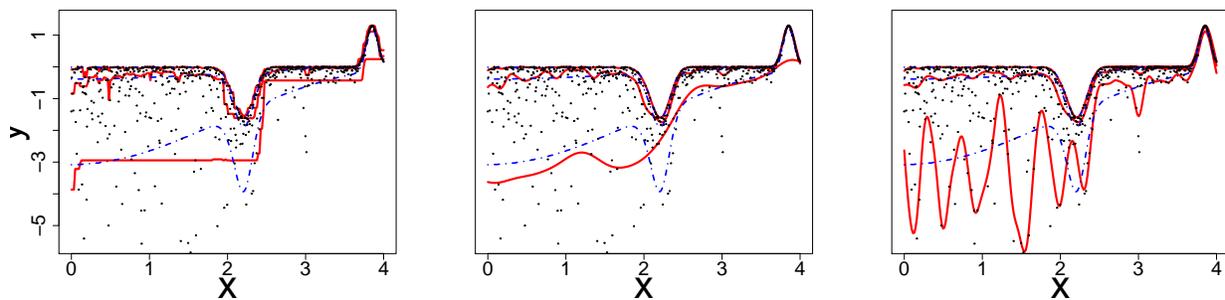


Figure 18: Quantiles estimates using RF (left), RK (middle), VB (right) for a 640-point sample of the toy problem 3. Dots: observations; plain red lines: metamodels for the 0.1, 0.5, 0.9 quantile estimates; dotted blue lines: actual quantiles.

9.4 On the non-crossing of the quantile functions

While the quantile functions (for different quantile levels) may obviously never cross, unfortunately, their estimators may not always satisfy this property. This is a well-known issue against which none of the methods presented here is immune.

The neighborhood approaches first estimate the CDF, then extract the quantiles. If the hyperparameters are the same for all quantiles, crossing is impossible. However in our setting, different neighborhood sizes were used for different quantiles.

With functional analysis approaches, crossing may happen even if each quantile is built with the same hyperparameters. In the literature, authors have produced methods to address this issue. It could be reduced by the introduction of

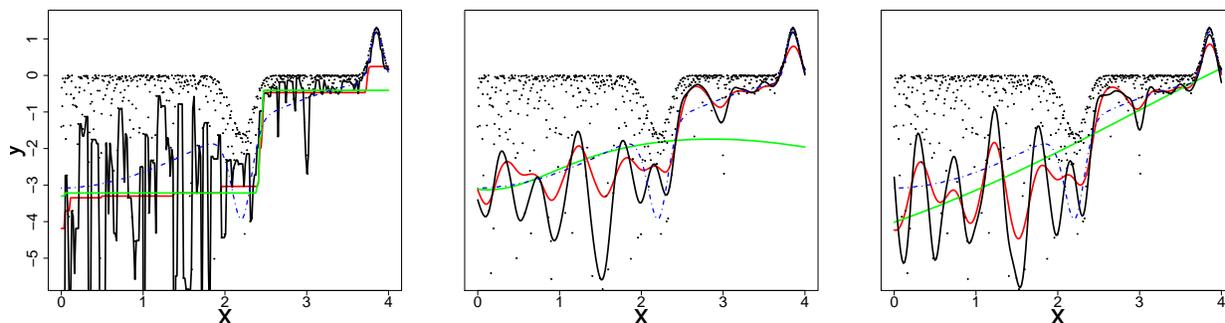


Figure 19: Metamodel responses for toy problem 3 and $\tau = 0.1$ with 640 (left: RF, center: RK, right: VB) for different values of hyperparameters. The true 01-quantile is presented in dotted blue lines. In green and black two extreme metamodels associated to two extreme hyperparameter values, while in red the oracle metamodels are represented.

additional constraints in the model [75] or by the construction of a new model that intrinsically produces non-crossing curves [60]. However in both cases, the dimension of the optimization problem then increases significantly.

