

Quantile Regression

Marcelo Castiel Ruas*, Henrique Helfer Hoeltgebaum†, Alexandre Street‡,
Cristiano Fernandes§
January 24, 2017

*Aluno de doutorado do Departamento de Engenharia Elétrica da PUC-RIO.

†Aluno de doutorado do Departamento de Engenharia Elétrica da PUC-RIO.

‡Professor do Departamento de Engenharia Elétrica da PUC-RIO.

§Professor do Departamento de Engenharia Elétrica da PUC-RIO.

1 Introduction

Quantile Regression is a powerful tool for measuring quantiles others than the median or predicting the mean. A quantile of a random variable is important in risk measuring, as we can measure the probability of occurrence of extreme events, and in many other fields. While working with energy forecasts, quantile regression can produce interesting results when working with both short term (hourly) or long term (monthly) data. As an example, we present a solar time series for the short term and a wind time series for long term. The first set of data is measured at the location of Tubarao (Brazil) on the year of 2014, while the latter is a dataset of mean power monthly observations from Icaraizinho (Brazil) between 1981 to 2011 of measured in Megawatts. Figures 1.1, 1.2 and Figures **Inserir figuras de dados solares.** illustrate the seasonality present in these datasets.

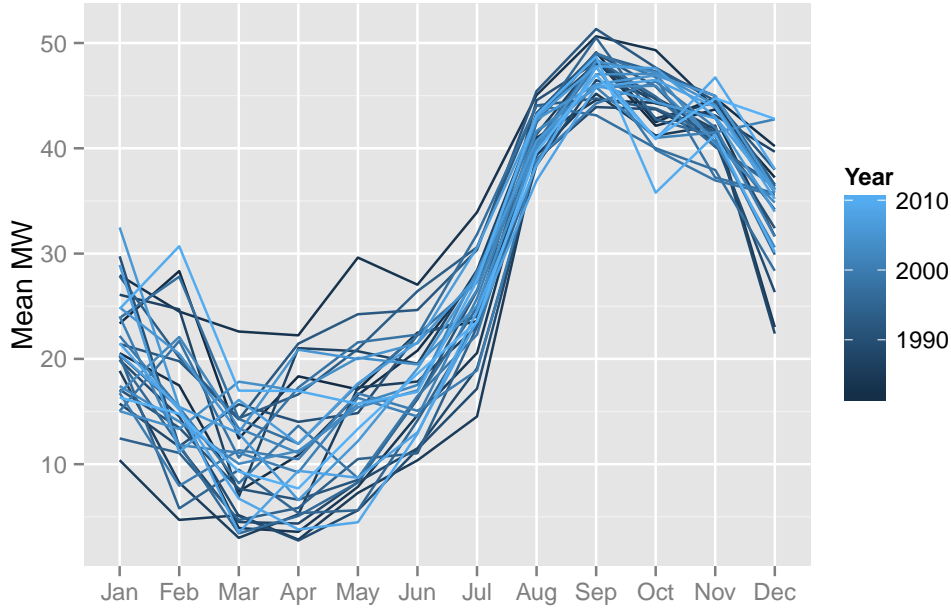


Figure 1.1: Icaraizinho yearly data. Each serie consists of monthly observations for each year.

In this work, we apply a few different techniques to forecast the quantile function a few steps ahead. The main frameworks we investigate are parametric linear models and a non-parametric regression. In all approaches we use the time series lags as the regression covariates. We also investigate how to apply quantile estimations to produce an empirical distribution for the k -step ahead forecasting by using a nonparametric approach.

To make good predictions of random variables, one must find good explanatory variables: it can be either autoregressive, exogenous terms or even a deterministic function that repeats itself. Figure 1.3 shows scatter plots relating y_t with some of its lags. We can see that in both sort and long term, past values are good explanatory variables to use for forecasting.

In contrast to the linear regression model through ordinary least squares (OLS), which provides only an estimation of the dependent variable conditional mean, quantile regression model yields a much more detailed information concerning the complex relationship about the dependent variable and its covariates. Here we denote as parametric linear model the well-known quantile regression model [3]. A Quantile Regression for the α -quantile is the solution of the following optimization problem:

$$\min_q \sum_{t=1}^n \alpha |y_t - q(x_t)|^+ + (1 - \alpha) |y_t - q(x_t)|^-, \quad (1.1)$$

where $q(x_t)$ is the estimated quantile value at a given time t and $|x|^+ = \max\{0, x\}$ and $|x|^- = -\min\{0, x\}$. To model this problem as a Linear Programming problem, thus being able to use a

