

# CONQUER THE QUANTILE: CONVOLUTION-SMOOTHED QUANTILE REGRESSION WITH NEURAL NETWORKS AND MINIMAX GUARANTEES

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Quantile regression provides a flexible approach to modeling heterogeneous effects and tail behaviors. This paper introduces the first quantile neural network estimator built upon the **convolution-smoothing quantile regression** (known as *conquer*) framework, which preserves both convexity and differentiability while retaining the robustness of the quantile loss. Extending the conquer estimator beyond linear models, we develop a non-parametric deep learning framework and establish sharp statistical guarantees. Specifically, we show that our estimator attains the minimax convergence rate over Besov spaces up to a logarithmic factor, matching the fundamental limits of nonparametric quantile estimation, and further derive general upper bounds for the estimation error in more general function classes. Empirical studies demonstrate that our method consistently outperforms existing quantile networks in both estimating accuracy and computational efficiency, underscoring the benefits of incorporating conquer into deep quantile learning.

## 1 INTRODUCTION

Quantile regression is a widely considered statistical tool for modeling heterogeneous effects and capturing the distributional structure of responses beyond the conditional mean. In many fields such as quantitative finance, survival analysis, and econometrics (Baur & Dimpfl, 2019; Horowitz, 1998; Chernozhukov & Hansen, 2005), quantile regression is used to understand the tail behaviors and provide robust outcomes faced with skewed or heavy-tailed data. Formally, given i.i.d. samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  from the random vector  $(X, Y)$  where  $X \in \mathbb{R}^d, Y \in \mathbb{R}$ , the conditional  $\tau$ -quantile function is defined as

$$f_{\tau}^*(\mathbf{x}_i) = F_{y_i|\mathbf{x}_i}^{-1}(\tau) = \inf\{y : P(y_i \leq y|\mathbf{x}_i) \geq \tau\}, \quad \tau \in (0, 1). \quad (1.1)$$

The standard approach estimates the conditional quantile function  $f_{\tau}^*$  by minimizing the empirical quantile loss

$$\hat{f}_{\tau} = \arg \min_f \sum_{i=1}^n \rho_{\tau}(y_i - f(\mathbf{x}_i)), \quad \rho_{\tau}(u) = \max\{\tau u, (\tau - 1)u\}. \quad (1.2)$$

We refer to Koenker et al. (2017) for an overview of quantile regression models. Although the estimator in equation 1.2 is theoretically well-founded, its loss function is non-differentiable. In models with a large parameter space, optimization becomes challenging so that efficient training is difficult. The neural network is a class of such models: even with a fixed input dimension  $d$ , the number of parameters can be very large. The flexibility of neural networks enables approximation of complex functions and, in theory, achieving minimax-optimal performance over Besov spaces (Suzuki, 2019). The success of neural networks depends heavily on generalization. Gradient-based optimizers such as stochastic gradient descent (SGD) often generalize well (Wu et al., 2022; Dziugaite & Roy, 2017), but sharp minima (Hochreiter & Schmidhuber, 1997)

047 and non-smooth loss surfaces (Huang et al., 2020; Foret et al., 2021) can degrade performance, especially  
048 with large-batch training (Keskar et al., 2016). This is because the high sensitivity to parameter changes  
049 caused by the sharp minima reduces the stability of optimization and model performance, and thus weakens  
050 generalization. Unfortunately, the quantile loss  $\rho_\tau(u)$  is non-differentiable at  $u = 0$  forming a sharp minima,  
051 and its gradient jumps abruptly from  $\tau$  to  $\tau - 1$  nearby.

052 To address the difficulties brought by the sharp minima and non-smooth loss functions, smoothing tech-  
053 niques have been developed, as seen in Berrada et al. (2018), which utilizes a smoothed loss function for  
054 classification tasks. For the quantile loss function, the conquer estimator (Fernandes et al., 2021) is particu-  
055 larly attractive, as it preserves both convexity and differentiability while retaining the robustness of quantile  
056 regression. Existing applications of conquer have been largely confined to linear models, and it remains un-  
057 clear whether its advantages carry over to more flexible function classes. In this paper, we demonstrate that  
058 integrating conquer smoothing with neural networks mitigate optimization difficulties, preserve statistical  
059 guarantees, while leveraging the expressive capacity of deep architectures.

## 061 1.1 RELATED WORK

062  
063 Our work builds upon smoothed quantile regression, neural networks, and minimax analysis in Besov spaces.  
064 For smoothed quantile regression, Horowitz (1998) introduced a kernel-based smoothing approach that al-  
065 leviates non-differentiability but sacrifices convexity, leading to challenging optimization problems. More  
066 recently, Fernandes et al. (2021) proposed the convolution smoothing method, known as the conquer estima-  
067 tor, which preserves both convexity and differentiability while gaining smoothness. The conquer framework  
068 has since been widely applied in quantile regression, and subsequent works have established its strong sta-  
069 tistical guarantees, including minimax optimality and asymptotic as well as nonasymptotic properties; see  
070 Kaplan & Sun (2017), Tan et al. (2022), and He et al. (2023). Apart from parametric models, Hu et al. (2025)  
071 proposed a local linear conquer estimator for time-varying coefficient models. However, the conquer esti-  
072 mator has primarily been developed in linear models, with limited exploration in nonlinear nonparametric  
073 settings. To the best of our knowledge, nonlinear extensions of the conquer estimator remain underdevel-  
074 oped, particularly in the context of modern neural networks. This motivates our focus, since Suzuki (2019)  
075 showed that deep neural networks can achieve minimax rate in Besov spaces where classical linear estima-  
076 tors such as kernel ridge regression, Nadaraya–Watson, or sieve methods are suboptimal.

077 Alongside these developments, a growing body of work has applied neural networks to quantile regression.  
078 Quantile networks have found applications in credit portfolio analysis, transportation, and survival analysis;  
079 see Feng et al. (2010), Rodrigues & Pereira (2020), and Pearce et al. (2022). On the theoretical side, statis-  
080 tical guarantees for quantile networks have been investigated, such as error bounds and minimax optimality;  
081 see Padilla et al. (2022), Shen et al. (2024), and Shen et al. (2025). Especially for minimax rates, Padilla  
082 et al. (2022) established that ReLU networks achieve near-minimax rates for quantile regression in Besov  
083 spaces. More recent work has also considered settings with covariate shift and noncrossing constraints (Feng  
084 et al., 2024; Shen et al., 2025), but these results remain restricted to Hölder classes. Together, this litera-  
085 ture suggests that bridging smoothed quantile regression with the expressive power of neural networks is a  
086 promising and unexplored direction.

## 087 1.2 CONTRIBUTIONS

088  
089 In this paper, we demonstrate that the convolution-type smoothed quantile regression technique (conquer)  
090 can be effectively integrated with neural networks to advance quantile regression. [Our study contributes](#)  
091 [to distributional learning through smoothed quantile objectives, which naturally relate to quantile-based](#)  
092 [modeling and theoretical guarantees—areas of growing interest within the deep learning community](#) (Suzuki,  
093 2019; Sun et al., 2022; Nishimura & Suzuki, 2024; Kelen et al., 2025). The main contributions are as follows:

- 094 (i) On the theoretical side, we first prove that the proposed conquer neural network estimator attains the  
 095 minimax convergence rate over Besov spaces, up to a logarithmic factor. In addition, we establish  
 096 general nonasymptotic error bounds that apply without assuming specific smoothness conditions  
 097 on the target function. These results demonstrate that the combination of conquer smoothing and  
 098 neural networks retains the desirable statistical guarantees of nonparametric quantile estimation  
 099 while leveraging the expressive capacity of deep learning models.
- 100 (ii) On the methodological and empirical side, we extend the conquer framework from linear non-  
 101 parametric models to neural networks, thereby enabling its application in highly nonlinear non-  
 102 parametric settings. Simulation studies are conducted to assess the performance of the proposed  
 103 method, and the results show consistent improvements over existing quantile networks in terms  
 104 of both estimation accuracy and computational efficiency. Together with the theoretical findings,  
 105 these empirical results highlight the effectiveness of applying conquer to modern neural network  
 106 architectures.

107  
 108 The rest of the article is organized as follows. We introduce the conquer framework with ReLU networks in  
 109 Section 2. In Section 3, we present the minimax rate for our estimation in Besov spaces and further develop  
 110 general upper bounds. Simulation studies are conducted in Section 4. All proofs and additional simulation  
 111 results are given in the Appendix.

## 113 2 PRELIMINARIES

114  
 115 Before starting our main result, we specify some notations regarding the ReLU neural networks with layer  
 116 parameter  $L$ , node parameter  $W$ , sparsity constraint  $S$ , and norm constraint  $B$ . In specific, we define the  
 117 class of sparse networks  $\mathcal{I}(L, W, S, B)$  with ReLU activation  $\sigma(x) = \max\{x, 0\}$  as

$$\begin{aligned}
 \mathcal{I}(L, W, S, B) := & \left\{ \left( A^{(L)}\sigma(\cdot) + b^{(L)} \right) \circ \dots \circ \left( A^{(1)}x + b^{(1)} \right) : A^{(1)} \in \mathbb{R}^{W \times d}, b^{(1)} \in \mathbb{R}^d, \right. \\
 & A^{(l)} \in \mathbb{R}^{1 \times W}, b^{(l)} \in \mathbb{R}, A^{(l)} \in \mathbb{R}^{W \times W}, b^{(l)} \in \mathbb{R}^W (1 < l < L), \\
 & \left. \sum_{l=1}^L \left( \|A^{(l)}\|_0 + \|b^{(l)}\|_0 \right) \leq S, \max_l \left( \|A^{(l)}\|_\infty \vee \|b^{(l)}\|_\infty \right) \leq B \right\}, \quad (2.1)
 \end{aligned}$$

125 where  $\circ$  denotes the composition of functions,  $\|A\|_0$  denotes the number of non-zero elements of the matrix  
 126  $A$ , and  $\|A\|_\infty$  denotes the maximum of the absolute values of the elements in matrix  $A$ . In our paper,  
 127 we also define  $\infty$ -norm for function  $f(\cdot)$  on the compact domain  $\mathcal{X}$ ,  $\|f\|_\infty = \sup_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x})$ . The sparse  
 128 network equation 2.1 has been widely investigated, by for example Suzuki (2019), Schmidt-Hieber (2020),  
 129 and Padilla et al. (2022). Based on  $\mathcal{I}(L, W, S, B)$ , our estimator  $\hat{f}_h$  is obtained from the ReLU networks by  
 130 minimizing the convolution-type smoothed quantile loss for  $\tau \in (0, 1)$ , i.e.,

$$\hat{f}_h := \arg \min_{f \in \mathcal{I}(L, W, S, B), \|f\|_\infty \leq F} \sum_{i=1}^n \ell_h(y_i - f(\mathbf{x}_i)), \quad \ell_h(u) := \int_{-\infty}^{\infty} \rho_\tau(v) K_h(v - u) dv, \quad (2.2)$$

135 where  $K_h(x) = K(x/h)/h$  with bandwidth  $h > 0$  and  $F > 0$  is a sufficiently large constant providing  
 136 technical convenience. The kernel function  $K(u)$  is required to be a bounded and nonnegative density  
 137 function such that  $\int uK(u)du = 0$ ,  $\int K(u)du = 1$ , and  $\int u^2K(u)du = \sigma_K^2$  where  $\sigma_K^2 > 0$  is a constant.  
 138  $K(u)$  can be chosen from commonly used kernel functions including: uniform kernel  $K(u) = (1/2)\mathbf{1}(|u| \leq$   
 139  $1)$ ; Gaussian kernel  $K(u) = (2\pi)^{-1/2}e^{-u^2/2}$ ; Epanechnikov kernel  $K(u) = (3/4)(1 - u^2)\mathbf{1}(|u| \leq 1)$ ,  
 140 etc.

For a given neural network  $f \in \mathcal{I}(L, W, S, B)$ ,  $\|f\|_\infty \leq F$ , we define the  $\|\cdot\|_\infty$ -projection of  $f_\tau^*$  onto  $\mathcal{I}(L, W, S, B)$  as

$$f_n := \arg \min_{f \in \mathcal{I}(L, W, S, B), \|f\|_\infty \leq F} \|f - f_\tau^*\|_\infty,$$

where  $f_\tau^*$  is assumed to be a function that belongs to Besov spaces defined below. Note that the architectural parameters  $(L, W, S, B)$  are usually chosen as functions of the sample size  $n$  (see Theorem 3.1), and therefore the projection  $f_n$  inherits this dependence through the network class  $\mathcal{I}(L, W, S, B)$ .

**Definition 2.1** (Besov space). For a function  $f \in L^p(\mathcal{X})$  and  $p \in (0, \infty]$ , denote the  $r$ -modulus of continuity as

$$w_{r,p}(f, t) = \sup_{\|u\|_2 \leq t} \|R_u^r(f)\|_p$$

where

$$R_u^r(f) = \begin{cases} \sum_{j=0}^r \frac{r!}{j!(r-j)!} (-1)^{r-j} f(x + uj) & , \text{ if } x \in \mathcal{X}, x + ru \in \mathcal{X} \\ 0 & , \text{ otherwise.} \end{cases}$$

For  $q \in (0, \infty]$  and  $\alpha > 0$ ,  $r = \lfloor \alpha \rfloor + 1$ , we define the Besov space  $B_{p,q}^\alpha(\mathcal{X})$  as

$$B_{p,q}^\alpha(\mathcal{X}) = \left\{ f \in L^p(\mathcal{X}) : \|f\|_{B_{p,q}^\alpha(\mathcal{X})} < \infty \right\}$$

where  $\|f\|_{B_{p,q}^\alpha(\mathcal{X})} = \|f\|_p + |f|_{B_{p,q}^\alpha(\mathcal{X})}$  and

$$|f|_{B_{p,q}^\alpha(\mathcal{X})} = \begin{cases} \left( \int_0^\infty (t^{-\alpha} w_{r,p}(f, t))^q t^{-1} dt \right)^{\frac{1}{q}} & \text{if } q < \infty, \\ \sup_{t>0} t^{-\alpha} w_{r,p}(f, t) & \text{if } q = \infty. \end{cases}$$

Throughout this paper, we write  $\lceil x \rceil$  for the smallest integer greater than or equal to  $x$ . For any two positive real sequences  $a_n$  and  $b_n$ , we write  $a_n \asymp b_n$  if there exist constants  $0 < c < C < \infty$  such that  $c \leq \liminf_{n \rightarrow \infty} a_n/b_n \leq \limsup_{n \rightarrow \infty} a_n/b_n \leq C$ . We write  $a_n \lesssim b_n$  ( $a_n \gtrsim b_n$ ) if there exists constant  $C > 0$  such that  $a_n \leq Cb_n$  ( $Ca_n \geq b_n$ ) for all  $n$ .

We introduce the performance metric in risk and empirical loss norms. For bounded functions  $f$  and  $g$ , we define

$$\Delta_n^2(f, g) := \frac{1}{n} \sum_{i=1}^n D^2(f(\mathbf{x}_i) - g(\mathbf{x}_i)),$$

where  $D^2(t) = \min\{|t|, t^2\}$ . Furthermore, we define  $\Delta^2(f, g) := \mathbb{E}(D^2(f(X) - g(X)))$ ,  $\Delta(f, g) := \sqrt{\Delta^2(f, g)}$ ,  $\|f - g\|_{\ell_2} := \sqrt{\mathbb{E}((f(X) - g(X))^2)}$ , and

$$\|f - g\|_n^2 := \frac{1}{n} \sum_{i=1}^n (f(\mathbf{x}_i) - g(\mathbf{x}_i))^2.$$

### 3 THEORY

In this section, we provide statistical guarantees for our conquer quantile ReLU neural networks. Firstly, we evaluate how well our methodology can estimate the functions in Besov spaces. Our main result in Theorem 3.1 shows that our conquer quantile regression with ReLU networks achieves the minimax error rate with only an additional logarithmic factor. Secondly, we develop a general risk bound on our estimator with respect to an arbitrary architecture of the neural networks, assuming the true quantile function  $f_\tau^*(\cdot)$  in a space more general than the Besov spaces, which imposes no smoothness conditions. This general risk bound accommodates a broad class of network architectures and modeling objectives, indicating that our theoretical development extends beyond the minimax analysis and offers an extensible foundation for future methodological advances.

### 3.1 MINIMAX RATE

In this subsection, we derive the convergence rate for our conquer estimator with ReLU networks when the quantile function belongs to Besov spaces. The Besov space is a very general function class which plays an important role in fields like statistical learning (Donoho & Johnstone, 1998; Giné & Nickl, 2021; Padilla et al., 2022) and approximation analysis (Temlyakov, 1993; Suzuki, 2019). As defined in Definition 2.1, Besov spaces unify and extend many classical smoothness spaces. In particular, the Sobolev space  $W^{\alpha,p}(\mathcal{X})$  coincides with the Besov space  $B_{p,p}^{\alpha}(\mathcal{X})$  and the Hölder space  $C^{\alpha}(\mathcal{X})$  corresponds to  $B_{\infty,\infty}^{\alpha}(\mathcal{X})$ . Thus, Besov space provides a more general framework that strictly contains both Sobolev and Hölder spaces that are commonly assumed for smoothness conditions in statistical guarantee analysis (Farrell et al., 2021; Montanelli, 2021; Schmidt-Hieber, 2020).

Based on the function space and network class specification, we impose the following assumptions on the data generation settings.

**Assumption 1.** We assume that  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  are i.i.d. samples from  $(X, Y)$ . Write  $F_{y_i|\mathbf{x}_i}$  is cumulative distribution function of  $y_i$  conditioning on  $\mathbf{x}_i$  for  $i = 1, \dots, n$ .

**Assumption 2.** There exists a constant  $\kappa > 0$  such that for  $\delta \in \mathbb{R}$  satisfying  $|\delta| \leq \kappa$  we have that, a.s.,

$$|F_{y_i|\mathbf{x}_i}(f_{\tau}^*(\mathbf{x}_i) + \delta) - F_{y_i|\mathbf{x}_i}(f_{\tau}^*(\mathbf{x}_i))| \geq \underline{p}|\delta| \quad (3.1)$$

for some constant  $\underline{p} > 0$ . Denote  $p_{y_i|\mathbf{x}_i}(\cdot)$  as the probability density function of  $y_i$  conditioning on  $\mathbf{x}_i$  uniformly for all  $i$ . We also require that  $p_{y_i|\mathbf{x}_i}(\cdot)$  is continuously differentiable and the derivative  $p'_{y_i|\mathbf{x}_i}(\cdot)$  satisfies almost surely that

$$|p'_{y_i|\mathbf{x}_i}(f_{\tau}^*(\mathbf{x}_i) + \delta) - p'_{y_i|\mathbf{x}_i}(f_{\tau}^*(\mathbf{x}_i))| \leq l_0|\delta|, \quad (3.2)$$

for some constant  $l_0 > 0$  uniformly for all  $i$ .

**Assumption 3.** Assume that  $X$  has a bounded probability density function  $g_X(\cdot) \leq c_2, c_2 > 0$  with support in  $[-H, H]^d$ .

**Assumption 4.** The quantile function satisfies  $f_{\tau}^* \in B_{p,q}^s([-H, H]^d)$ ,  $\|f_{\tau}^*\|_{\infty} \leq F$ , where for  $0 < p, q \leq \infty$ , and  $0 < s < \infty$  we have  $s \geq d/p$ . Furthermore, there exists  $m \in \mathbb{N}$  such that  $0 < s < \min\{m, m - 1 + 1/p\}$ . Here,  $B_{p,q}^s([-H, H]^d)$  is a Besov space in  $[-H, H]^d$  as in Definition 2.1.

**Remark 3.1.** Assumption 2 ensures the density around the target quantile is bounded away from zero. When the density around the target quantile is very small, the quantile becomes ill-conditioned, leading to inflated variance and unstable optimization.

1. Inflated variance. Standard quantile regression theory (e.g. Koenker & Xiao (2006)) shows that the asymptotic variance involves an inverse-density factor  $1/p_{Y|X}(q_{\tau}(x))$ . Hence, near-zero densities amplify the estimation error for any quantile estimator. In our neural network analysis, the constant appearing in equation A.1 of Lemma A.1 which bounds  $\Delta^2(f_n, f)$ , is proportional to the lower density bound  $1/\underline{p}$  where  $\underline{p}$  is imposed in Assumption 2. This reflects the same phenomenon: extremely small density around the target quantile leads to an unfavorable constant in the convergence rate.
2. Optimization instability. When the density is very small, the derivative of the quantile loss is nearly constant, resulting in weak curvature and slow convergence. Smoothing via the convolution kernel can mitigate this effect to some extent, but it cannot fully eliminate the ill-posedness caused by a vanishing density.

Condition 3.1 in Assumption 2 ensures local identifiability, which is a necessary condition that the quantile is estimatable. Similar conditions are widely imposed in quantile regression literature (Pollard, 1991; Belloni & Chernozhukov, 2011; Padilla et al., 2022). Meanwhile, condition 3.2 requires the conditional density  $p_{Y|X}(t)$  to be smooth in a neighborhood of the quantile. Such Lipschitz-type conditions are common in quantile regression (Koenker & Xiao, 2006; Zhou, 2010; He et al., 2023).

Assumptions 3 and 4 are mild and commonly assumed in both quantile regression (Wu & Zhou, 2017; He et al., 2023) and deep learning studies (Padilla et al., 2022; Suzuki, 2019). **Assumption 3 now only requires that  $g_X(\cdot)$  is bounded on  $[-H, H]^d$ .** The previous global lower bound  $c_1$ , as assumed in Padilla et al. (2022), is unnecessary for our results because regions with  $g_X(\mathbf{x}) = 0$  do not contribute to the  $\ell_2$  norm. This simplification does not affect identifiability or our main results, as the required local density condition is already guaranteed by Assumption 2. It is worth pointing out that Assumption 4 requires that the quantile function  $f_\tau^*$  belongs to a Besov space, which, as discussed, provides a more general smoothness framework than commonly assumed Sobolev or Hölder spaces, and is flexible enough to allow even discontinuous functions.

