

Forecasting Value-at-Risk Using Deep Neural Network Quantile Regression*

Ilias Chronopoulos ¹, Aristeidis Raftapostolos ², and George Kapetanios²

¹Essex Business School, University of Essex and ²King's Business School, King's College London
Address correspondence to Aristeidis Raftapostolos, King's Business School, King's College London, Bush House, 30 Aldwych, London WC2B 4BG, UK, or e-mail: aristeidis.1.raftapostolos@kcl.ac.uk.

Received April 5, 2021; revised February 6, 2023; editorial decision April 1, 2023; accepted April 13, 2023

Abstract

In this article, we use a *deep quantile* estimator, based on neural networks and their universal approximation property to examine a non-linear association between the conditional quantiles of a dependent variable and predictors. This methodology is versatile and allows both the use of different penalty functions, as well as high dimensional covariates. We present a Monte Carlo exercise where we examine the finite sample properties of the deep quantile estimator and show that it delivers good finite sample performance. We use the deep quantile estimator to forecast value-at-risk and find significant gains over linear quantile regression alternatives and other models, which are supported by various testing schemes. Further, we consider also an alternative architecture that allows the use of mixed frequency data in neural networks. This article also contributes to the interpretability of neural network output by making comparisons between the commonly used Shapley Additive Explanation values and an alternative method based on partial derivatives.

Key words: forecasting, machine learning, neural networks, quantile regression, value-at-risk

JEL classification: C45, C58, G17

Since the seminal work of [Koenker and Bassett \(1978\)](#) and [Koenker and Hallock \(2001\)](#), quantile regression has grown in popularity and has found applications in several disciplines both in academia and industry, see for example, [Chernozhukov and Umantsev \(2001\)](#); [Adams et al. \(2021\)](#); and [Koenker et al. \(2017\)](#). They generalize ordinary sample quantiles to the regression setting, that give more extensive information on the conditional distribution of a dependent variable, given the covariates, relative to the classical regression setting; that is, estimation of the conditional mean. This extension can be of great

* We thank the Editor and two anonymous referees for constructive comments and valuable suggestions, which led to significant improvements on the article. An earlier version of this article was circulated under the title “Deep Quantile Regression.”

importance under extreme events, where the conditional distribution of variables such as asset returns tends to exhibit skewness, or under the presence of outliers and/or asymmetries, see for example, [Baur and Schulze \(2005\)](#).

An assumption made in the early literature, was the linear association between the conditional quantile of the target variable and predictors. This was predominately an assumption that allowed for streamlined computation and theoretical inference, but was clearly restrictive. A more recent strand of the literature, relaxed the linearity assumption and considered non-parametric estimators for the conditional quantile, that is, based on different methods, see for example, [Belloni et al. \(2019\)](#) and references therein. Recent advances in machine learning (ML) literature, which is the focus of this article, show how modeling frameworks such as neural networks can be used to estimate general, non-linear, and potentially highly complicated associations.

Specifically, a large number of studies have shown that *feed-forward* neural networks can approximate arbitrarily well any continuous function of several real variables, see for example, [Hornik \(1991\)](#); [Hornik, Stinchcombe, and White \(1989\)](#); [Gallant and White \(1992\)](#); and [Park and Sandberg \(1991\)](#). Recent works by [Liang and Srikant \(2016\)](#) and [Yarotsky \(2017\)](#) extend this result for feed-forward neural networks with multiple layers, provided sufficiently many hidden neurons and layers are available. Notice that, besides neural networks, other non-parametric approaches, for example, splines, wavelets, the Fourier basis, as well as simple polynomial approximations, do have the universal approximation property, based on the Stone–Weierstrass theorem.

There is considerable empirical work identifying non-linearities and asymmetries in financial variables, see for example, [Gu, Kelly, and Xiu \(2020, 2021\)](#); [He and Krishnamurthy \(2013\)](#); and [Pohl, Schmedders, and Wilms \(2018\)](#), where they illustrate that ML offers richer functional form specifications that can capture potential non-linearities between dependent and independent variables. Some examples include [Gu, Kelly, and Xiu \(2020\)](#) in which, they evaluate the forecast accuracy of ML methods in measuring equity risk premia, and find that neural networks give substantial forecasting gains in asset pricing compared to linear models, and [Bucci \(2020\)](#), where a recurrent neural network is proposed, that approximates realized volatility well and outperforms other classic non-linear estimators in forecasting. In a similar fashion, [Smalter Hall and Cook \(2017\)](#) use several neural network architectures to predict unemployment in the United States and find that neural networks outperform forecasts from a linear benchmark model at short horizons. In addition, [Gu, Kelly, and Xiu \(2021\)](#) propose the use of a conditional Autoencoder,¹ and illustrate its superior performance relative to linear unsupervised learning methods.

Before we discuss the contributions of this article, we provide a succinct summary of the current ML literature on non-linear quantile and value-at-risk (VaR) estimation, but we note that the majority of this work, was not available during the writing of this article. [Keilbar and Wang \(2022\)](#) use neural networks to estimate a non-linear conditional VaR model introduced by [Tobias and Brunnermeier \(2016\)](#) and find that, it gives significant gains in modeling systemic risk. In addition, [Tambwekar et al. \(2022\)](#) estimate a non-linear binary quantile regression and develop confidence scores to assess the reliability of prediction. [Padilla, Tansey, and Chen \(2022\)](#) examine the performance of a quantile neural

1 Autoencoders are artificial neural networks that can be used as a dimensionality reduction technique.

network using the rectified linear unit (ReLU) as activation function and show that it has superior statistical performance relative to other quantile regression methods. [Chen et al. \(2020\)](#) propose a unified non-linear framework, based on feed-forward neural networks, that allows the estimation of treatment effects, for which they establish consistency and asymptotic normality. Their framework includes the quantile estimator and allows for high-dimensional covariates. ML-based estimators for quantiles have been proposed in other fields, see for example, [Meinshausen \(2006\)](#), where quantile random forests are introduced, and [Zhang, Quan, and Srinivasan \(2019\)](#) that propose a quantile neural network estimator.

In this article, we contribute to the expanding literature on the use of ML in Finance and use a deep quantile estimator that can capture non-linear associations between asset returns and predictors to forecast VaR. Note that this estimator also allows for high dimensional data. We further consider an alternative architecture that allows the use of mixed frequency data. We also contribute toward the explainable ML literature, by proposing the use of partial derivatives as a means to “peeking” inside the black box.

We first explore the small sample properties of the deep quantile estimator via Monte Carlo (MC) experiments, which show that the estimator delivers good finite sample performance. Then we examine the performance of the deep quantile estimator, in the context of one of the most widely examined problems in finance: that of measuring and subsequently forecasting the risk of a portfolio adequately, via VaR modeling. VaR is a popular model that was first introduced in the late 80s and since then, has become a standard toolkit in measuring market risk. It measures how much value a portfolio can lose within a given time period with some small probability, τ . VaR and quantiles are related in the following manner, let $\mathbf{r} = (r_1, \dots, r_T)'$ denote the returns of a portfolio, then, the τ th VaR is equivalent of computing the negative value of the τ th quantile of \mathbf{r} , $-q_\tau(\mathbf{r})$.

In this article, we argue, following the non-parametric literature, that the linear relationship between VaR and predictors can be restrictive and use the deep quantile neural network estimator that allows a non-linear association between covariates and VaR. This method appears particularly suitable for developing sound predictions for the past stock return losses in the United States over the sample period from September 1985 up to August 2020, the importance of which has been brought to the forefront by the recent COVID-19 pandemic. Specifically, our aim is to forecast 10-day ahead VaR produced from daily VaR forecasts. We use daily frequency returns in a fixed forecasting framework that is outlined below.

Under this forecasting framework, mixed frequency models become relevant benchmarks to the non-linear quantile estimator, see for example, [Ghysels, Plazzi, and Valkanov \(2016\)](#). Hence, we also include a linear MIXed DATA Sampling (MIDAS) model as a competitor and also a non-linear MIDAS model, which is an extension to the deep quantile estimator. Further, we consider 10-day compounded VaR forecasts that exhibit similar patterns, which we relegate to the [Supplementary Appendix](#).

We are not the first to use ML methods for VaR forecasting, see for example, [Du, Wang, and Xu \(2019\)](#), where they propose a recurrent neural network, as a forecasting methodology for the VaR model and exhibit an improved forecast performance relative to traditional methods. To the best of our knowledge though, there has been no application that uses a neural network quantile estimator in finance for forecasting VaR. Note that in this article, we consider a set of neural networks that allows for mixed frequency

estimation, following the current literature see for example, [Xu et al. \(2021\)](#); [Borup, Rapach, and Schütte \(2022\)](#); and [Babii, Ghysels, and Striaukas \(2022b\)](#).

Our empirical analysis shows that the deep quantile estimator outperforms the linear, MIDAS, and other non-parametric quantile models, in forecasting VaR. We assess the forecasting accuracy between models based on two statistical tests. The first is the [Diebold and Mariano \(1995\)](#) (DM) test with the [Harvey, Leybourne, and Newbold \(1997\)](#) adjustment, and the second is the [Giacomini and White \(2006\)](#) (GW) test. Results from both tests suggest that the neural network estimator has higher accuracy in forecasting VaR. We use the linear quantile method as a benchmark to assess whether the deep quantile estimator has predictive gains or not. This measure illustrates gains up to 98% relative to the linear one, for the deep quantile estimator and up to 84% for the non-linear MIDAS model. Further, we use the quantile score test that provides further evidence in favor of the neural network estimator.

We further examine whether the deep quantile estimator nests forecasts produced from the linear and other non-parametric models, using the encompassing test of [Giacomini and Komunjer \(2005\)](#). Overall, we find that forecasts from the deep quantile estimator encompass forecasts from competing models more times than vice versa. There are some cases where the test is inconclusive, suggesting that a forecast combination from a different pair of models would provide a better result, which is in line with the result of [Bates and Granger \(1969\)](#).

While ML methods show a great capacity at both approximating highly complicated non-linear functions and forecasting, they are routinely criticized as they lack interpretability and are considered a “black box”; in the sense that they do not offer simple summaries of relationships in the data. Recently though, there has been a number of studies that try to make ML output interpretable, see for example, [Athey and Imbens \(2017\)](#); [Wager and Athey \(2018\)](#); [Belloni, Chernozhukov, and Hansen \(2014\)](#); and [Joseph \(2019\)](#). In this article, we also try to understand in a semi-structural fashion, which variables impact the forecasting performance of the deep quantile estimator more. To this end, we first use Shapley Additive Explanation (SHAP) values as proposed by [Lundberg and Lee \(2017\)](#) and further developed by [Joseph \(2019\)](#), that have started to become a standard tool for interpretability in ML methods. Further, we use partial derivatives as a means of investigating the marginal contribution/influence of each variable to the output. We compare the partial derivatives and SHAP values over time, and our results can be summarized as follows. First, partial derivatives overall are more stable than SHAP values, and are able to produce interpretable results, at a fraction of the computational time of SHAP. Second, the partial derivatives of the deep quantile estimator fluctuate around the estimate of the conditional linear quantile and (i) exhibit time variation and (ii) can capture stressful events in the U.S. economy for instance the COVID-19 pandemic and the 2008 financial crisis.

The remainder of the article is organized as follows. Section 1 introduces the deep quantile estimator. Section 2 contains the MC exercise. Section 3 presents our empirical application. Section 4 presents the semi-structural analysis. Conclusions are set out in Section 5. We relegate to the [Supplementary Appendix](#) the specifications of the competing models, empirical results from a 1-day ahead VaR forecasting exercise, 10-day compounded VaR forecasts, GW test, and results from the quantile score test and predictive gains.

1 Empirical Methodology

In this section, we start by summarizing the underlying theory of a quantile regression as outlined by [Koenker and Bassett \(1978\)](#) and [Koenker \(2005\)](#) and argue that the linear relationship of the conditional quantile between a dependent variable given the covariates, can be restrictive. We illustrate how some fundamental results on the universal approximation property of neural networks can be used to approximate a non-linear relationship instead, and define the deep quantile estimator. We conclude with a discussion on how different penalization schemes can be used and further how hyper-parameters can be selected via time-series cross validation (CV).

1.1 Linear Quantile Regression

The standard goal in econometric analysis is to infer a relationship between a dependent variable and one or more covariates. We consider the following linear regression model:

$$y_t = \mathbf{x}_t' \boldsymbol{\beta} + u_t, \quad (1)$$

where y_t is the dependent variable at time t , $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$ is a vector of unobserved slope parameters, $\mathbf{x}_t = (x_{t1}, \dots, x_{tp})'$ is a vector of known covariates, and u_t is the random error of the regression which satisfies $E(u_t | \mathbf{x}_t) = 0$. Standard regression analysis tries to come up with an estimate of the conditional mean of y_t given \mathbf{x}_t , that minimizes the expected squared error loss:

$$\hat{\boldsymbol{\beta}} = \arg \min_{\boldsymbol{\beta}} \frac{1}{T} \sum_{t=1}^T (y_t - \mathbf{x}_t' \boldsymbol{\beta})^2. \quad (2)$$

This can be restrictive though, when (i) non-linearities and outliers exist and (ii) since it provides just an aspect of the conditional distribution of y_t , given \mathbf{x}_t by construction. These potential limitations led to the development of quantile regression. In their seminal work, [Koenker and Bassett \(1978\)](#) generalize ordinary sample quantiles to the regression setting, that give more complete information on the conditional distribution of y_t given \mathbf{x}_t , for which we now provide a succinct description.

