



**HAL**  
open science

# Inference for extremal regression with dependent heavy-tailed data

Abdelaati Daouia, Gilles Claude Stupfler, Antoine Usseglio-Carleve

► **To cite this version:**

Abdelaati Daouia, Gilles Claude Stupfler, Antoine Usseglio-Carleve. Inference for extremal regression with dependent heavy-tailed data. *Annals of Statistics*, 2023, 51 (5), pp.2040-2066. 10.1214/23-AOS2320 . hal-04554050v4

**HAL Id: hal-04554050**

**<https://hal.science/hal-04554050v4>**

Submitted on 22 Apr 2024

**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.

WORKING PAPERS

N° 1324

August 2023

“Inference for extremal regression with  
dependent heavy-tailed data”

Abdelaati Daouia, Gilles Stupfler and Antoine Usseglio-Carleve

# INFERENCE FOR EXTREMAL REGRESSION WITH DEPENDENT HEAVY-TAILED DATA

BY ABDELAATI DAOUIA<sup>1,a</sup>, GILLES STUPFLER<sup>2,b</sup> AND ANTOINE USSEGLIO-CARLEVE<sup>3,c</sup>

<sup>1</sup>Toulouse School of Economics, University of Toulouse Capitole, France, <sup>a</sup>abdelaati.daouia@tse-fr.eu

<sup>2</sup>Univ Angers, CNRS, LAREMA, SFR MATHSTIC, F-49000 Angers, France, <sup>b</sup>gilles.stupfler@univ-angers.fr

<sup>3</sup>Avignon Université, LMA UPR 2151, 84000 Avignon, France, <sup>c</sup>antoine.usseglio-carleve@univ-avignon.fr

Nonparametric inference on tail conditional quantiles and their least squares analogs, expectiles, remains limited to i.i.d. data. We develop a fully operational inferential theory for extreme conditional quantiles and expectiles in the challenging framework of  $\alpha$ -mixing, conditional heavy-tailed data whose tail index may vary with covariate values. This requires a dedicated treatment to deal with data sparsity in the far tail of the response, in addition to handling difficulties inherent to mixing, smoothing, and sparsity associated to covariate localization. We prove the pointwise asymptotic normality of our estimators and obtain optimal rates of convergence reminiscent of those found in the i.i.d. regression setting, but which had not been established in the conditional extreme value literature. Our assumptions hold in a wide range of models. We propose full bias and variance reduction procedures, and simple but effective data-based rules for selecting tuning hyperparameters. Our inference strategy is shown to perform well in finite samples and is showcased in applications to stock returns and tornado loss data.

## 1. Introduction.

1.1. *Motivation.* Quantile regression is a well-established statistical tool for the assessment of the impact of a vector of covariates  $\mathbf{X} \in \mathbb{R}^p$  upon a response variable  $Y \in \mathbb{R}$ . It fully describes the conditional distribution of  $Y$  given  $\mathbf{X}$  by considering the conditional quantiles

$$(1) \quad q(\tau|\mathbf{x}) = \arg \min_{\theta \in \mathbb{R}} \mathbb{E}([\varrho_\tau(Y - \theta) - \varrho_\tau(Y)]|\mathbf{X} = \mathbf{x}), \quad \tau \in (0, 1),$$

where  $\varrho_\tau(y) = |\tau - \mathbb{1}_{\{y \leq 0\}}| |y|$  denotes the quantile check function. Regression quantile estimators at the tails nonetheless typically suffer from instability and inconsistency due to data sparseness, especially when the underlying conditional distributions are heavy-tailed. This class of distributions is ubiquitous in, among others, insurance (large losses due to large claims), finance (large drops in stock indices) and natural sciences (large earthquake magnitudes, flood intensity, extreme rainfall). Existing approaches to extremal quantile regression in the heavy-tailed case fall into two main categories: linear quantile regression approaches, such as those of [6], [7] and [25], and at the opposite, nonparametric approaches such as those of [9] and [10]. Yet, tail quantile regression theory typically assumes that the data is independent and identically distributed (i.i.d.). Only [7] provides feasible inference tools for  $\alpha$ -mixing data, by assuming that conditional quantiles  $q(\tau|\mathbf{x})$  are linear in  $\mathbf{x}$ , that the extreme value index is constant, and using self-normalized quantile regression statistics rather

---

*MSC 2010 subject classifications:* Primary 62G32; secondary 62G05, 62G08, 62G15, 62G20, 62G30.

*Keywords and phrases:* Conditional quantiles, conditional expectiles, extreme value analysis, heavy tails, inference, mixing, nonparametric regression.

than simpler asymptotic Gaussian confidence intervals. More broadly, the problem of inference on nonlinear extremal quantile regression remains untouched under serial dependence.

In actuarial and financial risk management, the robustness of quantiles may constitute a weakness as they are not sensitive to the severity of extreme losses. Their failure to satisfy the coherence property (introduced in [1]) is also a serious drawback. A better alternative in these respects is expectile regression [22], which focuses on the conditional expectiles

$$(2) \quad e(\tau|\mathbf{x}) = \arg \min_{\theta \in \mathbb{R}} \mathbb{E}([\eta_\tau(Y - \theta) - \eta_\tau(Y)]|\mathbf{X} = \mathbf{x}), \quad \tau \in (0, 1),$$

where  $\eta_\tau(y) = |\tau - \mathbb{1}_{\{y \leq 0\}}|y^2$  is an asymmetric quadratic loss function. Expectiles extend the mean as quantiles extend the median, and induce the only coherent and elicitable risk measure [29]. As such, they come endowed with a natural backtesting methodology. Besides, the  $\tau$ th expectile is in fact the  $\tau$ th quantile of a transformed distribution function, that is,

$$(3) \quad e(\tau|\mathbf{x}) = \inf\{y \in \mathbb{R} \mid E(y|\mathbf{x}) \geq \tau\} \text{ with } E(y|\mathbf{x}) = \frac{\mathbb{E}[|Y - y|\mathbb{1}_{\{Y \leq y\}}|\mathbf{X} = \mathbf{x}]}{\mathbb{E}[|Y - y||\mathbf{X} = \mathbf{x}]}$$

Regression expectiles thus make an efficient use of the data, since they rely on the distance to all observations and not only on their probability, and they benefit from a transparent financial meaning in terms of their acceptance sets and the gain-loss ratio [3]. These properties and others have motivated the development of extremal expectile regression. The pioneering contribution of [24] is limited to elliptical heavy-tailed distributions, and the nonparametric approach of [17] hinges on the i.i.d. assumption. The approach of [16] can handle time series location-scale models, but is highly sensitive to model misspecification, makes the strong assumption of a constant tail index, and its bootstrap scheme is difficult to calibrate.

*1.2. Contribution and outline of the paper.* We propose a general and fully operational nonparametric inferential theory for conditional tail quantiles and expectiles when the data comes from an  $\alpha$ -mixing dependent sequence  $((\mathbf{X}_t, Y_t))_{t \geq 1}$ . We allow heavy-tailed data for both regression modes, although moment restrictions are inevitably required for conditional extreme expectile inference, since expectiles extend the mean. The conditional tail index is allowed to depend on covariate values. Our unifying argument is that any quantile can be estimated by inverting an estimator of the associated distribution function. This is very beneficial theoretically, as it reduces the problem of investigating the asymptotic normality of extreme conditional quantile and expectile estimators to proving the asymptotic normality of row sums of identically distributed and dependent random variables, while a solution via empirical versions of (1) and (2) would be technically involved (see [13] and [15] in the unconditional setting and under a stronger  $\beta$ -mixing assumption). When the level  $\tau$  is intermediate in the sense that  $\tau = \tau_n \uparrow 1$  slowly enough as the sample size  $n \rightarrow \infty$ , an asymptotically normal estimator is then obtained by inverting a kernel smoother of the associated conditional distribution function. The intermediate level  $\tau_n$  is assumed to satisfy  $nh_n^p(1 - \tau_n) \rightarrow \infty$ , where  $h_n > 0$  is the bandwidth sequence featuring in the kernel estimator. At properly extreme levels  $\tau = \tau'_n$  such that  $nh_n^p(1 - \tau'_n) = O(1)$ , we consider Weissman-type estimators (after [27]) that rely on a conditional tail index estimator based on quantiles or expectiles.

Our asymptotic theory is obtained through a standard ‘‘big blocks separated by small blocks’’ argument, for algebraically fast mixing and under reasonable technical conditions on the distributional behavior of  $\mathbf{X}$  and of  $Y$  given  $\mathbf{X}$ , marginally and through time; in particular, we make conditions on  $(\mathbf{X}_1, \mathbf{X}_{t+1})$  and on the joint extreme value behavior of  $(Y_1, Y_{t+1})$  given  $(\mathbf{X}_1, \mathbf{X}_{t+1})$  that are weaker than typical conditions in the literature. Our framework encompasses, among others, location-scale models with possible temporal misspecification, a general class of nonlinear regression models containing Generalized Linear Models, and

autoregressive models. Surprisingly perhaps, the asymptotic distribution of the intermediate estimators remains the same as in the nonparametric i.i.d. setting [10, 17], unlike in the absence of covariates [13, 15]. We nonetheless improve upon results of the i.i.d. literature by deriving the asymptotic normality of our estimators at faster rates of convergence than had been found so far. Obtaining these rates, reminiscent of the optimal rates found in classical nonparametric statistics under twice differentiability of the function to be estimated, requires developing an innovative and careful approach to the quantification of bias in smoothed extremal regression estimators. In particular, we find that the variation in conditional extremes induces bias that bears a link to the *design bias* in the terminology of [26]; this had not been appreciated in the literature before. From the inferential point of view, the Weissman-type structure makes it possible to come up with very accurate and novel bias-corrected versions and precise approximations to the empirical variance of the estimators. We thus construct asymptotic Gaussian confidence intervals that substantially improve upon the naive solutions existing in the nonparametric extreme value regression literature, in which the problem of accurate Gaussian inference had been mostly ignored. The method is applied here to inference about extreme conditional quantiles and expectiles, but it can handle any kind of Weissman-type estimator: the bias correction methodology revolves upon identifying bias sources due to the extrapolation procedure, while the variance correction relies on a precise evaluation of the correlation between intermediate estimators and tail index estimators. We propose rules-of-thumb for the choices of the bandwidth  $h_n$  and of the sample fraction  $1 - \tau_n$  needed for tail extrapolation, resulting in confidence intervals achieving excellent coverage already for moderately large sample sizes.

The outline of our paper is the following. Sections 2 and 3 focus on nonparametric extremal quantile and expectile regression, respectively. Section 4 investigates examples of regression models where our assumptions are satisfied. Section 5 develops a fully operational inferential methodology, showcased in a simulation study in Section 6 and on real data in Section 7. Our methods and data have been incorporated into the R package `Expectrem`<sup>1</sup>. Further details about our technical conditions, the proofs of all theoretical results and extra finite-sample results are deferred to an online Supplementary Material document [12]. Throughout we denote by  $x_+ = \max(x, 0)$  and  $x_- = \max(-x, 0)$  the positive and negative parts of a real number  $x$ . For a function  $f$  on  $\mathbb{R}^p$ ,  $\nabla f(\mathbf{x})$ ,  $Jf(\mathbf{x})$  and  $Hf(\mathbf{x})$  stand respectively for its gradient vector, Jacobian matrix, and Hessian matrix at the point  $\mathbf{x}$ . For a function  $f = f(\mathbf{x}, \mathbf{y})$  on  $\mathbb{R}^p \times \mathbb{R}^q$ , let  $\nabla_{\mathbf{x}} f$  and  $H_{\mathbf{x}} f$  be its partial gradient vector and Hessian matrix with respect to  $\mathbf{x}$  (*i.e.* the first  $p$  components of its gradient vector and the submatrix made of the first  $p$  rows and columns of its Hessian matrix, respectively). The symbols  $\mathbf{0}_p$  and  $\mathbf{1}_p$  denote vectors in  $\mathbb{R}^p$  with all components equal to 0 and 1, respectively.

## 2. Nonparametric extremal quantile regression.

2.1. *Framework.* Let  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  be a strictly stationary sequence of copies of  $(\mathbf{X}, Y) \in \mathbb{R}^p \times \mathbb{R}$ . Let  $F(\cdot|\mathbf{x})$  denote the distribution function of  $Y$  given  $\mathbf{X} = \mathbf{x}$ , that is,  $F(y|\mathbf{x}) = \mathbb{P}(Y \leq y|\mathbf{X} = \mathbf{x})$ . Assume that  $\mathbf{X}$  has a probability density function (p.d.f.)  $g$  on  $\mathbb{R}^p$  and fix  $\mathbf{x} \in \mathbb{R}^p$  with  $g(\mathbf{x}) > 0$ . Consider the following kernel estimator of  $F(\cdot|\mathbf{x})$ :

$$\widehat{F}_n(y|\mathbf{x}) = \frac{1}{nh_n^p \widehat{g}_n(\mathbf{x})} \sum_{t=1}^n \mathbb{1}_{\{Y_t \leq y\}} K\left(\frac{\mathbf{x} - \mathbf{X}_t}{h_n}\right) \quad \text{with} \quad \widehat{g}_n(\mathbf{x}) = \frac{1}{nh_n^p} \sum_{t=1}^n K\left(\frac{\mathbf{x} - \mathbf{X}_t}{h_n}\right).$$

<sup>1</sup>Available on GitHub at <https://github.com/AntoineUC/Expectrem>

Hereafter  $K$  is a kernel p.d.f. on  $\mathbb{R}^p$  and  $h_n \rightarrow 0$  is a (positive) bandwidth sequence, with  $\hat{g}_n$  being the associated classical Parzen-Rosenblatt estimator of  $g$ . A conditional quantile  $q(\tau|\mathbf{x}) \equiv \inf \{y \in \mathbb{R} \mid F(y|\mathbf{x}) \geq \tau\}$  can then be estimated by its empirical counterpart

$$\hat{q}_n(\tau|\mathbf{x}) = \inf \left\{ y \in \mathbb{R} \mid \hat{F}_n(y|\mathbf{x}) \geq \tau \right\}.$$

The standard roadmap for estimating extreme conditional quantiles at a level  $\tau = \tau_n \rightarrow 1$  as  $n \rightarrow \infty$  is to consider first intermediate (“extreme, but not too much”) levels  $\tau_n$ , for which  $\hat{q}_n(\tau_n|\mathbf{x})$  is a (relatively) consistent estimator of  $q(\tau_n|\mathbf{x})$ . Then, for the estimation of properly extreme quantiles  $q(\tau'_n|\mathbf{x})$  without any restriction on  $\tau'_n$ , one extrapolates such intermediate estimators using the shape of the tail of the underlying conditional distribution.

In order to assess the asymptotic behavior of this procedure, we introduce the following conditional, second-order regularly varying tails assumption about  $Y$  given  $\mathbf{X} = \mathbf{x}$ :

*Condition  $\mathcal{C}_2(\gamma(\mathbf{x}), \rho(\mathbf{x}), A(\cdot|\mathbf{x}))$*  There exist  $\gamma(\mathbf{x}) > 0$ ,  $\rho(\mathbf{x}) \leq 0$  and a positive or negative measurable function  $A(\cdot|\mathbf{x})$  converging to 0 at infinity such that for any  $y > 0$ ,

$$\lim_{s \rightarrow \infty} \frac{1}{A(1/\bar{F}(s|\mathbf{x})|\mathbf{x})} \left( \frac{\bar{F}(sy|\mathbf{x})}{\bar{F}(s|\mathbf{x})} - y^{-1/\gamma(\mathbf{x})} \right) = \begin{cases} y^{-1/\gamma(\mathbf{x})} \frac{y^{\rho(\mathbf{x})/\gamma(\mathbf{x})} - 1}{\gamma(\mathbf{x})\rho(\mathbf{x})} & \text{if } \rho(\mathbf{x}) < 0, \\ y^{-1/\gamma(\mathbf{x})} \frac{\log(y)}{\gamma^2(\mathbf{x})} & \text{if } \rho(\mathbf{x}) = 0. \end{cases}$$

This standard condition, wherein  $A(\cdot|\mathbf{x})$  is regularly varying with index  $\rho(\mathbf{x})$  [see 14, Theorem 2.3.3 p.44], controls the proximity between the extremes of the underlying conditional distribution and those of the ideal Pareto distribution with extreme value index  $\gamma(\mathbf{x})$ . We therefore make the fundamental modeling assumption that

*Condition  $\mathcal{M}$*   $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is a stationary  $\alpha$ -mixing sequence of copies of a random vector  $(\mathbf{X}, Y)$  satisfying the second-order regularly varying tails assumption  $\mathcal{C}_2(\gamma(\mathbf{x}), \rho(\mathbf{x}), A(\cdot|\mathbf{x}))$ .

The  $\alpha$ -mixing (or strong mixing) assumption we shall make is expressed as follows: let, for any two positive integers  $a \leq b$ ,  $\mathcal{F}_a^b = \sigma(\{(\mathbf{X}_j, Y_j), a \leq j \leq b\})$  be the  $\sigma$ -algebra generated by  $\{(\mathbf{X}_j, Y_j), a \leq j \leq b\}$ , and say that  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is  $\alpha$ -mixing if and only if

$$\alpha(n) = \sup_{k \geq 1} \sup_{A \in \mathcal{F}_1^k} \sup_{B \in \mathcal{F}_{k+n}^\infty} |\mathbb{P}(A \cap B) - \mathbb{P}(A)\mathbb{P}(B)| \rightarrow 0 \text{ as } n \rightarrow \infty.$$

More restrictive assumptions that will be employed below to relax certain regularity conditions involve the  $\beta$ -,  $\rho$ -,  $\phi$ - and  $\psi$ -mixing coefficients, see Section A.1 of the Supplementary Material document [12] for definitions and relationships between the different kinds of mixing. We introduce the following condition used to develop a “big blocks separated by small blocks” argument for the theory:

*Condition  $\mathcal{A}(l_n, r_n)$*  There exist sequences  $(l_n)$  and  $(r_n)$  such that  $l_n \rightarrow \infty$ ,  $r_n \rightarrow \infty$ ,  $l_n/r_n \rightarrow 0$ ,  $r_n/n \rightarrow 0$  and  $n\alpha(l_n)/r_n \rightarrow 0$  as  $n \rightarrow \infty$ .