**Theorem 3.1.** *Suppose that Assumptions 1-4 hold. Given  $N \asymp n^{\frac{d}{2s+d}}$ , let  $\epsilon \asymp N^{-s/d} + \log^{-1} N$  with  $v = sp/d - 1$ , suppose  $h^2 \lesssim n^{-\frac{s}{2s+d}}$  and the class  $\mathcal{I}(L, W, S, B)$  satisfies*

$$L = 3 + 2 \left\lceil \log_2 \left( \frac{3^{\max\{d, m\}}}{\epsilon c_{d, m}} \right) + 5 \right\rceil \lceil \log_2 \max\{d, m\} \rceil, \quad W = W_0 N, \quad (3.3)$$

$$S = (L - 1)W_0^2 N + N, \quad B = O\left(N^{v^{-1} + d^{-1}}\right),$$

for a constant  $c_{d, m} > 0$  that depends on  $d$  and  $m$ . Then there exists a constant  $C > 0$  such that

$$\mathbb{P} \left( \max \left\{ \left\| \hat{f}_h - f_\tau^* \right\|_{\ell_2}^2, \left\| \hat{f}_h - f_\tau^* \right\|_n^2 \right\} \leq C(\log n)^3 \max\{\delta n^{-1}, n^{-\frac{2s}{2s+d}}\} \right) \geq 1 - e^{-\delta} \log n.$$

**Remark 3.2.** *The expression for  $L$  in equation 3.3 is not meant to prescribe a single fixed value. Since the quantity  $\epsilon$  may be chosen within the asymptotic range  $\epsilon \asymp N^{-s/d} + (\log N)^{-1}$  with  $N \asymp n^{d/(2s+d)}$ , the resulting depth  $L$  can vary accordingly. Thus, equation 3.3 should be interpreted as describing an admissible range of depths that ensures the minimax rate, rather than requiring  $L$  to take a specific exact value.*

**Remark 3.3.** *The constant  $C$  in the upper bound depends on the regularity assumptions on the conditional density. In particular,  $C$  is proportional to  $c_2$  in Assumption 3 and inversely proportional to the lower bound  $p$  in Assumption 2.*

Note that the rate  $n^{-\frac{2s}{2s+d}}$  is known to be minimax optimal for function estimation in Besov spaces (Kerkycharian & Picard, 1992; Donoho & Johnstone, 1998; Suzuki, 2019). **Theorem 3.1 explicitly yields the minimax rate  $O_p(n^{-\frac{2s}{2s+d}} \log^3 n)$  when  $\delta \asymp \log n$  in Besov space up to constants and logarithmic factors.** Our estimator attains this minimax rate up to a logarithmic factor  $\log^3 n$ , implying that no estimator can substantially improve upon the performance of our conquer quantile ReLU network beyond this logarithmic term. The minimax rate also clarifies the role of the bandwidth parameter  $h$ . Specifically, the condition  $h^2 \lesssim n^{-\frac{s}{2s+d}}$  ensures that the smoothing bias remains negligible, so that the estimator still targets the true quantile function. Conversely, when  $h$  is extremely small, the conquer loss function  $\ell_h(y_i - f(\mathbf{x}_i))$  effectively reduces to the original quantile loss  $\rho_\tau(y_i - f(\mathbf{x}_i))$ ; see Figure 1. Smaller  $h$  leads to a sharper minimum and a smaller smoothing region, which makes the gradients change dramatically in this region. The excessive sharpness induced by very small  $h$  will increase sensitivity to parameter changes, which reduces optimization stability and leads to poor generalization.

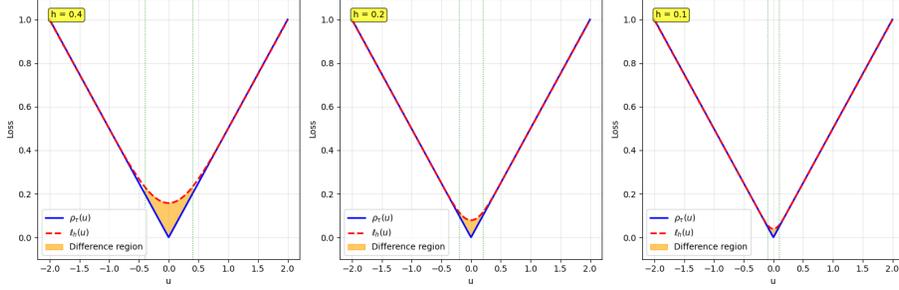


Figure 1: Relationship between loss functions  $\rho_\tau(u)$  (solid lines) and  $\ell_h(u)$  (dashed lines) for  $h = 0.4, 0.2, 0.1$ , with  $\tau = 0.5$ . Smaller  $h$  yields smaller difference region between  $\rho_\tau(u)$  and  $\ell_h(u)$ .

### 3.2 GENERAL UPPER BOUND

While the minimax analysis in the previous subsection relies on structural assumptions on the neural networks and target function, it is also of interest to study performance guarantees in a more general setting without such restrictions. To this end, we provide a general upper bound for our conquer neural network estimator without imposing constraints on the width, depth, magnitude of parameters, and sparsity of the network class, or on the smoothness assumptions (such as belonging to a Besov space) on the target quantile function  $f_\tau^*$ .

Consider a general neural network function with ReLU activation  $f \in \mathcal{F}(P, U, L)$  denoted as

$$\mathcal{F}(P, U, L) := \{f : \mathbb{R}^{p_0} \rightarrow \mathbb{R}^{p_{L+1}}, \quad \mathbf{x} \mapsto f(\mathbf{x}) = W_L \sigma_{\mathbf{v}_L} \circ W_{L-1} \sigma_{\mathbf{v}_{L-1}} \circ \cdots \circ W_1 \sigma_{\mathbf{v}_1} \circ W_0 \mathbf{x}\}, \quad (3.4)$$

where  $W_i$  is a  $p_{i+1} \times p_i$  weight matrix,  $\sigma_{\mathbf{v}_i}$  is a shifted activation function with shifting vector  $\mathbf{v}_i = (v_{i,1}, \dots, v_{i,p_i})^\top \in \mathbb{R}^{p_i}$ , i.e.,

$$\sigma_{\mathbf{v}_i} \begin{pmatrix} a_1 \\ \vdots \\ a_{p_i} \end{pmatrix} = \begin{pmatrix} \sigma(a_1 - v_{i,1}) \\ \vdots \\ \sigma(a_{p_i} - v_{i,p_i}) \end{pmatrix}.$$

Notice that the number of parameters is  $P = \sum_{l=0}^L (p_{l+1} p_l + p_{l+1})$  with  $p_0 = d, p_{L+1} = 1$ , the number of nodes is  $U = \sum_{l=1}^L p_l$ , and the number of layers is  $L$ . For  $f \in \mathcal{F}(P, U, L)$ , we redefine its conquer neural network estimator

$$\hat{f}_h := \arg \min_{f \in \mathcal{F}(P, U, L), \|f\|_\infty \leq F} \sum_{i=1}^n \ell_h(y_i - f(\mathbf{x}_i)), \quad (3.5)$$

and the approximation error is defined as

$$err_1 := \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \ell_h(y_i - \hat{f}_n(\mathbf{x}_i)) - \frac{1}{n} \sum_{i=1}^n \ell_h(y_i - f_\tau^*(\mathbf{x}_i)) \right] \quad (3.6)$$

where

$$f_n := \arg \min_{f \in \mathcal{F}(P, U, L), \|f\|_\infty \leq F} \mathbb{E} \left[ \sum_{i=1}^n \ell_h(y_i - f(\mathbf{x}_i)) \right]. \quad (3.7)$$

**Theorem 3.2.** *Suppose that Assumptions 1-2 hold and  $n \geq CLP \log(U)$  for a sufficiently large  $C > 0$ . Then with  $c_1 > 0$  a constant,  $\hat{f}_h$  defined in equation 3.5 satisfies*

$$\mathbb{E} \left[ \Delta_n^2 \left( \hat{f}_h, f_\tau^* \right) \right] \leq c_1 \left[ F \left( \frac{LP \log U \cdot \log n}{n} \right)^{1/2} + h^2 \mathbb{E} \|\hat{f}_h - f_\tau^*\|_{n,1} + err_1 \right], \quad (3.8)$$

where  $\|f - g\|_{n,1} = \frac{1}{n} \sum_{i=1}^n |f(\mathbf{x}_i) - g(\mathbf{x}_i)|$ . Furthermore, if  $h^2 = o\left(\sqrt{\mathbb{E}\|\hat{f}_h - f_\tau^*\|_n^2}\right)$ , it also holds that

$$\mathbb{E} \left[ \left\| \hat{f}_h - f_\tau^* \right\|_n^2 \right] \leq 2c_1 \max\{1, F\} \left[ F \left( \frac{LP \log U \cdot \log n}{n} \right)^{1/2} + err_1 \right]. \quad (3.9)$$

**Remark 3.4.** *Our general risk bound in Theorem 3.2 is quite flexible and does not rely on a very rigid architecture, but rather on general capacity measures (e.g., network class size, sparsity, norm bounds) and approximation properties. Therefore, in principle, it can accommodate more complex architectures, including residual (skip-connected) MLPs, as long as we can control the relevant complexity measures. However, to the best of our knowledge, there is no existing minimax-rate analysis for quantile regression specifically with residual-based (ResNet/skip-connection) networks. The most closely related work is the minimax-rate analysis for deep ReLU networks under quantile loss by Padilla et al. (2022), which considers plain feed-forward (non-skip) ReLU networks. Given the lack of minimax theory for residual architectures in quantile regression, we focused on the simpler network class to establish minimax optimality for our proposed conquer NN method, and leave for the possible extension to MLPs with skip connections as a promising future work.*

Theorem 3.2 drops out Assumptions 3-4 in Theorem 3.1 and yields a general error bound that depends on parameters  $P, U, L$ , the approximation error  $err_1$ , and the sample size  $n$ . As long as  $h^2 = o\left(\sqrt{\mathbb{E}\|\hat{f}_h - f_\tau^*\|_n^2}\right)$ , the risk bound in equation 3.9 remains essentially unaffected, implying that taking  $h$  sufficiently small does not deteriorate the estimator’s performance. However, from an optimization perspective, a very small  $h$  causes the smoothed loss  $\ell_h(\cdot)$  to behave almost identically to the original quantile loss  $\rho_\tau(\cdot)$ ; see Figure 1, which raises sharp minima concerns discussed in Section 1.

## 4 EMPIRICAL STUDY

In this section, we study the performance of our conquer framework with ReLU networks in several synthetic data scenarios, evaluated by mean squared error (MSE) and training time. We use the convolution-type smoothed quantile losses defined in equation 2.2 by three different kernel functions: Gaussian, uniform, and Epanechnikov kernels. We compare them against the baseline quantile ReLU network in Padilla et al. (2022) under different sample sizes and quantile levels. The experiment results show that our conquer networks perform better in MSE with less training time, especially for extreme quantile levels.

### 4.1 SCENARIO SETTINGS

Under the comparable number of parameters, we consider two networks with different shapes. One has 5 hidden layers of 70 nodes each, denoted by Model A, and the other one has 10 hidden layers of 50 nodes each, denoted by Model B. We consider both smooth and piecewise continuous quantile functions with heavy-tailed noises. Specifically, the data for the simulation are generated by the following mechanism,

$$y_i = g(\mathbf{x}_i) + \varepsilon_i, \quad i = 1, \dots, n,$$

where  $g(\cdot) : [0, 1]^d \rightarrow \mathbb{R}$ ,  $\{\mathbf{x}_i\}_{i=1}^n$  are independently sampled from uniform distribution on  $[0, 1]^d$  and  $d$  is the dimension of  $\mathbf{x}_i$ . We consider the following 3 scenarios of data generation:

**Scenario 1 (S1).** We set  $d = 2$ ,  $g(\mathbf{z}) = \cos(2\pi z_1^2) + \sin(\sqrt{z_1^2 + 2z_2} + 2)$ ,  $\mathbf{z} \in [0, 1]^2$ , and  $\varepsilon_i = \|\mathbf{x}_i - (1, 0)^\top\| t_i/2$ , where  $t_i \stackrel{\text{iid}}{\sim} t(2)$  for  $i = 1, \dots, n$  and  $t(2)$  is the t-distribution with 2 degrees of freedom.

**Scenario 2 (S2).** We set  $d = 5$ ,  $g(\mathbf{z}) = \sqrt{z_1 + 2z_2 + z_3 + 2z_4 + z_5}$ ,  $\mathbf{z} \in [0, 1]^5$ , and  $\varepsilon_i = \sqrt{\mathbf{x}_i^\top} \eta v_i$  for  $i = 1, \dots, n$ , where  $\eta = (1/2, 0, 1/2, 0, 1/2)^\top$  and  $v_i \stackrel{\text{iid}}{\sim} t(3)$ ,  $t(3)$  is the t-distribution with 3 degrees of freedom.

**Scenario 3 (S3).** We set  $d = 5$ ,  $g(\mathbf{z}) = g_2 \circ g_1(\mathbf{z})$ ,  $\mathbf{z} \in [0, 1]^5$  where  $g_1(\mathbf{z}) = (z_1 + 3z_2, \cos(2\pi(z_3 + z_4)), z_2 + \sqrt{z_3 + 2z_5})^\top$  and

$$g_2(\mathbf{z}) = \begin{cases} z_1 + \sqrt{z_2^2 + z_3} & \text{if } z_2 < 0, \\ \sqrt{z_1 + z_2} + 0.5z_3 & \text{otherwise,} \end{cases}$$

with  $\varepsilon_i \stackrel{\text{iid}}{\sim} \text{Laplace}(0, 2)$  for  $i = 1, \dots, n$ .

## 4.2 EXPERIMENT DETAILS

We evaluate the performance using the mean squared error (MSE) and training time. The training data set is denoted by  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$  and the test data set is by  $\{\tilde{\mathbf{x}}_i\}_{i=1}^n$ . For a fixed quantile  $\tau \in (0, 1)$ , we train the baseline network and conquer network from equation 2.2, and get the trained quantile estimate  $\hat{f}_h$ . By the data generating mechanism described above, we can calculate the true quantile  $f_\tau^*(\tilde{\mathbf{x}}_i) = g(\tilde{\mathbf{x}}_i) + F_{\varepsilon_i|\tilde{\mathbf{x}}_i}^{-1}(\tau)$ .

Then the MSE of quantile estimator is obtained by  $\sum_{i=1}^T (f_\tau^*(\tilde{\mathbf{x}}_i) - \hat{f}_h(\tilde{\mathbf{x}}_i))^2 / T$ , which is evaluated by using  $\{\tilde{\mathbf{x}}_i\}_{i=1}^T$  of size  $T = 10000$ . We set training sample sizes  $n \in \{1000, 5000, 10000\}$  and quantile levels  $\tau \in \{0.05, 0.25, 0.5, 0.75, 0.95\}$ . For each experiment setting, we run the experiment 50 times independently and get the MSE results 50 times. Table 1 shows the averaged MSE results over 50 trials under different experiment settings. **The bold fonts represent that the results of the conquer networks are better than those of the baseline models. In addition, to make the representation clearer, we box the baseline result if it outperforms our conquer networks.**

From Table 1, first, we can see a significant improvement in MSE as the sample size increases. A sample size of 1000 is inadequate for model training in our settings and is much smaller than the number of parameters of the networks. For sample size 10000, our conquer networks outperform the baseline model in most scenarios and quantile levels, regardless of which kernel is used, especially for low ( $\tau = 0.05$ ) and high ( $\tau = 0.95$ ) quantiles. It is reasonable because the original quantile loss is not differentiable at the origin point, which leads to biased subgradients for parameter update. Second, Scenario 3 has discontinuity points in the function  $g(\cdot)$ , while  $g(\cdot)$  is smooth in Scenario 1 and Scenario 2. The MSE results are consistent with Theorem 3.1 in the sense that MSE is smaller when the smoothness  $s$  increases. Third, Model B has more hidden layers than Model A, while the number of parameters is close. Meanwhile,  $d = 2$  in Scenario 1,  $d = 5$  in Scenario 2 and 3. The MSE results show that when  $d$  is small, shallow networks are more suitable. In contrast, for  $d = 5$ , deep networks perform better, which confirms the layer condition for the minimax rate in equation 3.3 of Theorem 3.1 to some extent.

We state the choice of bandwidth based on theoretical results in Section 3 and the experiments. As shown in Figure 1, the difference between smoothed loss  $\ell_h(u)$  and quantile loss  $\rho_\tau(u)$  grows larger with  $h$  increases. Furthermore, as  $h$  increases, the smoothed area becomes larger and the gradients become smaller, which affects the efficiency of SGD. On the other hand, too small  $h$  reduces to the original quantile loss and increases the sharpness of the minima, which tends to result in poor generalization. Therefore, a proper

Table 1: Mean squared error (MSE) performances for scenario 1-3, model A and B under different sample sizes, quantile levels, and smoothing kernels. The MSEs are averaged over 50 independent trials.

Method	$n=1000$					$n=5000$					$n=10000$					
	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
S1 Model A	Baseline	0.3820	0.0402	0.0278	0.0374	0.3784	0.0996	0.0120	0.0086	0.0108	0.0939	0.0618	0.0067	0.0047	0.0063	0.1366
	Gaussian	0.4354	<b>0.0383</b>	<b>0.0224</b>	0.0382	0.5035	0.1124	<b>0.0087</b>	<b>0.0055</b>	<b>0.0097</b>	0.1081	0.0678	<b>0.0064</b>	<b>0.0034</b>	<b>0.0055</b>	<b>0.0625</b>
	Uniform	0.4448	<b>0.0385</b>	<b>0.0224</b>	<b>0.0328</b>	<b>0.3721</b>	0.1079	<b>0.0087</b>	<b>0.0059</b>	<b>0.0104</b>	0.1312	0.0789	<b>0.0058</b>	<b>0.0036</b>	<b>0.0057</b>	<b>0.0543</b>
	Epanechnikov	0.6733	<b>0.0390</b>	<b>0.0222</b>	<b>0.0373</b>	0.5913	0.1251	<b>0.0093</b>	<b>0.0055</b>	<b>0.0095</b>	0.1169	0.0653	<b>0.0061</b>	<b>0.0037</b>	<b>0.0053</b>	<b>0.0665</b>
S1 Model B	Baseline	0.3842	0.0527	0.0319	0.0475	0.4222	0.1143	0.0149	0.0107	0.0151	0.1277	0.1691	0.0099	0.0066	0.0082	0.3882
	Gaussian	0.4158	0.0665	0.0383	0.0633	0.5202	0.1172	<b>0.0144</b>	<b>0.0097</b>	<b>0.0142</b>	<b>0.1066</b>	<b>0.1062</b>	<b>0.0095</b>	<b>0.0055</b>	0.0086	<b>0.0726</b>
	Uniform	0.5652	0.0612	0.0378	0.0614	0.5685	<b>0.1033</b>	<b>0.0145</b>	<b>0.0091</b>	0.0154	<b>0.1252</b>	<b>0.0974</b>	<b>0.0093</b>	<b>0.0055</b>	<b>0.0080</b>	<b>0.0736</b>
	Epanechnikov	0.4275	0.0582	0.0383	0.0626	0.6050	<b>0.1115</b>	<b>0.0146</b>	<b>0.0094</b>	<b>0.0145</b>	<b>0.1162</b>	<b>0.1294</b>	<b>0.0089</b>	<b>0.0057</b>	0.0092	<b>0.0663</b>
S2 Model A	Baseline	0.8292	0.0868	0.0619	0.0839	0.7874	0.2752	0.0275	0.0222	0.0308	0.2747	0.1704	0.0205	0.0145	0.0223	0.1670
	Gaussian	<b>0.7994</b>	<b>0.0711</b>	<b>0.0587</b>	<b>0.0778</b>	1.1466	<b>0.2202</b>	0.0276	<b>0.0169</b>	<b>0.0273</b>	<b>0.2598</b>	<b>0.1316</b>	<b>0.0176</b>	<b>0.0128</b>	<b>0.0193</b>	<b>0.1537</b>
	Uniform	0.8892	<b>0.0721</b>	<b>0.0471</b>	0.0964	0.8566	<b>0.2308</b>	0.0286	<b>0.0181</b>	<b>0.0275</b>	<b>0.2586</b>	<b>0.1457</b>	<b>0.0180</b>	<b>0.0128</b>	<b>0.0191</b>	<b>0.1467</b>
	Epanechnikov	0.9192	<b>0.0787</b>	<b>0.0522</b>	0.1048	0.9074	<b>0.2415</b>	0.0284	<b>0.0177</b>	<b>0.0265</b>	<b>0.2492</b>	<b>0.1430</b>	<b>0.0162</b>	<b>0.0126</b>	<b>0.0191</b>	<b>0.1388</b>
S2 Model B	Baseline	0.4930	0.0583	0.0493	0.0840	0.5898	0.1732	0.0257	0.0178	0.0300	0.1966	0.1323	0.0202	0.0129	0.0220	0.1358
	Gaussian	0.7367	0.0639	0.0500	<b>0.0759</b>	0.6273	0.1800	<b>0.0246</b>	<b>0.0155</b>	<b>0.0245</b>	0.2157	<b>0.1085</b>	<b>0.0160</b>	<b>0.0119</b>	<b>0.0178</b>	<b>0.1144</b>
	Uniform	<b>0.4515</b>	0.0717	0.0503	0.0935	0.6034	<b>0.1706</b>	<b>0.0216</b>	<b>0.0176</b>	<b>0.0241</b>	0.2239	<b>0.1172</b>	<b>0.0151</b>	<b>0.0118</b>	<b>0.0175</b>	<b>0.1342</b>
	Epanechnikov	0.5800	0.0723	0.0546	<b>0.0819</b>	0.7484	0.2151	0.0263	<b>0.0159</b>	<b>0.0289</b>	<b>0.1802</b>	<b>0.1038</b>	<b>0.0146</b>	<b>0.0107</b>	<b>0.0173</b>	<b>0.1302</b>
S3 Model A	Baseline	3.0766	0.7564	0.5350	0.7407	2.8969	1.2870	0.3854	0.2146	0.3387	1.3495	0.9232	0.2420	0.1459	0.2413	0.9960
	Gaussian	<b>2.8841</b>	0.7776	<b>0.4788</b>	<b>0.7056</b>	3.4222	1.3139	<b>0.3379</b>	<b>0.1993</b>	<b>0.3269</b>	<b>1.1897</b>	<b>0.8395</b>	<b>0.2040</b>	<b>0.1295</b>	<b>0.2145</b>	<b>0.7477</b>
	Uniform	3.3730	0.7735	<b>0.4876</b>	0.7832	3.3576	<b>1.1577</b>	<b>0.3449</b>	<b>0.1912</b>	<b>0.3082</b>	<b>1.2392</b>	<b>0.8790</b>	<b>0.2064</b>	<b>0.1300</b>	<b>0.2146</b>	<b>0.7824</b>
	Epanechnikov	3.4479	<b>0.7308</b>	<b>0.4679</b>	0.8153	3.3156	<b>1.1229</b>	<b>0.3615</b>	<b>0.1980</b>	<b>0.3112</b>	<b>1.3312</b>	<b>0.8230</b>	<b>0.2129</b>	<b>0.1349</b>	<b>0.2206</b>	<b>0.7998</b>
S3 Model B	Baseline	2.8166	0.7558	0.5175	0.7665	2.2839	1.0061	0.3596	0.2196	0.3551	1.0990	0.7786	0.2306	0.1391	0.2380	0.7249
	Gaussian	<b>2.3193</b>	0.8405	<b>0.5142</b>	<b>0.7543</b>	2.6520	1.0602	0.4263	0.2304	<b>0.3525</b>	<b>1.0681</b>	0.8256	0.2518	0.1416	0.2444	0.7536
	Uniform	<b>2.1946</b>	0.8316	0.5193	<b>0.7352</b>	2.5765	1.2056	0.4038	0.2320	<b>0.3373</b>	<b>1.0528</b>	<b>0.7167</b>	0.2409	<b>0.1372</b>	0.2493	<b>0.6920</b>
	Epanechnikov	2.9702	0.9008	<b>0.4892</b>	0.8702	<b>2.1331</b>	<b>1.0053</b>	0.3750	0.2297	<b>0.3473</b>	<b>1.0895</b>	<b>0.7502</b>	0.2454	<b>0.1390</b>	<b>0.2315</b>	<b>0.6683</b>

choice of bandwidth is necessary. By Theorem 3.1, we accept a bigger  $h$  when the sample size  $n$  is small and a smaller  $h$  when  $n$  is big. For example,  $h = 0.01/0.005/0.001$  for  $n = 1000/5000/10000$  in Scenario 2, Model A. We also find that when  $n = 10000$ , the performance remains outstanding with a wide range of bandwidth  $h$ , see Appendix B.

