

Nongaussian time series model via Quantile Regression

Marcelo Ruas and Alexandre Street, *Member, IEEE*

Abstract—Lorem ipsum dolor sit amet, consectetur adipiscing elit. Etiam lobortis facilisis sem. Nullam nec mi et neque pharetra sollicitudin. Praesent imperdiet mi nec ante. Donec ullamcorper, felis non sodales commodo, lectus velit ultrices augue, a dignissim nibh lectus placerat pede. Vivamus nunc nunc, molestie ut, ultricies vel, semper in, velit. Ut porttitor. Praesent in sapien. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Duis fringilla tristique neque. Sed interdum libero ut metus. Pellentesque placerat. Nam rutrum augue a leo. Morbi sed elit sit amet ante lobortis sollicitudin. Praesent blandit blandit mauris. Praesent lectus tellus, aliquet aliquam, luctus a, egestas a, turpis. Mauris lacinia lorem sit amet ipsum. Nunc quis urna dictum turpis accumsan semper.

Keywords—Quantile Regression, Model Identification, Non-gaussian time series model

TODO LIST

Mudar voz do início da introdução. Apenas apresentar, apenas mostrar a nossa escolha no objetivo . . .	1
revisar este paragrafo	1
colocar parte sobre nossa solução dos crossing quantiles	1
an optimal specification methodology regularization	2
terminar parágrafo de explicação sobre CV	5
acabar	5

I. INTRODUCTION

Mudar voz do início da introdução. Apenas apresentar, apenas mostrar a nossa escolha no objetivo

Renewable energy power is an emergent topic which is demanding attention from the academic community. The installed capacity of renewable energy plants has been increasing in a fast pace and projections point out that wind power alone will account to 18% of global power by 2050 [1]. In spite of its virtues, several new challenges are inherent when dealing with such power source. Many applications in Power Systems use renewable scenarios as input. In stochastic optimization problems such as Unit Commitment, economic dispatch, transmission expansion planning all use it. For robust optimization, bounds for probable ranges of coefficients are needed. For all the aforementioned applications, the knowledge of the time series conditional distribution can provide all that needed information. New statistical models capable of handling such difficulties are an emerging field in power systems literature. The main objective in such literature is to propose new models capable of generating scenarios of renewable energy source which are demanded in (i) energy trading, (ii) unit commitment, (iii) grid expansion planning, and (iv) investment decisions (see ([2]–[5]) and references

therein). To provide good scenarios from an array of potential influential factors, one has to properly select which features are relevant and create a good model for the conditional distribution. Notwithstanding, a little attention is devoted to addressing both at the same time.

revisar este paragrafo

Conventional statistical models are often focused on estimating the conditional mean of a given random variable. By reducing the outcome to a single statistic, we loose important informations about the series random behavior. In order to account for the process inherent variability it is important to consider probability forecasting. [6] reviews the commonly used methodologies regarding probabilistic forecasting models, splitting them in parametric and nonparametric classes. Main characteristics of **parametric models** are (i) assuming a distribution shape and (ii) low computational costs. ARIMA-GARCH, for example, model the renewable series by assuming the distribution *a priori*. On the other hand, **nonparametric models** (i) don't require a distribution to be specified, (ii) needs mode data to produce a good approximation and (iii) have a higher computational cost. Popular methods are Quantile Regression (QR), Kernel Density Estimation, Artificial Intelligence or a mix of them.

Most time series methods rely on the assumption of Gaussian errors. However, renewable series such as wind and solar are reported as non-Gaussian [7]–[10]. To circumvent this problem, the usage of nonparametric methods - which doesn't rely on assuming any previously assumed distribution - is adequate. Quantile Regression (QR) is a tool for constructing a methodology for non-gaussian time series, because of its facility to implement on commercial solvers and to extend the original model. However, when estimating a distribution function, as each quantile is estimated independently, the monotonicity of the distribution function may be violated. This issue can be addressed by constraining the sequence of quantiles to be in an increasing order. Other possibility is making a transformation afterwards, as shown in [11].

colocar parte sobre nossa solução dos crossing quantiles

The seminal work [12] defines QR as we use today. By this formulation, the conditional quantile is the solution of an optimization problem where we minimize the sum of the check function (defined formally in the next session). Instead of using the classical regression to estimate the conditional mean, the QR determines any quantile from the conditional distribution. Applications are enormous, ranging from risk measuring at financial funds (the Value-at-Risk) to a central measure robust to outliers. By estimating many quantiles on a thin grid of probabilities, one can have as many points as

desired from the estimated conditional distribution function. In [13], the application of QR is extended to time series, when the covariates are lagged values of y_t . In our work, beyond autoregressive terms, it is also considered other exogenous variables as covariates.