Finally the stochastic process approaches estimate each quantile in independent Gaussian processes, so crossing may occur. To our knowledge, this issue has not yet been addressed. Another approach available for all the methods presented here is related to the rearrangement of curves or isotonic regression [5, 2]. The idea is to perform many quantile regressions with a large number of different values of τ or a large set of bootstrapped versions of \mathcal{D}_n and then to rearrange the curves, thereby obtaining the whole distribution and then extracting the quantiles that by definition do not cross.

In theory, adding non-crossing constraints and predicting several quantiles simultaneously could improve the quality of the estimates (in particular as it might add some robustness). However, in practice, it also amounts to making the model more rigid (i.e. a single regularization hyperparameter for all quantiles), and preliminary experiments showed no gain in accuracy compared to independent predictions, despite a considerably higher computational cost. Hence, multi-quantile predictors were not considered in our study.

9.4.1 Concluding comments

In our benchmark, we generally followed the approaches as presented by their authors. However, most of them could be improved. The optimization scheme of NN is the computational bottleneck of the method, which makes it the most expensive method in our benchmark system. One possible improvement would be using the BFGS algorithm [41] or the ADAM algorithm [37] to perform the optimization. A faster scheme would allow more restarts, and hence improve robustness.

Another improvement concerns the splitting criterion (7) of RF, which is not designed for the quantile but for the expectation. This could lead to poor estimates for problems where quantiles are weakly correlated with expectations. Defining an appropriate splitting criterion could significantly improve the performance of this method.

In our experiments, QK used a predefined number of sampling points that were heuristically defined as a trade-off between space-filling and pointwise quantile estimation accuracy. The performance of QK could be significantly improved by optimally tuning the ratio between the number of points and repetitions, in the spirit of [7].

Finally, in practice, finding the best hyperparameters was the most difficult part of the proposed benchmark system. While this aspect is often toned down by authors, we believe is a key practical aspect that remains a challenging problem in quantile regression.

Acknowledgments

This work is part of a Ph.D. of L. Torossian funded by INRA and Région Occitanie. The authors would like to thank Edouard Pauwels for the discussions on the topic of convex and non-convex optimization and Pierre Casadebaig for its suggestions and advices concerning the crop model sunflo.