We also record the training time for each setting and for 50 trials. Due to the page limit, we present the result for  $\tau = 0.05$  in Figure 2. The rest results are shown in Figures 5-8 in Appendix B. We can conclude that our conquer networks are generally faster than the baseline, with only 80% of the training time consumed, and are stable enough.

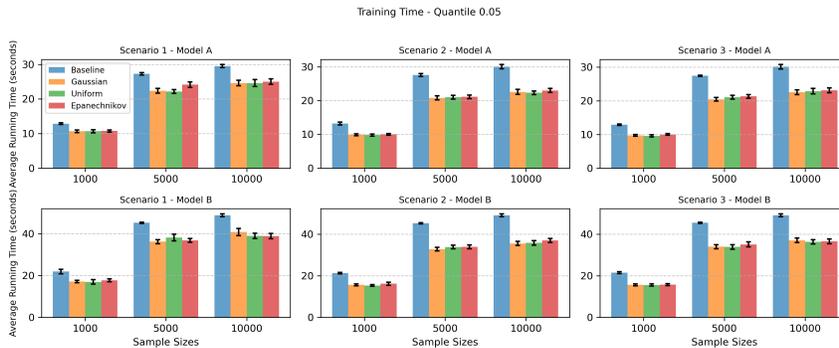


Figure 2: Bar chart with error bars with average training time over 50 trials under quantile level  $\tau = 0.05$ . The error bars represent 95% confidence intervals of training time for each setting.

470 ETHICS STATEMENT  
471

472 Our research does not pose any potential ethical issues. This work does not involve human subjects, personal  
473 data, or sensitive demographic information.  
474

475 REPRODUCIBILITY STATEMENT  
476

477 We have made efforts to ensure reproducibility of our results. All theoretical results are presented with  
478 complete assumptions in Section 3 and detailed proofs in Appendix A. The experimental setup, includ-  
479 ing simulation designs, model configurations, and hyperparameter choices, is described in Section 4 and  
480 Appendix B. Random seeds are set to ensure reproducibility of the experiment. Upon publication, we  
481 will release the full implementation, including training scripts and simulation setups, in a GitHub repos-  
482 itory. An anonymous version of the GitHub repository is available during the open review, see <https://anonymous.4open.science/r/conquernn-F625/>.  
483  
484

485 REFERENCES  
486

- 487 Peter L Bartlett, Olivier Bousquet, and Shahar Mendelson. Local rademacher complexities. *The Annals of*  
488 *Statistics*, 33(4):1497–1537, 2005.
- 489 Dirk G Baur and Thomas Dimpfl. A quantile regression approach to estimate the variance of financial  
490 returns. *Journal of Financial Econometrics*, 17(4):616–644, 2019.
- 491 Alexandre Belloni and Victor Chernozhukov.  $\ell_1$ -penalized quantile regression in high-dimensional sparse  
492 models. *The Annals of Statistics*, 39(1):82 – 130, 2011.
- 493 Leonard Berrada, Andrew Zisserman, and M Pawan Kumar. Smooth loss functions for deep top-k classifi-  
494 cation. *arXiv preprint arXiv:1802.07595*, 2018.
- 495 Jonathan Berrisch and Florian Ziel. Crps learning. *Journal of Econometrics*, 237(2):105221, 2023.
- 496 Victor Chernozhukov and Christian Hansen. An iv model of quantile treatment effects. *Econometrica*, 73  
497 (1):245–261, 2005.
- 498 Will Dabney, Mark Rowland, Marc Bellemare, and Rémi Munos. Distributional reinforcement learning with  
499 quantile regression. In *Proceedings of the AAAI conference on artificial intelligence*, volume 32, 2018.
- 500 David L Donoho and Iain M Johnstone. Minimax estimation via wavelet shrinkage. *The annals of Statistics*,  
501 26(3):879–921, 1998.
- 502 Gintare Karolina Dziugaite and Daniel M Roy. Computing nonvacuous generalization bounds for  
503 deep (stochastic) neural networks with many more parameters than training data. *arXiv preprint*  
504 *arXiv:1703.11008*, 2017.
- 505 Max H Farrell, Tengyuan Liang, and Sanjog Misra. Deep neural networks for estimation and inference.  
506 *Econometrica*, 89(1):181–213, 2021.
- 507 Xingdong Feng, Xin He, Yuling Jiao, Lican Kang, and Caixing Wang. Deep nonparametric quantile regres-  
508 sion under covariate shift. *Journal of Machine Learning Research*, 25(385):1–50, 2024.
- 509 Yijia Feng, Runze Li, Agus Sudjianto, and Yiyun Zhang. Robust neural network with applications to credit  
510 portfolio data analysis. *Statistics and its Interface*, 3(4):437, 2010.  
511  
512  
513  
514  
515  
516



- 564 Oscar Hernan Madrid Padilla and Sabyasachi Chatterjee. Risk bounds for quantile trend filtering.  
565 *Biometrika*, 109(3):751–768, 09 2021.  
566
- 567 Hadrien Montanelli. Deep relu networks overcome the curse of dimensionality for generalized bandlimited  
568 functions. *Journal of Computational Mathematics*, 39(6), 2021.
- 569 Yuto Nishimura and Taiji Suzuki. Minimax optimality of convolutional neural networks for infinite di-  
570 mensional input-output problems and separation from kernel methods. In *The Twelfth International*  
571 *Conference on Learning Representations*, 2024. URL [https://openreview.net/forum?id=](https://openreview.net/forum?id=EW8ZErxRzkJ)  
572 [EW8ZErxRzkJ](https://openreview.net/forum?id=EW8ZErxRzkJ).  
573
- 574 Oscar Hernan Madrid Padilla, Wesley Tansey, and Yanzhen Chen. Quantile regression with relu networks:  
575 Estimators and minimax rates. *Journal of Machine Learning Research*, 23(247):1–42, 2022.  
576
- 577 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen,  
578 Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary  
579 DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and  
580 Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in*  
581 *Neural Information Processing Systems*, volume 32, 2019.
- 582 Tim Pearce, Jong-Hyeon Jeong, Jun Zhu, et al. Censored quantile regression neural networks for  
583 distribution-free survival analysis. *Advances in neural information processing systems*, 35:7450–7461,  
584 2022.
- 585 David Pollard. Asymptotics for least absolute deviation regression estimators. *Econometric Theory*, 7(2):  
586 186–199, 1991.  
587
- 588 Filipe Rodrigues and Francisco C Pereira. Beyond expectation: Deep joint mean and quantile regression  
589 for spatiotemporal problems. *IEEE transactions on neural networks and learning systems*, 31(12):5377–  
590 5389, 2020.
- 591 Johannes Schmidt-Hieber. Nonparametric regression using deep neural networks with relu activation func-  
592 tion. *The Annals of Statistics*, 48(4):1875, 2020.  
593
- 594 Bodhisattva Sen. A gentle introduction to empirical process theory and applications, 2018.  
595
- 596 Guohao Shen, Yuling Jiao, Yuanyuan Lin, Joel L Horowitz, and Jian Huang. Nonparametric estimation of  
597 non-crossing quantile regression process with deep relu neural networks. *Journal of Machine Learning*  
598 *Research*, 25(88):1–75, 2024.
- 599 Guohao Shen, Runpeng Dai, Guojun Wu, Shikai Luo, Chengchun Shi, and Hongtu Zhu. Deep distributional  
600 learning with non-crossing quantile network. *arXiv preprint arXiv:2504.08215*, 2025.  
601
- 602 Jesse Sun, Dihong Jiang, and Yaoliang Yu. Conditional generative quantile networks via optimal transport.  
603 In *ICLR Workshop on Deep Generative Models for Highly Structured Data*, 2022. URL [https://](https://openreview.net/forum?id=BBxeo2Vuvbq)  
604 [openreview.net/forum?id=BBxeo2Vuvbq](https://openreview.net/forum?id=BBxeo2Vuvbq).
- 605 Taiji Suzuki. Adaptivity of deep relu network for learning in besov and mixed smooth besov spaces: optimal  
606 rate and curse of dimensionality. In *International Conference on Learning Representations*, 2019.  
607
- 608 Kean Ming Tan, Lan Wang, and Wen-Xin Zhou. High-dimensional quantile regression: Convolution  
609 smoothing and concave regularization. *Journal of the Royal Statistical Society Series B: Statistical*  
610 *Methodology*, 84(1):205–233, 2022.

611 V.N. Temlyakov. *Approximation of Periodic Functions*. Computational mathematics and analysis series.  
612 Nova Science Publishers, 1993.  
613

614 Lei Wu, Mingze Wang, and Weijie Su. The alignment property of sgd noise and how it helps select flat  
615 minima: A stability analysis. *Advances in Neural Information Processing Systems*, 35:4680–4693, 2022.  
616

617 Weichi Wu and Zhou Zhou. Nonparametric inference for time-varying coefficient quantile regression. *Jour-  
618 nal of Business & Economic Statistics*, 35(1):98–109, 2017.

619 Xuzhi Yang and Tengyao Wang. Multiple-output composite quantile regression through an optimal transport  
620 lens. In *Proceedings of Thirty Seventh Conference on Learning Theory*, volume 247 of *Proceedings of  
621 Machine Learning Research*, pp. 5076–5122. PMLR, 30 Jun–03 Jul 2024.

622 Rui Zhang, Christian Walder, Edwin V Bonilla, Marian-Andrei Rizoïu, and Lexing Xie. Quantile propaga-  
623 tion for wasserstein-approximate gaussian processes. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Bal-  
624 can, and H. Lin (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 21566–21578.  
625 Curran Associates, Inc., 2020. URL [https://proceedings.neurips.cc/paper\\_files/  
626 paper/2020/file/f5e62af885293cf4d511ceef31e61c80-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/f5e62af885293cf4d511ceef31e61c80-Paper.pdf).

627 Zhou Zhou. Nonparametric inference of quantile curves for nonstationary time series. *The Annals of Statis-  
628 tics*, 38(4):2187 – 2217, 2010.  
629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657

658 **A PROOFS**

659 We first state several definitions to develop our Empirical process theorems and auxiliary lemmas. Define  
660 the empirical loss function as

$$661 \hat{M}_n(f) = \sum_{i=1}^n \hat{M}_{n,i}(f), \quad \hat{M}_{n,i}(f) = \frac{1}{n}(\ell_h(y_i - f(\mathbf{x}_i)) - \ell_h(y_i - f_n(\mathbf{x}_i))),$$

662 and we set

$$663 M_n(f) = \mathbb{E} \left\{ \frac{1}{n} \sum_{i=1}^n [\ell_h(y_i - f(\mathbf{x}_i)) - \ell_h(y_i - f_n(\mathbf{x}_i))] \right\}.$$

664 For  $\epsilon > 0$  and a metric  $\text{dist}(\cdot, \cdot)$  on the class of functions  $\mathcal{F}$ , we define the covering number  $\mathcal{N}(\epsilon, \mathcal{F}, \text{dist}(\cdot, \cdot))$   
665 as the minimum number of balls of the form  $\{g : \text{dist}(g, f) \leq \epsilon\}$ , with  $f \in \mathcal{F}$ , needed to cover  $\mathcal{F}$  (see Defini-  
666 tion 2.2 in Sen (2018) for details). For simplicity, we consider  $[0, 1]^d$  instead of  $[-H, H]^d$  in Assumptions  
667 3-4.

673 **A.1 AUXILIARY LEMMAS**

674 **Lemma A.1.** *Suppose that  $\|f_n - f_\tau^*\|_\infty \leq c$  for a small enough constant  $c$  we have*

$$675 \Delta^2(f_n, f) \leq C \left[ \mathbb{E}(\ell_h(Y - f(X)) - \ell_h(Y - f_n(X))) + (\|f_n - f_\tau^*\|_\infty + h^2) \Delta(f, f_n) \sqrt{F} \right]. \quad (\text{A.1})$$

676 and

$$677 \|f - f_n\|_{\ell_2}^2 \leq C \max\{1, F\} \left[ \mathbb{E}(\ell_h(Y - f(X)) - \ell_h(Y - f_n(X))) + (\|f_n - f_\tau^*\|_\infty + h^2) \|f - f_n\|_{\ell_2} \sqrt{F} \right],$$

678 (A.2)

679 for any  $f \in \tilde{\mathcal{L}}(L, W, S, B)$  and for some constant  $C > 0$ .

680 *Proof.* By Knight identity (Knight, 1998),

$$681 \rho_\tau(Y - f(X) + s) - \rho_\tau(Y - f_n(X) + s)$$

$$682 = -(f(X) - f_n(X)) (\tau - \mathbf{1}\{Y + s \leq f_n(X)\}) + \int_0^{f(X) - f_n(X)} [\mathbf{1}\{Y + s \leq f_n(X) + z\} - \mathbf{1}\{Y + s \leq f_n(X)\}] dz,$$

$$683 = -(f(X) - f_n(X)) (\tau - \mathbf{1}\{Y + s \leq f_\tau^*(X)\}) - (f(X) - f_n(X)) (\mathbf{1}\{Y + s \leq f_\tau^*(X)\} - \mathbf{1}\{Y + s \leq f_n(X)\})$$

$$684 + \int_0^{f(X) - f_n(X)} [\mathbf{1}\{Y + s \leq f_n(X) + z\} - \mathbf{1}\{Y + s \leq f_n(X)\}] dz. \quad (\text{A.3})$$

685 By equation 3.1 in Assumption 2 and mean value expansion, applying Fubini's theorem and the fact  
686  $\int s K_h(s) ds = 0$ ,  $\int s^2 K_h(s) ds = \sigma_K^2 h^2$ , we have for some constant  $\underline{c}, c_\tau > 0$ ,

$$687 \mathbb{E} \left[ \int K_h(s) \int_0^{f(X) - f_n(X)} \mathbf{1}\{Y + s \leq f_n(X) + z\} - \mathbf{1}\{Y + s \leq f_n(X)\} dz ds \middle| X \right],$$

$$688 = \mathbb{E} \left[ \int K_h(s) \int_0^{f(X) - f_n(X)} F_{Y|X}(f_n(X) + z - s) - F_{Y|X}(f_n(X) - s) dz ds \right],$$

$$689 \geq \mathbb{E} \left[ \underline{p} \int_0^{f(X) - f_n(X)} \min\{z, \kappa\} dz \right] - \underline{c} h^2 \mathbb{E}|f(X) - f_n(X)|,$$

$$690 \geq c_\tau (\mathbb{E}[D^2(f(X) - f_n(X))] - h^2 \mathbb{E}|f(X) - f_n(X)|), \quad (\text{A.4})$$

Combining equation A.25, equation A.3, and equation A.4, by Fubini's theorem, we have

$$\begin{aligned}
& \mathbb{E} \{ \ell_h(Y - f(X)) - \ell_h(Y - f_n(X)) \} \\
&= \int K_h(s) \mathbb{E} \{ \rho_\tau(Y + s - f(X)) - \rho_\tau(Y + s - f_n(X)) \} ds, \\
&\geq -C \mathbb{E} \{ |f(X) - f_n(X)| \cdot (|f_\tau^*(X) - f_n(X)| + h^2) \} + c_\tau (\mathbb{E} [D^2(f(X) - f_n(X))]), \\
&\geq -C(\|f_\tau^* - f_n\|_\infty + h^2) \sqrt{F} \Delta^2(f, f_n) + c_\tau \Delta^2(f, f_n), \tag{A.5}
\end{aligned}$$

which yields equation A.1. Furthermore, note that  $\frac{\|f - f_n\|_{\ell_2}^2}{\max\{F, 1\}} \leq \Delta^2(f, f_n) \leq \|f - f_n\|_{\ell_2}^2$ , then equation A.2 holds.  $\square$

**Lemma A.2.** Suppose that  $f_n \in \tilde{\mathcal{I}}(L, W, S, B)$  and  $\|f_n - f_\tau^*\|_\infty \leq c$  for a sufficiently small constant  $c > 0$ . The estimator  $\hat{f}_h$  defined in equation 2.2 satisfies for some constant  $C > 0$ ,

$$\Delta^2(\hat{f}_h, f_n) \leq C \left[ M_n(\hat{f}_h) - \hat{M}_n(\hat{f}_h) + (\|f_n - f_\tau^*\|_\infty + h^2) \Delta(\hat{f}_h, f_n) \sqrt{F} \right]. \tag{A.6}$$

Furthermore,

$$\|\hat{f}_h - f_n\|_{\ell_2}^2 \leq C \max\{1, F\} \left[ M_n(\hat{f}_h) - \hat{M}_n(\hat{f}_h) + (\|f_n - f_\tau^*\|_\infty + h^2) \|\hat{f}_h - f_n\|_{\ell_2} \sqrt{F} \right]. \tag{A.7}$$

*Proof.* Since  $\hat{f}_h$  in equation 2.2 satisfies  $\hat{f}_h \in \tilde{\mathcal{I}}(L, W, S, B)$ , then by Lemma A.1,

$$\begin{aligned}
\Delta^2(\hat{f}_h, f_n) &\leq C \left[ \mathbb{E} \left( \ell_h(Y - \hat{f}_h(X)) - \ell_h(Y - f_n(X)) \right) + \|f_n - f_\tau^*\|_\infty \Delta(\hat{f}_h, f_n) \sqrt{F} + h^2 \right], \\
&\leq C \left[ M_n(\hat{f}_h) - \hat{M}_n(\hat{f}_h) + \|f_n - f_\tau^*\|_\infty \Delta(\hat{f}_h, f_n) \sqrt{F} + h^2 \right], \tag{A.8}
\end{aligned}$$

where the last inequality is obtained by the fact  $\hat{M}_n(\hat{f}_h) \leq 0$ . Similarly, equation A.7 can also hold by Lemma A.1 and the optimality of  $\hat{f}_h$ .  $\square$

**Lemma A.3.** Suppose that

$$3\mathbb{E} \left( \sup_{f \in \tilde{\mathcal{I}}(L, W, S, B), \|f - f_n\|_{\ell_2}^2 \leq r^2} \frac{1}{n} \sum_{i=1}^n \xi_i (f(\mathbf{x}_i) - f_n(\mathbf{x}_i))^2 \right) \leq r^2, \tag{A.9}$$

for  $\{\xi_i\}_{i=1}^n$  Rademacher variables independent of  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , and

$$2F \sqrt{\frac{7\gamma}{3n}} \leq r \tag{A.10}$$

Then with probability at least  $1 - e^{-\gamma}$ ,  $\|f - f_n\|_{\ell_2}^2 \leq r^2$  with  $f \in \tilde{\mathcal{I}}(L, W, S, B)$  implies

$$\|f - f_n\|_n^2 \leq (2r)^2.$$

*Proof.* Using the fact  $(a - b)^2 \leq 2(a^2 + b^2)$ , we have for  $f \in \tilde{\mathcal{I}}(L, W, S, B)$

$$\begin{aligned}
\text{Var} [(f(X) - f_n(X))^2] &\leq \mathbb{E}(f(X) - f_n(X))^4, \\
&\leq 2(F^2 + \|f_n\|_\infty^2) \mathbb{E}(f(X) - f_n(X))^2, \\
&\leq 4F^2 \|f - f_n\|_{\ell_2}^2. \tag{A.11}
\end{aligned}$$

Note that  $0 \leq (f(\mathbf{x}) - f_n(\mathbf{x}))^2 \leq 4F^2$ , then by Theorem 2.1 in Bartlett et al. (2005), for every  $\gamma > 0$ , with probability at least  $1 - e^{-\gamma}$ ,

$$\begin{aligned} & \sup_{f \in \tilde{\mathcal{I}}(L, W, S, B), \|f - f_n\|_{\ell_2}^2 \leq r^2} \left\{ \|f - f_n\|_n^2 - \|f - f_n\|_{\ell_2}^2 \right\} \\ & \leq 3\mathbb{E} \left( \sup_{f \in \tilde{\mathcal{I}}(L, W, S, B), \|f - f_n\|_{\ell_2}^2 \leq r^2} \frac{1}{n} \sum_{i=1}^n \xi_i (f(X_i) - f_n(X_i))^2 \right) + r \cdot 2F \sqrt{\frac{2\gamma}{n}} + \frac{28F^2\gamma}{3n}, \\ & \leq 3r^2, \end{aligned}$$

where the last inequality is obtained by equation A.9 and equation A.10. Then with probability at least  $1 - e^{-\gamma}$ ,  $\|f - f_n\|_{\ell_2}^2 \leq r^2$  with  $f \in \tilde{\mathcal{I}}(L, W, S, B)$  implies  $\|f - f_n\|_n^2 \leq 4r^2$ .  $\square$

**Lemma A.4.** Suppose that  $h^2 \lesssim n^{-\frac{s}{2s+d}}$  and  $\|\hat{f}_h - f_n\|_{\ell_2} \leq r_0$ , with  $r_0$  satisfying equation A.9, equation A.10, and Assumption 4 holds. Also, with the notation of Assumption 3, suppose that for the class  $\mathcal{I}(L, W, S, B)$  the parameters are chosen as

$$L = 3 + 2 \left\lceil \log_2 \left( \frac{3^{\max\{d, m\}}}{\epsilon c_{d, m}} \right) + 5 \right\rceil \lceil \log_2 \max\{d, m\} \rceil, \quad W = W_0 N \quad (\text{A.12})$$

$$S = (L - 1)W_0^2 N + N, \quad B = O\left(N^{(v^{-1} + d^{-1})(\max\{1, (d/p-s)+\})}\right) \quad (\text{A.13})$$

for a constant  $c_{d, m}$  that depends on  $d$  and  $m$ , a constant  $W_0$ , and where  $v = (s - \delta)/\delta$ ,

$$\delta = \frac{d}{p}, \quad N \asymp n^{\frac{d}{2s+d}}. \quad (\text{A.14})$$

Then there exists a universal constant  $C_0 > 0$  such that

$$\begin{aligned} \|\hat{f}_h - f_n\|_{\ell_2}^2 & \leq C_0 \left[ r_0 F^{5/2} \sqrt{\frac{\gamma}{n}} + \frac{F^{5/2}\gamma}{n} + \right. \\ & \left. r_0 F \sqrt{\frac{N(\log N)^2}{n}} + r_0 F \sqrt{\frac{N[(\log N)^2 + \log r_0^{-1} + \log n]}{n}} + N^{-s/d} r_0 F^{3/2} \right] \end{aligned}$$

with probability at least  $1 - 3e^{-\gamma}$ , where  $N \asymp n^{\frac{d}{2s+d}}$ .

