



# Nonparametric Quantile Regression with ReLU-Activated Recurrent Neural Networks



Hang Yu



Lyumin Wu



Wen-Xin Zhou



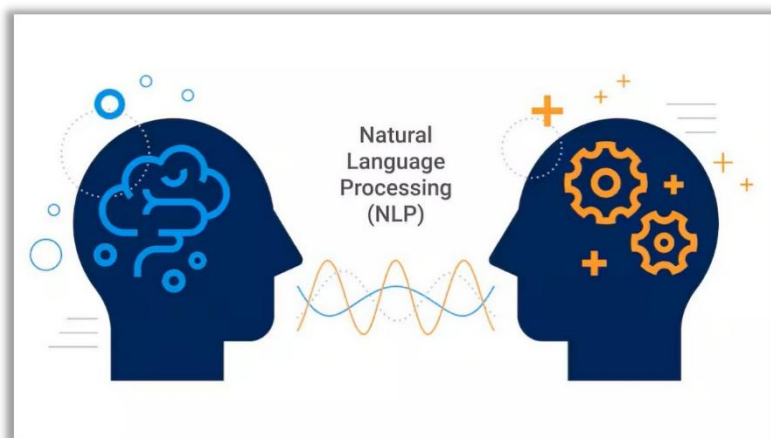
Zhao Ren

Presented by Hang Yu

2025.11.06

# Recurrent Neural Networks (RNNs)

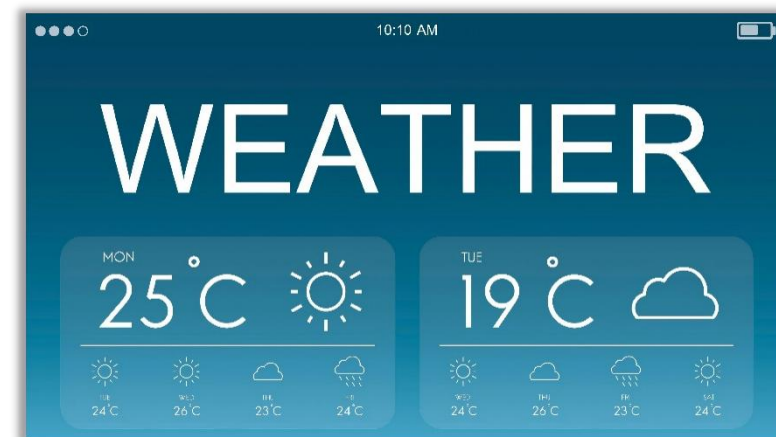
□ RNNs have achieved *remarkable success* in various applications.



*Natural Language Processing*



*Finance*



*Weather Forecasting*

However, the *theoretical* foundations of RNNs remain *incomplete*.

# RNNs and Sparse RNNs (SRNNs)

## □ RNN structure: Many-to-one setting

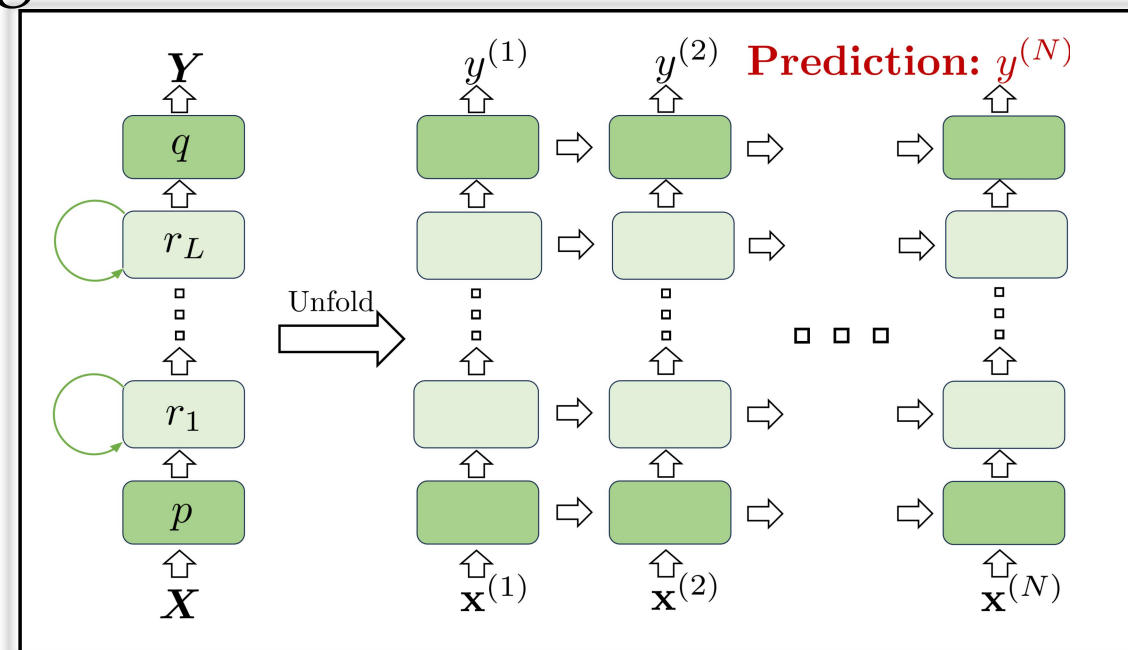
Given the width  $W$ , the length  $L$ , and the time horizon  $N$ , an RNN processes an input sequence  $\mathbf{X} = (\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)})$  sequentially through