References

- [1] Sachintha Abeywardana and Fabio Ramos. Variational inference for nonparametric bayesian quantile regression. In *AAAI*, pages 1686–1692, 2015.
- [2] Jason Abrevaya. Isotonic quantile regression: asymptotics and bootstrap. *Sankhyā: The Indian Journal of Statistics*, pages 187–199, 2005.
- [3] Philippe Artzner, Freddy Delbaen, Jean-Marc Eber, and David Heath. Coherent measures of risk. *Mathematical finance*, 9(3):203–228, 1999.
- [4] Sunil Arya, David M Mount, Nathan S Netanyahu, Ruth Silverman, and Angela Y Wu. An optimal algorithm for approximate nearest neighbor searching fixed dimensions. *Journal of the ACM (JACM)*, 45(6):891–923, 1998.
- [5] Alexandre Belloni, Victor Chernozhukov, Denis Chetverikov, and Iván Fernández-Val. Conditional quantile processes based on series or many regressors. *arXiv preprint arXiv:1105.6154*, 2011.
- [6] Pallab K Bhattacharya and Ashis K Gangopadhyay. Kernel and nearest-neighbor estimation of a conditional quantile. *The Annals of Statistics*, pages 1400–1415, 1990.
- [7] Mickael Binois, Robert B Gramacy, and Mike Ludkovski. Practical heteroskedastic gaussian process modeling for large simulation experiments. *Journal of Computational and Graphical Statistics*, (just-accepted):1–41, 2018.
- [8] Christopher M Bishop. *Neural networks for pattern recognition*. Oxford university press, 1995.
- [9] Alexis Boukouvalas, Remi Barillec, and Dan Cornford. Gaussian process quantile regression using expectation propagation. *arXiv preprint arXiv:1206.6391*, 2012.
- [10] George EP Box and George C Tiao. *Bayesian inference in statistical analysis*, volume 40. John Wiley & Sons, 2011.
- [11] Stephen Boyd and Lieven Vandenberghe. *Convex optimization*. Cambridge university press, 2004.
- [12] Mikio L Braun. Accurate error bounds for the eigenvalues of the kernel matrix. *Journal of Machine Learning Research*, 7(Nov):2303–2328, 2006.
- [13] Leo Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
- [14] Leo Breiman. *Classification and regression trees*. Routledge, 2017.
- [15] Alex J Cannon. Quantile regression neural networks: Implementation in r and application to precipitation down-scaling. *Computers & geosciences*, 37(9):1277–1284, 2011.
- [16] Pierre Casadebaig, Lydie Guilioni, Jérémie Lecoeur, Angélique Christophe, Luc Champolivier, and Philippe Debaeke. Sunflo, a model to simulate genotype-specific performance of the sunflower crop in contrasting environments. *Agricultural and forest meteorology*, 151(2):163–178, 2011.
- [17] Marco Cavazzuti. Design of experiments. In *Optimization Methods*, pages 13–42. Springer, 2013.
- [18] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):27, 2011.
- [19] Ching-Hua Chang and Chung-Wei Ha. On eigenvalues of differentiable positive definite kernels. *Integral Equations and Operator Theory*, 33(1):1–7, 1999.
- [20] Andreas Christmann and Ingo Steinwart. Consistency of kernel-based quantile regression. *Applied Stochastic Models in Business and Industry*, 24(2):171–183, 2008.
- [21] Andreas Christmann and Ingo Steinwart. How svms can estimate quantiles and the median. In *Advances in neural information processing systems*, pages 305–312, 2008.
- [22] Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine learning*, 20(3):273–297, 1995.
- [23] Noel Cressie. The origins of kriging. *Mathematical geology*, 22(3):239–252, 1990.
- [24] Arthur P Dempster, Nan M Laird, and Donald B Rubin. Maximum likelihood from incomplete data via the em algorithm. *Journal of the royal statistical society. Series B (methodological)*, pages 1–38, 1977.
- [25] Ramón Díaz-Uriarte and Sara Alvarez De Andres. Gene selection and classification of microarray data using random forest. *BMC bioinformatics*, 7(1):3, 2006.
- [26] L C W Dixon and G P Szego. The global optimization problem. an introduction. *Toward global optimization*, 2:1–15, 1978.
- [27] Harris Drucker, Christopher JC Burges, Linda Kaufman, Alex J Smola, and Vladimir Vapnik. Support vector regression machines. In *Advances in neural information processing systems*, pages 155–161, 1997.