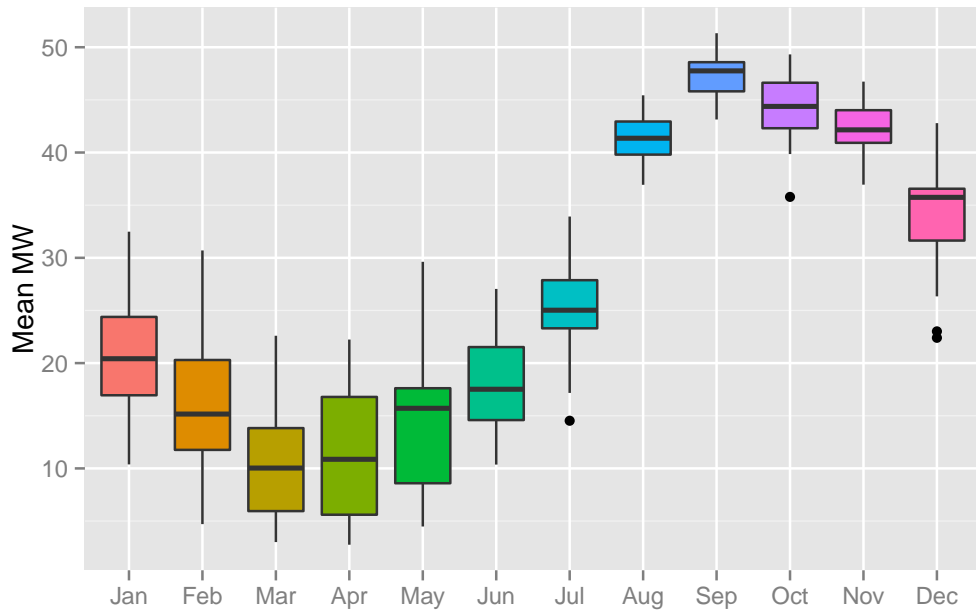


Figure 1.2: Boxplot for each month for the Icaraizinho dataset

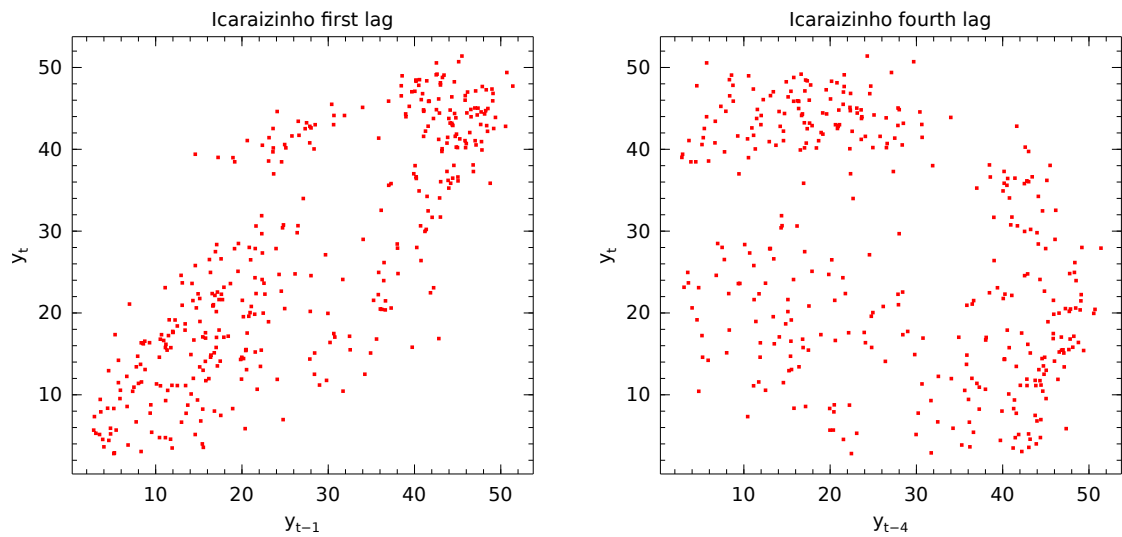


Figure 1.3: Relationship between y_t and some chosen lags.

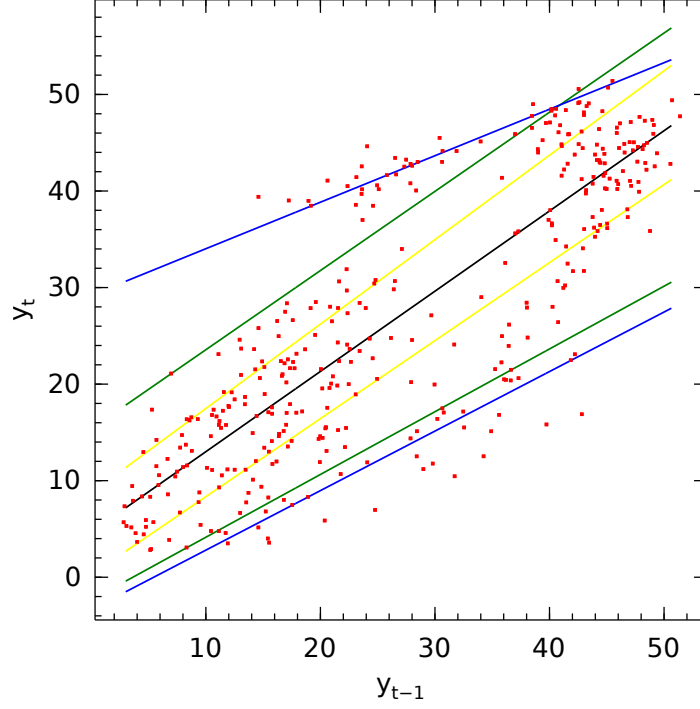


Figure 1.4: Linear quantile estimator with crossing quantiles for $\alpha = 0.95$ and $\alpha = 0.9$