Input layer  $p$ ,

Recurrent layers  $\{r_l\}_{l=1}^L$ ,

Output layer  $q$ ,

outputs  $\mathbf{Y} = (y^{(1)}, \dots, y^{(N)})$ , and obtains the **final prediction**  $y^{(N)}$ .



□ SRNN: The *sparsity*  $s$  of an RNN is defined as the number of its nonzero nodes.

We focus on establishing theoretical guarantees for RNNs/SRNNs  
in a fundamental area — **quantile regression**.

# Quantile Regression

## □ Problem setup

Consider the sequentially *stationary* observations  $\{(\mathbf{x}_t, y_t)\}_{t=1}^n$ , where any consecutive  $N$  observations share the same joint distribution as  $Z = ((X_1, Y_1), \dots, (X_N, Y_N))$ . Given a quantile level  $\tau \in (0, 1)$  of interest, we define the conditional  $\tau$ -th quantile of  $y_t$  (or  $Y_N$ ) given  $\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t$  (or  $X_1, \dots, X_N$ ) as

$$q_\tau(y_t | \mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t) = f_0(\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t), \quad N \leq t \leq n,$$

where  $f_0 : \mathbb{R}^{d_{\mathbf{x}} \times N} \rightarrow \mathbb{R}$  is the unknown conditional quantile function.

## □ Empirical risk minimization estimator:

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \mathcal{R}_n(f) := \frac{1}{n - N + 1} \sum_{t=N}^n \rho_\tau(y_t - f(\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t)),$$

where  $\rho_\tau(u) = (\tau - \mathbb{1}(u < 0))u$  is the check loss.

# Quantile Regression with **RNNs/SRNNs**

## □ Problem setup

*We consider the setting where  $\mathcal{F}$  is the class of RNNs/SRNNs.*

Consider the sequentially **stationary** observations  $\{(\mathbf{x}_t, y_t)\}_{t=1}^n$ , where any consecutive  $N$  observations share the same joint distribution as  $Z = ((X_1, Y_1), \dots, (X_N, Y_N))$ . Given a quantile level  $\tau \in (0, 1)$  of interest, we define the conditional  $\tau$ -th quantile of  $y_t$  (or  $Y_N$ ) given  $\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t$  (or  $X_1, \dots, X_N$ ) as

$$q_\tau(y_t | \mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t) = f_0(\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t), \quad N \leq t \leq n,$$

where  $f_0 : \mathbb{R}^{d_{\mathbf{x}} \times N} \rightarrow \mathbb{R}$  is the unknown conditional quantile function.

## □ Empirical risk minimization estimator:

$$\hat{f} \in \arg \min_{f \in \mathcal{F}} \mathcal{R}_n(f) := \frac{1}{n - N + 1} \sum_{t=N}^n \rho_\tau(y_t - f(\mathbf{x}_{t-N+1}, \dots, \mathbf{x}_t)),$$

where  $\rho_\tau(u) = (\tau - \mathbb{1}(u < 0))u$  is the check loss.

