

Sparse Composite Quantile Regression in Ultrahigh Dimensions with Tuning Parameter Calibration

Yuwen Gu, Hui Zou

Abstract

When estimating coefficients in a linear model, the (sparse) composite quantile regression was first proposed in Zou and Yuan (2008) as an efficient alternative to the (sparse) least squares regardless of the error distribution. The highly nonsmooth nature of the composite loss in the sparse composite quantile regression makes its theoretical analysis as well as numerical computation much more challenging than the least squares method. The theory in Zou and Yuan (2008) was proven under fixed-dimension asymptotics and the estimator was computed via linear programming. None of these can be extended to ultrahigh dimensions. In this paper, we study the sparse composite quantile regression under ultrahigh dimensionality and make three contributions. First, we provide a non-asymptotic analysis of both the lasso and the folded concave penalized composite quantile regression, which reveals a practical way of achieving the oracle estimator. Second, we construct a novel information criterion for selecting the regularization parameter in the folded concave penalized composite quantile regression and prove its selection consistency. Third, we exploit the structure of the composite loss and design a specialized optimization algorithm for computing the penalized composite quantile regression via the alternating direction method of multipliers. We conduct extensive simulations to illustrate the theoretical results. Our analysis provides a unified treatment of the concentration inequalities involving the composite loss. Those inequalities could be of independent interest.

I. INTRODUCTION

Coefficient estimation in linear models is routinely done via the least squares (LS) regression. Under Gaussian errors, the LS estimator has the likelihood interpretation and is most efficient. It is reasonably efficient under other light-tailed error distributions besides Gaussian. When the error distribution is heavy-tailed, the LS estimator may fail to be consistent. See numerical studies in Section VI for a clear demonstration. The quantile regression (QR, [1]) can consistently estimate the coefficients of a linear model under very heavy-tailed errors, like a Student's t (with few degrees of freedom) or Cauchy. The robustness of the QR estimator, a property often mentioned in the literature, comes from the fact that its asymptotic variance does not depend on the moments of the error distribution, upon which that of the LS estimator relies, however. In terms of efficiency, it is well known that the asymptotic

Y. Gu is with the Department of Statistics, University of Connecticut, Storrs, CT, USA (e-mail: yuwen.gu@uconn.edu).

H. Zou is with the School of Statistics, University of Minnesota, Minneapolis, MN, USA (e-mail: zouxx019@umn.edu).

Manuscript received XXXXXXXXXX; revised XXXXXXXXXX.

variance of a QR estimator is inversely proportional to (the square of) the error density evaluated at the true quantile of the error distribution ([2], [3]). Thus, under certain errors, it is expected that the QR estimator is more efficient than the LS estimator. Nevertheless, the quantile regression considers only one quantile at a time and may not fully grasp the distributional information to always produce efficient estimation. To its extreme, when the error density at a specified quantile approaches zero, the asymptotic variance of the corresponding QR estimator explodes to infinity, which results in an estimator having arbitrarily small efficiency. As an example, under the mixture normal error $\frac{1}{2}N(-3, 1) + \frac{1}{2}N(3, 1)$, the least absolute deviation estimator is 1272.8 times less efficient than the LS estimator.

To safeguard quantile regression against potential efficiency loss, methods based on the idea of combining quantile information across multiple levels have been proposed in the literature. The idea is natural: as more quantiles are used, we have more distributional information to dispense and can hence obtain more efficient estimation if we do it properly. One such approach named the composite quantile regression (CQR, [4], [5]) combines information over different quantiles via a mix of quantile loss functions. It was shown by [4] that the CQR estimator is much more (or arbitrarily more) efficient than the LS estimator under many heavy-tailed errors. Another notable approach by [6] seeks an optimal weighting scheme to combine QR estimators at given levels to achieve as much efficiency gain as possible. It was shown that as the number of quantiles increases, the asymptotic variance of their proposed estimator achieves the Cramér–Rao lower bound under certain regularity conditions.

When considering fitting a sparse CQR model, it is natural to adopt the sparse penalties used in the sparse LS. In [4], Zou and Yuan studied the sparse CQR using the adaptive lasso penalty [7] and proved its oracle properties under fixed-dimension asymptotics. We note that the approach by [6] cannot be easily regularized to obtain desired sparse solutions since their estimator is based on a weighted average of multiple estimators, each of whose sparsity patterns may be different.

Given the favorable theoretical properties of CQR under fixed dimensions, we expect the sparse CQR to also enjoy very competitive performance under the ultrahigh dimensional setting. However, there are few results in the literature to firmly establish such a claim, despite the massive literature on the sparse LS under ultrahigh dimensions. This is mainly caused by the severe nonsmoothness of the composite quantile loss. For example, even with a single quantile, the analysis of the lasso penalized QR was only recently done in [8], and the analysis therein is technically very different from the standard analysis for the least squares lasso. CQR uses the sum of many different quantile losses and hence makes it even more challenging to handle than the single QR estimator. The highly nonsmooth nature of the composite quantile loss is also a major obstacle for using standard algorithms for the penalized LS as its numeric solvers. In fact, coordinate descent, the most popular algorithm thus far for solving the least squares lasso, is not suitable for optimization problems involving a nonsmooth loss.

The contributions of this article are as follows. Firstly, we provide nonasymptotic analysis of both lasso and folded concave penalized CQR. Our analysis holds for very general fixed pair (n, p) . As (n, p) go to infinity, we prove that the lasso estimator is estimation consistent under ultrahigh dimensions. Moreover, we show that the lasso estimator is tuning free, meaning that the rate of convergence is achieved by using an explicit penalization parameter. Secondly, we establish the oracle property for the folded concave penalized CQR and construct a new information criterion for calibrating the tuning parameter therein to give consistent model selection. Our paper

demonstrates a unified treatment of the concentration inequalities involving the CQR loss. Those inequalities could be of independent interest to studies on other models involving the check loss. Lastly, we exploit the structure of the composite quantile loss and design a specialized ADMM algorithm for efficiently computing the sparse penalized CQR estimator. The results in this work make sparse CQR a real alternative to the sparse LS for real applications.

The rest of the article is organized as follows. In Section II, we introduce the framework for the penalized CQR, followed by a discussion of the theoretical properties of the lasso and folded concave penalized CQR in Section III. We propose a new information criterion for selecting the tuning parameter and investigate its selection consistency in Section IV. In Section V, we present the efficient algorithm to solve the penalized CQR. Numerical studies are conducted in Section VI to show the superior finite-sample performance of penalized CQR over penalized LS. All proofs are relegated to Section VIII.

II. PENALIZED COMPOSITE QUANTILE REGRESSION

Consider variable selection and coefficient estimation in the linear model

$$y = \beta_0 + \sum_{j=1}^p x_j \beta_j + \varepsilon, \quad (1)$$

where ε is independent of $\mathbf{x} = (x_1, \dots, x_p)^\top$. Suppose β_0^* and $\boldsymbol{\beta}^* = (\beta_1^*, \dots, \beta_p^*)^\top$ are the true coefficients in model (1) that generate our i.i.d. data $(\mathbf{x}_i, y_i)_{i=1}^n$, where $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$. Denote the response vector by $\mathbf{y} = (y_1, \dots, y_n)^\top$ and the design matrix by $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$. We also write $\mathbf{X} = (X_1, \dots, X_p)$, where $X_j = (x_{1j}, \dots, x_{nj})^\top, 1 \leq j \leq p$. Let $\mathbb{X} = (X_0, \mathbf{X})$ be the augmented design with $X_0 = \mathbf{1}_n$ (corresponding to an intercept term), where $\mathbf{1}_n$ stands for the n -dimensional vector of all ones.

As mentioned in Section I, we consider the CQR rather than the LS or QR to estimate $\boldsymbol{\beta}$ in model (1). Assume that the random error ε has cumulative distribution function $F(\cdot)$ and probability density function $f(\cdot)$. To ensure identifiability of β_0 , assume $F(0) = 1/2$. Given an ordered sequence of quantile levels $\tau_1 < \tau_2 < \dots < \tau_K \in (0, 1)$, let $\alpha_k^* = \beta_0^* + F^{-1}(\tau_k)$, where $F^{-1}(\tau_k) = \inf\{x : F(x) \geq \tau_k\}$ denotes the τ_k -th quantile of $\varepsilon, 1 \leq k \leq K$. The canonical composite quantile regression estimates $\boldsymbol{\beta}$ by minimizing

$$\sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^\top \boldsymbol{\beta}) \quad (2)$$

jointly over $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_K) \in \mathbb{R}^K$ and $\boldsymbol{\beta} \in \mathbb{R}^p$, where $\rho_{\tau_k}(u) = \{\tau_k - I(u < 0)\}u$ denotes the check loss for $1 \leq k \leq K$. A typical choice is to take equally spaced τ_k 's: $\tau_k = k/(K+1), 1 \leq k \leq K$. As $K \rightarrow \infty$, [4] showed that the asymptotic efficiency of the CQR estimator relative to the LS estimator has a universal lower bound, $12\text{var}(\varepsilon)\{\mathbb{E}_\varepsilon f(\varepsilon)\}^2$, which is at least 70.26% for an arbitrary error distribution and can be made arbitrarily large for non-normal distributions. The relative efficiency lower bound 0.7026 is further improved to 0.864 in [9]. Substantial efficiency gain can be achieved already with a relatively small K such as $K = 9$ or 19.

In the high-dimensional regime, the number of parameters p is typically large and may even exceed the number of observations n . Under the sparsity assumption on the model, assume that many components of $\boldsymbol{\beta}^*$ are zero. Let

$\mathcal{A} = \{1 \leq j \leq p: \beta_j^* \neq 0\}$ be the active set of β^* and denote the effective dimensionality of the model by $s = |\mathcal{A}|$.

To harness the sparsity structure of β^* , let us consider the sparse penalized CQR

$$\min_{\alpha, \beta} \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta) + P_\lambda(\beta), \quad (3)$$

where $P_\lambda(\cdot)$ is a penalty function with regularization parameter λ . For instance, $P_\lambda(\cdot)$ can be the lasso [10], SCAD [11], MCP [12], and so on. In [4], the adaptive lasso was used to show the oracle property of the corresponding penalized estimator under fixed dimensions. The techniques used therein cannot be used to handle the ultrahigh dimensionality setting. As for the computation, the adaptive lasso penalized CQR was formulated as a linear program and was solved by a standard linear programming solver. However, such an approach does not scale well with dimensions. Other efficient alternatives are needed when p is large.

III. ANALYSIS OF PENALIZED COMPOSITE QUANTILE REGRESSION

In this section, we show the theoretical properties of the penalized CQR under ultrahigh dimensionality with both lasso and folded concave penalties. All our results are nonasymptotic and holds for general (n, p) . From these results, we establish the rate of convergence of the lasso penalized CQR and show its tuning free property. We also establish the strong oracle property for a feasible solution of the folded concave penalized CQR by incorporating the lasso estimator as initial estimation. For ease of exposition, we introduce the following notation.

For $u \in \mathbb{R}$, let $u_+ = uI(u > 0)$ and $u_- = -uI(u < 0)$ be the positive and negative parts of u , respectively. Moreover, let $\text{sgn}(u) = I(u > 0) - I(u < 0)$ be the sign function. The largest and smallest eigenvalues of a symmetric matrix \mathbf{A} are denoted by $\Lambda_{\max}(\mathbf{A})$ and $\Lambda_{\min}(\mathbf{A})$, respectively. We also let ∂g be the subdifferential of a convex function g . For two matrices $\mathbf{A}_1, \mathbf{A}_2 \in \mathbb{R}^{m \times n}$, let $\langle \mathbf{A}_1, \mathbf{A}_2 \rangle = \text{tr}(\mathbf{A}_1^T \mathbf{A}_2)$ be their trace inner product and $\|\mathbf{A}_1\|_F = \langle \mathbf{A}_1, \mathbf{A}_1 \rangle^{1/2}$ be the Frobenius norm of \mathbf{A}_1 . For any vector $\mathbf{v} = (v_1, \dots, v_p)^T \in \mathbb{R}^p$ and an arbitrary index set $I \subset \{1, \dots, p\}$, we write $\mathbf{v}_I = (v_j, j \in I)^T$ and denote by $\mathbf{X}_I = (X_j, j \in I)$ the submatrix consisting of the columns of \mathbf{X} whose indices are in I . The complement of I is denoted by $I^c = \{1, \dots, p\} \setminus I$. For $q \in [1, \infty]$, the L_q -norm of \mathbf{v} is denoted by $\|\mathbf{v}\|_q$.

A. Lasso penalized composite quantile regression

For $\lambda > 0$, we define the lasso penalized CQR estimator as

$$(\hat{\alpha}_\lambda, \hat{\beta}_\lambda) := \arg \min_{\alpha, \beta} Q_n(\alpha, \beta) + \lambda \sum_{j=1}^p |\beta_j|, \quad (4)$$

where $Q_n(\alpha, \beta) = (nK)^{-1} \sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta)$. In the sequel, we refer to $(\hat{\alpha}_\lambda, \hat{\beta}_\lambda)$ as the CQR lasso estimator.

For $\Delta \in \mathbb{R}^p$ and integer $m \geq 0$, let $\overline{\mathcal{A}}(\Delta, m) \subset \mathcal{A}^c$ be the support of the m largest in absolute value components of $\Delta_{\mathcal{A}^c}$. When $m = 0$, we take $\overline{\mathcal{A}}(\Delta, m)$ to be the empty set. The following assumption is imposed on the data and error distribution, which is typical in the QR literature.

(C0). The observations $(\mathbf{x}_i, y_i)_{i=1}^n$ are i.i.d. with $\min(n, p) \geq 3$. The density function is continuously differentiable and satisfies $f(u) \leq \bar{f} < \infty$ and $f'(u) \leq \bar{f}' \in (0, \infty)$ for all u in the support of ε . Moreover, there exists a

constant $\mathcal{U}_0 > 0$ such that $f(\alpha_k^* + u) \geq \underline{f} > 0$ for all $1 \leq k \leq K$ and $|u| \leq \mathcal{U}_0$. Also, $(\mathbf{x}_{i,\mathcal{A}}, y_i)_{i=1}^n$ are in general positions (Section 2.2, [3]) and there is at least one continuous covariate in the true model.

Note that we do not impose any moment or light tail assumptions on the error distribution and the assumptions on the error density are mild and can be satisfied by many commonly seen distributions, including heavy-tailed distributions like Cauchy. We also assume that \bar{f} , \underline{f}' and \underline{f} are all positive constants. The assumptions on $(\mathbf{x}_{i,\mathcal{A}}, y_i)$'s ensure that the CQR oracle estimator (5) is unique. This is a fairly common assumption in the QR literature (see [3], [13]). More discussions of the CQR oracle estimator can be found in Section III-B and Section B of the appendix.

We assume two additional conditions to establish the estimation consistency of the CQR lasso estimator. For the sake of brevity, only fixed design is considered. Define the restricted set $\mathcal{C} = \{(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p : \|\Delta_{\mathcal{A}^c}\|_1 \leq 3\|\Delta_{\mathcal{A}}\|_1 + \frac{3}{K}\|\delta\|_1\}$. The two assumptions are both imposed on the design matrix.

(C1). The design matrix \mathbf{X} satisfies

$$\kappa_m = \inf_{(\delta, \Delta) \in \mathcal{C}, (\delta, \Delta) \neq \mathbf{0}} \frac{\sum_{k=1}^K \sum_{i=1}^n (\delta_k + \mathbf{x}_i^T \Delta)^2}{n(K\|\Delta_{\mathcal{A} \cup \overline{\mathcal{A}}}(\Delta, m)\|_2^2 + \|\delta\|_2^2)} > 0.$$

(C2). The design matrix \mathbf{X} satisfies

$$q = \frac{3}{8} \frac{\underline{f}^{3/2}}{\bar{f}'} \inf_{(\delta, \Delta) \in \mathcal{C}, (\delta, \Delta) \neq \mathbf{0}} \frac{[n^{-1} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^T \Delta)^2]^{3/2}}{n^{-1} \sum_{i=1}^n \sum_{k=1}^K |\delta_k + \mathbf{x}_i^T \Delta|^3} > 0.$$

Condition (C1) is an extension of the restricted identifiability property (RIP), also known as the restricted eigenvalue (RE) condition, to the case of the penalized CQR. RIP is a common assumption in the literature for sparse penalized regressions. For example, it is assumed in the penalized LS, Dantzig selector [14], [15] and penalized QR [8]. Condition (C2) is similar to the restricted nonlinearity assumption in [8]. The quantity q , referred to as the restricted nonlinear impact (RNI) coefficient by those authors, describes how well the CQR empirical loss function can be minorized by a quadratic function over the restricted set \mathcal{C} . We present in the following theorem the L_2 -risk bound for the CQR lasso estimator, from which the estimation consistency of the estimator follows.

Theorem 1. Under conditions (C0), (C1) and (C2), with probability at least $1 - p_1(\lambda)$, where

$$p_1(\lambda) = 2K \exp\left(-\frac{9n\lambda^2}{2}\right) + 2p \exp\left(-\frac{n\lambda^2}{2M_0}\right) + \exp\left\{-2M_0 \frac{s(1 + \log p)}{\kappa_0}\right\},$$

the CQR lasso estimator $(\hat{\alpha}_\lambda, \hat{\beta}_\lambda)$ satisfies

$$\|\hat{\alpha}_\lambda - \alpha^*\|_2 \leq \frac{8}{\underline{f}} \sqrt{\frac{K}{\kappa_m}} \left\{ 32 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1) + \lambda \sqrt{\frac{s}{\kappa_0}} \right\}$$

and for integer $m > 0$,

$$\|\hat{\beta}_\lambda - \beta^*\|_2 \leq \frac{8}{\underline{f} \sqrt{\kappa_m}} \sqrt{1 + \frac{18s}{m} + \frac{18}{m}} \cdot \left\{ 32 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1) + \lambda \sqrt{\frac{s}{\kappa_0}} \right\},$$

provided that the growth condition

$$64 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1) + 2\lambda \sqrt{\frac{s}{\kappa_0}} \leq q \sqrt{\frac{\underline{f}}{K}}$$

holds, where $M_0 = \max_{0 \leq j \leq p} \|X_j\|_2^2/n$.

Remark 1. By Theorem 1, one can typically choose the tuning parameter $\lambda = C\sqrt{\log p/n}$ for the CQR lasso estimator, where $C > \sqrt{2M_0}$ is some constant. For example, one can choose $C = 2\sqrt{M_0}$. Note that given the design \mathbf{X} , M_0 can be readily obtained. Therefore, in principle, the parameter λ in the lasso penalized CQR is tuning free. This is in similar spirit to the square-root lasso [16]. With such choice of λ , we can see that $p_1(\lambda) = o(1)$ as $n, p \rightarrow \infty$, which leads to

$$\|\widehat{\beta}_\lambda - \beta^*\|_2 = \mathcal{O}_P\left(\frac{1}{\sqrt{\kappa_0 \kappa_s}} \sqrt{\frac{s \log p}{n}}\right)$$

provided $q^{-1} \sqrt{s \log p / (n \kappa_0)} = o(1)$ and $\kappa_0 (s \log p)^{-1} = o(1)$, by taking $m = s$. When κ_0 and κ_s are both positive constants, the CQR lasso estimator achieves the near-optimal rate $\sqrt{s \log p / n}$, which implies that it is a consistent estimator even when p is of exponential order of n , i.e., $\log p = \mathcal{O}(n^\gamma)$ for some $0 < \gamma < 1$, provided $s \log p = o(n)$.

B. Folded concave penalized composite quantile regression

Folded concave penalized regression has been widely adopted in the statistical analysis of high-dimensional data due to its strong oracle optimality [17], [18]. In order to establish the oracle property of the folded concave penalized CQR estimator, let us first define the CQR oracle estimator,

$$(\widehat{\alpha}^o, \widehat{\beta}^o) := \underset{\alpha, \beta: \beta_{\mathcal{A}^c} = \mathbf{0}}{\operatorname{argmin}} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta). \quad (5)$$

The oracle estimator $(\widehat{\alpha}^o, \widehat{\beta}^o)$ is the ideal estimator one could possibly get using the CQR. It is not feasible in practice since \mathcal{A} is unknown, but it serves as a benchmark estimator to which one can compare a penalized CQR estimator. In the following lemma, we show the rate of convergence of the CQR oracle estimator under the growing-dimension regime, i.e., the true dimensionality s is allowed to grow with n .

Let $\mathcal{A}_0 = \{0\} \cup \mathcal{A}$ and $\mathbb{X}_{\mathcal{A}_0} = (\mathbf{1}_n, \mathbf{X}_{\mathcal{A}})$. Denote $\underline{\mu} = \Lambda_{\min}(n^{-1} \mathbb{X}_{\mathcal{A}_0}^T \mathbb{X}_{\mathcal{A}_0})$ and $\bar{\mu} = \Lambda_{\max}(n^{-1} \mathbb{X}_{\mathcal{A}_0}^T \mathbb{X}_{\mathcal{A}_0})$. Moreover, let $M_{\mathcal{A}} = \max_{1 \leq i \leq n} (s+1)^{-1} (1 + \|\mathbf{x}_{i\mathcal{A}}\|_2^2)$ and $M_{\mathcal{A}^c} = \max_{1 \leq i \leq n, j \in \mathcal{A}^c} |x_{ij}|$. In this article, we assume that $M_{\mathcal{A}}$ and $M_{\mathcal{A}^c}$ are both positive constants.