In [14]–[18], QR is employed to model the conditional distribution of Wind Power Time Series. An updating quantile regression model is presented by [15]. The authors present a modified version of the simplex algorithm to incorporate new observations without restarting the optimization procedure. In [16], the authors build a quantile model from already existent independent Wind Power forecasts. The approach by [14] is to use QR with a nonparametric methodology. The authors add a penalty term based on the Reproducing Kernel Hilbert Space, which allows a nonlinear relationship between the explanatory variables and the output. This paper also develops an on-line learning technique, where the model is easily updated after each new observation. In [18], wind power probabilistic forecasts are made by using QR with a special type of Neural Network (NN) with one hidden layer, called extreme learning machine. In this setup, each quantile is a different linear combination of the features of the hidden layer. The authors of [19] use the weighted Nadaraya-Watson to estimate the conditional function in the time series.

Regularization is a topic already explored in previous QR papers. The work by [20] defines the proprieties and convergence rates for QR when adding a penalty proportional to the ℓ_1 -norm to perform variable selection, using the same idea as the LASSO [21]. The ADALASSO equivalent to QR is proposed by [22]. In this variant, the penalty for each variable has a different weight, and this modification ensures that the oracle propriety is being respected.

an optimal specification methodology regularization

We propose using Quantile Autoregression (QAR) to create a methodology capable of estimating and simulating a nongaussian time series, such as renewable energy source. By estimating a regularized QAR we model the conditional quantile function. For the best of the authors knowledge, no other work has developed a methodology where regularization and estimation of the conditional distribution using QR is carried on at the same time. We propose to attack both problems simultaneously by using either Mixed Integer Linear Programming (MILP) or a LASSO penalization. On the LASSO formulation, regularization is performed for an individual quantile as described in [20], with the difference that all quantiles are estimated at the same time. In [23], the best subset with size K is selected by solving a MILP problem to minimize the sum of squared errors. The idea is straightforward: integer variables are used to count whether a variable is included or not in the model; a total number of K variables is allowed. Model selection for QR is performed using this same approach. The advantage we highlight on using the latter methodology is that the solution provided is optimal in the sense of minimizing the check function for a given number K of variables.

The objective of this paper is, then, to propose a new methodology to address nonparametric time-series focused on

renewable energy. This may be seen as a multiple quantile regression that specifies a time series model based on the empirical conditional distribution. The main contributions are:

- A nonparametric methodology to model the conditional distribution of time series.
- We propose a parsimonious regularized based methodology that selects the global optimal solution.
- Regularization techniques applied to an ensemble of quantile functions to estimate the conditional distribution, solving the issue of non-crossing quantiles.

The remaining of the paper is organized as follows. In section II, we present both the linear parametric and the nonlinear QR based time series models. In section III, we discuss the estimation procedures for them. The regularization strategies are also presented on this section. Finally, in section IV, a case study using real data from both solar and wind power is presented in order to test our methodology. Section V will conclude this article.

II. QUANTILE REGRESSION BASED TIME SERIES MODEL

Let the α -conditional quantile function of Y for a given value x of the d -dimensional random variable X , i.e., $Q_{Y|X} : [0, 1] \times \mathbb{R}^d \rightarrow \mathbb{R}$, can be defined as

$$Q_{Y|X}(\alpha, x) = F_{Y|X}^{-1}(\alpha, x) = \inf\{y : F_{Y|X}(y, x) \geq \alpha\}. \quad (1)$$

Let a dataset be composed of n observations of $\{y_t, x_t\}_{t=1}^n$. The sample quantile function is based on a finite number of observations and is the solution to the following optimization problem:

$$\hat{Q}_{Y|X}(\alpha, \cdot) \in \arg \min_{q_\alpha} \sum_{t \in T} \rho_\alpha(y_t - q_\alpha(x_t)), \quad (2)$$

$$q_\alpha \in \mathcal{Q}. \quad (3)$$

where ρ is the check function, defined as

$$\rho_\alpha(x) = \begin{cases} \alpha x & \text{if } x \geq 0 \\ (1 - \alpha)x & \text{if } x < 0 \end{cases}. \quad (4)$$

Quantile q_α belongs to a function space \mathcal{Q} . We might have different assumptions for space \mathcal{Q} , depending on the type of function we want to find for q_α . A few properties, however, must be achieved by our choice of space, such as being continuous and having limited first derivative. In this paper, we consider the case where \mathcal{Q} is a linear function's space.