*Proof.* Let

$$\mathcal{G} = \left\{ g : g(\mathbf{x}, y) = \ell_h(y - f(\mathbf{x})) - \ell_h(y - f_n(\mathbf{x})), \quad f \in \tilde{\mathcal{I}}(L, W, S, B), \|f - f_n\|_{\ell_2} \leq r_0 \right\}.$$

Then for  $\xi_1, \dots, \xi_n$  independent Rademacher variables independent of  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , by Theorem 2.1 in Bartlett et al. (2005), with probability at least  $1 - 2e^{-\gamma}$ , we have that

$$\begin{aligned} & M_n(\hat{f}_h) - \hat{M}_n(\hat{f}_h) + (\|f_n - f_\tau^*\|_\infty + h^2) \|\hat{f}_h - f_n\|_{\ell_2} \sqrt{F} \\ & \lesssim \mathbb{E} \left( \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \xi_i g(\mathbf{x}_i, y_i) \mid (\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n) \right) \\ & \quad + 4r_0 F^{3/2} \sqrt{\frac{\gamma}{n}} + \frac{100F^{3/2}\gamma}{3n} + (\|f_n - f_\tau^*\|_\infty + h^2) \|\hat{f}_h - f_n\|_{\ell_2} F^{1/2}. \quad (\text{A.15}) \end{aligned}$$

Denote  $\mathbb{E}_\xi$  as the expectation with respect to  $\xi_1, \dots, \xi_n$ . Let

$$\varphi_{f,i}(t_i) = \ell_h(y_i - (t_i + f_n(\mathbf{x}_i))) - \ell_h(y_i - f_n(\mathbf{x}_i)),$$

where  $t_i = f(\mathbf{x}_i) - f_n(\mathbf{x}_i)$ . Note that  $\ell_h(\cdot)$  is 1-Lipschitz continuous and  $\varphi_{f,i}(0) = 0$ , by Talagrand's inequality (Ledoux & Talagrand (2013)), Lemma A.3, with probability at least  $1 - e^{-\gamma}$ , we have

$$\begin{aligned} \mathbb{E}_\xi \left( \sup_{g \in \mathcal{G}} \frac{1}{n} \sum_{i=1}^n \xi_i g(\mathbf{x}_i, y_i) \right) &= \mathbb{E}_\xi \left( \sup_{f \in \mathcal{I}(L, W, S, B), \|f\|_\infty \leq F, \|f - f_n\|_{\ell_2} \leq r_0} \frac{1}{n} \sum_{i=1}^n \xi_i \varphi_{f,i}(t_i) \right), \\ &\leq \mathbb{E}_\xi \left( \sup_{f \in \mathcal{I}(L, W, S, B), \|f\|_\infty \leq F, \|f - f_n\|_{\ell_2} \leq r_0} \frac{1}{n} \sum_{i=1}^n \xi_i (f(\mathbf{x}_i) - f_n(\mathbf{x}_i)) \right), \\ &\leq \mathbb{E}_\xi \left( \sup_{f \in \mathcal{I}(L, W, S, B), \|f\|_\infty \leq F, \|f_n - f\|_n \leq 2r_0} \frac{1}{n} \sum_{i=1}^n \xi_i (f(\mathbf{x}_i) - f_n(\mathbf{x}_i)) \right). \end{aligned} \tag{A.16}$$

By Dudley's chaining inequality and arguments in the proof of Theorem 2 in Suzuki (2019), we further have for some constant  $C > 0$ ,

$$\begin{aligned} &\mathbb{E}_\xi \left( \sup_{f \in \mathcal{I}(L, W, S, B), \|f\|_\infty \leq F, \|f_n - f\|_n \leq 2r_0} \frac{1}{n} \sum_{i=1}^n \xi_i (f(\mathbf{x}_i) - f_n(\mathbf{x}_i)) \right) \\ &\leq \inf_{0 < \alpha < r_0} \left\{ 4\alpha + \frac{24r_0}{\sqrt{n}} \sqrt{\log \mathcal{N} \left( \alpha, \tilde{\mathcal{I}}(L, W, S, B), \|\cdot\|_\infty \right)} \right\} \\ &\leq C \inf_{0 < \alpha < r_0} \left\{ \alpha + r_0 \sqrt{\frac{N \log(N) [\log^2 N + \log \alpha^{-1}]}{n}} \right\} \end{aligned} \tag{A.17}$$

By Proposition 1 in Suzuki (2019) and  $N \asymp n^{\frac{d}{2s+d}}$ ,  $h^2 \lesssim n^{-\frac{s}{2s+d}}$ , we have  $\|f_n - f_\tau^*\|_\infty + h^2 \lesssim N^{-s/d}$ . Let

$$\alpha = r_0 \sqrt{\frac{N(\log N)^2}{n}},$$

together with Lemma A.2, equation A.17, equation A.15, and equation A.16, we have for some constant  $C > 0$ ,

$$\begin{aligned} \left\| \hat{f}_h - f_\tau^* \right\|_{\ell_2}^2 &\leq C \left[ 4r_0 F^{5/2} \sqrt{\frac{\gamma}{n}} + \frac{100F^{5/2}\gamma}{3n} + \right. \\ &\quad \left. r_0 F \sqrt{\frac{N(\log N)^2}{n}} + r_0 F \sqrt{\frac{N \log(N) [(\log N)^2 + \log r_0^{-1} + \log n]}{n}} + r_0 N^{-s/d} F^{3/2} \right] \end{aligned}$$

with probability at least  $1 - 3e^{-\gamma}$ .  $\square$

**Lemma A.5** (Lemma 20 in Padilla et al. (2022)). *Let  $r^*$  be defined as*

$$r^* = \inf \left\{ r > 0 : 3\mathbb{E} \left( \sup_{f \in \tilde{\mathcal{I}}(L, W, S, B), \|f - f_n\|_{\ell_2} \leq s} \frac{1}{n} \sum_{i=1}^n \xi_i (f(x_i) - f_n(x_i))^2 \right) < s^2, \forall s \geq r \right\},$$

846 for  $\{\xi_i\}_{i=1}^n$  Rademacher variables independent of  $\{(x_i, y_i)\}_{i=1}^n$ . Then under the conditions A.12, A.13, and  
 847 A.14 of Lemma A.4,

$$848 \quad r^* \leq \tilde{C} \left[ \sqrt{\frac{N(\log N)^2}{n}} + \sqrt{\frac{N \log(N) [(\log N)^2 + \log n]}{n}} \right], \quad (\text{A.18})$$

851 for a constant  $\tilde{C} > 0$  and with  $N$  satisfying  $N \asymp n^{\frac{d}{2s+d}}$ .

852 Let  $\mathcal{H}$  be a class of functions from  $\mathcal{X}$  to  $\mathbb{R}$ . We define the pseudodimension of  $\mathcal{H}$ , denoted as  $\text{Pdim}(\mathcal{H})$ ,  
 853 as the largest integer  $m$  for which there exist  $(a_1, b_1) \dots, (a_m, b_m) \in \mathcal{X} \times \mathbb{R}$  such that for all  $\eta \in \{0, 1\}^m$   
 854 there exists  $f \in \mathcal{H}$  such that

$$855 \quad f(a_i) > b_i \iff \eta_i,$$

856 for  $i = 1, \dots, m$ .

857 **Lemma A.6** (Theorem 12 from Padilla et al. (2022)). *With the notation from before, for the neural network*  
 858 *function class  $\mathcal{F}(P, U, L)$ , we have*

$$859 \quad \text{Pdim}(\mathcal{F}(P, U, L)) = O(LP \log(U)).$$

## 860 A.2 PROOF OF THEOREM 3.1

861 *Proof.* Write the ball centered in  $f$  of radius  $r$

$$862 \quad \text{B}(f, \|\cdot\|_{\ell_2}, r) = \{g : \|f - g\|_{\ell_2} \leq r\}.$$

863 We divide the space  $\tilde{\mathcal{I}}(L, W, S, B) = \{f : f \in \mathcal{I}(L, W, S, B), \|f\|_{\infty} \leq F\}$  into sets of increasing radius

$$864 \quad \text{B}(f_n, \|\cdot\|_{\ell_2}, \bar{r}), \text{B}(f_n, \|\cdot\|_{\ell_2}, 2\bar{r}) \setminus \text{B}(f_n, \|\cdot\|_{\ell_2}, \bar{r}), \dots, \text{B}(f_n, \|\cdot\|_{\ell_2}, 2^l \bar{r}) \setminus \text{B}(f_n, \|\cdot\|_{\ell_2}, 2^{l-1} \bar{r}),$$

865 where

$$866 \quad l = \left\lceil \log_2 \left( \frac{2F}{\sqrt{\log n/n}} \right) \right\rceil.$$

867 If for some  $j \leq l$ ,

$$868 \quad \hat{f}_h \in \text{B}(f_n, \|\cdot\|_{\ell_2}, 2^j \bar{r}),$$

869 then by Lemma A.4, with probability at least  $1 - 3e^{-\gamma}$ , we have for some constant  $\tilde{C} > 0$ ,

$$870 \quad \left\| \hat{f}_h - f_n \right\|_{\ell_2}^2 \leq \tilde{C} \left( 2^j \bar{r} F^{5/2} \sqrt{\frac{\gamma}{n}} + \frac{F^{5/2} \gamma}{n} + \right. \\ 871 \quad \left. 2^j \bar{r} F \sqrt{\frac{N(\log N)^2}{n}} + 2^j \bar{r} F \sqrt{\frac{N \log(N) [(\log N)^2 + 2 \log n]}{n}} + N^{-s/d} 2^j \bar{r} F^{3/2} \right). \quad (\text{A.19})$$

872 Recall the  $r^*$  defined in Lemma A.5. We set

$$873 \quad \bar{r} = 8\tilde{C} \left[ F^{5/2} \sqrt{\frac{\gamma}{n}} + F \sqrt{\frac{N(\log N)^2}{n}} + F \sqrt{\frac{N \log(N) [(\log N)^2 + 2 \log n]}{n}} + N^{-s/d} F^{3/2} \right] \\ 874 \quad + 2\sqrt{2\tilde{C}} \cdot \sqrt{\frac{F^{5/2} \gamma}{n}} + r^*. \quad (\text{A.20})$$

By Lemma A.5 and  $N \asymp n^{d/(2s+d)}$ , setting  $\gamma \asymp \log n$ , it follows that  $\tilde{\mathcal{I}}(L, W, S, B) \subset \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, 2^l \bar{r})$  when  $n$  is sufficiently large. Therefore,  $\hat{f}_h \in \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, 2^l \bar{r})$  with probability 1 if  $n$  is sufficiently large.

Elementary calculation yields  $\bar{r} > r^*$  and for all  $0 \leq j \leq l$ ,

$$\frac{2^j \bar{r}}{8} \geq \tilde{C} \left[ F^{5/2} \sqrt{\frac{\gamma}{n}} + F \sqrt{\frac{N(\log N)^2}{n}} + F \sqrt{\frac{N \log(N) [(\log N)^2 + 2 \log n]}{n}} + N^{-s/d} F^{3/2} \right], \quad (\text{A.21})$$

and

$$\frac{2^{2j} \bar{r}^2}{8} \geq \tilde{C} \left( \frac{F^{5/2} \gamma}{n} \right). \quad (\text{A.22})$$

Combining equation A.21, equation A.22, and equation A.19, we have with probability at least  $1 - 3e^{-\gamma}$ ,

$$\|\hat{f}_h - f_n\|_{\ell_2} \leq 2^{j-1} \bar{r}. \quad (\text{A.23})$$

Now we begin from the first step of our localization procedure. Note that  $\bar{r} > r^*$ , by Lemmas A.3 and A.4, with probability at least  $1 - e^{-\gamma}$ , we have that

$$\|\hat{f}_h - f_n\|_{\ell_2} \leq 2^l \bar{r} \quad \text{implies} \quad \|\hat{f}_h - f_n\|_n \leq 2^{l+1} \bar{r}.$$

Then by arguments above, with probability at least  $1 - 4e^{-\gamma}$ ,  $\|\hat{f}_h - f_n\|_{\ell_2} \leq 2^l \bar{r}$  implies that

$$\|\hat{f}_h - f_n\|_{\ell_2} \leq 2^{l-1} \bar{r}, \quad \text{and} \quad \|\hat{f}_h - f_n\|_n \leq 2^l \bar{r}.$$

Continue recursively, we arrive at

$$\|\hat{f}_h - f_n\|_{\ell_2} \leq \bar{r}, \quad \text{and} \quad \|\hat{f}_h - f_n\|_n \leq 2\bar{r},$$

with probability approaching one.

Specifically, this procedure can be formulated as

$$\begin{aligned} \mathbb{P}(\hat{f}_h \in \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, \bar{r})) &= \mathbb{P}(\hat{f}_h \in \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, 2\bar{r})) - \mathbb{P}(\hat{f}_h \in \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, 2\bar{r}) \setminus \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, \bar{r})), \\ &\geq \mathbb{P}(\hat{f}_h \in \mathcal{B}(f_n, \|\cdot\|_{\ell_2}, 2\bar{r})) - 4e^{-\gamma}, \\ &\geq \dots, \\ &\geq 1 - 4(l+1)e^{-\gamma} = 1 - o(1), \end{aligned} \quad (\text{A.24})$$

with setting  $\gamma \asymp \log n$ .

By Lemma A.5 and equation A.20, we have

$$\begin{aligned} \bar{r} &\leq 8\tilde{C} \left[ F^{5/2} \sqrt{\frac{\gamma}{n}} + F \sqrt{\frac{N(\log N)^2}{n}} + F \sqrt{\frac{N \log(N) [(\log N)^2 + 2 \log n]}{n}} + N^{-s/d} F^{3/2} \right] \\ &\quad + 2\sqrt{2} \cdot \sqrt{\frac{\tilde{C} F^{5/2} \gamma}{n}} + \tilde{C} \left[ \sqrt{\frac{N(\log N)^2}{n}} + \sqrt{\frac{N \log(N) [(\log N)^2 + \log n]}{n}} \right]. \end{aligned}$$

Then the claim follows by  $N \asymp n^{d/(2s+d)}$  and  $\|f_n - f_\tau^*\|_\infty \lesssim N^{-s/d}$  in Proposition 1 of Suzuki (2019).  $\square$

## A.3 PROOF OF THEOREM 3.2

*Proof.* Throughout this proof, the covariates  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are fixed. For simplicity, we write  $f \in \mathcal{F}$  instead of  $f \in \mathcal{F}(P, U, L)$  within the proof. Let  $\hat{\delta}_i = \hat{f}_h(\mathbf{x}_i) - f_\tau^*(\mathbf{x}_i)$  for  $i = 1, \dots, n$ . Note that  $\ell_h(u) = \int_{-\infty}^{\infty} \rho_\tau(v) K_h(v-u) dv = \int_{-\infty}^{\infty} \rho_\tau(u+s) K_h(s) ds$ , then for any  $x, y \in \mathbb{R}$ ,

$$\ell_h(x) - \ell_h(y) = \int_{-\infty}^{\infty} K_h(s) (\rho_\tau(x+s) - \rho_\tau(y+s)) ds. \quad (\text{A.25})$$

By Knight identity (Knight, 1998), for any  $\delta \in \mathbb{R}$ ,

$$\begin{aligned} & \rho_\tau(y_i - (f_\tau^*(\mathbf{x}_i) + \hat{\delta}_i) + s) - \rho_\tau(y_i - f_\tau^*(\mathbf{x}_i) + s) \\ &= -\hat{\delta}_i (\tau - \mathbf{1}\{y_i + s \leq f_\tau^*(\mathbf{x}_i)\}) + \int_0^{\hat{\delta}_i} (\mathbf{1}\{y_i \leq f_\tau^*(\mathbf{x}_i) + z - s\} - \mathbf{1}\{y_i \leq f_\tau^*(\mathbf{x}_i) - s\}) dz. \end{aligned} \quad (\text{A.26})$$

By Assumption 1, Fubini's theorem and mean value expansion, using the fact  $\int s K_h(s) ds = 0$ ,  $\int s^2 K_h(s) ds = \sigma_K^2 h^2$ , we have

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \int K_h(s) \int_0^{\hat{\delta}_i} (\mathbf{1}\{y_i \leq f_\tau^*(\mathbf{x}_i) + z - s\} - \mathbf{1}\{y_i \leq f_\tau^*(\mathbf{x}_i) - s\}) dz ds \middle| \mathbf{x}_1, \dots, \mathbf{x}_n \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \int K_h(s) \int_0^{\hat{\delta}_i} F_{y_i|\mathbf{x}_i}(f_\tau^*(\mathbf{x}_i) + z - s) - F_{y_i|\mathbf{x}_i}(f_\tau^*(\mathbf{x}_i) - s) dz ds \middle| \mathbf{x}_1, \dots, \mathbf{x}_n \right], \\ &\geq \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n \left( p \int_0^{|\hat{\delta}_i|} \min\{z, \kappa\} dz - \underline{c} |\hat{\delta}_i| h^2 \right) \middle| \mathbf{x}_1, \dots, \mathbf{x}_n \right], \end{aligned} \quad (\text{A.27})$$

where the constants  $p, \underline{c} > 0$  are uniform for all  $i = 1, \dots, n$ .

Similarly, we also have

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n -\hat{\delta}_i \int K_h(s) (\tau - \mathbf{1}\{y_i + s \leq f_\tau^*(\mathbf{x}_i)\}) ds \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n -\hat{\delta}_i \int K_h(s) (F_{y_i|\mathbf{x}_i}(f_\tau^*(\mathbf{x}_i)) - F_{y_i|\mathbf{x}_i}(f_\tau^*(\mathbf{x}_i) - s)) ds \right] \\ &\geq -\underline{c} \mathbb{E} \left( \frac{1}{n} \sum_{i=1}^n |\hat{\delta}_i| h^2 \right). \end{aligned} \quad (\text{A.28})$$

Combining equation A.25, equation A.26, equation A.27, and equation A.28, for  $D_h(t) = \min\{|t|, t^2\} - h^2|t|$ , we have

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n D_h^2 \left\{ f_\tau^*(\mathbf{x}_i) - \hat{f}_h(\mathbf{x}_i) \right\} \right] \\ &\leq \frac{1}{c_\tau n} \mathbb{E} \left( \sum_{i=1}^n \mathbb{E} \left[ \ell_h \left\{ y_i - f_\tau^*(\mathbf{x}_i) - \hat{\delta}_i \right\} \right] - \sum_{i=1}^n \mathbb{E} \left[ \ell_h \left\{ y_i - f_\tau^*(\mathbf{x}_i) \right\} \right] \right), \\ &= \frac{1}{c_\tau} \mathbb{E} \left\{ M_n(\hat{f}_h) \right\} + O(\text{err}_1). \end{aligned} \quad (\text{A.29})$$

Next, we only need to bound  $E \left\{ M_n(\hat{f}_h) \right\}$ . By symmetrization Lemma 10 in Madrid Padilla & Chatterjee (2021), and Talagrand’s inequality (Ledoux & Talagrand, 2013) using the fact  $\ell_h(\cdot)$  is 1-Lipschitz continuous, for i.i.d. Rademacher variables  $\xi_i, i = 1, \dots, n$ ,

$$\begin{aligned} \mathbb{E} \left\{ M_n(\hat{f}_h) \right\} &\leq \mathbb{E} \left\{ \sup_{f \in \mathcal{F}, \|f\|_\infty \leq F} \left[ M_n(f) - \hat{M}_n(f) \right] \right\}, \\ &\leq 2 \mathbb{E} \left\{ \sup_{f \in \mathcal{F}, \|f\|_\infty \leq F} \sum_{i=1}^n \xi_i \hat{M}_{n,i}(f) \right\}, \\ &\leq 2 \mathbb{E} \left\{ \sup_{f \in \mathcal{F}, \|f\|_\infty \leq 2F} \frac{1}{n} \sum_{i=1}^n \xi_i f(\mathbf{x}_i) \right\}. \end{aligned} \quad (\text{A.30})$$

By Dudley’s theorem and Lemma 4 in Farrell et al. (2021), we further have

$$\begin{aligned} F \mathbb{E} \left\{ \sup_{f \in \mathcal{F}, \|f\|_\infty \leq F} \frac{1}{n} \sum_{i=1}^n \xi_i \frac{f(\mathbf{x}_i)}{F} \right\} &\leq \frac{CF}{\sqrt{n}} \int_0^2 \sqrt{\log \mathcal{N}(\mu, \mathcal{F}/F, \|\cdot\|_n)} d\mu, \\ &\leq \frac{CF}{\sqrt{n}} \int_0^2 \sqrt{\log \left( \left( \frac{2 \cdot e \cdot n}{\mu \cdot \text{Pdim}(\mathcal{F})} \right)^{\text{Pdim}(\mathcal{F})} \right)} d\mu. \end{aligned} \quad (\text{A.31})$$

Combining Lemma A.6, equation A.30 and equation A.31, for some constant  $\tilde{C} > 0$ , we have

$$\mathbb{E} \left\{ M_n(\hat{f}_h) \right\} \leq \tilde{C} F \sqrt{\frac{LP \log U \cdot \log n}{n}}, \quad (\text{A.32})$$

which shows equation 3.8 combining with equation A.29.