The quantile regression model can be defined as

$$Q_y(\tau | \mathbf{x}_t) = \mathbf{x}_t' \boldsymbol{\beta}(\tau), \quad \tau \in (0, 1), \quad (3)$$

such that y_t satisfies the quantile constraint $Pr[y_t \leq \mathbf{x}_t' \boldsymbol{\beta}(\tau) | \mathbf{x}_t] = \tau$, where $\boldsymbol{\beta}(\tau)$ are regression coefficients that depend on τ . Quantile regression tries to come up with an estimate for the τ th conditional quantile, $\hat{Q}_y(\tau, \mathbf{x}_t) := \hat{\boldsymbol{\beta}}(\tau)$, by minimizing the following function

$$\hat{\boldsymbol{\beta}}(\tau) = \arg \min_{\boldsymbol{\beta}} \frac{1}{T} \sum_{t=1}^T \rho_{\tau}(y_t - \mathbf{x}_t' \boldsymbol{\beta}(\tau)), \quad (4)$$

where $\rho_{\tau}(\cdot)$ is the quantile loss function defined as

$$\rho_{\tau}(u_t) = \begin{cases} \tau u_t(\tau), & \text{if } u_t(\tau) \geq 0 \\ (1 - \tau) u_t(\tau), & \text{if } u_t(\tau) < 0 \end{cases}$$

and $u_t(\tau) = y_t - \mathbf{x}_t' \boldsymbol{\beta}(\tau)$. The quantile estimator in [Equation \(4\)](#), provides (i) much richer information on the whole conditional distribution of y_t as function of the \mathbf{x}_t , and (ii) more

robust estimates under the presence of outliers and non-linearities, when compared with the ordinary least squares estimator.

Notice that the linear association assumption, $Q_y(\tau|\mathbf{x}_t) = \mathbf{x}_t' \boldsymbol{\beta}(\tau)$, can be generally restrictive. Instead, we consider the case of the following non-linear association,

$$Q_y(\tau|\mathbf{x}_t) = h_\tau(\mathbf{x}_t),$$

where $h_\tau(\cdot)$ is some unknown, (potentially highly) non-linear function. In this article, we use an estimation strategy to approximate $h_\tau(\mathbf{x}_t)$ with neural networks using their universal approximation property. Specifically, we assume that there exists a neural network with a function $G_\tau(\mathbf{x}_t, \boldsymbol{w})$, to be defined below, that can approximate $h_\tau(\mathbf{x}_t)$ well. Before we illustrate how this methodology is implemented, we provide a discussion on how neural networks can approximate $h_\tau(\mathbf{x}_t)$.

1.2 Neural Networks

In this article, we limit our attention to feed-forward neural networks, to approximate $h_\tau(\mathbf{x}_t)$. This architecture consists of an input layer of covariates, the hidden layer(s) where non-linear transformations of the covariates occur, and the output layer that gives the final prediction. Each hidden layer has several interconnected neurons relating it to both the previous and next ones. Specifically, information flows from one layer to the other, via neurons only in one direction, and the connections correspond to weights. Optimizing a loss function w.r.t. these weights makes neural networks capable of learning.

Throughout our exposition, L denotes the total number of hidden layers, a measure for the depth of a neural network, and $J^{(l)}$ denotes the total number of neurons at layer l , a measure of its width. We start by presenting a general definition of a deep (multi-layer) feed-forward neural network. Let $\sigma_l(\cdot)$, $l = 0, \dots, L$ be the activation function used at the l th layer, that is applied element-wise and induces non-linearity. We use the ReLU activation function, $\sigma_l(\cdot) = \max(\cdot, 0)$, for $l = 1, \dots, L - 1$, applied element-wise and a linear one for the output layer, $l = L$. We denote by $g^{(l)}$ the output of the l th layer which is a vector of length equal to the number of the $J^{(l)}$ neurons in that layer, such that $g^{(0)} = \mathbf{x}_t$. Then, the overall structure of the network is equal to:

$$G_\tau(\mathbf{x}_t, \boldsymbol{w}) = g^{(L)}\left(g^{(L-1)}\left(\dots\left(g^{(1)}(\cdot)\right)\right)\right), \quad (5)$$

where

$$g^{(l)}(\mathbf{x}_t) = \sigma_l\left(\mathbf{W}^{(l-1)}g^{(l-1)} + \mathbf{b}^{(l)}\right), \quad l = 1, \dots, L, \quad (6)$$

$\mathbf{W}^{(l)}$ is a $J^{(l)} \times J^{(l-1)}$ matrix of weights, $\mathbf{b}^{(l)}$ is a $J^{(l)} \times 1$ vector of biases giving an overall vector $\boldsymbol{w} = (\text{vec}(\mathbf{W}^{(0)})', \dots, \text{vec}(\mathbf{W}^{(L)})', \mathbf{b}^{(1)'}, \dots, \mathbf{b}^{(L)'})'$ of trainable parameters of dimensions $J^{(l)}(1 + J^{(l-1)})$ total number of parameters in each hidden layer l , $J^{(0)} = p$ and $J^{(L)} = 1$.

According to various universal approximation theorems [see e.g., the theoretical results in Hornik (1991), Hornik, Stinchcombe, and White (1989), Gallant and White (1992), Kapanios and Blake (2010), Liang and Srikant (2016), and Yarotsky (2017)], $G_\tau(\mathbf{x}_t, \boldsymbol{w})$ can approximate arbitrarily well $h_\tau(\mathbf{x}_t)$, such that, for any $\epsilon > 0$,

$$\sup_t |G_\tau(\mathbf{x}_t, \boldsymbol{w}) - h_\tau(\mathbf{x}_t)| < \epsilon. \quad (7)$$

In this sense, the above (ϵ) -approximation can be seen as a sieve-type non-parametric estimation bound, where ϵ can become arbitrarily small by increasing the complexity of $G_\tau(\mathbf{x}_t, \boldsymbol{w})$.

The increase in complexity can occur, either by letting $L \rightarrow \infty$, which stands for deep learning, or by letting $J^{(l)} \rightarrow \infty$. While asymptotically, both ways deliver the same results (see e.g., Farrell, Liang, and Misra (2021) and references therein), the approximation error has been shown to decline exponentially with L , see for example, Babii et al. (2020) but only polynomially with $J^{(l)}$, providing some evidence for the prevalent use of deep learning. Notice that there also exists an alternative approximation theory for sparse deep learning, see for example, the work of Schmidt-Hieber (2020). As an illustration, in the Supplementary Appendix, we depict a simple feed-forward neural network with two inputs, two hidden layers, a total of five neurons and one output layer.

1.3 Non-linear Quantile Regression

We assume that the conditional quantile follows a non-linear relationship $Q_y(\tau|\mathbf{x}_t) = h_\tau(\mathbf{x}_t)$ and there exists a function $G_\tau(\mathbf{x}_t, \boldsymbol{w})$, that can (ϵ) -approximate $h_\tau(\mathbf{x}_t)$, see the bound in Equation (7). Using this assumption, we can formally define the conditional quantile function as the following approximation

$$Q_y(\tau|\mathbf{x}_t) = G_\tau(\mathbf{x}_t, \boldsymbol{w}) + O(\epsilon),$$

where $G_\tau(\mathbf{x}_t, \boldsymbol{w})$ is the unknown non-linear function we want to estimate in order to approximate $h_\tau(\mathbf{x}_t)$. We obtain the deep neural network conditional quantile estimate from the solution of the following minimization problem:

$$Q_y(\tau|\mathbf{x}_t) = \underset{\boldsymbol{w}}{\operatorname{argmin}} \frac{1}{T} \sum_{t=1}^T \rho_\tau(y_t - G_\tau(\mathbf{x}_t, \boldsymbol{w})), \tag{8}$$

where $\boldsymbol{w} = (\operatorname{vec}(\mathbf{W}^{(0)})', \dots, \operatorname{vec}(\mathbf{W}^{(L)})', \mathbf{b}^{(1)'}', \dots, \mathbf{b}^{(L)'}')$ contains all model parameters, and $G_\tau(\mathbf{x}_t, \boldsymbol{w})$ denotes the overall non-linear mapping, described in Equations (5) and (6). Notice that the choice of $G_\tau(\mathbf{x}_t, \boldsymbol{w})$ will govern whether the model is parametric or non-parametric. If the number of neurons and layers is small, then the model is parametric, if the above number becomes large, then the model becomes non-parametric, since the number of estimated parameters increases with the sample size, similar to sieve non-parametric approximations.

To allow the use of mixed frequency data, we can make the following changes to the structure of the network $G_\tau(\mathbf{x}_t, \boldsymbol{w})$:

In the input layer, we implement frequency alignment on each input variable \mathbf{x}_t according to the corresponding maximum lag order K . Thus, each high-frequency predictor \mathbf{x}_t is transformed into a low-frequency vector $\mathbf{x}_t^* = B(L_\varphi; \boldsymbol{\vartheta})\mathbf{x}_t$,

$$B(L_\varphi; \boldsymbol{\vartheta}) = \sum_{k=0}^K B(k; \boldsymbol{\vartheta})L_\varphi^k, \quad B(k; \boldsymbol{\vartheta}) = \frac{\exp(\vartheta_1 k + \vartheta_2 k^2)}{\sum_{k=1}^K \exp(\vartheta_1 k + \vartheta_2 k^2)}, \tag{9}$$

where $B(k; \boldsymbol{\vartheta})$ is the normalized Almon polynomial, L_φ^k is a lag operator such that $L_\varphi^k \mathbf{x}_t^\varphi = \mathbf{x}_{t-k}^\varphi$; the lag coefficients in $B(k; \boldsymbol{\vartheta})$ of the corresponding lag operator L_φ^k are parameterized as a function of a small dimensional vector of parameters $\boldsymbol{\vartheta}$. We use this weight

function on the frequency alignment vector to reduce the number of parameters and ensure a parsimonious specification. As a consequence, the low-frequency variable \mathbf{x}_t^* which has the same frequency as the output y_t is obtained. The rest of the architecture of the *deep MIDAS* follows the architecture of the deep quantile estimator, but instead of using \mathbf{x}_t in Equation (6), we use \mathbf{x}_t^* .

1.4 Regularized Non-linear Quantile Regression

Neural networks have a great capacity to estimate non-linear relationships from the data, but this comes at a cost, since they are prone to overfitting. This can lead to a severe drop in their forecasting performance, especially in small samples. There is a variety of commonly used techniques in ML, see for example, Gu, Kelly, and Xiu (2021) for a good summary, that can be used to ease this impact, originally coming from the high-dimensional statistical literature. The reader is also referred to Goodfellow, Bengio, and Courville (2016) for an excellent summary of different topics about the implementation of neural networks, including regularization.

1.4.1 Regularization

A common solution to this caveat is regularization, where a penalty term is imposed on the weights of the neural network and is appended in the loss function. Regularization, generally improves the out-of-sample performance of the network by decreasing the in-sample noise from over-parameterization, utilizing the bias-variance trade-off. Further, another benefit of regularization is that it provides computational gains in the optimization algorithm. The penalized loss function, for a given quantile τ , can be written as:

$$L(G_\tau(\mathbf{x}_t, \mathbf{w}), y_t) = \frac{1}{T} \sum_{t=1}^T \rho_\tau(y_t - \widehat{G}_\tau(\mathbf{x}_t, \mathbf{w})) + \phi(\mathbf{w}), \quad (10)$$

where the penalty term is

$$\phi(\mathbf{w}) = \begin{cases} \lambda \|\mathbf{w}\|_1, & \text{LASSO} \\ \lambda \|\mathbf{w}\|_2^2, & \text{Ridge} \\ \lambda(1 - \alpha) \|\mathbf{w}\|_1 + \lambda\alpha \|\mathbf{w}\|_2^2, & \text{Elastic Net} \\ 0, & \text{otherwise} \end{cases},$$

and λ and α are tuning parameters, for which we discuss their selection below. Generally, there is a plethora of loss functions, and the choice among them, depends mainly on the task at hand. In this article, we use the quantile loss function. The different penalization schemes on $\phi(\mathbf{w})$ work as follows: *deep LASSO* or l_1 -norm penalization, is a regularization method that shrinks uniformly all the weights to zero, and some at exactly zero. The latter is referred to as the variable selection property of the deep LASSO. *Deep Ridge* works in a similar manner to the deep LASSO, by shrinking the weights, uniformly to zero, but not at exactly zero. Finally, the *deep Elnet*² is a combination of *deep LASSO* and *deep Ridge*, which has been shown to retain good features from both methods, see for example, Zou and Hastie (2005).

Throughout the MC experiments and empirical application, we set $\kappa = 20$ and $\delta = 5$ as in Babii, Ghysels, and Striaukas (2022a). Finally, we use the optimal parameters and hyper-parameters from the time-series CV and evaluate the out-of-sample performance of the network.

1.4.3 Optimization

The estimation of neural networks is generally a computational cumbersome optimization problem due to non-linearities and non-convexities. The most commonly used solution utilizes stochastic gradient descent (SGD) to train a neural network. SGD uses a batch of a specific size, that is, a small subset of the data at each epoch (iteration) of the optimization to evaluate the gradient, to alleviate the computation hurdle. The step of the derivative at each epoch is controlled by the learning rate, γ . We use the adaptive moment estimation algorithm (ADAM) proposed by Kingma and Ba (2014),⁴ which is a more efficient version of SGD. Finally, we set the number of epochs to 5000 and use early stopping, following Gu, Kelly, and Xiu (2020) to avoid any potential overfitting.

2 Monte Carlo

2.1 Setup

In this section, we present MC experiments, in order to study the finite sample performance of the deep quantile estimator as outlined in Section 1, for the different penalization schemes. We generate artificial data $\{y_t\}$ using a single predictor $\{x_t\}$, according to the following model:

$$y_t = h_\tau(x_t) + u_t, \quad (13)$$

where u_t is the realization of a random variable u distributed as, $u_t \sim \text{i.i.d.}N(-\sigma\Phi^{-1}(\tau), \sigma^2)$, $\sigma = 0.1$ and Φ^{-1} is the quantile function of the standard normal distribution. $h_\tau(\cdot)$ is the general non-linear function that we wish to approximate via the deep quantile estimator.