The problem (2)-(3) can be rewritten as a Linear Programming problem as in (5)-(8), thus being able to use a modern solver to fit our model. Variables ε_t^+ e ε_t^- represent the quantities $|y - q(\cdot)|^+$ and $|y - q(\cdot)|^-$, respectively. A is the set containing a sequence of probabilities α_i such that $0 < \alpha_1 < \alpha_2 < \dots < \alpha_Q < 1$. This set represents a finite

discretization of the interval $[0, 1]$.

$$\min_{\beta_{0\alpha}, \beta_{\alpha}, \varepsilon_{t\alpha}^+, \varepsilon_{t\alpha}^-} \sum_{\alpha \in A} \sum_{t \in T} (\alpha \varepsilon_{t\alpha}^+ + (1 - \alpha) \varepsilon_{t\alpha}^-) \quad (5)$$

s.t.

$$\varepsilon_{t\alpha}^+ - \varepsilon_{t\alpha}^- = y_t - \beta_{0\alpha} - \beta_{\alpha}^T x_t, \quad \forall t \in T, \forall \alpha \in A, \quad (6)$$

$$\varepsilon_{t\alpha}^+, \varepsilon_{t\alpha}^- \geq 0, \quad \forall t \in T, \forall \alpha \in A, \quad (7)$$

$$\beta_{0\alpha} + \beta_{\alpha}^T x_t \leq \beta_{0\alpha'} + \beta_{\alpha'}^T x_t, \quad \forall t \in T, \forall (\alpha, \alpha') \in A \times A, \alpha < \alpha', \quad (8)$$

After solving the problem, the sequence $\{q_{\alpha}\}_{\alpha \in A}$ is fully defined by the optimum values $\beta_{0\alpha}^*$ and β_{α}^* , for every α .

We apply QR to estimate the conditional distribution $\hat{Q}_{Y_{t+k}|X_{t+k}, Y_t, Y_{t-1}, \dots}(\alpha, \cdot)$ for a k -step ahead forecast of time series $\{y_t\}$, where X_{t+k} is a vector of exogenous variables at the time we want to forecast. Once the conditional distribution is estimated, we are able to simulate and generate scenarios. In the next session, regularization techniques are presented, in order to choose parsimoniously which variables will be input for \hat{Q} .

III. REGULARIZATION

When dealing with many candidates to use as covariates, one has to deal with the problem of selecting a subset of variables to use in constructing the model. This means that the vector of coefficients $\beta_{\alpha} = [\beta_{1\alpha} \cdots \beta_{P\alpha}]$ should not have all nonzero values. There are many ways of selecting a subset of variables among the available options. Classical approaches for this problem are the Stepwise algorithm [24], [25], [21], which includes variables in sequence.

Two approaches will be employed. At first, we use a Mixed Integer Linear Programming optimization problem (MILP) to find the best subset among all choices of covariates. The second way is by using a LASSO-type technique, which consists in penalizing the ℓ_1 -norm of regressors, thus shrinking the size of estimated coefficients towards zero.

A. Best subset selection via MILP

We use MILP to select variables by including constraints which limits their number in K . Only K coefficients $\beta_{p\alpha}$ may have nonzero values, for each α . Binary variable $z_{p\alpha}$ indicates whether $\beta_{p\alpha}$ has a nonzero value. The optimization problem

that incorporates this idea is described below:

$$\min_{\beta_{0\alpha}, \beta_{\alpha}, z_{p\alpha}, \varepsilon_{t\alpha}^+, \varepsilon_{t\alpha}^-} \sum_{\alpha \in A} \sum_{t \in T} (\alpha \varepsilon_{t\alpha}^+ + (1 - \alpha) \varepsilon_{t\alpha}^-) \quad (9)$$

s.t.