Lemma 1. *Under condition (C0), if $32K(s+1)M_{\mathcal{A}}/(\sqrt{n}f\underline{\mu}) \leq \mathcal{U}_0$, with probability at least $1 - \exp[-(s+1)M_{\mathcal{A}}/(2\bar{\mu})]$, the CQR oracle estimator satisfies*

$$\|\widehat{\alpha}^o - \alpha^*\|_2^2 + \|\widehat{\beta}^o - \beta^*\|_2^2 \leq \frac{1024K^2(s+1)M_{\mathcal{A}}}{nf^2\underline{\mu}^2}.$$

Remark 2. Assuming $s/(\underline{\mu}\sqrt{n}) = o(1)$ and $\underline{\mu}/s = o(1)$, the CQR oracle estimator has the following rate of convergence

$$\|\widehat{\alpha}^o - \alpha^*\|_2 = \mathcal{O}_P\left(\frac{1}{\underline{\mu}} \sqrt{\frac{s}{n}}\right), \quad \|\widehat{\beta}^o - \beta^*\|_2 = \mathcal{O}_P\left(\frac{1}{\underline{\mu}} \sqrt{\frac{s}{n}}\right)$$

as $n \rightarrow \infty$. This implies that $\widehat{\beta}^o$ is $\sqrt{n/s}$ -consistent when s diverges with n , if we assume that $\underline{\mu} > 0$ is a fixed constant. Then it is required that $s = \mathcal{O}(n^\gamma)$ for some $0 < \gamma < \frac{1}{2}$.

Next, we introduce the details of the folded concave penalized CQR and show that the CQR oracle estimator is attainable via the folded concave penalized CQR. The folded concave penalized CQR at penalty level $\lambda > 0$ solves the following minimization problem

$$\min_{\alpha, \beta} Q_n(\alpha, \beta) + \sum_{j=1}^p p_\lambda(|\beta_j|), \quad (6)$$

where $p_\lambda(t), t \geq 0$ belongs to a class of folded concave penalties that satisfy the following properties:

- (P1) $p_\lambda(t)$ is nondecreasing and concave in $t \geq 0$ and $p_\lambda(0) = 0$;
- (P2) $p_\lambda(t)$ is differentiable in $t > 0$;
- (P3) $p'_\lambda(t) \geq a_1\lambda, 0 < t \leq a_2\lambda$ and $p'_\lambda(0) := p'_\lambda(0+) \geq a_1\lambda$, where $a_1, a_2 > 0$ are fixed constants;
- (P4) $p'_\lambda(t) = 0, t \geq a\lambda$ for a fixed constant $a > a_2$.

It can be shown that both the SCAD penalty and MCP belong to this class (see, e.g., [19]). For the analysis of the minimizer, we consider the local linear approximation (LLA, [20]) algorithm, where the initial estimator is chosen to be the CQR lasso estimator.

- 1) Initialize α and β with respectively $\hat{\alpha}^{(0)}$ and $\hat{\beta}^{(0)}$. Compute weights

$$\hat{w}_j^{(0)} = p'_\lambda(|\hat{\beta}_j^{(0)}|), j = 1, \dots, p.$$

- 2) For $m = 1, 2, \dots$, repeat the LLA iterations in (2.a) and (2.b).

- (2.a) Solve the following convex optimization problem for $\hat{\alpha}^{(m)}$ and $\hat{\beta}^{(m)}$

$$(\hat{\alpha}^{(m)}, \hat{\beta}^{(m)}) := \arg \min_{\alpha, \beta} Q_n(\alpha, \beta) + \sum_{j=1}^p \hat{w}_j^{(m-1)} |\beta_j|.$$

- (2.b) Calculate the weights

$$\hat{w}_j^{(m)} = p'_\lambda(|\hat{\beta}_j^{(m)}|), j = 1, \dots, p.$$

In order to establish the oracle property, we assume a “beta-min” condition, i.e., the true coefficient exhibit sufficient signal:

$$(C3) \min_{j \in \mathcal{A}} |\beta_j^*| > (a+1)\lambda.$$

The “beta-min” condition is always assumed for non-convexly penalized regressions and is almost a necessary condition for establishing consistency results. See, e.g., [11], [13], [18].

Theorem 2. *Suppose the folded concave penalized CQR (6) is solved with the LLA algorithm that is initialized with the CQR lasso estimator (4) at penalty level $\lambda_0 = c \sqrt{\log p/n}$ for some constant $c > \sqrt{2M_0}$. Let $r_0 = \min_{j \in \mathcal{A}} |\beta_j^*| - a\lambda$ and $r_* = \sqrt{(s+1)M_{\mathcal{A}} \log n/n}$. Assume the folded concave penalty $p_\lambda(\cdot)$ satisfies properties (P1) – (P4), where for integer $m > 0$, λ is taken such that*

$$\lambda \geq \frac{8}{a_0 f \sqrt{\kappa_m}} \sqrt{1 + \frac{18s}{m} + \frac{18}{m}} \cdot \left\{ 32 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1) + \lambda_0 \sqrt{\frac{s}{\kappa_0}} \right\}. \quad (7)$$

Under conditions (C0) – (C3) and the assumptions that $r_0 \sqrt{(s+1)M_{\mathcal{A}}^c} \leq \mathcal{U}_0$, $r_* \sqrt{(s+1)M_{\mathcal{A}}^c} \leq \mathcal{U}_0$ and $\lambda > 8K(\underline{f}\underline{\mu})^{-1} \sqrt{(s+1)M_{\mathcal{A}}^c/n}$, with probability at least $p_0 = 1 - p_1(\lambda_0) - p_2(r_0) - p_2(r_*) - p_3$, the LLA algorithm converges to the oracle estimator $(\hat{\alpha}^o, \hat{\beta}^o)$ in two iterations, where $p_1(\cdot)$ is given in Theorem 1, $p_2(\cdot)$ is defined as

$$p_2(r) = \exp\left\{-\frac{n(t(r))^2}{32\bar{\mu}r^2}\right\}, \text{ where } t(r) = \frac{\underline{f}}{4K}\bar{\mu}r^2 - 2r\sqrt{\frac{(s+1)M_{\mathcal{A}}^c}{n}},$$

and

$$p_3 = 2(p-s)\exp\left(-\frac{2nB^2}{M_0}\right) + 2(p-s)n^{2(K+s)}\exp\left(-\frac{3nB^2}{24\bar{f}M_{\mathcal{A}}^c\bar{\mu}^{1/2}r_* + 8M_{\mathcal{A}}^cB}\right) \\ + 2(p-s)n^{2(K+s)}\exp\left(-\frac{3nB_0^2}{24\bar{f}M_{\mathcal{A}}^c(s+1)^{1/2}M_{\mathcal{A}}^{1/2}n^{-2}r_* + 4M_{\mathcal{A}}^cB_0}\right)$$

with $B = \frac{1}{2}[a_1\lambda - \frac{K+s}{n}M_{\mathcal{A}}^c - \bar{f}M_{\mathcal{A}}^c\bar{\mu}^{1/2}r_*]_+$ and $B_0 = [\frac{B}{2} - \frac{8\bar{f}M_{\mathcal{A}}^c\sqrt{(s+1)M_{\mathcal{A}}^c}}{n^2}r_*]_+$.

It is easy to translate Theorem 2 into an asymptotic statement that the folded concave penalized CQR estimator finds the oracle CQR estimator with overwhelming probability. For brevity, we omit such discussions. We emphasize that unlike the lasso penalized CQR, the regularization parameter in the folded concave penalized CQR involves unknown quantities. In order to apply Theorem 2 in applications, we need a data-driven choice of the regularization parameter. To this end, we construct a new information criterion for selecting the tuning parameter in the next section.

IV. TUNING PARAMETER CALIBRATION

For the folded concave penalized CQR, there exists a tuning sequence λ_n (we write λ_n here to signify its dependence on n) such that the LLA algorithm yields the CQR oracle estimator in two iterations with probability approaching one (Theorem 2). However, as we pointed out already, there is no direct way to use such λ_n as given in Theorem 2, since it relies on unknown quantities. We thus pursue data-driven approach to the selection of λ . Consider the following high-dimensional Bayesian information criterion:

$$\text{BIC}^H(\lambda) = \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \rho_{v_k}(y_i - \hat{\alpha}_k^\lambda - \mathbf{x}_i^\top \hat{\beta}^\lambda) + |\hat{A}_\lambda| \frac{C_n \log(p)}{n}, \quad (8)$$

where $(\hat{\alpha}^\lambda, \hat{\beta}^\lambda)$ is the two-step estimator from the LLA algorithm initialized by the CQR lasso estimator with regularization parameter λ_0 (Theorem 2), $\hat{A}_\lambda = \{1 \leq j \leq p : \hat{\beta}_j^\lambda \neq 0\}$ and C_n is a positive number depending on n (allowed to grow with n). We compare the values of $\text{BIC}^H(\lambda)$ for $\lambda \in \Xi_n = \{\lambda : |\hat{A}_\lambda| \leq J_n\}$, where $J_n > s$ represents a rough estimate of the upper bound of the model sparsity and is allowed to (slowly) diverge as $n \rightarrow \infty$. Typically, J_n is much smaller than p , so that one can avoid searching over a notoriously large model space. The tuning parameter selected via BIC is given by

$$\hat{\lambda}_n = \arg \min_{\lambda \in \Xi_n} \text{BIC}^H(\lambda).$$

We call the information criterion BIC because it can be shown in the following Theorem that $\Pr(\hat{A}_{\hat{\lambda}_n} = \mathcal{A}) \rightarrow 1$ as $n \rightarrow \infty$. Note that model selection consistency is the signature property of BIC in the fixed-dimension setting.

Let $M = \max_{1 \leq i \leq n, 0 \leq j \leq p} |x_{ij}|$ and assume M is a positive constant. Also, define

$$\underline{\zeta} = \inf_{A \supset \mathcal{A}, |A| \leq 2J_n} \Lambda_{\min}(n^{-1}(\mathbf{1}_n, \mathbf{X}_A)^T(\mathbf{1}_n, \mathbf{X}_A))$$

and

$$\bar{\zeta} = \sup_{A \supset \mathcal{A}, |A| \leq 2J_n} \Lambda_{\max}(n^{-1}(\mathbf{1}_n, \mathbf{X}_A)^T(\mathbf{1}_n, \mathbf{X}_A)),$$

and assume that both $\underline{\zeta}$ and $\bar{\zeta}$ are positive constants, and so are $\kappa_0, \kappa_s, \underline{\mu}$ and $\bar{\mu}$.

Theorem 3. *Under the conditions of Theorem 2, and as $n \rightarrow \infty$, assuming that $s = \mathcal{O}(1), n = \mathcal{O}(p), J_n^2 \log p = \mathcal{O}(n), C_n \rightarrow \infty$ and $\sqrt{\max(J_n, C_n) \log p / n} = \mathcal{O}(\min_{j \in \mathcal{A}} |\beta_j|)$, the criterion $\text{BIC}^H(\lambda)$ is selection consistent, i.e., $\Pr(\hat{A}_{\hat{\lambda}_n} = \mathcal{A}) \rightarrow 1$ as $n \rightarrow \infty$.*

Remark 3. The sequence C_n is often taken to slowly diverge to infinity, e.g., $C_n = \log \log n$. Under fixed model sparsity, $s = \mathcal{O}(1)$, it is implied from Theorem 3 that $\text{BIC}^H(\lambda)$ is consistent when $\log p = \mathcal{O}(n^{\gamma_1})$ and $J_n = \mathcal{O}(n^{\gamma_2})$ for some positive constants γ_1 and γ_2 such that $\gamma_1 + 2\gamma_2 < 1$. It is worth mentioning that the fixed model sparsity is assumed in order to achieve ultrahigh dimensionality due to technical difficulties with the check loss (see, e.g., [21]). If instead we allow the model sparsity to grow, using current technique, we must assume p can be of at most polynomial order of n (see, e.g., [22]).

One problem with the criterion in (8) is that it is not scale invariant. In practice, one needs to standardize the variables beforehand. Therefore, we also consider a scale invariant version of the high-dimensional Bayesian information criterion:

$$\text{BIC}^{\text{HL}}(\lambda) = \log \left(\frac{1}{K} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \hat{\alpha}_k^\lambda - \mathbf{x}_i^T \hat{\beta}^\lambda) \right) + |\hat{A}_\lambda| \frac{C_n \log(p)}{n}. \quad (9)$$

Theorem 4. *In model (1), assume that $\mathbb{E}(|\varepsilon|) < \infty$. Under the conditions of Theorem 2, as $n \rightarrow \infty$, assume moreover that $s = \mathcal{O}(1), n = \mathcal{O}(p), J_n^2 \log p = \mathcal{O}(n), C_n \rightarrow \infty$ and $\sqrt{\max(J_n, C_n) \log p / n} = \mathcal{O}(\min_{j \in \mathcal{A}} |\beta_j|)$. The criterion $\text{BIC}^{\text{HL}}(\lambda)$ is selection consistent, i.e., $\Pr(\hat{A}_{\hat{\lambda}_n} = \mathcal{A}) \rightarrow 1$ as $n \rightarrow \infty$.*

Remark 4. The selection consistency of $\text{BIC}^{\text{HL}}(\lambda)$ requires additionally that $\mathbb{E}(|\varepsilon|) < \infty$. Our empirical study in Section VI suggests that this might be a necessary condition. Indeed, in the numerical comparison there, we see that $\text{BIC}^{\text{HL}}(\lambda)$ does not perform well under the Cauchy error.

V. OPTIMIZATION

Note that both lasso and folded concave penalized CQR can be solved with one or more runs of the following weighted lasso penalized CQR:

$$\min_{\alpha, \beta} \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta) + \lambda_1 \sum_{j=1}^p d_j |\beta_j|, \quad (10)$$

where $\lambda_1 > 0$ and $d_j \geq 0$ for $j = 1, \dots, p$. Specifically, the CQR lasso estimator can be achieved by letting $d_j = 1$ for all $j = 1, \dots, p$, while the folded concave penalized CQR estimator can be obtained by iteratively solving (10)

with $d_j = \hat{w}_j^{(m-1)}$ in the m th LLA iteration. Hence, in the sequel, we only need to focus on developing the algorithm for solving (10).

Traditionally, (10) can be solved by linear programming if n and p are moderate. However, linear programming does not scale well when p is large [23]. We hence propose an efficient alternating direction method of multipliers (ADMM) algorithm for solving (10). The algorithm is based on a reformulation that turns the original problem into one that can harness the power of ADMM. We point out that there are multiple ways to formulate (10) into problems that are solvable by ADMM. For instance, in a relevant context, two ADMM versions are proposed in [23] to efficiently solve the penalized QR and they can be readily modified to solve (10). However, the formulation we present here is different from the ideas in [23] and it results in a more stable ADMM algorithm. To elaborate on our algorithm, let $z_{ik} = y_i - \alpha_k - \mathbf{x}_i^T \boldsymbol{\beta}$ for $i = 1, \dots, n$ and $k = 1, \dots, K$, and define matrix $\mathbf{Z} = (z_{ik})_{n \times K}$ in terms of the z_{ik} 's. By convexity, it can be immediately seen that (10) is equivalent to

$$\begin{aligned} & \text{minimize} && \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(z_{ik}) + \lambda_1 \sum_{j=1}^p d_j |\gamma_j| \\ & \text{subject to} && \mathbf{Z} = \mathbf{1}_K^T \otimes \mathbf{y} - \mathbf{1}_n \otimes \boldsymbol{\alpha}^T - \mathbf{1}_K^T \otimes (\mathbf{X}\boldsymbol{\beta}) \\ & && \boldsymbol{\gamma} = \boldsymbol{\beta}, \end{aligned} \tag{11}$$

where $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)^T$ and \otimes denotes the Kronecker product. As will be seen, the introduction of the γ_j 's renders a more stable ADMM algorithm where only the dual updates involve the non-smooth functions. For ease of notation, let $\mathbf{Y} = \mathbf{1}_K \otimes \mathbf{y}$, $\boldsymbol{\varphi} = (\boldsymbol{\alpha}^T, \boldsymbol{\beta}^T)^T$,

$$\mathbb{X}_1 = \begin{pmatrix} \mathbf{1}_n & \cdots & \mathbf{0} & \mathbf{X} \\ \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \cdots & \mathbf{1}_n & \mathbf{X} \end{pmatrix}_{(nK) \times (p+K)}, \text{ and } \mathbb{X}_2 = (\mathbf{0}_{p \times K} \quad \mathbf{I}_p)_{p \times (p+K)}.$$

Then, (11) can be equivalently written as

$$\begin{aligned} & \text{minimize} && \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(z_{ik}) + \lambda_1 \sum_{j=1}^p d_j |\gamma_j| \\ & \text{subject to} && \begin{pmatrix} \mathbb{X}_1 \\ -\mathbb{X}_2 \end{pmatrix} \boldsymbol{\varphi} + \begin{pmatrix} \text{vec}(\mathbf{Z}) \\ \boldsymbol{\gamma} \end{pmatrix} = \begin{pmatrix} \mathbf{Y} \\ \mathbf{0} \end{pmatrix}, \end{aligned} \tag{12}$$

where vec stands for the vectorization operator that stacks the columns of a matrix one underneath the other to form a single vector. The augmented Lagrangian of problem (12) is

$$\begin{aligned} L_{\sigma}(\boldsymbol{\varphi}, \mathbf{Z}, \boldsymbol{\gamma}, \mathbf{U}, \mathbf{v}) & := \frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n \rho_{\tau_k}(z_{ik}) + \lambda_1 \sum_{j=1}^p d_j |\gamma_j| \\ & + \langle \text{vec}(\mathbf{U}), \text{vec}(\mathbf{Z}) + \mathbb{X}_1 \boldsymbol{\varphi} - \mathbf{Y} \rangle + \langle \mathbf{v}, \boldsymbol{\gamma} - \mathbb{X}_2 \boldsymbol{\varphi} \rangle \\ & + \frac{\sigma}{2} \|\text{vec}(\mathbf{Z}) + \mathbb{X}_1 \boldsymbol{\varphi} - \mathbf{Y}\|_F^2 + \frac{\sigma}{2} \|\boldsymbol{\gamma} - \mathbb{X}_2 \boldsymbol{\varphi}\|_2^2, \end{aligned} \tag{13}$$

where $\mathbf{U} = (u_{ik})_{n \times K}$ and $\mathbf{v} = (v_1, \dots, v_p)^T$ are the Lagrangian multipliers and $\sigma > 0$. Let $\boldsymbol{\varphi}^r, \mathbf{Z}^r, \boldsymbol{\gamma}^r, \mathbf{U}^r$ and \mathbf{v}^r be the iterate after the r th iteration of the algorithm, where $r \geq 0$. The ADMM has the following updates in the $(r+1)$ st

iteration

$$\begin{cases} \boldsymbol{\varphi}^{r+1} := \arg \min_{\boldsymbol{\varphi}} L_{\sigma}(\boldsymbol{\varphi}, \mathbf{Z}^r, \boldsymbol{\gamma}^r, \mathbf{U}^r, \mathbf{v}^r), \\ (\mathbf{Z}^{r+1}, \boldsymbol{\gamma}^{r+1}) := \arg \min_{\mathbf{Z}, \boldsymbol{\gamma}} L_{\sigma}(\boldsymbol{\varphi}^{r+1}, \mathbf{Z}, \boldsymbol{\gamma}, \mathbf{U}^r, \mathbf{v}^r), \\ \text{vec}(\mathbf{U}^{r+1}) := \text{vec}(\mathbf{U}^r) + \boldsymbol{\sigma}\{\text{vec}(\mathbf{Z}^{r+1}) + \mathbb{X}_1 \boldsymbol{\varphi}^{r+1} - \mathbf{Y}\}, \\ \mathbf{v}^{r+1} := \mathbf{v}^r + \boldsymbol{\sigma}\{\boldsymbol{\gamma}^{r+1} - \mathbb{X}_2 \boldsymbol{\varphi}^{r+1}\}. \end{cases} \quad (14)$$

It follows from (13) that

$$\boldsymbol{\varphi}^{r+1} = \frac{1}{\boldsymbol{\sigma}} (\mathbb{X}_1^T \mathbb{X}_1 + \mathbb{X}_2^T \mathbb{X}_2)^{-1} \{ \mathbb{X}_1^T (\boldsymbol{\sigma} \mathbf{Y} - \boldsymbol{\sigma} \text{vec}(\mathbf{Z}^r) - \text{vec}(\mathbf{U}^r)) + \mathbb{X}_2^T (\boldsymbol{\sigma} \boldsymbol{\gamma}^r + \mathbf{v}^r) \}.$$

Note that

$$\mathbb{X}_1^T \mathbb{X}_1 + \mathbb{X}_2^T \mathbb{X}_2 = \begin{pmatrix} n \mathbf{I}_K & \mathbf{1}_K \mathbf{1}_n^T \mathbf{X} \\ \mathbf{X}^T \mathbf{1}_n \mathbf{1}_K^T & \mathbf{I}_p + K \mathbf{X}^T \mathbf{X} \end{pmatrix}.$$

Let the Schur complement of $n \mathbf{I}_K$ in the above matrix be

$$\mathbf{S} = \mathbf{I}_p + K \mathbf{X}^T \mathbf{X} - \frac{1}{n} (\mathbf{X}^T \mathbf{1}_n \mathbf{1}_K^T) (\mathbf{1}_K \mathbf{1}_n^T \mathbf{X}) = \mathbf{I}_p + K \mathbf{X}_0^T \mathbf{X}_0,$$

where $\mathbf{X}_0 = (\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}_n^T) \mathbf{X}$ is the centered design matrix. Then, we have

$$(\mathbb{X}_1^T \mathbb{X}_1 + \mathbb{X}_2^T \mathbb{X}_2)^{-1} = \begin{pmatrix} \frac{1}{n} \mathbf{I}_K + \frac{1}{n^2} \mathbf{1}_K \mathbf{1}_n^T \mathbf{X} \mathbf{S}^{-1} \mathbf{X}^T \mathbf{1}_n \mathbf{1}_K^T & -\frac{1}{n} \mathbf{1}_K \mathbf{1}_n^T \mathbf{X} \mathbf{S}^{-1} \\ -\frac{1}{n} \mathbf{S}^{-1} \mathbf{X}^T \mathbf{1}_n \mathbf{1}_K^T & \mathbf{S}^{-1} \end{pmatrix}.$$

When p is large, the computation of \mathbf{S}^{-1} can be expensive. We can apply the Sherman–Morrison–Woodbury formula to get

$$\mathbf{S}^{-1} = \mathbf{I}_p - K \mathbf{X}_0^T (\mathbf{I}_n + K \mathbf{X}_0 \mathbf{X}_0^T)^{-1} \mathbf{X}_0,$$

where we only need to evaluate the inverse of an $n \times n$ matrix. When n is relatively small compared to p , this formula can be very helpful.