All the experiments are based on the following values: $\tau \in (1\%, 2.5\%, 5\%, 10\%, 20\%)$, $T \in (100, 300, 500, 1000)$ and the number of MC replications is 1000. We consider the following five *data-generating mechanisms* to assess the finite sample properties of the deep quantile estimator:

Case I: We consider the case of a $N(0, 1)$ simulated single predictor that is generated as

$$y_t = h_\tau(x_t) + u_t, \quad h_\tau(x_t) = \sin(2\pi x_t), \quad x_t \sim_{\text{i.i.d.}} N(0, 1).$$

This is the simplest design in our MC experiments. We use this simple case to showcase that linear methods, as expected, cannot produce reasonable performance under a sigmoid type of a non-linear function $h_\tau(\cdot)$.

4 ADAM is using estimates for the first and second moments of the gradient to calculate the learning rate.

Case II: We consider an AR(1) simulated single predictor as follows:

$$y_t = h_\tau(x_t) + u_t, \quad h_\tau(x_t) = \sin(2\pi x_t),$$

where x_t is simulated as

$$x_t = 0.8x_{t-1} + \varepsilon_t, \quad \varepsilon_t \sim_{\text{i.i.d.}} N(0, 1).$$

In this design, we increase the complexity by introducing a correlated predictor.

Case III: We consider the case of a single predictor generated via a GARCH(1,1) model

$$y_t = h_\tau(x_t) + u_t, \quad h_\tau(x_t) = \sin(2\pi x_t),$$

where x_t is simulated as:

$$x_t = \sigma_t \varepsilon_t, \quad \sigma_t^2 = 1 + 0.7x_{t-1}^2 + 0.2\sigma_{t-1}^2, \quad \varepsilon_t \sim_{\text{i.i.d.}} N(0, 1).$$

In this design, we wish to examine, how the deep quantile estimator fares, when the regressor is conditionally heteroskedastic, following a GARCH(1, 1) model. A GARCH type of assumption on the distribution of asset returns is one commonly used in the literature.

Case IV: We consider the case of a single predictor that is generated as follows:

$$y_t = h_\tau(x_t) + u_t, \quad h_\tau(x_t) = G_\tau(x_t, \mathbf{w}), \quad x_t \sim_{\text{i.i.d.}} N(0, 1).$$

In this case, we simulate $h_\tau(x_t)$ to reflect a function composition, commonly used in neural networks. We simulate it with three hidden layers and a specific number of neurons, such as

$$G_\tau(x_t, \mathbf{w}) = \left(\mathbf{W}^{(3)} \left(\sin \left(\mathbf{W}^{(2)} \left(\sin \left(\mathbf{W}^{(1)} \left(\sin \left(\mathbf{W}^{(0)} x_t' + \mathbf{b}^{(1)} \right) \right) + \mathbf{b}^{(2)} \right) \right) + \mathbf{b}^{(3)} \right) \right) \right)',$$

where $\mathbf{w} = (\text{vec}(\mathbf{W}^{(0)})', \dots, \text{vec}(\mathbf{W}^{(3)})', \mathbf{b}^{(1)'}, \dots, \mathbf{b}^{(3)'})'$, $\mathbf{W}^{(0)}$ is 50×1 , $\mathbf{W}^{(1)}$ is 10×50 , $\mathbf{W}^{(2)}$ is 8×10 and $\mathbf{W}^{(3)}$ is 1×8 . Further, we simulate the weights, \mathbf{w} , so that, every entry $w_{i,j}$ is simulated as, $w_{i,j} = \delta_{i,j} 1(\delta_{i,j} > 0.5)$, where $\delta_{i,j} \sim U(0, 1)$, allowing for some sparsity.

Case V: We consider an AR(1) simulated error as follows:

$$y_t = h_\tau(x_t) + \varepsilon_t, \quad h_\tau(x_t) = \sin(2\pi x_t), \quad x_t \sim_{\text{i.i.d.}} N(0, 1),$$

where ε_t is simulated as

$$\varepsilon_t = 0.6\varepsilon_{t-1} + u_t.$$

We use this design to examine whether correlated errors impact the deep quantile estimator.

Across all cases, we estimate $h_\tau(x_t)$ using the deep quantile estimator with different penalization schemes. Let $\hat{h}_{\tau, \text{pen}} = \hat{G}_{\tau, \text{pen}}(x_t, \mathbf{w})$ denotes the estimate, where pen corresponds to no regularization, deep LASSO, deep Ridge, and deep Elnet. We use the following metrics in order to evaluate the small sample properties, of the deep quantile estimator across $R=1000$, MC replications: (i) the average mean-squared error of the true residuals,

$AMSE_{u_i} = \frac{1}{R} \frac{1}{T} \sum_{i=1}^R \left(\sum_{t=1}^T u_t^2 \right)_i$, (ii) the average mean-squared error of the estimated residuals, $AMSE_{\hat{u}_{t,pen}} = \frac{1}{R} \frac{1}{T} \sum_{i=1}^R \left(\sum_{t=1}^T (y_t - \hat{y}_{t,pen})^2 \right)_i$ and finally, (iii) the average absolute bias $ABIAS_{\hat{b}_{\tau,pen}} = \frac{1}{R} \frac{1}{T} \sum_{i=1}^R \left(\sum_{t=1}^T |(b_{\tau}(x_t) - \hat{G}_{\tau,pen}(x_t, \mathbf{w}))| \right)_i$. We report results only for $AMSE_{\hat{u}_{t,pen}}$ below, since results for the alternative metrics exhibit similar patterns and are available upon request.

2.2 Results

We consider a linear quantile estimator, other sieve-type have estimators such as polynomials and B-splines (defined in the [Supplementary Appendix](#)), shallow and deep neural network estimators with their depth being described in Section 1. Note, that we also use the time-series CV to select the optimal number of knots of the B-splines estimator from the following grid of six discrete values in the interval [5, 10]. Our MC aims to answer which estimator has the best asymptotic properties under a non-linear setup. We find that both shallow and deep networks deliver good finite sample properties across quantiles, with shallow learning being the best between the two in most cases. Further, neural networks perform better than the linear quantile and other sieve estimators. We present our MC results for Cases I–V in [Tables 1–5](#), respectively.

In [Table 1](#), we can see that the linear quantile estimator, under a non-linear setup does not work, as expected, and the MSE remains constant as the sample size increases. Next, we present the asymptotic properties for the deep quantile estimator across different penalization schemes, namely deep quantile, deep LASSO, deep Ridge, and *deep Elastic Net*, and find that the deep quantile non-linear estimators have good finite sample properties.

When $\tau = 1\%$ it appears that the deep quantile estimator works well for sample sizes larger than $T = 300$, but in comparison with the linear one it generally works better. In Case II, the non-linear estimators depict fine finite sample properties and their performance is better than the linear one. In this case, the non-regularized estimator performs better than the regularized ones. Next, similar behavior appears in Case III. In Case IV, where we allow for some sparsity in the weights, we find, as expected, that the linear quantile regression estimator, does not work under non-linearity, while the non-linear one works as expected. In Case V, where we consider serial correlated errors, we find that adding a penalty term in the non-linear estimators improves the performance of the deep quantile estimator in extreme quantiles.

We further find, as expected, that the linear quantile regression, second-order quantile polynomial, and cubic splines estimators, do not work under non-linearity. Finally, in very few occasions, we find that splines estimator performs better than shallow networks for small sample sizes and extreme quantiles.

Overall, our MC results suggest that the deep quantile estimator using both deep and shallow learning has good finite sample properties, and can approximate non-linear functions. We also find evidence in favor of the penalization schemes described in Section 1. Specifically, the penalized deep quantile estimators also have good finite sample properties, and in some cases, perform better than the non-regularized one; a finding in favor of weight regularization.

Table 1. MC results for Case I. Model: $y_t = h_\tau(x_t) + u_t$, $h_\tau(x_t) = \sin(2\pi x_t)$, $x_t \sim \text{i.i.d.}N(0, 1)$, $u_t \sim \text{i.i.d.}N(-\sigma\Phi^{-1}(\tau), \sigma^2)$, $\sigma = 0.1$ and Φ^{-1} is the quantile function of the standard normal distribution

τ	T	deep Quantile	deep LASSO	deep Ridge	deep Elnet	shallow Quantile	shallow LASSO	shallow Ridge	shallow Elnet	Linear	Polynomial	Splines
1%	100	1.610	1.837	1.860	1.685	1.563	1.839	1.790	1.715	1.754	1.710	1.543
	300	1.667	1.735	1.737	1.694	1.411	1.703	1.500	1.332	1.749	1.743	1.727
	500	1.687	1.651	1.677	1.676	1.285	1.430	1.172	1.094	1.757	1.757	1.690
	1000	1.723	1.381	1.470	1.186	1.150	0.745	1.023	0.849	1.768	1.767	1.710
2.5%	100	1.432	1.561	1.548	1.493	1.380	1.441	1.397	1.407	1.645	1.620	1.547
	300	1.315	1.253	1.375	1.281	1.057	1.218	0.899	0.848	1.637	1.635	1.619
	500	1.335	0.875	0.968	0.989	0.755	0.872	0.712	0.680	1.640	1.638	1.583
	1000	1.301	0.547	0.823	0.674	0.558	0.459	0.580	0.450	1.651	1.650	1.607
5%	100	1.257	1.308	1.312	1.280	1.184	1.243	1.127	1.096	1.532	1.524	1.477
	300	0.895	0.699	0.873	0.802	0.860	0.799	0.580	0.624	1.526	1.523	1.506
	500	0.928	0.469	0.642	0.647	0.562	0.494	0.501	0.438	1.530	1.528	1.473
	1000	0.919	0.279	0.616	0.399	0.300	0.215	0.380	0.259	1.542	1.542	1.500
10%	100	0.958	0.939	0.971	0.987	1.002	0.955	0.848	0.831	1.387	1.378	1.321
	300	0.498	0.369	0.466	0.473	0.584	0.503	0.423	0.471	1.378	1.376	1.356
	500	0.702	0.277	0.387	0.386	0.311	0.281	0.373	0.298	1.387	1.386	1.331
	1000	0.542	0.158	0.360	0.181	0.143	0.108	0.141	0.178	1.394	1.394	1.350
20%	100	0.694	0.616	0.666	0.665	0.665	0.631	0.615	0.617	1.124	1.118	1.002
	300	0.325	0.225	0.307	0.294	0.336	0.279	0.290	0.251	1.126	1.121	1.076
	500	0.404	0.148	0.265	0.199	0.178	0.150	0.178	0.177	1.131	1.129	1.068
	1000	0.179	0.099	0.225	0.136	0.068	0.058	0.060	0.097	1.142	1.141	1.088

Notes: The figure presents the average mean-squared error of the estimated residuals, $AMSE_{\hat{u}_t, \text{pen}}$ for the different penalization schemes (Model), $T = 100, 300, 500, 1000$ and different quantiles, $\tau = (1\%, 2.5\%, 5\%, 10\%, 20\%)$. Entries in bold indicate the model with minimum AMSE value for each quantile and different sample size.

Table 2. MC results for Case II. Model: $y_t = h_\tau(x_t) + u_t$, $h_\tau(x_t) = \sin(2\pi x_t)$, $x_t = 0.8x_{t-1} + \varepsilon_t$, $\varepsilon_t \sim \text{i.i.d.}N(0, 1)$, $u_t \sim \text{i.i.d.}N(-\sigma\Phi^{-1}(\tau), \sigma^2)$, $\sigma = 0.1$ and Φ^{-1} is the quantile function of the standard normal distribution

τ	T	deep Quantile	deep LASSO	deep Ridge	deep Elnet	shallow Quantile	shallow LASSO	shallow Ridge	shallow Elnet	Linear	Polynomial	Splines
1%	100	1.647	1.873	1.848	1.686	1.645	2.200	1.942	1.816	1.774	1.716	1.553
	300	1.761	1.798	1.803	1.784	1.558	1.791	1.765	1.689	1.781	1.773	1.744
	500	1.755	1.763	1.758	1.768	1.546	1.688	1.551	1.511	1.756	1.753	1.733
	1000	1.779	1.735	1.747	1.670	1.435	1.235	1.395	1.291	1.768	1.767	1.759
2.5%	100	1.512	1.624	1.592	1.554	1.541	1.647	1.634	1.567	1.651	1.640	1.556
	300	1.564	1.569	1.611	1.565	1.357	1.522	1.326	1.266	1.662	1.659	1.638
	500	1.551	1.411	1.441	1.433	1.190	1.356	1.147	1.072	1.639	1.638	1.627
	1000	1.555	1.133	1.297	1.114	0.902	0.966	0.974	0.904	1.653	1.652	1.645
5%	100	1.388	1.441	1.442	1.408	1.349	1.462	1.359	1.378	1.543	1.526	1.458
	300	1.329	1.212	1.289	1.244	1.136	1.260	0.985	0.958	1.552	1.549	1.530
	500	1.300	0.943	1.045	1.033	0.907	1.080	0.918	0.845	1.532	1.531	1.520
	1000	1.198	0.650	0.863	0.790	0.569	0.690	0.726	0.594	1.544	1.543	1.537
10%	100	1.171	1.192	1.217	1.183	1.152	1.190	1.068	1.092	1.384	1.376	1.304
	300	0.940	0.769	0.851	0.842	0.897	0.918	0.677	0.687	1.404	1.400	1.377
	500	0.935	0.564	0.712	0.660	0.657	0.685	0.656	0.591	1.383	1.381	1.367
	1000	0.813	0.367	0.616	0.449	0.357	0.342	0.455	0.382	1.395	1.394	1.387
20%	100	0.904	0.850	0.877	0.879	0.878	0.865	0.822	0.810	1.135	1.128	1.035
	300	0.624	0.501	0.546	0.567	0.556	0.538	0.466	0.434	1.142	1.139	1.108
	500	0.593	0.337	0.446	0.403	0.387	0.372	0.391	0.368	1.130	1.128	1.105
	1000	0.440	0.211	0.414	0.228	0.180	0.124	0.224	0.224	1.140	1.139	1.128

Notes: The figure presents the average mean-squared error of the estimated residuals, $AMSE_{\hat{u}_{t,pen}}$ for the different penalization schemes (Model), $T = 100, 300, 500, 1000$ and different quantiles, $\tau = (1\%, 2.5\%, 5\%, 10\%, 20\%)$. Entries in bold indicate the model with minimum AMSE value for each quantile and different sample size.