modern solver to fit our model, we can create variables ε_t^+ e ε_t^- to represent $|y - q(x_t)|^+$ and $|y - q(x_t)|^-$, respectively. So we have:

$$\begin{aligned} \min_{q, \varepsilon_t^+, \varepsilon_t^-} \quad & \sum_{t=1}^n (\alpha \varepsilon_t^+ + (1 - \alpha) \varepsilon_t^-) \\ \text{s.t.} \quad & \varepsilon_t^+ - \varepsilon_t^- = y_t - q^\alpha(x_t), \quad \forall t \in \{1, \dots, n\}, \\ & \varepsilon_t^+, \varepsilon_t^- \geq 0, \quad \forall t \in \{1, \dots, n\}. \end{aligned} \quad (1.2)$$

One of our goals with quantile regression is to estimate a distribution function F_X of a given random variable X from a sequence of quantiles $q_t^{\alpha_1} \leq q_t^{\alpha_2} \leq \dots \leq q_t^{\alpha_Q}$, with $0 < \alpha_1 < \alpha_2 < \dots < \alpha_Q < 1$. The process of fitting \hat{F}_X is by mapping every α_i with its estimated quantile \hat{q}^{α_i} . When this sequence of chosen α_i is thin enough, we can approximate well the distribution function of X , as is shown in Figure 1.5. Thus, the distribution found for X is nonparametric, as no previous assumptions are made about its shape, and its form is fully recovered by the data we have.

A typical problem, however, arises when working with quantile regression. When quantiles are estimated independently, it is possible to find $q_t^{\alpha_1} > q_t^{\alpha_2}$, for a given t , when $\alpha_1 < \alpha_2$. An example can be seen on Figure 1.4, where quantiles $\alpha = 0.95$ and $\alpha = 0.9$ cross. This problem, called *crossing quantiles*, can be prevented by estimating all quantiles with a single maximization problem.

Let A be the set containing all quantiles α_i such that $0 < \alpha_1 < \alpha_2 < \dots < \alpha_Q < 1$ and $T_\tau = \{1, 2, \dots, \tau\}$. In order to estimate all quantiles simultaneously, the new objective function will be the sum of all individual objective functions, as well as include all constraints from all individual problems. The only difference is the inclusion of an equation to guarantee that quantiles won't cross. When modifying problem 1.2 to account for all quantiles, we have the following new problem:

$$\min_{q, \varepsilon_{t,\alpha}^+, \varepsilon_{t,\alpha}^-} \quad \sum_{\alpha \in A} \sum_{t \in T} (\alpha \varepsilon_{t,\alpha}^+ + (1 - \alpha) \varepsilon_{t,\alpha}^-) \quad (1.3)$$

$$\text{s.t.} \quad \varepsilon_{t,\alpha}^+ - \varepsilon_{t,\alpha}^- = y_t - q^\alpha(x_t), \quad \forall t \in T_\tau, \forall \alpha \in A, \quad (1.4)$$

$$\varepsilon_{t,\alpha}^+, \varepsilon_{t,\alpha}^- \geq 0, \quad \forall t \in T_\tau, \forall \alpha \in A, \quad (1.5)$$

$$q_t^\alpha \leq q_t^{\alpha'}, \quad \forall t \in T_\tau, \forall \alpha, \alpha' \in A, \alpha < \alpha', \quad (1.6)$$

where constraint 1.6 assures that no lower quantile will have a bigger value than a higher quantile.