Remark 5. In the actual implementation, we often center the design matrix before fitting the model. Then, $\mathbb{X}_1^T \mathbb{X}_1 + \mathbb{X}_2^T \mathbb{X}_2$ is block diagonal since $\mathbf{X}^T \mathbf{1}_n = \mathbf{0}$ and its inverse can be readily obtained.

The update of \mathbf{Z}^{r+1} and $\boldsymbol{\gamma}^{r+1}$ can be carried out component-wisely. This pertains to the application of the proximity operator of the check loss $\rho_{\tau}(\cdot)$ and the absolute value function $|\cdot|$, respectively. For $v \in \mathbb{R}$, the proximity operator of $\rho_{\tau}(\cdot)$ with respect to a parameter $a > 0$ is defined as

$$\text{Prox}_{\rho_{\tau}}(v, a) := \arg \min_{u \in \mathbb{R}} \rho_{\tau}(u) + \frac{a}{2} (u - v)^2.$$

The following lemma gives the closed form expression of $\text{Prox}_{\rho_{\tau}}$.

Lemma 2. For $v \in \mathbb{R}$ and $a > 0$, the proximity operator of the check loss $\rho_{\tau}(\cdot)$ with respect to parameter a is given by

$$\text{Prox}_{\rho_{\tau}}(v, a) = v - \max\left(\frac{\tau - 1}{a}, \min\left(v, \frac{\tau}{a}\right)\right).$$

Now by Lemma 2, we obtain

$$z_{ik}^{r+1} = \text{Prox}_{\rho_{\tau_k}} \left(y_i - \alpha_k^{r+1} - \mathbf{x}_i^T \boldsymbol{\beta}^{r+1} - \frac{u_{ik}^r}{\sigma}, nK\sigma \right), 1 \leq i \leq n, 1 \leq k \leq K.$$

The proximity operator of $|\cdot|$ is the soft-thresholding operator and thus

$$\gamma_j^{r+1} = \text{Shrink} \left(\beta_j^{r+1} - \frac{v_j^r}{\sigma}, \frac{\lambda_1 d_j}{\sigma} \right),$$

where $\text{Shrink}(v, a) = \text{sgn}(v)(|v| - a)_+$.

We summarize the above ADMM algorithm in Algorithm 1. A discussion of the convergence criterion for this algorithm can be found in the appendix.

Algorithm 1: The ADMM algorithm for solving the weighted lasso penalized composite quantile regression

- 1) Initialize the algorithm with $(\boldsymbol{\varphi}^0, \mathbf{Z}^0, \boldsymbol{\gamma}^0, \mathbf{U}^0, \mathbf{v}^0)$, where $\boldsymbol{\varphi}^0 = ((\boldsymbol{\alpha}^0)^T, (\boldsymbol{\beta}^0)^T)^T$.
- 2) For $r = 0, 1, 2, \dots$, repeat steps (2.1) – (2.3) until convergence.

(2.1) Update

$$\begin{aligned} \boldsymbol{\varphi}^{r+1} = ((\boldsymbol{\alpha}^{r+1})^T, (\boldsymbol{\beta}^{r+1})^T)^T \leftarrow & \frac{1}{\sigma} (\mathbb{X}_1^T \mathbb{X}_1 + \mathbb{X}_2^T \mathbb{X}_2)^{-1} \\ & \cdot \{ \mathbb{X}_1^T (\boldsymbol{\sigma} \mathbf{Y} - \boldsymbol{\sigma} \text{vec}(\mathbf{Z}^r) - \text{vec}(\mathbf{U}^r)) + \mathbb{X}_2^T (\boldsymbol{\sigma} \boldsymbol{\gamma}^r + \mathbf{v}^r) \}. \end{aligned}$$

(2.2) Update

$$z_{ik}^{r+1} \leftarrow \text{Prox}_{\rho_{\tau_k}} \left(y_i - \alpha_k^{r+1} - \mathbf{x}_i^T \boldsymbol{\beta}^{r+1} - \frac{u_{ik}^r}{\sigma}, nK\sigma \right), 1 \leq i \leq n, 1 \leq k \leq K,$$

and

$$\gamma_j^{r+1} \leftarrow \text{Shrink} \left(\beta_j^{r+1} - \frac{v_j^r}{\sigma}, \frac{\lambda_1 d_j}{\sigma} \right), 1 \leq j \leq p.$$

(2.3) Update

$$\text{vec}(\mathbf{U}^{r+1}) \leftarrow \text{vec}(\mathbf{U}^r) + \boldsymbol{\sigma} \{ \text{vec}(\mathbf{Z}^{r+1}) + \mathbb{X}_1 \boldsymbol{\varphi}^{r+1} - \mathbf{Y} \}$$

and

$$\mathbf{v}^{r+1} \leftarrow \mathbf{v}^r + \boldsymbol{\sigma} \{ \boldsymbol{\gamma}^{r+1} - \mathbb{X}_2 \boldsymbol{\varphi}^{r+1} \}.$$

VI. NUMERICAL EXPERIMENTS

We conduct Monte Carlo studies to assess the finite sample performance of the proposed method as well as the tuning criterion. First, we compare the estimators from the penalized LS, the penalized CQR, the ideal oracle LS, and the oracle CQR. Recall that the oracle estimators are obtained through applying the canonical LS and CQR to the true underlying model. Second, we compare the tuned penalized CQR estimation by using cross-validation (CV) and by using the proposed information criteria.

Our simulated data are from the linear model

$$y = \beta_0^* + \mathbf{x}^T \boldsymbol{\beta}^* + \varepsilon, \quad (15)$$

where $\beta_0^* = 0$ and $\boldsymbol{\beta}^* = (3, 1.5, 0, 0, 2, \mathbf{0}_{p-5})^T$. The covariates are drawn from the multivariate normal distribution, $\mathbf{x} \sim N_p(\mathbf{0}, \boldsymbol{\Sigma})$, where two different covariance matrices $\boldsymbol{\Sigma} = (0.5^{|i-j|})$ and $\boldsymbol{\Sigma} = (0.8^{|i-j|})$ are considered. For the error distribution, we refer to [4] and consider five different shapes:

- (a) the normal distribution, $\varepsilon \sim N(0, 3)$;
- (b) the mixture normal distribution (MN), $\varepsilon \sim \sqrt{6} \times \varepsilon^*$, where $\varepsilon^* \sim 0.5N(0, 1) + 0.5N(0, 0.5^6)$;
- (c) the mixture double gamma distribution (MDG), $\varepsilon \sim \frac{1}{9}\varepsilon^*$, where $\varepsilon^* \sim f(\varepsilon) = e^{-14} \cdot \frac{1}{2}e^{-|\varepsilon|} + (1 - e^{-14}) \cdot \frac{1}{\Gamma(15)}|\varepsilon|^{14}e^{-|\varepsilon|}$;
- (d) the t -distribution with 3 degrees of freedom $\varepsilon \sim t_3$; and,
- (e) the Cauchy distribution, $\varepsilon \sim \text{Cauchy}$.

In the simulation study, our training data are composed of n observations $(\mathbf{x}_i, y_i)_{i=1}^n$, independently generated from model (15). An independent set of n observations is also simulated from the same model for parameter tuning of the training model. We evaluate the variable selection performance of the estimated coefficients $\hat{\beta}$ by the number of false positives $\text{FP} = |\hat{A} \setminus A^*|$ and the number of false negatives $\text{FN} = |A^* \setminus \hat{A}|$, where $A^* = \{1 \leq j \leq p: \beta_j^* \neq 0\}$ and $\hat{A} = \{1 \leq j \leq p: \hat{\beta}_j \neq 0\}$. The estimation accuracy of $\hat{\beta}$ is measured by the model error $(\hat{\beta} - \beta^*)^T \Sigma (\hat{\beta} - \beta^*)$. Two sets of data dimensions $(n, p) = (100, 600)$ and $(n, p) = (200, 1200)$ are used in our simulations. In all settings, we use $K = 19$ quantile levels $\tau_k = 0.05k$, $k = 1, \dots, 19$. The simulation results are summarized in Tables I and II.

It can be seen from the tables that the CQR oracle estimator has very similar model error to the LS estimator under normal error, while is more efficient under the other error distributions. In particular, the model error of the LS estimator explodes under the Cauchy error. In theory, it can be arbitrarily large. SCAD penalized CQR estimators have very close model errors to the CQR oracle estimator under most error distributions and outperform the penalized LS estimators. In terms of model selection accuracy, the SCAD penalized CQR estimator also outperform all the other penalized estimators.

The comparison between CV, BIC^{H} and BIC^{HL} for tuning parameter selection is shown in Tables III and IV. Note that BIC^{HL} does not perform well under the Cauchy error. This confirms its requirement for the first moment of the error distribution. The information criterion is computationally more efficient than CV and also delivers better results.

VII. DISCUSSION

In this article, we studied the sparse penalized CQR under various forms of regularization. In particular, we established the estimation consistency of the CQR lasso estimator. Through the LLA algorithm, we showed that the CQR oracle estimator could be achieved via folded concave penalized CQR. Our theoretical analysis remains valid even when the dimensionality is ultrahigh in the sense that $p = \mathcal{O}(n^\nu)$ with $0 < \nu < \frac{1}{2}$.

We also developed a fast and stable ADMM algorithm for solving the weighted L_1 -penalized CQR. Numerical studies proved the efficiency of the algorithm. The methodologies and numerical solvers proposed in this article make the sparse CQR a real alternative to the sparse LS in practice. It can be applied whenever the estimation efficiency of the coefficients are concerned.

TABLE I

SIMULATION RESULTS FOR THE NUMERICAL COMPARISON OF FOUR METHODS: LS-LASSO, LS-SCAD, CQR-LASSO AND CQR-SCAD, UNDER MODEL (15) WITH $n = 100$ AND $p = 600$. THE LS-ORACLE AND CQR-ORACLE SERVE AS THE BENCHMARK. TWO COVARIANCE STRUCTURES $\Sigma = (0.5^{i-j})$ AND $\Sigma = (0.8^{i-j})$ ARE SHOWN, UNDER EACH OF WHICH FIVE ERROR DISTRIBUTIONS ARE CONSIDERED: $N(0,3)$, MIXTURE NORMAL, MIXTURE DOUBLE GAMMA, t_3 AND CAUCHY. THE ESTIMATION ACCURACY IS REPORTED AS “MODEL ERROR” AND THE SELECTION ACCURACY IS REPORTED AS “FP, FN”. NUMBERS LISTED ARE AVERAGES OVER 100 INDEPENDENT RUNS, WITH STANDARD ERRORS REPORTED IN THE PARENTHESES

		$N(0,3)$	MN	MDG	t_3	Cauchy
$\Sigma = (0.5^{i-j})$						
Model error	LS-oracle	0.093 (0.008)	0.093 (0.007)	0.084 (0.007)	0.098 (0.009)	9350.072 (6837.507)
	CQR-oracle	0.105 (0.008)	0.004 (0.002)	0.025 (0.003)	0.047 (0.004)	0.094 (0.011)
	LS-lasso	0.664 (0.035)	0.620 (0.025)	0.588 (0.031)	0.663 (0.054)	18.963 (1.513)
	LS-SCAD	0.671 (0.038)	0.646 (0.036)	0.523 (0.033)	0.578 (0.036)	31.738 (7.959)
	CQR-lasso	0.792 (0.041)	0.272 (0.029)	0.465 (0.034)	0.374 (0.022)	1.672 (0.144)
	CQR-SCAD	0.122 (0.019)	0.006 (0.002)	0.032 (0.004)	0.064 (0.006)	0.438 (0.098)
	FP, FN	LS-lasso	16.55, 0 (1.28), (0)	16.91, 0 (0.92), (0)	16.53, 0 (1.07), (0)	15.83, 0 (1.08), (0)
LS-SCAD		18.04, 0 (1.54), (0)	18.00, 0 (1.44), (0)	17.04, 0 (1.58), (0)	16.81, 0 (1.10), (0)	17.11, 1.78 (2.93), (0.13)
CQR-lasso		15.33, 0 (0.67), (0)	15.29, 0 (0.67), (0)	14.49, 0 (0.53), (0)	12.89, 0 (0.55), (0)	38.75, 0.01 (2.93), (0.01)
CQR-SCAD		1.68, 0.01 (0.25), (0.01)	1.62, 0 (0.29), (0)	2.27, 0 (0.32), (0)	2.33, 0 (0.34), (0)	2.18, 0.01 (0.38), (0.01)
$\Sigma = (0.8^{i-j})$						
Model error	LS-oracle	0.097 (0.008)	0.092 (0.008)	0.079 (0.006)	0.088 (0.008)	184.788 (91.476)
	CQR-oracle	0.097 (0.008)	0.005 (0.002)	0.023 (0.002)	0.046 (0.004)	0.134 (0.011)
	LS-lasso	0.488 (0.025)	0.493 (0.028)	0.441 (0.021)	0.422 (0.031)	19.649 (1.623)
	LS-SCAD	0.524 (0.026)	0.443 (0.020)	0.422 (0.019)	0.442 (0.029)	27.706 (5.775)
	CQR-lasso	0.498 (0.029)	0.152 (0.023)	0.259 (0.029)	0.269 (0.014)	0.993 (0.087)
	CQR-SCAD	0.123 (0.013)	0.005 (0.001)	0.032 (0.003)	0.060 (0.005)	0.355 (0.055)
	FP, FN	LS-lasso	13.06, 0 (0.85), (0)	12.34, 0 (0.87), (0)	11.97, 0 (0.91), (0)	13.59, 0 (0.99), (0)
LS-SCAD		13.89, 0 (0.90), (0)	11.97, 0 (0.78), (0)	11.28, 0 (0.64), (0)	14.04, 0 (0.99), (0)	13.53, 1.62 (2.27), (0.11)
CQR-lasso		12.15, 0 (0.68), (0)	13.28, 0 (0.62), (0)	12.09, 0 (0.57), (0)	13.06, 0 (0.66), (0)	28.26, 0.03 (2.46), (0.02)
CQR-SCAD		2.09, 0.02 (0.40), (0.01)	1.29, 0 (0.24), (0)	1.83, 0 (0.32), (0)	1.97, 0 (0.30), (0)	2.38, 0.06 (0.32), (0.03)

TABLE II

SIMULATION RESULTS FOR THE NUMERICAL COMPARISON OF FOUR METHODS: LS-LASSO, LS-SCAD, CQR-LASSO AND CQR-SCAD, UNDER MODEL (15) WITH $n = 200$ AND $p = 1200$. THE LS-ORACLE AND CQR-ORACLE SERVE AS THE BENCHMARK. TWO COVARIANCE STRUCTURES $\Sigma = (0.5^{i-j})$ AND $\Sigma = (0.8^{i-j})$ ARE SHOWN, UNDER EACH OF WHICH FIVE ERROR DISTRIBUTIONS ARE CONSIDERED: $N(0,3)$, MIXTURE NORMAL, MIXTURE DOUBLE GAMMA, t_3 AND CAUCHY. THE ESTIMATION ACCURACY IS REPORTED AS “MODEL ERROR” AND THE SELECTION ACCURACY IS REPORTED AS “FP, FN”. NUMBERS LISTED ARE AVERAGES OVER 100 INDEPENDENT RUNS, WITH STANDARD ERRORS REPORTED IN THE PARENTHESES

		$N(0,3)$	MN	MDG	t_3	Cauchy
$\Sigma = (0.5^{i-j})$						
Model error	LS-oracle	0.051 (0.004)	0.045 (0.004)	0.041 (0.003)	0.048 (0.004)	1136.066 (965.520)
	CQR-oracle	0.047 (0.005)	0.001 (0)	0.011 (0.001)	0.023 (0.002)	0.060 (0.005)
	LS-lasso	0.340 (0.015)	0.337 (0.014)	0.281 (0.011)	0.284 (0.013)	28.450 (7.987)
	LS-SCAD	0.061 (0.006)	0.061 (0.005)	0.055 (0.005)	0.062 (0.005)	41.685 (24.654)
	CQR-lasso	0.394 (0.018)	0.072 (0.011)	0.180 (0.014)	0.239 (0.013)	0.830 (0.073)
	CQR-SCAD	0.046 (0.004)	0.001 (0)	0.011 (0.001)	0.023 (0.002)	0.137 (0.030)
	FP, FN	LS-lasso	19.62, 0 (1.41), (0)	20.09, 0 (1.25), (0)	20.15, 0 (1.22), (0)	19.59, 0 (1.27), (0)
LS-SCAD		5.24, 0 (1.08), (0)	5.76, 0 (0.94), (0)	4.56, 0 (0.92), (0)	6.53, 0 (1.10), (0)	25.49, 1.54 (4.01), (0.12)
CQR-lasso		18.76, 0 (0.95), (0)	19.05, 0 (0.77), (0)	19.11, 0 (0.78), (0)	18.59, 0 (0.95), (0)	60.56, 0 (6.35), (0)
CQR-SCAD		2.21, 0 (0.30), (0)	2.29, 0 (0.48), (0)	2.99, 0 (0.44), (0)	2.31, 0 (0.36), (0)	1.50, 0 (0.30), (0)
$\Sigma = (0.8^{i-j})$						
Model error	LS-oracle	0.042 (0.004)	0.047 (0.005)	0.034 (0.003)	0.046 (0.004)	71435.826 (68875.639)
	CQR-oracle	0.049 (0.004)	0.001 (0)	0.011 (0.001)	0.022 (0.002)	0.055 (0.005)
	LS-lasso	0.252 (0.012)	0.235 (0.011)	0.219 (0.009)	0.208 (0.013)	22.598 (2.334)
	LS-SCAD	0.071 (0.006)	0.073 (0.007)	0.050 (0.004)	0.065 (0.009)	23.726 (2.856)
	CQR-lasso	0.255 (0.014)	0.030 (0.004)	0.099 (0.008)	0.153 (0.009)	0.730 (0.070)
	CQR-SCAD	0.086 (0.008)	0.001 (0)	0.023 (0.003)	0.048 (0.004)	0.539 (0.099)
	FP, FN	LS-lasso	14.83, 0 (1.03), (0)	16.17, 0 (1.1), (0)	15.24, 0 (1.09), (0)	14.80, 0 (1.23), (0)
LS-SCAD		6.36, 0 (0.99), (0)	5.68, 0 (0.72), (0)	5.64, 0 (0.83), (0)	4.74, 0.01 (0.76), (0.01)	15.22, 1.84 (2.92), (0.10)
CQR-lasso		16.86, 0 (1.02), (0)	14.89, 0 (0.7), (0)	15.83, 0 (0.78), (0)	16.47, 0 (0.98), (0)	42.75, 0 (4.13), (0)
CQR-SCAD		2.19, 0 (0.32), (0)	1.93, 0 (0.36), (0)	2.70, 0 (0.50), (0)	2.59, 0 (0.44), (0)	1.85, 0 (0.24), (0)