We also require reasonable regularity assumptions on the kernel function  $K$  and the probabilistic behavior of the covariate sequence  $(\mathbf{X}_t)_{t \geq 1}$ . Let  $\|\cdot\|$  denote the Euclidean norm on  $\mathbb{R}^p$  and  $B(\mathbf{x}, r)$  be the open  $\|\cdot\|$ -ball with center  $\mathbf{x}$  and radius  $r > 0$ . Full details for the rationale behind our conditions and their interpretation are given in Section A.2 of the Supplementary Material document [12].

*Condition  $\mathcal{K}$*  The p.d.f.  $K$  is bounded with a support contained in the unit closed  $\|\cdot\|$ -ball.

*Condition  $\mathcal{L}_g$*  The p.d.f.  $g$  satisfies  $g(\mathbf{x}) > 0$  and is Lipschitz continuous at  $\mathbf{x}$ : there exist  $c, r > 0$  such that for any  $\mathbf{x}' \in B(\mathbf{x}, r)$ ,  $|g(\mathbf{x}) - g(\mathbf{x}')| \leq c\|\mathbf{x} - \mathbf{x}'\|$ .

*Condition  $\mathcal{B}_p$*  There exists an integer  $t_0 \geq 1$  such that

$$1 \leq t < t_0 \Rightarrow \lim_{r \rightarrow 0} r^{-p} \mathbb{P}(\mathbf{X}_1 \in B(\mathbf{x}, r), \mathbf{X}_{t+1} \in B(\mathbf{x}, r)) = 0$$

$$\text{and } \limsup_{r \rightarrow 0} \sup_{t \geq t_0} r^{-2p} \mathbb{P}(\mathbf{X}_1 \in B(\mathbf{x}, r), \mathbf{X}_{t+1} \in B(\mathbf{x}, r)) < \infty.$$

Assumptions  $\mathcal{K}$ ,  $\mathcal{L}_g$  and  $\mathcal{B}_p$  (whose first half is trivially true if  $t_0$  can be chosen equal to 1) are in particular imposed to control the asymptotic behavior of  $\widehat{g}_n(\mathbf{x})$ . Assumption  $\mathcal{B}_p$  will be satisfied especially if, for all  $t \geq 1$ , the pair  $(\mathbf{X}_1, \mathbf{X}_{t+1})$  has a p.d.f.  $g_t$  such that  $\sup_{t \geq 1} g_t$  is bounded on  $B(\mathbf{x}, r) \times B(\mathbf{x}, r)$  for some  $r > 0$ . Under condition  $\mathcal{L}_g$ , this local boundedness condition is automatically true if  $(\mathbf{X}_t)$  is  $\beta$ -mixing. When  $p \geq 2$ , the causal and invertible  $\text{AR}(p)$  process does not satisfy this boundedness condition, but does satisfy assumption  $\mathcal{B}_p$ .

To control the variation in conditional extreme value behavior across the covariate space, we make a Lipschitz assumption on the log-conditional survival function at extreme levels.

*Condition  $\mathcal{L}_\omega$*  There exists  $r > 0$  such that

$$\limsup_{y \rightarrow \infty} \sup_{\substack{\mathbf{x}' \in B(\mathbf{x}, r) \\ \mathbf{x}' \neq \mathbf{x}}} \frac{1}{\|\mathbf{x}' - \mathbf{x}\|} \left| \frac{1}{\log(y)} \log \frac{\overline{F}(y|\mathbf{x}')}{\overline{F}(y|\mathbf{x})} \right| < \infty.$$

We also impose an anti-clustering condition that translates into assuming that a joint conditional extreme value of  $(Y_1, Y_{t+1})$  cannot be much more likely than a marginal conditional extreme of  $Y_1$ , uniformly across time and locally uniformly across the covariate space. Let

$$\Omega_h(z|\mathbf{x}) = \sup_{t \geq 1} \sup_{y, y' \geq z} \sup_{\mathbf{x}', \mathbf{x}'' \in B(\mathbf{x}, h)} \frac{\mathbb{P}(Y_1 > y, Y_{t+1} > y' | \mathbf{X}_1 = \mathbf{x}', \mathbf{X}_{t+1} = \mathbf{x}'')}{\sqrt{\overline{F}(y|\mathbf{x}')\overline{F}(y'|\mathbf{x}'')}}.$$

*Condition  $\mathcal{B}_\Omega$*  There exist  $h, z > 0$  such that  $\Omega_h(z|\mathbf{x}) < \infty$ .

This should be considered as a weak requirement compared with the existence of a (conditional) tail copula, as assumed in *e.g.* [13, 15] in the unconditional setting. It is an important ingredient in the quantification of the correlation between tail empirical moments.

These conditions will ensure the pointwise asymptotic normality of our estimators at rates of convergence that have hitherto been standard in the conditional extreme value framework. Achieving better rates of convergence requires, similarly to classical nonparametric estimation, stronger regularity conditions: it is well-known that the optimal rate of convergence  $n^{-2/(p+4)}$  of  $\widehat{g}_n(\mathbf{x})$  to  $g(\mathbf{x})$  is obtained by solving the bias-variance tradeoff if  $K$  is symmetric and  $g$  is twice differentiable at  $\mathbf{x}$ . This motivates the following additional assumptions.

*Condition  $\mathcal{KS}$*  The p.d.f.  $K$  is bounded and symmetric (*i.e.*  $K(\mathbf{u}) = K(-\mathbf{u})$ ) with a support contained in the unit closed  $\|\cdot\|$ -ball.

*Condition  $\mathcal{D}_g$*  The p.d.f.  $g$  satisfies  $g(\mathbf{x}) > 0$ , is continuously differentiable in a neighborhood of  $\mathbf{x}$  and its gradient is Lipschitz continuous at  $\mathbf{x}$ .

*Condition  $\mathcal{D}_\omega$*  For  $y$  large enough, the function  $\overline{F}(y|\cdot)$  is differentiable at  $\mathbf{x}$ , the function  $y \mapsto \nabla_{\mathbf{x}} \log \overline{F}(y|\mathbf{x}) / \log(y)$  has a limit  $\boldsymbol{\mu}(\mathbf{x}) \in \mathbb{R}^p$  as  $y \rightarrow \infty$ , and there exists  $r > 0$  with

$$\limsup_{y \rightarrow \infty} \sup_{\substack{\mathbf{x}' \in B(\mathbf{x}, r) \\ \mathbf{x}' \neq \mathbf{x}}} \frac{1}{\|\mathbf{x}' - \mathbf{x}\|^2} \left| \frac{1}{\log(y)} \log \frac{\overline{F}(y|\mathbf{x}')}{\overline{F}(y|\mathbf{x})} - (\mathbf{x}' - \mathbf{x})^\top \frac{\nabla_{\mathbf{x}} \log \overline{F}(y|\mathbf{x})}{\log(y)} \right| < \infty.$$

The motivation for the assumption that  $\nabla_{\mathbf{x}} \log \overline{F}(y|\mathbf{x}) / \log(y)$  converges as  $y \rightarrow \infty$  is the fact that, in the setup of conditional heavy tails,  $\log \overline{F}(y|\mathbf{x}) / \log(y) = -1/\gamma(\mathbf{x}) + \log L(y|\mathbf{x}) / \log(y)$ , where  $L(\cdot|\mathbf{x})$  is slowly varying. In particular,  $\log L(y|\mathbf{x}) / \log(y) \rightarrow 0$

as  $y \rightarrow \infty$  [see 14, Proposition B.1.9.1 p.36]. The assumption asks for this convergence to hold also when taking the gradient with respect to  $\mathbf{x}$ , *i.e.* the function  $L(\cdot|\mathbf{x})$  should not vary too wildly in  $\mathbf{x}$  when  $y$  is large. The finite limit of  $\nabla_{\mathbf{x}} \log \bar{F}(y|\mathbf{x})/\log(y)$  as  $y \rightarrow \infty$  will then be  $\boldsymbol{\mu}(\mathbf{x}) = \nabla \gamma(\mathbf{x})/\gamma^2(\mathbf{x}) \in \mathbb{R}^p$ . In summary, while condition  $\mathcal{L}_\omega$  is a Lipschitz condition on the log-conditional survival function at extreme levels, condition  $\mathcal{D}_\omega$  is essentially an appropriate analog about its gradient with respect to  $\mathbf{x}$ .

**2.2. Intermediate quantile estimation.** Let  $\tau = \tau_n \uparrow 1$  as  $n \rightarrow \infty$ . We show that  $\hat{q}_n(\tau_n|\mathbf{x})$  is asymptotically normal if  $nh_n^p(1 - \tau_n) \rightarrow \infty$ , *i.e.*  $\tau_n \uparrow 1$  slowly enough. Hereafter, under condition  $\mathcal{KS}$  and  $\mathcal{D}_\omega$ , we let  $\Upsilon_K : (\Delta, \mathbf{x}) \in \mathbb{R} \times \mathbb{R}^p \mapsto \frac{\gamma(\mathbf{x})\Delta}{2} \int_{\mathbb{R}^p} K(\mathbf{u})(\mathbf{u}^\top \boldsymbol{\mu}(\mathbf{x}))^2 d\mathbf{u}$ .

**THEOREM 2.1.** *Assume that conditions  $\mathcal{M}$ ,  $\mathcal{A}(l_n, r_n)$ ,  $\mathcal{K}$ ,  $\mathcal{L}_g$ ,  $\mathcal{L}_\omega$ ,  $\mathcal{B}_p$  and  $\mathcal{B}_\Omega$  hold with  $\sum_{j=1}^\infty j^\eta \alpha(j) < \infty$  for some  $\eta > 1$ . Let  $\tau_n \uparrow 1$ , fix  $J \geq 1$ , pick distinct  $c_j \in (0, 1]$  and let  $\tau_{n,j}$  be such that  $1 - \tau_{n,j} = c_j(1 - \tau_n)(1 + o(1))$  as  $n \rightarrow \infty$  (for  $1 \leq j \leq J$ ). Assume further that  $h_n \rightarrow 0$  is such that  $nh_n^p(1 - \tau_n) \rightarrow \infty$ ,  $nh_n^{p+2}(1 - \tau_n) \log^2(1 - \tau_n) \rightarrow 0$  and  $\sqrt{nh_n^p(1 - \tau_n)}A((1 - \tau_n)^{-1}|\mathbf{x}) = O(1)$ , and that there is  $\delta > 0$  such that  $r_n(r_n/\sqrt{nh_n^p(1 - \tau_n)})^\delta \rightarrow 0$ . Then*

$$\sqrt{nh_n^p(1 - \tau_n)} \left( \frac{\hat{q}_n(\tau_{n,j}|\mathbf{x})}{q(\tau_{n,j}|\mathbf{x})} - 1 \right)_{1 \leq j \leq J} \xrightarrow{d} \mathcal{N} \left( \mathbf{0}_J, \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} \gamma^2(\mathbf{x}) \mathbf{M} \right),$$

where  $\mathbf{M} = [1/\max(c_j, c_l)]_{1 \leq j, l \leq J}$ . If conditions  $\mathcal{KS}$ ,  $\mathcal{D}_g$  and  $\mathcal{D}_\omega$  hold, then condition  $nh_n^{p+2}(1 - \tau_n) \log^2(1 - \tau_n) \rightarrow 0$  can be replaced by the weaker bias assumption  $\sqrt{nh_n^p(1 - \tau_n)} \times h_n^2 \log^2(1 - \tau_n) \rightarrow \Delta \in [0, \infty)$ , in which case, provided  $r_n h_n^p \rightarrow 0$ , the asymptotic mean  $\mathbf{0}_J$  of the above Gaussian limit is replaced by  $\Upsilon_K(\Delta, \mathbf{x}) \mathbf{1}_J$ .

If  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is  $\rho$ -mixing, condition  $\sum_{j=1}^\infty j^\eta \alpha(j) < \infty$  can be replaced by summability of the  $\rho$ -mixing coefficient series. If  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is also  $\psi$ -mixing with  $\sum_{j=1}^\infty \psi(j) < \infty$  (instead of  $\sum_{j=1}^\infty j^\eta \alpha(j) < \infty$  for some  $\eta > 1$ , or summability of  $\rho$ -mixing coefficients), all conditions on  $(l_n)$  and  $(r_n)$  (including  $\mathcal{A}(l_n, r_n)$ ) as well as  $\mathcal{B}_p$  and  $\mathcal{B}_\Omega$  can also be dropped.

It is remarkable that the asymptotic distribution in Theorem 2.1 is exactly the one obtained in the i.i.d. setting by [10] under an unnecessary regularity assumption on conditional tails. This is not true in the unconditional setting, see [15]. The essential difference is that the kernel estimator only takes into account those pairs  $(\mathbf{X}_t, Y_t)$  such that  $\mathbf{X}_t$  are close enough to  $\mathbf{x}$ , and the mixing and stationarity assumptions ensure that such data points are far enough apart in time and hence asymptotically independent. This phenomenon has been observed in other contexts, such as nonparametric conditional Expected Shortfall estimation [21, p.784].

Theorem 2.1 actually holds under weaker bias assumptions than those of [10]. Indeed, when  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is geometrically  $\alpha$ -mixing, and if  $A(t|\mathbf{x}) \propto t^{\rho(\mathbf{x})}$ , as is the case in a wide range of heavy-tailed models used in practice [see *e.g.* 2, Table 2.1 p.59], the optimal rate of convergence is  $n^{\rho(\mathbf{x})/(1-(p+2)\rho(\mathbf{x}))}$  when  $\mathcal{K}$ ,  $\mathcal{L}_g$  and  $\mathcal{L}_\omega$  hold, while it is equal to  $n^{2\rho(\mathbf{x})/(2-(p+4)\rho(\mathbf{x}))}$  if  $\mathcal{KS}$ ,  $\mathcal{D}_g$  and  $\mathcal{D}_\omega$  are satisfied, see Section A.3 of the Supplementary Material document [12] for details. In the latter setting, note that  $p = 0$  yields the optimal convergence rate  $n^{\rho(\mathbf{x})/(1-2\rho(\mathbf{x}))}$  of unconditional extreme value estimators in heavy-tailed models [see 14, p.77], while the case  $\rho(\mathbf{x}) \rightarrow -\infty$ , corresponding to the ideal but unrealistic scenario where all the  $Y_t$  such that  $\mathbf{X}_t \in B(\mathbf{x}, h_n)$  can be used, yields the optimal convergence rate  $n^{-2/(p+4)}$ , *i.e.* the optimal convergence rate of nonparametric estimators of a twice continuously differentiable *central* conditional quantile, see [5]. Of course, the nonparametric nature of the methodology coupled with the double sparsity of relevant data (in  $\mathbf{x}$ , due to kernel smoothing, and in  $y$ , due to taking only extreme values of the response into account) limits



the applicability of the method to low dimensions. The nonparametric approach is, however, indifferent to the structure of the conditional distribution, unlike techniques adapted to taking large dimensions into account that rely on strong model specifications such as linear quantile structures, see for instance [6, 7, 25].

The bias component  $\Upsilon_K(\Delta, \mathbf{x})$  appearing under conditions  $\mathcal{D}_g$  and  $\mathcal{D}_\omega$  is an analog of the bias component in kernel regression for the mean. It is linked to the gradient of the target extreme conditional quantile through the vector  $\boldsymbol{\mu}(\mathbf{x}) = \nabla\gamma(\mathbf{x})/\gamma^2(\mathbf{x})$ . A key factor in extremal regression is therefore the *variation in conditional extremes*, while in standard regression it is also important to account for *changes of shape* in the regression function through its second derivatives [see *e.g.* 26, p.73]. We finally note that  $\gamma(\mathbf{x})\boldsymbol{\mu}(\mathbf{x}) = \nabla\gamma(\mathbf{x})/\gamma(\mathbf{x}) = \nabla(\log\gamma)(\mathbf{x})$  is reminiscent of the *design bias* for classical regression in the terminology of [26]: the higher and less variable  $\gamma$  around  $\mathbf{x}$ , the bigger and more stable the local number of extreme observations, and the easier the conditional extreme value estimation problem.

**2.3. Extreme quantile estimation.** Consider now  $\tau'_n$  such that  $nh_n^p(1 - \tau'_n) \rightarrow c < \infty$ . Then, in a neighborhood of  $\mathbf{x}$ , very few or no top observations  $Y_t$  will be close to the extreme value  $q(\tau'_n|\mathbf{x})$ , so the empirical estimator  $\widehat{q}_n(\tau'_n|\mathbf{x})$  will not be consistent. However, the conditional heavy tail assumption suggests the extrapolation formula  $q(\tau'_n|\mathbf{x}) \approx ((1 - \tau'_n)/(1 - \tau_n))^{-\gamma(\mathbf{x})}q(\tau_n|\mathbf{x})$  for  $n$  large. Plugging in consistent estimators  $\widehat{\gamma}(\mathbf{x})$  of  $\gamma(\mathbf{x})$  and  $\bar{q}_n(\tau_n|\mathbf{x})$  of  $q(\tau_n|\mathbf{x})$  yields a conditional Weissman-type estimator (see [27]) of  $q(\tau'_n|\mathbf{x})$ :

$$\widehat{q}_{n,\tau_n}^W(\tau'_n|\mathbf{x}) = \left(\frac{1 - \tau'_n}{1 - \tau_n}\right)^{-\widehat{\gamma}(\mathbf{x})} \bar{q}_n(\tau_n|\mathbf{x}).$$

We prove below that the choice of  $\widehat{\gamma}(\mathbf{x})$  is crucial since  $\widehat{q}_{n,\tau_n}^W(\tau'_n|\mathbf{x})$  inherits its asymptotics.