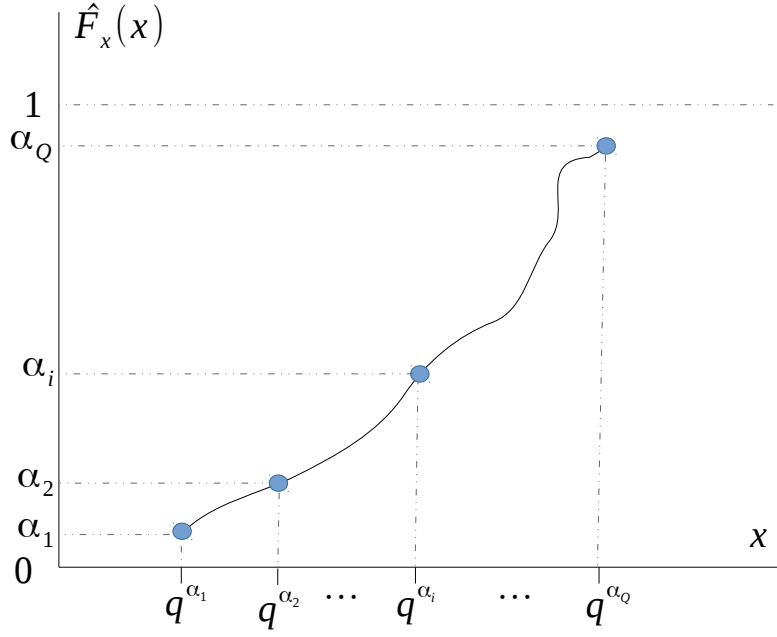


Figure 1.5: Fitting a distribution function from quantile estimations

The next section discusses with bigger details how to fit a distribution function from a sequence of estimated quantiles, as well as showing two different strategies to estimate them: linear models and nonparametric models. In the former, q is a linear function of the series past values, up to a maximum number of lags p , such as:

$$q(y_t, \alpha; \beta) = \beta_0(\alpha) + \beta_1(\alpha)y_{t-1} + \beta_2(\alpha)y_{t-2} + \dots + \beta_p(\alpha)y_{t-p}. \quad (1.7)$$

In the latter, we let $q(x_t)$ assume any functional form. To prevent overfitting, however, we penalize the function's roughness by incorporating a penalty on the second derivative.

In section 3 we investigate how to simulate S scenarios of y_t , considering a linear model and errors ε_t for which the distribution is unknown. To address this issue, we use quantile linear regression to calculate a thin grid of quantiles and fit a distribution function \hat{F}_{y_t} . This function will be used to simulate the innovations on the model.

2 Estimating distribution function from quantile regressions

In many applications where a time series model is employed, we often consider the innovations' distribution as known. Take, for example, the ARMA(p,q) model:

$$y_t = c + \varepsilon_t + \sum_{i=1}^p \phi_i y_{t-i} + \sum_{i=1}^q \theta_i \varepsilon_{t-i}.$$

In this model, errors ε_t are assumed to have normal distribution with zero mean.

When we are dealing with meteorological time series, however, we can't always assume normality. In these cases, one can either find a distribution that has a better fit to the data or have a nonparametric method to estimate the distribution directly from the available data.

In a time series framework, where a time series y_t is given by a linear model of its regressors x_t

$$y_t = \beta^T x_t + \varepsilon_t,$$

we propose to estimate the k -step ahead distribution of y_t with a nonparametric approach. Let a quantile be given by

$$q_\alpha(x_t) = \hat{F}_{y_t|X_t=x_t}^{-1}(\alpha), \quad (2.1)$$

when joining a thin grid of quantiles we can approximate the distribution function F_{y_t} .

In any given t , by choosing many different values of α_i , we can estimate a sequence of quantiles $q_t^{\alpha_1} \leq q_t^{\alpha_2} \leq \dots \leq q_t^{\alpha_Q}$ with $0 < \alpha_1 < \alpha_2 < \dots < \alpha_Q < 1$. This is done by solving the problem defined on equations (1.3)-(1.6). Let F_{y_t} be the estimated distribution function of y_t . The process of fitting \hat{F}_{y_t} is by mapping every α_i with its estimated quantile \hat{q}^{α_i} . A problem arises for the distribution extremities, because when $\alpha = 0$ or $\alpha = 1$, the optimization problem becomes unbounded. In order to find good estimates for y_t when F_{y_t} approaches 0 or 1, we can either use a kernel smoothing function, splines, linear approximation, or any other method. **This will be developed later.**

In the sections 2.1 and 2.2 we investigate two ways of estimating quantiles. The next session explores estimating the quantile q_t^α is given by the following linear model:

$$q_t^\alpha = \beta_0^\alpha + x_t^T \beta^\alpha + \varepsilon_t, \quad (2.2)$$

where β^α is a vector of coefficients for the explanatory variables. In section 2.2 we introduce a Nonparametric Quantile Autoregressive model with a ℓ_1 -penalty term, in order to properly simulate densities for several α -quantiles. In this nonparametric approach we don't assume any form for $q(x_t)$, but rather let the function adjust to the data. To prevent overfitting, the ℓ_1 penalty for the second derivative (approximated by the second difference of the ordered observations) is included in the objective function.

2.1 Linear Models for the Quantile Autoregression

Given a time series $\{y_t\}$, we investigate how to select which lags will be included in the Quantile Autoregression. We won't be choosing the full model because this normally leads to a bigger variance in our estimators, which is often linked with bad performance in forecasting applications. So our strategy will be to use some sort of regularization method in order to improve performance. We investigate two ways of accomplishing this goal. The first of them consists of selecting the best subset of variables through Mixed Integer Programming, given that K variables are included in the model. Using MIP to select the best subset of variables is investigated in [1]. The second way is including a ℓ_1 penalty on the linear quantile regression, as in [2], and let the model select which and how many variables will have nonzero coefficients. Both of them will be built over the standard Quantile Linear Regression model. In the end of the section, we discuss a information criteria to be used for quantile regression and verify how close are the solutions in the eyes of this criteria.