Table 3. MC results for Case III. Model: $y_t = h_\tau(x_t) + u_t$, $h_\tau(x_t) = \sin(2\pi x_t)$, $x_t = \sigma_t \varepsilon_t$, $\sigma_t^2 = 1 + 0.7x_{t-1}^2 + 0.2\sigma_{t-1}^2$, $u_t \sim \text{i.i.d.}N(-\sigma\Phi^{-1}(\tau), \sigma^2)$, $\sigma = 0.1$ and Φ^{-1} is the quantile function of the standard normal distribution

τ	T	deep Quantile	deep LASSO	deep Ridge	deep Elnet	shallow Quantile	shallow LASSO	shallow Ridge	shallow Elnet	Linear	Polynomial	Splines
1%	100	1.617	1.872	1.810	1.658	1.804	2.545	2.400	1.990	1.732	1.682	1.518
	300	1.751	1.780	1.779	1.772	1.645	1.842	1.860	1.783	1.750	1.742	1.716
	500	1.778	1.797	1.795	1.795	1.680	1.793	1.794	1.761	1.764	1.762	1.741
	1000	1.779	1.792	1.786	1.778	1.660	1.639	1.705	1.686	1.762	1.760	1.750
2.5%	100	1.516	1.605	1.581	1.542	1.594	1.735	1.675	1.621	1.613	1.600	1.516
	300	1.610	1.619	1.624	1.619	1.541	1.607	1.569	1.525	1.635	1.632	1.607
	500	1.636	1.611	1.625	1.627	1.520	1.578	1.502	1.457	1.649	1.648	1.634
	1000	1.639	1.532	1.593	1.540	1.444	1.460	1.429	1.378	1.644	1.644	1.637
5%	100	1.411	1.461	1.456	1.432	1.435	1.520	1.450	1.455	1.510	1.505	1.426
	300	1.465	1.439	1.454	1.447	1.393	1.443	1.355	1.299	1.522	1.521	1.499
	500	1.482	1.371	1.410	1.391	1.346	1.407	1.313	1.236	1.539	1.536	1.523
	1000	1.478	1.183	1.274	1.216	1.176	1.236	1.170	1.123	1.535	1.534	1.528
10%	100	1.264	1.273	1.278	1.272	1.271	1.294	1.276	1.255	1.365	1.357	1.280
	300	1.241	1.154	1.199	1.206	1.189	1.197	1.089	1.081	1.375	1.374	1.350
	500	1.234	1.041	1.108	1.093	1.091	1.168	1.060	0.993	1.391	1.389	1.376
	1000	1.179	0.825	0.972	0.913	0.804	0.889	0.855	0.840	1.390	1.389	1.382
20%	100	0.985	0.963	0.979	0.980	0.976	0.990	0.965	0.964	1.100	1.094	1.008
	300	0.934	0.812	0.865	0.883	0.877	0.880	0.801	0.778	1.120	1.118	1.087
	500	0.869	0.690	0.762	0.733	0.703	0.734	0.681	0.654	1.134	1.132	1.111
	1000	0.770	0.524	0.680	0.554	0.478	0.466	0.574	0.543	1.136	1.135	1.125

Notes: The figure presents the average mean-squared error of the estimated residuals, $AMSE_{\hat{u}_t, \text{pen}}$ for the different penalization schemes (Model), $T = 100, 300, 500, 1000$ and different quantiles, $\tau = (1\%, 2.5\%, 5\%, 10\%, 20\%)$. Entries in bold indicate the model with minimum AMSE value for each quantile and different sample size.

Table 4. MC results for Case IV. Model: $y_t = h_\tau(x_t) + u_t$, $h_\tau(x_t) = G(x_t, w)$, $x_t \sim \text{i.i.d.}N(0, 1)$, $w_{i,j} = \delta_{ij}1(\delta_{ij} > 0.1)$, $\delta_{ij} \sim U(0, 1)$, $u_t \sim \text{i.i.d.}N(-\sigma\Phi^{-1}(\tau), \sigma^2)$, $\sigma = 0.1$ and Φ^{-1} is the quantile function of the standard normal distribution

τ	T	deep Quantile	deep LASSO	deep Ridge	deep Elnet	shallow Quantile	shallow LASSO	shallow Ridge	shallow Elnet	Linear	Polynomial	Splines
1%	100	0.567	0.733	0.740	0.627	0.580	0.595	0.626	0.629	0.749	0.775	0.263
	300	0.544	0.558	0.603	0.584	0.399	0.323	0.410	0.385	0.779	0.780	0.410
	500	0.492	0.396	0.405	0.406	0.255	0.217	0.255	0.252	0.710	0.705	0.330
	1000	0.542	0.315	0.339	0.281	0.257	0.169	0.266	0.202	0.794	0.777	0.496
2.5%	100	0.402	0.456	0.429	0.428	0.388	0.392	0.379	0.413	0.629	0.629	0.270
	300	0.301	0.280	0.325	0.299	0.256	0.178	0.199	0.214	0.688	0.717	0.303
	500	0.306	0.192	0.224	0.201	0.178	0.146	0.162	0.156	0.680	0.681	0.350
	1000	0.322	0.141	0.208	0.157	0.127	0.101	0.154	0.120	0.715	0.724	0.371
5%	100	0.304	0.313	0.317	0.348	0.306	0.283	0.290	0.288	0.618	0.663	0.213
	300	0.157	0.124	0.170	0.161	0.165	0.123	0.128	0.146	0.619	0.661	0.282
	500	0.189	0.094	0.129	0.124	0.111	0.076	0.122	0.100	0.644	0.674	0.286
	1000	0.182	0.077	0.146	0.095	0.079	0.064	0.075	0.077	0.636	0.657	0.340
10%	100	0.190	0.169	0.183	0.200	0.193	0.191	0.197	0.200	0.514	0.546	0.148
	300	0.101	0.092	0.101	0.120	0.098	0.084	0.092	0.093	0.558	0.593	0.267
	500	0.136	0.065	0.093	0.082	0.055	0.055	0.060	0.065	0.536	0.566	0.274
	1000	0.093	0.049	0.077	0.055	0.039	0.041	0.042	0.044	0.577	0.612	0.339
20%	100	0.133	0.112	0.128	0.127	0.134	0.124	0.119	0.113	0.425	0.465	0.126
	300	0.081	0.053	0.073	0.070	0.063	0.048	0.051	0.050	0.420	0.443	0.210
	500	0.060	0.038	0.059	0.050	0.036	0.032	0.034	0.035	0.419	0.450	0.267
	1000	0.050	0.034	0.038	0.033	0.023	0.026	0.023	0.028	0.436	0.469	0.239

Notes: The figure presents the average mean-squared error of the estimated residuals, $AMSE_{\hat{u}_t, \text{pen}}$ for the different penalization schemes (Model), $T = 100, 300, 500, 1000$ and different quantiles, $\tau = (1\%, 2.5\%, 5\%, 10\%, 20\%)$. Entries in bold indicate the model with minimum AMSE value for each quantile and different sample size.

Table 5. MC results for Case V. Model: $y_t = h_\tau(x_t) + \varepsilon_t$, $h_\tau(x_t) = \sin(2\pi x_t)$, $x_t \sim \text{i.i.d. } N(0, 1)$, $\varepsilon_t = 0.6\varepsilon_{t-1} + u_t$, $u_t \sim \text{i.i.d. } N(-\sigma\Phi^{-1}(\tau), \sigma^2)$, $\sigma = 0.1$ and Φ^{-1} is the quantile function of the standard normal distribution

τ	T	deep Quantile	deep LASSO	deep Ridge	deep Elnet	shallow Quantile	shallow LASSO	shallow Ridge	shallow Elnet	Linear	Polynomial	Splines
1%	100	1.687	1.811	1.769	1.690	1.615	1.707	1.840	1.752	1.827	1.780	1.571
	300	1.805	1.780	1.805	1.797	1.519	1.321	1.609	1.481	1.854	1.849	1.834
	500	1.798	1.688	1.758	1.714	1.444	1.029	1.321	1.213	1.852	1.849	1.775
	1000	1.852	1.524	1.572	1.456	1.254	0.789	0.999	0.904	1.852	1.849	1.781
2.5%	100	1.527	1.569	1.580	1.535	1.433	1.389	1.512	1.464	1.682	1.660	1.565
	300	1.556	1.312	1.441	1.406	1.139	0.800	1.058	1.025	1.699	1.699	1.686
	500	1.551	1.048	1.246	1.166	0.972	0.572	0.867	0.746	1.716	1.715	1.647
	1000	1.517	0.768	1.013	0.826	0.764	0.348	0.567	0.505	1.708	1.707	1.636
5%	100	1.367	1.371	1.385	1.349	1.261	1.102	1.224	1.180	1.565	1.562	1.454
	300	1.187	0.818	0.952	0.975	0.885	0.559	0.674	0.669	1.574	1.569	1.526
	500	1.156	0.604	0.832	0.738	0.683	0.363	0.638	0.503	1.577	1.576	1.523
	1000	1.091	0.405	0.741	0.443	0.414	0.205	0.349	0.289	1.572	1.571	1.509
10%	100	1.055	0.998	1.082	1.018	0.979	0.805	0.907	0.891	1.384	1.377	1.311
	300	0.733	0.444	0.582	0.626	0.579	0.358	0.479	0.425	1.402	1.400	1.359
	500	0.752	0.330	0.547	0.449	0.404	0.252	0.402	0.343	1.408	1.407	1.350
	1000	0.562	0.203	0.460	0.219	0.172	0.125	0.168	0.187	1.405	1.405	1.347
20%	100	0.753	0.678	0.715	0.727	0.716	0.612	0.674	0.677	1.112	1.111	0.997
	300	0.429	0.262	0.399	0.325	0.351	0.234	0.322	0.257	1.132	1.130	1.089
	500	0.381	0.201	0.348	0.206	0.229	0.155	0.191	0.209	1.135	1.133	1.082
	1000	0.230	0.120	0.194	0.146	0.062	0.071	0.071	0.112	1.133	1.132	1.074

Notes: The figure presents the average mean-squared error of the estimated residuals, $AMSE_{\hat{u}_t, \text{pen}}$ for the different penalization schemes (Model), $T = 100, 300, 500, 1000$ and different quantiles, $\tau = (1\%, 2.5\%, 5\%, 10\%, 20\%)$. Entries in bold indicate the model with minimum AMSE value for each quantile and different sample size.

3 Empirical Setup

In this section, we outline our empirical application setup, where we use the deep quantile estimator to forecast VaR. We examine the predictive ability of the deep quantile estimator and other non-parametric models, relative to the linear one, using the quantile-encompassing test of [Giacomini and Komunjer \(2005\)](#). We further examine the predictive performance of the different methods by testing their forecasting accuracy, using the DM, GW, and quantile score tests.

3.1 Deep Quantile VaR Forecasting

The data used in our empirical application consist of around 36 years of daily prices on the S&P500 index (source: Bloomberg), from September 1985 to August 2020 ($T=9053$ observations). We use daily log returns, defined as $r_t = \log(P_t/P_{t-1})$ for our forecasting analysis. We use four different classes of VaR models and produce forecasts for $\tau = (1\%, 5\%, 10\%)$ empirical conditional quantiles, using the deep quantile estimator.

The first VaR specification we consider is the GARCH(1,1) model that has been proposed by [Bollerslev \(1986\)](#), in which $\sigma_{1,t}^2 = \omega_0 + \omega_1\sigma_{1,t-1}^2 + \omega_2r_{t-1}^2$, see [Equation \(14\)](#). The second VaR specification we consider, is RiskMetrics, proposed by [Morgan \(1996\)](#), which assumes $\sigma_{2,t}^2 = \lambda\sigma_{2,t-1}^2 + (1-\lambda)r_{t-1}^2$, where for daily returns, $\lambda = 0.94$, see [Equation \(15\)](#).