For equation 3.9, when  $h^2 = o(\sqrt{\mathbb{E} \|\hat{f}_h - f_\tau^*\|_n^2})$ , by  $\sum_i^n |a_i|/n \leq \sqrt{\sum_{i=1}^n a_i^2/n}$  and Jensen’s inequality,

$$h^2 \mathbb{E} \|\hat{f}_h - f_\tau^*\|_{n,1} \leq h^2 \mathbb{E} \sqrt{\|\hat{f}_h - f_\tau^*\|_n^2} \leq h^2 \sqrt{\mathbb{E} \|\hat{f}_h - f_\tau^*\|_n^2} = o(\mathbb{E} \|\hat{f}_h - f_\tau^*\|_n^2). \quad (\text{A.33})$$

Then combining with the fact  $\|\hat{f}_h - f_\tau^*\|_n^2 \leq \max\{1, F\} \Delta_n^2(\hat{f}_h, f_\tau^*)$ , equation 3.9 can also hold.  $\square$

## B ADDITIONAL EXPERIMENTS

We state additional details of the experiment in this section. The whole experiments are implemented in PyTorch (Paszke et al., 2019). We use stochastic gradient descent (SGD) with the Nesterov method of momentum factor 0.9. For each sample size, we keep 1/10 of the data for validation. We start the learning rate at 0.1, and use scheduler ReduceLRonPlateau with factor 0.5 and patience 5 to adjust the learning rate dynamically. We also implement gradient calculation for three different conquer loss functions, see Remark 3.1 in He et al. (2023). We manually calculate gradients of the loss functions at the backpropagation step, and keep automatic differentiation for the rest of the computational graph. We also show MSE results without a stop criterion, that is, to train the models with full 100 epochs. Compared to the results in Table 1, there is no obvious advantage to training without stopping. On the contrary, using a stop criterion can improve generalization ability and reduce training time.

## B.1 TABLES

In this section, we first present MSE results for our conquer neural networks with a sample size of 10000 and different bandwidth choices, as shown in Table 2. Bandwidths  $h = \{0.001, 0.005, 0.01, 0.05, 0.1\}$  are considered. The results show that our conquer neural networks outperform the baseline model over a wide range of bandwidths.

Table 2: MSE performances for scenario 1-3, model A and B, sample size 10000 under different bandwidths, quantile levels, and smoothing kernels. The MSEs are averaged over 50 independent trials.

$h$	Method	Scenario 1, Model A					Scenario 1, Model B					
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
0.001	Baseline	0.0618	0.0067	0.0047	0.0063	0.1366	0.1691	0.0099	0.0066	0.0082	0.3882	
	Gaussian	0.0678	<b>0.0064</b>	<b>0.0034</b>	<b>0.0055</b>	<b>0.0625</b>	<b>0.0654</b>	0.0101	<b>0.0061</b>	0.0083	<b>0.0708</b>	
	Uniform	0.0789	<b>0.0058</b>	<b>0.0036</b>	<b>0.0057</b>	<b>0.0543</b>	<b>0.1356</b>	<b>0.0094</b>	<b>0.0060</b>	0.0089	<b>0.0679</b>	
	Epanechnikov	0.0653	<b>0.0061</b>	<b>0.0037</b>	<b>0.0053</b>	<b>0.0665</b>	<b>0.1225</b>	<b>0.0098</b>	<b>0.0056</b>	0.0098	<b>0.0699</b>	
	0.005	Gaussian	0.0716	<b>0.0061</b>	<b>0.0033</b>	<b>0.0054</b>	<b>0.0580</b>	<b>0.1062</b>	<b>0.0095</b>	<b>0.0055</b>	0.0086	<b>0.0726</b>
		Uniform	0.0645	<b>0.0056</b>	<b>0.0038</b>	<b>0.0056</b>	<b>0.0549</b>	<b>0.0974</b>	<b>0.0093</b>	<b>0.0055</b>	<b>0.0080</b>	<b>0.0736</b>
		Epanechnikov	0.0676	<b>0.0062</b>	<b>0.0036</b>	<b>0.0053</b>	<b>0.0587</b>	<b>0.1294</b>	<b>0.0089</b>	<b>0.0057</b>	0.0092	<b>0.0663</b>
	0.01	Gaussian	0.0624	<b>0.0062</b>	<b>0.0037</b>	<b>0.0060</b>	<b>0.0662</b>	<b>0.0755</b>	<b>0.0091</b>	<b>0.0058</b>	0.0087	<b>0.0722</b>
		Uniform	0.0687	<b>0.0057</b>	<b>0.0037</b>	<b>0.0056</b>	<b>0.0805</b>	<b>0.0717</b>	<b>0.0090</b>	<b>0.0053</b>	0.0088	<b>0.0617</b>
Epanechnikov		0.0800	<b>0.0059</b>	<b>0.0035</b>	<b>0.0056</b>	<b>0.0749</b>	<b>0.0846</b>	<b>0.0093</b>	<b>0.0056</b>	<b>0.0081</b>	<b>0.0737</b>	
0.05	Gaussian	0.0913	0.0081	<b>0.0034</b>	0.0074	<b>0.0623</b>	<b>0.1055</b>	0.0120	<b>0.0059</b>	0.0107	<b>0.0748</b>	
	Uniform	0.0718	<b>0.0066</b>	<b>0.0036</b>	<b>0.0060</b>	<b>0.0582</b>	<b>0.0746</b>	<b>0.0094</b>	<b>0.0055</b>	0.0090	<b>0.0878</b>	
	Epanechnikov	0.1112	<b>0.0063</b>	<b>0.0033</b>	<b>0.0054</b>	<b>0.0744</b>	<b>0.0801</b>	<b>0.0093</b>	<b>0.0056</b>	0.0096	<b>0.0607</b>	
0.1	Gaussian	0.1153	0.0183	<b>0.0035</b>	0.0174	<b>0.1076</b>	<b>0.1593</b>	0.0215	<b>0.0057</b>	0.0203	<b>0.1167</b>	
	Uniform	0.0687	0.0101	<b>0.0037</b>	0.0090	<b>0.0758</b>	<b>0.1179</b>	0.0135	<b>0.0060</b>	0.0130	<b>0.0767</b>	
	Epanechnikov	0.0690	0.0085	<b>0.0036</b>	0.0071	<b>0.0681</b>	<b>0.1006</b>	0.0122	<b>0.0054</b>	0.0107	<b>0.0872</b>	
$h$	Method	Scenario 2, Model A					Scenario 2, Model B					
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
0.001	Baseline	0.1704	0.0205	0.0145	0.0223	0.1670	0.1323	0.0202	0.0129	0.0220	0.1358	
	Gaussian	<b>0.1316</b>	<b>0.0176</b>	<b>0.0128</b>	<b>0.0193</b>	<b>0.1537</b>	<b>0.1085</b>	<b>0.0160</b>	<b>0.0119</b>	<b>0.0178</b>	<b>0.1144</b>	
	Uniform	<b>0.1457</b>	<b>0.0180</b>	<b>0.0128</b>	<b>0.0191</b>	<b>0.1467</b>	<b>0.1172</b>	<b>0.0151</b>	<b>0.0118</b>	<b>0.0175</b>	<b>0.1342</b>	
	Epanechnikov	<b>0.1430</b>	<b>0.0162</b>	<b>0.0126</b>	<b>0.0191</b>	<b>0.1388</b>	<b>0.1038</b>	<b>0.0146</b>	<b>0.0107</b>	<b>0.0173</b>	<b>0.1302</b>	
	0.005	Gaussian	<b>0.1428</b>	<b>0.0184</b>	<b>0.0134</b>	<b>0.0182</b>	<b>0.1371</b>	<b>0.1029</b>	<b>0.0148</b>	<b>0.0118</b>	<b>0.0164</b>	<b>0.1264</b>
		Uniform	<b>0.1457</b>	<b>0.0184</b>	<b>0.0128</b>	<b>0.0185</b>	<b>0.1423</b>	<b>0.1250</b>	<b>0.0151</b>	<b>0.0115</b>	<b>0.0171</b>	<b>0.1275</b>
		Epanechnikov	<b>0.1551</b>	<b>0.0167</b>	<b>0.0121</b>	<b>0.0185</b>	<b>0.1262</b>	<b>0.1262</b>	<b>0.0186</b>	0.0129	<b>0.0172</b>	<b>0.1326</b>
	0.01	Gaussian	<b>0.1315</b>	<b>0.0165</b>	<b>0.0123</b>	<b>0.0172</b>	<b>0.1413</b>	<b>0.1166</b>	<b>0.0154</b>	<b>0.0112</b>	<b>0.0177</b>	<b>0.1117</b>
		Uniform	<b>0.1590</b>	<b>0.0169</b>	<b>0.0109</b>	<b>0.0180</b>	<b>0.1493</b>	<b>0.1079</b>	<b>0.0128</b>	<b>0.0108</b>	<b>0.0187</b>	<b>0.1142</b>
Epanechnikov		<b>0.1364</b>	<b>0.0184</b>	<b>0.0118</b>	<b>0.0197</b>	<b>0.1366</b>	<b>0.1118</b>	<b>0.0152</b>	<b>0.0111</b>	<b>0.0184</b>	<b>0.1119</b>	
0.05	Gaussian	<b>0.1478</b>	<b>0.0178</b>	<b>0.0122</b>	<b>0.0191</b>	<b>0.1553</b>	<b>0.1151</b>	<b>0.0146</b>	<b>0.0100</b>	<b>0.0193</b>	<b>0.1333</b>	
	Uniform	<b>0.1328</b>	<b>0.0171</b>	<b>0.0131</b>	<b>0.0195</b>	<b>0.1400</b>	<b>0.1302</b>	<b>0.0135</b>	<b>0.0120</b>	<b>0.0174</b>	<b>0.1071</b>	
	Epanechnikov	<b>0.1650</b>	<b>0.0162</b>	<b>0.0125</b>	<b>0.0202</b>	<b>0.1363</b>	<b>0.1093</b>	<b>0.0167</b>	<b>0.0109</b>	<b>0.0212</b>	<b>0.1247</b>	
0.1	Gaussian	<b>0.1558</b>	<b>0.0179</b>	<b>0.0121</b>	<b>0.0193</b>	<b>0.1341</b>	<b>0.0980</b>	<b>0.0157</b>	<b>0.0106</b>	<b>0.0190</b>	<b>0.1256</b>	
	Uniform	<b>0.1474</b>	<b>0.0176</b>	<b>0.0121</b>	<b>0.0188</b>	<b>0.1347</b>	<b>0.1198</b>	<b>0.0153</b>	<b>0.0126</b>	<b>0.0192</b>	<b>0.1169</b>	
	Epanechnikov	<b>0.1533</b>	<b>0.0181</b>	<b>0.0120</b>	<b>0.0196</b>	<b>0.1421</b>	<b>0.1170</b>	<b>0.0162</b>	<b>0.0118</b>	<b>0.0201</b>	<b>0.1302</b>	
$h$	Method	Scenario 3, Model A					Scenario 3, Model B					
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
0.001	Baseline	0.9232	0.2420	0.1459	0.2413	0.9960	0.7786	0.2306	0.1391	0.2380	0.7249	
	Gaussian	<b>0.8395</b>	<b>0.2040</b>	<b>0.1295</b>	<b>0.2145</b>	<b>0.7477</b>	0.8256	0.2518	0.1416	0.2444	0.7536	
	Uniform	<b>0.8790</b>	<b>0.2064</b>	<b>0.1300</b>	<b>0.2146</b>	<b>0.7824</b>	<b>0.7167</b>	0.2409	<b>0.1372</b>	0.2493	<b>0.6920</b>	
	Epanechnikov	<b>0.8230</b>	<b>0.2129</b>	<b>0.1349</b>	<b>0.2206</b>	<b>0.7998</b>	<b>0.7502</b>	0.2454	<b>0.1390</b>	<b>0.2315</b>	<b>0.6683</b>	
	0.005	Gaussian	<b>0.8707</b>	<b>0.2077</b>	<b>0.1308</b>	<b>0.2263</b>	<b>0.8742</b>	<b>0.7740</b>	0.2385	0.1413	0.2528	<b>0.6860</b>
		Uniform	<b>0.8739</b>	<b>0.2021</b>	<b>0.1276</b>	<b>0.2249</b>	<b>0.8332</b>	<b>0.7142</b>	<b>0.2279</b>	<b>0.1374</b>	<b>0.2292</b>	0.7574
		Epanechnikov	<b>0.8138</b>	<b>0.2084</b>	<b>0.1307</b>	<b>0.2166</b>	<b>0.8117</b>	<b>0.7051</b>	0.2376	<b>0.1344</b>	0.2529	0.7390
	0.01	Gaussian	<b>0.8995</b>	<b>0.1998</b>	<b>0.1275</b>	<b>0.2313</b>	<b>0.7526</b>	<b>0.7052</b>	0.2473	<b>0.1362</b>	0.2381	<b>0.7225</b>
		Uniform	<b>0.8739</b>	<b>0.2058</b>	<b>0.1289</b>	<b>0.2074</b>	<b>0.8256</b>	<b>0.7404</b>	0.2411	<b>0.1384</b>	0.2390	0.7370
Epanechnikov		<b>0.8079</b>	<b>0.2167</b>	<b>0.1317</b>	<b>0.2294</b>	<b>0.7893</b>	0.8089	0.2330	0.1393	0.2497	<b>0.6786</b>	
0.05	Gaussian	<b>0.8527</b>	<b>0.2083</b>	<b>0.1356</b>	<b>0.2211</b>	<b>0.7853</b>	<b>0.7434</b>	0.2357	<b>0.1332</b>	<b>0.2350</b>	<b>0.6987</b>	
	Uniform	<b>0.8250</b>	<b>0.2164</b>	<b>0.1352</b>	<b>0.2280</b>	<b>0.7966</b>	<b>0.7513</b>	0.2485	<b>0.1349</b>	0.2485	<b>0.6703</b>	
	Epanechnikov	<b>0.7879</b>	<b>0.2079</b>	<b>0.1310</b>	<b>0.2128</b>	<b>0.7526</b>	<b>0.6812</b>	0.2372	<b>0.1379</b>	0.2439	<b>0.7231</b>	
0.1	Gaussian	<b>0.8220</b>	<b>0.2125</b>	<b>0.1348</b>	<b>0.2114</b>	<b>0.8293</b>	<b>0.7644</b>	<b>0.2302</b>	<b>0.1326</b>	0.2467	<b>0.6604</b>	
	Uniform	<b>0.8100</b>	<b>0.2118</b>	<b>0.1371</b>	<b>0.2253</b>	<b>0.8218</b>	<b>0.7501</b>	<b>0.2295</b>	<b>0.1391</b>	0.2393	<b>0.7201</b>	
	Epanechnikov	<b>0.9015</b>	<b>0.2132</b>	<b>0.1310</b>	<b>0.2204</b>	<b>0.7852</b>	<b>0.6904</b>	0.2471	<b>0.1302</b>	<b>0.2262</b>	<b>0.6741</b>	

In the main text, Table 1 and Figure 2 show that our conquer networks obtain better performance and higher training efficiency. The faster training mainly comes from our stopping strategy, that is to stop training if the learning rate is lower than a threshold and the validation loss does not decrease for several consecutive epochs, while the baseline models are trained using all epochs according to the codes provided in Padilla et al. (2022). Therefore, we also train the conquer networks without the stop criterion and study the MSE results, which are shown in 3. We can conclude that our conquer networks still outperform the baseline models.

Table 3: MSE performances without stop criterion for scenario 1-3, Model A and B under different sample sizes, quantile levels and smoothing kernels. The MSEs are averaged over 50 independent trials.

Method	n=1000					n=5000					n=10000					
	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
S1 Model A	Baseline	<u>0.3820</u>	0.0402	0.0278	0.0374	0.3784	<u>0.0996</u>	0.0120	0.0086	0.0108	<u>0.0939</u>	<u>0.0618</u>	0.0067	0.0047	0.0063	0.1366
	Gaussian	0.4271	<b>0.0379</b>	<b>0.0224</b>	<b>0.0374</b>	0.5027	0.1118	<b>0.0088</b>	<b>0.0053</b>	<b>0.0095</b>	0.1081	0.0669	<b>0.0063</b>	<b>0.0033</b>	<b>0.0054</b>	<b>0.0626</b>
	Uniform	0.4413	<b>0.0378</b>	<b>0.0224</b>	<b>0.0326</b>	<b>0.3704</b>	0.1072	<b>0.0087</b>	<b>0.0059</b>	<b>0.0101</b>	0.1312	0.0789	<b>0.0058</b>	<b>0.0035</b>	<b>0.0055</b>	<b>0.0530</b>
Epanechnikov	0.6718	<b>0.0390</b>	<b>0.0218</b>	<b>0.0371</b>	0.5913	0.1249	<b>0.0091</b>	<b>0.0055</b>	<b>0.0095</b>	0.1153	0.0648	<b>0.0060</b>	<b>0.0037</b>	<b>0.0053</b>	<b>0.0662</b>	
S1 Model B	Baseline	<u>0.3842</u>	<u>0.0527</u>	<u>0.0319</u>	<u>0.0475</u>	<u>0.4222</u>	0.1143	0.0149	0.0107	0.0151	0.1277	0.1691	0.0099	0.0066	0.0082	0.3882
	Gaussian	0.4158	0.0666	0.0383	0.0632	0.5080	0.1172	<b>0.0143</b>	<b>0.0098</b>	<b>0.0141</b>	<b>0.1067</b>	<b>0.1063</b>	<b>0.0094</b>	<b>0.0055</b>	0.0084	<b>0.0729</b>
	Uniform	0.5645	0.0610	0.0376	0.0616	0.5685	<b>0.1034</b>	<b>0.0143</b>	<b>0.0091</b>	0.0153	<b>0.1255</b>	<b>0.0975</b>	<b>0.0093</b>	<b>0.0055</b>	<b>0.0080</b>	<b>0.0738</b>
Epanechnikov	0.4247	0.0576	0.0383	0.0623	0.6050	<b>0.1109</b>	<b>0.0144</b>	<b>0.0092</b>	<b>0.0146</b>	<b>0.1164</b>	<b>0.1291</b>	<b>0.0089</b>	<b>0.0058</b>	0.0092	<b>0.0662</b>	
S2 Model A	Baseline	0.8292	0.0868	0.0619	0.0839	<u>0.7874</u>	0.2752	<u>0.0275</u>	0.0222	0.0308	0.2747	0.1704	0.0205	0.0145	0.0223	0.1670
	Gaussian	<b>0.7994</b>	<b>0.0711</b>	<b>0.0587</b>	<b>0.0778</b>	1.1450	<b>0.2202</b>	0.0276	<b>0.0169</b>	<b>0.0273</b>	<b>0.2598</b>	<b>0.1316</b>	<b>0.0176</b>	<b>0.0128</b>	<b>0.0193</b>	<b>0.1537</b>
	Uniform	0.8890	<b>0.0722</b>	<b>0.0472</b>	0.0968	0.8566	<b>0.2320</b>	0.0286	<b>0.0181</b>	<b>0.0275</b>	<b>0.2586</b>	<b>0.1457</b>	<b>0.0180</b>	<b>0.0128</b>	<b>0.0191</b>	<b>0.1467</b>
Epanechnikov	0.9192	<b>0.0787</b>	<b>0.0522</b>	0.1048	0.9065	<b>0.2415</b>	0.0284	<b>0.0177</b>	<b>0.0265</b>	<b>0.2492</b>	<b>0.1430</b>	<b>0.0162</b>	<b>0.0125</b>	<b>0.0191</b>	<b>0.1388</b>	
S2 Model B	Baseline	0.4930	<u>0.0583</u>	<u>0.0493</u>	0.0840	<u>0.5898</u>	0.1732	0.0257	0.0178	0.0300	0.1966	0.1323	0.0202	0.0129	0.0220	0.1358
	Gaussian	0.7345	0.0633	0.0499	<b>0.0759</b>	0.6273	0.1800	<b>0.0246</b>	<b>0.0155</b>	<b>0.0245</b>	0.2157	<b>0.1085</b>	<b>0.0160</b>	<b>0.0119</b>	<b>0.0178</b>	<b>0.1144</b>
	Uniform	<b>0.4506</b>	0.0716	0.0504	0.0935	0.6030	<b>0.1707</b>	<b>0.0216</b>	<b>0.0176</b>	<b>0.0241</b>	0.2233	<b>0.1172</b>	<b>0.0151</b>	<b>0.0118</b>	<b>0.0175</b>	<b>0.1342</b>
Epanechnikov	0.5800	0.0723	0.0543	<b>0.0806</b>	0.7471	0.2151	0.0263	<b>0.0159</b>	<b>0.0289</b>	<b>0.1792</b>	<b>0.1038</b>	<b>0.0146</b>	<b>0.0107</b>	<b>0.0173</b>	<b>0.1302</b>	
S3 Model A	Baseline	3.0766	0.7564	0.5350	0.7407	<u>2.8969</u>	1.2870	0.3854	0.2146	0.3387	1.3495	0.9232	0.2420	0.1459	0.2413	0.9960
	Gaussian	<b>2.9080</b>	0.7769	<b>0.4784</b>	<b>0.7056</b>	3.4181	1.3139	<b>0.3379</b>	<b>0.1994</b>	<b>0.3229</b>	<b>1.1897</b>	<b>0.8395</b>	<b>0.2038</b>	<b>0.1295</b>	<b>0.2145</b>	<b>0.7477</b>
	Uniform	3.3730	0.7690	<b>0.4876</b>	0.7832	3.3576	<b>1.1593</b>	<b>0.3449</b>	<b>0.1912</b>	<b>0.3088</b>	<b>1.2392</b>	<b>0.8790</b>	<b>0.2064</b>	<b>0.1300</b>	<b>0.2146</b>	<b>0.7824</b>
Epanechnikov	3.4479	<b>0.7308</b>	<b>0.4679</b>	0.8076	3.3156	<b>1.1229</b>	<b>0.3616</b>	<b>0.1980</b>	<b>0.3099</b>	<b>1.3312</b>	<b>0.8230</b>	<b>0.2123</b>	<b>0.1337</b>	<b>0.2206</b>	<b>0.7998</b>	
S3 Model B	Baseline	2.8166	<u>0.7558</u>	0.5175	0.7665	2.2839	1.0061	<u>0.3596</u>	<u>0.2196</u>	0.3551	1.0990	0.7786	<u>0.2306</u>	0.1391	0.2380	0.7249
	Gaussian	<b>2.3173</b>	0.8406	<b>0.5142</b>	<b>0.7536</b>	2.6542	1.0602	0.4255	0.2302	<b>0.3526</b>	<b>1.0681</b>	0.8244	0.2516	0.1416	0.2443	0.7536
	Uniform	<b>2.1946</b>	0.8319	0.5193	<b>0.7352</b>	2.5765	1.2056	0.4038	0.2292	<b>0.3372</b>	<b>1.0528</b>	<b>0.7167</b>	0.2411	<b>0.1359</b>	0.2493	<b>0.6890</b>
Epanechnikov	2.9684	0.9008	<b>0.4892</b>	0.8704	<b>2.1302</b>	<b>1.0053</b>	0.3741	0.2297	<b>0.3473</b>	<b>1.0896</b>	<b>0.7474</b>	0.2453	<b>0.1390</b>	<b>0.2315</b>	<b>0.6683</b>	

We also study the effect of residual-based neural networks. We add one residual block for every hidden layer in the baseline and conquer networks to study the MSE performance. We compare the baseline model performance without residual blocks to that using residual-based structures. The results are shown in Table 4. The numbers with brackets in the table represent the standard error of the MSEs over the 50 trials, and the bold font means that the model has a smaller MSE. We find from the results that the residual-based neural networks have no advantage in reducing MSE in many scenarios. We also make the comparison for conquer networks, as shown in Table 5 for the Gaussian kernel, for example. Residual blocks also fail to significantly reduce the MSE in the Gaussian conquer model. However, we find that for high ( $\tau = 0.95$ ) or low ( $\tau = 0.05$ ) quantile levels, the models with residual structures tend to have smaller MSEs and standard errors, which increases the training stability. We provide the table of the MSE performances for baseline models and conquer networks with residual blocks, see Table 6. We can see that for residual-based networks, the conquer networks still have better performance than the baseline models.