TABLE III

SIMULATION RESULTS FOR NUMERICAL COMPARISON BETWEEN CV, BIC^H AND BIC^{HL} IN TERMS OF TUNING PARAMETER SELECTION, UNDER MODEL (15) WITH $n = 100$ AND $p = 600$. TWO COVARIANCE STRUCTURES $\Sigma = (0.5^{i-j})$ AND $\Sigma = (0.8^{i-j})$ ARE SHOWN, UNDER EACH OF WHICH FIVE ERROR DISTRIBUTIONS ARE CONSIDERED: $N(0,3)$, MIXTURE NORMAL, MIXTURE DOUBLE GAMMA, t_3 AND CAUCHY. THE ESTIMATION ACCURACY IS REPORTED AS “MODEL ERROR” AND THE SELECTION ACCURACY IS REPORTED AS “FP, FN”. NUMBERS LISTED ARE AVERAGES OVER 100 INDEPENDENT RUNS, WITH STANDARD ERRORS REPORTED IN THE PARENTHESES

		$N(0,3)$	MN	MDG	t_3	Cauchy
$\Sigma = (0.5^{i-j})$						
Model error	CV	0.087 (0.008)	0.005 (0.002)	0.021 (0.002)	0.043 (0.004)	0.312 (0.083)
	BIC^H	0.115 (0.019)	0.017 (0.011)	0.053 (0.017)	0.058 (0.012)	0.988 (0.374)
	BIC^{HL}	0.094 (0.011)	0.003 (0.001)	0.022 (0.003)	0.048 (0.004)	6.287 (0.824)
FP, FN	CV	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0.01, 0.13 (0.01), (0.04)
	BIC^H	0, 0.03 (0), (0.02)	0, 0.01 (0), (0.01)	0, 0.03 (0), (0.02)	0, 0.01 (0), (0.01)	0.01, 0.27 (0.01), (0.05)
	BIC^{HL}	0, 0.01 (0), (0.01)	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0, 1.17 (0), (0.09)
$\Sigma = (0.8^{i-j})$						
Model error	CV	0.273 (0.047)	0.008 (0.003)	0.088 (0.025)	0.118 (0.025)	0.718 (0.116)
	BIC^H	0.393 (0.066)	0.077 (0.035)	0.200 (0.051)	0.295 (0.059)	2.169 (0.338)
	BIC^{HL}	0.607 (0.091)	0.024 (0.020)	0.153 (0.043)	0.114 (0.020)	5.411 (0.517)
FP, FN	CV	0, 0.07 (0), (0.03)	0, 0 (0), (0)	0.02, 0.02 (0.02), (0.01)	0, 0.01 (0), (0.01)	0.1, 0.2 (0.03), (0.04)
	BIC^H	0.01, 0.15 (0.01), (0.04)	0, 0.04 (0), (0.02)	0, 0.09 (0), (0.03)	0, 0.12 (0), (0.03)	0.14, 0.53 (0.06), (0.07)
	BIC^{HL}	0.05, 0.23 (0.02), (0.04)	0, 0.01 (0), (0.01)	0, 0.06 (0), (0.02)	0, 0.01 (0), (0.01)	0.12, 1.07 (0.04), (0.08)

TABLE IV

SIMULATION RESULTS FOR NUMERICAL COMPARISON BETWEEN CV, BIC^H AND BIC^{HL} IN TERMS OF TUNING PARAMETER SELECTION, UNDER MODEL (15) WITH $n = 200$ AND $p = 1200$. TWO COVARIANCE STRUCTURES $\Sigma = (0.5^{|i-j|})$ AND $\Sigma = (0.8^{|i-j|})$ ARE SHOWN, UNDER EACH OF WHICH FIVE ERROR DISTRIBUTIONS ARE CONSIDERED: $N(0,3)$, MIXTURE NORMAL, MIXTURE DOUBLE GAMMA, t_3 AND CAUCHY. THE ESTIMATION ACCURACY IS REPORTED AS “MODEL ERROR” AND THE SELECTION ACCURACY IS REPORTED AS “FP, FN”. NUMBERS LISTED ARE AVERAGES OVER 100 INDEPENDENT RUNS, WITH STANDARD ERRORS REPORTED IN THE PARENTHESES

		$N(0,3)$	MN	MDG	t_3	Cauchy
$\Sigma = (0.5^{ i-j })$						
Model error	CV	0.121 (0.019)	0.006 (0.002)	0.026 (0.003)	0.043 (0.006)	0.477 (0.224)
	BIC^H	0.041 (0.004)	0.001 (0.000)	0.012 (0.001)	0.020 (0.002)	0.915 (0.423)
	BIC^{HL}	0.042 (0.004)	0.001 (0.000)	0.009 (0.001)	0.023 (0.002)	3.632 (0.757)
FP, FN	CV	0, 0.02 (0), (0.01)	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0.01, 0.11 (0.01), (0.04)
	BIC^H	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0.26, 0.05 (0.13), (0.03)
	BIC^{HL}	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0.01, 0.61 (0.01), (0.08)
$\Sigma = (0.8^{ i-j })$						
Model error	CV	0.372 (0.059)	0.042 (0.025)	0.081 (0.023)	0.206 (0.044)	0.736 (0.141)
	BIC^H	0.235 (0.05)	0.001 (0.000)	0.070 (0.03)	0.055 (0.017)	0.538 (0.086)
	BIC^{HL}	0.132 (0.032)	0.002 (0.001)	0.019 (0.002)	0.041 (0.004)	3.28 (0.518)
FP, FN	CV	0.01, 0.12 (0.01), (0.03)	0, 0.02 (0), (0.01)	0.01, 0.02 (0.01), (0.01)	0, 0.06 (0), (0.02)	0.05, 0.16 (0.02), (0.04)
	BIC^H	0, 0.09 (0), (0.03)	0, 0 (0), (0)	0, 0.03 (0), (0.02)	0, 0.01 (0), (0.01)	0.21, 0.12 (0.06), (0.03)
	BIC^{HL}	0, 0.03 (0), (0.02)	0, 0 (0), (0)	0, 0 (0), (0)	0, 0 (0), (0)	0.10, 0.55 (0.04), (0.07)

VIII. PROOFS

We provide proofs of all previously stated results in this section. For the sake of brevity, some auxiliary results are relegated to the appendix.

Assume without loss of generality $\beta_0^* = 0$ such that $\alpha_k^* = F^{-1}(\tau_k)$ for $1 \leq k \leq K$. Also recall $M_0 = \max_{0 \leq j \leq p} \|X_j\|_2^2/n$ from Theorem 1.

Lemma 3. *Under condition (C0), with probability at least*

$$1 - 2K \exp\left(-\frac{9}{2}n\lambda^2\right) - 2p \exp\left(-\frac{n\lambda^2}{2M_0}\right),$$

the CQR lasso estimator $(\widehat{\alpha}_\lambda, \widehat{\beta}_\lambda)$ satisfies

$$(\widehat{\delta}^\lambda, \widehat{\Delta}^\lambda) \in \mathcal{C} = \{(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p : \|\Delta_{\mathcal{A}^c}\|_1 \leq 3\|\Delta_{\mathcal{A}}\|_1 + \frac{3}{K}\|\delta\|_1\},$$

where $\widehat{\delta}^\lambda = \widehat{\alpha}_\lambda - \alpha^*$ and $\widehat{\Delta}^\lambda = \widehat{\beta}_\lambda - \beta^*$.

Proof of Lemma 3. See Section A of the appendix. □

Now let $v_n(\alpha, \beta) = Q_n(\alpha, \beta) - Q_n(\alpha^*, \beta^*) - \mathbb{E}[Q_n(\alpha, \beta) - Q_n(\alpha^*, \beta^*)]$. For $r > 0$, set $\mathcal{C}_r = \{(\delta, \Delta) \in \mathcal{C} : (nK)^{-1} \sum_{k=1}^K \sum_{i=1}^n (\delta_k + \mathbf{x}_i^T \Delta)^2 \leq r^2\}$ and define $e(r) = \sup_{(\delta, \Delta) \in \mathcal{C}_r} |v_n(\alpha^* + \delta, \beta^* + \Delta)|$.

Lemma 4. *For $r, t > 0$, under conditions (C0) and (C1), with probability at least $1 - \exp[-nt^2/(32r^2)]$, we have*

$$e(r) \leq 16 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1)r + t.$$

It follows immediately that, if one takes

$$t = 16 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1)r,$$

then with probability at least $1 - \exp[-2M_0\kappa_0^{-1}s(1 + \log p)]$, we have

$$e(r) \leq 32 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1)r.$$

Proof of Lemma 4. See Section A of the appendix. □

Lemma 5. *Under conditions (C0) and (C2), for any $(\delta, \Delta) \in \mathcal{C}$, we have*

$$\mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \geq \min\{\underline{f}r^2/4, q(\underline{f}/K)^{1/2}r\},$$

where $r^2 = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^T \Delta)^2$.

Proof of Lemma 5. See Section A of the appendix. □

Proof of Theorem 1. Let

$$r_* = 8\underline{f}^{-1} \left[32 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1 + \log p}{n}} (\sqrt{s} + 1) + \lambda \sqrt{\frac{s}{\kappa_0}} \right]$$

and set $\mathcal{C}^* = \{(\delta, \Delta) \in \mathcal{C} : (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^T \Delta)^2 = r_*^2\}$. Moreover, define $\widehat{\delta}^\lambda = \widehat{\alpha}_\lambda - \alpha^*$ and $\widehat{\Delta}^\lambda = \widehat{\beta}_\lambda - \beta^*$. Under event $\mathcal{E}_1 = \{(\widehat{\delta}^\lambda, \widehat{\Delta}^\lambda) \in \mathcal{C}\}$, if

$$\inf_{(\delta, \Delta) \in \mathcal{C}^*} Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*) + \lambda(\|\beta^* + \Delta\|_1 - \|\beta^*\|_1) > 0, \quad (16)$$

then by convexity of Q_n , this implies that $(\widehat{\delta}^\lambda, \widehat{\Delta}^\lambda) \in \mathcal{C}_{r_*}^*$. To show (16), first note that for all $(\delta, \Delta) \in \mathcal{C}^*$,

$$\begin{aligned} & Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*) + \lambda(\|\beta^* + \Delta\|_1 - \|\beta^*\|_1) \\ & \geq \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] - e(r_*) \\ & \quad + \lambda(\|\Delta_{\mathcal{A}^c}\|_1 - \|\Delta_{\mathcal{A}}\|_1). \end{aligned} \quad (17)$$

Now let $\mathcal{E}_2 = \{e(r_*) \leq 32\sqrt{2M_0(1+\log p)}/(n\kappa_0)(\sqrt{s}+1)r_*\}$. It follows from Lemma 4 that $\Pr(\mathcal{E}_2) \geq 1 - \exp[-2M_0\kappa_0^{-1}s(1+\log p)]$. By Lemma 5, for any $(\delta, \Delta) \in \mathcal{C}^*$, we have

$$\mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \geq \min\{fr_*^2/4, q(\underline{f}/K)^{1/2}r_*\}.$$

Also, by condition (C1) and (A.5), for $(\delta, \Delta) \in \mathcal{C}^*$, we have $\|\Delta_{\mathcal{A}^c}\|_1 \leq r_*\sqrt{s/\kappa_0}$. Thus, under event $\mathcal{E}_1 \cap \mathcal{E}_2$, for any $(\delta, \Delta) \in \mathcal{C}^*$, it follows from (17) and the growth condition that

$$\begin{aligned} & Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*) + \lambda(\|\beta^* + \Delta\|_1 - \|\beta^*\|_1) \\ & \geq \frac{f}{4}r_*^2 - \left[32\sqrt{\frac{2M_0}{\kappa_0}}\sqrt{\frac{1+\log p}{n}}(\sqrt{s}+1) + \lambda\sqrt{\frac{s}{\kappa_0}}\right]r_* > 0 \end{aligned}$$

by our choice of r_* . Therefore, by Lemma 3 and 4, with probability at least

$$\Pr(\mathcal{E}_1 \cap \mathcal{E}_2) \geq 1 - \Pr(\mathcal{E}_1^c) - \Pr(\mathcal{E}_2^c) \geq 1 - p_1(\lambda),$$

we have $(\widehat{\delta}^\lambda, \widehat{\Delta}^\lambda) \in \mathcal{C}_{r_*}^*$. This, by condition (C1), further implies that

$$\begin{aligned} r_*^2 & \geq \frac{\kappa_m}{K} \left[\|\widehat{\delta}^\lambda\|_2^2 + K \|\widehat{\Delta}_{\mathcal{A} \cup \overline{\mathcal{A}}}^\lambda(\widehat{\Delta}^\lambda, m)\|_2^2 \right] \\ & \geq \frac{\kappa_m}{K} \|\widehat{\delta}^\lambda\|_2^2 + \kappa_m \|\widehat{\Delta}_{\mathcal{A} \cup \overline{\mathcal{A}}}^\lambda(\widehat{\Delta}^\lambda, m)\|_2^2. \end{aligned}$$

As a result, we obtain $\|\widehat{\delta}^\lambda\|_2 \leq r_*\sqrt{K/\kappa_m}$ and

$$\|\widehat{\Delta}_{\mathcal{A} \cup \overline{\mathcal{A}}}^\lambda(\widehat{\Delta}^\lambda, m)\|_2 \leq r_*/\sqrt{\kappa_m}. \quad (18)$$

Note that the j th largest in absolute value component of $\widehat{\Delta}_{\mathcal{A}^c}$ is bounded by $\|\widehat{\Delta}_{\mathcal{A}^c}\|_1/j$. Therefore, it follows that

$$\begin{aligned} \left\| \widehat{\Delta}_{(\mathcal{A} \cup \overline{\mathcal{A}})(\widehat{\Delta}^\lambda, m)^c} \right\|_2^2 & \leq \sum_{j=m+1}^p \frac{\|\widehat{\Delta}_{\mathcal{A}^c}\|_1^2}{j^2} \leq \frac{1}{m} \|\widehat{\Delta}_{\mathcal{A}^c}\|_1^2 \\ & \leq \frac{1}{m} [3\|\widehat{\Delta}_{\mathcal{A}}\|_1 + 3K^{-1}\|\widehat{\delta}^\lambda\|_1]^2 \leq \frac{18s}{m} \|\widehat{\Delta}_{\mathcal{A}}\|_2^2 + \frac{18}{mK} \|\widehat{\delta}^\lambda\|_2^2 \\ & \leq \frac{18s}{m} \|\widehat{\Delta}_{\mathcal{A} \cup \overline{\mathcal{A}}}^\lambda(\widehat{\Delta}^\lambda, m)\|_2^2 + \frac{18}{mK} \|\widehat{\delta}^\lambda\|_2^2, \end{aligned}$$

which, together with (18), implies that

$$\begin{aligned} \|\widehat{\Delta}^\lambda\|_2^2 & \leq \left(1 + \frac{18s}{m}\right) \|\widehat{\Delta}_{\mathcal{A} \cup \overline{\mathcal{A}}}^\lambda(\widehat{\Delta}^\lambda, m)\|_2^2 + \frac{18}{mK} \|\widehat{\delta}^\lambda\|_2^2 \\ & \leq \frac{r_*^2}{\kappa_m} \left(1 + \frac{18s}{m} + \frac{18}{m}\right). \end{aligned}$$

This completes the proof of Theorem 1. \square

For $r > 0$, define $B_{\mathcal{A}}(r) = \{(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p: \|\delta\|_2^2 + \|\Delta_{\mathcal{A}}\|_2^2 \leq r^2, \Delta_{\mathcal{A}^c} = \mathbf{0}\}$ and $S_{\mathcal{A}}(r) = \{(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p: \|\delta\|_2^2 + \|\Delta_{\mathcal{A}}\|_2^2 = r^2, \Delta_{\mathcal{A}^c} = \mathbf{0}\}$. Moreover, let $z(r) = \sup_{(\delta, \Delta) \in B_{\mathcal{A}}(r)} |\mathbf{v}_n(\alpha^* + \delta, \beta^* + \Delta)|$.

Lemma 6. *Under condition (C0), for any $r, t > 0$ such that $r\sqrt{(s+1)M_{\mathcal{A}}} \leq \mathcal{U}_0$, with probability at least $1 - \exp[-nt^2/(32\bar{\mu}r^2)]$, we have*

$$z(r) \leq 4r\sqrt{\frac{(s+1)M_{\mathcal{A}}}{n}} + t$$

and

$$\inf_{(\delta, \Delta) \in S_{\mathcal{A}}(r)} [Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \geq \frac{f}{2K}\bar{\mu}r^2 - 4r\sqrt{\frac{(s+1)M_{\mathcal{A}}}{n}} - t.$$

Proof of Lemma 6. See Section A of the appendix. \square

Proof of Lemma 1. For ease of notation, let $\hat{\delta}^{\circ} = \hat{\alpha}^{\circ} - \alpha^*$ and $\hat{\Delta}^{\circ} = \hat{\beta}^{\circ} - \beta^*$. In Lemma 6, let $r = r^* = 32K(f\bar{\mu})^{-1} \times \sqrt{M_{\mathcal{A}}(s+1)/n}$ and $t = 4r^* \sqrt{M_{\mathcal{A}}(s+1)/n}$. By assumption, with the choice of r^* , we have $r^* \sqrt{(s+1)M_{\mathcal{A}}} \leq \mathcal{U}_0$. It follows immediately that with probability at least $1 - \exp[-(s+1)M_{\mathcal{A}}/(2\bar{\mu})]$, we have

$$\inf_{(\delta, \Delta) \in S_{\mathcal{A}}(r^*)} [Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \geq \frac{f\bar{\mu}(r^*)^2}{2K} - 8r^* \sqrt{\frac{(s+1)M_{\mathcal{A}}}{n}} > 0.$$

By convexity of Q_n and optimality of $(\hat{\alpha}^{\circ}, \hat{\beta}^{\circ})$, this implies that

$$\|\hat{\delta}^{\circ}\|_2^2 + \|\hat{\Delta}^{\circ}\|_2^2 \leq (r^*)^2,$$

which completes the lemma. \square

Lemma 7. *Suppose the folded concave penalized CQR (6) is solved with the LLA algorithm. Let $a_0 = \min(a_2, 1)$ and define*

$$\begin{aligned} \mathcal{E}_1 &= \{\|\hat{\beta}^{(0)} - \beta^*\|_{\infty} \leq a_0\lambda\}, \\ \mathcal{E}_2 &= \{\|\nabla_{\mathcal{A}^c} Q_n(\hat{\alpha}^{\circ}, \hat{\beta}^{\circ})\|_{\infty} < a_1\lambda\}, \\ \mathcal{E}_3 &= \left\{ \min_{j \in \mathcal{A}} |\hat{\beta}_j^{\circ}| > a\lambda \right\}, \end{aligned}$$

where $\nabla_{\mathcal{A}^c} Q_n(\hat{\alpha}^{\circ}, \hat{\beta}^{\circ}) = (\nabla_j Q_n(\hat{\alpha}^{\circ}, \hat{\beta}^{\circ}), j \in \mathcal{A}^c)$ with

$$\nabla_j Q_n(\hat{\alpha}^{\circ}, \hat{\beta}^{\circ}) = \frac{1}{2n} \sum_{i=1}^n x_{ij} \left(1 - \frac{2}{K} \sum_{k=1}^K \tau_k\right) - \frac{1}{2nK} \sum_{i=1}^n \sum_{k=1}^K \text{Sgn}(\hat{r}_{ik}) x_{ij},$$

$\hat{r}_{ik} = y_i - \hat{\alpha}_k^{\circ} - \mathbf{x}_i^T \hat{\beta}^{\circ}$, $1 \leq i \leq n$, $1 \leq k \leq K$, and

$$\text{Sgn}(u) = \begin{cases} 1, & \text{if } u > 0 \\ [-1, 1], & \text{if } u = 0 \\ -1, & \text{if } u < 0. \end{cases}$$

Then under $\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3$ and condition (C0), the LLA algorithm converges to the CQR oracle estimator.

Proof of Lemma 7. See Section A of the appendix. \square

For each $j \in \mathcal{A}^c$, define $S_j^n(\alpha, \beta) = \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K [I(y_i - \alpha_k - \mathbf{x}_i^\top \beta \leq 0) - \tau_k] x_{ij}$, and for $r > 0$, let

$$\begin{aligned} \gamma_j(r) = & \sup_{(\delta, \Delta) \in B_{\mathcal{A}^c}(r)} |S_j^n(\alpha^* + \delta, \beta^* + \Delta) - S_j^n(\alpha^*, \beta^*) \\ & - \mathbb{E}[S_j^n(\alpha^* + \delta, \beta^* + \Delta) - S_j^n(\alpha^*, \beta^*)]|. \end{aligned}$$

Lemma 8. For $r, t > 0$, $0 < \psi < r$ and $j \in \mathcal{A}^c$, under condition (C0) we have

$$\begin{aligned} \Pr(\gamma_j(r) > t) \leq & 2N_\psi \exp\left(-\frac{nt^2}{8\bar{f}M_{\mathcal{A}^c}^2 \bar{\mu}^{1/2} r + \frac{8}{3}M_{\mathcal{A}^c} t}\right) \\ & + 2N_\psi \exp\left(-\frac{nt_0^2}{2\bar{f}M_{\mathcal{A}^c}^2 ((s+1)M_{\mathcal{A}^c})^{1/2} \psi + \frac{4}{3}M_{\mathcal{A}^c} t_0}\right), \end{aligned}$$

where N_ψ is the ψ -covering number (see, e.g., [24]) of $B_{\mathcal{A}^c}(r)$ and $t_0 = \lceil t/2 - 2\bar{f}M_{\mathcal{A}^c}((s+1)M_{\mathcal{A}^c})^{1/2} \psi \rceil_+$.