The last two specifications we consider follow the *Conditional Autoregressive VaR* (CAViAR) model, proposed by [Engle and Manganelli \(2004\)](#), where a specific quantile is analyzed, rather than the whole distribution. Specifically, the CAViAR model corrects the past $\text{VaR}_{j,t-1}$ estimates in the following way: it increases $\text{VaR}_{j,t}$ when $\text{VaR}_{j,t-1}$ is above the τ th quantile, while, when the $\text{VaR}_{j,t-1}$ is less than the τ th quantile, it reduces $\text{VaR}_{j,t}$. Thus, the third VaR we examine is the symmetric absolute value (SV) that responds symmetrically to past returns, see [Equation \(16\)](#) and lastly, we consider the asymmetric slope value (ASV) as it offers a different response to positive and negative returns, see [Equation \(17\)](#). For ease of exposition, we refer to the above specification as $\text{VaR}_{1,t}, \dots, \text{VaR}_{4,t}$, respectively. We summarize their specifications below:

$$\text{VaR}_{1,t} = \beta_0 + \beta_1\sigma_{1,t} \quad (14)$$

$$\text{VaR}_{2,t} = \beta_0 + \beta_1\sigma_{2,t} \quad (15)$$

$$\text{VaR}_{3,t} = \beta_0 + \beta_1\text{VaR}_{3,t-1} + \beta_2|r_{t-1}| \quad (16)$$

$$\text{VaR}_{4,t} = \beta_0 + \beta_1\text{VaR}_{4,t-1} + \beta_2r_{t-1}^+ - \beta_3r_{t-1}^-, \quad (17)$$

where $\beta_i, i = 0, \dots, 3$ are parameters to be estimated. We use these specifications following [Giacomini and Komunjer \(2005\)](#). Under the mixed frequency setup, we consider the following equation

$$\text{VaR}_{i,t}^{(\text{MIDAS})} = B(L_\varphi; \boldsymbol{\vartheta})\text{VaR}_{i,t}, \quad (18)$$

where $B(L_\varphi; \boldsymbol{\vartheta})$ is defined in [Equation \(9\)](#), $i = 1, \dots, 4$ and $\boldsymbol{\vartheta}$ are parameters to be estimated. For a more detailed summary of MIDAS we refer the reader to [Ghysels, Santa-Clara, and Valkanov \(2004\)](#). As discussed in Section 1, the linear association between VaR and the covariates can be restrictive. Instead, we assume that the relationship between the response variable, VaR, and the covariates has an unknown non-linear form for a given τ , that we wish to approximate with the deep quantile estimator as

$$\text{VaR}_{1,t} = G_\tau(\sigma_{1,t}, \mathbf{w}) \quad (19)$$

$$\text{VaR}_{2,t} = G_\tau(\sigma_{2,t}, \mathbf{w}) \quad (20)$$

$$\text{VaR}_{3,t} = G_\tau(\text{VaR}_{3,t-1}, |r_{t-1}|, \mathbf{w}) \quad (21)$$

$$\text{VaR}_{4,t} = G_\tau(\text{VaR}_{4,t-1}, r_{t-1}^+, r_{t-1}^-, \mathbf{w}), \quad (22)$$

where $\text{VaR}_{j,t}$, $j = 1, \dots, 4$ is indexed at (day) $t = 1, \dots, T$. The dimension p of covariates that we use in our analysis depends on the specification chosen for VaR. Specifically, if $j = 1, 2$ then $p = 1$, if $j = 3$, $p = 2$ and finally if $j = 4$ then $p = 3$.

In the [Supplementary Appendix](#), we briefly delineate the model specifications for the quantile B-splines, quantile polynomial, and quantile MIDAS estimators.

3.2 Forecasting Exercise Design

This section presents our forecasting exercise design. We reserve the last 2000 observations to evaluate the out-of-sample performance using various tests and use the remaining 7053 observations to tune parameters via time-series CV as described in Section 1. This specific split is used because we follow [Giacomini and Komunjer \(2005\)](#) and want the power of the conditional quantile forecast encompassing (CQFE) test to be comparable with their exercise. Generally, a forecasting exercise is performed either via a recursive or rolling window, see for example, [Ghysels et al. \(2019\)](#), yet in either setting to produce all h -step ahead forecasts for the last 2000 observations and to tune the hyper-parameters can be computationally challenging. Instead, we follow [Giacomini and Komunjer \(2005\)](#) and perform a fixed forecast window exercise, in which we estimate our models once.

For our forecasting design, we use a fixed forecast window exercise and predict the ten-day-ahead VaR as:

$$\widehat{\text{VaR}}_{1,t+10} | \mathcal{F}_t = G_\tau(\sigma_{1,t}, \mathbf{w}^*), \quad (23)$$

where \mathcal{F}_t denotes the information set up to time t , \mathbf{w}^* denotes the optimal weights obtained from the time-series CV. [Equation \(23\)](#) illustrates how forecasts for the first VaR specification were obtained via the deep quantile estimator. In a similar manner, forecasts can be obtained for other VaR specifications and alternative models, using [Equations \(14\)–\(22\)](#).

We evaluate the forecasting performance of VaR models with the deep quantile estimator as in Section 1. Further, we consider ten-day compounded VaR forecasts, which we relegate to the [Supplementary Appendix](#).

3.3 Forecast Evaluation

In this section, we discuss the various tests we have considered, in order to evaluate the predictive ability of the deep quantile estimator and present the testing results. In general, *root mean-squared forecast error* (RMSFE) is used to measure the accuracy of point estimates and is defined as

$$\text{RMSFE}_j = \sqrt{\frac{\sum_{t=1}^T (y_{j,t+b} - \widehat{G}_\tau(\mathbf{x}_{j,t+b}, \mathbf{w}))^2}{T}},$$

where $j = 1, \dots, 4$ depends on the specification chosen for VaR_j , b denotes the forecasting horizon and $\hat{G}_\tau(x_{j,t+b}, \mathbf{w})$ is the solution of Equation (8) after selecting the optimal \mathbf{w} via CV at the τ th quantile.

Table 6 reports relative $\text{RMSFE}_j = \frac{\text{RMSFE}_j^M}{\text{RMSFE}_j^B}$, where RMSFE^M comes from the competing estimator and RMSFE^B from the linear quantile estimator for each VaR specification. In most cases, the linear quantile estimator outperforms polynomial and splines estimators across competing models and quantiles. In all cases, deep quantile and deep quantile MIDAS estimators outperform the linear quantile estimator. The forecast gains of the deep quantile vary from 50% to 98%, while the gains from deep quantile MIDAS fluctuate between 11% and 84%. We find that neural network models improve VaR forecasts in all VaR models across the different quantiles we consider.

3.3.1 DM test

We perform a quantitative forecast comparison across different methods and test their statistical significance. To do so, we calculate the RMSFE for each method and perform the DM test, with the Harvey, Leybourne, and Newbold (1997) adjustment to gauge the statistical significance of the forecasts. As our empirical application entails quantiles, we compute the DM statistics based on the comparison of empirical quantile losses rather than the MSE loss. With the DM test, we assess the forecasting accuracy of the deep quantile estimator relative to the benchmark linear quantile regression model. In this exercise, we set τ equal to 1%, 5%, and 10%.

Results from the DM test are reported in Table 6, where asterisks denote the statistical significance of rejecting the null hypothesis of the test at 1%, 5%, and 10% level of significance, for all quantiles and models we consider. These results suggest that forecasts produced from the non-linear estimator outperform, for the majority of cases, forecasts obtained from the linear and non-parametric quantile regression estimators.

3.3.2 GW test

In a similar manner and to complement the DM test, we follow Carriero, Kapetanios, and Marcellino (2009) and further calculate the GW test of equal forecasting accuracy, that can handle forecasts based on both nested and non-nested models, regardless of the estimation procedures used for the derivation of the forecasts, including the deep quantile estimator. As in the DM test, we compute the GW statistics based on the empirical quantile losses rather than the MSE one. Table 2 in the Supplementary Appendix illustrates the results for GW test, where daggers denote the statistical significance of rejecting the null hypothesis of the test at 1%, 5% and 10% level of significance, for all quantiles and different models we consider. Similar to the DM forecasting accuracy test, the GW test is again significant at 1% in most cases, with the following exceptions.

Quantile polynomial regression forecasts are only significant at the 1% level of significance for ASV model at $\tau = 5\%$. In quantile splines, forecasts for the RM specification at $\tau = 10\%$ are significant at the 5% level of significance and under SV at $\tau = 10\%$ are significant at the 10% level of significance. Forecasts from the linear MIDAS, under the SV specification, at $\tau = 1\%$ are insignificant and at $\tau = 5\%$ are significant at the 5% level of significance and under the ASV specification, at $\tau = 1\%$ and $\tau = 5\%$, are significant at the 10% significance level.

Results for the SV with *Deep* MIDAS estimator are not significant at $\tau = 1\%$ and $\tau = 5\%$ for the GW test, while results for the ASV with Deep MIDAS estimator are significant at 5% for all quantiles we consider. Results for the SV with Deep LASSO MIDAS estimator at $\tau = 1\%$ are not significant, while at $\tau = 5\%$ are significant at the 5% level of significance.

Results for the SV with Deep Ridge MIDAS estimator are only significant at the 5% level of significance for the GW test at $\tau = 10\%$. Results for the ASV with Deep Ridge MIDAS estimator are only significant at the 5% level of significance across quantiles. Forecasts from deep Elnet MIDAS model under SV specification at $\tau = 10\%$ are significant at 5% level of significance. Finally, forecasts from deep Elnet under ASV specification across quantiles are significant at the 5% level of significance.

Overall, results from both the DM and GW tests suggest that the non-linear estimators outperform, for the majority of times, competing linear and non-parametric estimators in VaR forecasting.

3.3.3 CQFE

We present the implementation of the CQFE test as proposed by [Giacomini and Komunjer \(2005\)](#) and the generalized method of moments (GMM) estimation as proposed by [Hansen \(1982\)](#). Let $\hat{q}_{1,t}$ be a vector of the τ th quantile forecasts produced from model 1 and $\hat{q}_{2,t}$ be the competing forecasts produced from model 2. The basic principle of CQFE is to test whether $\hat{q}_{1,t}$ conditionally encompasses $\hat{q}_{2,t}$. Encompassing occurs when the second set of forecasts fails to add new information to the first set of quantile forecasts (or vice versa) in which case the first (second) quantile forecast is said to encompass the second (first).

The aim of the CQFE test is to test the null hypothesis, that $\hat{q}_{1,t}$ performs better than any linear combination of $\hat{q}_{1,t}$ and $\hat{q}_{2,t}$. Under the null hypothesis, it holds

$$E_t\left(\rho_\tau(y_{t+1} - \hat{q}_{1,t})\right) \leq E_t\left(\rho_\tau(y_{t+1} - \theta_0 - \theta_1\hat{q}_{1,t} - \theta_2\hat{q}_{2,t})\right), \tag{24}$$

that is satisfied if and only if the weights (θ_1, θ_2) are equal to $(1, 0)$. The objective function of the GMM is:

$$J_T = g_T(\theta)' \mathbf{W}_T g_T(\theta).$$

The optimal weights are computed as:

$$\theta^* = \arg \min_{\theta} g_T(\theta)' \mathbf{W}_T g_T(\theta), \quad g_T(\theta) = \frac{\sum_{t=1}^T (\tau - \mathbf{1}_{\tau}\{y_{t+1} - \theta' \mathbf{q}_t < 0\}) \mathbf{z}_T}{T},$$

where \mathbf{W}_T is a positive definite matrix, $g_T(\theta)$ is the sample moment condition, $\theta = (\theta_0, \theta_1, \theta_2)'$ is a set of weights, $\theta^* = (\theta_0^*, \theta_1^*, \theta_2^*)'$ denotes the optimal weights, and $\hat{\mathbf{q}}_t = (1, \hat{q}_{1,t}, \hat{q}_{2,t})'$ is a vector with the forecasted values based on the pairwise models 1, and 2 in the CQFE test, m denotes the out-of-sample size and \mathbf{z}_T is a vector of instruments. [Hansen \(1982\)](#) showed that by setting $\mathbf{W}_T = \mathbf{S}_T^{-1}$ that is, the inverse of an asymptotic covariance matrix, is optimal as it estimates θ^* with as small as possible asymptotic variance. \mathbf{S} is also known as the spectral density matrix of g_T . We follow [Newey and West \(1987\)](#) and use a heteroskedasticity robust estimate $\hat{\mathbf{S}}_T$, of \mathbf{S} defined as:

Table 6. Comparison of the forecasting methods

Model	τ	Polynomial	Splines	MIDAS	Deep QR	Deep MIDAS	Deep LASSO	Deep LASSO MIDAS	Deep Ridge	Deep Ridge MIDAS	Deep Elnet	Deep Elnet MIDAS
GARCH	1%	0.909	0.791***	1.083***	0.346***	0.346***	0.236***	0.335***	0.177***	0.453***	0.146***	0.469***
	5%	1.256	1.182***	1.364***	0.331***	0.430***	0.190***	0.364***	0.223***	0.331***	0.190***	0.529***
	10%	1.486	1.554***	1.486***	0.189***	0.459***	0.122***	0.419***	0.270***	0.405***	0.189***	0.541**
RM	1%	0.780	0.570***	1.066***	0.289***	0.246***	0.141***	0.164***	0.207***	0.220***	0.239***	0.410***
	5%	1.000	0.806***	1.232***	0.348***	0.174***	0.142***	0.174***	0.29***	0.213***	0.168***	0.310***
	10%	1.282	1.176***	1.447***	0.400***	0.259***	0.129***	0.224***	0.329***	0.235***	0.176***	0.200***
SV	1%	1.000	1.511***	0.838	0.198***	0.850	0.024***	0.871	0.225***	0.850	0.498***	0.895
	5%	1.000	1.716***	0.795	0.148***	0.824	0.062***	0.795	0.097***	0.830	0.102***	0.864
	10%	1.000	1.264***	0.736*	0.142***	0.745*	0.047***	0.708**	0.151***	0.764	0.104***	0.792
ASV	1%	1.000	1.040***	4.918***	0.235***	0.601*	0.096***	0.586*	0.014***	0.584*	0.037***	0.586*
	5%	1.000***	1.034***	4.559**	0.080***	0.584*	0.122***	0.576*	0.067***	0.567*	0.063***	0.584*
	10%	1.000	1.275***	1.209***	0.033***	0.582*	0.013***	0.542**	0.098***	0.542**	0.052***	0.536**

Notes: The table reports relative RMSFE. The smaller the entry (<1) the better the forecast. *, **, and *** denote results from DM test with the [Harvey, Leybourne, and Newbold \(1997\)](#) adjustment for predictive accuracy, indicating rejection of the null hypothesis that the models have the same predictive accuracy at the 10%, 5%, and 1% levels of significance, respectively.

$$\widehat{S}_T = \widehat{S}_0 + \sum_{j=1}^m \left(1 - \frac{j}{m+1}\right) (\widehat{S}_j + \widehat{S}'_j), \quad \text{where} \quad \widehat{S}_j = \frac{1}{T} \sum_{t=j+1}^T g_t(\widehat{\theta}) g_{t-j}(\widehat{\theta}).$$

\widehat{S}_0 is the estimated spectral density matrix evaluated at frequency zero. The GMM estimation is performed recursively, that is, (i) minimize J_T using an identity weighting matrix to get θ^* , which gives \mathbf{W}_T via \widehat{S}_T and (ii) minimize J_T using $\mathbf{W}_T = \widehat{S}_T^{-1}$ from step (i).