$$\varepsilon_{t\alpha}^+ - \varepsilon_{t\alpha}^- = y_t - \beta_{0\alpha} - \sum_{p=1}^P \beta_{p\alpha} x_{t,p}, \quad \forall t \in T, \forall \alpha \in A, \quad (10)$$

$$\varepsilon_{t\alpha}^+, \varepsilon_{t\alpha}^- \geq 0, \quad \forall t \in T, \forall \alpha \in A, \quad (11)$$

$$-M z_{p\alpha} \leq \beta_{p\alpha} \leq M z_{p\alpha}, \quad \forall \alpha \in A, \forall p \in P, \quad (12)$$

$$\sum_{p=1}^P z_{p\alpha} \leq K, \quad \forall \alpha \in A, \quad (13)$$

$$z_{p\alpha} \in \{0, 1\}, \quad \forall \alpha \in A, \forall p \in P, \quad (14)$$

$$\beta_{0\alpha} + \beta_{\alpha}^T x_t \leq \beta_{0\alpha'} + \beta_{\alpha'}^T x_t, \quad \forall t \in T, \forall (\alpha, \alpha') \in A \times A, \alpha < \alpha', \quad (15)$$

The objective function and constraints (10), (11) and (15) are the same from standard linear quantile regression. By constraint (12), variable $z_{p\alpha}$ is a binary that assumes 1 when coefficient $\beta_{p\alpha}$ is included, while (13) guarantees that at most K of them are nonzero. The value of M is chosen in order to guarantee that $M \geq \|\hat{\beta}_{\alpha}\|_{\infty}$. The solution given by $\beta_{0\alpha}^*$ and $\beta_{\alpha}^* = [\beta_{1\alpha}^* \cdots \beta_{P\alpha}^*]$ will be the best linear α -quantile regression with K nonzero coefficients.

Defining groups for variables: Consider the optimization problem defined on (9)-(15). The choice of the best subset is independent for different values of α . This means that the best subset may include two completely different sets of regressors for two probabilities α and α' close to each other. Take $K = 2$ for the example, selecting $\beta_{1\alpha}$ and $\beta_{4\alpha}$ for α while $\beta_{2\alpha'}$ and $\beta_{5\alpha'}$ is possible, but unlikely to be true.

To address this issue, we propose to divide all $\alpha \in A$ in groups. The collection G of all groups g form a partition of A , and each α belongs to exactly one group g . The subset of selected covariates must be the same for all α in the same group g . To model these properties as constraints on problem (9)-(15), we substitute constraint (12) for the following equations:

$$z_{p\alpha g} := 2 - (1 - z_{pg}) - I_{g\alpha} \quad (16)$$

$$\sum_{g \in G} I_{g\alpha} = 1, \quad \forall \alpha \in A, \quad (17)$$

$$-M z_{p\alpha g} \leq \beta_{p\alpha g} \leq M z_{p\alpha g}, \quad \forall p \in P, \forall \alpha \in A, \forall g \in G, \quad (18)$$

$$I_{g\alpha}, z_{pg} \in \{0, 1\}, \quad \forall p \in P, \forall g \in G, \quad (19)$$

on problem (9)-(15). where G is a set of group index and z_{pg} is a binary variable that equals 1 iff covariate p is included on group g and $I_{g\alpha}$ equals 1 iff probability α belongs to group g . Constraint (18) forces that

$$\text{if } z_{pg} = 0 \text{ and } I_{g\alpha} = 1 \text{ then } \beta_{p\alpha} = 0.$$

Hence, if covariate p belongs to group g , this covariate is not among group's g subset of variables, than its coefficient must

be equal to 0, for that α . Note that variable $z_{p\alpha}$ behaves differently than when we are not considering groups. This means that if probability α belongs to group g but variable p is not selected to be among the ones of group g , then $\beta_{p\alpha}$ is zero. Equation (16) defines $z_{p\alpha}$ to simplify writing.

B. Variable selection via LASSO

Another way of doing regularization is including the ℓ_1 -norm of the coefficients on the objective function. In [20], the reader can find properties and convergence rate when using the LASSO to select variables in a quantile regression setting. The ADALASSO variant is presented in [22]. The advantage of this method is that coefficients are shrunk towards zero by changing a continuous parameter λ , which penalizes the size of the ℓ_1 -norm. When the value of λ gets bigger, fewer variables are selected to be used. This is the same strategy of the LASSO methodology, and its usage for the quantile regression is discussed in [26]. The proposed optimization problem to be solved is:

$$\min_{\beta_{0\alpha}, \beta_{p\alpha}} \sum_{t \in T} \alpha |y_t - q_\alpha(x_t)|^+ + \sum_{t \in T} (1-\alpha) |y_t - q_\alpha(x_t)|^- + \lambda \|\beta_\alpha\|_1, \quad (20)$$

$$q_\alpha(x_t) = \beta_0 - \sum_{p=1}^P \beta_p x_{t,p}.$$

For such estimation to be coherent, however, each covariate must have the same relative weight in comparison with one another. So, before solving the optimization problem, we perform a linear transformation such that all variables have mean $\mu = 0$ and variance $\sigma^2 = 1$. We apply the transformation $\tilde{x}_{t,p} = (x_{t,p} - \bar{x}_{t,p}) / \hat{\sigma}_{x_{t,p}}$, where $\bar{x}_{t,p}$ and $\hat{\sigma}_{x_{t,p}}$ are respectively the sample's unconditional mean and standard deviation. The $\tilde{y}_{t-p,i}$ series will be used to estimate the coefficients, as this series has the desired properties.