Proof of Lemma 8. See Section A of the appendix. \square

Proof of Theorem 2. Let $\hat{\delta}^o = \hat{\alpha}^o - \alpha^*$ and $\hat{\Delta}^o = \hat{\beta}^o - \beta^*$. For $1 \leq i \leq n$, $1 \leq k \leq K$, write $\hat{r}_{ik} = y_i - \hat{\alpha}_k^o - \mathbf{x}_i^\top \hat{\beta}^o$ and $r_{ik}^* = y_i - \alpha_k^* - \mathbf{x}_i^\top \beta^*$. For ease of notation, let $F(\delta, \Delta) = Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)$ for $(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p$.

According to Lemma 7, with probability at least

$$\Pr(\mathcal{E}_1 \cap \mathcal{E}_2 \cap \mathcal{E}_3) \geq 1 - \Pr(\mathcal{E}_1^c) - \Pr(\mathcal{E}_2^c) - \Pr(\mathcal{E}_3^c),$$

the LLA algorithm will converge to the oracle estimator in two iterations. In the sequel, we will split the proof into three parts and provide the upper bound on each of $\Pr(\mathcal{E}_1^c)$, $\Pr(\mathcal{E}_2^c)$ and $\Pr(\mathcal{E}_3^c)$, separately.

(i) First, we deal with $\Pr(\mathcal{E}_1^c) = \Pr(\|\hat{\beta}^{(0)} - \beta^*\|_\infty > a_0 \lambda)$. Since in the LLA algorithm, we take $(\hat{\alpha}^{(0)}, \hat{\beta}^{(0)})$ to be the lasso estimator $(\hat{\alpha}_{\lambda_0}, \hat{\beta}_{\lambda_0})$, by Theorem 1, we have

$$\Pr(\mathcal{E}_1) = \Pr(\|\hat{\beta}_{\lambda_0} - \beta^*\|_\infty \leq a_0 \lambda) \geq \Pr(\|\hat{\beta}_{\lambda_0} - \beta^*\|_2 \leq a_0 \lambda) \geq 1 - p_1(\lambda_0),$$

which implies that $\Pr(\mathcal{E}_1^c) \leq p_1(\lambda_0)$.

(ii) We next derive the upper bound on $\Pr(\mathcal{E}_3^c) = \Pr(\min_{j \in \mathcal{A}} |\hat{\beta}_j^o| \leq a\lambda)$. Let $r_0 = \min_{j \in \mathcal{A}} |\beta_j^*| - a\lambda$. It can be seen that $\Pr(\mathcal{E}_3^c) \leq \Pr(\|\hat{\Delta}^o\|_\infty > r_0)$. Note that by convexity of Q_n , $\|\hat{\Delta}^o\|_2 \leq r_0$ is implied by the event that $\inf_{(\delta, \Delta) \in S_{\mathcal{A}}(r_0)} F(\delta, \Delta) > 0$. Since $r_0 \sqrt{(s+1)M_{\mathcal{A}}} \leq \mathcal{U}_0$, it follows from Lemma 6 that for any $t > 0$,

$$\inf_{(\delta, \Delta) \in S_{\mathcal{A}}(r_0)} F(\delta, \Delta) \geq \frac{f}{2K} \mu r_0^2 - 4r_0 \sqrt{\frac{(s+1)M_{\mathcal{A}}}{n}} - t$$

holds with probability at least $1 - \exp[-nt^2/(32\bar{\mu}r_0^2)]$. By condition (C3), it can be seen that $r_0 > \lambda > 8K(f\bar{\mu})^{-1} \times \sqrt{(s+1)M_{\mathcal{A}}/n}$. Now take $t = f\bar{\mu}r_0^2/(4K) - 2r_0 \sqrt{(s+1)M_{\mathcal{A}}/n}$. Then, we can see that $t > 0$. It follows immediately that $\inf_{(\delta, \Delta) \in S_{\mathcal{A}}(r_0)} F(\delta, \Delta) \geq t > 0$. With this specific choice of t , we get

$$\Pr(\|\hat{\Delta}^o\|_2 \leq r_0) \geq 1 - \exp[-nt^2/(32\bar{\mu}r_0^2)],$$

which implies that

$$\Pr(\mathcal{E}_3^c) \leq \Pr(\|\hat{\Delta}^o\|_\infty > r_0) \leq \Pr(\|\hat{\Delta}^o\|_2 > r_0) \leq \exp[-nt^2/(32\bar{\mu}r_0^2)].$$

(iii) Finally, we look at $\Pr(\mathcal{E}_2^c) = \Pr(\|\nabla_{\mathcal{A}^c} \mathcal{Q}_n(\hat{\alpha}^0, \hat{\beta}^0)\|_\infty \geq a_1 \lambda)$. To this end, we set $r_* = \sqrt{(s+1)M_{\mathcal{A}^c} \log n/n}$ and let $\mathcal{R} = \{(i, k) : \hat{r}_{ik} = 0, 1 \leq i \leq n, 1 \leq k \leq K\}$ be the index set of zero residuals. From Section B of the appendix, we have $|\mathcal{R}| \leq K(K+s)$. It follows that

$$\begin{aligned} \nabla_j \mathcal{Q}_n(\hat{\alpha}^0, \hat{\beta}^0) &= \frac{1}{2nK} \sum_{i=1}^n \sum_{k=1}^K [(1-2\tau_k) - \text{Sgn}(\hat{r}_{ik})] x_{ij} \\ &= \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K [I(\hat{r}_{ik} \leq 0) - \tau_k] x_{ij} - \frac{1}{2nK} \sum_{(i,k) \in \mathcal{R}} [\text{Sgn}(\hat{r}_{ik}) + 1] x_{ij}, \end{aligned}$$

where we have

$$\max_{j \in \mathcal{A}^c} \left| \frac{1}{2nK} \sum_{(i,k) \in \mathcal{R}} [\text{Sgn}(\hat{r}_{ik}) + 1] x_{ij} \right| \leq \frac{(K+s)M_{\mathcal{A}^c}}{n} := B_1.$$

Now define event $\mathcal{E}_0 = \{(\hat{\delta}^0, \hat{\Delta}^0) \in B_{\mathcal{A}^c}(r_*)\}$. Under \mathcal{E}_0 , by the triangular inequality, we have

$$\begin{aligned} \max_{j \in \mathcal{A}^c} \left| \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K [I(\hat{r}_{ik} \leq 0) - \tau_k] x_{ij} \right| &\leq \max_{j \in \mathcal{A}^c} \gamma_j(r_*) + \max_{j \in \mathcal{A}^c} |S_j^n(\alpha^*, \beta^*)| \\ &+ \max_{j \in \mathcal{A}^c} \sup_{(\delta, \Delta) \in B_{\mathcal{A}^c}(r_*)} \left| \mathbb{E}[S_j^n(\alpha^* + \delta, \beta^* + \Delta) - S_j^n(\alpha^*, \beta^*)] \right|. \end{aligned}$$

By the mean value theorem, it can be seen that

$$\begin{aligned} \max_{j \in \mathcal{A}^c} \sup_{(\delta, \Delta) \in B_{\mathcal{A}^c}(r_*)} \left| \mathbb{E}[S_j^n(\alpha^* + \delta, \beta^* + \Delta) - S_j^n(\alpha^*, \beta^*)] \right| \\ \leq \frac{1}{nK} \bar{f} M_{\mathcal{A}^c} \sup_{(\delta, \Delta) \in B_{\mathcal{A}^c}(r_*)} \sum_{i=1}^n \sum_{k=1}^K |\delta_k + \mathbf{x}_{i, \mathcal{A}^c}^\top \Delta_{\mathcal{A}^c}| \leq \bar{f} M_{\mathcal{A}^c} \bar{\mu}^{1/2} r_* := B_2. \end{aligned}$$

Note that if $B \geq 0$, then $2B = a_1 \lambda - B_1 - B_2$. It follows that

$$\begin{aligned} \Pr(\mathcal{E}_2^c) &\leq \Pr((\hat{\delta}^0, \hat{\Delta}^0) \notin B_{\mathcal{A}^c}(r_*)) + \Pr\left(\max_{i \in \mathcal{A}^c} \gamma_j(r_*) \geq B\right) \\ &+ \Pr\left(\max_{j \in \mathcal{A}^c} |S_j^n(\alpha^*, \beta^*)| \geq B\right). \end{aligned}$$

Note that $r_* \sqrt{(s+1)M_{\mathcal{A}^c}} \leq \mathcal{U}_0$. By similar arguments in (ii), it can be shown that

$$\Pr((\hat{\delta}^0, \hat{\Delta}^0) \notin B_{\mathcal{A}^c}(r_*)) \leq \exp\left(-\frac{m_*^2}{32\bar{\mu}r_*^2}\right),$$

where $t_* = \underline{f} \mu r_*^2 / (4K) - 2r_* \sqrt{(s+1)M_{\mathcal{A}^c}/n}$. Applying Hoeffding's inequality, we obtain

$$\Pr\left(\max_{j \in \mathcal{A}^c} |S_j^n(\alpha^*, \beta^*)| \geq B\right) \leq 2(p-s) \exp\left(-\frac{2nB^2}{M_0}\right).$$

Lastly, we apply Lemma 8 to obtain the bound on $\Pr\left(\max_{j \in \mathcal{A}^c} \gamma_j(r_*) \geq B\right)$. Let $\psi = 4r_*/n^2$. It can be shown that the ψ -covering number of $B_{\mathcal{A}^c}(r_*)$ satisfies

$$\left(\frac{r_*}{\psi}\right)^{K+s} \leq N_\psi \leq \left(\frac{2r_* + \psi}{\psi}\right)^{K+s} \leq n^{2(K+s)}, n \geq 2.$$

By Lemma 8, we have

$$\begin{aligned} \Pr\left(\max_{j \in \mathcal{A}^c} \gamma_j(r_*) \geq B\right) &\leq 2(p-s)N_\psi \exp\left(-\frac{nB^2}{8\bar{f}M_{\mathcal{A}^c}^2 \bar{\mu}^{1/2} r_* + \frac{8}{3}M_{\mathcal{A}^c} B}\right) \\ &+ 2(p-s)N_\psi \exp\left(-\frac{nB_0^2}{2\bar{f}M_{\mathcal{A}^c}^2 ((s+1)M_{\mathcal{A}^c})^{1/2} \psi + \frac{4}{3}M_{\mathcal{A}^c} B_0}\right), \end{aligned}$$

where $B_0 = [B/2 - 2\bar{f}M_{\mathcal{A}^c}((s+1)M_{\mathcal{A}})^{1/2}\boldsymbol{\psi}]_+$. This completes the proof. \square

Lemma 9. For any $A \in \{S : S \supset \mathcal{A}, |S| \leq 2J_n\}$ and $r, t > 0$, let $z(A, r) = \sup_{(\boldsymbol{\delta}, \Delta) \in B_A(r)} |\mathbf{v}_n(\boldsymbol{\alpha}^* + \boldsymbol{\delta}, \boldsymbol{\beta}^* + \Delta)|$. With probability at least $1 - \exp\{-nt^2/(32\bar{\zeta}r^2)\}$, we have

$$z(A, r) \leq 4Mr \sqrt{\frac{|A|+1}{n}} + t.$$

Proof of Lemma 9. See Section A of the appendix. \square

Proof of Theorem 3. Split all models (denoted by their index sets) under consideration, $\{\hat{A}_\lambda : \lambda \in \Xi_n\}$, into three groups: $\{\hat{A}_\lambda : \lambda \in \Xi_n^-\}$, $\{\hat{A}_\lambda : \lambda \in \Xi_n^0\}$, and $\{\hat{A}_\lambda : \lambda \in \Xi_n^+\}$, where $\Xi_n^- = \{\lambda \in \Xi_n : \mathcal{A} \not\subset \hat{A}_\lambda\}$ (underfitted models), $\Xi_n^0 = \{\lambda \in \Xi_n : \hat{A}_\lambda = \mathcal{A}\}$ and $\Xi_n^+ = \{\lambda \in \Xi_n : \mathcal{A} \subset \hat{A}_\lambda, \hat{A}_\lambda \neq \mathcal{A}\}$ (overfitted models). The proof then boils down to show that $P\left(\inf_{\lambda \in \Xi_n^-} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0\right) \rightarrow 1$ and $P\left(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0\right) \rightarrow 1$ as $n \rightarrow \infty$. Note that BIC contains both a goodness-of-fit part and a model complexity part. The basic idea is to show that for underfitted models, the goodness-of-fit part dominates, while for overfitted models, the model complexity part dominates.

Let $\hat{Q}_n^\lambda = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \hat{\alpha}_k^\lambda - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}^\lambda)$, where $(\hat{\alpha}^\lambda, \hat{\boldsymbol{\beta}}^\lambda)$ is the two-step LLA estimator to the folded concave penalized CQR (6) with lasso initialization. Also, let $Q_n^* = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k^* - \mathbf{x}_i^\top \boldsymbol{\beta}^*)$. For any $A \subset \{1, 2, \dots, p\}$, let $(\hat{\alpha}^A, \hat{\boldsymbol{\beta}}^A)$ be the estimator obtained by fitting the canonical CQR to model A , i.e.,

$$(\hat{\alpha}^A, \hat{\boldsymbol{\beta}}^A) = \arg \min_{\alpha, \boldsymbol{\beta} : \beta_{A^c} = \mathbf{0}} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^\top \boldsymbol{\beta}). \quad (19)$$

Define $\hat{Q}_n^A = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \hat{\alpha}_k^A - \mathbf{x}_i^\top \hat{\boldsymbol{\beta}}^A)$. For any $\lambda \in \Xi_n$, recall that $\hat{A}_\lambda = \{1 \leq j \leq p : \hat{\beta}_j^\lambda \neq 0\}$ corresponds to the active set of the two-step LLA estimator $(\hat{\alpha}^\lambda, \hat{\boldsymbol{\beta}}^\lambda)$. By optimality of $(\hat{\alpha}^{\hat{A}_\lambda}, \hat{\boldsymbol{\beta}}^{\hat{A}_\lambda})$ in (39), we have $\hat{Q}_n^{\hat{A}_\lambda} \leq \hat{Q}_n^\lambda$.

Let $\mathcal{G}_n^+ = \{A : A \supset \mathcal{A}, A \neq \mathcal{A}, |A| \leq J_n\}$. It can be seen that $\{\hat{A}_\lambda : \lambda \in \Xi_n^+\} \subset \mathcal{G}_n^+$. For $r > 0$ and $A \in \mathcal{G}_n^+$, let $B_A(r) = \{(\boldsymbol{\delta}, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p : \|\boldsymbol{\delta}\|^2 + \|\Delta_A\|^2 \leq r^2, \Delta_{A^c} = \mathbf{0}\}$ and $S_A(r) = \{(\boldsymbol{\delta}, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p : \|\boldsymbol{\delta}\|^2 + \|\Delta_A\|^2 = r^2, \Delta_{A^c} = \mathbf{0}\}$.

Case I: overfitted models. By Theorem 2, under the assumptions of this theorem, we have

$$P(\hat{A}_{\lambda_n} \neq \mathcal{A}) \rightarrow o(1) \quad \text{as } n \rightarrow \infty.$$

Therefore, for any $\lambda \in \Xi_n^+$, we have

$$\begin{aligned} & \Pr\left(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0\right) \\ &= \Pr\left(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0, \hat{A}_{\lambda_n} = \mathcal{A}\right) \\ & \quad + \Pr\left(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0, \hat{A}_{\lambda_n} \neq \mathcal{A}\right) \\ &= \Pr\left(\inf_{\lambda \in \Xi_n^+} \left[(\hat{Q}_n^\lambda - \hat{Q}_n^{\mathcal{A}}) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n}\right] > 0\right) + o(1) \\ &\geq \Pr\left(\inf_{\lambda \in \Xi_n^+} \left[(\hat{Q}_n^{\hat{A}_\lambda} - \hat{Q}_n^{\mathcal{A}}) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n}\right] > 0\right) + o(1), \end{aligned}$$

where the last inequality follows from the fact that $\widehat{Q}_n^\lambda \geq \widehat{Q}_n^{\hat{\lambda}}$. Moreover, note that $\widehat{Q}_n^{\hat{\lambda}} \leq \widehat{Q}_n^{\mathcal{A}} \leq Q_n^*$ due to inclusion $\mathcal{A} \subset \hat{A}_\lambda$.

Let $F(\delta, \Delta) = Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)$ for $(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p$. For each model $A \in \mathcal{G}_n^+$, let $r_A^* = 16K(M + \bar{\xi}^{1/2})(f\underline{\zeta})^{-1} \sqrt{(|A|+1)\log(p)/n}$. If we can show that $\inf_{(\delta, \Delta) \in S_A(r_A^*)} F(\delta, \Delta) > 0$, then by convexity of $\rho_\tau(\cdot)$, we must have $\|\hat{\alpha}^A - \alpha^*\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \leq (r_A^*)^2$. Indeed, by Knight's identity (see (28) of the appendix) and the mean value theorem, we have

$$\begin{aligned} & \inf_{(\delta, \Delta) \in S_A(r_A^*)} \mathbb{E}[F(\delta, \Delta)] \\ &= \inf_{(\delta, \Delta) \in S_A(r_A^*)} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^\top \Delta} [F(\alpha_k^* + t) - F(\alpha_k^*)] dt \\ &= \inf_{(\delta, \Delta) \in S_A(r_A^*)} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^\top \Delta} [tf(\alpha_k^* + \bar{u}_{ik,t})] dt. \end{aligned} \quad (20)$$

Note that J_n satisfies $J_n^2 \log(p)/n = o(1)$. Therefore, for any $1 \leq i \leq n, 1 \leq k \leq K$ and $(\delta, \Delta) \in S_A(r_A^*)$, we have

$$\begin{aligned} |\delta_k + \mathbf{x}_i^\top \Delta| &\leq \sqrt{1 + \|\mathbf{x}_{iA}\|^2} \cdot \sqrt{\delta_k^2 + \|\Delta_A\|_2^2} \\ &\leq Mr_A^* \sqrt{|A|+1} = 16KM(M + \bar{\xi}^{1/2})(f\underline{\zeta})^{-1} (|A|+1) \sqrt{\log(p)/n} \\ &\leq 16KM(M + \bar{\xi}^{1/2})(f\underline{\zeta})^{-1} (J_n + 1) \sqrt{\log(p)/n} = o(\mathcal{Q}_0). \end{aligned}$$

It follows from condition (C0) and (40) that

$$\inf_{(\delta, \Delta) \in S_A(r_A^*)} \mathbb{E}[F(\delta, \Delta)] \geq \inf_{(\delta, \Delta) \in S_A(r_A^*)} \frac{f}{2nK} \sum_{k=1}^K \sum_{i=1}^n (\delta_k + \mathbf{x}_i^\top \Delta)^2 \geq \frac{f}{2K} \underline{\zeta} (r_A^*)^2.$$

Therefore, by Lemma 9, with probability at least $1 - \exp\{-nt^2/[32\bar{\zeta}(r_A^*)^2]\}$, we have

$$\begin{aligned} \inf_{(\delta, \Delta) \in S_A(r_A^*)} F(\delta, \Delta) &\geq \inf_{(\delta, \Delta) \in S_A(r_A^*)} \mathbb{E}[F(\delta, \Delta)] - z(A, r_A^*) \\ &\geq \frac{f}{2K} \underline{\zeta} (r_A^*)^2 - 4Mr_A^* \sqrt{\frac{|A|+1}{n}} - t. \end{aligned}$$

Now take $t = 8r_A^* \sqrt{\bar{\zeta}(|A|+1)\log(p)/n}$. It follows that for each $A \in \mathcal{G}_n^+$, with probability at least $b_n^A = 1 - \exp\{-2(|A|+1)\log(p)\}$, we have

$$\inf_{(\delta, \Delta) \in S_A(r_A^*)} F(\delta, \Delta) \geq \frac{f\underline{\zeta}(r_A^*)^2}{2K} - r_A^* \sqrt{\frac{|A|+1}{n}} (8\sqrt{\bar{\zeta}\log(p)} + 4M) > 0,$$

which immediately implies that

$$\|\hat{\alpha}^A - \alpha^*\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \leq (r_A^*)^2.$$

Now by the Bonferroni inequality, we have $\|\hat{\alpha}^A - \alpha\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \leq (r_A^*)^2$ for all $A \in \mathcal{G}_n^+$ simultaneously with probability at least

$$\begin{aligned}
 b_n &= 1 - \sum_{|A|=s+1}^{J_n} \binom{p-s}{|A|-s} (1 - b_n^A) \\
 &= 1 - \sum_{k=1}^{J_n-s} \binom{p-s}{k} \exp\{-2(k+s+1)\log(p)\} \\
 &\geq 1 - p^{-2(s+1)} \sum_{k=1}^{p-s} \binom{p-s}{k} \left(\frac{1}{p^2}\right)^k \\
 &= 1 - p^{-2(s+1)} \left[\left(1 + \frac{1}{p^2}\right)^{p-s} - 1 \right] \rightarrow 1 \quad \text{as } p \rightarrow \infty.
 \end{aligned} \tag{21}$$