**THEOREM 2.2.** *Assume that condition  $\mathcal{M}$  holds with  $\rho(\mathbf{x}) < 0$ . Let  $\tau_n, \tau'_n \uparrow 1$  be such that  $(1 - \tau'_n)/(1 - \tau_n) \rightarrow 0$  and assume that  $v_n(\bar{q}_n(\tau_n|\mathbf{x})/q(\tau_n|\mathbf{x}) - 1) = O_{\mathbb{P}}(1)$  and  $v_n(\widehat{\gamma}(\mathbf{x}) - \gamma(\mathbf{x})) \xrightarrow{d} \Gamma$ , where  $\Gamma$  is a nondegenerate distribution and  $v_n \rightarrow \infty$ . If  $v_n A((1 - \tau_n)^{-1}|\mathbf{x}) = O(1)$  and  $v_n/\log[(1 - \tau_n)/(1 - \tau'_n)] \rightarrow \infty$ , then*

$$\frac{v_n}{\log[(1 - \tau_n)/(1 - \tau'_n)]} \left( \frac{\widehat{q}_{n,\tau_n}^W(\tau'_n|\mathbf{x})}{q(\tau'_n|\mathbf{x})} - 1 \right) \xrightarrow{d} \Gamma.$$

In our context,  $v_n$  is typically  $\sqrt{nh_n^p(1 - \tau_n)}$  and  $\bar{q}_n(\tau_n|\mathbf{x}) = \widehat{q}_n(\tau_n|\mathbf{x})$ . As a conditional tail index estimator, we use the so-called “kernel version of the Hill estimator” from [10]:

$$\widehat{\gamma}_{\tau_n}^{(J)}(\mathbf{x}) = \frac{1}{\log(J!)} \sum_{j=1}^J \log \left( \frac{\widehat{q}_n(1 - (1 - \tau_n)/j|\mathbf{x})}{\widehat{q}_n(\tau_n|\mathbf{x})} \right), \text{ for a fixed } J \geq 2.$$

The asymptotic distribution of  $\widehat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})$  can be deduced from Theorem 2.1. Note that the number  $J$  of high quantiles is fixed; the case  $J = J_n \rightarrow \infty$ , which would correspond to the classical established theory of the Hill estimator [see *e.g.* 14, Section 3.2], cannot be handled using Theorem 2.1. Obtaining asymptotic results for a growing number of summands  $J = J_n$  would involve the difficult study of conditional tail empirical quantile processes.

**THEOREM 2.3.** *Work under the conditions of Theorem 2.1, and assume in addition that  $\sqrt{nh_n^p(1 - \tau_n)}A((1 - \tau_n)^{-1}|\mathbf{x}) \rightarrow \lambda(\mathbf{x}) \in \mathbb{R}$ . Then,*

$$\sqrt{nh_n^p(1 - \tau_n)} \left( \widehat{\gamma}_{\tau_n}^{(J)}(\mathbf{x}) - \gamma(\mathbf{x}), \frac{\widehat{q}_n(\tau_n|\mathbf{x})}{q(\tau_n|\mathbf{x})} - 1 \right)$$

$$\xrightarrow{d} \mathcal{N} \left( \left( \frac{\lambda(\mathbf{x})}{\log(J!)} \sum_{j=2}^J \frac{j^{\rho(\mathbf{x})} - 1}{\rho(\mathbf{x})}, 0 \right), \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} \gamma^2(\mathbf{x}) \begin{pmatrix} \frac{J(J-1)(2J-1)}{6 \log^2(J!)} & 0 \\ 0 & 1 \end{pmatrix} \right).$$

If conditions  $\mathcal{KS}$ ,  $\mathcal{D}_g$  and  $\mathcal{D}_\omega$  hold, under the weaker bias assumption  $\sqrt{nh_n^p(1-\tau_n)} \times h_n^2 \log^2(1-\tau_n) \rightarrow \Delta \in [0, \infty)$  and if  $r_n h_n^p \rightarrow 0$ , the second component of the asymptotic mean of the above Gaussian limit is replaced by  $\Upsilon_K(\Delta, \mathbf{x})$ .

The asymptotic variance of  $\widehat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})$  is minimal when  $J = 9$ . We shall adopt this choice in our finite-sample experiments. Theorem 2.3 improves upon Corollary 2 in [10], by removing unnecessary assumptions about the right conditional tail, and by weakening a bias assumption, see the discussion below Theorem 2.1. The bias component  $\Delta$  does not appear in the limiting distribution of  $\widehat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})$ , which suggests that the local variation of conditional extremes is not as important in the estimation of the conditional shape parameter  $\gamma(\mathbf{x})$  of the Pareto approximating distribution as it is for its scale  $q(\tau_n|\mathbf{x})$ .

### 3. Nonparametric extremal expectile regression.

**3.1. Framework.** Rewrite the conditional distribution function in (3) as  $E(y|\mathbf{x}) = 1 - \varphi^{(1)}(y|\mathbf{x}) / (2\varphi^{(1)}(y|\mathbf{x}) + y - m(\mathbf{x}))$  where  $\varphi^{(a)}(y|\mathbf{x}) = \mathbb{E}((Y - y)^a \mathbb{1}_{\{Y > y\}} | \mathbf{X} = \mathbf{x})$  and  $m(\mathbf{x}) = \mathbb{E}(Y | \mathbf{X} = \mathbf{x})$ . Consider the following kernel smoother for  $E(y|\mathbf{x})$ :

$$\widehat{E}_n(y|\mathbf{x}) = 1 - \frac{\widehat{\varphi}_n^{(1)}(y|\mathbf{x})}{2\widehat{\varphi}_n^{(1)}(y|\mathbf{x}) + (y - \widehat{m}_n(\mathbf{x}))} \text{ with } \widehat{m}_n(\mathbf{x}) = \frac{1}{nh_n^p \widehat{g}_n(\mathbf{x})} \sum_{t=1}^n Y_t K\left(\frac{\mathbf{x} - \mathbf{X}_t}{h_n}\right)$$

$$\text{and } \widehat{\varphi}_n^{(1)}(y|\mathbf{x}) = \frac{1}{nh_n^p \widehat{g}_n(\mathbf{x})} \sum_{t=1}^n (Y_t - y) \mathbb{1}_{\{Y_t > y\}} K\left(\frac{\mathbf{x} - \mathbf{X}_t}{h_n}\right).$$

The estimator  $\widehat{m}_n$  is the Nadaraya-Watson estimator of the regression function  $m$ . The characterization of conditional expectiles as  $e(\tau|\mathbf{x}) \equiv \inf \{y \in \mathbb{R} | E(y|\mathbf{x}) \geq \tau\}$  (see [20]) implies that they can be estimated by their empirical counterparts

$$\widehat{e}_n(\tau|\mathbf{x}) = \inf \left\{ y \in \mathbb{R} \mid \widehat{E}_n(y|\mathbf{x}) \geq \tau \right\}.$$

This is in fact exactly the *asymmetric least squares estimator*  $\widehat{e}_n(\tau|\mathbf{x}) = \arg \min_{\theta \in \mathbb{R}} \int_{\mathbb{R}} [\eta_\tau(y - \theta) - \eta_\tau(y)] d\widehat{F}_n(y|\mathbf{x})$  obtained by smoothing up the loss function defining  $e(\tau|\mathbf{x})$  in (2).

The definition of conditional expectiles in (2) and (3) requires  $\mathbb{E}(|Y| | \mathbf{X} = \mathbf{x}) < \infty$ . To obtain the asymptotic normality of  $\widehat{e}_n(\tau|\mathbf{x})$  at intermediate levels, we make the following additional assumptions on conditional tail heaviness and regularity of conditional moments.

**Condition  $\mathcal{H}_\delta$**  One has  $\gamma(\mathbf{x}) < 1/(2 + \delta)$  and there exists  $r > 0$  such that the function  $\mathbb{E}(Y_-^{2+\delta} | \mathbf{X} = \cdot)$  is bounded on  $B(\mathbf{x}, r)$ .

**Condition  $\mathcal{L}_m$**  The response  $Y$  has a finite second moment given  $\mathbf{X} = \mathbf{x}$ , and the conditional mean functions  $\mathbb{E}(Y | \mathbf{X} = \cdot)$  and  $\mathbb{E}(Y^2 | \mathbf{X} = \cdot)$  are Lipschitz continuous at  $\mathbf{x}$ .

**Condition  $\mathcal{B}_m$**  There exists  $r > 0$  such that

$$\sup_{t \geq 1} \sup_{\mathbf{x}_1, \mathbf{x}_{t+1} \in B(\mathbf{x}, r)} \mathbb{E}(Y_1^2 + Y_{t+1}^2 | \mathbf{X}_1 = \mathbf{x}_1, \mathbf{X}_{t+1} = \mathbf{x}_{t+1}) < \infty.$$

Condition  $\mathcal{H}_\delta$  (in which  $\delta > 0$ ) guarantees a finite conditional moment of order  $(2 + \delta)$  in a neighborhood of  $\mathbf{x}$ ; in the unconditional framework, Theorem 2 in [11] requires the analogous  $\mathbb{E}(Y_-^{2+\delta}) < \infty$ . Conditions  $\mathcal{L}_m$  and  $\mathcal{B}_m$  ensure the convergence of the regression mean

at a standard rate, see Proposition A.1 in Section A.4 of the Supplementary Material document [12]. Finally, in the spirit of Section 2, a stronger version of condition  $\mathcal{L}_m$  will be imposed to obtain better rates of convergence:

*Condition  $\mathcal{D}_m$*  The response  $Y$  has a finite second moment given  $\mathbf{X} = \mathbf{x}$ , and the conditional mean functions  $\mathbb{E}(Y|\mathbf{X} = \cdot)$  and  $\mathbb{E}(Y^2|\mathbf{X} = \cdot)$  are continuously differentiable in a neighborhood of  $\mathbf{x}$  and have Lipschitz continuous gradients at  $\mathbf{x}$ .

An alternative approach uses the asymptotic proportionality between extreme quantiles and expectiles of heavy-tailed distributions, namely,  $e(\tau|\mathbf{x})/q(\tau|\mathbf{x}) \rightarrow (1/\gamma(\mathbf{x}) - 1)^{-\gamma(\mathbf{x})}$  as  $\tau \uparrow 1$  [see Proposition 2.3 in 3]. Plugging in the estimators  $\hat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})$  and  $\hat{q}_n(\tau_n|\mathbf{x})$  results in the *quantile-based estimator*

$$\check{e}_n(\tau_n|\mathbf{x}) \equiv \check{e}_n^{(J)}(\tau_n|\mathbf{x}) = (1/\hat{\gamma}_{\tau_n}^{(J)}(\mathbf{x}) - 1)^{-\hat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})} \hat{q}_n(\tau_n|\mathbf{x}).$$

**3.2. Intermediate expectile estimation.** We first derive the asymptotic distribution of the intermediate expectile estimator  $\hat{e}_n(\tau_n|\mathbf{x})$  jointly with an intermediate quantile estimator. This will be key to the construction of an expectile-based estimator for the conditional tail index, and therefore to our extrapolation of conditional expectiles to extreme levels.

**THEOREM 3.1.** *Assume that conditions  $\mathcal{M}$ ,  $\mathcal{A}(l_n, r_n)$ ,  $\mathcal{K}$ ,  $\mathcal{H}_\delta$ ,  $\mathcal{L}_g$ ,  $\mathcal{L}_m$ ,  $\mathcal{L}_\omega$ ,  $\mathcal{B}_p$ ,  $\mathcal{B}_m$  and  $\mathcal{B}_\Omega$  hold with  $\sum_{j=1}^{\infty} j^\eta [\alpha(j)]^{\delta/(2+\delta)} < \infty$  for some  $\eta > \delta/(2 + \delta)$ . Let  $\tau_n \uparrow 1$  and  $\kappa > 0$  be given, and let  $\beta_n$  be such that  $1 - \beta_n = \kappa(1 - \tau_n)(1 + o(1))$  as  $n \rightarrow \infty$ . Assume further that  $h_n \rightarrow 0$  is such that  $nh_n^p(1 - \tau_n) \rightarrow \infty$ ,  $nh_n^{p+2}(1 - \tau_n) \log^2(1 - \tau_n) \rightarrow 0$ ,  $\sqrt{nh_n^p(1 - \tau_n)}A((1 - \tau_n)^{-1}|\mathbf{x}) = O(1)$  and  $r_n(r_n/\sqrt{nh_n^p(1 - \tau_n)})^\delta \rightarrow 0$ . Then*

$$\sqrt{nh_n^p(1 - \tau_n)} \left( \frac{\hat{e}_n(\tau_n|\mathbf{x})}{e(\tau_n|\mathbf{x})} - 1, \frac{\hat{q}_n(\beta_n|\mathbf{x})}{q(\beta_n|\mathbf{x})} - 1 \right) \xrightarrow{d} \mathcal{N} \left( (0, 0), \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} \gamma^2(\mathbf{x}) \boldsymbol{\Sigma}(\mathbf{x}) \right),$$

where the  $2 \times 2$  symmetric matrix  $\boldsymbol{\Sigma}(\mathbf{x})$  has entries  $\Sigma_{1,1}(\mathbf{x}) = 2\gamma(\mathbf{x})/(1 - 2\gamma(\mathbf{x}))$ ,  $\Sigma_{2,2}(\mathbf{x}) = \kappa^{-1}$  and

$$\Sigma_{1,2}(\mathbf{x}) = \begin{cases} \kappa^{-1} & \text{if } \kappa \geq 1/\gamma(\mathbf{x}) - 1, \\ \left( \frac{1}{\gamma(\mathbf{x})} - 1 \right)^{\gamma(\mathbf{x})} \frac{\kappa^{-\gamma(\mathbf{x})}}{1 - \gamma(\mathbf{x})} - 1 & \text{if } \kappa < 1/\gamma(\mathbf{x}) - 1. \end{cases}$$

If conditions  $\mathcal{KS}$ ,  $\mathcal{D}_g$ ,  $\mathcal{D}_m$  and  $\mathcal{D}_\omega$  hold, then condition  $nh_n^{p+2}(1 - \tau_n) \log^2(1 - \tau_n) \rightarrow 0$  can be replaced by the weaker bias assumption  $\sqrt{nh_n^p(1 - \tau_n)} \times h_n^2 \log^2(1 - \tau_n) \rightarrow \Delta \in [0, \infty)$ , and if moreover  $r_n h_n^p \rightarrow 0$ , then the asymptotic mean  $(0, 0)$  of the above Gaussian limit is replaced by  $\Upsilon_K(\Delta, \mathbf{x}) \times (1, 1)$ .

If  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is  $\rho$ -mixing with summability of the  $\rho$ -mixing coefficient series (instead of  $\sum_{j=1}^{\infty} j^\eta [\alpha(j)]^{\delta/(2+\delta)} < \infty$  for some  $\eta > \delta/(2 + \delta)$ ), condition  $\mathcal{B}_m$  can be dropped and condition  $\mathcal{H}_\delta$  can be replaced by  $\gamma(\mathbf{x}) < 1/(2 + \delta)$ . If  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is in fact also  $\psi$ -mixing with  $\sum_{j=1}^{\infty} \psi(j) < \infty$  (instead of  $\sum_{j=1}^{\infty} j^\eta [\alpha(j)]^{\delta/(2+\delta)} < \infty$  for some  $\eta > \delta/(2 + \delta)$ ), or summability of the  $\rho$ -mixing coefficients), all conditions on  $(l_n)$  and  $(r_n)$  (including  $\mathcal{A}(l_n, r_n)$ ) as well as conditions  $\mathcal{B}_p$  and  $\mathcal{B}_\Omega$  can also be dropped, and condition  $\mathcal{H}_\delta$  can be replaced by the weaker requirement that  $0 < \gamma(\mathbf{x}) < 1/2$ .

For  $\hat{e}_n(\tau_n|\mathbf{x})$ , we obtain again the same asymptotic distribution as in the i.i.d. setting under weaker moment and regularity assumptions and with a faster optimal convergence rate, see Theorem 1 in [17] and the discussion below our Theorem 2.1. By contrast, mixing changes the asymptotic distribution of unconditional empirical intermediate expectiles, see [13].

We turn to the asymptotic normality of the quantile-based expectile estimator  $\check{e}_n(\tau_n|\mathbf{x})$ .

**THEOREM 3.2.** *Work under the conditions of Theorem 2.1 with  $\gamma(\mathbf{x}) < 1$  and  $\mathbb{E}(Y_- | \mathbf{X} = \mathbf{x}) < \infty$ , and assume in addition that  $\sqrt{nh_n^p(1-\tau_n)}A((1-\tau_n)^{-1}|\mathbf{x}) \rightarrow \lambda_1(\mathbf{x}) \in \mathbb{R}$  and  $\sqrt{nh_n^p(1-\tau_n)}/q(\tau_n|\mathbf{x}) \rightarrow \lambda_2(\mathbf{x}) \in \mathbb{R}$ . Then*

$$\sqrt{nh_n^p(1-\tau_n)} \left( \frac{\check{e}_n(\tau_n|\mathbf{x})}{e(\tau_n|\mathbf{x})} - 1 \right) \xrightarrow{d} \mathcal{N} \left( b_1^{(J)}(\gamma(\mathbf{x}), \rho(\mathbf{x}))\lambda_1(\mathbf{x}) + b_2(\gamma(\mathbf{x}), m(\mathbf{x}))\lambda_2(\mathbf{x}), \right. \\ \left. \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} \gamma^2(\mathbf{x}) \left( 1 + [(1-\gamma(\mathbf{x}))^{-1} - \log(1/\gamma(\mathbf{x}) - 1)]^2 \frac{J(J-1)(2J-1)}{6 \log^2(J!)} \right) \right),$$

where

$$b_1^{(J)}(\gamma(\mathbf{x}), \rho(\mathbf{x})) = \frac{(1-\gamma(\mathbf{x}))^{-1} - \log(1/\gamma(\mathbf{x}) - 1)}{\log(J!)} \sum_{j=2}^J \frac{j^{\rho(\mathbf{x})} - 1}{\rho(\mathbf{x})} \\ - \left( \frac{(1/\gamma(\mathbf{x}) - 1)^{-\rho(\mathbf{x})}}{1 - \gamma(\mathbf{x}) - \rho(\mathbf{x})} + \frac{(1/\gamma(\mathbf{x}) - 1)^{-\rho(\mathbf{x})} - 1}{\rho(\mathbf{x})} \right)$$

and  $b_2(\gamma(\mathbf{x}), m(\mathbf{x})) = -\gamma(\mathbf{x})(1/\gamma(\mathbf{x}) - 1)^{\gamma(\mathbf{x})}m(\mathbf{x})$ .

If conditions  $\mathcal{KS}$ ,  $\mathcal{D}_g$  and  $\mathcal{D}_\omega$  hold, under the weaker bias assumption  $\sqrt{nh_n^p(1-\tau_n)} \times h_n^2 \log^2(1-\tau_n) \rightarrow \Delta \in [0, \infty)$  and if  $r_n h_n^p \rightarrow 0$ , the quantity  $\Upsilon_K(\Delta, \mathbf{x})$  is added to the asymptotic mean of the above Gaussian limit.

As is the case for  $\hat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})$ , the asymptotic variance of  $\check{e}_n(\tau_n|\mathbf{x}) \equiv \check{e}_n^{(J)}(\tau_n|\mathbf{x})$  is minimal when  $J = 9$ . An important benefit of this quantile-based estimator is that its asymptotic normality does not require conditions  $\mathcal{H}_\delta$ ,  $\mathcal{L}_m$  and  $\mathcal{B}_m$ . In particular, it applies as soon as a finite conditional first moment exists. The estimator  $\check{e}_n(\tau_n|\mathbf{x})$  is biased, however, and its variance is higher than the variance of  $\hat{e}_n(\tau_n|\mathbf{x})$  when  $\gamma(\mathbf{x})$  is only moderately large.