We also implement the simulations for the joint estimation of multiple quantile levels under non-crossing constraints, see Padilla et al. (2022) for the estimation of baseline networks. For the conquer networks, for multiple quantile levels are given in the set  $\Gamma \subset (0, 1)$ , we estimate the quantile functions by solving

$$\begin{aligned} \{\hat{f}_\tau\}_{\tau \in \Gamma} = \arg \min_{\{f_\tau\}_{\tau \in \Gamma}} \sum_{\tau \in \Gamma} \sum_{i=1}^n \ell_\tau(y_i - f_\tau(\mathbf{x}_i)) \\ \text{subject to } f_\tau(\mathbf{x}_i) \leq f_{\tau'}(\mathbf{x}_i) \quad \forall \tau < \tau', \tau, \tau' \in \Gamma, i = 1, \dots, n. \end{aligned}$$

Table 4: Mean squared error (MSE) and its standard error performances for baseline, scenario 1-3, model A and B with residual blocks under different sample sizes, quantile levels. ‘‘Original’’ means the original baseline network without residual blocks, and ‘‘+Res’’ means the baseline network with residual blocks. The MSEs are averaged over 50 independent trials. The bracket means the standard error of the MSEs over the trials.

Baseline	$n=1000$					$n=5000$					$n=10000$							
	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$			
S1	Model A	Original	0.3820 (0.0916)	<b>0.0402</b> (0.0025)	<b>0.0278</b> (0.0019)	<b>0.0374</b> (0.0024)	0.3784 (0.0524)	<b>0.0996</b> (0.0172)	<b>0.0120</b> (0.0005)	<b>0.0086</b> (0.0007)	<b>0.0108</b> (0.0005)	0.0939 (0.0207)	<b>0.0618</b> (0.0048)	<b>0.0067</b> (0.0003)	<b>0.0047</b> (0.0002)	<b>0.0063</b> (0.0004)	0.1366 (0.0585)	
		+ Res	<b>0.3316</b> (0.0343)	0.0575 (0.0029)	0.0324 (0.0018)	0.0376 (0.0025)	<b>0.3508</b> (0.0729)	0.1051 (0.0071)	0.0221 (0.0010)	0.0120 (0.0006)	0.0121 (0.0008)	<b>0.0797</b> (0.0057)	0.0661 (0.0037)	0.0150 (0.0006)	0.0088 (0.0004)	0.0076 (0.0005)	<b>0.0596</b> (0.0070)	
	Model B	Original	0.3842 (0.0705)	<b>0.0527</b> (0.0036)	<b>0.0319</b> (0.0019)	0.0475 (0.0032)	0.4222 (0.0709)	<b>0.1143</b> (0.0160)	0.0149 (0.0006)	<b>0.0107</b> (0.0004)	<b>0.0151</b> (0.0008)	0.1277 (0.0234)	0.1691 (0.0656)	<b>0.0099</b> (0.0006)	<b>0.0066</b> (0.0003)	<b>0.0082</b> (0.0004)	<b>0.3882</b> (0.3125)	
		+ Res	<b>0.3457</b> (0.0365)	0.0716 (0.0040)	0.0373 (0.0017)	<b>0.0424</b> (0.0031)	<b>0.3231</b> (0.0515)	0.1249 (0.0081)	0.0293 (0.0011)	0.0191 (0.0008)	0.0162 (0.0011)	<b>0.0620</b> (0.0048)	<b>0.0783</b> (0.0057)	0.0200 (0.0007)	0.0121 (0.0005)	0.0118 (0.0007)	0.0455 (0.0037)	
S2	Model A	Original	0.8292 (0.1170)	0.0868 (0.0162)	<b>0.0619</b> (0.0046)	<b>0.0839</b> (0.0061)	0.7874 (0.1149)	0.2752 (0.0264)	<b>0.0275</b> (0.0021)	<b>0.0222</b> (0.0016)	<b>0.0308</b> (0.0020)	0.2747 (0.0216)	0.1704 (0.0121)	<b>0.0205</b> (0.0013)	<b>0.0145</b> (0.0006)	<b>0.0223</b> (0.0013)	0.1670 (0.0076)	
		+ Res	<b>0.7892</b> (0.1507)	<b>0.0801</b> (0.0061)	0.0655 (0.0059)	0.1017 (0.0060)	<b>0.7819</b> (0.1265)	<b>0.2290</b> (0.0220)	0.0315 (0.0022)	0.0244 (0.0017)	0.0354 (0.0019)	<b>0.2551</b> (0.0164)	<b>0.1285</b> (0.0106)	0.0213 (0.0015)	0.0157 (0.0008)	0.0231 (0.0012)	0.1723 (0.0096)	
	Model B	Original	<b>0.4930</b> (0.0742)	<b>0.0583</b> (0.0038)	<b>0.0493</b> (0.0043)	<b>0.0840</b> (0.0062)	<b>0.5898</b> (0.0712)	<b>0.1732</b> (0.0114)	<b>0.0257</b> (0.0015)	<b>0.0178</b> (0.0009)	<b>0.0300</b> (0.0018)	<b>0.1966</b> (0.0164)	0.1323 (0.0146)	0.0202 (0.0017)	<b>0.0129</b> (0.0006)	0.0220 (0.0013)	0.1358 (0.0070)	
		+ Res	0.5210 (0.1177)	0.0683 (0.0056)	0.0651 (0.0084)	0.0869 (0.0054)	0.7145 (0.1019)	0.1847 (0.0214)	0.0261 (0.0015)	0.0205 (0.0008)	0.0335 (0.0023)	0.2108 (0.0160)	<b>0.1254</b> (0.0102)	<b>0.0174</b> (0.0010)	0.0142 (0.0007)	<b>0.0204</b> (0.0009)	<b>0.1302</b> (0.0090)	
S3	Model A	Original	<b>3.0766</b> (0.2728)	<b>0.7564</b> (0.0388)	<b>0.5350</b> (0.0248)	<b>0.7407</b> (0.0330)	<b>2.8969</b> (0.2179)	1.2870 (0.0844)	0.3854 (0.0216)	<b>0.2146</b> (0.0051)	<b>0.3387</b> (0.0109)	1.3495 (0.0976)	0.9232 (0.0348)	<b>0.2420</b> (0.0082)	<b>0.1459</b> (0.0028)	<b>0.2413</b> (0.0067)	0.9960 (0.0636)	
		+ Res	3.1729 (0.4222)	0.8820 (0.0415)	0.5587 (0.0139)	0.8456 (0.0403)	3.0478 (0.2550)	<b>1.1872</b> (0.0536)	<b>0.3741</b> (0.0106)	0.2298 (0.0057)	0.3626 (0.0107)	<b>1.1419</b> (0.0664)	<b>0.7821</b> (0.0382)	0.2505 (0.0066)	0.1708 (0.0035)	0.2539 (0.0062)	<b>0.7933</b> (0.0389)	
	Model B	Original	2.8166 (0.2967)	<b>0.7558</b> (0.0383)	0.5175 (0.0205)	0.7665 (0.0509)	2.2839 (0.2600)	1.0061 (0.0579)	0.3596 (0.0131)	0.2196 (0.0054)	0.3551 (0.0109)	1.0990 (0.0780)	0.7786 (0.0315)	<b>0.2306</b> (0.0075)	<b>0.1391</b> (0.0036)	<b>0.2380</b> (0.0084)	0.7249 (0.0330)	
		+ Res	<b>2.2566</b> (0.2357)	0.7801 (0.0379)	<b>0.5042</b> (0.0171)	<b>0.7446</b> (0.0301)	<b>2.2770</b> (0.1769)	<b>0.9813</b> (0.0587)	<b>0.3330</b> (0.0103)	<b>0.2148</b> (0.0043)	<b>0.3146</b> (0.0083)	<b>0.9632</b> (0.0474)	<b>0.7292</b> (0.0358)	0.2358 (0.0068)	0.1554 (0.0025)	0.2451 (0.0061)	<b>0.6422</b> (0.0267)	

Table 5: Mean squared error (MSE) and its standard error performances for Gaussian conquer, scenario 1-3, model A and B with residual blocks under different sample sizes, quantile levels. ‘‘Original’’ means the original Gaussian conquer network without residual blocks, and ‘‘+Res’’ means the Gaussian conquer model with residual blocks. The MSEs are averaged over 50 independent trials. The bracket means the standard error of the MSEs over the trials.

Gaussian	$n=1000$					$n=5000$					$n=10000$						
	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$		
S1	Model A	Original	0.4354 (0.0767)	<b>0.0383</b> (0.0025)	0.0224 (0.0014)	0.0382 (0.0030)	0.5035 (0.0858)	0.1124 (0.0120)	<b>0.0087</b> (0.0004)	<b>0.0055</b> (0.0003)	0.0097 (0.0007)	0.1081 (0.0148)	0.0678 (0.0088)	<b>0.0064</b> (0.0005)	<b>0.0034</b> (0.0001)	<b>0.0055</b> (0.0004)	0.0625 (0.0060)
		+ Res	<b>0.3787</b> (0.0573)	0.0620 (0.0030)	<b>0.0217</b> (0.0014)	<b>0.0262</b> (0.0021)	<b>0.4817</b> (0.0825)	<b>0.1006</b> (0.0058)	0.0249 (0.0011)	0.0118 (0.0006)	<b>0.0090</b> (0.0006)	<b>0.0889</b> (0.0105)	<b>0.0641</b> (0.0037)	0.0245 (0.0009)	0.0122 (0.0006)	0.0100 (0.0008)	<b>0.0465</b> (0.0044)
	Model B	Original	0.4158 (0.0341)	<b>0.0665</b> (0.0043)	0.0383 (0.0014)	0.0633 (0.0024)	0.5202 (0.0427)	<b>0.1172</b> (0.0113)	<b>0.0144</b> (0.0006)	<b>0.0097</b> (0.0004)	0.0142 (0.0005)	0.1066 (0.0122)	0.1062 (0.0307)	<b>0.0095</b> (0.0005)	<b>0.0055</b> (0.0003)	<b>0.0086</b> (0.0005)	0.0726 (0.0124)
		+ Res	<b>0.3525</b> (0.0380)	0.0797 (0.0036)	<b>0.0306</b> (0.0018)	<b>0.0308</b> (0.0021)	<b>0.5173</b> (0.0915)	0.1360 (0.0094)	0.0386 (0.0016)	0.0187 (0.0011)	<b>0.0131</b> (0.0009)	<b>0.0709</b> (0.0077)	<b>0.0829</b> (0.0052)	0.0276 (0.0011)	0.0156 (0.0006)	0.0110 (0.0007)	<b>0.0443</b> (0.0044)
S2	Model A	Original	0.7994 (0.1005)	0.0711 (0.0068)	0.0587 (0.0088)	<b>0.0778</b> (0.0053)	1.1466 (0.2029)	0.2202 (0.0174)	0.0276 (0.0025)	<b>0.0169</b> (0.0010)	<b>0.0273</b> (0.0017)	0.2598 (0.0284)	0.1316 (0.0091)	0.0176 (0.0011)	<b>0.0128</b> (0.0007)	0.0193 (0.0016)	0.1537 (0.0102)
		+ Res	<b>0.7481</b> (0.1358)	<b>0.0704</b> (0.0082)	<b>0.0577</b> (0.0043)	0.0906 (0.0064)	<b>0.7275</b> (0.0984)	<b>0.2076</b> (0.0219)	<b>0.0260</b> (0.0015)	0.0223 (0.0028)	0.0308 (0.0017)	<b>0.2394</b> (0.0315)	<b>0.1302</b> (0.0075)	<b>0.0152</b> (0.0009)	0.0134 (0.0006)	<b>0.0187</b> (0.0010)	<b>0.1354</b> (0.0092)
	Model B	Original	0.7367 (0.1652)	<b>0.0639</b> (0.0057)	<b>0.0500</b> (0.0049)	<b>0.0759</b> (0.0053)	0.6273 (0.0972)	0.1800 (0.0171)	0.0246 (0.0020)	<b>0.0155</b> (0.0009)	<b>0.0245</b> (0.0024)	0.2157 (0.0236)	<b>0.1085</b> (0.0087)	0.0160 (0.0017)	<b>0.0119</b> (0.0007)	<b>0.0178</b> (0.0013)	0.1144 (0.0078)
		+ Res	<b>0.7154</b> (0.1862)	0.0700 (0.0083)	0.0553 (0.0043)	0.0827 (0.0053)	<b>0.5483</b> (0.1012)	<b>0.1430</b> (0.0134)	<b>0.0237</b> (0.0011)	0.0195 (0.0013)	0.0284 (0.0014)	<b>0.1835</b> (0.0138)	0.1008 (0.0093)	<b>0.0152</b> (0.0010)	0.0120 (0.0006)	0.0187 (0.0011)	<b>0.1036</b> (0.0070)
S3	Model A	Original	2.8841 (0.3062)	<b>0.7776</b> (0.0746)	<b>0.4788</b> (0.0206)	<b>0.7056</b> (0.0500)	3.4222 (0.3403)	1.3139 (0.1159)	0.3379 (0.0135)	<b>0.1993</b> (0.0067)	0.3269 (0.0100)	1.1897 (0.0783)	0.8395 (0.0565)	<b>0.2040</b> (0.0055)	<b>0.1295</b> (0.0027)	<b>0.2145</b> (0.0055)	0.7477 (0.0402)
		+ Res	<b>2.8543</b> (0.2822)	0.9541 (0.0714)	0.5828 (0.0455)	0.7702 (0.0347)	<b>2.7395</b> (0.2494)	<b>1.1222</b> (0.0464)	<b>0.3339</b> (0.0117)	0.2143 (0.0045)	<b>0.3045</b> (0.0084)	<b>0.9996</b> (0.0518)	<b>0.7556</b> (0.0360)	0.2403 (0.0084)	0.1515 (0.0027)	0.2281 (0.0056)	<b>0.6931</b> (0.0286)
	Model B	Original	<b>2.3193</b> (0.2530)	0.8405 (0.0573)	<b>0.5142</b> (0.0182)	<b>0.7543</b> (0.0441)	2.6520 (0.5055)	1.0602 (0.0791)	0.4263 (0.0202)	0.2304 (0.0070)	0.3525 (0.0122)	1.0681 (0.0822)	0.8256 (0.0460)	0.2518 (0.0132)	<b>0.1416</b> (0.0047)	0.2444 (0.0085)	0.7536 (0.0450)
		+ Res	2.5481 (0.2880)	<b>0.7569</b> (0.0426)	0.5322 (0.0230)	0.8249 (0.0643)	<b>2.0108</b> (0.1294)	<b>0.9965</b> (0.0655)	<b>0.3196</b> (0.0112)	<b>0.2130</b> (0.0083)	<b>0.3027</b> (0.0108)	<b>0.8176</b> (0.0462)	<b>0.6785</b> (0.0374)	<b>0.2169</b> (0.0068)	0.1470 (0.0035)	<b>0.2303</b> (0.0064)	<b>0.6120</b> (0.0310)

Table 6: Mean squared error (MSE) performances for scenario 1-3, model A and B with residual blocks under different sample sizes, quantile levels, and smoothing kernels. The MSEs are averaged over 50 independent trials. The baseline model is boxed if it outperforms our conquer networks.

Method	n=1000					n=5000					n=10000					
	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
S1 Model A	Baseline	<b>0.3316</b>	<b>0.0575</b>	0.0324	0.0376	<b>0.3508</b>	0.1051	<b>0.0221</b>	0.0120	0.0121	<b>0.0797</b>	0.0661	<b>0.0150</b>	<b>0.0088</b>	<b>0.0076</b>	0.0596
	Gaussian	0.3787	0.0620	<b>0.0217</b>	<b>0.0262</b>	0.4817	<b>0.1006</b>	0.0249	<b>0.0118</b>	<b>0.0090</b>	0.0889	<b>0.0641</b>	0.0245	0.0122	0.0100	<b>0.0465</b>
	Uniform	0.3719	0.0581	<b>0.0218</b>	<b>0.0234</b>	0.4008	0.1139	0.0255	<b>0.0113</b>	<b>0.0092</b>	0.0881	0.0689	0.0239	0.0127	0.0099	<b>0.0483</b>
	Epanechnikov	0.3356	0.0629	<b>0.0225</b>	<b>0.0259</b>	0.3973	0.1149	0.0259	<b>0.0115</b>	<b>0.0090</b>	0.0928	<b>0.0627</b>	0.0254	0.0127	0.0094	<b>0.0434</b>
S1 Model B	Baseline	<b>0.3457</b>	<b>0.0716</b>	0.0373	0.0424	<b>0.3231</b>	0.1249	<b>0.0293</b>	0.0191	0.0162	<b>0.0620</b>	0.0783	<b>0.0200</b>	<b>0.0121</b>	0.0118	0.0455
	Gaussian	0.3525	0.0797	<b>0.0306</b>	<b>0.0308</b>	0.5173	0.1360	0.0386	<b>0.0187</b>	<b>0.0131</b>	0.0709	0.0829	0.0276	0.0156	<b>0.0110</b>	<b>0.0443</b>
	Uniform	0.3596	0.0795	<b>0.0322</b>	<b>0.0352</b>	0.4551	<b>0.1240</b>	0.0389	0.0196	<b>0.0137</b>	0.0698	<b>0.0769</b>	0.0275	0.0153	0.0123	<b>0.0440</b>
	Epanechnikov	0.3680	0.0871	<b>0.0348</b>	<b>0.0304</b>	0.4675	0.1272	0.0379	<b>0.0190</b>	<b>0.0131</b>	0.0669	0.0810	0.0269	0.0157	<b>0.0113</b>	<b>0.0379</b>
S2 Model A	Baseline	0.7892	0.0801	0.0655	0.1017	0.7819	0.2290	0.0315	0.0244	0.0354	0.2551	0.1285	0.0213	0.0157	0.0231	0.1723
	Gaussian	<b>0.7481</b>	<b>0.0704</b>	<b>0.0577</b>	<b>0.0906</b>	<b>0.7275</b>	<b>0.2076</b>	<b>0.0260</b>	<b>0.0223</b>	<b>0.0308</b>	<b>0.2394</b>	1.1302	<b>0.0152</b>	<b>0.0134</b>	<b>0.0187</b>	<b>0.1354</b>
	Uniform	0.8190	<b>0.0734</b>	<b>0.0647</b>	0.1053	0.8890	<b>0.2000</b>	<b>0.0240</b>	<b>0.0209</b>	<b>0.0353</b>	<b>0.2544</b>	<b>0.1136</b>	<b>0.0168</b>	<b>0.0137</b>	<b>0.0188</b>	<b>0.1168</b>
	Epanechnikov	<b>0.6887</b>	0.0827	0.0716	0.1185	<b>0.6827</b>	<b>0.1907</b>	<b>0.0265</b>	<b>0.0232</b>	<b>0.0302</b>	0.2563	<b>0.1227</b>	<b>0.0168</b>	<b>0.0137</b>	<b>0.0197</b>	<b>0.1233</b>
S2 Model B	Baseline	<b>0.5210</b>	0.0683	0.0651	0.0869	0.7145	0.1847	0.0261	0.0205	0.0335	0.2108	0.1254	0.0174	0.0142	0.0204	0.1302
	Gaussian	0.7154	0.0700	<b>0.0553</b>	<b>0.0827</b>	<b>0.5483</b>	<b>0.1430</b>	<b>0.0237</b>	<b>0.0195</b>	<b>0.0284</b>	<b>0.1835</b>	<b>0.1008</b>	<b>0.0152</b>	<b>0.0120</b>	<b>0.0187</b>	<b>0.1036</b>
	Uniform	0.6440	<b>0.0605</b>	<b>0.0523</b>	<b>0.0808</b>	<b>0.5217</b>	<b>0.1445</b>	<b>0.0215</b>	<b>0.0180</b>	<b>0.0318</b>	<b>0.1884</b>	<b>0.1043</b>	<b>0.0149</b>	<b>0.0113</b>	<b>0.0182</b>	<b>0.1074</b>
	Epanechnikov	0.6414	0.0879	<b>0.0580</b>	<b>0.0868</b>	<b>0.5577</b>	<b>0.1682</b>	<b>0.0235</b>	<b>0.0179</b>	<b>0.0268</b>	<b>0.1823</b>	<b>0.1049</b>	<b>0.0143</b>	<b>0.0116</b>	<b>0.0203</b>	<b>0.1154</b>
S3 Model A	Baseline	3.1729	0.8820	<b>0.5587</b>	0.8456	3.0478	1.1872	0.3741	0.2298	0.3626	1.1419	0.7821	0.2505	0.1708	0.2539	0.7933
	Gaussian	<b>2.8543</b>	0.9541	0.5828	<b>0.7702</b>	<b>2.7395</b>	<b>1.1222</b>	<b>0.3339</b>	<b>0.2143</b>	<b>0.3045</b>	<b>0.9996</b>	<b>0.7556</b>	<b>0.2403</b>	<b>0.1515</b>	<b>0.2281</b>	<b>0.6931</b>
	Uniform	3.9049	<b>0.8399</b>	0.5665	<b>0.7472</b>	<b>2.8420</b>	<b>1.1147</b>	<b>0.3521</b>	<b>0.2104</b>	<b>0.3096</b>	<b>1.1170</b>	0.8147	<b>0.2228</b>	<b>0.1468</b>	<b>0.2301</b>	<b>0.7264</b>
	Epanechnikov	<b>3.0467</b>	<b>0.8131</b>	0.5638	<b>0.7606</b>	<b>2.6823</b>	1.2444	<b>0.3520</b>	<b>0.2145</b>	<b>0.3247</b>	1.2093	<b>0.7679</b>	<b>0.2242</b>	<b>0.1531</b>	<b>0.2232</b>	<b>0.7798</b>
S3 Model B	Baseline	2.2566	0.7801	0.5042	0.7446	2.2770	0.9813	0.3330	0.2148	0.3146	0.9632	0.7292	0.2358	0.1554	0.2451	0.6422
	Gaussian	2.5481	<b>0.7569</b>	0.5322	0.8249	<b>2.0108</b>	0.9965	<b>0.3196</b>	<b>0.2130</b>	<b>0.3027</b>	<b>0.8176</b>	<b>0.6785</b>	<b>0.2169</b>	<b>0.1470</b>	<b>0.2303</b>	<b>0.6120</b>
	Uniform	<b>1.9425</b>	0.8181	<b>0.4824</b>	0.7644	<b>1.9375</b>	<b>0.9431</b>	0.3394	<b>0.2002</b>	<b>0.3013</b>	<b>0.8535</b>	<b>0.6307</b>	<b>0.2177</b>	<b>0.1462</b>	<b>0.2086</b>	<b>0.6415</b>
	Epanechnikov	2.4130	0.7807	0.5403	<b>0.7333</b>	<b>1.9167</b>	<b>0.8839</b>	<b>0.3262</b>	<b>0.2112</b>	0.3211	<b>0.9264</b>	<b>0.7278</b>	<b>0.2058</b>	<b>0.1471</b>	<b>0.2277</b>	<b>0.5911</b>