The process is done in two stages: variable selection and coefficients estimation. At first, all covariates are input on the following optimization problem:

$$\begin{aligned} \tilde{\beta}_\lambda^{*LASSO} = \arg \min_{\beta_0, \beta, \varepsilon_{t\alpha}^+, \varepsilon_{t\alpha}^-} \sum_{\alpha \in A} \sum_{t \in T} (\alpha \varepsilon_{t\alpha}^+ + (1-\alpha) \varepsilon_{t\alpha}^-) \\ + \lambda \sum_{p=1}^P \xi_{p\alpha} \end{aligned} \quad (21)$$

subject to

$$\varepsilon_{t\alpha}^+ - \varepsilon_{t\alpha}^- = y_t - \beta_{0\alpha} - \sum_{p=1}^P \beta_{p\alpha} \tilde{x}_{t,p}, \quad \forall t \in T, \forall \alpha \in A, \quad (22)$$

$$\varepsilon_{t\alpha}^+, \varepsilon_{t\alpha}^- \geq 0, \quad \forall t \in T, \forall \alpha \in A, \quad (23)$$

$$\xi_{p\alpha} \geq \beta_{p\alpha}, \quad \forall p \in P, \forall \alpha \in A, \quad (24)$$

$$\xi_{p\alpha} \geq -\beta_{p\alpha}, \quad \forall p \in P, \forall \alpha \in A. \quad (25)$$

This model is built upon the standard linear programming model for the quantile regression (5)-(8). On the above formulation, the ℓ_1 norm of equation (20) is substituted by the sum of ξ_p , which represents the absolute value of $\beta_{p\alpha}$. The

link between variables ξ_p and $\beta_{p\alpha}$ is made by constraints (24) and (25). Note that the linear coefficient $\beta_{0\alpha}$ is not included in the penalization, as the sum of penalties on the objective function 21.

For low values of λ , the penalty over the size of coefficients is small. Because of that, the output of problem (21)-(25) is a model where most coefficients have nonzero value. On the other hand, when the penalty on $\|\beta_\alpha\|_1$ is big, many covariates will have zero valued coefficients. When λ approaches infinity, one has a constant model. For instance, the penalty isn't applied to the linear coefficient $\beta_{0\alpha}$.

In fact, the LASSO coefficients are biased, so it is employed only as a variable selector. The optimum vector of coefficients $\tilde{\beta}_\lambda^{*LASSO}$ for a given λ may be composed by both nonzero and zero coefficients. We then define L_λ as the set of indexes of selected variables given by

$$L_\lambda = \{p \in \{1, \dots, P\} \mid |\beta_{\lambda,p}^{*LASSO}| \neq 0\}.$$

Hence, we have that, for each $p \in \{1, \dots, P\}$,

$$\beta_{\lambda,p}^{*LASSO} = 0 \implies \beta_{\lambda,p}^* = 0.$$

Note that problem (5)-(8) is employed to act as variable selection only. On the second stage, the optimal coefficient vector β_λ^{*LASSO} is estimated by the non-regularized QR, where only variables that belongs to L_λ are input:

$$(obj_\lambda^*, \beta_\lambda^*) \xleftarrow{(obj, var)} \min_{\beta_0, \beta, \varepsilon_t^+, \varepsilon_t^-} \sum_{t \in T} (\alpha \varepsilon_t^+ + (1-\alpha) \varepsilon_t^-) \quad (26)$$

subject to

$$\varepsilon_t^+ - \varepsilon_t^- = y_t - \beta_0 - \sum_{p \in L_\lambda} \beta_p x_{t,p}, \quad \forall t \in T, \quad (27)$$

$$\varepsilon_t^+, \varepsilon_t^- \geq 0, \quad \forall t \in T. \quad (28)$$

The variable obj_λ^* receives the value of the objective function on its optimal solution. In summary, the optimization in equation 20 acts as a variable selection for the subsequent estimation, which is normally called the post-LASSO estimation [27].