**3.3. Extreme expectile estimation.** The asymptotic proportionality of extreme quantiles and expectiles entails  $e(\tau'_n|\mathbf{x}) \approx ((1-\tau'_n)/(1-\tau_n))^{-\gamma(\mathbf{x})}e(\tau_n|\mathbf{x})$ , for  $n$  large. Plugging in consistent estimators  $\hat{\gamma}(\mathbf{x})$  of  $\gamma(\mathbf{x})$  and  $\bar{e}_n(\tau_n|\mathbf{x})$  of  $e(\tau_n|\mathbf{x})$  (such as  $\hat{e}_n(\tau_n|\mathbf{x})$  or  $\check{e}_n(\tau_n|\mathbf{x})$ ) yields a Weissman-type estimator of  $e(\tau'_n|\mathbf{x})$ :

$$\hat{e}_{n,\tau_n}^W(\tau'_n|\mathbf{x}) = \left( \frac{1-\tau'_n}{1-\tau_n} \right)^{-\hat{\gamma}(\mathbf{x})} \bar{e}_n(\tau_n|\mathbf{x}).$$

The next result states that  $\hat{e}_{n,\tau_n}^W(\tau'_n|\mathbf{x})$  also inherits the asymptotic distribution of  $\hat{\gamma}(\mathbf{x})$ .

**THEOREM 3.3.** *Assume that condition  $\mathcal{M}$  holds, with  $\gamma(\mathbf{x}) < 1$ ,  $\rho(\mathbf{x}) < 0$  and  $\mathbb{E}(Y_- | \mathbf{X} = \mathbf{x}) < \infty$ . Let  $\tau_n, \tau'_n \uparrow 1$  be such that  $(1-\tau'_n)/(1-\tau_n) \rightarrow 0$  and assume that  $v_n(\bar{e}_n(\tau_n|\mathbf{x})/e(\tau_n|\mathbf{x}) - 1) = O_{\mathbb{P}}(1)$  and  $v_n(\hat{\gamma}(\mathbf{x}) - \gamma(\mathbf{x})) \xrightarrow{d} \Gamma$ , where  $\Gamma$  is a nondegenerate distribution and  $v_n \rightarrow \infty$ . If  $v_n A((1-\tau_n)^{-1}|\mathbf{x}) = O(1)$ ,  $v_n/q(\tau_n|\mathbf{x}) = O(1)$  and  $v_n/\log[(1-\tau_n)/(1-\tau'_n)] \rightarrow \infty$ , then*

$$\frac{v_n}{\log[(1-\tau_n)/(1-\tau'_n)]} \left( \frac{\hat{e}_{n,\tau_n}^W(\tau'_n|\mathbf{x})}{e(\tau'_n|\mathbf{x})} - 1 \right) \xrightarrow{d} \Gamma.$$

Let now  $(\tau_n)$  be an intermediate sequence. Define an expectile-based estimator of  $\gamma(\mathbf{x})$  as

$$\hat{\gamma}_{\tau_n}^E(\mathbf{x}) = \left( 1 + \frac{\hat{F}_n(\hat{e}_n(\tau_n|\mathbf{x})|\mathbf{x})}{1-\tau_n} \right)^{-1},$$

where  $\widehat{F}_n(\cdot|\mathbf{x}) = 1 - \widehat{F}_n(\cdot|\mathbf{x})$ . We derive the joint asymptotic normality of  $(\widehat{\gamma}_{\tau_n}^E(\mathbf{x}), \widehat{e}_n(\tau_n|\mathbf{x}))$ .

**THEOREM 3.4.** *Work under the conditions of Theorem 3.1. Assume further that  $\sqrt{nh_n^p(1-\tau_n)}A((1-\tau_n)^{-1}|\mathbf{x}) \rightarrow \lambda_1(\mathbf{x}) \in \mathbb{R}$  and  $\sqrt{nh_n^p(1-\tau_n)}/q(\tau_n|\mathbf{x}) \rightarrow \lambda_2(\mathbf{x}) \in \mathbb{R}$ . Then*

$$\sqrt{nh_n^p(1-\tau_n)} \left( \widehat{\gamma}_{\tau_n}^E(\mathbf{x}) - \gamma(\mathbf{x}), \frac{\widehat{e}_n(\tau_n|\mathbf{x})}{e(\tau_n|\mathbf{x})} - 1 \right) \xrightarrow{d} \mathcal{N} \left( (b_E(\mathbf{x}), 0), \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} v_E(\mathbf{x}) \begin{pmatrix} 1 - \gamma(\mathbf{x}) & 1 \\ 1 & 2 \end{pmatrix} \right),$$

where  $v_E(\mathbf{x}) = \gamma^3(\mathbf{x})/(1 - 2\gamma(\mathbf{x}))$  and

$$b_E(\mathbf{x}) = \frac{\gamma(\mathbf{x})(1/\gamma(\mathbf{x}) - 1)^{1-\rho(\mathbf{x})}}{1 - \gamma(\mathbf{x}) - \rho(\mathbf{x})} \lambda_1(\mathbf{x}) + \gamma^2(\mathbf{x})(1/\gamma(\mathbf{x}) - 1)^{\gamma(\mathbf{x})+1} m(\mathbf{x}) \lambda_2(\mathbf{x}).$$

If conditions  $\mathcal{KS}$ ,  $\mathcal{D}_g$ ,  $\mathcal{D}_m$  and  $\mathcal{D}_\omega$  hold, under the weaker bias assumption  $\sqrt{nh_n^p(1-\tau_n)} \times h_n^2 \log^2(1-\tau_n) \rightarrow \Delta \in [0, \infty)$  and if  $r_n h_n^p \rightarrow 0$ , the second component of the asymptotic mean of the above Gaussian limit is replaced by  $\Upsilon_K(\Delta, \mathbf{x})$ .

Again, mixing does not impact the asymptotic distribution of  $\widehat{\gamma}_{\tau_n}^E(\mathbf{x})$  obtained in the i.i.d. setting [see Theorem 4 in 17] and the bias component  $\Delta$  does not appear in the limit. The advantage of using  $\widehat{\gamma}_{\tau_n}^E(\mathbf{x})$  is its low asymptotic variance when  $\gamma(\mathbf{x})$  is moderately high, at the expense of a strong bias. We discuss a bias-correction methodology in Section 5.2.

**4. Regression models covered by our framework.** We now draw a non-exhaustive list of examples satisfying our assumptions. The conditions involving the mixing coefficients or the sequences  $(l_n)$  and  $(r_n)$  hold automatically when the stochastic process  $((\mathbf{X}_t, Y_t))_{t \geq 1}$  is geometrically  $\alpha$ -mixing, namely, there exists  $a \in (0, 1)$  such that  $\alpha(n) = O(a^n)$ . Besides, in the typical extreme value models where  $A(t|\mathbf{x}) \propto t^{\rho(\mathbf{x})}$ , our assumptions linking  $\tau_n$  and  $h_n$  will be satisfied if  $h_n = C_1 n^{-h}$  and  $\tau_n = 1 - C_2 n^{-\tau}$ , for any  $C_1, C_2 > 0$  and suitably chosen  $h, \tau > 0$ , see the discussion below Theorem 2.1. We therefore focus in this section on the validity of assumptions  $\mathcal{M}$ ,  $\mathcal{H}_\delta$  (for a given  $\delta > 0$ ),  $\mathcal{L}_g$ ,  $\mathcal{L}_m$ ,  $\mathcal{L}_\omega$ ,  $\mathcal{B}_p$ ,  $\mathcal{B}_m$ ,  $\mathcal{B}_\Omega$ ,  $\mathcal{D}_g$ ,  $\mathcal{D}_m$  and  $\mathcal{D}_\omega$ . We provide an extended discussion with more insight, and sometimes weaker conditions, in Section B of the Supplementary Material document [12]. We also give a full treatment therein of the instructive case of  $m$ -dependent (including i.i.d.) observations.

*Location-scale model with possible temporal misspecification.* Suppose that  $Y_t = m(\mathbf{X}_t) + \sigma(\mathbf{X}_t)\varepsilon_t$  where  $m$  and  $\sigma > 0$  are location and scale components, and  $(\varepsilon_t)$  is a stationary and centered sequence of unobserved heavy-tailed innovations independent from the sequence  $(\mathbf{X}_t)$ . The  $\varepsilon_t$  can be dependent, yielding a possibly misspecified regression model in the sense that relevant, serially correlated covariates can be missing in  $\mathbf{X}_t$  but left in the error  $\varepsilon_t$ .

**PROPOSITION 4.1 (Location-scale model).** *Assume that  $(\mathbf{X}_t)_{t \geq 1}$  is  $\beta$ -mixing (i.e. absolutely regular) and  $(\varepsilon_t)_{t \geq 1}$  is strongly mixing. Suppose further that:*

- For any  $t \geq 1$ , the random pairs  $(\mathbf{X}_1, \mathbf{X}_{t+1})$  have absolutely continuous distributions whose first marginal has a p.d.f.  $g$  such that  $g(\mathbf{x}) > 0$ , and the functions  $g$ ,  $m$  and  $\sigma$  are continuously differentiable in a neighborhood of  $\mathbf{x}$ .

• The errors  $\varepsilon_t$  have a common p.d.f.  $f_\varepsilon$  with respect to the Lebesgue measure such that  $f_\varepsilon(z) = c_0 z^{-1/\gamma-1}(1 + d_0 z^{-a} + d'_0 z^{-a-b}(1 + o(1)))$  as  $z \rightarrow \infty$ , where  $\gamma > 0$ ,  $c_0 > 0$ ,  $d_0, d'_0 \neq 0$  and  $a, b > 0$  are such that either  $a \neq 1$ ,  $a = 1 \neq b$ , or  $a = b = 1$  with then  $2d'_0(1 + \gamma) \neq d_0^2(1 + 2\gamma)$ .

Then conditions  $\mathcal{M}$ ,  $\mathcal{L}_g$ ,  $\mathcal{L}_\omega$ ,  $\mathcal{B}_p$  and  $\mathcal{B}_\Omega$  hold. If moreover  $\gamma < 1/(2 + \delta)$  and  $\mathbb{E}(\varepsilon_-^{2+\delta}) < \infty$ , then conditions  $\mathcal{H}_\delta$ ,  $\mathcal{L}_m$  and  $\mathcal{B}_m$  hold as well. If in addition  $g$ ,  $m$  and  $\sigma$  are twice continuously differentiable in a neighborhood of  $\mathbf{x}$ , then conditions  $\mathcal{D}_g$  and  $\mathcal{D}_m$  also hold. If moreover  $f_\varepsilon$  is continuously differentiable in a neighborhood of infinity and satisfies the second-order von Mises condition  $-z f'_\varepsilon(z)/f_\varepsilon(z) \rightarrow 1/\gamma + 1$ , then condition  $\mathcal{D}_\omega$  holds, with  $\lim_{y \rightarrow \infty} \nabla_{\mathbf{x}} \log \bar{F}(y|\mathbf{x})/\log(y) = 0$ .

The class of  $\beta$ -mixing processes covers many important cases such as Harris recurrent aperiodic Markov chains [4, Corollary 3.6]. Our conditions on  $m$  and  $\sigma$  cover all standard location-scale regression models, such as heteroscedastic linear regression with smooth conditional variance function, single-index models with smooth link functions, and additive models. The assumptions on  $f_\varepsilon$  hold in the vast majority of standard heavy-tailed models, such as the Fréchet, Burr, Student and Fisher distributions. We note that this setting cannot cover time series models such as autoregressive models, which require a specific treatment, because a key assumption is that the series  $(\mathbf{X}_t)$  must be independent from the series of errors  $(\varepsilon_t)$ .

*Nonlinear regression model.* Let  $F(\cdot, \boldsymbol{\theta})$  be a parametric family of heavy-tailed distribution functions on  $\mathbb{R}$ , where  $\boldsymbol{\theta} \in \Theta \subset \mathbb{R}^d$  is a finite-dimensional, convex and open set of parameters, and let  $q(\cdot, \boldsymbol{\theta})$  be the associated quantile function (the left-continuous inverse of  $y \mapsto F(y, \boldsymbol{\theta})$ ). We abuse notation and let  $\boldsymbol{\theta} = \boldsymbol{\theta}(\cdot) : \mathbf{x} \in \mathbb{R}^p \mapsto \boldsymbol{\theta}(\mathbf{x}) \in \Theta$  be a smooth mapping, and we consider the model  $Y_t = q(U_t, \boldsymbol{\theta}(\mathbf{X}_t))$  where  $(U_t)$  is a stationary sequence of unobserved, uniformly distributed innovations independent from the series  $(\mathbf{X}_t)$ . Then  $Y$  given  $\mathbf{X} = \mathbf{x}$  has distribution function  $F(\cdot, \boldsymbol{\theta}(\mathbf{x}))$ : when  $\boldsymbol{\theta}(\cdot)$  is linear, this setting covers the (V)GLM model of [28] for a univariate response variable.

**PROPOSITION 4.2 (Nonlinear regression model).** *Assume that  $(\mathbf{X}_t)_{t \geq 1}$  is  $\beta$ -mixing and  $(U_t)_{t \geq 1}$  is strongly mixing. Suppose further that:*

- For any  $\boldsymbol{\theta}$ , the survival function  $\bar{F}(\cdot, \boldsymbol{\theta})$  is second-order regularly varying, with tail index  $\gamma = \underline{\gamma}(\boldsymbol{\theta})$ , second-order parameter  $\rho = \rho(\boldsymbol{\theta})$  and auxiliary function  $A = A(\cdot|\boldsymbol{\theta})$ . The function  $\boldsymbol{\theta} \mapsto \bar{F}(y, \boldsymbol{\theta})$  is continuously differentiable for  $y$  large enough.
- For any  $t \geq 1$ , the random pairs  $(\mathbf{X}_1, \mathbf{X}_{t+1})$  have absolutely continuous distributions whose first marginal has a p.d.f.  $g$  such that  $g(\mathbf{x}) > 0$  and  $g$  is continuously differentiable in a neighborhood of  $\mathbf{x}$ .
- The parameter mapping  $\boldsymbol{\theta} : \mathbb{R}^p \rightarrow \Theta$  is continuously differentiable in a neighborhood of  $\mathbf{x}$ . There are  $y_0 > 0$  and a continuous function  $\kappa$  on  $\Theta$  with

$$\forall \boldsymbol{\theta} \in \Theta, \sup_{y \geq y_0} \left\| \frac{\nabla_{\boldsymbol{\theta}} \log \bar{F}(y, \boldsymbol{\theta})}{\log(y)} \right\| = \sup_{y \geq y_0} \frac{1}{\log(y)} \frac{\|\nabla_{\boldsymbol{\theta}} \bar{F}(y, \boldsymbol{\theta})\|}{\bar{F}(y, \boldsymbol{\theta})} \leq \kappa(\boldsymbol{\theta}).$$

Then conditions  $\mathcal{M}$ ,  $\mathcal{L}_g$ ,  $\mathcal{L}_\omega$ ,  $\mathcal{B}_p$  and  $\mathcal{B}_\Omega$  hold. Assume moreover that:

- The survival function  $\bar{F}(\cdot, \boldsymbol{\theta})$  only puts mass on  $[0, \infty)$ .
- At the point  $\mathbf{x}$ ,  $\gamma(\boldsymbol{\theta}(\mathbf{x})) < 1/(2 + \delta)$ .
- The first and second moments of the distribution  $F(\cdot, \boldsymbol{\theta})$  define continuously differentiable functions  $m_1$  and  $m_2$  of  $\boldsymbol{\theta}$  wherever they are defined.

Then conditions  $\mathcal{H}_\delta$ ,  $\mathcal{L}_m$  and  $\mathcal{B}_m$  hold. If  $m_1$  and  $m_2$  are twice continuously differentiable functions of  $\boldsymbol{\theta}$  wherever they are defined and  $g$  and  $\boldsymbol{\theta}(\cdot)$  are twice continuously differentiable in a neighborhood of  $\boldsymbol{x}$ , then conditions  $\mathcal{D}_g$  and  $\mathcal{D}_m$  hold. If  $\boldsymbol{\theta} \mapsto \gamma(\boldsymbol{\theta})$  is also twice continuously differentiable on  $\Theta$  and there exist  $y_0 > 0$  and a continuous function  $\kappa$  on  $\Theta$  with

$$\forall \boldsymbol{\theta} \in \Theta, \sup_{y \geq y_0} \left\{ \frac{1}{\log(y)} (\|\nabla_{\boldsymbol{\theta}} \log \bar{F}(y, \boldsymbol{\theta})\| + \|H_{\boldsymbol{\theta}} \log \bar{F}(y, \boldsymbol{\theta})\|) \right\} \leq \kappa(\boldsymbol{\theta})$$

then condition  $\mathcal{D}_\omega$  holds, with  $\lim_{y \rightarrow \infty} \nabla_{\boldsymbol{x}} \log \bar{F}(y|\boldsymbol{x}) / \log(y) = [\nabla \gamma(\boldsymbol{\theta}(\boldsymbol{x}))]^\top J_{\boldsymbol{\theta}}(\boldsymbol{x}) / \gamma^2(\boldsymbol{\theta}(\boldsymbol{x}))$ .

Our conditions on the statistical model  $F(\cdot, \boldsymbol{\theta})$  are mild and readily checked in typical heavy-tailed models such as the Fréchet, Burr, Generalized Pareto and half- $t$  models. If  $\bar{F}(\cdot, \boldsymbol{\theta})$  puts mass on a neighborhood of  $-\infty$ , then extra assumptions on the left conditional tail (such as symmetry) are required to ensure that conditions  $\mathcal{H}_\delta$ ,  $\mathcal{L}_m$ ,  $\mathcal{B}_m$  and  $\mathcal{D}_m$  hold.

*Autoregressive model.* Consider the causal and invertible AR( $p$ ) model  $Y_t = \sum_{j=1}^p \phi_j Y_{t-j} + \varepsilon_t$ ,  $t \in \mathbb{Z}$ , where the polynomial  $P(z) = 1 - \sum_{j=1}^p \phi_j z^j$  has no root inside the closed unit disk in  $\mathbb{C}$ , and  $(\varepsilon_t)$  is an i.i.d. sequence of innovations. Here  $Y_t$  should be understood as the stationary solution of the AR( $p$ ) equations, and  $\mathbf{X}_t = (Y_{t-1}, Y_{t-2}, \dots, Y_{t-p})^\top$ .