To solve the problem, we let  $\tau_0 < \dots < \tau_m$  be the elements of  $\Gamma$  and solve

$$\{\hat{g}_\tau\}_{\tau \in \Gamma} = \arg \min_{\{g_\tau\}_{\tau \in \Gamma}} \sum_{i=1}^n \ell_h(y_i - g_{\tau_0}(\mathbf{x}_i)) + \sum_{j=1}^m \sum_{i=1}^n \ell_h \left\{ y_i - g_{\tau_0}(\mathbf{x}_i) - \sum_{l=1}^j \log \left( 1 + e^{g_{\tau_l}(\mathbf{x}_i)} \right) \right\}$$

and set

$$\hat{f}_{\tau_0}(\mathbf{x}) = \hat{g}_{\tau_0}(\mathbf{x}), \quad \text{and} \quad \hat{f}_{\tau_j}(\mathbf{x}) = \hat{g}_{\tau_0}(\mathbf{x}) + \sum_{l=1}^j \log \left( 1 + e^{\hat{g}_{\tau_l}(\mathbf{x})} \right) \quad \text{for } j = 1, \dots, m,$$

where we recall that  $\ell_h(\cdot)$  is the convolution-type smoothed quantile loss for the quantile  $\tau$  in equation 2.2. For the simulation data, we study the MSE performance, see Table 7 below. In contrast to Table 1 in the paper, we can see that the joint estimation of multiple quantile levels performs much better for high ( $\tau=0.95$ ) and low ( $\tau=0.05$ ) quantiles, while the performance for  $\tau = 0.25/0.5/0.75$  declines a little bit in exchange.

We also make it clear to the data-driven rules for bandwidth selection. In detail, we propose the K-fold cross-validation algorithm for the bandwidth selection. For a candidate list of bandwidths, we train the models on the training set and calculate the pinball loss on the validation set for K times. Select the bandwidth with the minimum mean validation loss. Simulation studies have been implemented for Scenario 2 and the MSE results are shown in Table 8. By the cross-validation algorithm, our conquer networks still outperform the baseline models, especially for large sample sizes, which shows the stability of our method. We also refer to Table 2 to show the results' stability for different bandwidth choices.

In addition, we implemented the MAE loss and the pinball loss in the experiment part. The MAE results for multiple quantile levels joint estimation are shown in Table 9. The MAE results for 5-fold cross-validation of Scenario 2 are shown in Table 10. The results for MAE show that our conquer networks outperform the baseline networks in most cases, which remains consistent compared to the MSE results in Tables 7 and 8,

Table 7: Mean squared error (MSE) performances for scenario 1-3, model A and B under different sample sizes, quantile levels, and smoothing kernels. Multiple quantile levels are trained jointly under the non-crossing constraint. The MSEs are averaged over 50 independent trials.

	Method	$n=1000$					$n=5000$					$n=10000$					
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
S1	Model A	Baseline	0.1312	0.0514	0.0278	0.0417	0.1221	0.0494	0.0147	0.0100	0.0146	0.0397	0.0305	0.0088	0.0056	0.0078	0.0234
		Gaussian	<b>0.1258</b>	<b>0.0427</b>	<b>0.0224</b>	<b>0.0395</b>	0.1235	<b>0.0477</b>	<b>0.0111</b>	<b>0.0070</b>	<b>0.0120</b>	<b>0.0361</b>	0.0338	<b>0.0076</b>	<b>0.0040</b>	<b>0.0060</b>	<b>0.0215</b>
		Uniform	<b>0.1212</b>	<b>0.0424</b>	<b>0.0240</b>	<b>0.0444</b>	0.1331	<b>0.0471</b>	<b>0.0111</b>	<b>0.0067</b>	<b>0.0113</b>	<b>0.0355</b>	0.0340	<b>0.0077</b>	<b>0.0039</b>	<b>0.0059</b>	<b>0.0222</b>
		Epanechnikov	<b>0.1204</b>	<b>0.0453</b>	<b>0.0235</b>	<b>0.0408</b>	0.1289	0.0495	<b>0.0123</b>	<b>0.0079</b>	<b>0.0128</b>	<b>0.0372</b>	0.0361	<b>0.0080</b>	<b>0.0043</b>	<b>0.0063</b>	<b>0.0217</b>
		Baseline	0.1911	0.0623	0.0352	0.0591	0.1821	0.0767	0.0192	0.0119	0.0191	0.0621	0.0547	0.0123	0.0074	0.0106	0.0362
		Model B	Gaussian	0.2105	0.0711	0.0378	0.0681	0.2041	0.0886	<b>0.0178</b>	<b>0.0100</b>	<b>0.0169</b>	<b>0.0617</b>	0.0674	0.0126	<b>0.0063</b>	<b>0.0096</b>
		Uniform	0.2145	0.070	0.0378	0.0662	0.2066	0.0864	<b>0.0171</b>	<b>0.0097</b>	<b>0.0159</b>	<b>0.0597</b>	0.0711	0.0134	<b>0.0067</b>	<b>0.0096</b>	0.0375
		Epanechnikov	0.2136	0.073	0.0412	0.0685	0.2034	0.0883	<b>0.0183</b>	<b>0.0103</b>	<b>0.0168</b>	<b>0.0605</b>	0.0726	0.0126	<b>0.0064</b>	<b>0.0099</b>	0.0391
S2	Model A	Baseline	0.1830	0.0829	0.0712	0.0910	0.2613	0.0650	0.0272	0.0248	0.0321	0.0888	0.0469	0.0199	0.0182	0.0234	0.0641
		Gaussian	<b>0.1736</b>	<b>0.0718</b>	<b>0.0645</b>	0.0948	<b>0.2526</b>	<b>0.0542</b>	<b>0.0222</b>	<b>0.0202</b>	<b>0.0262</b>	<b>0.0810</b>	<b>0.0395</b>	<b>0.0140</b>	<b>0.0129</b>	<b>0.0181</b>	<b>0.0539</b>
		Uniform	<b>0.1669</b>	<b>0.0641</b>	<b>0.0592</b>	<b>0.0827</b>	<b>0.2362</b>	<b>0.0549</b>	<b>0.0233</b>	<b>0.0211</b>	<b>0.0283</b>	<b>0.0870</b>	<b>0.0377</b>	<b>0.0132</b>	<b>0.0122</b>	<b>0.0165</b>	<b>0.0477</b>
		Epanechnikov	0.1901	<b>0.0786</b>	0.0749	0.1066	0.2634	<b>0.0535</b>	<b>0.0222</b>	<b>0.0192</b>	<b>0.0267</b>	<b>0.0803</b>	<b>0.0396</b>	<b>0.0148</b>	<b>0.0130</b>	<b>0.0172</b>	<b>0.0524</b>
		Baseline	0.1358	0.0494	0.0481	0.0690	0.2026	0.0645	0.0222	0.0191	0.0276	0.0890	0.0499	0.0164	0.0143	0.0198	0.0560
		Model B	Gaussian	0.1449	0.0572	0.0537	0.0780	0.2060	0.0737	<b>0.0207</b>	<b>0.0171</b>	<b>0.0261</b>	0.0899	<b>0.0480</b>	<b>0.0126</b>	<b>0.0108</b>	<b>0.0163</b>
		Uniform	0.1489	0.0559	0.0546	0.0835	0.2252	0.0728	0.0237	<b>0.0190</b>	<b>0.0266</b>	0.0906	<b>0.0492</b>	<b>0.0132</b>	<b>0.0112</b>	<b>0.0165</b>	<b>0.0542</b>
		Epanechnikov	0.1485	0.0567	0.0514	0.0800	0.2220	0.0777	0.0241	0.0208	0.0298	0.0937	<b>0.0466</b>	<b>0.0135</b>	<b>0.0110</b>	<b>0.0159</b>	<b>0.0541</b>
S3	Model A	Baseline	0.9823	0.7126	0.6113	0.7058	1.1469	0.3581	0.2787	0.2524	0.2740	0.4018	0.2699	0.1865	0.1778	0.2012	0.2859
		Gaussian	<b>0.8659</b>	<b>0.6814</b>	<b>0.5797</b>	<b>0.6803</b>	<b>1.0587</b>	<b>0.3368</b>	<b>0.2476</b>	<b>0.2292</b>	<b>0.2494</b>	<b>0.3801</b>	<b>0.2301</b>	<b>0.1597</b>	<b>0.1502</b>	<b>0.1654</b>	<b>0.2361</b>
		Uniform	<b>0.9403</b>	0.7268	0.6202	<b>0.6916</b>	<b>1.0572</b>	<b>0.3504</b>	<b>0.2589</b>	<b>0.2351</b>	<b>0.2550</b>	<b>0.3745</b>	<b>0.2297</b>	<b>0.1672</b>	<b>0.1572</b>	<b>0.1720</b>	<b>0.2371</b>
		Epanechnikov	<b>0.9368</b>	0.7153	<b>0.5725</b>	<b>0.6167</b>	<b>1.0125</b>	<b>0.3208</b>	<b>0.2507</b>	<b>0.2305</b>	<b>0.2518</b>	<b>0.3541</b>	<b>0.2305</b>	<b>0.1668</b>	<b>0.1554</b>	<b>0.1670</b>	<b>0.2398</b>
		Baseline	0.6933	0.5887	0.5255	0.5704	0.7379	0.2942	0.2617	0.2486	0.2605	0.3159	0.1990	0.1743	0.1704	0.1867	0.2247
		Model B	Gaussian	0.7480	0.6832	0.5910	0.6510	0.7794	0.3240	0.2964	0.2767	0.2867	0.3347	0.2008	0.1824	0.1748	<b>0.1826</b>
		Uniform	0.7214	0.6326	0.5468	0.6127	0.8332	0.3078	0.2791	0.2623	0.2682	0.3108	<b>0.1916</b>	<b>0.1700</b>	<b>0.1620</b>	<b>0.1698</b>	<b>0.1973</b>
		Epanechnikov	0.8064	0.7124	0.6185	0.6936	0.8509	0.3292	0.2900	0.2717	0.2821	0.3409	<b>0.1916</b>	<b>0.1712</b>	<b>0.1650</b>	<b>0.1725</b>	<b>0.2014</b>

Table 8: Mean squared error (MSE) performances for scenario 2, model A and B with 5-fold cross-validation under different sample sizes, quantile levels, and smoothing kernels. The MSEs are averaged over 50 independent trials.

$n$	Method	Scenario 2, Model A					Scenario 2, Model B				
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$
1000	Baseline	0.9195	0.0714	0.0586	0.0894	0.8167	0.5480	0.0674	0.0462	0.0838	0.4792
	Gaussian	<b>0.7859</b>	0.0763	<b>0.0498</b>	<b>0.0890</b>	0.8510	0.5849	0.0751	0.0555	<b>0.0814</b>	0.6314
	Uniform	0.9437	0.0827	<b>0.0569</b>	<b>0.0744</b>	<b>0.7308</b>	<b>0.5434</b>	0.0740	0.0499	<b>0.0802</b>	1.0782
	Epanechnikov	1.0141	0.0759	<b>0.0462</b>	0.0919	<b>0.6248</b>	0.5938	0.0688	0.0561	<b>0.0787</b>	0.6151
5000	Baseline	0.3446	0.0258	0.0241	0.0335	0.3439	0.1629	0.0265	0.0188	0.0272	0.2205
	Gaussian	<b>0.2358</b>	0.0275	<b>0.0180</b>	<b>0.0302</b>	<b>0.3064</b>	0.1724	<b>0.0199</b>	<b>0.0161</b>	<b>0.0261</b>	0.2244
	Uniform	<b>0.2477</b>	0.0310	<b>0.0184</b>	<b>0.0268</b>	<b>0.2732</b>	0.1804	<b>0.0256</b>	<b>0.0182</b>	<b>0.0258</b>	<b>0.2045</b>
	Epanechnikov	<b>0.2556</b>	<b>0.0237</b>	<b>0.0169</b>	<b>0.0273</b>	<b>0.2657</b>	0.2104	<b>0.0227</b>	<b>0.0161</b>	0.0295	0.2357
10000	Baseline	0.1511	0.0193	0.0164	0.0229	0.1853	0.1218	0.0152	0.0128	0.0223	0.1522
	Gaussian	0.1577	<b>0.0162</b>	<b>0.0139</b>	<b>0.0184</b>	<b>0.1401</b>	<b>0.1192</b>	<b>0.0138</b>	<b>0.0113</b>	<b>0.0189</b>	<b>0.1375</b>
	Uniform	<b>0.1489</b>	<b>0.0153</b>	<b>0.0120</b>	<b>0.0210</b>	<b>0.1508</b>	0.1255	0.0164	<b>0.0112</b>	<b>0.0171</b>	0.1592
	Epanechnikov	<b>0.1431</b>	<b>0.0165</b>	<b>0.0121</b>	<b>0.0187</b>	<b>0.1852</b>	<b>0.1091</b>	0.0167	<b>0.0113</b>	<b>0.0166</b>	<b>0.1122</b>

respectively. The pinball loss results for real data analysis are presented in Table 11. These table results indicate the stability of our conquer networks with respect to different loss metrics.

Table 9: Mean absolute error (MAE) performances for scenario 1-3, model A and B under different sample sizes, quantile levels, and smoothing kernels. Multiple quantile levels are trained jointly under the non-crossing constraint. The MAEs are averaged over 50 independent trials.

	Method	$n=1000$					$n=5000$					$n=10000$					
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	
S1	Baseline	0.2805	0.1798	0.1310	0.1622	0.2792	0.1718	0.0960	0.0781	0.0948	0.1583	0.1314	0.0723	0.0573	0.0694	0.1202	
	Model A	Gaussian	<b>0.2775</b>	<b>0.1645</b>	<b>0.1204</b>	<b>0.1587</b>	<b>0.2780</b>	<b>0.1678</b>	<b>0.0826</b>	<b>0.0647</b>	<b>0.0859</b>	<b>0.1516</b>	0.1396	<b>0.0675</b>	<b>0.0483</b>	<b>0.0607</b>	<b>0.1144</b>
	Uniform	<b>0.2713</b>	<b>0.1654</b>	<b>0.1246</b>	0.1671	0.2926	<b>0.1663</b>	<b>0.0829</b>	<b>0.0633</b>	<b>0.0828</b>	<b>0.1487</b>	0.1400	<b>0.0679</b>	<b>0.0487</b>	<b>0.0605</b>	<b>0.1165</b>	
	Epanechnikov	<b>0.2707</b>	<b>0.1715</b>	<b>0.1230</b>	<b>0.1607</b>	0.2870	<b>0.1703</b>	<b>0.0849</b>	<b>0.0673</b>	<b>0.0883</b>	<b>0.1510</b>	0.1456	<b>0.0695</b>	<b>0.0507</b>	<b>0.0622</b>	<b>0.1157</b>	
	Model B	Baseline	0.3414	0.1974	0.1484	0.1948	0.3444	0.2151	0.1084	0.0856	0.1094	0.1983	0.1807	0.0876	0.0675	0.0816	0.1524
	Gaussian	0.3552	0.2121	0.1551	0.2063	0.3620	0.2304	<b>0.1039</b>	<b>0.0786</b>	<b>0.1032</b>	0.1993	0.1971	0.0880	<b>0.0617</b>	<b>0.0775</b>	<b>0.1516</b>	
Uniform	0.3616	0.2121	0.1546	0.2017	0.3659	0.2281	<b>0.1016</b>	<b>0.0772</b>	<b>0.1003</b>	<b>0.1976</b>	0.2027	0.0910	<b>0.0639</b>	<b>0.0780</b>	0.1550		
Epanechnikov	0.3611	0.2133	0.1610	0.2073	0.3654	0.2312	<b>0.1056</b>	<b>0.0797</b>	<b>0.1027</b>	<b>0.1975</b>	0.2062	<b>0.0876</b>	<b>0.0621</b>	<b>0.0791</b>	0.1589		
S2	Baseline	0.3352	0.2137	0.1999	0.2265	0.3765	0.1955	0.1250	0.1182	0.1344	0.2236	0.1674	0.1060	0.1013	0.1159	0.1899	
	Model A	Gaussian	<b>0.3271</b>	<b>0.1968</b>	<b>0.1882</b>	<b>0.2235</b>	<b>0.3686</b>	<b>0.1794</b>	<b>0.1141</b>	<b>0.1086</b>	<b>0.1240</b>	<b>0.2125</b>	<b>0.1534</b>	<b>0.0885</b>	<b>0.0843</b>	<b>0.0986</b>	<b>0.1724</b>
	Uniform	<b>0.3256</b>	<b>0.1913</b>	<b>0.1851</b>	<b>0.2173</b>	<b>0.3611</b>	<b>0.1815</b>	<b>0.1143</b>	<b>0.1088</b>	<b>0.1256</b>	<b>0.2153</b>	<b>0.1494</b>	<b>0.0877</b>	<b>0.0832</b>	<b>0.0971</b>	<b>0.1661</b>	
	Epanechnikov	0.3432	<b>0.2054</b>	0.2012	0.2380	0.3781	<b>0.1798</b>	<b>0.1117</b>	<b>0.1057</b>	<b>0.1232</b>	<b>0.2117</b>	<b>0.1529</b>	<b>0.0900</b>	<b>0.0851</b>	<b>0.0987</b>	<b>0.1732</b>	
	Model B	Baseline	0.2947	0.1700	0.1680	0.2033	0.3517	0.1998	0.1162	0.1073	0.1289	0.2301	0.1733	0.0971	0.0919	0.1086	0.1814
	Gaussian	0.3027	0.1820	0.1770	0.2153	<b>0.3513</b>	0.2124	<b>0.1109</b>	<b>0.0990</b>	<b>0.1223</b>	<b>0.2294</b>	<b>0.1693</b>	<b>0.0869</b>	<b>0.0794</b>	<b>0.0977</b>	<b>0.1781</b>	
Uniform	0.3058	0.1780	0.1763	0.2184	0.3633	0.2112	<b>0.1150</b>	<b>0.1022</b>	<b>0.1233</b>	0.2320	<b>0.1726</b>	<b>0.0890</b>	<b>0.0807</b>	<b>0.0986</b>	<b>0.1801</b>		
Epanechnikov	0.3035	0.1741	0.1713	0.2136	0.3596	0.2188	0.1194	0.1076	0.1304	0.2368	<b>0.1672</b>	<b>0.0890</b>	<b>0.0801</b>	<b>0.0969</b>	<b>0.1791</b>		
S3	Baseline	0.7794	0.6503	0.6017	0.6386	0.8088	0.4631	0.4050	0.3822	0.3952	0.4868	0.4030	0.3326	0.3220	0.3402	0.4076	
	Model A	Gaussian	<b>0.7381</b>	<b>0.6396</b>	<b>0.5890</b>	<b>0.6228</b>	<b>0.7545</b>	<b>0.4453</b>	<b>0.3817</b>	<b>0.3651</b>	<b>0.3784</b>	<b>0.4728</b>	<b>0.3708</b>	<b>0.3071</b>	<b>0.2947</b>	<b>0.3067</b>	<b>0.3738</b>
	Uniform	<b>0.7678</b>	0.6559	0.6132	0.6444	<b>0.7743</b>	<b>0.4575</b>	<b>0.3907</b>	<b>0.3680</b>	<b>0.3797</b>	<b>0.4688</b>	<b>0.3712</b>	<b>0.3133</b>	<b>0.3017</b>	<b>0.3134</b>	<b>0.3742</b>	
	Epanechnikov	<b>0.7650</b>	<b>0.6417</b>	<b>0.5790</b>	<b>0.6037</b>	<b>0.7491</b>	<b>0.4391</b>	<b>0.3849</b>	<b>0.3672</b>	<b>0.3805</b>	<b>0.4604</b>	<b>0.3700</b>	<b>0.3130</b>	<b>0.2988</b>	<b>0.3085</b>	<b>0.3771</b>	
	Model B	Baseline	0.6603	0.5983	0.5712	0.5957	0.6654	0.4162	0.3933	0.3816	0.3874	0.4271	0.3380	0.3167	0.3113	0.3227	0.3520
	Gaussian	0.6923	0.6354	0.5995	0.6336	0.6862	0.4385	0.4147	0.3997	0.4032	0.4360	0.3422	0.3229	0.3138	<b>0.3173</b>	<b>0.3404</b>	
Uniform	0.6805	0.6133	0.5824	0.6204	0.6978	0.4292	0.4047	0.3915	0.3917	<b>0.4207</b>	0.3353	<b>0.3142</b>	<b>0.3030</b>	<b>0.3081</b>	<b>0.3357</b>		
Epanechnikov	0.7175	0.6433	0.6115	0.6542	0.7074	0.4433	0.4119	0.3987	0.4003	0.4336	0.3335	<b>0.3143</b>	<b>0.3040</b>	<b>0.3074</b>	<b>0.3351</b>		

Table 10: Mean absolute error (MAE) performances for scenario 2, model A and B with 5-fold cross-validation under different sample sizes, quantile levels, and smoothing kernels. The MAEs are averaged over 50 independent trials.