□ Classically, we assume the true function  $f_0$  belongs to a **Hölder class**.

**Definition 1** (Hölder Class of Functions  $\mathcal{C}_d^\beta(\mathcal{X}, K)$ ). Given a domain  $\mathcal{X} \subseteq \mathbb{R}^d$ , a positive Hölder smoothness parameter  $\beta$ , and a constant  $K > 0$ , the  $\beta$ -Hölder function class is defined as

$$\mathcal{C}_d^\beta(\mathcal{X}, K) = \left\{ f : \mathcal{X} \rightarrow \mathbb{R} \mid \sum_{\alpha: \|\alpha\|_1 < \beta} \|\partial^\alpha f\|_\infty + \sum_{\alpha: \|\alpha\|_1 = r} \sup_{\substack{\mathbf{x}, \mathbf{y} \in \mathcal{X} \\ \mathbf{x} \neq \mathbf{y}}} \frac{|\partial^\alpha f(\mathbf{x}) - \partial^\alpha f(\mathbf{y})|}{\|\mathbf{x} - \mathbf{y}\|_2^s} \leq K \right\},$$

where  $r = \lfloor \beta \rfloor$ ,  $s = \beta - r$ ,  $\partial^\alpha = \partial^{\alpha_1} \dots \partial^{\alpha_d}$  with  $\alpha = (\alpha_1, \dots, \alpha_d) \in \mathbb{N}_0^d$  and  $\|\alpha\|_1 = \sum_{i=1}^d \alpha_i$ . Moreover, we refer to  $\gamma = \beta/d$  as the dimension-adjusted degree of smoothness of  $\mathcal{C}_d^\beta(\mathcal{X}, K)$ .

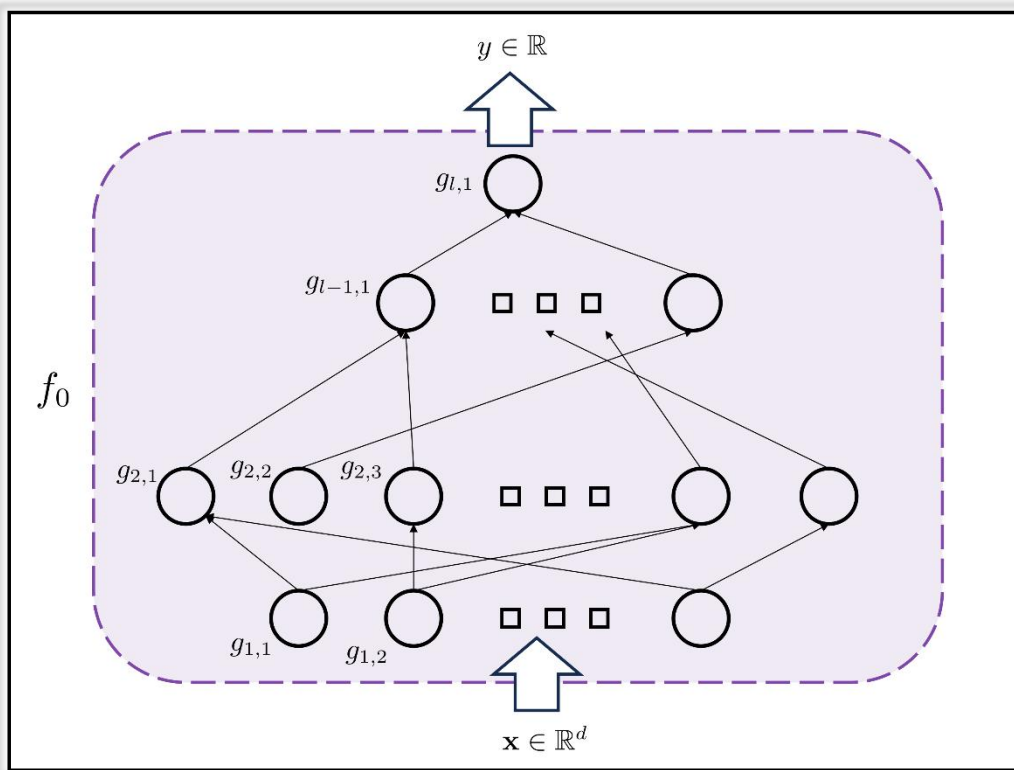
Stone (1982) established that the **minimax convergence rate** for estimating a regression function under the  $L_2$  norm over the Hölder class  $\mathcal{C}_d^\beta(\mathcal{X}, K)$  is  $n^{-\gamma/(2\gamma+1)}$ .