**PROPOSITION 4.3 (Autoregressive model).** *Assume that the common distribution of the  $\varepsilon_t$  is centered, square-integrable, and has a Lipschitz continuous, everywhere strictly positive p.d.f.  $f_\varepsilon$  with respect to the Lebesgue measure that satisfies  $f_\varepsilon(z) = c_0 z^{-1/\gamma-1} (1 + d_0 z^{-a} + d'_0 z^{-a-b} (1 + o(1)))$  as  $z \rightarrow \infty$ , where  $\gamma > 0$ ,  $c_0 > 0$ ,  $d_0, d'_0 \neq 0$  and  $a, b > 0$  are such that either  $a \neq 1$ ,  $a = 1 \neq b$ , or  $a = b = 1$  with then  $2d'_0(1 + \gamma) \neq d_0^2(1 + 2\gamma)$ .*

*Then conditions  $\mathcal{M}$ ,  $\mathcal{L}_g$ ,  $\mathcal{D}_m$ ,  $\mathcal{L}_\omega$ ,  $\mathcal{B}_p$  and  $\mathcal{B}_\Omega$  hold, and the process  $(Y_t)$  is geometrically  $\beta$ - and  $\rho$ -mixing. If  $\varepsilon$  also has a finite moment of order  $(2 + \delta)$  and  $\gamma < 1/(2 + \delta)$ , then condition  $\mathcal{H}_\delta$  holds. If moreover  $f_\varepsilon$  is continuously differentiable, with a uniformly bounded Lipschitz continuous derivative  $f'_\varepsilon$ , then condition  $\mathcal{D}_g$  holds. Finally, if  $f_\varepsilon$  also satisfies the second-order von Mises condition  $-z f'_\varepsilon(z) / f_\varepsilon(z) \rightarrow 1/\gamma + 1$ , then condition  $\mathcal{D}_\omega$  holds, with  $\lim_{y \rightarrow \infty} \nabla_{\boldsymbol{x}} \log \bar{F}(y|\boldsymbol{x}) / \log(y) = 0$ .*

Condition  $\mathcal{B}_m$  is unnecessary because  $(Y_t)$  is  $\rho$ -mixing. Unlike in our other examples, checking condition  $\mathcal{B}_\Omega$  is nontrivial, because the sequences  $(\mathbf{X}_t) = ((Y_{t-1}, Y_{t-2}, \dots, Y_{t-p})^\top)$  and  $(\varepsilon_t)$  are not independent. This is done by noting that  $(Y_t)$  is a Markov chain of order  $p$  and then by checking the conditions of Lemma A.1(iii) of the Supplementary Material document [12] (with  $t_0 = p$ ). Handling models whose natural covariate is infinite-dimensional, such as the ARMA and GARCH models containing lagged values of the innovation, or moving maxima processes such as the m4 process of [23], in which joint distributions of pairs of responses at different time points are never absolutely continuous, should be viewed as substantially harder and worthy of future research.

**5. Practical implementation.** We discuss hyperparameter selection, bias and variance correction, and asymptotic confidence interval construction for both tail quantities of interest.

*5.1. Selection of tuning parameters.* The value  $k_n = n(1 - \tau_n)$ , rounded to the next integer, can be viewed as the effective sample size for tail extrapolation. A larger  $k_n$  leads to larger bias, while smaller  $k_n$  results in more variance. In the regression case, one should also determine the bandwidth  $h_n$ . Results in finite samples indicate that it is often enough

to use the global bandwidth obtained by minimizing the mean integrated squared error  $\mathbb{E}(\int_{\mathcal{X}}(\widehat{g}_n(\mathbf{x}) - g(\mathbf{x}))^2 d\mathbf{x})$  of the density estimator over the support  $\mathcal{X}$  of  $\mathbf{X}$ , that is,

$$(4) \quad h_{n,\star} = \left( \frac{p \int_{\mathbb{R}^p} K^2}{\int_{\mathcal{X}} \left( \int_{\mathbb{R}^p} (\mathbf{u}^\top H_g(\mathbf{x}) \mathbf{u}) K(\mathbf{u}) d\mathbf{u} \right)^2 d\mathbf{x}} \right)^{1/(p+4)} n^{-1/(p+4)}.$$

For example, when  $p = 1$ , the classical normal scale rule derived from (4), assuming a Gaussian p.d.f.  $g$  and a naive kernel  $K(u) = 1/2$  on  $[-1, 1]$  yields  $\widehat{h}_{n,\star} = (12\sqrt{\pi})^{1/5} \widehat{\sigma}_n n^{-1/5}$ , where  $\widehat{\sigma}_n$  is the empirical standard deviation of the  $X_t$ . We adopt this version in our examples with  $p = 1$  and abuse notation by denoting it  $h_{n,\star}$ .

We turn to a choice of  $k_n$  optimizing the bias-variance tradeoff in extreme value estimation: if the chosen tail index estimator has an asymptotic variance  $V(\mathbf{x})$  and an asymptotic bias component  $\lambda(\mathbf{x})B(\mathbf{x})$ , where  $\lambda(\mathbf{x}) = \lim_{n \rightarrow \infty} \sqrt{k_n h_{n,\star}^p} A(n/k_n | \mathbf{x})$ , we define

$$k_{n,\star}(\mathbf{x}) = \arg \min_{1 \leq k \leq n} \{ (k h_{n,\star}^p)^{-1} V(\mathbf{x}) + B^2(\mathbf{x}) A^2(n/k | \mathbf{x}) \}.$$

A common practice is to consider the very general case when  $A(t | \mathbf{x}) = b(\mathbf{x}) \gamma(\mathbf{x}) t^{\rho(\mathbf{x})}$ , for  $b(\mathbf{x}) \in \mathbb{R}$  and  $\rho(\mathbf{x}) < 0$ , see among others [19]. This yields the closed form expression

$$k_{n,\star}(\mathbf{x}) = \left( \frac{1}{-2\rho(\mathbf{x}) b^2(\mathbf{x}) \gamma^2(\mathbf{x})} \frac{V(\mathbf{x})}{B^2(\mathbf{x})} \right)^{1/(1-2\rho(\mathbf{x}))} h_{n,\star}^{-p/(1-2\rho(\mathbf{x}))} n^{-2\rho(\mathbf{x})/(1-2\rho(\mathbf{x}))}.$$

The quantities  $b(\mathbf{x})$  and  $\rho(\mathbf{x})$  are estimated using naive kernel regression versions  $\bar{b}(\mathbf{x})$  and  $\bar{\rho}(\mathbf{x})$  of the moment estimators provided by the R function `mop` from the `Expectrem` package, see Section C.1 of the Supplementary Material document [12]. To estimate  $B(\mathbf{x})$  and  $V(\mathbf{x})$ , note that when the tail index estimator is chosen as  $\widehat{\gamma}_{1-k_n/n}^{(J)}(\mathbf{x})$ , we have  $V(\mathbf{x}) = (\int_{\mathbb{R}^p} K^2/g(\mathbf{x})) v_q^{(J)}(\mathbf{x})$  and  $B(\mathbf{x}) = B^{(J)}(\mathbf{x})$ , with

$$v_q^{(J)}(\mathbf{x}) = \gamma^2(\mathbf{x}) \frac{J(J-1)(2J-1)}{6 \log^2(J!)} \quad \text{and} \quad B^{(J)}(\mathbf{x}) = \frac{1}{\log(J!)} \sum_{j=1}^J \frac{j^{\rho(\mathbf{x})} - 1}{\rho(\mathbf{x})}.$$

When  $\widehat{\gamma}_{1-k_n/n}^E(\mathbf{x})$  is chosen, then  $V(\mathbf{x}) = (\int_{\mathbb{R}^p} K^2/g(\mathbf{x})) v_E(\mathbf{x})$  and  $B(\mathbf{x}) = B_E(\mathbf{x})$  where

$$v_E(\mathbf{x}) = \frac{\gamma^3(\mathbf{x})(1-\gamma(\mathbf{x}))}{1-2\gamma(\mathbf{x})} \quad \text{and} \quad B_E(\mathbf{x}) = \frac{\gamma(\mathbf{x})(1/\gamma(\mathbf{x})-1)^{1-\rho(\mathbf{x})}}{1-\gamma(\mathbf{x})-\rho(\mathbf{x})}.$$

These quantities are estimated by plugging in  $\widehat{g}_n(\mathbf{x})$  (with  $h = h_{n,\star}$ ) and  $\bar{\rho}(\mathbf{x})$  in place of  $g(\mathbf{x})$  and  $\rho(\mathbf{x})$ , and a naive kernel regression version  $\bar{\gamma}(\mathbf{x})$  of the Hill estimator in place of  $\gamma(\mathbf{x})$ , with  $h = h_{n,\star}$  and effective sample size  $k$  corresponding to 25% of the local sample size  $N_h(\mathbf{x}) = \sum_{t=1}^n \mathbb{1}_{\{\|\mathbf{X}_t - \mathbf{x}\| \leq h\}}$  (see Section C.1 of [12] for the expression of  $\bar{\gamma}(\mathbf{x})$ ).

This results in data-driven choices  $\widehat{k}_{n,\star}^{(J)}(\mathbf{x})$  and  $\widehat{k}_{n,\star}^E(\mathbf{x})$ , depending on the tail index estimation technique. We omit the dependence of this selected value upon the tail index estimator and we denote it again by  $k_{n,\star}$  for the sake of brevity. Our choices may not be optimal in certain difficult cases, but they afford effective data-based rules on our simulated and real data.

**5.2. Bias correction guidelines.** The quality of the pure Pareto approximation to conditional heavy tails deteriorates as  $\rho(\mathbf{x})$  gets closer to 0, with the resulting bias being possibly very substantial then. Bias-reduced versions of  $\widehat{\gamma}_{1-k_n/n}^{(J)}(\mathbf{x})$  and  $\widehat{q}_{n,\tau_n}^W(\tau'_n | \mathbf{x})$  are

$$\widehat{\gamma}_{1-k_n/n}^{(J, BR)}(\mathbf{x}) = \widehat{\gamma}_{1-k_n/n}^{(J)}(\mathbf{x}) \left( 1 - \frac{1}{\log(J!)} \sum_{j=1}^J \frac{j^{\bar{\rho}(\mathbf{x})} - 1}{\bar{\rho}(\mathbf{x})} \bar{b}(\mathbf{x}) \left( \frac{n}{k_n} \right)^{\bar{\rho}(\mathbf{x})} \right)$$



$$\text{and } \hat{q}_{n,\tau_n}^{W,BR}(\tau'_n|\mathbf{x}) = \hat{q}_{n,\tau_n}^W(\tau'_n|\mathbf{x}) \left( 1 + \frac{\left(\frac{k_n}{n(1-\tau'_n)}\right)^{\bar{\rho}(\mathbf{x})} - 1}{\bar{\rho}(\mathbf{x})} \bar{b}(\mathbf{x}) \hat{\gamma}_{1-k_n/n}^{(J,BR)}(\mathbf{x}) \left(\frac{n}{k_n}\right)^{\bar{\rho}(\mathbf{x})} \right).$$

These estimators are respectively inspired by Theorem 2.3 and by the approach of [19]. We use  $\hat{\gamma}_{1-k_n/n}^{(J,BR)}(\mathbf{x})$  within  $\hat{q}_{n,\tau_n}^W(\tau'_n|\mathbf{x})$  in our subsequent implementation of  $\hat{q}_{n,\tau_n}^W(\tau'_n|\mathbf{x})$ .

As for extreme conditional expectile-based estimation, we note that bias reduction is even more crucial, due to the fact that, in view of Theorems 3.3 and 3.4, the estimators contain two sources of bias: one coming from the second-order regular variation framework, and the other stemming from the asymptotic proportionality between quantiles and expectiles. We suggest using the following bias-reduced version of  $\hat{\gamma}_{1-k_n/n}^E(\mathbf{x})$ :

$$\begin{aligned} \hat{\gamma}_{1-k_n/n}^{E,BR}(\mathbf{x}) &= \left( 1 + \frac{\hat{F}_n(\hat{e}_n(1-k_n/n|\mathbf{x})|\mathbf{x})}{k_n/n} \frac{1}{1 + \hat{r}(1-k_n/n|\mathbf{x})} \right)^{-1} \text{ where } \hat{r}(1-k_n/n|\mathbf{x}) \\ &= \left( 1 - \frac{\hat{m}_n(\mathbf{x})}{\hat{e}_n(1-k_n/n|\mathbf{x})} \right) \frac{1}{1-2k_n/n} \left( 1 + \frac{\bar{b}(\mathbf{x})[\hat{F}_n(\hat{e}_n(1-k_n/n|\mathbf{x})|\mathbf{x})]^{-\bar{\rho}(\mathbf{x})}}{1 - \hat{\gamma}_{1-k_n/n}^E(\mathbf{x}) - \bar{\rho}(\mathbf{x})} \right)^{-1} - 1. \end{aligned}$$

This is a kernel regression version of the work in [18] in the unconditional setting. Here  $\hat{\gamma}_{1-k_n/n}^{E,BR}(\mathbf{x})$  is computed using the R function `tindexp` with argument `br=TRUE` (from the R package `Expectrem`) applied to the  $Y_t$  such that  $\|\mathbf{X}_t - \mathbf{x}\| \leq h_{n,*}$ . Bias-reduced versions  $\hat{e}_{n,\tau_n}^{W,BR}(\tau'_n|\mathbf{x})$  and  $\check{e}_{n,\tau_n}^{W,BR}(\tau'_n|\mathbf{x})$  of the extreme conditional expectile estimators  $\hat{e}_{n,\tau_n}^W(\tau'_n|\mathbf{x})$  (extrapolating  $\hat{e}_n(\tau_n|\mathbf{x})$ ) and  $\check{e}_{n,\tau_n}^W(\tau'_n|\mathbf{x})$  (extrapolating  $\check{e}_n(\tau_n|\mathbf{x})$ ), computed using the R function `extExpect` with `br=TRUE` from that package, are described in Section C.2 of the Supplementary Material document [12].

### 5.3. Pointwise asymptotic confidence intervals.

5.3.1. *Extremal conditional quantiles.* We use the equivalent version of Theorem 2.2 on the log-scale, which tends to be more accurate in practice. Combined with Theorem 2.3, this suggests that  $\log(\hat{q}_{n,1-k_{n,*}/n}^{W,BR}(\tau'_n|\mathbf{x})/q(\tau'_n|\mathbf{x}))$  is approximately Gaussian centered with variance  $(\int_{\mathbb{R}^p} K^2/g(\mathbf{x}))v_q^{(J)}(\mathbf{x})$ , with  $v_q^{(J)}(\mathbf{x}) = \gamma^2(\mathbf{x}) \times J(J-1)(2J-1)/(6\log^2(J!))$ . A first 95% asymptotic Gaussian confidence interval for  $q(\tau'_n|\mathbf{x})$  is then

$$\begin{aligned} \hat{I}_{q,1}(\tau'_n|\mathbf{x}) &= \left[ \hat{q}_{n,1-k_{n,*}/n}^{W,BR}(\tau'_n|\mathbf{x}) \exp \left( \pm \frac{\sqrt{\int_{\mathbb{R}^p} K^2 \hat{v}_q^{(J)}(\mathbf{x})}}{\sqrt{k_{n,*} h_{n,*}^p}} \log \left( \frac{k_{n,*}}{n(1-\tau'_n)} \right) z_{0.975} \right) \right] \\ &\text{with } \hat{v}_q^{(J)}(\mathbf{x}) = \frac{J(J-1)(2J-1)}{6\log^2(J!)} (\hat{\gamma}_{1-k_{n,*}/n}^{(J,BR)}(\mathbf{x}))^2, \end{aligned}$$

and where  $z_\tau$  is the  $\tau$ th quantile of the standard normal distribution. This relies exclusively on the asymptotic distribution of the tail index estimator used in the extrapolation. For small sample sizes, the variability of the intermediate quantile estimator also has an impact on the variance of the extrapolated estimator. By Theorem 2.3,  $\hat{\gamma}_{\tau_n}^{(J)}(\mathbf{x})$  is asymptotically independent of  $\hat{q}_n(\tau_n|\mathbf{x})$ , so a straightforward calculation provides a refined version  $\tilde{v}_q^{(J)}(\mathbf{x})$  of the asymptotic variance estimate of  $\log(\hat{q}_{n,1-k_{n,*}/n}^{W,BR}(\tau'_n|\mathbf{x})/q(\tau'_n|\mathbf{x}))$ , and thus a corrected 95%

asymptotic Gaussian confidence interval as

$$\widehat{I}_{q,2}(\tau'_n|\mathbf{x}) = \left[ \widehat{q}_{n,1-k_{n,\star}/n}^{W,BR}(\tau'_n|\mathbf{x}) \exp \left( \pm \frac{\sqrt{\int_{\mathbb{R}^p} K^2 \widehat{v}_q^{(J)}(\mathbf{x})}}{\sqrt{k_{n,\star} h_{n,\star}^p}} \log \left( \frac{k_{n,\star}}{n(1-\tau'_n)} \right) z_{0.975} \right) \right],$$

$$\text{with } \widehat{v}_q^{(J)}(\mathbf{x}) = \left( \frac{J(J-1)(2J-1)}{6 \log^2(J!)} + \frac{1}{\log^2(k_{n,\star}/[n(1-\tau'_n)])} \right) (\widehat{\gamma}_{1-k_{n,\star}/n}^{(J,BR)}(\mathbf{x}))^2.$$

Simulation evidence shows that this correction improves coverage for low sample sizes.