$n$	Method	Scenario 2, Model A					Scenario 2, Model B				
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$
1000	Baseline	0.6325	0.2014	0.1820	0.2285	0.6348	0.5071	0.1994	0.1684	0.2264	0.5078
	Gaussian	<b>0.5974</b>	0.2043	<b>0.1710</b>	<b>0.2199</b>	<b>0.6320</b>	<b>0.4969</b>	0.2017	0.1770	<b>0.2173</b>	0.5303
	Uniform	0.6401	0.2062	<b>0.1794</b>	<b>0.2024</b>	<b>0.5837</b>	<b>0.5030</b>	<b>0.1993</b>	0.1694	<b>0.2105</b>	0.5735
	Epanechnikov	0.6348	<b>0.1928</b>	<b>0.1663</b>	<b>0.2211</b>	<b>0.5642</b>	0.5129	<b>0.1922</b>	0.1778	<b>0.2142</b>	0.5409
5000	Baseline	0.3797	0.1214	0.1158	0.1371	0.4034	0.2954	0.1237	0.1053	0.1275	0.3513
	Gaussian	<b>0.3376</b>	0.1220	<b>0.1007</b>	<b>0.1302</b>	<b>0.4002</b>	<b>0.2905</b>	<b>0.1085</b>	<b>0.0973</b>	<b>0.1234</b>	<b>0.3410</b>
	Uniform	<b>0.3501</b>	0.1276	<b>0.1019</b>	<b>0.1249</b>	<b>0.3751</b>	0.3059	<b>0.1220</b>	<b>0.1020</b>	<b>0.1201</b>	<b>0.3337</b>
	Epanechnikov	<b>0.3513</b>	<b>0.1154</b>	<b>0.0982</b>	<b>0.1239</b>	<b>0.3680</b>	0.3228	<b>0.1151</b>	<b>0.0969</b>	0.1297	<b>0.3435</b>
10000	Baseline	0.2858	0.1034	0.0951	0.1156	0.3194	0.2515	0.0949	0.0866	0.1124	0.2863
	Gaussian	<b>0.2768</b>	<b>0.0937</b>	<b>0.0858</b>	<b>0.1014</b>	<b>0.2762</b>	<b>0.2373</b>	<b>0.0882</b>	<b>0.0800</b>	<b>0.1022</b>	<b>0.2774</b>
	Uniform	<b>0.2738</b>	<b>0.0937</b>	<b>0.0818</b>	<b>0.1064</b>	<b>0.2852</b>	0.2579	0.0951	<b>0.0802</b>	<b>0.0988</b>	0.2916
	Epanechnikov	<b>0.2744</b>	<b>0.0959</b>	<b>0.0818</b>	<b>0.1020</b>	<b>0.3157</b>	<b>0.2383</b>	<b>0.0943</b>	<b>0.0794</b>	<b>0.0974</b>	<b>0.2553</b>

## B.2 REAL DATA ANALYSIS

We also implement real data analysis to study the performance of the conquer networks. We study the BMI (body mass index) data (<https://www.kaggle.com/datasets/rehan497/health-lifestyle-dataset>) and attempt to predict it with age, daily steps, sleep hours, water intake, and calories consumed. We consider the samples without the family history, group the samples by gender and train our conquer networks with 5-fold cross-validation on the training set, which contains 28000 samples. We then evaluate the pinball losses on the testing set, which consists of 7000 samples. The results are shown in Table 11. Our conquer networks achieve better performance in most cases, showing the good generalizability of the conquer neural network.

Table 11: Pinball loss performance for BMI (body mass index) prediction. “Single Quantile” represents that the models are trained with every single quantile level. “Multiple Quantiles” represents that the models are trained with multiple quantile levels jointly under non-crossing constrains.

Gender	Method	Single Quantile					Multiple Quantiles				
		$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$	$\tau=0.05$	$\tau=0.25$	$\tau=0.5$	$\tau=0.75$	$\tau=0.95$
Male	Baseline	0.5231	2.0638	2.7387	2.0534	0.5190	0.5218	2.0629	2.7415	2.0523	0.5198
	Gaussian	<b>0.5221</b>	<b>2.0609</b>	<b>2.7384</b>	<b>2.0513</b>	<b>0.5189</b>	<b>0.5214</b>	<b>2.0625</b>	<b>2.7407</b>	<b>2.0520</b>	<b>0.5192</b>
	Uniform	<b>0.5217</b>	<b>2.0618</b>	2.7403	<b>2.0521</b>	0.5192	<b>0.5218</b>	<b>2.0619</b>	<b>2.7395</b>	<b>2.0517</b>	0.5201
	Epanechnikov	<b>0.5218</b>	<b>2.0636</b>	2.7389	<b>2.0518</b>	0.5205	<b>0.5215</b>	<b>2.0621</b>	<b>2.7399</b>	<b>2.0521</b>	<b>0.5197</b>
Female	Baseline	0.5218	2.0596	2.7366	2.0555	0.5249	0.5204	2.0449	2.7291	2.0531	0.5251
	Gaussian	<b>0.5202</b>	<b>2.0446</b>	<b>2.7323</b>	<b>2.0531</b>	0.5259	<b>0.5203</b>	2.0450	2.7305	<b>2.0529</b>	<b>0.5250</b>
	Uniform	<b>0.5211</b>	<b>2.0497</b>	<b>2.7311</b>	2.0642	0.5251	<b>0.5202</b>	<b>2.0449</b>	2.7295	<b>2.0526</b>	<b>0.5248</b>
	Epanechnikov	<b>0.5205</b>	<b>2.0454</b>	<b>2.7276</b>	<b>2.0523</b>	<b>0.5248</b>	0.5212	2.0460	2.7325	2.0579	0.5290

## B.3 PLOTS OF MSEs BY SAMPLE SIZES

We study the plot of MSEs as a function of sample sizes to directly corroborate Theorem 3.2. We plotted the curves of log MSEs by log sample sizes. For scenario 2 and model A, we trained the baseline model and the conquer network with the Gaussian smoothing kernel. We took sample sizes of {1000, 3000, 5000, 7000, 10000} and evaluated the mean MSEs of test sets over 50 trials. The plots are shown in Figure 3. The curves have linear shapes, and in equation 3.9 of Theorem 3.2, the log MSE also has a nearly linear upper bound with respect to log sample size, which corroborates the theoretical results in Theorem 3.2. Besides, the slopes of both baseline and conquer networks are larger than  $-1$ , which is the asymptotic slope of the minimax rate  $n^{-2s/(2s+d)}$ . We also fit linear models for the curves in the figure and find that the slope of the conquer networks is smaller than that of the baseline models, indicating that our upper bound is better than the baseline from the simulation.

## B.4 PLOTS OF LOSS LANDSCAPE

To express the benefit of the conquer network more intuitively, we tried to plot the loss landscape of the baseline network and the conquer network based on Li et al. (2018), which are shown in Figure 4. The networks have the same structures, consisting of 20 layers, each with 35 nodes. We generated two random directions and ensured the models shared the same directions. Then we used the filter-wise normalization method in Li et al. (2018) and calculated the loss surfaces. The loss landscapes plot shows that the baseline network (see the subplot on the left) has at least three local minima. In contrast, the conquer network in the right subplot has one minimum, which implies that the conquer network improves the training dynamics.

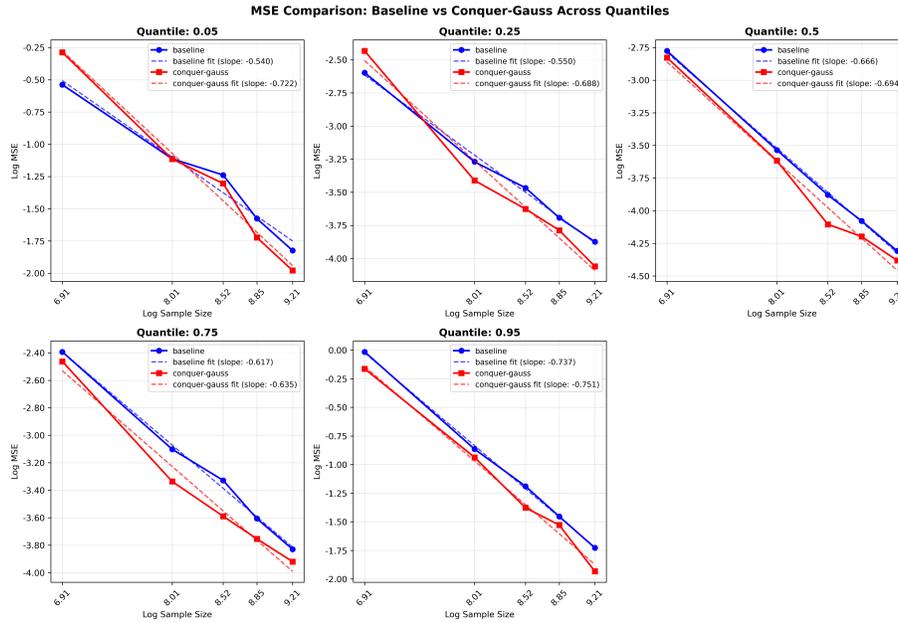


Figure 3: Plots of log MSEs by log sample sizes for scenario 2, model A. The red lines represent the conquer network smoothed by the Gaussian kernel with  $h = 0.1$ . The blue lines represent the baseline network. The solid lines with points represent the line charts of simulation results. The dashed lines represent the fitted linear models. The MSEs are averaged over 50 independent trials.

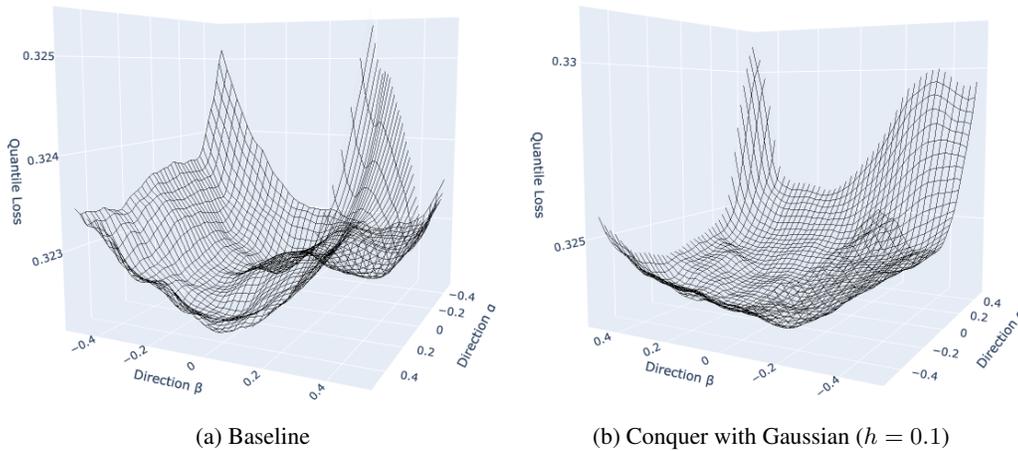


Figure 4: Plot of 3D loss landscape of scenario 2,  $\tau = 0.5$ . The subplot on the left represents the baseline model, and the right subplot represents the conquer smoothed by the Gaussian kernel with  $h = 0.1$ . Both networks consist of 20 layers, each with 35 nodes.

B.5 PLOTS OF TRAINING TIME

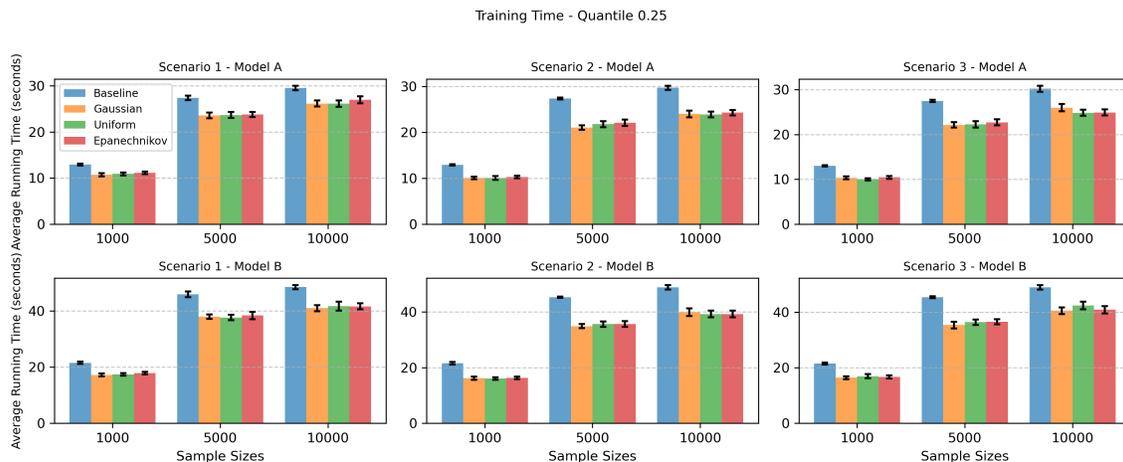


Figure 5: Bar chart with error bars with average training time over 50 trials under quantile level  $\tau = 0.25$ . The error bars represent 95% confidence intervals of training time for each setting.

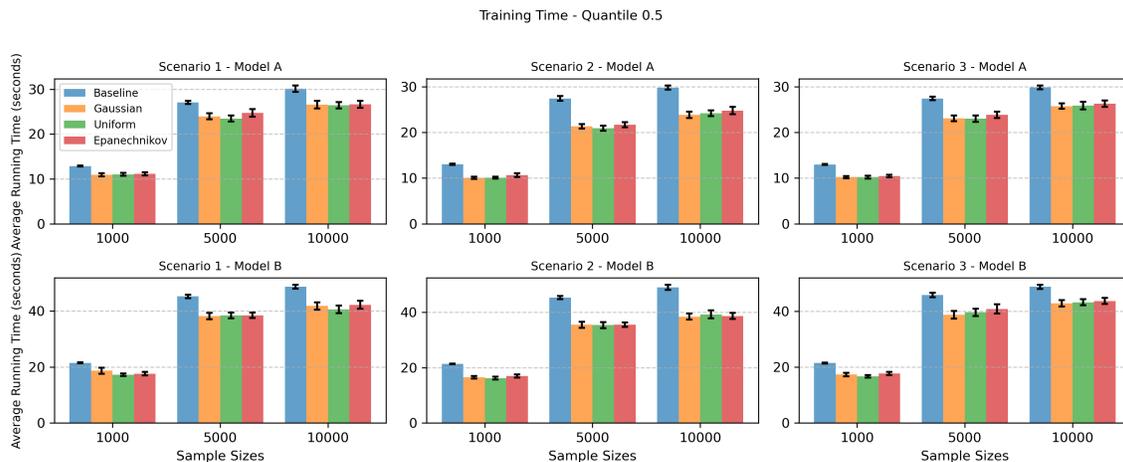


Figure 6: Bar chart with error bars with average training time over 50 trials under quantile level  $\tau = 0.5$ . The error bars represent 95% confidence intervals of training time for each setting.

1457  
 1458  
 1459  
 1460  
 1461  
 1462  
 1463  
 1464  
 1465  
 1466  
 1467  
 1468  
 1469  
 1470  
 1471  
 1472  
 1473  
 1474  
 1475  
 1476  
 1477  
 1478  
 1479  
 1480  
 1481  
 1482  
 1483  
 1484  
 1485  
 1486  
 1487  
 1488  
 1489  
 1490  
 1491  
 1492  
 1493  
 1494  
 1495  
 1496  
 1497  
 1498  
 1499  
 1500  
 1501  
 1502  
 1503

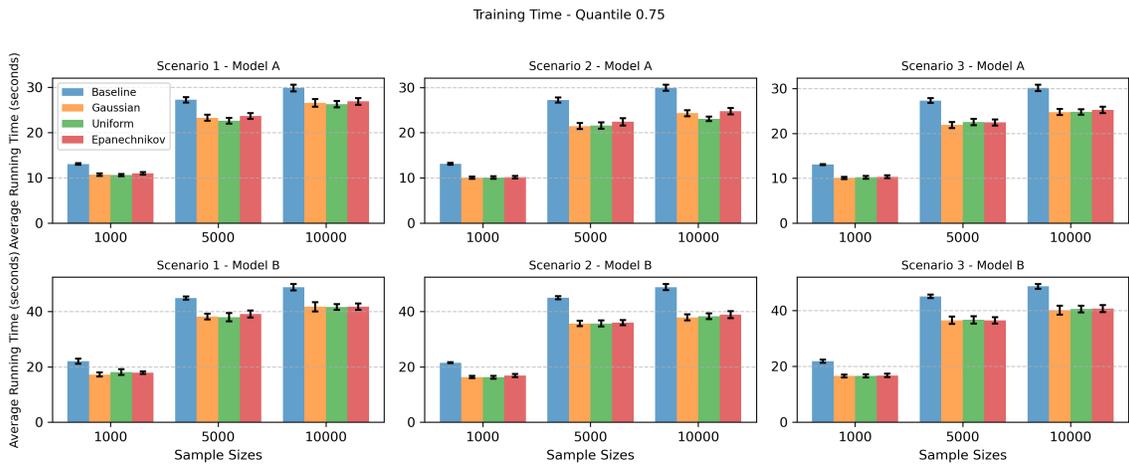


Figure 7: Bar chart with error bars with average training time over 50 trials under quantile level  $\tau = 0.75$ . The error bars represent 95% confidence intervals of training time for each setting.

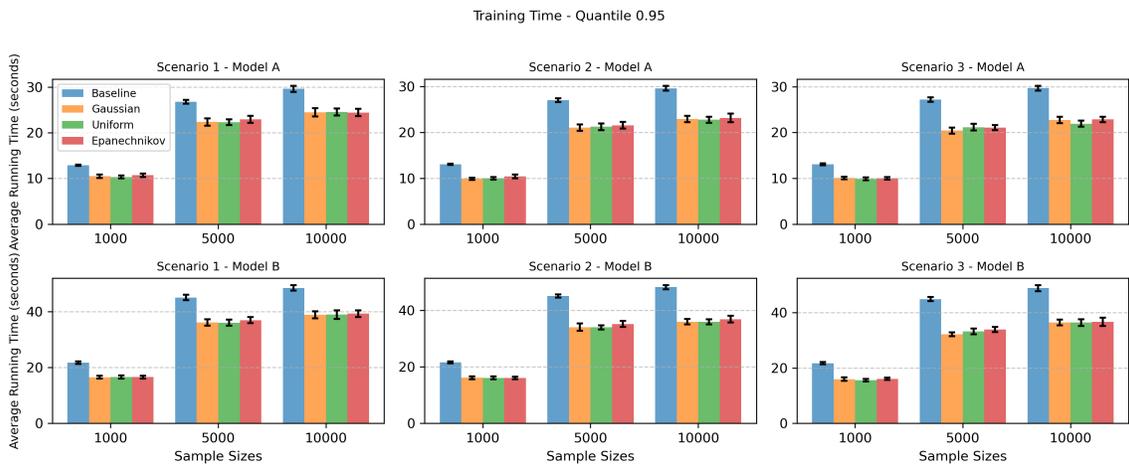


Figure 8: Bar chart with error bars with average training time over 50 trials under quantile level  $\tau = 0.95$ . The error bars represent 95% confidence intervals of training time for each setting.

Figures 2 and 5-8 show that our conquer networks require 20% less training time compared to the baseline method. Within the same scenario, model and sample size, the training times between 3 different kernels and between 5 different quantile levels are close to each other. We can also conclude that shallower networks have shorter training times when the number of parameters is similar.

## C DISCUSSION

Our smoothing principle can potentially be extended beyond quantile regression to other distributional objectives such as CRPS and Wasserstein distances. As discussed in Berrisch & Ziel (2023), CRPS can be expressed as an integral of quantile losses over  $\tau \in [0, 1]$ . Therefore, one can consider the smoothed objective

$$\min_{f_\tau \in \mathcal{F}_{DNN}, \tau \in [0,1]} \frac{1}{n} \sum_{i=1}^n \int_0^1 \ell_h(y_i - f_\tau(x_i)) d\tau,$$

where  $\ell_h$  is our conquer smoothed quantile loss.

Regarding the potential extension to Wasserstein-based objectives, consider the empirical Wasserstein- $p$  distance between the predictive CDF  $F_f(\cdot|x_i)$  ( $F_f^{-1}(\tau|x_i) = f_\tau(x_i)$ ) and the target CDF  $G(\cdot|x_i)$

$$W_p(F_f, G) = \sum_{i=1}^n \left( \int_0^1 |f_\tau(x_i) - G^{-1}(\tau|x_i)|^p d\tau \right)^{1/p}.$$

It is important to note that directly minimizing a Wasserstein distance is not straightforward in our conditional quantile regression, since the target  $G^{-1}(\tau|x_i)$  is unknown. One alternative is to use the empirical distribution  $G_n$  as the target, i.e.,

$$\hat{W}_p(F_f, G_n) = \sum_{i=1}^n \left( \int_0^1 |f_\tau(x_i) - G_n^{-1}(\tau|x_i)|^p d\tau \right)^{1/p},$$

where

$$G_n^{-1}(\tau|x_i) = \inf \left\{ y : \frac{1}{|\mathcal{N}_i|} \sum_{j \in \mathcal{N}_i} \mathbf{1}_{\{y_j \leq y\}} \geq \tau \right\}, \quad \mathcal{N}_i = \{j : x_j = x_i\}.$$

However, this requires repeated observations  $(y_j, x_j)$  at the same covariate  $x_j = x_i$ . Such a setting arises in reinforcement learning (Dabney et al., 2018), where multiple returns  $y_i$  for a fixed state  $x_i$  can be observed under different rollouts. In Bayesian settings, Zhang et al. (2020) chain quantile regression with Wasserstein-based objectives for posterior inference. Some existing literature provides partial connections between the quantile loss and Wasserstein objectives. Lheritier & Bondoux (2022) discussed that optimization based on the quantile loss can be interpreted as a 1-Wasserstein projection in certain settings, while Yang & Wang (2024) introduces a Wasserstein-improved objective for composite quantile regression. While a direct application is beyond the scope of our current work, these studies highlight potential directions for future research.

## D THE USE OF LARGE LANGUAGE MODELS (LLMs)

We used large language models (LLMs) solely to polish the writing and improve grammar and readability. All theoretical results, proofs, methodological developments, and simulation studies were carried out entirely by the authors without the aid of LLMs.