**Curse of Dimensionality:** In RNN applications,  $d$  is often **large**, resulting in a **small** value of the dimension-adjusted smoothness  $\gamma$ , which in turn leads to **slow** convergence rates.



# Hierarchical Interaction Model

- Assume  $f_0$  belongs to the hierarchical interaction model  $\mathcal{H}_d^l(\mathcal{P}, K)$ . Then  $f_0$  exhibits a *compositional structure* as follows:



- Each function  $g_{i,j}$  belongs to a Hölder class.
- We define the intrinsic smoothness of  $\mathcal{H}_d^l(\mathcal{P}, K)$  by
 
$$\gamma^* = \beta^*/t^*, \quad \text{where } (\beta^*, t^*) = \operatorname{argmin}_{(\beta, t) \in \mathcal{P}} \beta/t.$$
- $\gamma^*$  does not depend on the ambient input dimension, thereby *mitigating* the curse of dimensionality.

# Stationary $\beta$ -mixing

- Since RNNs naturally handle dependent sequences, we assume the data are **stationary** and  **$\beta$ -mixing** instead of **i.i.d.**

**Definition 2** ( $\beta$ -mixing (Bradley, 1983)). Let  $\{\mathbf{z}_t\}_{t=-\infty}^{\infty}$  be a sequence of random vectors. For any  $i, j \in \mathbb{Z} \cup \{-\infty, +\infty\}$ , define  $\sigma_i^j = \sigma(\mathbf{z}_i, \mathbf{z}_{i+1}, \dots, \mathbf{z}_j)$  as the  $\sigma$ -algebra generated by  $\mathbf{z}_k, i \leq k \leq j$ . For any  $a \in \mathbb{N}$ , the  $\beta$ -mixing coefficient of the stochastic process  $\{\mathbf{z}_t\}_{t=-\infty}^{\infty}$  is defined as

$$\beta(a) = \sup_{k \geq 1} \mathbb{E} \left[ \sup_{B \in \sigma_{-\infty}^k} \left| \mathbb{P}(A|B) - \mathbb{P}(A) \right| \right].$$

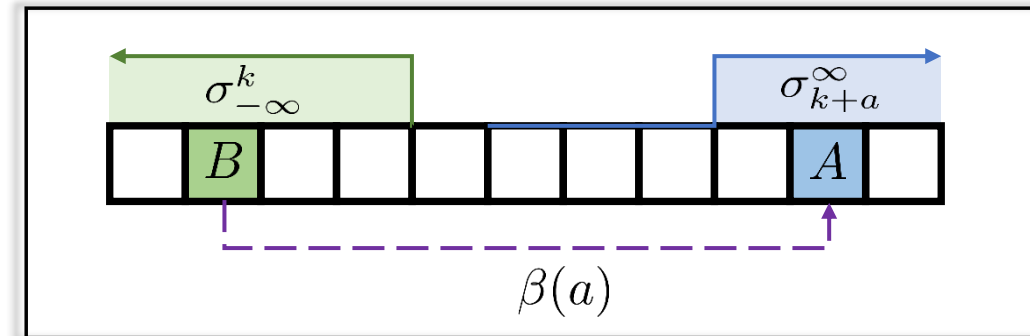
*quantifies maximum dependence between past and future observations*