Consequently, we consider two separate test $H_{10} : (\theta_1^*, \theta_2^*) = (1, 0)$ versus $H_{1a} : (\theta_1^*, \theta_2^*) \neq (1, 0)$ and $H_{20} : (\theta_1^*, \theta_2^*) = (0, 1)$ versus $H_{2a} : (\theta_1^*, \theta_2^*) \neq (0, 1)$, which correspond to testing whether forecast $\widehat{q}_{1,t}$ encompasses $\widehat{q}_{2,t}$ or $\widehat{q}_{2,t}$ encompasses $\widehat{q}_{1,t}$. Then the CQFE statistics are defined as:

$$\begin{aligned} \text{ENC}_1 &= T \left((\theta_1^*, \theta_2^*) - (1, 0) \right)' \widehat{\Omega} \left((\theta_1^*, \theta_2^*) - (1, 0) \right) \\ \text{ENC}_2 &= T \left((\theta_1^*, \theta_2^*) - (0, 1) \right)' \widehat{\Omega} \left((\theta_1^*, \theta_2^*) - (0, 1) \right), \end{aligned}$$

where $\widehat{\Omega} = g_T(\theta)' S^{-1} g_T(\theta)$. The asymptotic distribution of the GMM estimates of θ requires the moment conditions to be once differentiable. To satisfy this requirement, we follow [Giacomini and Komunjer \(2005\)](#) and replace the moment condition with the following smooth approximation:

$$g_\tau(\theta) = \frac{\sum_{t=1}^T [\tau - (1 - \exp((y_{t+1} - \theta' \widehat{q}_t)/\eta))] \mathbf{1}\{y_{t+1} - \theta' \widehat{q}_t < 0\} z_T}{T},$$

where η is the smoothing parameter. We choose the critical values, c_{crit} of the test from a χ^2 distribution, in which $\widehat{q}_{i,t}$ encompasses $\widehat{q}_{j,t}$, if $\text{ENC}_i \leq c_{\text{crit}} \forall i \neq j = 1, 2$. In the empirical application, the vector of instruments, z_T , is $(1, r_t, \text{VaR}_{i,t}, \text{VaR}_{j,t}), \forall i \neq j = 1, 2$.

We select η to be 0.005, following the CQFE test rejection probabilities in [Giacomini and Komunjer \(2005\)](#), since our POOS size is 2000 observations. We consider the following five blocks: (i) the non-parametric, (ii) the non-linear, (iii) the non-linear MIDAS, (iv) the linear, and (v) the linear MIDAS blocks. The non-parametric block consists of the quantile polynomial and quantile splines estimators, the non-linear block consists of the deep quantile estimators for the different regularization schemes and the non-linear MIDAS block consists of the deep MIDAS estimators for the different regularization schemes. Finally, the linear and linear MIDAS blocks consist of the linear quantile and linear quantile MIDAS estimators, respectively.

We examine each block of models across different quantiles. Specifically, we consider how many times the models within a specific block outperform models from other blocks and present these results in [Table 7](#). Under this setting, a *win* denotes that the prevailing model encompasses the competing benchmark model, while a *loss* means that the competing model encompasses the prevailing one. Precisely, we consider a win when the computed p -value of the CQFE test fails to reject the null hypothesis, that is, H_{10} or H_{20} . On the contrary, in the case where the CQFE test suggests that there is no encompassing between the forecasts, we consider this as a loss, that is, the null hypothesis is rejected. Furthermore, the CQFE test has a gray zone in which the test can fail to reject both null hypotheses (H_{10} and H_{20}), hence the test is inconclusive. Below we summarize the CQFE testing results for the different quantiles when $\eta = 0.005$.

For the 10th quantile, the non-linear block encompasses 713 times the competing blocks, in comparison to the linear block, which encompasses the competing blocks 173 times and the non-parametric block that encompasses the others 322 times. The linear block does not encompass other blocks less than 15 times and the non-linear block for 86 times. Additionally, the test is inconclusive 696 times for the non-linear block and 159 times for the linear one. Thus, the non-linear block is ranked first in terms of how many times it encompasses the other blocks and the non-linear MIDAS block is ranked second.

For the 5th quantile, the non-linear block encompasses 739 times other blocks, 342 times the non-parametric, and the linear 170 times. Further, the linear block does not encompass the other blocks 18 times and the non-linear 60 times. Finally, for the non-linear block, the CQFE test is inconclusive 726 times and 166 times for the linear block. The ranking of the first two blocks is the same as in the 10th quantile.

Finally, we examine the first quantile. In this case, the non-linear block encompasses 757 times the other blocks, 336 times the non-parametric, and the linear block 173 times. Furthermore, the linear block does not encompass 15 times the other blocks and the non-linear 42 times. The test is inconclusive 751 times for the non-linear block and 170 times for the linear one. The ranking remains the same as above. Results for different smoothing parameters η suggest similar patterns and are available upon request.

4 Semi-structural Analysis

A general issue in ML is the trade-off between accuracy and interpretability; where the output of a highly complicated model, for example, a deep neural network, can have great accuracy or forecasting performance, but cannot be easily interpreted. In this section, we first discuss the details of two methods that can be used to make ML methods interpretable. The first one is the SHAP values, which has received a lot of attention recently, and the second is partial derivatives. Further, we make a formal comparison on the output of both methods, based on the output of the deep quantile estimator that illustrates, (i) that both methods can be used to make the impact of each covariate in neural networks interpretable and (ii) perhaps surprisingly that the use of partial derivatives, offers more stable results at a fraction of the computational cost.

4.1 Shapley Values

SHAP values are a general class of additive attribution methods, based on the initial work of [Shapley \(1953\)](#) where the goal was to determine how to fairly split a pay-off among players in a cooperative game. In the context of ML, the goal of SHAP values is to explain the prediction of the dependent variable by estimating the contribution of each covariate to the prediction. SHAP values, following the exposition in [Lundberg and Lee \(2017\)](#) and [Lundberg, Erion, and Lee \(2018\)](#) can be constructed as follows.

Let $f(\mathbf{x}_t) = \widehat{G}(\mathbf{x}_t, \boldsymbol{w})$ be the output of the estimated model we wish to interpret, given a $p \times 1$ vector of covariates \mathbf{x}_t , and \widehat{f} the explanation model, to be defined below. Further, let \mathbf{x}_t^\dagger be the $M \times 1$ subset (vector) of \mathbf{x}_t that contains simplified covariates. These simplified covariates, can be mapped to the original through a mapping function $h_{x_i}(\cdot)$, such that $\mathbf{x}_t = h_{x_i}(\mathbf{x}_t^\dagger)$. Then under the local accuracy property of [Lundberg and Lee \(2017\)](#), if there exists a vector, \mathbf{z}_t^\dagger , with binary inputs, such that $\mathbf{z}_t^\dagger \approx \mathbf{x}_t^\dagger$, then $\widehat{f}(\mathbf{z}_t^\dagger) \approx f(h_{x_i}(\mathbf{z}_t^\dagger))$, where the explanation model (i.e., the additive attribution function) is

Table 7. Entries of the table present the number of times a block encompasses (wins), does not encompass (losses) and is inconclusive, according to the CQFE test for different quantiles τ

$\eta =$ 0.005	Block	$\tau = 1\%$			$\tau = 5\%$			$\tau = 10\%$		
		wins	losses	inconclusive	wins	losses	inconclusive	wins	losses	Inconclusive
	Linear	173	15	170	170	18	166	173	15	159
	Non-parametric	336	40	334	342	34	332	322	54	306
	Non-linear	757	42	751	739	60	726	713	86	696
	MIDAS	175	13	173	174	14	172	179	9	169
	Non-linear MIDAS	709	43	704	700	52	687	657	95	638

Notes: Results are reported for $\eta = 0.005$.

$$\widehat{f}(z_t^\dagger) = \phi_0 + \sum_{i=1}^M \phi_i z_{t,i}^\dagger, \tag{25}$$

and $f(z_t^\dagger)$ represents the linear decomposition of the original ML model, where ϕ_0 is the intercept, $\phi_i \in \mathbb{R}$ is the effect to each dependent variable $z_i^\dagger \in (0, 1)$, that provides local and global inference at the same time. If $z_{t,i} = 1$ then the covariate is observed, on the contrary, if $z_{t,i} = 0$ then the covariate is unknown. Under the following three properties: (i) local accuracy that is, the explanation function should match the original model, (ii) missingness, which ensures that input variable have no attributed effect and (iii) consistency, under which, if an input variables is important, then the effect to each dependent variable should not decline, the SHAP value is

$$\phi_i = \sum_{M \subseteq p \setminus \{i\}} \frac{|M|!(p - |M| - 1)!}{p!} [f_{M \cup \{i\}}(x_{M \cup \{i\}}) - f_M(x_M)], \tag{26}$$

where p is the set of all predictors, $|M|$ is the number of non-zero elements in x_t^\dagger , $f_M(x_M)$ is the model's output using except from the i th covariate, and $f_{M \cup \{i\}}(x_{M \cup \{i\}})$ is the output of the model, when $\{i\}$ is included in the covariate set.

The calculation of SHAP values can be computationally expensive, as it requires 2^N possible permutations of the predictors. For the case of deep neural networks Lundberg and Lee (2017), and Shrikumar, Greenside, and Kundaje (2017), have shown that *DeepLIFT* can be used as an approximation of the deep SHAP that is computationally feasible,⁵ preserving the three properties above. *DeepLIFT* is a recursive prediction explanation method for deep learning. The Additive feature attribution methods analogy of *DeepLIFT* is called the summation-to-delta property is

$$\sum_{i=1}^p C_{\Delta x_i, \Delta o} = \Delta o. \tag{27}$$

Then, the SHAP values can be obtained as

5 There are other methods that can be used to achieve this, such as Tree Explainer, Kernel Explainer, Linear Explainer, and Gradient Explainer.

$$\phi_i = C_{\Delta x_{t,i} \Delta o},$$

where $C_{\Delta x_{t,i} \Delta o}$, represents the impact of a covariate to a reference value relative to the initial value, is assigned to each $x_{t,i}$ covariate, $o = f(\cdot)$ is the output of the model, $\Delta o = f(x) - f(r)$, $\Delta x_{t,i} = f(x_{t,i}) - r_{t,i}$ and r the reference value. Equation (27) matches Equation (25), if in Δo we set $\phi_0 = f(r_{t,i})$ and $\phi_i = C_{\Delta x_{t,i} \Delta o}$.

4.2 Partial Derivatives

The use of partial derivatives for the interpretation of a model is straight forward in econometrics, with various uses, ranging from the simple linear regression model to impulse response analysis. In this section, we show how partial derivatives can be used even in highly non-linear deep neural networks. Before we start the analysis, note that while the deep neural networks are highly non-linear, their solution/output via SGD optimization methods can be treated as differentiable function, as the majority of activation functions are differentiable. Let's consider the case of ReLU, that is, not differentiable at 0, whereas it is in every other point. From the point of gradient descent, heuristically, it works well enough to treat it as a differentiable function. Further, Goodfellow, Bengio, and Courville (2016) argue that this issue is negligible and ML softwares are prone to rounding errors, which make it very unlikely to compute the gradient at a singularity point. Note that even in this extreme case, both SGD and ADAM, will use the right subgradient at 0.

For a general $x_t \in \mathbb{R}^p$, let

$$d_{j,i,t} = \frac{\partial \widehat{G}_{j,\tau}(x_t, w)}{\partial x_{j,i,t-1}}, \quad (28)$$

denote the partial derivative of covariate $x_i = x_{it}$, for $i = 1, \dots, p$ at time $t = 1, \dots, T$, $\widehat{G}_{j,\tau}(x_t, w)$ is the forecasted VaR $_{j,t}$, across the j different VaR specifications we consider. We assess the partial derivative in time, since, following Kapetanios (2007), we expect it to vary in time, due to the inherent non-linearity of the neural network. Our covariate(s) x_t is the conditional volatility for GARCH and RM, VaR lagged values, the absolute S&P500 daily return, and the positive and negative S&P500 daily returns for SV and ASV, respectively. It is evident that under the classic linear regression problem, or linear quantile regression model, the effect of the covariates x_t to the dependent variable y_t is constant, time invariant, and corresponds to $\widehat{\beta}(\tau)$.

4.3 Results

In this application, we use the whole sample size, that is, around 36 years of daily returns on the S&P500 index to provide an accurate interpretation of the deep quantile estimator. Figures 1–4 illustrate the partial derivatives and SHAP values evaluated in time on the output of the deep quantile⁶ estimator, for a specific quantile τ . Further, we compare the partial derivatives of the deep quantile estimator relative to the linear quantile regression partial derivative, that is, the $\beta(\tau)$ coefficient. Both partial derivatives and SHAP values

6 In this section, we limit our attention in the output of the best performing model, in terms of its forecasting capacity, as reflected by the forecast gains measure in Section 3, for each model, based on the different penalization schemes. Results from all the different penalization schemes suggest similar patterns to the ones discussed above and are available upon request.