When we choose $q(x_t)$ to be a linear function, as on equation 1.1 (that we reproduce below for convenience):

$$\min_{q_t} \sum_{t \in T_\tau}^n \alpha |y_t - q(x_t)|^+ + (1 - \alpha) |y_t - q(x_t)|^-, \quad (2.3)$$

we can substitute it on problem 1.2, getting the following LP problem:

$$\begin{aligned} \min_{\beta_0, \beta, \varepsilon_t^+, \varepsilon_t^-} \quad & \sum_{t=1}^n (\alpha \varepsilon_t^+ + (1 - \alpha) \varepsilon_t^-) \\ \text{s.t.} \quad & \varepsilon_t^+ - \varepsilon_t^- = y_t - \beta_0 - \beta^T x_t, \quad \forall t \in \{1, \dots, n\}, \\ & \varepsilon_t^+, \varepsilon_t^- \geq 0, \quad \forall t \in \{1, \dots, n\}. \end{aligned} \quad (2.4)$$

$$\min_{q, \varepsilon_{t,\alpha}^+, \varepsilon_{t,\alpha}^-} \quad \sum_{\alpha \in A} \sum_{t \in T} (\alpha \varepsilon_{t,\alpha}^+ + (1 - \alpha) \varepsilon_{t,\alpha}^-) \quad (2.5)$$

$$\text{s.t.} \quad \varepsilon_{t,\alpha}^+ - \varepsilon_{t,\alpha}^- = y_t - \beta_0^\alpha - (\beta^\alpha)^T x_t, \quad \forall t \in T_\tau, \forall \alpha \in A, \quad (2.6)$$

$$\varepsilon_{t,\alpha}^+, \varepsilon_{t,\alpha}^- \geq 0, \quad \forall t \in T_\tau, \forall \alpha \in A, \quad (2.7)$$

$$q_t^\alpha \leq q_t^{\alpha'}, \quad \forall t \in T_\tau, \forall \alpha, \alpha' \in A, \alpha < \alpha', \quad (2.8)$$

In this work, we didn't explore the addition of terms other than the terms y_t past lags. For example, we could include functions of y_{t-p} , such as $\log(y_{t-p})$ or $\exp(y_{t-p})$. We leave such inclusion for further works.

2.2 Quantile Autoregression with a nonparametric approach

Fitting a linear estimator for the Quantile Auto Regression isn't appropriate when nonlinearity is present in the data. This nonlinearity may produce a linear estimator that underestimates the quantile for a chunk of data while overestimating for the other chunk (for example, scatter plot of y_t versus y_{t-1} that is seen on the upper left of figure 1.3). To prevent this issue from occurring we propose a modification which we let the prediction $\mathcal{Q}_{y_t|y_{t-1}}(\alpha)$ adjust freely to the data and its nonlinearities. To prevent overfitting and smoothen our predictor, we include a penalty on its roughness by including the ℓ_1 norm of its second derivative. For more information on the ℓ_1 norm acting as a filter, one can refer to [2].

Let $\{\tilde{y}_t\}_{t=1}^n$ be the sequence of observations in time t . Now, let \tilde{x}_t be the p -lagged time series of \tilde{y}_t , such that $\tilde{x}_t = L^p(\tilde{y}_t)$, where L is the lag operator. Matching each observation \tilde{y}_t with its p -lagged correspondent \tilde{x}_t will produce $n - p$ pairs $\{(\tilde{y}_t, \tilde{x}_t)\}_{t=p+1}^n$ (note that the first p observations of y_t must be discarded). When we order the observation of x in such way that they are in growing order

$$\tilde{x}^{(p+1)} \leq \tilde{x}^{(p+2)} \leq \dots \leq \tilde{x}^{(n)},$$

we can then define $\{x_i\}_{i=1}^{n-p} = \{\tilde{x}^{(t)}\}_{t=p+1}^n$ and $\{y_i\}_{i=1}^{n-p} = \{\tilde{y}^{(t)}\}_{t=p+1}^n$ and $T' = \{2, \dots, n - p - 1\}$. As we need the second difference of q_i , I has to be shortened by two elements.

Our optimization model to estimate the nonparametric quantile is as follows:

$$\begin{aligned} \mathcal{Q}_{y_\tau|y_{\tau-1}}^\alpha(\tau) = \arg \min_{q_t} \quad & \sum_{t \in T'} (\alpha |y_t - q_t|^+ + (1 - \alpha) |y_t - q_t|^-) \\ & + \lambda \sum_{t \in T'} |D_{x_t}^2 q_t|, \end{aligned} \quad (2.9)$$

where $D^2 q_t$ is the second derivative of the q_t function, calculated as follows:

$$D_{x_t}^2 q_t = \frac{\left(\frac{q_{t+1} - q_t}{x_{t+1} - x_t} \right) - \left(\frac{q_t - q_{t-1}}{x_t - x_{t-1}} \right)}{x_{t+1} - 2x_t + x_{t-1}}.$$

The first part on the objective function is the usual quantile regression condition for $\{q_t\}$. The second part is the ℓ_1 -filter. The purpose of a filter is to control the amount of variation for our estimator q_t . When no penalty is employed we would always get $q_t = y_t$. On the other hand, when $\lambda \rightarrow \infty$, our estimator approaches the linear quantile regression.