- [28] Sahibsingh A Dudani. The distance-weighted k-nearest-neighbor rule. *IEEE Transactions on Systems, Man, and Cybernetics*, (4):325–327, 1976.
- [29] Dani Gamerman and Hedibert F Lopes. *Markov chain Monte Carlo: stochastic simulation for Bayesian inference*. Chapman and Hall/CRC, 2006.
- [30] Vincent Garcia, Eric Debreuve, and Michel Barlaud. Fast k nearest neighbor search using gpu. In *Computer Vision and Pattern Recognition Workshops, 2008. CVPRW'08. IEEE Computer Society Conference on*, pages 1–6. IEEE, 2008.
- [31] Walter R Gilks, Sylvia Richardson, and David Spiegelhalter. *Markov chain Monte Carlo in practice*. CRC press, 1995.
- [32] Xuming He. Quantile curves without crossing. *The American Statistician*, 51(2):186–192, 1997.
- [33] Marti A. Hearst, Susan T Dumais, Edgar Osuna, John Platt, and Bernhard Scholkopf. Support vector machines. *IEEE Intelligent Systems and their applications*, 13(4):18–28, 1998.
- [34] Hemant Ishwaran. The effect of splitting on random forests. *Machine Learning*, 99(1):75–118, 2015.
- [35] Michael I Jordan, Zoubin Ghahramani, Tommi S Jaakkola, and Lawrence K Saul. An introduction to variational methods for graphical models. *Machine learning*, 37(2):183–233, 1999.
- [36] Kristian Kersting, Christian Plagemann, Patrick Pfaff, and Wolfram Burgard. Most likely heteroscedastic gaussian process regression. In *Proceedings of the 24th international conference on Machine learning*, pages 393–400. ACM, 2007.
- [37] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [38] Roger Koenker and Gilbert Bassett Jr. Regression quantiles. *Econometrica: journal of the Econometric Society*, pages 33–50, 1978.
- [39] Samuel Kotz, Tomasz Kozubowski, and Krzysztof Podgorski. *The Laplace distribution and generalizations: a revisit with applications to communications, economics, engineering, and finance*. Springer Science & Business Media, 2012.
- [40] Miguel Lázaro-Gredilla and Michalis Titsias. Variational heteroscedastic gaussian process regression. 2011.
- [41] Adrian S Lewis and Michael L Overton. Nonsmooth optimization via bfgs. *Submitted to SIAM J. Optimiz*, pages 1–35, 2009.
- [42] Youjuan Li, Yufeng Liu, and Ji Zhu. Quantile regression in reproducing kernel hilbert spaces. *Journal of the American Statistical Association*, 102(477):255–268, 2007.
- [43] Nikolas List and Hans Ulrich Simon. Svm-optimization and steepest-descent line search. In *Proceedings of the 22nd Annual Conference on Computational Learning Theory*, 2009.
- [44] Gilles Louppe. Understanding random forests: From theory to practice. *arXiv preprint arXiv:1407.7502*, 2014.
- [45] Kristian Lum, Alan E Gelfand, et al. Spatial quantile multiple regression using the asymmetric laplace process. *Bayesian Analysis*, 7(2):235–258, 2012.
- [46] Lawrence C Marsh and David R Cormier. *Spline regression models*, volume 137. Sage, 2001.
- [47] Michael D McKay, Richard J Beckman, and William J Conover. Comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics*, 21(2):239–245, 1979.
- [48] Geoffrey McLachlan and Thiriyambakam Krishnan. *The EM algorithm and extensions*, volume 382. John Wiley & Sons, 2007.
- [49] Nicolai Meinshausen. Quantile regression forests. *Journal of Machine Learning Research*, 7(Jun):983–999, 2006.
- [50] Nicolai Meinshausen and Maintainer Nicolai Meinshausen. The quantregforest package. 2007.
- [51] Thomas P Minka. Expectation propagation for approximate bayesian inference. In *Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence*, pages 362–369. Morgan Kaufmann Publishers Inc., 2001.
- [52] Vincent Moutoussamy, Simon Nanty, and Benoît Pauwels. Emulators for stochastic simulation codes. *ESAIM: Proceedings and Surveys*, 48:116–155, 2015.
- [53] Rémi Munos. Optimistic optimization of a deterministic function without the knowledge of its smoothness. In *Advances in neural information processing systems*, pages 783–791, 2011.