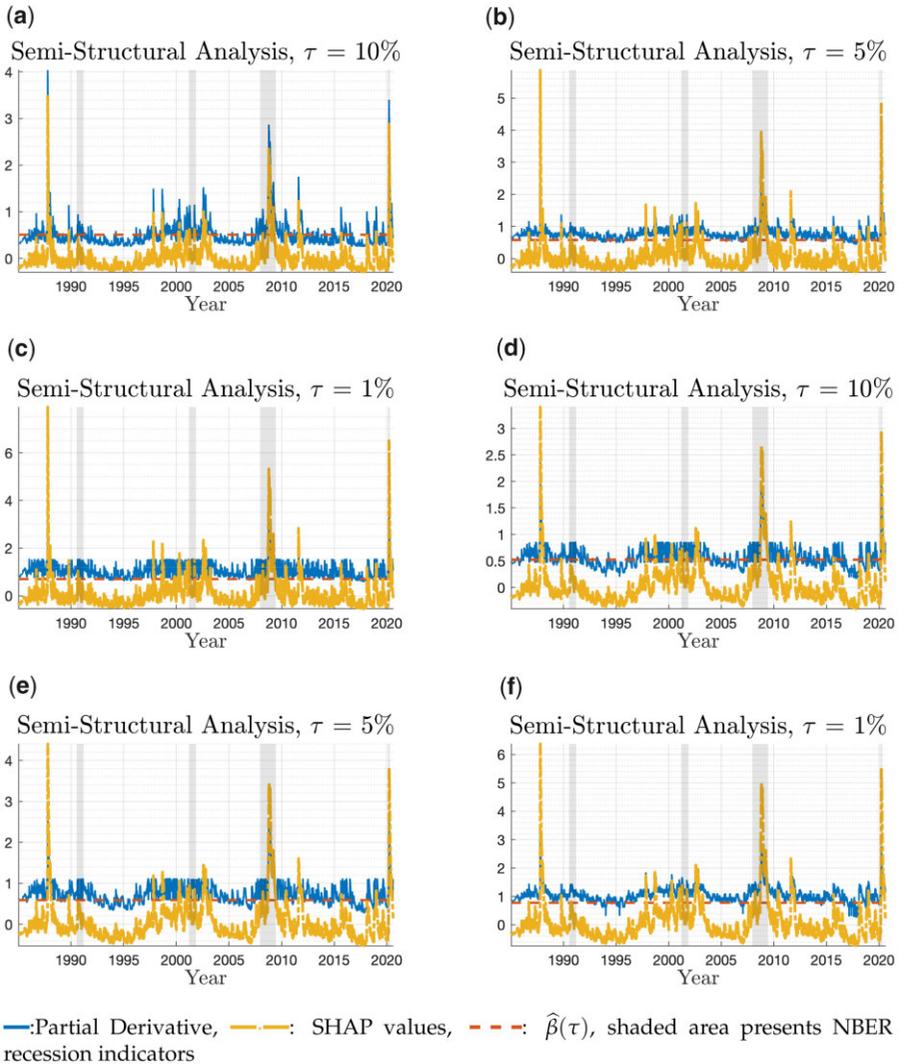


Figure 1. Partial derivative, SHAP, and $\hat{\beta}(\tau)$ for GARCH and RM models. (a) GARCH without penalty. (b) GARCH with Elnet penalty. (c) GARCH with Elnet penalty. (d) RM without penalty. (e) RM without penalty. (f) RM with Ridge penalty.

seem to identify interesting patterns that can be linked to some well-known events. Below we discuss our results for all models we have considered in our empirical application.

The results for the first two models, that is, GARCH and RM can be summarized together, since in both models there is only one covariate, that is the conditional volatility, but with a different specification. The results from this model are illustrated in Figure 1. We find that the partial derivative appears to be more stable over time, fluctuating around the constant partial derivative, $\beta(\tau)$, of the linear quantile estimator. When there is a crisis or a stressful event in the financial markets, they increase. As an example, we see significant

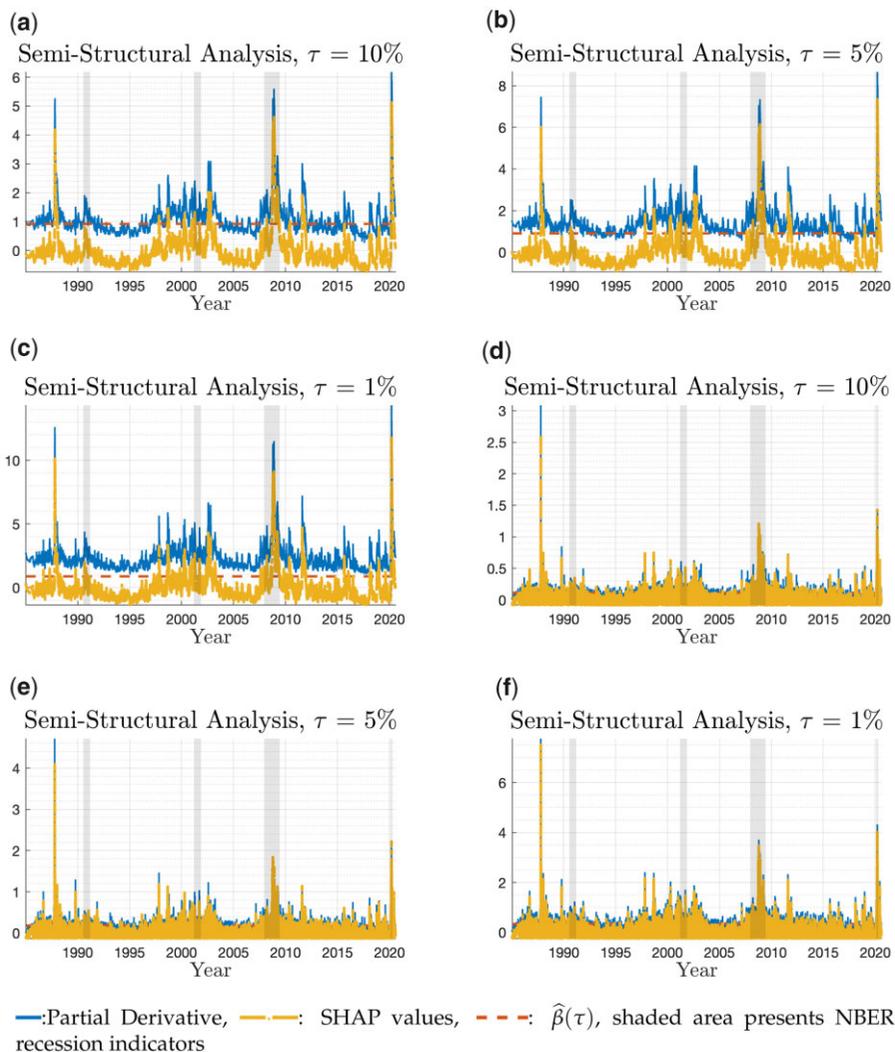


Figure 2. Partial derivative, SHAP, and $\hat{\beta}(\tau)$ for SV model. (a) SV without penalty. (b) SV without penalty. (c) SV with LASSO penalty. (d) VaR lagged values without penalty. (e) VaR lagged values without penalty. (f) VaR lagged values LASSO penalty.

spikes in the partial derivatives, both in March 2020 as well as in 2008, which stand for the onset of the COVID-19 pandemic and the *Great Recession*, respectively. We also find that the biggest increase occurs in 1987, the year when *Black Monday* happened, and also significant variation during the U.S. government shutdown in 2019. The values for the partial derivatives generally increase, as τ decreases. SHAP values have a similar behavior with the partial derivatives, but are more volatile across time. For the first two models, there are some events, for example, during the 1991, where the values for both SHAP and partial

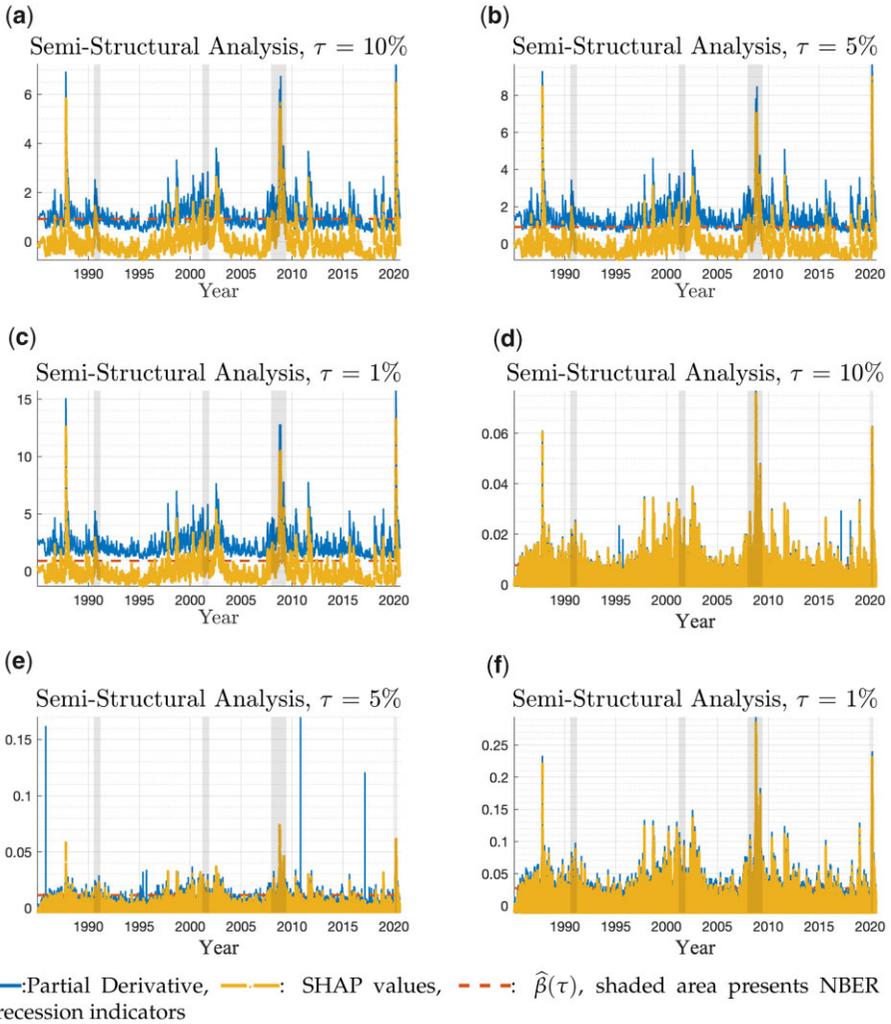


Figure 3. Partial derivative, SHAP, and $\hat{\beta}(\tau)$ for ASV model. (a) ASV with Ridge penalty. (b) ASV with Ridge penalty. (c) ASV without penalty. (d) S&P500 positive values with Ridge penalty. (e) S&P500 positive values with Ridge penalty. (f) S&P500 positive values without penalty.

derivatives do not increase a lot. We view this finding as an inability of these two models, to properly account for this crisis.

In the last two models, the merit of SHAP values and partial derivatives becomes clear, since in these models we have more than one covariate and both methods can provide an indication on the effect of each covariate on the final output. Overall, we find that increasing the number of covariates, allow the models to account for all crises within the sample. For the case of the SV model, we find that the important covariate is the lagged values of VaR, rather than the absolute values of S&P500. Similar to the one covariate models, we

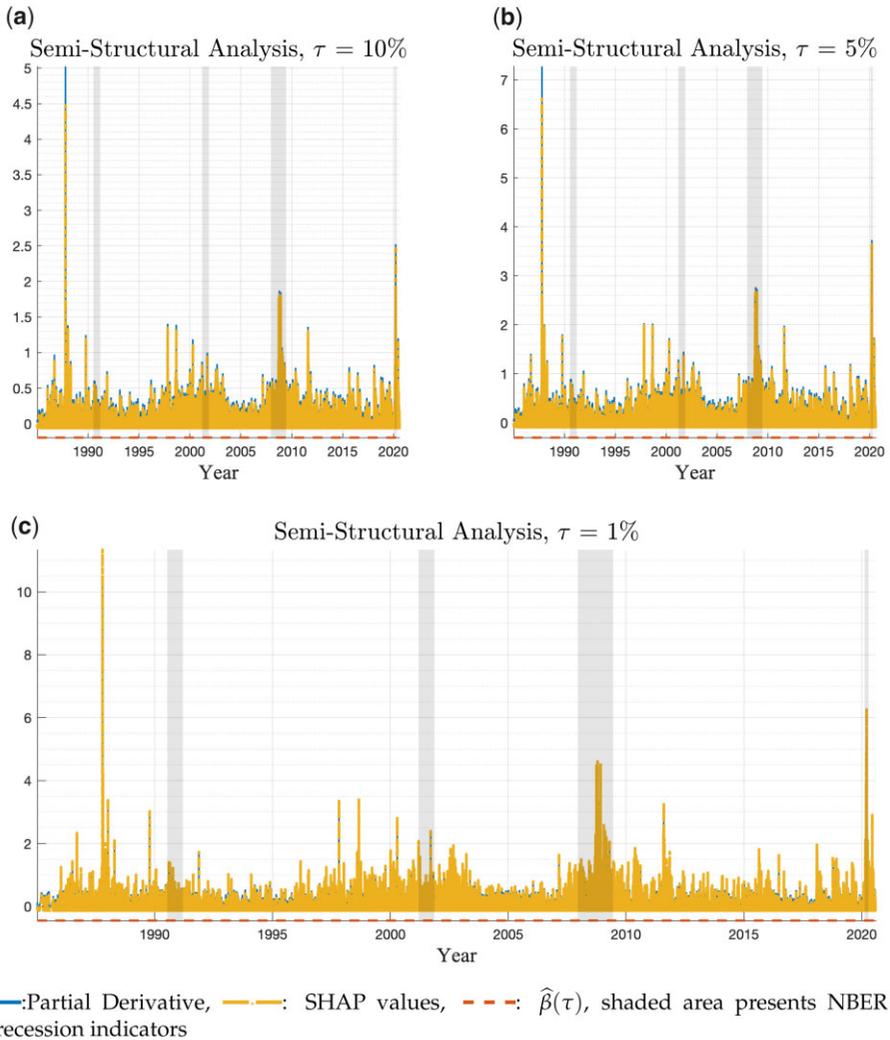


Figure 4. Partial derivative, SHAP, and $\hat{\beta}(\tau)$ for ASV model. (a) S&P500 negative values with Ridge penalty. (b) S&P500 negative values with Ridge penalty. (c) S&P500 negative values without penalty.

find that the partial derivatives are more stable than SHAP values, fluctuating closely around $\hat{\beta}(\tau)$ and picking up when there are crisis or distress in the economy or financial markets. The SHAP values again appear to be more volatile with a wider range. Similar to the findings of the one covariate models, the higher the values for the partial derivative and SHAP, the lower the τ quantile.