IV. ESTIMATION

A. Time-series cross validation

Sections III-A and III-B presented two different methods to estimate the conditional distribution in a parsimonious way. However, as presented, the aforementioned methods don't provide a unique solution, but a set of solutions for a range of tuning parameters. For instance, on the MILP method, the quantity K of nonzero coefficients is an input of the problem. Similarly, the LASSO needs a penalization parameter λ , that tunes how much penalty the ℓ_1 -norm receives.

In statistics and machine learning, a popular technique is using Cross-validation (CV) to select the best model from this range of possibilities. It is a technique used to have an estimate of the model's quality of prediction in an independent testing set. The best model that minimizes the CV error is the model which presumably will have the best performance on out of sample data.

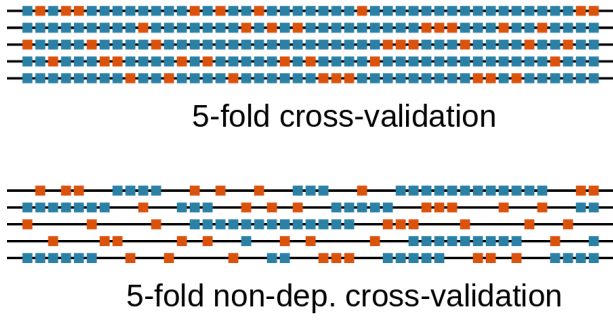


Fig. 1. K -fold CV and K -fold with non-dependent data. Observations in blue are used to estimation and in orange for evaluation. Note that non-dependent data doesn't use all dataset in each fold.

The usage of CV is not straightforward when data is dependent, which is the case when working with time series. As the data is time dependent, one can be interested in using either all observations or to take the dependency away. The works [30] and [31] deals specifically with the usage of CV in a time series context. They provide tests with both K -fold CV and K -fold with non-dependent data. Both schemes are shown of Figure 1. In both settings, the training data is randomly split into a collection of sets J_k , forming a K size partition. Each of these J_k is used as test set, while

terminar parágrafo de explicação sobre CV

B. Information Criteria for Quantile Regression

Sometimes, using CV can be computationally expensive, as the full estimation is done several times for each tuning parameter - in this case, either K or λ . Other form of deciding the quantity of variables that provides a good equilibrium between in-sample prediction and parsimony is the Information Criteria.

Information criteria summarizes two aspects. One of them refers to how well the model fits the in-sample observations and the other part penalizes the quantity of covariates used in the model. By penalizing how big our model is, we prevent overfitting from happening. So, in order for a covariate to be included in the model, it must supply enough goodness of fit. In [29], it is presented a variation of the Schwarz criteria for M-estimators that includes quantile regression. The Schwarz Information Criteria (SIC), adapted to the quantile autoregression case, is presented below:

$$SIC(m) = n \log(obj^*) + \frac{1}{2} K \log n, \quad (29)$$

where K is the model's dimension. This procedure leads to a consistent model selection if the model is well specified.

By minimizing the SIC

acabar

C. Model selection distance

Solving a LP problem such as the LASSO is many times faster than a similar-sized MILP problem, for introducing

binary variables breaks the problem's convexity. On the other hand, in our case the MILP solution is the exact best solution in minimizing the QR objective function, while the LASSO is an approximation of that.

One of our goals is to test how far from the optimal solution is the LASSO. For each number K of total nonzero coefficients, there will be a penalty λ_K^* which minimizes the errors from the quantile regression's objective function (given on equation (26)):

$$\lambda_K^* = \arg \min_{\lambda} \{obj_{\lambda}^* \mid \|\beta_{\lambda}^*\|_0 = K\}, \quad (30)$$

where the quantity $\|\beta_{\lambda}^*\|_0$ is the 0-norm, which gives the total of nonzero coefficients, for a given lambda of the LASSO estimations.

We, then, define the sets L_K^{LASSO} and L_K^{MILP} , which contains all nonzero indexes, for a given K , when using methods LASSO and MILP for regularization, respectively. Thus, we can compare the best LASSO fit where exactly K variables are selected with the best fit given by the MILP problem, also with K variables selected.