- [54] Victor Picheny, Ronan Trépos, and Pierre Casadebaig. Optimization of black-box models with uncertain climatic inputs application to sunflower ideotype design. *PloS one*, 12(5):e0176815, 2017.
- [55] Matthew Plumlee and Rui Tuo. Building accurate emulators for stochastic simulations via quantile kriging. *Technometrics*, 56(4):466–473, 2014.
- [56] Adrian E Raftery. Bayesian model selection in social research. *Sociological methodology*, pages 111–163, 1995.
- [57] Carl Edward Rasmussen and Christopher KI Williams. Gaussian processes for machine learning. 2006. *The MIT Press, Cambridge, MA, USA*, 38:715–719, 2006.
- [58] Roman Rosipal and Leonard J Trejo. Kernel partial least squares regression in reproducing kernel hilbert space. *Journal of machine learning research*, 2(Dec):97–123, 2001.
- [59] Olivier Roustant, David Ginsbourger, and Yves Deville. Dicekriging, diceoptim: Two r packages for the analysis of computer experiments by kriging-based metamodeling and optimization. 2012.
- [60] Maxime Sangnier, Olivier Fercoq, and Florence d’Alché Buc. Joint quantile regression in vector-valued rkhs. In *Advances in Neural Information Processing Systems*, pages 3693–3701, 2016.
- [61] Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation function. *arXiv preprint arXiv:1708.06633*, 2017.
- [62] Bernhard Schölkopf. The kernel trick for distances. In *Advances in neural information processing systems*, pages 301–307, 2001.
- [63] George AF Seber and Alan J Lee. *Linear regression analysis*, volume 329. John Wiley & Sons, 2012.
- [64] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. Taking the human out of the loop: A review of bayesian optimization. *Proceedings of the IEEE*, 104(1):148–175, 2016.
- [65] Jooyong Shim, Changha Hwang, and Kyung Ha Seok. Non-crossing quantile regression via doubly penalized kernel machine. *Computational Statistics*, 24(1):83–94, 2009.
- [66] Joseph Sill, Gábor Takács, Lester Mackey, and David Lin. Feature-weighted linear stacking. *arXiv preprint arXiv:0911.0460*, 2009.
- [67] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958, 2014.
- [68] Michael L Stein. *Interpolation of spatial data: some theory for kriging*. Springer Science & Business Media, 2012.
- [69] Ingo Steinwart and Andreas Christmann. *Support vector machines*. Springer Science & Business Media, 2008.
- [70] Ingo Steinwart, Andreas Christmann, et al. Estimating conditional quantiles with the help of the pinball loss. *Bernoulli*, 17(1):211–225, 2011.
- [71] Charles J Stone. Nearest neighbor estimators of a nonlinear regression function. In *Computer Science and Statistics: 8th Annual Symposium on the Interface*, pages 413–418, 1975.
- [72] Charles J Stone. Consistent nonparametric regression. *The annals of statistics*, pages 595–620, 1977.
- [73] C.B. Storlie and J.C. Helton. Multiple predictor smoothing methods for sensitivity analysis: Description of techniques. *Reliability Engineering and System Safety*, 93:28–54, 2008.
- [74] Matthew A Taddy and Athanasios Kottas. A bayesian nonparametric approach to inference for quantile regression. *Journal of Business & Economic Statistics*, 28(3):357–369, 2010.
- [75] Ichiro Takeuchi, Quoc V Le, Timothy D Sears, and Alexander J Smola. Nonparametric quantile estimation. *Journal of Machine Learning Research*, 7(Jul):1231–1264, 2006.
- [76] Ole Tange. Gnu parallel 2018. 2018.
- [77] Berwin A Turlach and Andreas Weingessel. quadprog: Functions to solve quadratic programming problems. *CRAN-Package quadprog*, 2007.
- [78] Dimitris G Tzikas, Aristidis C Likas, and Nikolaos P Galatsanos. The variational approximation for bayesian inference. *IEEE Signal Processing Magazine*, 25(6):131–146, 2008.
- [79] Vladimir Vapnik. *The nature of statistical learning theory*. Springer science & business media, 2013.
- [80] Ghislain Verdier and Ariane Ferreira. Adaptive mahalanobis distance and k -nearest neighbor rule for fault detection in semiconductor manufacturing. *IEEE Transactions on Semiconductor Manufacturing*, 24(1):59–68, 2011.

- [81] N. Villa-Vialaneix, M. Follador, M. Ratto, and A. Leip. A comparison of eight metamodeling techniques for the simulation of N_2O fluxes and N leaching from corn crops. *Environmental Modelling & Software*, 34:51–66, 2012.
- [82] Ian H Witten, Eibe Frank, Mark A Hall, and Christopher J Pal. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, 2016.
- [83] Keming Yu and Rana A Moyeed. Bayesian quantile regression. *Statistics & Probability Letters*, 54(4):437–447, 2001.
- [84] Keming Yu and Jin Zhang. A three-parameter asymmetric laplace distribution and its extension. *Communications in Statistics-Theory and Methods*, 34(9-10):1867–1879, 2005.
- [85] Peng Zhao and Bin Yu. On model selection consistency of lasso. *Journal of Machine learning research*, 7(Nov):2541–2563, 2006.
- [86] Hui Zou and Trevor Hastie. Regularization and variable selection via the elastic net. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 67(2):301–320, 2005.