The full model can be rewritten as a LP problem as bellow:

$$\min_{q_t} \sum_{t=1}^n (\alpha \delta_t^+ + (1 - \alpha) \delta_t^-) + \lambda \sum_{t \in T'} \xi_t \quad (2.10)$$

$$s.t. \quad \delta_t^+ - \delta_t^- = y_t - q_t, \quad \forall t \in \{3, \dots, n-1\}, \quad (2.11)$$

$$D_t = \left(\frac{q_{t+1} - q_t}{x_{t+1} - x_t} \right) - \left(\frac{q_t - q_{t-1}}{x_t - x_{t-1}} \right) \quad \forall t \in \{3, \dots, n-1\}, \quad (2.12)$$

$$\xi_t \geq D_t, \quad \forall t \in \{3, \dots, n-1\}, \quad (2.13)$$

$$\xi_t \geq -D_t, \quad \forall t \in \{3, \dots, n-1\}, \quad (2.14)$$

$$\delta_t^+, \delta_t^-, \xi_t \geq 0, \quad \forall t \in \{3, \dots, n-1\}. \quad (2.15)$$

The output of our optimization problem is a sequence of ordered points $\{(x_t, q_t)\}_{t \in T}$. The next step is to interpolate these points in order to provide an estimation for any other value of x . To address this issue, we propose using a B-splines interpolation, that will be developed in another study.

The quantile estimation is done for different values of λ . By using different levels of penalization on the second difference, the estimation can be more or less adaptive to the fluctuation. It is important to notice that the usage of the ℓ_1 -norm as penalty leads to a piecewise linear solution q_t . Figure 2.1 shows the quantile estimation for a few different values of λ .

When estimating quantiles for a few different values of α , however, sometimes we find them overlapping each other, which we call crossing quantiles. This effect can be seen in figure 2.1f, where the 95%-quantile crosses over the 90%-quantile. To prevent this, we can include a non-crossing constraint:

$$q_i^\alpha \leq q_i^{\alpha'}, \quad \forall i \in I, \alpha < \alpha'. \quad (2.16)$$

This means that when α' is a higher quantile than α , then the values from the α' -quantile must be bigger than those of the α -quantile for each and every point.

As a result of this nonparametric estimation, we are able to establish a relation between y_t and y_{t-p} in a way that the model adjusts itself automatically to the present nonlinearities. For this, we only have to supply a numeric value for λ . This approach, however, have yet some issues do be discussed.

The first issue is how to select an appropriate value for λ . A simple way is to do it by inspection, which means to test many different values and pick the one that suits best our needs by looking at them. The other alternative is to use a metric to which we can select the best tune. We can achieve this by using a cross-validation method, for example.

The other issue occurs when we try to add more than one lag to the analysis at the same time. This happens because the problem solution is a set of points that we need to interpolate. This multivariate interpolation, however, is not easily solved, in the sense that we can either choose using a very naive estimator such as the K-nearest neighbors or just find another method that is not yet adopted for a wide range of applications.