**5.3.2. Extremal conditional expectiles.** Our extreme conditional expectile estimators are constructed using the same Weissman extrapolation argument: we build upon Theorems 3.3 and 3.4 to deduce a first 95% asymptotic Gaussian confidence interval for  $e(\tau'_n|\mathbf{x})$  as

$$\widehat{I}_{E,1}(\tau'_n|\mathbf{x}) = \left[ \widehat{e}_{n,1-k_{n,\star}/n}^{W,BR}(\tau'_n|\mathbf{x}) \exp \left( \pm \frac{\sqrt{\int_{\mathbb{R}^p} K^2 \widehat{v}_E(\mathbf{x})}}{\sqrt{k_{n,\star} h_{n,\star}^p}} \log \left( \frac{k_{n,\star}}{n(1-\tau'_n)} \right) z_{0.975} \right) \right]$$

$$\text{with } \widehat{v}_E(\mathbf{x}) = \frac{(\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}))^3 (1 - \widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}))}{1 - 2\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x})}.$$

For a low-to-moderate sample size  $n$ , the empirical variance of the estimates tends to be very far from the asymptotic variance. This is not only due to neglecting the correlation between the conditional tail index estimator and the empirical expectiles, but also to the use of the asymptotic proportionality between extreme quantiles and expectiles in the derivation of the asymptotic results. Calculating the errors incurred in using these two approximations requires, first of all, an accurate quantification of the variance matrix of the random vector

$$\sqrt{k_{n,\star} h_{n,\star}^p} \left( \frac{n\widehat{F}_n(\widehat{e}_n(1-k_{n,\star}/n|\mathbf{x})|\mathbf{x})}{k_{n,\star}} - \left( \frac{1}{\gamma(\mathbf{x})} - 1 \right), \frac{\widehat{e}_n(1-k_{n,\star}/n|\mathbf{x})}{e(1-k_{n,\star}/n|\mathbf{x})} - 1 \right)$$

since, up to order  $1/\sqrt{k_{n,\star} h_{n,\star}^p}$ , we have  $\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}) \approx \mathcal{G}(n\widehat{F}_n(\widehat{e}_n(1-k_{n,\star}/n|\mathbf{x})|\mathbf{x})/k_{n,\star})$ , where  $\mathcal{G}(u) = 1/(1+u)$ . An inspection of the proofs of Theorems 3.1 and 3.4 (see Section C.3 of the Supplementary Material document [12]) suggests that a corrected asymptotic variance matrix for this random vector is  $(\int_{\mathbb{R}^p} K^2/g(\mathbf{x}))\mathbf{T}_n(\mathbf{x})$  where the  $2 \times 2$  symmetric matrix  $\mathbf{T}_n(\mathbf{x})$  has components

$$T_{n,11}(\mathbf{x}) = \frac{2(1-\gamma(\mathbf{x}))^2}{\gamma(\mathbf{x})(1-2\gamma(\mathbf{x}))} \frac{\kappa_{1,n}(\mathbf{x})}{\kappa_{2,n}^2(\mathbf{x})} - 2 \frac{1-\gamma(\mathbf{x})}{\gamma(\mathbf{x})} \frac{\sqrt{\kappa_{1,n}(\mathbf{x})}}{\kappa_{2,n}(\mathbf{x})} + \frac{1-\gamma(\mathbf{x})}{\gamma(\mathbf{x})},$$

$$T_{n,12}(\mathbf{x}) = -\frac{2\gamma(\mathbf{x})(1-\gamma(\mathbf{x}))}{1-2\gamma(\mathbf{x})} \frac{\kappa_{1,n}(\mathbf{x})}{\kappa_{2,n}(\mathbf{x})} + \gamma(\mathbf{x})\sqrt{\kappa_{1,n}(\mathbf{x})}, \quad T_{n,22}(\mathbf{x}) = \frac{2\gamma^3(\mathbf{x})}{1-2\gamma(\mathbf{x})} \kappa_{1,n}(\mathbf{x}),$$

$$\text{with } \kappa_{1,n}(\mathbf{x}) = \frac{1-2k_{n,\star}/n}{1-m(\mathbf{x})/e(1-k_{n,\star}/n|\mathbf{x})} \quad \text{and} \quad \kappa_{2,n}(\mathbf{x}) = 1 - \frac{\gamma(\mathbf{x})m(\mathbf{x})}{e(1-k_{n,\star}/n|\mathbf{x})}.$$

Replacing  $\gamma(\mathbf{x})$ ,  $m(\mathbf{x})$  and  $e(1-k_{n,\star}/n|\mathbf{x})$  with  $\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x})$ ,  $\widehat{m}_n(\mathbf{x})$  and  $\widehat{e}_n(1-k_{n,\star}/n|\mathbf{x})$  yields estimators  $\widehat{\kappa}_{1,n}(\mathbf{x})$  and  $\widehat{\kappa}_{2,n}(\mathbf{x})$  and therefore  $\widehat{T}_{n,11}(\mathbf{x})$ ,  $\widehat{T}_{n,12}(\mathbf{x})$  and  $\widehat{T}_{n,22}(\mathbf{x})$ . Then

$$\sqrt{k_{n,\star} h_{n,\star}^p} (\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}) - \gamma(\mathbf{x}))$$

$$\begin{aligned} &\approx \sqrt{k_{n,\star} h_{n,\star}^p} \left( \mathcal{G}(n\widehat{F}_n(\widehat{e}_n(1 - k_{n,\star}/n|\mathbf{x})|\mathbf{x})/k_{n,\star}) - \mathcal{G}(1/\gamma(\mathbf{x}) - 1) \right) \\ &= \sum_{k=1}^{\infty} \frac{(-1)^k \gamma^{k+1}(\mathbf{x})}{(k_{n,\star} h_{n,\star}^p)^{(k-1)/2}} \left( \sqrt{k_{n,\star} h_{n,\star}^p} \left( \frac{n\widehat{F}_n(\widehat{e}_n(1 - k_{n,\star}/n|\mathbf{x})|\mathbf{x})}{k_{n,\star}} - \left( \frac{1}{\gamma(\mathbf{x})} - 1 \right) \right) \right)^k \end{aligned}$$

through a power series expansion. [This is conceptually simpler than an Edgeworth expansion, which would approximate the p.d.f. of  $\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x})$ .] Taking only the first term above leads to the delta-method for  $\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x})$ ; to obtain a finer approximation, we use all terms up to order 4. Based on our calculation of  $\mathbf{T}_n(\mathbf{x})$ , an asymptotic approximation of the covariance matrix of  $\sqrt{k_{n,\star} h_{n,\star}^p}(\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}) - \gamma(\mathbf{x}), \widehat{e}_n(1 - k_{n,\star}/n|\mathbf{x})/e(1 - k_{n,\star}/n|\mathbf{x}) - 1)$  is  $(\int_{\mathbb{R}^p} K^2/g(\mathbf{x}))\mathbf{S}_n(\mathbf{x})$ , where the symmetric matrix  $\mathbf{S}_n(\mathbf{x})$  has components

$$\begin{aligned} S_{n,11}(\mathbf{x}) &= \gamma^4(\mathbf{x})T_{n,11}(\mathbf{x}) \left( 1 + 8 \frac{\gamma^2(\mathbf{x})T_{n,11}(\mathbf{x})}{k_{n,\star} h_{n,\star}} \times \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} \right), \\ S_{n,12}(\mathbf{x}) &= -\gamma^2(\mathbf{x})T_{n,12}(\mathbf{x}) \left( 1 + 3 \frac{\gamma^2(\mathbf{x})T_{n,11}(\mathbf{x})}{k_{n,\star} h_{n,\star}} \times \frac{\int_{\mathbb{R}^p} K^2}{g(\mathbf{x})} \right) \text{ and } S_{n,22}(\mathbf{x}) = T_{n,22}(\mathbf{x}). \end{aligned}$$

We denote by  $\widehat{\mathbf{S}}_n(\mathbf{x})$  the associated estimator, which hinges upon the estimators  $\widehat{\mathbf{T}}_n(\mathbf{x})$  and  $\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x})$  previously introduced. The final step is to recall that (see Section C.2 of [12])

$$\begin{aligned} \log \frac{\widehat{e}_{n,1-k_{n,\star}/n}^{W,BR}(\tau'_n|\mathbf{x})}{e(\tau'_n|\mathbf{x})} &\approx \log \frac{\widehat{e}_{n,1-k_{n,\star}/n}^W(\tau'_n|\mathbf{x})}{e(\tau'_n|\mathbf{x})} - (\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}) - \gamma(\mathbf{x})) \log(\kappa_{1,n}(\mathbf{x})) \\ &\approx \left( \log \left( \frac{k_{n,\star}}{n(1 - \tau'_n)} \right) - \log(\kappa_{1,n}(\mathbf{x})) \right) (\widehat{\gamma}_{1-k_{n,\star}/n}^{E,BR}(\mathbf{x}) - \gamma(\mathbf{x})) + \log \frac{\widehat{e}_n(1 - k_{n,\star}/n|\mathbf{x})}{e(1 - k_{n,\star}/n|\mathbf{x})}. \end{aligned}$$

The variance of  $(\sqrt{k_{n,\star} h_{n,\star}^p} / \log(k_{n,\star}/(n(1 - \tau'_n)))) \times \log(\widehat{e}_{n,1-k_{n,\star}/n}^{W,BR}(\tau'_n|\mathbf{x})/e(\tau'_n|\mathbf{x}))$  is then estimated by  $(\int_{\mathbb{R}^p} K^2/\widehat{g}_n(\mathbf{x}))\widetilde{v}_E(\mathbf{x})$ , where

$$\widetilde{v}_E(\mathbf{x}) = \frac{\widehat{\mathbf{S}}_{n,11}(\mathbf{x})\widehat{\mathcal{L}}_n^2 + 2\widehat{\mathbf{S}}_{n,12}(\mathbf{x})\widehat{\mathcal{L}}_n + \widehat{\mathbf{S}}_{n,22}(\mathbf{x})}{\log^2(k_{n,\star}/(n(1 - \tau'_n)))} \text{ and } \widehat{\mathcal{L}}_n = \log(k_{n,\star}/(n(1 - \tau'_n)\widehat{\kappa}_{1,n}(\mathbf{x}))).$$

This results in the corrected 95% asymptotic Gaussian confidence interval

$$\widehat{I}_{E,2}(\tau'_n|\mathbf{x}) = \left[ \widehat{e}_{n,1-k_{n,\star}/n}^{W,BR}(\tau'_n|\mathbf{x}) \exp \left( \pm \frac{\sqrt{\int_{\mathbb{R}^p} K^2/\widehat{g}_n(\mathbf{x})} \widetilde{v}_E(\mathbf{x})}{\sqrt{k_{n,\star} h_{n,\star}^p}} \log \left( \frac{k_{n,\star}}{n(1 - \tau'_n)} \right) z_{0.975} \right) \right].$$

It is shown below to have a greatly improved coverage probability compared to  $\widehat{I}_{E,1}(\tau'_n|\mathbf{x})$ .

## 6. Simulation study.

6.1. *Models and setup.* For the sake of brevity, we only report here results in the non-linear (Burr) regression model, with covariate dimension  $p = 1$ . Other cases spanning our list of worked-out examples, in dimensions  $p = 1$  and 2, are presented in Section C.4 of the Supplementary Material document [12].

We consider a nonlinear Burr process  $Y_t = ((1 - U_t)^{\rho(X_t)} - 1)^{-\gamma(X_t)/\rho(X_t)}$  where:

- $X_t = \Phi(Z_t)$ , where  $\Phi$  is the standard Gaussian distribution function and  $(Z_t)$  (simulated using the `garch.sim` routine from the R package TSA) is a GARCH(1,1) process with  $\omega = 0.25$ ,  $\alpha = 0.75$ ,  $\beta = 0.2$ , i.e.  $Z_{t+1} = \Sigma_{t+1}\delta_{t+1}$ , where the  $\delta_t$  are i.i.d. standard Gaussian and  $\Sigma_{t+1}$  is defined recursively as  $\Sigma_{t+1} = (\omega + \alpha Z_t^2 + \beta \Sigma_t^2)^{1/2}$ .

- $(U_t)$  is defined recursively as  $U_0 \sim \text{Uniform}[0, 1]$  and, for  $t \geq 1$ ,  $U_t = \frac{1}{r}U_{t-1} + \eta_t$ , where the  $\eta_t$  are i.i.d. uniformly drawn over  $\{0, 1/r, \dots, (r-1)/r\}$ , and  $r = 5$ . The  $U_t$  are standard uniform and  $\alpha$ -mixing, see Section C.4 of [12] for further details.

We fix  $\rho(x) = -1$  for all  $x \in [0, 1]$  and consider three different models for  $\gamma(x)$ ,  $x \in [0, 1]$ :

- (NL-P) The polynomial model  $\gamma(x) = 0.15 + 0.5x(1-x)$ ;
- (NL-S) The sinusoidal model  $\gamma(x) = 0.2 + 0.05 \sin(2\pi x)$ ;
- (NL-C) The constant model  $\gamma(x) = 0.2$ .

In these three cases,  $\gamma(x) \in (0, 1/2)$  for any  $x \in [0, 1]$ . The true value of the conditional quantile is  $q(\tau|x) = ((1-\tau)^{-1} - 1)^{\gamma(x)}$ . The theoretical conditional expectile  $e(\tau|x)$  is computed numerically using the function `eburr` in the R package `Expectrem`.

We simulate  $N = 1,000$  replications of size  $n = 10,000$  of each of these models (results for  $n \in \{1,000, 5,000\}$  are found in Section C.4 of [12]). We estimate extreme conditional quantiles and expectiles at level  $\tau'_n = 1 - 10/n = 0.999$ . We let  $K(u) = 0.5 \mathbb{1}_{\{|u| \leq 1\}}$  be the uniform kernel, we select  $h_n = h_{n,\star}$  and  $\tau_n = \tau_{n,\star} = 1 - k_{n,\star}/n$  as in Section 5.1, and take  $J = 9$  in the conditional tail index estimator  $\hat{\gamma}_{1-k_{n,\star}/n}^{(J, BR)}(x)$ . Our quantile (resp. expectile) estimators are compared with the non-extrapolated quantile estimator  $\hat{q}_n(\tau'_n|x)$  (resp. the non-extrapolated expectile estimator  $\hat{e}_n(\tau'_n|x)$ ) and the simple extrapolated version  $\hat{q}_{n,1-k_{n,\star}/n}^W(\tau'_n|x)$  based on  $\hat{\gamma}_{1-k_{n,\star}/n}^{(J)}(x)$  (resp. the simple extrapolated version  $\hat{e}_{n,1-k_{n,\star}/n}^W(\tau'_n|x)$  based on  $\hat{\gamma}_{1-k_{n,\star}/n}^E(x)$ ). Our proposed 95% asymptotic Gaussian confidence intervals for  $q(\tau'_n|x)$  are compared with the following competing intervals:

$$\hat{I}_{q,3}(\tau'_n|\mathbf{x}) = \left[ \hat{q}_{n,1-k_{n,\star}/n}^W(\tau'_n|\mathbf{x}) \exp \left( \pm \frac{\sqrt{\frac{\int_{\mathbb{R}^p} K^2}{\hat{g}_n(\mathbf{x})} \bar{v}_q^{(J)}(\mathbf{x})}}{\sqrt{k_{n,\star} h_{n,\star}^p}} \log \left( \frac{k_{n,\star}}{n(1-\tau'_n)} \right) z_{0.975} \right) \right],$$

with  $\bar{v}_q^{(J)}(\mathbf{x}) = \frac{J(J-1)(2J-1)}{6 \log^2(J!)} (\hat{\gamma}_{1-k_{n,\star}/n}^{(J)}(\mathbf{x}))^2$ , similar to  $\hat{I}_{q,1}(\tau'_n|\mathbf{x})$  (without correction), and

$$\hat{I}_{q,4}(\tau'_n|\mathbf{x}) = \left[ \hat{q}_n(\tau'_n|\mathbf{x}) \exp \left( \pm \hat{\gamma}_{1-k_{n,\star}/n}^{(J)}(\mathbf{x}) \frac{\sqrt{\frac{\int_{\mathbb{R}^p} K^2}{\hat{g}_n(\mathbf{x})}}}{\sqrt{n h_{n,\star}^p (1-\tau'_n)}} z_{0.975} \right) \right],$$

based on an application of Theorem 2.1 written on the log-scale, which does not feature any extrapolation at all. Likewise, for  $e(\tau'_n|\mathbf{x})$ , we compute

$$\hat{I}_{E,3}(\tau'_n|\mathbf{x}) = \left[ \hat{e}_{n,1-k_{n,\star}/n}^W(\tau'_n|\mathbf{x}) \exp \left( \pm \frac{\sqrt{\frac{\int_{\mathbb{R}^p} K^2}{\hat{g}_n(\mathbf{x})} \bar{v}_E(\mathbf{x})}}{\sqrt{k_{n,\star} h_{n,\star}^p}} \log \left( \frac{k_{n,\star}}{n(1-\tau'_n)} \right) z_{0.975} \right) \right],$$

with  $\bar{v}_E(\mathbf{x}) = \frac{(\hat{\gamma}_{1-k_{n,\star}/n}^E(\mathbf{x}))^3 (1 - \hat{\gamma}_{1-k_{n,\star}/n}^E(\mathbf{x}))}{1 - 2\hat{\gamma}_{1-k_{n,\star}/n}^E(\mathbf{x})}$ , similar to  $\hat{I}_{E,1}(\tau'_n|\mathbf{x})$  (without correction), and

$$\hat{I}_{E,4}(\tau'_n|\mathbf{x}) = \left[ \hat{e}_n(\tau'_n|\mathbf{x}) \exp \left( \pm \sqrt{\frac{2(\hat{\gamma}_{1-k_{n,\star}/n}^E(\mathbf{x}))^3}{1 - 2\hat{\gamma}_{1-k_{n,\star}/n}^E(\mathbf{x})}} \times \frac{\sqrt{\frac{\int_{\mathbb{R}^p} K^2}{\hat{g}_n(\mathbf{x})}}}{\sqrt{n h_{n,\star}^p (1-\tau'_n)}} z_{0.975} \right) \right],$$

suggested by an equivalent version of Theorem 3.1 using the log-scale. This will make it possible to assess the benefits of the extrapolation procedure and bias correction scheme.

6.2. *Results.* We represent in Figure 3 boxplots of the extreme conditional quantile and asymmetric least squares expectile estimates, as well as the coverage probabilities of the intervals  $\widehat{I}_{q,j}(\tau'_n|x)$  and  $\widehat{I}_{E,j}(\tau'_n|x)$ , for  $1 \leq j \leq 4$ , in models (NL-P), (NL-S) and (NL-C). Quantitative results at selected values of  $x$  are provided in Tables C.1, C.2 and C.3 in Section C.4 of the Supplementary Material document [12]. We discuss the conclusions from a full set of results for sample sizes  $n \in \{1,000, 5,000, 10,000\}$  in an expanded list of models found in Section C.4 of [12], see Figures C.1–C.8 therein. Comparing quantile-based expectile estimates with their bias-reduced asymmetric least squares counterparts (see Figure C.1) reveals that  $\widehat{e}_{n,\tau_n}^{W,BR}(\tau'_n|x)$  appears to have the smallest bias overall, although it has a slightly higher variance than the extrapolated, bias-reduced quantile-based estimator  $\check{e}_{n,\tau_n}^{W,BR}(\tau'_n|x)$ . This makes it difficult to decide which one of the asymmetric least squares or quantile-based expectile estimators performs best. We note that the main advantage of the quantile-based approach is to not rely upon the finite second moment assumption. Since our simulation setting considers models where the conditional second moment is finite, we focus throughout this section on the proposed extreme conditional quantile estimators as well as the asymmetric least squares conditional expectile estimators and their associated confidence intervals.