- Algebraically  $\beta$ -mixing:

$$\beta(a) \leq \beta_0 / a^r, \quad \forall a.$$

- Exponentially  $\beta$ -mixing:

$$\beta(a) \leq \beta_0 \exp(-\beta_1 a^r), \quad \forall a.$$

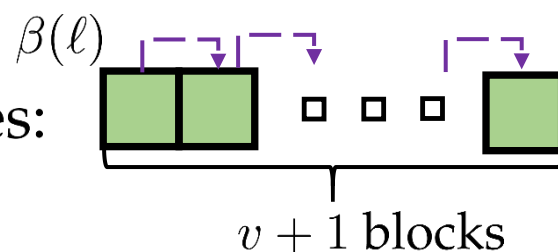
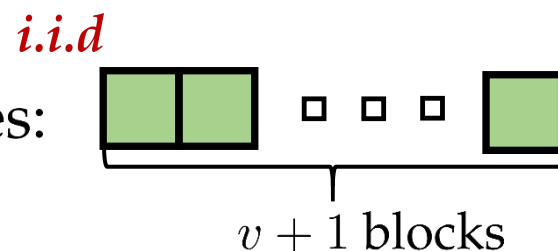




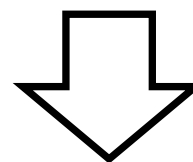
$m$  blocks:

Diagram illustrating the block structure of the proposed method. The sequence consists of:

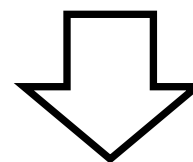
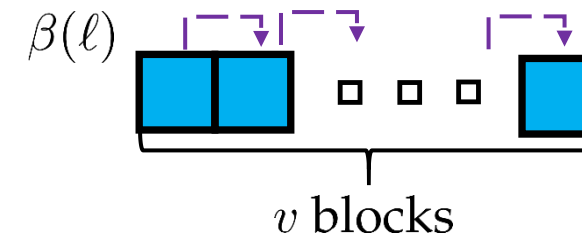
- $\ell$  blocks per group (light blue)
- $\ell$  blocks (green)
- $\ell$  blocks (grey)
- $k$  blocks *retained* (yellow)

 $\ell$  subsequences: $\ell$  subsequences:

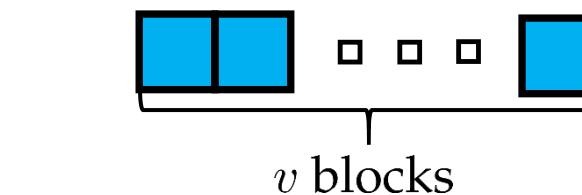
## Partition



### Lemma 5


$$\beta(\ell)$$


*i.i.d*



# Donut-set Decomposition

□ In analysis, an important quantity is

$$\mathbb{P} \left( \|\hat{f} - f_0\|_2 > \delta_\star \right),$$

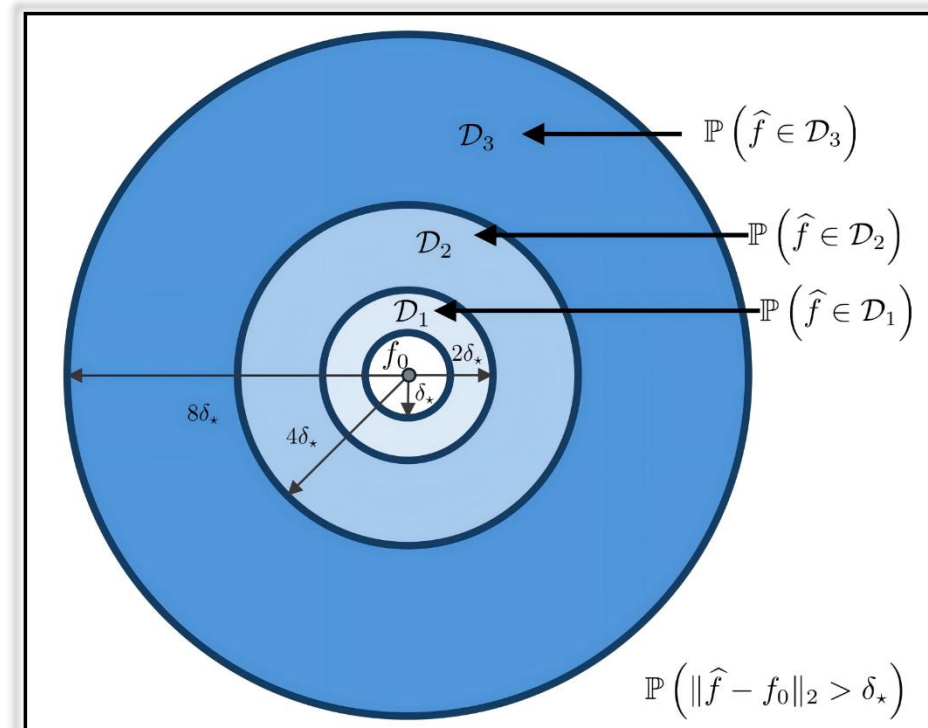
where  $\delta_\star = c(\delta_a + \delta_b + \delta_\beta)$ ,  $\delta_a$  is the approximation error,  $\delta_b$  is the stochastic error, and  $\delta_\beta$  is the dependence error.