As the MILP solution is the exact solution for the problem, while the LASSO solution is an approximation, we use the former as a *benchmarking* for the quality of the latter solution. It is desirable that the LASSO solution be as related with the MILP solution as possible. The difference in performance is given by a similarity metric d , which measures distance from solutions weighted by the correlation between variables. The similarity is calculated as the solution of the following optimization problem

$$d(\beta_{MILP(K)}^*, \beta_{\lambda_K^*}^*) = \min_{0 \leq \delta_{ij} \leq 1} \sum_{i,j=1}^K \delta_{ij} (1 - |\rho_{ij}|) \quad (31)$$

subject to

$$\sum_{j=1}^K \delta_{ij} = 1, \quad i = 1, \dots, K, \quad (32)$$

$$\sum_{i=1}^K \delta_{ij} = 1, \quad j = 1, \dots, K, \quad (33)$$

where ρ_{ij} is the correlation between the i -th and j -th independent variables in sets L_k^{MILP} and L_k^{LASSO} , respectively. The optimal value for the decision variables of this problem provides us with an assignment between selected covariates from both methods, namely, MILP and LASSO, that minimizes the overall "index of uncorrelation" between selected covariates. If $\delta_{ij}^* = 1$, the i -th selected variable in L_k^{MILP} is associated with the j -th variable in L_k^{LASSO} . For instance, if $d(\beta_{MILP(K)}^*, \beta_{\lambda_K^*}^*) = 0$, it means that there are K perfectly correlated pair of variables, even though not being the same subset.

D. Evaluation criteria

The full dataset is split between the test set - which evaluates our methodology's forecasting performance - and the training

set - which we use to estimate parameters. This setting mimics real world applications, where the future is unknown.

As conditional distribution is the focus in this paper, we use a performance measurement which emphasizes the correctness of each quantile. For each observation y_t and probability $\alpha \in A$, a score function is defined by

$$L(A) = \sum_{t \in T} \rho_\alpha(y_t - q_\alpha(x_t))$$

The error measure is defined as the average of the score function for all observations and target quantiles.