Now we derive the upper bound for $\sup_{A \in \mathcal{G}_n^+} |\hat{Q}_n^A - Q_n^*|$. Let $\hat{\delta}^A = \hat{\alpha}^A - \alpha^*$ and $\hat{\Delta}^A = \hat{\beta}^A - \beta^*$. Observe that

$$\begin{aligned}
 |\hat{Q}_n^A - Q_n^*| &= \left| \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \{ \rho_{\tau_k}(y_i - \hat{\alpha}_k^A - \mathbf{x}_i^T \hat{\beta}^A) - \rho_{\tau_k}(y_i - \alpha_k^* - \mathbf{x}_i^T \beta^*) \} \right| \\
 &\leq |\mathbb{E}\{F(\hat{\delta}^A, \hat{\beta}^A)\}| + z(A, r_A^*).
 \end{aligned}$$

Similarly, we have

$$\begin{aligned}
 |\mathbb{E}\{F(\hat{\delta}^A, \hat{\beta}^A)\}| &= \left| \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\hat{\delta}_k^A + \mathbf{x}_i^T \hat{\Delta}^A} [F(\alpha_k^* + t) - F(\alpha_k^*)] dt \right| \\
 &\leq \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{|\hat{\delta}_k^A + \mathbf{x}_i^T \hat{\Delta}^A|} [tf(\alpha_k^* + \bar{u}_{ik,t})] dt \\
 &\leq \frac{\bar{f}}{2nK} \sum_{i=1}^n \sum_{k=1}^K (\hat{\delta}^A + \mathbf{x}_i^T \hat{\Delta}^A)^2 \\
 &\leq \frac{1}{2} \bar{f} \bar{\zeta} (\|\hat{\delta}^A\|_2^2 + \|\hat{\Delta}^A\|_2^2) \leq \frac{1}{2} \bar{f} \bar{\zeta} (r_A^*)^2.
 \end{aligned}$$

It follows that with probability at least b_n ,

$$\begin{aligned}
 |\hat{Q}_n^A - Q_n^*| &\leq \frac{1}{2} \bar{f} \bar{\zeta} (r_A^*)^2 + 4Mr_A^* \sqrt{\frac{|A|+1}{n}} + 8r_A^* \sqrt{\frac{\bar{\zeta}(|A|+1)\log(p)}{n}} \\
 &\leq \frac{128K^2(M + \bar{\zeta}^{1/2})^2}{\underline{f}\bar{\zeta}} \frac{(|A|+1)\log(p)}{n}
 \end{aligned}$$

holds for all $A \in \mathcal{G}_n^+$. Now going back to BIC, we have

$$\begin{aligned}
 &\Pr\left(\inf_{\lambda \in \Xi_n^+} \left[(\hat{Q}_n^{\hat{\lambda}} - \hat{Q}_n^{\mathcal{A}}) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n} \right] > 0 \right) \\
 &\geq \Pr\left(\frac{C_n \log(p)}{n} - \sup_{A \in \mathcal{G}_n^+} \frac{\hat{Q}_n^{\mathcal{A}} - \hat{Q}_n^A}{(|A| - s)} > 0 \right).
 \end{aligned}$$

Therefore, with probability at least b_n , we have

$$\sup_{A \in \mathcal{G}_n^+} \frac{\hat{Q}_n^{\mathcal{A}} - \hat{Q}_n^A}{(|A| - s)} \leq \sup_{A \in \mathcal{G}_n^+} \frac{\hat{Q}_n^* - \hat{Q}_n^A}{(|A| - s)} = \mathcal{O}_P\left(\frac{s \log(p)}{n}\right).$$

Since $s = \mathcal{O}(1)$ and C_n diverges with n , we have $s \log(p)/n = \mathcal{O}(C_n \log(p)/n)$. It follows that

$$\Pr\left(\frac{C_n \log(p)}{n} - \sup_{A \in \mathcal{G}_n^+} \frac{\hat{Q}_n^{\mathcal{A}} - \hat{Q}_n^A}{(|A| - s)} > 0 \right) \rightarrow 1 \quad \text{as } n \rightarrow \infty,$$

which implies that $\Pr(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0) \rightarrow 1$ as $n \rightarrow \infty$.

Case II: underfitted models. For any $\lambda \in \Xi_n^-$, similar to Case I, we have

$$\begin{aligned} & \Pr\left(\inf_{\lambda \in \Xi_n^-} [\text{BIC}^H(\lambda) - \text{BIC}^H(\lambda_n)] > 0\right) \\ & \geq \Pr\left(\inf_{\lambda \in \Xi_n^-} \left[(\widehat{Q}_n^{\hat{\lambda}} - \widehat{Q}_n^{\mathcal{A}}) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n} \right] > 0\right) + o(1). \end{aligned}$$

Define $\text{BIC}^H(A) = \widehat{Q}_n^A + |A|C_n \log(p)/n$ and let $\mathcal{G}_n^- = \{A : |A| \leq J_n, \mathcal{A} \not\subset A\}$. We can see that $\{\hat{A}_\lambda : \lambda \in \Xi_n^-\} \subset \mathcal{G}_n^-$. It suffices to show $\inf_{A \in \mathcal{G}_n^-} \text{BIC}^H(A) > \text{BIC}^H(\mathcal{A})$ with probability tending to one as $n \rightarrow \infty$. For any $A \in \mathcal{G}_n^-$, let $\bar{A} = A \cup \mathcal{A}$. Let $\theta = \min_{j \in \mathcal{A}} |\beta_j^*|$. Since $A \not\supset \mathcal{A}$, we must have $\|\hat{\alpha}^A - \alpha^*\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \geq \theta^2$. However, since $\bar{A} \supset \mathcal{A}$ and $|\bar{A}| \leq 2J_n$, using Lemma 9, we can similarly show as in Case I that $\|\hat{\alpha}^{\bar{A}} - \alpha^*\|_2^2 + \|\hat{\beta}^{\bar{A}} - \beta^*\|_2^2 \leq \theta^2$ with probability at least $b_n^{\bar{A}} = 1 - \exp\{-2(|\bar{A}| + 1) \log(p)\}$ as long as

$$\theta > 8K(f\underline{\zeta})^{-1} (2\sqrt{\bar{\zeta} \log(p)} + M) \sqrt{\frac{2J_n + 1}{n}},$$

which is implied by the assumption $\sqrt{J_n \log(p)/n} = o(\theta)$. It then follows that $\|\hat{\alpha}^{\bar{A}} - \alpha^*\|_2^2 + \|\hat{\beta}^{\bar{A}} - \beta^*\|_2^2 \leq \theta^2$ holds for all $A \in \mathcal{G}_n^-$ with probability at least $\tilde{b}_n \rightarrow 1$ as $n \rightarrow \infty$, where

$$\tilde{b}_n = 1 - \sum_{|\bar{A}|=s+1}^{2J_n} \binom{p-s}{|\bar{A}|-s} (1 - b_n^{\bar{A}}) \geq 1 - p^{-2(s+1)} \left[\left(1 + \frac{1}{p^2}\right)^{p-s} - 1 \right].$$

Therefore, there exists $a \in [0, 1]$, $\bar{\alpha}^{\bar{A}} = a\hat{\alpha}^A + (1-a)\hat{\alpha}^{\bar{A}}$ and $\bar{\beta}^{\bar{A}} = a\hat{\beta}^A + (1-a)\hat{\beta}^{\bar{A}}$ such that $\|\bar{\alpha}^{\bar{A}} - \alpha^*\|_2^2 + \|\bar{\beta}^{\bar{A}} - \beta^*\|_2^2 = \theta^2$. By convexity of ρ_τ and the fact that $\widehat{Q}_n^A \geq \widehat{Q}_n^{\bar{A}}$, we have $\bar{Q}_n^{\bar{A}} = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \bar{\alpha}_k^{\bar{A}} - \mathbf{x}_i^T \bar{\beta}^{\bar{A}}) \leq \widehat{Q}_n^A$. Note that $\widehat{Q}_n^{\bar{A}} \leq \widehat{Q}_n^{\mathcal{A}} \leq Q_n^*$. It follows that $\bar{Q}_n^{\bar{A}} - \widehat{Q}_n^{\bar{A}} \geq \bar{Q}_n^{\bar{A}} - Q_n^*$. For ease of notation, let $\bar{\delta}^{\bar{A}} = \bar{\alpha}^{\bar{A}} - \alpha^*$ and $\bar{\Delta}^{\bar{A}} = \bar{\beta}^{\bar{A}} - \beta^*$. It can be seen that

$$\bar{Q}_n^{\bar{A}} - Q_n^* \geq \mathbb{E}[F(\bar{\delta}^{\bar{A}}, \bar{\Delta}^{\bar{A}})] - z(\bar{A}, \theta).$$

Following similar arguments from Case I and noting that the support of $\bar{\beta}^{\bar{A}}$ is a subset of \bar{A} , we can show that with probability at least \tilde{b}_n , for all $A \in \mathcal{G}_n^-$, we have

$$\bar{Q}_n^{\bar{A}} - Q_n^* \geq \frac{f}{2K} \underline{\zeta} \theta^2 - 4\theta \sqrt{\frac{|\bar{A}| + 1}{n}} (2\sqrt{\bar{\zeta} \log(p)} + M).$$

Now we have

$$\begin{aligned} \text{BIC}^H(A) - \text{BIC}^H(\bar{A}) &= (\widehat{Q}_n^A - \widehat{Q}_n^{\bar{A}}) + (|A| - |\bar{A}|) \frac{C_n \log(p)}{n} \\ &\geq (\bar{Q}_n^{\bar{A}} - Q_n^*) - \frac{C_n s \log(p)}{n}. \end{aligned}$$

Since $\sqrt{C_n s \log(p)/n} = o(\theta)$, it can be seen that with probability tending to one, we have $\inf_{A \in \mathcal{G}_n^-} \text{BIC}^H(A) - \text{BIC}^H(\bar{A}) > 0$. Following similar arguments as in Case I, we can show $\text{BIC}^H(\bar{A}) \geq \inf_{S \supset \mathcal{A}, |S| \leq 2J_n} \text{BIC}^H(S) \geq \text{BIC}^H(\mathcal{A})$ with probability tending to one. Case II then follows by noting that

$$\begin{aligned} & \inf_{A \in \mathcal{G}_n^-} [\text{BIC}^H(A) - \text{BIC}^H(\mathcal{A})] \\ &= \inf_{A \in \mathcal{G}_n^-} [\text{BIC}^H(A) - \text{BIC}^H(\bar{A}) + \text{BIC}^H(\bar{A}) - \text{BIC}^H(\mathcal{A})] \\ &\geq \inf_{A \in \mathcal{G}_n^-} [\text{BIC}^H(A) - \text{BIC}^H(\bar{A})]. \end{aligned}$$

□

Proof of Theorem 4. See Section A of the appendix. □

Proof of Lemma 2. See Lemma 1 of [23]. □

APPENDIX A

PROOFS

Proof of Lemma 3. Let $\zeta = (\zeta_1, \dots, \zeta_K)^\top$ and $\xi = (\xi_1, \dots, \xi_p)^\top$, where

$$\zeta_k = -\frac{1}{nK} \sum_{i=1}^n [\tau_k - I(\varepsilon_i \leq \alpha_k^*)], \quad 1 \leq k \leq K,$$

and

$$\xi_j = -\frac{1}{nK} \sum_{k=1}^K \sum_{i=1}^n [\tau_k - I(\varepsilon_i \leq \alpha_k^*)] x_{ij}, \quad 1 \leq j \leq p.$$

Note that $(\zeta^\top, \xi^\top)^\top \in \partial Q_n(\alpha^*, \beta^*)$, where the subdifferential is taken with respect to α and β . By convexity of $Q_n(\alpha, \beta)$ and optimality of $(\widehat{\alpha}_\lambda, \widehat{\beta}_\lambda)$, we have

$$\begin{aligned} 0 &\geq Q_n(\widehat{\alpha}_\lambda, \widehat{\beta}_\lambda) - Q_n(\alpha^*, \beta^*) + \lambda(\|\widehat{\beta}_\lambda\|_1 - \|\beta^*\|_1) \\ &\geq \zeta^\top(\widehat{\alpha}_\lambda - \alpha^*) + \xi^\top(\widehat{\beta}_\lambda - \beta^*) + \lambda(\|\widehat{\beta}_\lambda\|_1 - \|\beta^*\|_1) \\ &\geq -\|\zeta\|_\infty \cdot \|\widehat{\alpha}_\lambda - \alpha^*\|_1 - \|\xi\|_\infty \cdot \|\widehat{\beta}_\lambda - \beta^*\|_1 \\ &\quad + \lambda(\|\widehat{\beta}_{\lambda, \mathcal{A}^c} - \beta_{\mathcal{A}^c}^*\|_1 - \|\widehat{\beta}_{\lambda, \mathcal{A}} - \beta_{\mathcal{A}}^*\|_1), \end{aligned}$$

which implies that

$$\begin{aligned} (\lambda - \|\xi\|_\infty) \|\widehat{\beta}_{\lambda, \mathcal{A}^c} - \beta_{\mathcal{A}^c}^*\|_1 &\leq (\lambda + \|\xi\|_\infty) \|\widehat{\beta}_{\lambda, \mathcal{A}} - \beta_{\mathcal{A}}^*\|_1 \\ &\quad + \|\zeta\|_\infty \cdot \|\widehat{\alpha}_\lambda - \alpha^*\|_1. \end{aligned} \tag{22}$$

Under event $\mathcal{E} = \{\|\zeta\|_\infty \leq 3\lambda/(2K), \|\xi\|_\infty \leq \lambda/2\}$, it follows from (22) that

$$\|\widehat{\Delta}_{\mathcal{A}^c}^\lambda\|_1 \leq 3\|\widehat{\Delta}_{\mathcal{A}}^\lambda\|_1 + \frac{3}{K}\|\widehat{\delta}^\lambda\|_1.$$

The lemma then follows from Hoeffding's inequality

$$\begin{aligned} \Pr(\mathcal{E}) &\geq 1 - \Pr\left(\|\zeta\|_\infty > \frac{3\lambda}{2K}\right) - \Pr\left(\|\xi\|_\infty > \frac{\lambda}{2}\right) \\ &\geq 1 - \sum_{k=1}^K \Pr\left(\left|-\frac{1}{nK} \sum_{i=1}^n [\tau_k - I(\varepsilon_i \leq \alpha_k^*)]\right| > \frac{3\lambda}{2K}\right) \\ &\quad - \sum_{j=1}^p \Pr\left(\left|-\frac{1}{nK} \sum_{i=1}^n x_{ij} \sum_{k=1}^K [\tau_k - I(\varepsilon_i \leq \alpha_k^*)]\right| > \frac{\lambda}{2}\right) \\ &\geq 1 - 2K \exp\left(-\frac{9n\lambda^2}{2}\right) - 2p \exp\left(-\frac{n\lambda^2}{2M_0}\right). \end{aligned}$$

This proves the lemma. □

Proof of Lemma 4. First, let us show that the check loss $\rho_\tau(\cdot)$ is Lipschitz continuous with Lipschitz constant $\max(\tau, 1 - \tau)$. To see it, note that for any $u_1, u_2 \in \mathbb{R}$, we have

$$\begin{aligned} |\rho_\tau(u_1) - \rho_\tau(u_2)| &= |(\tau - 0.5)(u_1 - u_2) + 0.5(|u_1| - |u_2|)| \\ &\leq (|\tau - 0.5| + 0.5)|u_1 - u_2| = \max(\tau, 1 - \tau)|u_1 - u_2|. \end{aligned}$$

Now let $\delta = \alpha - \alpha^*$, $\Delta = \beta - \beta^*$, and define

$$\begin{aligned} U_i(\delta, \Delta) &= \frac{1}{K} \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta) - \frac{1}{K} \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k^* - \mathbf{x}_i^T \beta^*) \\ &= \frac{1}{K} \sum_{k=1}^K \rho_{\tau_k}(r_{ik}^* - \delta_k - \mathbf{x}_i^T \Delta) - \frac{1}{K} \sum_{k=1}^K \rho_{\tau_k}(r_{ik}^*), \end{aligned}$$

where $r_{ik}^* = y_i - \alpha_k^* - \mathbf{x}_i^T \beta^* = \varepsilon_i - \alpha_k^*$, $1 \leq i \leq n$, $1 \leq k \leq K$. It follows immediately that

$$e(r) = \sup_{(\delta, \Delta) \in \mathcal{C}_r} \left| \frac{1}{n} \sum_{i=1}^n [U_i(\delta, \Delta) - \mathbb{E}U_i(\delta, \Delta)] \right|.$$

By Lipschitz continuity of the check loss, it follows that

$$\begin{aligned} |U_i(\delta, \Delta)| &\leq \frac{1}{K} \sum_{k=1}^K |\rho_{\tau_k}(r_{ik}^* - \delta_k - \mathbf{x}_i^T \Delta) - \rho_{\tau_k}(r_{ik}^*)| \\ &\leq \frac{1}{K} \sum_{k=1}^K \max(\tau_k, 1 - \tau_k) |\delta_k + \mathbf{x}_i^T \Delta| \leq \frac{1}{K} \sum_{k=1}^K |\delta_k + \mathbf{x}_i^T \Delta|, \quad 1 \leq i \leq n. \end{aligned} \quad (23)$$

Now applying Massart's concentration inequality (Theorem 14.2, [25]), we obtain

$$\Pr(e(r) \geq \mathbb{E}[e(r)] + t) \leq \exp\left(-\frac{n^2 t^2}{8b_n^2(r)}\right), \quad (24)$$

where $b_n^2(r) = \sup_{(\delta, \Delta) \in \mathcal{C}_r} \sum_{i=1}^n \text{var}(U_i(\delta, \Delta))$. First, we derive the upper bound on $b_n^2(r)$. Note that by (23) and Cauchy–Schwarz inequality

$$\begin{aligned} b_n^2(r) &= \sup_{(\delta, \Delta) \in \mathcal{C}_r} \sum_{i=1}^n \mathbb{E}[U_i(\delta, \Delta) - \mathbb{E}U_i(\delta, \Delta)]^2 \leq 4 \sup_{(\delta, \Delta) \in \mathcal{C}_r} \sum_{i=1}^n \left[\sum_{k=1}^K \frac{1}{K} |\delta_k + \mathbf{x}_i^T \Delta| \right]^2 \\ &\leq 4 \sup_{(\delta, \Delta) \in \mathcal{C}_r} \sum_{i=1}^n \left(\sum_{k=1}^K \frac{1}{K} \right) \left[\sum_{k=1}^K \frac{1}{K} (\delta_k + \mathbf{x}_i^T \Delta)^2 \right] \leq 4nr^2. \end{aligned}$$

Next, we derive the upper bound on $\mathbb{E}[e(r)]$. By applying the symmetrization procedure [26] and the contraction principle [27], we have

$$\begin{aligned} \mathbb{E}[e(r)] &\leq 2\mathbb{E} \left[\sup_{(\delta, \Delta) \in \mathcal{C}_r} \frac{1}{n} \sum_{i=1}^n \xi_i U_i(\delta, \Delta) \right] \\ &\leq \frac{2}{nK} \sum_{k=1}^K \mathbb{E} \left[\sup_{(\delta, \Delta) \in \mathcal{C}_r} \left| \sum_{i=1}^n \xi_i \{ \rho_{\tau_k}(r_{ik}^* - \delta_k - \mathbf{x}_i^T \Delta) - \rho_{\tau_k}(r_{ik}^*) \} \right| \right] \\ &\leq \frac{4}{nK} \sum_{k=1}^K \mathbb{E} \left[\sup_{(\delta, \Delta) \in \mathcal{C}_r} \left| \sum_{i=1}^n \xi_i (\delta_k + \mathbf{x}_i^T \Delta) \right| \right], \end{aligned} \quad (25)$$

where ξ_1, \dots, ξ_n are i.i.d. Rademacher random variables that satisfy $\Pr(\xi_i = -1) = \Pr(\xi_i = 1) = 1/2$ and that are independent of $\varepsilon_1, \dots, \varepsilon_n$.