Even though empirical point estimates without extrapolation seem at first glance to perform respectably, their variance is substantially larger than that of extrapolated estimates, and most importantly the lack of extrapolation makes accurate inference impossible. This is obvious from the coverage probabilities of  $\widehat{I}_{E,4}(\tau'_n|x)$ , often very close to 1. While bias reduction and variance correction for extrapolated estimates are arguably not crucial for quantile estimation at large values of  $n$ , and in fact bias correction can slightly deteriorate coverage probabilities of the Gaussian confidence interval, bias reduction is however valuable when the second-order parameter  $\rho(x)$  gets close to 0 (see the top panels of Figure C.6), and variance correction brings noticeable improvements when  $n$  is moderately large ( $n = 1,000$ ), with coverage probabilities typically improving by about 5% when using our proposed interval  $\widehat{I}_{q,2}(\tau'_n|x)$  instead of  $\widehat{I}_{q,1}(\tau'_n|x)$ , see the third row of Figures C.2–C.7. Bias and variance correction are, however, of prime importance when estimating and inferring extreme conditional expectiles: standard extrapolated estimates are heavily biased, with the asymptotic consistency property not visibly evidenced even for large  $n$ , see the second row of Figures C.2–C.7. This is due to the presence of a typically very large bias term proportional to the reciprocal of the extreme conditional quantile in the expectile extrapolation formula that is key for the asymptotic normality of extrapolated expectile estimators, see Lemma A.10 in Section A.4 of [12]. In addition, the uncorrected expression of the asymptotic variance of the extrapolated expectile estimates yields values very far from their sample variance, while our corrected proposal gets very close to this observed variance, see the rightmost panels in Figure 3. These are the main reasons why the intervals  $\widehat{I}_{E,1}(\tau'_n|x)$  and  $\widehat{I}_{E,3}(\tau'_n|x)$  have very poor coverage, while  $\widehat{I}_{E,2}(\tau'_n|x)$  achieves a coverage remarkably close to the nominal rate when  $n \geq 5,000$  and performs much better than its competitors when  $n = 1,000$ , see the third row of Figures C.2–C.7. A coverage close to the nominal level is also achieved in the more challenging dimension  $p = 2$ , see the third column of Figure C.8, although misspecified values of the second-order parameters  $\rho$  and  $b$  were used in this difficult situation in order to better handle the finite-sample variability of the resulting quantile and expectile estimates. It should be highlighted that in a substantial number of cases, the non-corrected variance estimate  $\widehat{v}_E(x)$  was not positive and therefore the associated confidence interval  $\widehat{I}_{E,3}(\tau'_n|x)$  could not be calculated; by contrast, the interval  $\widehat{I}_{E,2}(\tau'_n|x)$  was always well-defined. As evidenced in the bottom panels of Figures C.2–C.7, the median length of the interval  $\widehat{I}_{q,2}(\tau'_n|x)$  (resp.  $\widehat{I}_{E,2}(\tau'_n|x)$ ) was found to be close to that of  $\widehat{I}_{q,1}(\tau'_n|x)$  and  $\widehat{I}_{q,3}(\tau'_n|x)$  (resp.  $\widehat{I}_{E,1}(\tau'_n|x)$  and  $\widehat{I}_{E,3}(\tau'_n|x)$ ) and always substantially shorter than that of  $\widehat{I}_{q,4}(\tau'_n|x)$  (resp.  $\widehat{I}_{E,4}(\tau'_n|x)$ ),

thus reflecting the added value of the extrapolation procedure combined with the bias and variance correction schemes.

We conclude that the Gaussian asymptotic theory, when properly corrected in order to account for the presence of conditional heavy tails, provides reasonably accurate inference. It should not come as a surprise that good inferential results require a fairly large sample size, of the order of several thousands, due to the four main difficulties of the problem: (i) Temporal mixing, which impacts the trustworthiness of the asymptotic Gaussian limits, (ii) Nonparametric smoothing, inducing an approximation bias due to local variation of the underlying distribution, (iii) Sparsity in  $x$  due to the nonparametric regression framework, and (iv) Sparsity in  $y$  due to the extreme value context. Sparsity in  $x$  is a particular concern in the difficult model (AR), where the covariate, being a lagged response, takes unbounded values and where a correct inference can only be expected where observations  $Y_t$  concentrate.

## 7. Real data analysis.

**7.1. Stock returns data.** The first dataset, available from the R package HRW, contains the  $n = 2,363$  values of the excess daily stock returns (daily log-returns minus risk-free interest rate) on General Electric and the S&P 500 index, from 1 November 1993 to 31 March 2003. Let  $X_t$  (resp.  $Y_t$ ) be the negative of the excess daily log-return of the S&P 500 index (resp. General Electric), whose large values represent large losses. In a full-fledged data analysis one would of course be interested in understanding the extreme value behavior of  $Y_t$  given the full history of the S&P 500 index; in this illustrative example we focus on predicting the conditional extreme quantiles and expectiles of  $Y_t$  given  $X_t = x$ , a fixed value. About half of the data is negative, which is detrimental to the calculation of our bias-reduced estimators that require only positive values, see Section C.1 of [12]. We thus shift most of the data above 0 by subtracting to the  $Y_t$  their empirical unconditional 10% quantile, approximately equal to  $-0.023$ , *i.e.* we apply our procedures to  $Y'_t = Y_t + 0.023$  before shifting back to the original position.

As in our simulation study, we infer extremal regression quantiles and expectiles at level  $\tau'_n = 1 - 10/n \approx 0.995$  by following the methodology described in Section 5.1. We represent in the middle panel of Figure 1 (restricted to the interval  $X_t \in [-0.015, 0.015]$  which contains more than 80% of the data) our bias-reduced extrapolated estimates  $\hat{q}_{n,1-k_{n,*}/n}^{W,BR}(\tau'_n|x)$  and  $\hat{e}_{n,1-k_{n,*}/n}^{W,BR}(\tau'_n|x)$ , along with their respective asymptotic 95% confidence intervals  $\hat{I}_{q,2}(\tau'_n|x)$  and  $\hat{I}_{E,2}(\tau'_n|x)$ . The purely empirical estimates  $\hat{q}_n(\tau'_n|x)$  and  $\hat{e}_n(\tau'_n|x)$  are graphed in the left panel along with their respective asymptotic 95% confidence intervals  $\hat{I}_{q,4}(\tau'_n|x)$  and  $\hat{I}_{E,4}(\tau'_n|x)$ . In agreement to what was observed in Section 6.2, the intervals  $\hat{I}_{q,4}(\tau'_n|x)$  and  $\hat{I}_{E,4}(\tau'_n|x)$  are almost overall wider than  $\hat{I}_{q,2}(\tau'_n|x)$  and  $\hat{I}_{E,2}(\tau'_n|x)$ , respectively.

The estimated curves, confidence intervals, and regression mean all point towards a linear trend. This motivated us to perform a residual-based extreme value estimation from the ordinary linear regression model  $Y_t = -0.00030 + 1.24X_t + \varepsilon_t$ . We construct the residuals  $\hat{\varepsilon}_t$  and we calculate corresponding residual-based, bias-corrected extreme quantile estimates  $\hat{q}_{n,\tau_n,\varepsilon}(\tau'_n)$ , following [19], and bias-corrected extreme expectile estimates  $\hat{e}_{n,\tau_n,\varepsilon}(\tau'_n)$ , following [18], in which we chose  $\tau_n = 1 - 200/n = 0.915$  after a graphical inspection of the Hill plot of residuals, with associated asymptotic 95% Gaussian confidence intervals on the log-scale. This yields conditional extreme point quantile and expectile estimates of the  $Y_t$  as  $\bar{q}_{n,\tau_n}(\tau'_n|x) = -0.00030 + 1.24x + \hat{q}_{n,\tau_n,\varepsilon}(\tau'_n)$  and  $\bar{e}_{n,\tau_n}(\tau'_n|x) = -0.00030 + 1.24x + \hat{e}_{n,\tau_n,\varepsilon}(\tau'_n)$ , and their corresponding confidence intervals through the same linear translation, all of them graphed in the right panel of Figure 1. In this panel, we also represented

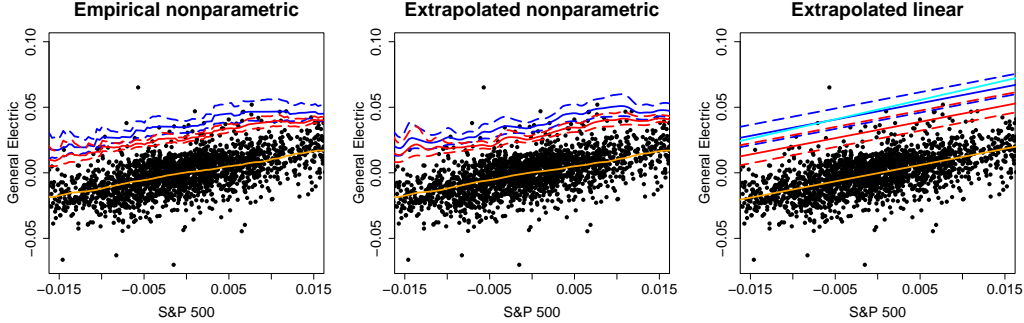


FIG 1. Stock market data. Blue line: Estimated conditional quantile at level  $\tau_n^l = 0.995$  (with 95% confidence interval in dashed line), red line: Estimated conditional expectile at level  $\tau_n^l = 0.995$  (with 95% confidence interval in dashed line). Left panel: Nonparametric empirical estimates, middle panel: Extrapolated nonparametric estimates, right panel: Extrapolated estimates based on the linear regression model. In the left and middle panels, the orange line is the Nadaraya-Watson estimate. In the right panel, the orange line is the ordinary least squares line, and the cyan line represents the extremal quantile regression estimate of [8].

the corresponding (extrapolated) extremal linear quantile regression estimate of [8], produced using their `extrap.rq` routine in which the preliminary estimate of the tail index is calculated using their `summary.rq.hill` routine at the intermediate level 0.95, chosen according to their recommendations. Although the middle and right panels yield a broadly similar message, the confidence intervals from the linear model are wider than our bias- and variance-corrected nonparametric confidence intervals. This is probably due to the unavoidable assumption of constant tail index in the linear model, whose validity is unclear here. The extremal linear regression estimates are thus inevitably driven by the few largest observations in the data cloud. By contrast, our nonparametric method is able to finely differentiate conditional extreme value behavior when  $x$  varies. Accurate inference on both extremal regression modes, without recourse to the *a priori* assumptions of linearity and common tail, is crucial in order to produce correct tail risk appraisal. This is especially important for conditional expectiles that typically result in more liberal assessments of risk than conditional quantiles because they satisfy the diversification principle, and that here appear indeed to induce less conservative risk measurements.

7.2. *Tornado losses data.* This dataset<sup>2</sup> records, for each tornado that has occurred in the United States between 1 January 2010 and 31 December 2019, the associated monetary loss (`loss`), its starting and ending latitude and longitude (`slat`, `slon`, `elat` and `elon`), and the length and width of the area traveled over by the tornado (`len` and `wid`). We focus on the loss per surface unit  $Y$  (in USD) in terms of the tornado’s (average) geographical location  $\mathbf{X} = (X_1, X_2)$ , that is,  $Y = \text{loss}/(\text{len} \times \text{wid})$ ,  $X_1 = (\text{slon} + \text{elon})/2$ , and  $X_2 = (\text{slat} + \text{elat})/2$ . To keep the analysis simple and illustrative of our results, we do not attempt to conduct a conditional extreme value analysis of losses given the random path of the tornado, and we do not incorporate other covariates (such as population density) that can contribute to the monetary loss. This results in a sample  $(\mathbf{X}_t, Y_t)$  of size  $n = 6,360$  across the whole of the US (excluding Alaska and Hawaii), including the major Joplin, Missouri, tornado which caused a total loss of 2.8 billion USD on 22 May 2011. We focus on the part of the US to the east of the 100th meridian west, due to sparsity of recorded tornadoes to the west of this geographical limit. The data, over the studied area, is represented in Figure 2 (a).

<sup>2</sup>Available at <https://www.spc.noaa.gov/wcm/\#data>

The heavy tail model assumption was checked using local Generalized Pareto QQ-plots [see pp. 90-91 in 14] omitted for the sake of brevity. Tail index estimates were found to exceed 1 in many locations  $\mathbf{x}$ , which prevents the use of the expectile risk measure.

Our target is the extremal conditional risk measure  $q(\tau'_n|\mathbf{x})$  at each location  $\mathbf{x}$ , where  $\tau'_n = 0.995$  corresponds to a catastrophic loss exceeded (on average) only once every 200 cases at the location of interest. This is a reasonable choice given that, for instance, the state of Florida has recorded around 200 tornado events over the 10-year period we examine. A major hurdle to address that goal in this two-dimensional setting is the practical calculation of the optimal bandwidth  $h_{n,\star}$ . The usual rule-of-thumb calculations for the evaluation of  $\int_{\mathbb{R}^2} (\mathbf{u}^\top H_g(\mathbf{x})\mathbf{u}) K(\mathbf{u}) d\mathbf{u}$ , involved in (4), based on a bivariate Gaussian assumption lead to tedious calculations and unappealing results, and a uniform distribution over the covariate space cannot be chosen as it would have an identically zero Hessian matrix. Instead, the crude observation that the latitude  $X_2 \in [18, 49]$  of the data points concentrates around its median while the longitude  $X_1 \in [-100, -66.5]$  appears more uniformly scattered suggests, for the specific purpose of calculating  $h_{n,\star}$  only, to make the simplifying assumption that

$$g(x_1, x_2) \propto \mathbb{1}_{[-100, -66.5]}(x_1) \times (49 - x_2)(x_2 - 18) \mathbb{1}_{[18, 49]}(x_2).$$

This leads to a diagonal nonzero Hessian matrix  $H_g(\mathbf{x})$ . Letting  $\kappa(u) = (15/16)(1 - u^2)^2$ , for  $u \in [-1, 1]$ , be the one-dimensional quartic kernel and  $K(\mathbf{u}) = (16/(5\pi))\kappa(\|\mathbf{u}\|)$  be its isotropic version on  $\mathbb{R}^2$ , leads to the global spatial bandwidth  $h_{n,\star} \approx 5.47$ . Then, for each geographic location  $\mathbf{x}$ , representing one of the 21,935 cities located east of the 100th meridian picked in the United States Cities Database<sup>3</sup>, we chose the corresponding local optimal hyperparameter  $k_{n,\star} = \hat{k}_{n,\star}(\mathbf{x})$  as described above in Section 5.1, with the only difference that we set  $\bar{\rho}(\mathbf{x}) \equiv -1$  and  $\bar{b}(\mathbf{x}) \equiv 1$ , as recommended in Section C.4 of [12] when the covariate has dimension 2. Each city  $\mathbf{x}$  has in its  $h_{n,\star}$ -vicinity 1,100 observations on average, with 90% of locations reporting at least 400 observations. Figure 2 (b) displays the historical frequency of tornadoes, showing that the area most often hit mainly comprises the states of Alabama, Mississippi, Louisiana, Arkansas, Missouri, Kentucky, and Tennessee.

The Nadaraya-Watson estimates of the conditional mean of losses per squared yard and the extrapolated bias-reduced conditional quantile estimates are shown in Figures 2 (c) and (d), respectively. The first conclusion is that the area most exposed to tornadoes is actually not the riskiest in terms of average and/or extremal conditional losses per tornado. By contrast, tornadoes in the states of Florida, Texas, Oklahoma, Nebraska, South Dakota, Iowa and their surroundings are found to carry the most extreme risk, with a 99.5%-regression Value-at-Risk exceeding 80 USD per squared yard, even though the frequency of tornadoes there is substantially lower. The large difference in order of magnitude between the regression mean and extremal quantile reflects the great variability and tail heaviness of the conditional loss distribution; an important benefit of the nonparametric approach is its ability to accurately identify conditional extreme value behavior, without recourse to any strong *a priori* spatio-temporal model specification. According to the results obtained at 9 selected cities, reported in Table 1, the conditional tail index varies with  $\mathbf{x}$ , with a minimum of 0.79 achieved at Harrisville, MI, as opposed to 1.53 in Woodson, TX, which is the location with maximal estimated extreme quantile risk. We also note that, among others, Charleston, SC and New Orleans, LA have very similar tail index estimates but completely different tail quantile estimates, owing to a strong geographical heterogeneity in the scale parameter of the loss distribution. Most importantly, the confidence intervals  $\hat{I}_{q,4}(0.995|\mathbf{x})$  based on purely empirical estimates are much wider than those provided by our bias- and variance-corrected proposal  $\hat{I}_{q,2}(0.995|\mathbf{x})$  based

<sup>3</sup>Available at <https://simplemaps.com/data/us-cities>



Location $\mathbf{x}$ (State)	$N_{h_{n,*}}(\mathbf{x})$	$\hat{\gamma}_{1-k_{n,*}/n}^{(J, BR)}(\mathbf{x})$	$\hat{q}_n(0.995 \mathbf{x})$ $\hat{I}_{q,4}(0.995 \mathbf{x})$	$\hat{q}_{n,1-k_{n,*}/n}^{W, BR}(0.995 \mathbf{x})$ $\hat{I}_{q,2}(0.995 \mathbf{x})$
New York (NY)	413	1.11	16.71 [2.46, 113.35]	63.53 [18.13, 222.70]
Charleston (SC)	839	0.96	48.08 [12.54, 184.28]	56.15 [22.58, 139.64]
Nashville (TN)	2,317	0.95	16.57 [8.31, 33.04]	24.51 [14.62, 41.08]
Captiva (FL)	205	0.93	236.74 [36.08, 1553.36]	144.67 [43.50, 481.13]
New Orleans (LA)	1,427	0.98	27.46 [11.44, 65.88]	31.01 [16.42, 58.58]
Woodson (TX)	958	1.53	118.37 [16.53, 847.72]	390.59 [100.95, 1511.26]
Kansas City (MO)	1,326	1.30	45.99 [11.29, 187.43]	69.00 [25.57, 186.21]
Minneapolis (MN)	620	0.93	34.09 [9.54, 121.79]	40.49 [17.01, 96.40]
Harrisville (MI)	472	0.79	24.86 [4.76, 129.71]	29.32 [10.28, 83.67]

TABLE 1

*Tornado losses data. Results at selected cities (first column), with the number of neighboring observations (second column), conditional tail index estimate (third column), empirical conditional quantile estimate at level 0.995 (fourth column) and extrapolated bias-reduced conditional quantile estimate at the same level (fifth column), along with the 95% asymptotic confidence interval corresponding to each quantile estimate in brackets. Captiva, FL is the city with maximal estimated average loss; Woodson, TX is the city with maximal extrapolated conditional quantile estimate; Harrisville, MI is the city with minimal estimated conditional tail index.*

on extrapolated quantile estimates, and overrepresent the uncertainty about high regression quantiles, in agreement with our conclusions based on simulated data.