For the case of the ASV model, we find that again the lagged values of VaR are the most significant covariate, the negative S&P500 returns have some impact and the positive S&P500 returns are almost insignificant. Similar to the cases above, we find that the partial derivative is more stable than SHAP values, fluctuating closely around $\hat{\beta}(\tau)$ and picking up when there is a crisis or distress in the economy or financial markets. The SHAP values

again appear to be more volatile with a wider range. Again and same as before, lower quantiles have higher partial derivatives. The results for these two models are illustrated in Figures 2–4.

Different penalization schemes maintain the aforementioned results, with a lower magnitude. Overall, we observe that the linear quantile regression shows a fixed pattern across time, and is evident that this model does not anticipate shocks in the economy. Our analysis suggests that it is higher during stressful events. As Engle and Manganelli (2004) suggest, SV and ASV react more to negative shocks, and in stressful events, their spike is larger than the GARCH and RM models. Finally, covariates with the minimum contribution on the forecasted values, such as the positive S&P500 returns have negligible impact on both SHAP and partial derivatives values.

5 Conclusion

In this article, we contribute to the expanding literature on the use of ML in finance and use the deep quantile estimator that has the potential to capture the non-linear association between asset returns and predictors. In Section 1, we lay out the exact workings of the deep quantile estimator, and illustrate how it generalizes linear quantile regression.

In the MC exercise in Section 2, we study the finite sample properties of the deep quantile estimator, based on a number of data-generating processes. We present extensive evidence the estimator gives good finite sample performance, that is a function of T , uniformly across different regularization schemes.

We use the deep quantile estimator, with various penalization schemes, to forecast VaR. We find that the deep quantile estimator gives considerable predictive gains, up to 98%, relative to the VaR forecasts produced by the linear quantile regression. This result is backed by the forecasting accuracy tests, that is, the DM, the GW, and the quantile score tests. Further, results from the CQFE test of Giacomini and Komunjer (2005) suggest that forecasts obtained from the non-linear estimators encompass forecasts from the linear and non-parametric models with a higher frequency. These findings are in support of the non-linear association between the conditional quantile of asset returns and covariates, hence suggesting a new avenue in forecasting in finance and in macroeconomics during extreme events.

In addition, we do a semi-structural analysis to examine the contribution of the predictors in VaR over time. We consider, following the ML literature, SHAP values, and further partial derivatives. Our findings suggest that the non-linear estimator reacts more in stressful events and exhibits time variation, while the linear quantile estimator presents, as expected, a constant time-invariant behavior. We conclude that financial variables are characterized by non-linearities, which the deep quantile estimator can approximate quite well.

Finally, we make a formal comparison between SHAP and partial derivatives, and interestingly find that partial derivatives can be used to make ML methods interpretable, are less volatile, easier to interpret and can be computed at a fraction of time used in the calculation of SHAP values.

Supplemental Data

Supplemental data are available at <https://www.datahostingsite.com>.

References

- Adams, P. A., T. Adrian, N. Boyarchenko, and D. Giannone. 2021. Forecasting Macroeconomic Risks. *International Journal of Forecasting* 37: 1173–1191.
- Athey, S., and G. W. Imbens. 2017. The State of Applied Econometrics: Causality and Policy Evaluation. *Journal of Economic Perspectives* 31: 3–32.
- Babii, A., X. Chen, E. Ghysels, and R. Kumar. 2020. Binary Choice with Asymmetric Loss in a Data-rich Environment: Theory and an Application to Racial Justice, *arXiv preprint arXiv: 2010.08463*.
- Babii, A., E. Ghysels, and J. Striaukas. 2022a. High-Dimensional Granger Causality Tests with an Application to VIX and News. *Journal of Financial Econometrics* nbac023.
- Babii, A., E. Ghysels, and J. Striaukas. 2022b. Machine Learning Time Series Regressions with an Application to Nowcasting. *Journal of Business & Economic Statistics* 40: 1094–1106.
- Bates, J. M., and C. W. Granger. 1969. The Combination of Forecasts. *Journal of the Operational Research Society* 20: 451–468.
- Baur, D., and N. Schulze. 2005. Coexceedances in Financial Markets—A Quantile Regression Analysis of Contagion. *Emerging Markets Review* 6: 21–43.
- Belloni, A., V. Chernozhukov, D. Chetverikov, and I. Fernández-Val. 2019. Conditional Quantile Processes Based on Series or Many Regressors. *Journal of Econometrics* 213: 4–29.
- Belloni, A., V. Chernozhukov, and C. Hansen. 2014. Inference on Treatment Effects after Selection among High-Dimensional Controls. *The Review of Economic Studies* 81: 608–650.
- Bollerslev, T. 1986. Generalized Autoregressive Conditional Heteroskedasticity. *Journal of Econometrics* 31: 307–327.
- Borup, D., D. Rapach, and E. C. M. Schütte. 2022. Mixed-Frequency Machine Learning: Nowcasting and Backcasting Weekly Initial Claims with Daily Internet Search-Volume Data. *International Journal of Forecasting*.
- Bucci, A. 2020. Realized Volatility Forecasting with Neural Networks. *Journal of Financial Econometrics* 18: 502–531.
- Carriero, A., G. Kapetanios, and M. Marcellino. 2009. Forecasting Exchange Rates with a Large Bayesian VAR. *International Journal of Forecasting* 25: 400–417.
- Chen, X, Y. Liu, S. Ma, and Z. Zhang. 2020. Efficient Estimation of General Treatment Effects Using Neural Networks with a Diverging Number of Confounders, *arXiv preprint arXiv: 2009.07055*.
- Chernozhukov, V., and L. Umantsev. 2001. Conditional Value-at-Risk: Aspects of Modeling and Estimation. *Empirical Economics* 26: 271–292.
- Diebold, F. X., and R. S. Mariano. 1995. Comparing Predictive Accuracy. *Journal of Business & Economic Statistics* 13: 253–263.
- Du, Z., M. Wang, and Z. Xu. 2019. “On Estimation of Value-at-Risk with Recurrent Neural Network.” *2019 Second International Conference on Artificial Intelligence for Industries (AI4I)*, Laguna Hills, CA, USA, 103–106.
- Engle, R. F., and S. Manganelli. 2004. CAViaR: Conditional Autoregressive Value at Risk by Regression Quantiles. *Journal of Business & Economic Statistics* 22: 367–381.
- Farrell, M. H., T. Liang, and S. Misra. 2021. Deep Neural Networks for Estimation and Inference. *Econometrica* 89: 181–213.
- Gallant, A. R., and H. White. 1992. On Learning the Derivatives of an Unknown Mapping with Multilayer Feedforward Networks. *Neural Networks* 5: 129–138.
- Ghysels, E., A. Plazzi, and R. Valkanov. 2016. Why Invest in Emerging Markets? The Role of Conditional Return Asymmetry. *The Journal of Finance* 71: 2145–2192.
- Ghysels, E., A. Plazzi, R. Valkanov, A. Rubia, and A. Dossani. 2019. Direct versus Iterated Multiperiod Volatility Forecasts. *Annual Review of Financial Economics* 11: 173–195.

- Ghysels, E., P. Santa-Clara, and R. Valkanov. 2004. *The MIDAS Touch: Mixed Data Sampling Regression Models*.
- Giacomini, R., and I. Komunjer. 2005. Evaluation and Combination of Conditional Quantile Forecasts. *Journal of Business & Economic Statistics* 23: 416–431.
- Giacomini, R., and H. White. 2006. Tests of Conditional Predictive Ability. *Econometrica* 74: 1545–1578.
- Goodfellow, I., Y. Bengio, and A. Courville. 2016. *Deep Learning*. Cambridge, MA: MIT Press.
- Gu, S., B. Kelly, and D. Xiu. 2020. Empirical Asset Pricing via Machine Learning. *The Review of Financial Studies* 33: 2223–2273.
- Gu, S., B. Kelly, and D. Xiu. 2021. Autoencoder Asset Pricing Models. *Journal of Econometrics* 222: 429–450.
- Hansen, L. P. 1982. Large Sample Properties of Generalized Method of Moments Estimators. *Econometrica* 50: 1029–1054.
- Harvey, D., S. Leybourne, and P. Newbold. 1997. Testing the Equality of Prediction Mean Squared Errors. *International Journal of Forecasting* 13: 281–291.
- He, Z., and A. Krishnamurthy. 2013. Intermediary Asset Pricing. *American Economic Review* 103: 732–770.
- Hornik, K. 1991. Approximation Capabilities of Multilayer Feedforward Networks. *Neural Networks* 4: 251–257.
- Hornik, K., M. Stinchcombe, and H. White. 1989. Multilayer Feedforward Networks Are Universal Approximators. *Neural Networks* 2: 359–366.
- Joseph, A. 2019. “Shapley regressions: a framework for statistical inference on machine learning models.” Discussion paper, Bank of England.
- Kapetanios, G. 2007. Measuring Conditional Persistence in Nonlinear Time Series. *Oxford Bulletin of Economics and Statistics* 69: 363–386.
- Kapetanios, G., and A. P. Blake. 2010. Tests of the Martingale Difference Hypothesis Using Boosting and RBF Neural Network Approximations. *Econometric Theory* 26: 1363–1397.
- Keilbar, G., and W. Wang. 2022. Modelling Systemic Risk Using Neural Network Quantile Regression. *Empirical Economics* 62: 93–118.
- Kingma, D. P., and J. Ba. 2014. Adam: A Method for Stochastic Optimization, *arXiv preprint arXiv:1412.6980*.
- Koenker, R. 2005. *Quantile Regression, Econometric Society Monographs*. Cambridge: Cambridge University Press.
- Koenker, R., and G. Bassett. 1978. Regression Quantiles. *Econometrica* 46: 33–50.
- Koenker, R., V. Chernozhukov, X. He, and L. Peng. 2017. *Handbook of Quantile Regression*. Boca Raton, FL: CRC Press.
- Koenker, R., and K. F. Hallock. 2001. Quantile Regression. *Journal of Economic Perspectives* 15: 143–156.
- Liang, S., and R. Srikant. 2016. Why Deep Neural Networks for Function Approximation, *arXiv preprint arXiv:1610.04161*
- Lundberg, S. M., G. G. Erion, and S. I. Lee. 2018. Consistent Individualized Feature Attribution for Tree Ensembles, *arXiv preprint arXiv:1802.03888*.
- Lundberg, S. M., and S. I. Lee. 2017. “A Unified Approach to Interpreting Model Predictions.” In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan and R. Garnett (eds.), *Advances in Neural Information Processing Systems*. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf
- Meinshausen, N. 2006. Quantile Regression Forests. *Journal of Machine Learning Research* 7: 983–999.

- Morgan, J. P. 1996. *Reuters (1996) RiskMetrics-Technical Document*. New York, NY: JP Morgan.
- Newey, W. K., and K. D. West. 1987. Hypothesis Testing with Efficient Method of Moments Estimation. *International Economic Review* 28: 777–787.
- Padilla, O. H. M., W. Tansey, and Y. Chen. 2022. Quantile Regression with ReLU Networks: Estimators and Minimax Rates. *Journal of Machine Learning Research* 23: 1–42.
- Park, J., and I. W. Sandberg. 1991. Universal Approximation Using Radial-Basis-Function Networks. *Neural Computation* 3: 246–257.
- Pohl, W., K. Schmedders, and O. Wilms. 2018. Higher Order Effects in Asset Pricing Models with Long-Run Risks. *The Journal of Finance* 73: 1061–1111.
- Schmidt-Hieber, J. 2020. Nonparametric Regression Using Deep Neural Networks with ReLU Activation Function. *The Annals of Statistics* 48: 1875–1897.
- Shapley, L. S. 1953. A Value for n-Person Games. *Contributions to the Theory of Games* 2: 307–317.
- Shrikumar, A., P. Greenside, and A. Kundaje. 2017. “Learning important features through propagating activation differences.” *International Conference on Machine Learning*, Sydney NSW Australia, 3145–3153.
- Smalter Hall, A., and T. R. Cook. 2017. “Macroeconomic indicator forecasting with deep neural networks.” Working paper, Federal Reserve Bank of Kansas City.
- Tambwekar, A., A. Maiya, S. Dhavala, and S. Saha. 2022. Estimation and Applications of Quantiles in Deep Binary Classification. *IEEE Transactions on Artificial Intelligence* 3: 275–286.
- Tobias, A., and M. K. Brunnermeier. 2016. CoVaR. *American Economic Review* 106: 1705–1741.
- Wager, S., and S. Athey. 2018. Estimation and Inference of Heterogeneous Treatment Effects Using Random Forests. *Journal of the American Statistical Association* 113: 1228–1242.
- Xu, Q., S. Liu, C. Jiang, and X. Zhuo. 2021. QRNN-MIDAS: A Novel Quantile Regression Neural Network for Mixed Sampling Frequency Data. *Neurocomputing* 457: 84–105.
- Yarotsky, D. 2017. Error Bounds for Approximations with Deep ReLU Networks. *Neural Networks: The Official Journal of the International Neural Network Society* 94: 103–114.
- Zhang, W., H. Quan, and D. Srinivasan. 2019. An Improved Quantile Regression Neural Network for Probabilistic Load Forecasting. *IEEE Transactions on Smart Grid* 10: 4425–4434.
- Zou, H., and T. Hastie. 2005. Regularization and Variable Selection via the Elastic Net. *Journal of the Royal Statistical Society Series B (Statistical Methodology)* 67: 301–320.