For $(\delta, \Delta) \in \mathcal{C}_r$, by condition (C1) and the Cauchy–Schwarz inequality, we have

$$r^2 \geq \frac{\kappa_0}{K} (\|\delta\|_2^2 + K\|\Delta_{\mathcal{A}}\|_2^2) \geq \frac{\kappa_0}{K^2} \|\delta\|_1^2 + \frac{\kappa_0}{s} \|\Delta_{\mathcal{A}}\|_1^2, \quad (26)$$

which implies that $\|\delta\|_1 \leq rK/\sqrt{\kappa_0}$ and $\|\Delta_{\mathcal{A}}\|_1 \leq r\sqrt{s/\kappa_0}$. Now let $\xi = (\xi_1, \dots, \xi_n)^T$. Note that for any $t \in \mathbb{R}$, we have by Taylor expansion

$$\begin{aligned} \mathbb{E}[\exp(tX_j^T \xi)] &= \prod_{i=1}^n \left[\frac{1}{2} (e^{tx_{ij}} + e^{-tx_{ij}}) \right] \\ &\leq \prod_{i=1}^n \exp\left(\frac{1}{2} t^2 x_{ij}^2\right) = \exp\left(\frac{t^2}{2} \sum_{i=1}^n x_{ij}^2\right), \quad 0 \leq j \leq p. \end{aligned}$$

Letting $t > 0$, by Jensen's inequality, we have

$$\begin{aligned}
\exp(t\mathbb{E}[\|\mathbb{X}^T \xi\|_\infty]) &= \exp\left(t\mathbb{E}\max_{0 \leq j \leq p} |X_j^T \xi|\right) \leq \mathbb{E}\exp\left(t\max_{0 \leq j \leq p} |X_j^T \xi|\right) \\
&= \mathbb{E}\left[\max_{0 \leq j \leq p} \exp(t|X_j^T \xi|)\right] \leq \mathbb{E}\max_{0 \leq j \leq p} \left(e^{tX_j^T \xi} + e^{-tX_j^T \xi}\right) \\
&\leq \sum_{j=0}^p \mathbb{E}\left(e^{tX_j^T \xi} + e^{-tX_j^T \xi}\right) \leq 2 \sum_{j=0}^p \exp\left(\frac{t^2}{2}\|X_j\|_2^2\right) \\
&\leq 2(1+p) \exp\left(\frac{t^2}{2}\max_{0 \leq j \leq p} \|X_j\|_2^2\right) = 2(1+p) \exp\left(\frac{1}{2}nM_0 t^2\right),
\end{aligned}$$

which implies that

$$\mathbb{E}(\|\mathbb{X}^T \xi\|_\infty) \leq \frac{1}{t} [\log 2 + \log(1+p)] + \frac{nM_0}{2} t, t > 0.$$

Taking $t = \sqrt{2[\log 2 + \log(1+p)]/(nM_0)}$ and noting that $p \geq 3$ by condition (C0), we obtain

$$\mathbb{E}(\|\mathbb{X}^T \xi\|_\infty) \leq \sqrt{2nM_0[\log 2 + \log(1+p)]} \leq \sqrt{2M_0} \cdot \sqrt{n(1+\log p)}. \quad (27)$$

It then follows from (25), (27) and Hölder's inequality that

$$\begin{aligned}
\mathbb{E}[e(r)] &\leq \frac{4}{nK} \mathbb{E}(\|\mathbb{X}^T \xi\|_\infty) \cdot \sup_{(\delta, \Delta) \in \mathcal{C}_r} \sum_{k=1}^K (|\delta_k| + \|\Delta\|_1) \\
&\leq \frac{4\sqrt{2M_0}}{n} \sqrt{n(1+\log p)} \sup_{(\delta, \Delta) \in \mathcal{C}_r} (K^{-1}\|\delta\|_1 + \|\Delta\|_1) \\
&\leq \frac{4\sqrt{2M_0}}{n} \sqrt{n(1+\log p)} \sup_{(\delta, \Delta) \in \mathcal{C}_r} [4K^{-1}\|\delta\|_1 + 4\|\Delta_{\mathcal{A}}\|_1] \\
&\leq 16 \sqrt{\frac{2M_0}{\kappa_0}} \sqrt{\frac{1+\log p}{n}} (\sqrt{s} + 1)r.
\end{aligned}$$

The lemma then follows from (24). \square

Proof of Lemma 5. By Knight's identity [2], we have for any two scalars $r \neq 0$ and s ,

$$|r-s| - |r| = -s[I(r > 0) - I(r < 0)] + 2 \int_0^s [I(r \leq t) - I(r \leq 0)] dt.$$

It follows that for any $\tau \in (0, 1)$, when $r \neq 0$,

$$\begin{aligned}
\rho_\tau(r-s) - \rho_\tau(r) &= (\tau - 0.5)[(r-s) - r] + 0.5[|r-s| - |r|] \\
&= (0.5 - \tau)s - 0.5s[I(r > 0) - I(r < 0)] + \int_0^s [I(r \leq t) - I(r \leq 0)] dt \\
&= s[I(r < 0) - \tau] + \int_0^s [I(r \leq t) - I(r \leq 0)] dt.
\end{aligned} \quad (28)$$

Let $r_{ik}^* = y_i - \alpha_k^* - \mathbf{x}_i^\top \beta^* = \varepsilon_i - \alpha_k^*$, $1 \leq i \leq n$, $1 \leq k \leq K$. Recall that ε has a density with respect to the Lebesgue measure. By condition (C0), identity (28) and the mean value theorem, we have for some $\bar{u}_{ik,t}$ between 0 and t ,

$$\begin{aligned} & \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \\ &= \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^\top \Delta} [F(\alpha_k^* + t) - F(\alpha_k^*)] dt \\ &= \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^\top \Delta} \left[t f(\alpha_k^*) + \frac{t^2}{2} f'(\alpha_k^* + \bar{u}_{ik,t}) \right] dt \\ &\geq \frac{f}{2nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2 - \frac{\bar{f}'}{6nK} \sum_{i=1}^n \sum_{k=1}^K |\delta_k + \mathbf{x}_i^\top \Delta|^3. \end{aligned} \quad (29)$$

For $(\delta, \Delta) \in \mathcal{C}$, note that if

$$\left[\frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2 \right]^{1/2} \leq \frac{4q}{K^{1/2} \underline{f}^{1/2}}, \quad (30)$$

then by condition (C2), this implies that

$$\frac{\bar{f}'}{6nK} \sum_{i=1}^n \sum_{k=1}^K |\delta_k + \mathbf{x}_i^\top \Delta|^3 \leq \frac{f}{4nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2,$$

which, together with (29), implies that for all $(\delta, \Delta) \in \mathcal{C}_{4q(K\underline{f})^{-1/2}}$,

$$\mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \geq \frac{f}{4nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2.$$

To show that the lemma holds for all $(\delta, \Delta) \in \mathcal{C}$, define

$$\begin{aligned} r_{\mathcal{C}} &= \sup_{r>0} \left\{ r: \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \right. \\ &\quad \left. \geq \frac{f}{4nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2, \forall (\delta, \Delta) \in \mathcal{C}_r \right\}. \end{aligned}$$

By previous arguments, we must have $r_{\mathcal{C}} \geq 4q(K\underline{f})^{-1/2}$. Now for any $(\delta, \Delta) \in \mathcal{C}$, let $r^2 = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2$. If $r < r_{\mathcal{C}}$, then by the definition of $r_{\mathcal{C}}$, we have

$$\mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \geq \frac{f}{4nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2. \quad (31)$$

If instead $r \geq r_{\mathcal{C}}$, let $\delta' = r_{\mathcal{C}} \delta / r$ and $\Delta' = r_{\mathcal{C}} \Delta / r$. It can be seen immediately that $(nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K (\delta'_k + \mathbf{x}_i^\top \Delta')^2 = r_{\mathcal{C}}^2$.

By convexity of Q_n , we have

$$\begin{aligned} & \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \\ &\geq \frac{r}{r_{\mathcal{C}}} \mathbb{E}[Q_n(\alpha^* + \delta', \beta^* + \Delta') - Q_n(\alpha^*, \beta^*)] \\ &\geq \frac{r}{r_{\mathcal{C}}} \frac{f}{4} r_{\mathcal{C}}^2 \geq q \left(\frac{f}{K} \right)^{1/2} \left[\frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^\top \Delta)^2 \right]^{1/2}. \end{aligned} \quad (32)$$

The lemma then follows from (31) and (32). \square

Proof of Lemma 6. As with the proof of Lemma 4, define

$$U_i(\delta, \Delta) = \frac{1}{K} \sum_{k=1}^K \rho_{\tau_k}(r_{ik}^* - \delta_k - \mathbf{x}_i^\top \Delta) - \frac{1}{K} \sum_{k=1}^K \rho_{\tau_k}(r_{ik}^*),$$

where $r_{ik}^* = y_i - \alpha_k^* - \mathbf{x}_i^T \beta^* = \varepsilon_i - \alpha_k^*$, $1 \leq i \leq n$, $1 \leq k \leq K$. Now applying Massart's concentration inequality, we get

$$\Pr(z(r) \geq \mathbb{E}[z(r)] + t) \leq \exp\left(-\frac{n^2 t^2}{8b_n^2(r)}\right), \quad (33)$$

where $b_n^2(r) = \sup_{(\delta, \Delta) \in B_{\mathcal{S}}(r)} \sum_{i=1}^n \text{var}(U_i(\delta, \Delta))$. For ease of notation, let $\Delta_{\mathcal{S}}^k = (\delta_k, \Delta_{\mathcal{S}}^T)^T$, $1 \leq k \leq K$. It follows from Lipschitz continuity of the check loss that

$$\begin{aligned} b_n^2(r) &\leq \frac{4}{K} \sup_{(\delta, \Delta) \in B_{\mathcal{S}}(r)} \sum_{k=1}^K \sum_{i=1}^n (\delta_k + \mathbf{x}_i^T \Delta)^2 \\ &= \frac{4}{K} \sup_{(\delta, \Delta) \in B_{\mathcal{S}}(r)} \sum_{k=1}^K (\Delta_{\mathcal{S}}^k)^T \mathbb{X}_{\mathcal{S}_0}^T \mathbb{X}_{\mathcal{S}_0} \Delta_{\mathcal{S}}^k \\ &\leq \frac{4n}{K} \sup_{(\delta, \Delta) \in B_{\mathcal{S}}(r)} \sum_{k=1}^K \bar{\mu} [\delta_k^2 + \|\Delta_{\mathcal{S}}\|_2^2] \leq 4n\bar{\mu}r^2. \end{aligned}$$

Moreover, by the symmetrization procedure and the contraction principle, we obtain

$$\begin{aligned} \mathbb{E}[z(r)] &\leq \frac{4}{nK} \sum_{k=1}^K \mathbb{E} \left[\sup_{(\delta, \Delta) \in B_{\mathcal{S}}(r)} \left| \sum_{i=1}^n \xi_i (\delta_k + \mathbf{x}_i^T \Delta_{\mathcal{S}}) \right| \right] \\ &\leq \frac{4}{nK} \mathbb{E}(\|\mathbb{X}_{\mathcal{S}_0}^T \xi\|_2) \cdot \sup_{(\delta, \Delta) \in B_{\mathcal{S}}(r)} \sum_{k=1}^K \|\Delta_{\mathcal{S}}^k\|_2 \leq \frac{4r}{n} \mathbb{E}(\|\mathbb{X}_{\mathcal{S}_0}^T \xi\|_2), \end{aligned} \quad (34)$$

where $\xi = (\xi_1, \dots, \xi_n)^T$ is a random vector of i.i.d. Rademacher variables that is independent of $\{\varepsilon_1, \dots, \varepsilon_n\}$. By Jensen's and Khintchine inequalities [28], we have

$$\begin{aligned} \mathbb{E}(\|\mathbb{X}_{\mathcal{S}_0}^T \xi\|_2) &\leq [\mathbb{E}(\xi^T \mathbb{X}_{\mathcal{S}_0} \mathbb{X}_{\mathcal{S}_0}^T \xi)]^{1/2} = \left[\sum_{j \in \mathcal{S}_0} \mathbb{E} \left(\sum_{i=1}^n \xi_i x_{ij} \right)^2 \right]^{1/2} \\ &\leq \left(\sum_{j \in \mathcal{S}_0} \sum_{i=1}^n x_{ij}^2 \right)^{1/2} = \left(\sum_{i=1}^n \sum_{j \in \mathcal{S}_0} x_{ij}^2 \right)^{1/2} \leq \sqrt{n(s+1)M_{\mathcal{S}}}. \end{aligned}$$

It follows from (34) that $\mathbb{E}[z(r)] \leq 4r \sqrt{(s+1)M_{\mathcal{S}}/n}$. The first part of the lemma then follows from (33).

Let $F(\delta, \Delta) = Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)$. To prove the second inequality of the lemma, it suffices to note that

$$\inf_{(\delta, \Delta) \in S_{\mathcal{S}}(r)} F(\delta, \Delta) \geq \inf_{(\delta, \Delta) \in S_{\mathcal{S}}(r)} \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] - z(r),$$

and that by (28) and the mean value theorem, we have for some $\bar{u}_{ik,t}$ between 0 and t such that

$$\begin{aligned} &\inf_{(\delta, \Delta) \in S_{\mathcal{S}}(r)} \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \\ &= \inf_{(\delta, \Delta) \in S_{\mathcal{S}}(r)} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^T \Delta} [F(\alpha_k^* + t) - F(\alpha_k^*)] dt \\ &= \inf_{(\delta, \Delta) \in S_{\mathcal{S}}(r)} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^T \Delta} [t f(\alpha_k^* + \bar{u}_{ik,t})] dt. \end{aligned} \quad (35)$$

Now for any $1 \leq i \leq n$, $1 \leq k \leq K$ and $(\delta, \Delta) \in S_{\mathcal{S}}(r)$, we have

$$|\delta_k + \mathbf{x}_i^T \Delta| \leq \sqrt{1 + \|\mathbf{x}_{i,\mathcal{S}}\|_2^2} \cdot \sqrt{\delta_k^2 + \|\Delta_{\mathcal{S}}\|_2^2} \leq r \sqrt{(s+1)M_{\mathcal{S}}} \leq \mathcal{U}_0.$$

It then follows from condition (C0) and (35) that

$$\begin{aligned} & \inf_{(\delta, \Delta) \in \mathcal{S}_{\mathcal{A}}(r)} \mathbb{E}[Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)] \\ & \geq \inf_{(\delta, \Delta) \in \mathcal{S}_{\mathcal{A}}(r)} \frac{f}{2nK} \sum_{i=1}^n \sum_{k=1}^K (\delta_k + \mathbf{x}_i^T \Delta)^2 \geq \frac{f}{2K} \mu r^2. \end{aligned}$$

This completes the proof. \square

Proof of Lemma 7. Note that Q_n is convex, but not differentiable. Denote the subdifferential of $Q_n(\alpha, \beta)$ at $(\hat{\alpha}^\circ, \hat{\beta}^\circ)$ by

$$\begin{aligned} \partial Q_n(\hat{\alpha}^\circ, \hat{\beta}^\circ) = & \left\{ (\zeta, \xi) : \zeta_k = \frac{1-2\tau_k}{2K} - \frac{1}{2nK} \sum_{i=1}^n \text{Sgn}(\hat{r}_{ik}), 1 \leq k \leq K, \right. \\ & \left. \xi_j = \frac{1}{2n} \sum_{i=1}^n x_{ij} \left(1 - \frac{2}{K} \sum_{l=1}^K \tau_l\right) - \frac{1}{2nK} \sum_{i=1}^n \sum_{l=1}^K \text{Sgn}(\hat{r}_{il}) x_{ij}, 1 \leq j \leq p \right\}. \end{aligned}$$

By convexity of Q_n , for any $(\zeta, \xi) \in \partial Q_n(\hat{\alpha}^\circ, \hat{\beta}^\circ)$ and (α, β) , we have

$$Q_n(\alpha, \beta) - Q_n(\hat{\alpha}^\circ, \hat{\beta}^\circ) \geq \zeta^T(\alpha - \hat{\alpha}^\circ) + \xi^T(\beta - \hat{\beta}^\circ).$$

Now by optimality of $(\hat{\alpha}^\circ, \hat{\beta}^\circ)$, we can take $\zeta = \mathbf{0}$ and $\xi_{\mathcal{A}} = \mathbf{0}$. It follows that

$$Q_n(\alpha, \beta) \geq Q_n(\hat{\alpha}^\circ, \hat{\beta}^\circ) + \sum_{j \in \mathcal{A}^c} \xi_j (\beta_j - \hat{\beta}_j^\circ). \quad (36)$$

Under event \mathcal{E}_1 , we have $\max_{j \in \mathcal{A}^c} |\hat{\beta}_j^{(0)}| \leq a_0 \lambda \leq a_2 \lambda$. Moreover, by condition (C3), we have

$$\min_{j \in \mathcal{A}} |\hat{\beta}_j^{(0)}| \geq \min_{j \in \mathcal{A}} |\beta_j^*| - \max_{j \in \mathcal{A}} |\hat{\beta}_j^{(0)} - \beta_j^*| \geq (a+1-a_0)\lambda \geq a\lambda.$$

Thus, under event \mathcal{E}_1 , it follows from properties (P3) and (P4) of $p_\lambda(\cdot)$ that

$$p'_\lambda(|\hat{\beta}_j^{(0)}|) \geq a_1 \lambda, \forall j \in \mathcal{A}^c \quad \text{and} \quad p'_\lambda(|\hat{\beta}_j^{(0)}|) = 0, \forall j \in \mathcal{A}.$$

Similarly, under event \mathcal{E}_3 and by the fact that $\hat{\beta}_{\mathcal{A}^c}^\circ = \mathbf{0}$, it can be shown that

$$p'_\lambda(|\hat{\beta}_j^\circ|) = 0, \forall j \in \mathcal{A} \quad \text{and} \quad p'_\lambda(|\hat{\beta}_j^\circ|) \geq a_1 \lambda, \forall j \in \mathcal{A}^c.$$

To this end, it can be seen from step (2.a) of the LLA algorithm that

$$(\hat{\alpha}^{(1)}, \hat{\beta}^{(1)}) = \arg \min_{\alpha, \beta} Q_n(\alpha, \beta) + \sum_{j \in \mathcal{A}^c} p'_\lambda(|\hat{\beta}_j^{(0)}|) |\beta_j|.$$

Now under $\mathcal{E}_2 = \{\|\xi_{\mathcal{A}^c}\|_\infty < a_1 \lambda\}$, it follows from (36) that for any (α, β) ,

$$\begin{aligned} & \left[Q_n(\alpha, \beta) + \sum_{j \in \mathcal{A}^c} p'_\lambda(|\hat{\beta}_j^{(0)}|) |\beta_j| \right] - \left[Q_n(\hat{\alpha}^\circ, \hat{\beta}^\circ) + \sum_{j \in \mathcal{A}^c} p'_\lambda(|\hat{\beta}_j^{(0)}|) |\hat{\beta}_j^\circ| \right] \\ & \geq \sum_{j \in \mathcal{A}^c} \xi_j (\beta_j - \hat{\beta}_j^\circ) + \sum_{j \in \mathcal{A}^c} p'_\lambda(|\hat{\beta}_j^{(0)}|) |\beta_j| \\ & \geq \sum_{j \in \mathcal{A}^c} [p'_\lambda(|\hat{\beta}_j^{(0)}|) - |\xi_j|] |\beta_j| \geq 0. \end{aligned} \quad (37)$$

The leftmost hand side of the above inequality is strictly positive unless $\beta_{\mathcal{A}^c} = \mathbf{0}$. Note that condition (C0) implies the uniqueness of the oracle estimator (See Section B of this appendix). It can be then seen that $(\hat{\alpha}^{(1)}, \hat{\beta}^{(1)})$ coincides

with the oracle estimator. Now given that $(\widehat{\alpha}^{(1)}, \widehat{\beta}^{(1)})$ is the oracle estimator, we show that $(\widehat{\alpha}^{(2)}, \widehat{\beta}^{(2)})$ yielded by the LLA algorithm will still be the oracle estimator. To see it, note that under event \mathcal{E}_2 ,

$$p'_\lambda(|\widehat{\beta}_j^{(1)}|) = 0, \forall j \in \mathcal{A} \quad \text{and} \quad p'_\lambda(|\widehat{\beta}_j^{(1)}|) \geq a_1 \lambda, \forall j \in \mathcal{A}^c.$$

By the LLA iteration, we have

$$(\widehat{\alpha}^{(2)}, \widehat{\beta}^{(2)}) = \arg \min_{\alpha, \beta} Q_n(\alpha, \beta) + \sum_{j \in \mathcal{A}^c} p'_\lambda(|\widehat{\beta}_j^{(1)}|) |\beta_j|.$$

Thus, we can follow similar arguments from (37) to show that under event \mathcal{E}_3 , $(\widehat{\alpha}^{(2)}, \widehat{\beta}^{(2)})$ is still the oracle estimator.

This proves the lemma. \square

Note that the above proof is slightly different from the general result (Theorems 1 and 2) in [18] since we need to deal with the intercept terms additionally.