**Donut-set Decomposition:**

$$\mathbb{P} \left( \|\hat{f} - f_0\|_2 > \delta_\star \right) \leq \sum_{i=1}^{\lfloor \log_2(2K/\delta_\star) \rfloor} \mathbb{P} \left( \hat{f} \in \mathcal{D}_i \right),$$

where

$$\mathcal{D}_i = \left\{ f \in \mathcal{F} : 2^{i-1}\delta_\star < \|f - f_0\|_2 \leq 2^i\delta_\star \right\}.$$



## □ Convergence Rate

**Theorem 1.** Let  $\mathcal{RNN}_{d_{\mathbf{x}},1}(W, L, K)$  be the hypothesis class  $\mathcal{F}$  and suppose regularity conditions, the hierarchical assumption, and the continuous probability measure assumption.

(i) Under the *exponentially* mixing assumption, Let  $W_0, L_0 \geq 3$  satisfy  $W_0 L_0 \asymp (n/(\log n)^{(6+1/r)})^{1/(4\gamma^*+2)}$ , and choose  $W = cW_0 \log W_0$  and  $L = cL_0 \log L_0$ . Then, the ERM estimator  $\hat{f}$  satisfies

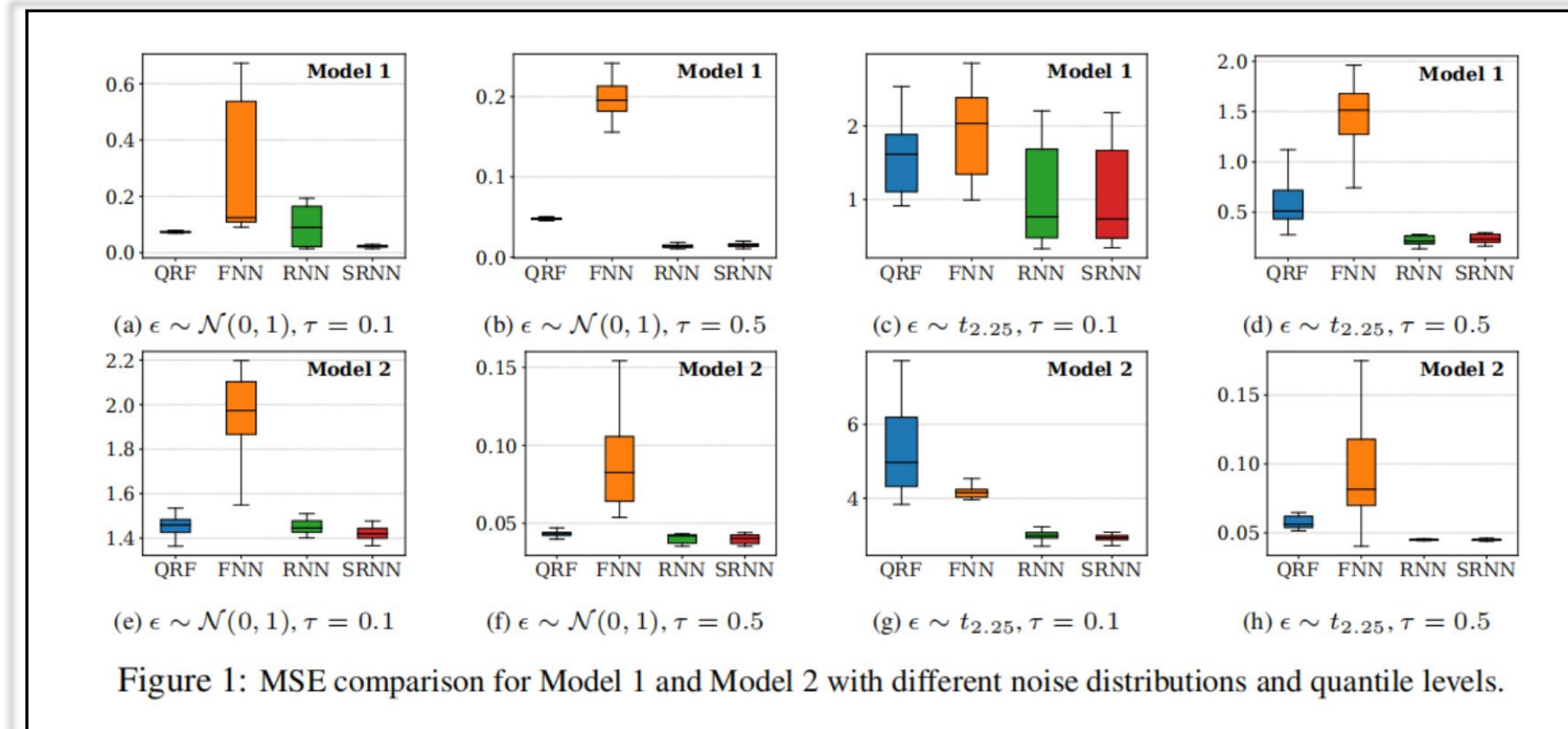
$$\|\hat{f} - f_0\|_2 = \mathcal{O}_p \left( n^{-\gamma^*/(2\gamma^*+1)} (\log n)^{(6+1/r)\gamma^*/(2\gamma^*+1)} \right). \quad \textbf{Minimax optimal!}$$