Our results are compared with the straightforward unconditional approach estimating, in each state, the mean and Value-at-Risk at level 99.5% from the univariate sample of losses in this state only, see Figures 2 (e) and (f). This approach yields quantile point estimates in the riskiest areas that are up to 67% lower than those of the regression method, and produces unrealistic discontinuities in the estimates, see the examples of Texas-Oklahoma-Kansas, Nebraska-South Dakota and Georgia-South Carolina, while the local nature of our proposed approach eliminates these discontinuities between neighboring states by combining regional tail information. A similar countywide analysis in the spirit of risk assessment exercises such as those reported by the US Federal Emergency Management Agency (FEMA)<sup>4</sup> cannot be carried out here, because certain counties did not report any observation.

<sup>4</sup> Available at <https://hazards.fema.gov/nri/tornado>

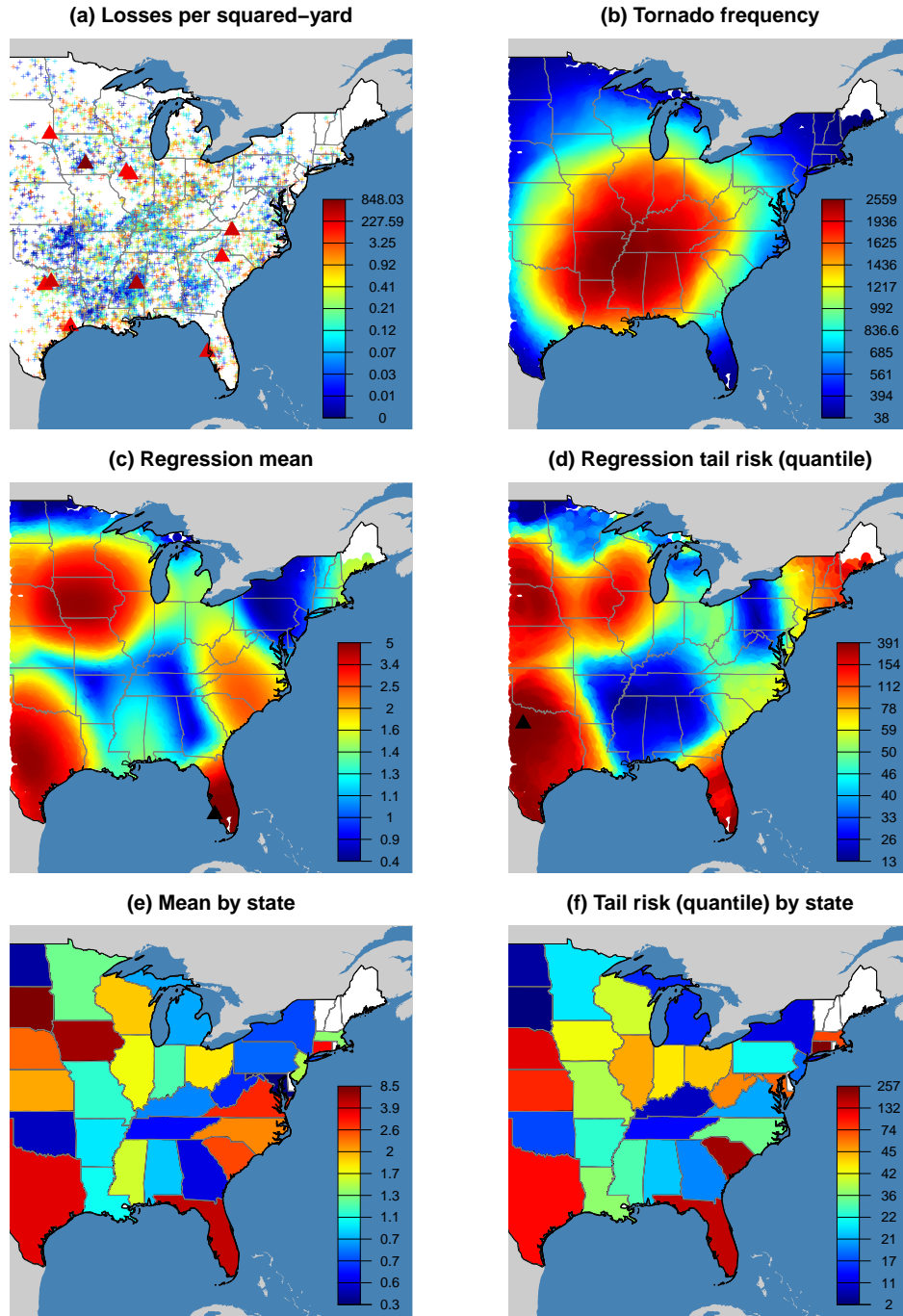


FIG 2. Tornado losses data. Top row, left: Data across the eastern half of the US, right: Local number of observations in the  $h_{n,\star}$ -ball. Middle row, left: Estimated conditional mean of losses per squared yard, right: Extrapolated conditional quantile estimate of those losses at level  $\tau'_n = 0.995$ . Cities with the highest estimated conditional average loss and extreme loss are marked with a black triangle in the left and right panels, respectively. Bottom row, left: Unconditional statewide estimation, using the sample average, right: Using the bias-reduced extreme quantile estimator of [19]. Losses (all panels except (b), in USD) and tornado frequency (panel (b)) are indicated by a color scheme, ranging from dark blue (lowest) to dark red (highest).

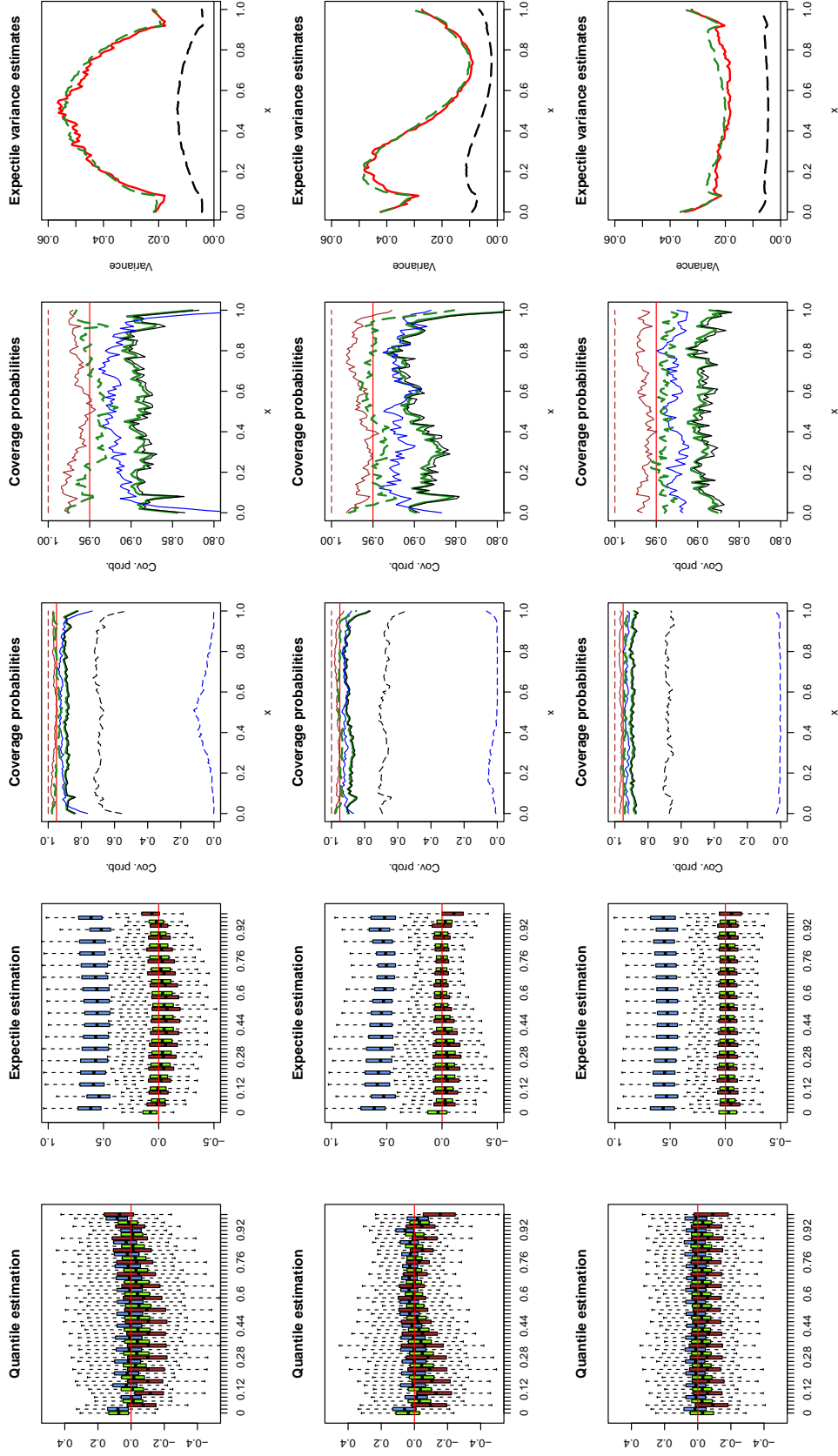


FIG 3. Simulation results in dimension  $p = 1$ . From left to right, boxplots of  $\log(\hat{q}_{1-k_{n,*}/n}^{W, BR}(\tau'_n|x)/q(\tau'_n|x))$  (green),  $\log(\hat{q}_{1-k_{n,*}/n}^W(\tau'_n|x)/q(\tau'_n|x))$  (blue) and  $\log(\hat{q}_n(\tau'_n|x)/q(\tau'_n|x))$  (brown); boxplots of  $\log(\hat{e}_{1-k_{n,*}/n}^{W, BR}(\tau'_n|x)/e(\tau'_n|x))$  (green),  $\log(\hat{e}_{1-k_{n,*}/n}^W(\tau'_n|x)/e(\tau'_n|x))$  (blue) and  $\log(\hat{e}_n(\tau'_n|x)/e(\tau'_n|x))$  (brown); empirical pointwise coverage probabilities of the asymptotic 95% confidence intervals  $\hat{I}_{q,1}(\tau'_n|x)$  (full black line),  $\hat{I}_{q,2}(\tau'_n|x)$  (dashed black line),  $\hat{I}_{q,3}(\tau'_n|x)$  (full green line),  $\hat{I}_{q,4}(\tau'_n|x)$  (dashed green line),  $\hat{I}_{q,2}(\tau'_n|x)$  (dashed blue line) and  $\hat{I}_{E,4}(\tau'_n|x)$  (dashed brown line), with the target 95% nominal level in full red line; same figure, zooming around nominal coverage; comparison of the pointwise empirical variances of the  $(\sqrt{k_{n,*}b_{n,*}}/\log(k_{n,*}/(n(1-\tau'_n)))) \log(\hat{e}_{1-k_{n,*}/n}^{W, BR}(\tau'_n|x)/e(\tau'_n|x))$  (full red line), the pointwise median of the corrected variance estimates  $(\int_{\mathbb{R}} K^2/\hat{g}_n(x))\tilde{v}_E(x)$  (dashed green line) and the pointwise median of the naive variance estimates  $(\int_{\mathbb{R}} K^2/\hat{g}_n(x))\tilde{v}_E(x)$  (dashed black line). Top row: model (NL-P), middle row: model (NL-S), bottom row: model (NL-C), all with  $n = 10,000$  and estimates at level  $\tau'_n = 0.999$ .

**Acknowledgments.** The authors acknowledge an anonymous Associate Editor and three anonymous reviewers for their helpful comments that led to a much improved article.

**Funding.** Support from the ANR (grants ANR-19-CE40-0013 and ANR-17-EURE-0010) and the Centre Henri Lebesgue (ANR-11-LABX-0020-01) is gratefully acknowledged. A. Daouia and G. Stupfler acknowledge support from the TSE-HEC ACPR Chair and an AXA Research Fund Award.

## SUPPLEMENTARY MATERIAL

**Supplementary material for “Inference for extremal regression with dependent heavy-tailed data”** The supplementary material document [12] contains further details about our technical conditions and an expanded discussion of the rates of pointwise convergence of our estimators. We then provide the proofs of all theoretical results in the main paper and a full analysis of our worked-out regression examples, preceded by auxiliary results and their proofs. We finally provide further details about our bias and variance correction procedures, and extra finite-sample results.

## REFERENCES

- [1] ARTZNER, P., DELBAEN, F., EBER, J. M. and HEATH, D. (1999). Coherent measures of risk. *Mathematical Finance* **9** 203–228.
- [2] BEIRLANT, J., GOEGEBEUR, Y., SEGERS, J. and TEUGELS, J. (2004). *Statistics of Extremes: Theory and Applications*. Wiley.
- [3] BELLINI, F. and DI BERNARDINO, E. (2017). Risk management with expectiles. *The European Journal of Finance* **23** 487–506.
- [4] BRADLEY, R. C. (2005). Basic properties of strong mixing conditions. A survey and some open questions. *Probability Surveys* **2** 107–144.
- [5] CHAUDHURI, P. (1991). Global nonparametric estimation of conditional quantile functions and their derivatives. *Journal of Multivariate Analysis* **39** 246–269.
- [6] CHERNOZHUKOV, V. (2005). Extremal quantile regression. *Annals of Statistics* **33** 806–839.
- [7] CHERNOZHUKOV, V. and FERNÁNDEZ-VAL, I. (2011). Inference for extremal conditional quantile models, with an application to market and birthweight risks. *Review of Economic Studies* **78** 559–589.
- [8] CHERNOZHUKOV, V., FERNÁNDEZ-VAL, I. and KAJI, T. (2017). Extremal Quantile Regression. In *Handbook of Quantile Regression* (R. Koenker, V. Chernozhukov, X. He and L. Peng, eds.) Chapman and Hall/CRC.
- [9] DAOUIA, A., GARDES, L. and GIRARD, S. (2013). On kernel smoothing for extremal quantile regression. *Bernoulli* **19** 2557–2589.
- [10] DAOUIA, A., GARDES, L., GIRARD, S. and LEKINA, A. (2011). Kernel estimators of extreme level curves. *TEST* **20** 311–333.
- [11] DAOUIA, A., GIRARD, S. and STUPFLER, G. (2018). Estimation of tail risk based on extreme expectiles. *Journal of the Royal Statistical Society: Series B* **80** 263–292.
- [12] DAOUIA, A., STUPFLER, G. and USSEGLIO-CARLEVE, A. (2023). Supplement to “Inference for extremal regression with dependent heavy-tailed data”. DOI: 10.1214/[provided by typesetter].
- [13] DAVISON, A. C., PADOAN, S. A. and STUPFLER, G. (2023). Tail risk inference via expectiles in heavy-tailed time series. *Journal of Business and Economic Statistics* **41** 876–889.
- [14] DE HAAN, L. and FERREIRA, A. (2006). *Extreme Value Theory: An Introduction*. Springer, New York.
- [15] DREES, H. (2003). Extreme quantile estimation for dependent data, with applications to finance. *Bernoulli* **9** 617–657.
- [16] GIRARD, S., STUPFLER, G. and USSEGLIO-CARLEVE, A. (2021). Extreme conditional expectile estimation in heavy-tailed heteroscedastic regression models. *Annals of Statistics* **49** 3358–3382.
- [17] GIRARD, S., STUPFLER, G. and USSEGLIO-CARLEVE, A. (2022). Nonparametric extreme conditional expectile estimation. *Scandinavian Journal of Statistics* **49** 78–115.
- [18] GIRARD, S., STUPFLER, G. and USSEGLIO-CARLEVE, A. (2022). On automatic bias reduction for extreme expectile estimation. *Statistics & Computing* **32** 64.

- [19] GOMES, M. I. and PESTANA, D. (2007). A sturdy reduced-bias extreme quantile (VaR) estimator. *Journal of the American Statistical Association* **102** 280–292.
- [20] JONES, M. C. (1994). Expectiles and M-quantiles are quantiles. *Statistics & Probability Letters* **20** 149–153.
- [21] LINTON, O. and XIAO, Z. (2013). Estimation of and inference about the Expected Shortfall for time series with infinite variance. *Econometric Theory* **29** 771–807.
- [22] NEWEY, W. K. and POWELL, J. L. (1987). Asymmetric least squares estimation and testing. *Econometrica* **55** 819–847.
- [23] SMITH, R. L. and WEISSMAN, I. (1996). Characterization and Estimation of the Multivariate Extremal Index. Technical Report, University of North Carolina.
- [24] USSEGLIO-CARLEVE, A. (2018). Estimation of conditional extreme risk measures from heavy-tailed elliptical random vectors. *Electronic Journal of Statistics* **12** 4057–4093.
- [25] WANG, H. J., LI, D. and HE, X. (2012). Estimation of high conditional quantiles for heavy-tailed distributions. *Journal of the American Statistical Association* **107** 1453–1464.
- [26] WASSERMAN, L. (2006). *All of Nonparametric Statistics*. Springer.
- [27] WEISSMAN, I. (1978). Estimation of parameters and large quantiles based on the  $k$  largest observations. *Journal of the American Statistical Association* **73** 812–815.
- [28] YEE, T. W. and WILD, C. J. (1996). Vector generalized additive models. *Journal of the Royal Statistical Society: Series B* **58** 481–493.
- [29] ZIEGEL, J. F. (2016). Coherence and elicibility. *Mathematical Finance* **26** 901–918.