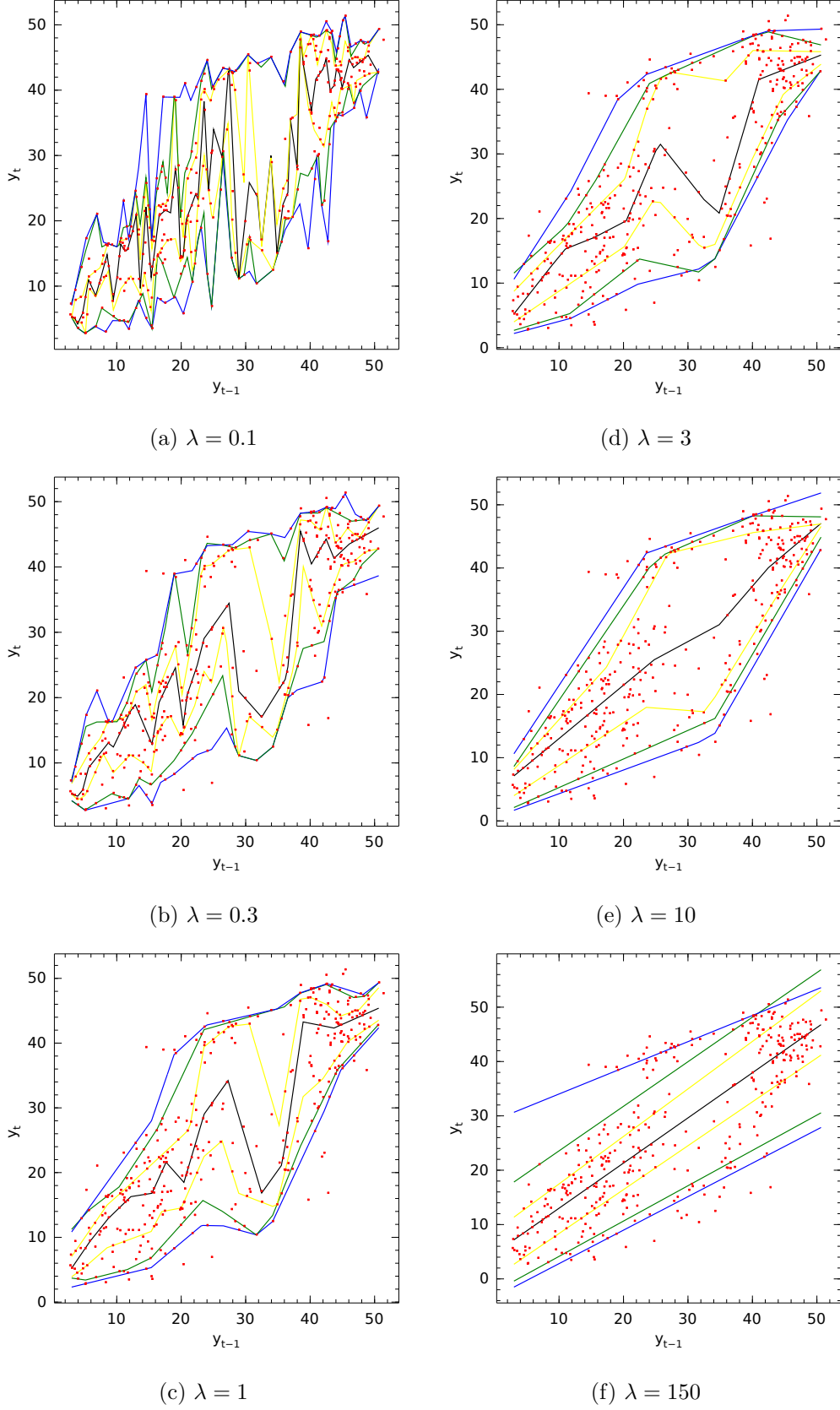


Figure 2.1: Quantile estimations for a few different values of λ . The quantiles represented here are $\alpha = (5\%, 10\%, 25\%, 50\%, 75\%, 90\%, 95\%)$. When $\lambda = 0.1$, on the upper left, we clearly see an overfitting on the estimations. The other extreme case is also shown, when $\lambda = 200$ the nonparametric estimator converges to the linear model.

3 Simulation

In this section, we investigate how to simulate future paths of the time series y_t . Let n be the total number of observations of y_t . We produce S different paths with size K for each. We have n observations of y_t and we want to produce . Given a vector of explanatory variables x_t , let q_t^α be given by the following linear model:

$$q_t^\alpha = \beta_0^\alpha + x_t^T \beta^\alpha + \varepsilon_t, \quad (3.1)$$

where β^α is a vector of coefficients for the explanatory variables. The variables chosen to compose x_t can be either exogenous variables, autoregressive components of y_t or both. As the distribution of ε_t is unknown, we have to use a nonparametric approach in order to estimate its one-step ahead density.

The coefficients β_0^α and β^α are the solution of the minimization problem given in the problem defined in (1.3)-(1.6), reproduced here for convenience:

$$\min_{q, \varepsilon_{t,\alpha}^+, \varepsilon_{t,\alpha}^-} \sum_{\alpha \in A} \sum_{t \in T} (\alpha \varepsilon_{t,\alpha}^+ + (1 - \alpha) \varepsilon_{t,\alpha}^-) \quad (3.2)$$

$$\text{s.t.} \quad \varepsilon_{t,\alpha}^+ - \varepsilon_{t,\alpha}^- = y_t - q^\alpha(x_t), \quad \forall t \in T_\tau, \forall \alpha \in A, \quad (3.3)$$

$$\varepsilon_{t,\alpha}^+, \varepsilon_{t,\alpha}^- \geq 0, \quad \forall t \in T_\tau, \forall \alpha \in A, \quad (3.4)$$

$$q_t^\alpha \leq q_t^{\alpha'}, \quad \forall t \in T_\tau, \forall \alpha, \alpha' \in A, \alpha < \alpha', \quad (3.5)$$

To produce S different paths of $\{\hat{y}_t\}_{t=n+1}^{n+K}$, we use the following procedure:

Procedure for simulating S scenarios of y_t

1. At first, let $\tau = n + 1$.
2. In any given period τ , for a sequence $0 < \alpha_1 < \alpha_2 < \dots < \alpha_Q < 1$, we use the problem defined on (1.3)-(1.6) to predict quantiles $q_t^{\alpha_1} \leq q_t^{\alpha_2} \leq \dots \leq q_t^{\alpha_Q}$. Note that x_τ is supposed to be known at time τ . In the presence of exogenous variables that are unknown, it is advisable to incorporate its uncertainty by considering different scenarios. In each scenario, though, x_τ must be considered fully known.
3. Let F_{y_τ} be the estimated distribution function of y_τ . The process of fitting \hat{F}_{y_τ} is by mapping every α_i with its estimated quantile \hat{q}^{α_i} . A problem arises for the distribution extremities, because when $\alpha = 0$ or $\alpha = 1$, the optimization problem becomes unbounded. In order to find good estimates for y_τ when F_{y_τ} approaches 0 or 1, we can either use a kernel smoothing function, splines, linear approximation, or any other method. **This will be developed later.** When this sequence of chosen α_i is thin enough, we can approximate well the distribution function of y_τ , as is shown in Figure 3.1. Thus, the distribution found for \hat{y}_τ is nonparametric, as no previous assumptions are made about its shape, and its form is fully recovered by the data we have.
4. Once we have a distribution for y_{n+1} , we can generate S different simulated values, drawn from the distribution $\hat{F}_{y_{n+1}}$ found by doing steps 2 and 3. Let X be a random variable with uniform distribution over the interval $[0, 1]$. By using results from the Probability Integral Transform, we know that the random variable $F_{y_{n+1}}^{-1}(X)$ has the same distribution as y_{n+1} . So, by drawing a sample of size S from X and applying the inverse function of $F_{y_{n+1}}$, we have our sample of size K for y_{n+1} .
5. Each one of the S different values for y_{n+1} will be the starting point of a different path. Now, for each $\tau \in [n + 2, n + K]$ and $s \in S$, we have to estimate the quantiles $q_{\tau,s}^{\alpha_i}$ and find a distribution function for $\hat{F}_{y_{\tau,s}}$ just like it was done on steps 2 and 3. Note that when $\tau > n + 2$, every estimate will be scenario dependent, hence there will be S distribution functions estimated for each period τ . From now on, in each path just one new value will be drawn randomly from the one-step ahead distribution function - as opposed to what was carried on step 3, when S values were simulated. As there will be S distribution functions - one for each path, in each period τ it

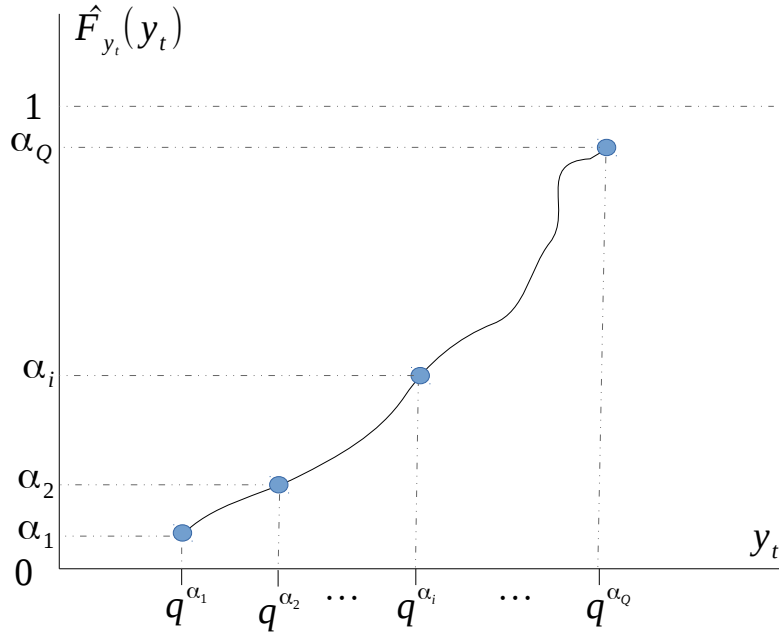


Figure 3.1: Fitting a distribution function from quantile estimations

will be produced exact S values for y_τ , one for its own path. Repeating this step until all values of τ and s are simulated will give us the full simulations that we are looking for.

References

- [1] Dimitris Bertsimas, Angela King, and Rahul Mazumder. Best subset selection via a modern optimization lens. *arXiv preprint arXiv:1507.03133*, 2015.
- [2] Seung-Jean Kim, Kwangmoo Koh, Stephen Boyd, and Dmitry Gorinevsky. ℓ_1 trend filtering. *SIAM review*, 51(2):339–360, 2009.
- [3] Roger Koenker. *Quantile regression*. Number 38. Cambridge university press, 2005.
- [4] Fernando Porrua, Bernardo Bezerra, Luiz Augusto Barroso, Priscila Lino, Francisco Ralston, and Mario Pereira. Wind power insertion through energy auctions in brazil. In *Power and Energy Society General Meeting, 2010 IEEE*, pages 1–8. IEEE, 2010.