(ii) Under the *algebraically* mixing assumption, Let  $W_0, L_0 \geq 3$  satisfy  $W_0 L_0 \asymp (n^{(1-1/r)}/(\log n)^7)^{1/(4\gamma^*+2)}$ , and choose  $W = cW_0 \log W_0$  and  $L = cL_0 \log L_0$ . Then, the ERM estimator  $\hat{f}$  satisfies

$$\|\hat{f} - f_0\|_2 = \mathcal{O}_p \left( n^{-(1-1/r)\gamma^*/(2\gamma^*+1)} (\log n)^{(7\gamma^*/(2\gamma^*+1))} \right).$$

# Experiments

## □ Simulations



# Real data analysis

## □ Dow Jones Industrial Average (DJIA) analysis: *stationary*

Model	$\tau = 0.1$	$\tau = 0.25$	$\tau = 0.5$	$\tau = 0.75$	$\tau = 0.9$
QRF	0.456	0.698	0.817	0.616	0.365
FNN	0.538	0.735	0.810	0.662	0.404
RNN	0.410	<b>0.640</b>	0.760	0.562	0.306
SRNN	<b>0.406</b>	0.647	<b>0.759</b>	<b>0.561</b>	<b>0.305</b>

Table 1: Out-of-sample prediction errors at different quantiles for DJIA growth analysis.

## □ GDP analysis: *non-stationary*

Model	$\tau = 0.1$	$\tau = 0.25$	$\tau = 0.5$	$\tau = 0.75$	$\tau = 0.9$
QRF	0.849	1.225	1.410	1.246	0.911
FNN	0.867	1.180	1.773	2.505	2.657
RNN	<b>0.835</b>	<b>1.113</b>	1.349	<b>1.154</b>	0.904
SRNN	0.837	1.700	<b>1.211</b>	1.200	<b>0.898</b>

Table 2: Out-of-sample prediction errors at different quantiles for GDP growth analysis.

# Summary

---

- ❑ Problem: *Quantile regression* with *RNNs/SRNNs*
- ❑ Underlying function class: **Hierarchical Interaction Model**
- ❑ Data assumption: *Stationary and  $\beta$ -mixing*
- ❑ Two technique: *Blocking technique and donut-set decomposition*
- ❑ Experiments: Simulations, DJIA analysis, and GDP analysis

---

*Thanks!*