REFERENCES

- [1] International energy agency. [Online]. Available: <https://www.iea.org/newsroom/news/2013/october/wind-power-seen-generating-up-to-18-of-global-power-by-2050.html>
- [2] A. Moreira, D. Pozo, A. Street, and E. Sauma, "Reliable renewable generation and transmission expansion planning: Co-optimizing system's resources for meeting renewable targets," *IEEE Transactions on Power Systems*, vol. PP, no. 99, pp. 1–1, 2016.
- [3] R. Jabr, "Robust transmission network expansion planning with uncertain renewable generation and loads," *IEEE Transactions on Power Systems*, vol. 28, no. 4, pp. 4558–4567, 2013.
- [4] C. Zhao and Y. Guan, "Data-driven stochastic unit commitment for integrating wind generation," *IEEE Transactions on Power Systems*, vol. 31, no. 4, pp. 2587–2596, July 2016.
- [5] A. C. Passos, A. Street, and L. A. Barroso, "A dynamic real option-based investment model for renewable energy portfolios," *IEEE Transactions on Power Systems*, vol. 32, no. 2, pp. 883–895, March 2017.
- [6] Y. Zhang, J. Wang, and X. Wang, "Review on probabilistic forecasting of wind power generation," *Renewable and Sustainable Energy Reviews*, vol. 32, pp. 255–270, Apr. 2014. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S1364032114000446>
- [7] R. J. Bessa, V. Miranda, A. Botterud, J. Wang, and M. Constantinescu, "Time adaptive conditional kernel density estimation for wind power forecasting," *Sustainable Energy, IEEE Transactions on*, vol. 3, no. 4, pp. 660–669, 2012.
- [8] J. Jeon and J. W. Taylor, "Using conditional kernel density estimation for wind power density forecasting," *Journal of the American Statistical Association*, vol. 107, no. 497, pp. 66–79, 2012.
- [9] J. W. Taylor and J. Jeon, "Forecasting wind power quantiles using conditional kernel estimation," *Renewable Energy*, vol. 80, pp. 370–379, 2015.
- [10] C. Wan, J. Lin, J. Wang, Y. Song, and Z. Y. Dong, "Direct quantile regression for nonparametric probabilistic forecasting of wind power generation," *IEEE Transactions on Power Systems*, vol. 32, no. 4, pp. 2767–2778, July 2017.
- [11] V. Chernozhukov, I. Fernández-Val, and A. Galichon, "Quantile and Probability Curves Without Crossing," *Econometrica*, vol. 78, no. 3, pp. 1093–1125, May 2010. [Online]. Available: <http://onlinelibrary.wiley.com/doi/10.3982/ECTA7880/abstract>
- [12] R. Koenker and G. Bassett Jr, "Regression quantiles," *Econometrica: journal of the Econometric Society*, pp. 33–50, 1978.
- [13] R. Koenker, Z. Xiao, J. Fan, Y. Fan, M. Knight, M. Hallin, B. J. M. Werker, C. M. Hafner, O. B. Linton, and P. M. Robinson, "Quantile Autoregression [with Comments, Rejoinder]," *Journal of the American Statistical Association*, vol. 101, no. 475, pp. 980–1006, 2006. [Online]. Available: <http://www.jstor.org/stable/27590777>
- [14] C. Gallego-Castillo, R. Bessa, L. Cavalcante, and O. Lopez-Garcia, "On-line quantile regression in the rkhs (reproducing kernel hilbert space) for operational probabilistic forecasting of wind power," *Energy*, vol. 113, pp. 355–365, 2016.
- [15] J. K. Møller, H. A. Nielsen, and H. Madsen, "Time-adaptive quantile regression," *Computational Statistics & Data Analysis*, vol. 52, no. 3, pp. 1292–1303, Jan. 2008. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0167947307002502>
- [16] H. A. Nielsen, H. Madsen, and T. S. Nielsen, "Using quantile regression to extend an existing wind power forecasting system with probabilistic forecasts," *Wind Energy*, vol. 9, no. 1-2, pp. 95–108, 2006.
- [17] J. B. Bremnes, "Probabilistic wind power forecasts using local quantile regression," *Wind Energy*, vol. 7, no. 1, pp. 47–54, Jan. 2004. [Online]. Available: <http://onlinelibrary.wiley.com/doi/10.1002/we.107/abstract>
- [18] C. Wan, J. Lin, J. Wang, Y. Song, and Z. Y. Dong, "Direct Quantile Regression for Nonparametric Probabilistic Forecasting of Wind Power Generation," *IEEE Transactions on Power Systems*, vol. 32, no. 4, pp. 2767–2778, Jul. 2017.
- [19] Z. Cai, "Regression Quantiles for Time Series," *Econometric Theory*, vol. 18, no. 1, pp. 169–192, 2002. [Online]. Available: <http://www.jstor.org/stable/3533031>
- [20] A. Belloni and V. Chernozhukov, "L1-Penalized Quantile Regression in High-Dimensional Sparse Models," *arXiv:0904.2931 [math, stat]*, Apr. 2009, arXiv: 0904.2931. [Online]. Available: <http://arxiv.org/abs/0904.2931>
- [21] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 267–288, 1996.
- [22] G. Ciuperca, "Adaptive LASSO model selection in a multiphase quantile regression," *Statistics*, vol. 50, no. 5, pp. 1100–1131, Sep. 2016. [Online]. Available: <http://dx.doi.org/10.1080/02331888.2016.1151427>
- [23] D. Bertsimas, A. King, and R. Mazumder, "Best Subset Selection via a Modern Optimization Lens," *arXiv:1507.03133 [math, stat]*, Jul. 2015, arXiv: 1507.03133. [Online]. Available: <http://arxiv.org/abs/1507.03133>
- [24] M. Efronson, "Multiple regression analysis," *Mathematical methods for digital computers*, vol. 1, pp. 191–203, 1960.
- [25] R. R. Hocking and R. N. Leslie, "Selection of the Best Subset in Regression Analysis," *Technometrics*, vol. 9, no. 4, pp. 531–540, Nov. 1967. [Online]. Available: <http://www.tandfonline.com/doi/abs/10.1080/00401706.1967.10490502>
- [26] Y. Li and J. Zhu, "L1-norm quantile regression," *Journal of Computational and Graphical Statistics*, 2012.
- [27] A. Belloni and V. Chernozhukov, "Least squares after model selection in high-dimensional sparse models," 2009.
- [28] B. Sherwood, A. Maidman, M. B. Sherwood, and T. ByteCompile, "Package 'rqpen'," *Penalized Quantile Regression. In*, vol. 1, 2017.
- [29] J. A. Machado, "Robust model selection and m-estimation," *Econometric Theory*, vol. 9, pp. 478–493, 1993.
- [30] C. Bergmeir, R. J. Hyndman, and B. Koo, "A Note on the Validity of Cross-Validation for Evaluating Time Series Prediction," Monash University, Department of Econometrics and Business Statistics, Tech. Rep. 10/15, 2017. [Online]. Available: <https://ideas.repec.org/p/msh/ebswps/2015-10.html>
- [31] C. Bergmeir and J. M. Benítez, "On the use of cross-validation for time series predictor evaluation," *Information Sciences*, vol. 191, pp. 192–213, May 2012. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0020025511006773>