Proof of Lemma 8. Consider a minimal ψ -cover of $B_{\mathcal{A}}(r)$ and denote this covering net by $\{(\delta^\ell, \Delta^\ell), \ell = 1, \dots, N_\psi\} \subset B_{\mathcal{A}}(r)$. For $j \in \mathcal{A}^c$, define

$$U_{ij}(\delta, \Delta) = \frac{1}{K} \sum_{k=1}^K [I(r_{ik}^* \leq \delta_k + \mathbf{x}_i^T \Delta) - \tau_k] x_{ij} - \frac{1}{K} \sum_{k=1}^K [I(r_{ik}^* \leq 0) - \tau_k] x_{ij},$$

where $r_{ik}^* = y_i - \alpha_k^* - \mathbf{x}_i^T \beta^* = \varepsilon_i - \alpha_k^*$, $1 \leq i \leq n$, $1 \leq k \leq K$. Then it can be seen that

$$\gamma_j(r) = \sup_{(\delta, \Delta) \in B_{\mathcal{A}}(r)} \left| \frac{1}{n} \sum_{i=1}^n [U_{ij}(\delta, \Delta) - \mathbb{E}U_{ij}(\delta, \Delta)] \right|.$$

For any $(\delta, \Delta) \in B_{\mathcal{A}}(r)$ and $j \in \mathcal{A}^c$, note that

$$|U_{ij}(\delta, \Delta)| \leq \frac{1}{K} \sum_{k=1}^K |I(r_{ik}^* \leq \delta_k + \mathbf{x}_i^T \Delta) - I(r_{ik}^* \leq 0)| \cdot |x_{ij}| \leq M_{\mathcal{A}^c}.$$

Let $p_{ik} = \Pr(-(\delta_k + \mathbf{x}_i^T \Delta)_- < r_{ik}^* \leq (\delta_k + \mathbf{x}_i^T \Delta)_+)$, $1 \leq i \leq n$, $1 \leq k \leq K$. It follows from the mean value theorem and condition (C0) that

$$p_{ik} = F(\alpha_k^* + (\delta_k + \mathbf{x}_i^T \Delta)_+) - F(\alpha_k^* - (\delta_k + \mathbf{x}_i^T \Delta)_-) \leq \bar{f} |\delta_k + \mathbf{x}_i^T \Delta|.$$

By Cauchy–Schwarz inequality and the mean value theorem, we have

$$\begin{aligned} \text{var}[U_{ij}(\delta, \Delta)] &\leq \frac{x_{ij}^2}{K} \sum_{k=1}^K \text{var}[I(-(\delta_k + \mathbf{x}_i^T \Delta)_- < r_{ik}^* \leq (\delta_k + \mathbf{x}_i^T \Delta)_+)] \\ &= \frac{x_{ij}^2}{K} \sum_{k=1}^K p_{ik}(1 - p_{ik}) \leq \frac{\bar{f} x_{ij}^2}{K} \sum_{k=1}^K |\delta_k + \mathbf{x}_i^T \Delta| \leq \frac{\bar{f} M_{\mathcal{A}^c}^2}{K} \sum_{k=1}^K |\delta_k + \mathbf{x}_i^T \Delta_{\mathcal{A}^c}|. \end{aligned}$$

Let $\Delta_{\mathcal{A}^c}^k = (\delta_k, \Delta_{\mathcal{A}^c}^T)^T$, $1 \leq k \leq K$. By Cauchy–Schwarz inequality again, we get

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \text{var}[U_{ij}(\delta, \Delta)] &\leq \frac{1}{nK} \bar{f} M_{\mathcal{A}^c}^2 \sum_{i=1}^n \sum_{k=1}^K |\delta_k + \mathbf{x}_i^T \Delta_{\mathcal{A}^c}| \\ &\leq \frac{\bar{f} M_{\mathcal{A}^c}^2}{K} \sum_{k=1}^K \left[\frac{1}{n} (\Delta_{\mathcal{A}^c}^k)^T \mathbb{X}_{\mathcal{A}^c}^T \mathbb{X}_{\mathcal{A}^c} \Delta_{\mathcal{A}^c}^k \right]^{1/2} \leq \bar{f} M_{\mathcal{A}^c}^2 \bar{\mu}^{1/2} r. \end{aligned}$$

Now applying Bernstein inequality, we have for any $(\delta, \Delta) \in B_{\mathcal{A}}(r)$ and $t > 0$,

$$\Pr\left(\left|\frac{1}{n} \sum_{i=1}^n [U_{ij}(\delta, \Delta) - \mathbb{E}U_{ij}(\delta, \Delta)]\right| > t\right) \leq 2 \exp\left(-\frac{nt^2}{2\bar{f} M_{\mathcal{A}^c}^2 \bar{\mu}^{1/2} r + \frac{4}{3} M_{\mathcal{A}^c} t}\right).$$

Now for $1 \leq \ell \leq N_\psi$, let $B_\ell(\boldsymbol{\psi}) = \{(\boldsymbol{\delta}, \Delta) : \|\boldsymbol{\delta} - \boldsymbol{\delta}^\ell\|_2^2 + \|\Delta - \Delta^\ell\|_2^2 \leq \boldsymbol{\psi}^2, \Delta_{\mathcal{A}^c} = \mathbf{0}\}$ be the ball centered at $(\boldsymbol{\delta}^\ell, \Delta^\ell) \in B_{\mathcal{A}}(r)$ with radius $\boldsymbol{\psi}$. For any $1 \leq i \leq n$, $1 \leq k \leq K$ and $(\boldsymbol{\delta}, \Delta) \in B_\ell(\boldsymbol{\psi})$, note that

$$|(\boldsymbol{\delta}_k + \mathbf{x}_i^T \Delta) - (\boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell)| \leq (1 + \|\mathbf{x}_{i\mathcal{A}}\|_2^2)^{1/2} \boldsymbol{\psi} \leq [(s+1)M_{\mathcal{A}}]^{1/2} \boldsymbol{\psi}.$$

For $1 \leq i \leq n$ and $j \in \mathcal{A}^c$, let

$$V_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell) = \frac{|x_{ij}|}{K} \sum_{k=1}^K [I(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell + ((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}) - I(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell)].$$

Since the indicator function $I(u \leq t)$ is nondecreasing in t , we have

$$\begin{aligned} & \sup_{(\boldsymbol{\delta}, \Delta) \in B_\ell(\boldsymbol{\psi})} \left| \frac{1}{n} \sum_{i=1}^n (U_{ij}(\boldsymbol{\delta}, \Delta) - U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)) - \mathbb{E}[U_{ij}(\boldsymbol{\delta}, \Delta) - U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)] \right| \\ & \leq \frac{1}{nK} \sum_{i=1}^n |x_{ij}| \sum_{k=1}^K [I(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell + ((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}) - I(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell)] \\ & \quad - \Pr(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell - ((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}) + \Pr(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell)] \\ & := I_1 + \frac{1}{n} \sum_{i=1}^n [V_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell) - \mathbb{E}V_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)], \end{aligned} \tag{38}$$

where

$$\begin{aligned} I_1 = \frac{1}{nK} \sum_{i=1}^n |x_{ij}| \sum_{k=1}^K & [\Pr(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell + ((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}) \\ & - \Pr(r_{ik}^* \leq \boldsymbol{\delta}_k^\ell + \mathbf{x}_i^T \Delta^\ell - ((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi})]. \end{aligned}$$

By the mean value theorem, we have

$$I_1 \leq \frac{1}{nK} \sum_{i=1}^n |x_{ij}| \sum_{k=1}^K \cdot 2((s+1)M_{\mathcal{A}})^{1/2} \bar{f} \boldsymbol{\psi} \leq 2\bar{f}M_{\mathcal{A}^c}((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}.$$

Similarly, it can be shown that $|V_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)| \leq M_{\mathcal{A}^c}$ and

$$\frac{1}{n} \sum_{i=1}^n \text{var}(V_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)) \leq \bar{f}M_{\mathcal{A}^c}^2((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}.$$

It then follows from (38) and Bernstein inequality that for $1 \leq \ell \leq N_\psi$,

$$\begin{aligned} & \Pr\left(\sup_{(\boldsymbol{\delta}, \Delta) \in B_\ell(\boldsymbol{\psi})} \left| \frac{1}{n} \sum_{i=1}^n (U_{ij}(\boldsymbol{\delta}, \Delta) - U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)) - \mathbb{E}[U_{ij}(\boldsymbol{\delta}, \Delta) - U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)] \right| > t \right) \\ & \leq 2 \exp\left(-\frac{nt_1^2}{2\bar{f}M_{\mathcal{A}^c}^2((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi} + \frac{4}{3}M_{\mathcal{A}^c}t_1} \right), \end{aligned}$$

where $t_1 = [t - 2\bar{f}M_{\mathcal{A}^c}((s+1)M_{\mathcal{A}})^{1/2} \boldsymbol{\psi}]_+$. The lemma then follows by noting that

$$\begin{aligned} \Pr(\gamma_j(r) > t) & = \Pr\left(\sup_{(\boldsymbol{\delta}, \Delta) \in B_{\mathcal{A}}(r)} \left| \frac{1}{n} \sum_{i=1}^n [U_{ij}(\boldsymbol{\delta}, \Delta) - \mathbb{E}U_{ij}(\boldsymbol{\delta}, \Delta)] \right| > t \right) \\ & \leq \sum_{\ell=1}^{N_\psi} \Pr\left(\sup_{(\boldsymbol{\delta}, \Delta) \in B_\ell(\boldsymbol{\psi})} \left| \frac{1}{n} \sum_{i=1}^n (U_{ij}(\boldsymbol{\delta}, \Delta) - U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)) \right. \right. \\ & \quad \left. \left. - \mathbb{E}[U_{ij}(\boldsymbol{\delta}, \Delta) - U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)] \right| > \frac{t}{2} \right) \\ & \quad + \sum_{\ell=1}^{N_\psi} \Pr\left(\left| \frac{1}{n} \sum_{i=1}^n [U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell) - \mathbb{E}U_{ij}(\boldsymbol{\delta}^\ell, \Delta^\ell)] \right| > \frac{t}{2} \right). \end{aligned}$$

This completes the proof. \square

Proof of Lemma 9. Let $U_i(\delta, \Delta) = \frac{1}{K} \sum_{k=1}^K \{\rho_{\tau_k}(r_{ik}^* - \delta_k - \mathbf{x}_i^T \Delta) - \rho_{\tau_k}(r_{ik}^*)\}$, where $r_{ik}^* = y_i - \alpha_k^* - \mathbf{x}_i^T \beta^* = \varepsilon_i - \alpha_k^*$, $1 \leq i \leq n, 1 \leq k \leq K$. By Massart's concentration inequality, we have

$$\Pr(z(A, r) \geq \mathbb{E}[z(A, r)] + t) \leq \exp\left(-\frac{n^2 t^2}{8b_n^2(A, r)}\right),$$

where $b_n^2(A, r) = \sup_{(\delta, \Delta) \in B_A(r)} \sum_{i=1}^n \text{var}(U_i(\delta, \Delta))$. It follows from Lipschitz continuity of the check loss that

$$\begin{aligned} b_n^2(A, r) &\leq \frac{4}{K} \sup_{(\delta, \Delta) \in B_A(r)} \sum_{k=1}^K \sum_{i=1}^n (\delta_k + \mathbf{x}_i^T \Delta)^2 \\ &= \frac{4}{K} \sup_{(\delta, \Delta) \in B_A(r)} \sum_{k=1}^K (\delta_k, \Delta_A^T) (\mathbf{1}_n, \mathbf{X}_A)^T (\mathbf{1}_n, \mathbf{X}_A) (\delta_k, \Delta_A^T)^T \\ &\leq \frac{4n}{K} \sup_{(\delta, \Delta) \in B_A(r)} \sum_{k=1}^K \bar{\xi} [\delta_k^2 + \|\Delta_A\|^2] \leq 4n \bar{\xi} r^2. \end{aligned}$$

Moreover, by the symmetrization procedure and the contraction principle, we obtain

$$\begin{aligned} \mathbb{E}[z(A, r)] &\leq \frac{4}{nK} \sum_{k=1}^K \mathbb{E} \left[\sup_{(\delta, \Delta) \in B_A(r)} \left| \sum_{i=1}^n \xi_i (\delta_k + \mathbf{x}_i^T \Delta) \right| \right] \\ &\leq \frac{4}{nK} \mathbb{E}(\|(\mathbf{1}_n, \mathbf{X}_A)^T \xi\|_2) \cdot \sup_{(\delta, \Delta) \in B_A(r)} \sum_{k=1}^K \|(\delta_k, \Delta_A^T)^T\|_2 \\ &\leq \frac{4r}{n} \mathbb{E}(\|(\mathbf{1}_n, \mathbf{X}_A)^T \xi\|_2), \end{aligned}$$

where $\xi = (\xi_1, \dots, \xi_n)^T$ is a random vector of i.i.d. Rademacher variables, independent of $\varepsilon_1, \dots, \varepsilon_n$. By Jensen's and Khintchine inequalities, we have

$$\begin{aligned} \mathbb{E}(\|(\mathbf{1}_n, \mathbf{X}_A)^T \xi\|_2) &\leq \{\mathbb{E}[\xi^T (\mathbf{1}_n, \mathbf{X}_A) (\mathbf{1}_n, \mathbf{X}_A)^T \xi]\}^{1/2} \\ &= \left[\sum_{j \in \{0\} \cup A} \mathbb{E} \left(\sum_{i=1}^n \xi_i x_{ij} \right)^2 \right]^{1/2} \leq \left(\sum_{j \in \{0\} \cup A} \sum_{i=1}^n x_{ij}^2 \right)^{1/2} \\ &= \left(\sum_{i=1}^n \sum_{j \in \{0\} \cup A} x_{ij}^2 \right)^{1/2} \leq M \sqrt{n(|A|+1)}. \end{aligned}$$

It follows that $\mathbb{E}[z(A, r)] \leq 4Mr \sqrt{(|A|+1)/n}$. The lemma then follows. \square

Proof of Theorem 4. As in the proof of Theorem 3, split all models under consideration, $\{\hat{A}_\lambda : \lambda \in \Xi_n\}$, into three groups: $\{\hat{A}_\lambda : \lambda \in \Xi_n^-\}$, $\{\hat{A}_\lambda : \lambda \in \Xi_n^0\}$, and $\{\hat{A}_\lambda : \lambda \in \Xi_n^+\}$, where $\Xi_n^- = \{\lambda \in \Xi_n : \mathcal{A} \not\subset \hat{A}_\lambda\}$, $\Xi_n^0 = \{\lambda \in \Xi_n : \hat{A}_\lambda = \mathcal{A}\}$ and $\Xi_n^+ = \{\lambda \in \Xi_n : \mathcal{A} \subset \hat{A}_\lambda, \hat{A}_\lambda \neq \mathcal{A}\}$.

Let $\hat{Q}_n^\lambda = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \hat{\alpha}_k^\lambda - \mathbf{x}_i^T \hat{\beta}^\lambda)$, where $(\hat{\alpha}^\lambda, \hat{\beta}^\lambda)$ is the two-step LLA estimator to the folded concave penalized CQR (6) with lasso initialization. Also, let $Q_n^* = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k^* - \mathbf{x}_i^T \beta^*)$. For any $A \subset \{1, 2, \dots, p\}$, let $(\hat{\alpha}^A, \hat{\beta}^A)$ be the estimator obtained by fitting the canonical CQR to model A , i.e.,

$$(\hat{\alpha}^A, \hat{\beta}^A) = \arg \min_{\alpha, \beta: \beta_{A^c} = \mathbf{0}} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta). \quad (39)$$

Define $\hat{Q}_n^A = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \hat{\alpha}_k^A - \mathbf{x}_i^T \hat{\beta}^A)$. For any $\lambda \in \Xi_n$, recall that $\hat{A}_\lambda = \{1 \leq j \leq p : \hat{\beta}_j^\lambda \neq 0\}$ corresponds to the active set of the two-step LLA estimator $(\hat{\alpha}^\lambda, \hat{\beta}^\lambda)$. By optimality of $(\hat{\alpha}^{\hat{A}_\lambda}, \hat{\beta}^{\hat{A}_\lambda})$ in (39), we have $\hat{Q}_n^{\hat{A}_\lambda} \leq \hat{Q}_n^\lambda$.

Let $\mathcal{G}_n^+ = \{A : A \supset \mathcal{A}, A \neq \mathcal{A}, |A| \leq J_n\}$. It can be seen that $\{\hat{A}_\lambda : \lambda \in \Xi_n^+\} \subset \mathcal{G}_n^+$. For $r > 0$ and $A \in \mathcal{G}_n^+$, let $B_A(r) = \{(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p : \|\delta\|^2 + \|\Delta_A\|^2 \leq r^2, \Delta_{A^c} = \mathbf{0}\}$ and $S_A(r) = \{(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p : \|\delta\|^2 + \|\Delta_A\|^2 = r^2, \Delta_{A^c} = \mathbf{0}\}$.

Case I: overfitted models. By Theorem 2, under the assumptions of this theorem, we have

$$P(\hat{A}_{\lambda_n} \neq \mathcal{A}) \rightarrow \mathcal{O}(1) \quad \text{as } n \rightarrow \infty.$$

Therefore, for any $\lambda \in \Xi_n^+$, we have

$$\begin{aligned} & \Pr\left(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^{\text{HL}}(\lambda) - \text{BIC}^{\text{HL}}(\lambda_n)] > 0\right) \\ &= \Pr\left(\inf_{\lambda \in \Xi_n^+} \left[\log\left(\frac{\hat{Q}_n^\lambda}{\hat{Q}_n^{\mathcal{A}}}\right) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n}\right] > 0\right) + \mathcal{O}(1) \\ &\geq \Pr\left(\inf_{\lambda \in \Xi_n^+} \left[\log\left(\frac{\hat{Q}_n^{\hat{A}_\lambda}}{\hat{Q}_n^{\mathcal{A}}}\right) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n}\right] > 0\right) + \mathcal{O}(1), \end{aligned}$$

where the last inequality follows from the fact that $\hat{Q}_n^\lambda \geq \hat{Q}_n^{\hat{A}_\lambda}$. Moreover, note that $\hat{Q}_n^{\hat{A}_\lambda} \leq \hat{Q}_n^{\mathcal{A}} \leq Q_n^*$ due to inclusion $\mathcal{A} \subset \hat{A}_\lambda$. We then apply the inequality $\log(1+x) \leq x, \forall x \geq 0$ to get

$$\log\left(\frac{\hat{Q}_n^{\hat{A}_\lambda}}{\hat{Q}_n^{\mathcal{A}}}\right) = -\log\left(\frac{\hat{Q}_n^{\mathcal{A}}}{\hat{Q}_n^{\hat{A}_\lambda}}\right) = -\log\left(1 + \frac{\hat{Q}_n^{\mathcal{A}} - \hat{Q}_n^{\hat{A}_\lambda}}{\hat{Q}_n^{\hat{A}_\lambda}}\right) \geq -\frac{\hat{Q}_n^{\mathcal{A}} - \hat{Q}_n^{\hat{A}_\lambda}}{\hat{Q}_n^{\hat{A}_\lambda}}.$$

Let $F(\delta, \Delta) = Q_n(\alpha^* + \delta, \beta^* + \Delta) - Q_n(\alpha^*, \beta^*)$ for $(\delta, \Delta) \in \mathbb{R}^K \times \mathbb{R}^p$. For each model $A \in \mathcal{G}_n^+$, let $r_A^* = 16K(M + \bar{\zeta}^{1/2})(f\underline{\zeta})^{-1} \sqrt{(|A|+1)\log(p)/n}$. If we can show that $\inf_{(\delta, \Delta) \in S_A(r_A^*)} F(\delta, \Delta) > 0$, then by convexity of $\rho_\tau(\cdot)$, we must have $\|\hat{\alpha}^A - \alpha^*\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \leq (r_A^*)^2$. Indeed, by Knight's identity (see (28) of the appendix) and the mean value theorem, we have

$$\begin{aligned} & \inf_{(\delta, \Delta) \in S_A(r_A^*)} \mathbb{E}[F(\delta, \Delta)] \\ &= \inf_{(\delta, \Delta) \in S_A(r_A^*)} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^T \Delta} [F(\alpha_k^* + t) - F(\alpha_k^*)] dt \\ &= \inf_{(\delta, \Delta) \in S_A(r_A^*)} \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\delta_k + \mathbf{x}_i^T \Delta} [tf(\alpha_k^* + \bar{u}_{ik,t})] dt. \end{aligned} \quad (40)$$

Note that J_n satisfies $J_n^2 \log(p)/n = \mathcal{O}(1)$. Therefore, for any $1 \leq i \leq n, 1 \leq k \leq K$ and $(\delta, \Delta) \in S_A(r_A^*)$, we have

$$\begin{aligned} |\delta_k + \mathbf{x}_i^T \Delta| &\leq \sqrt{1 + \|\mathbf{x}_{iA}\|^2} \cdot \sqrt{\delta_k^2 + \|\Delta_A\|_2^2} \\ &\leq Mr_A^* \sqrt{|A|+1} = 16KM(M + \bar{\zeta}^{1/2})(f\underline{\zeta})^{-1} (|A|+1) \sqrt{\log(p)/n} \\ &\leq 16KM(M + \bar{\zeta}^{1/2})(f\underline{\zeta})^{-1} (J_n + 1) \sqrt{\log(p)/n} = \mathcal{O}(\mathcal{U}_0). \end{aligned}$$

It follows from condition (C0) and (40) that

$$\inf_{(\delta, \Delta) \in S_A(r_A^*)} \mathbb{E}[F(\delta, \Delta)] \geq \inf_{(\delta, \Delta) \in S_A(r_A^*)} \frac{f}{2nK} \sum_{k=1}^K \sum_{i=1}^n (\delta_k + \mathbf{x}_i^T \Delta)^2 \geq \frac{f}{2K} \bar{\zeta} (r_A^*)^2.$$

Therefore, by Lemma 9, with probability at least $1 - \exp\{-nt^2/[32\bar{\zeta}(r_A^*)^2]\}$, we have

$$\begin{aligned} \inf_{(\delta, \Delta) \in S_A(r_A^*)} F(\delta, \Delta) &\geq \inf_{(\delta, \Delta) \in S_A(r_A^*)} \mathbb{E}[F(\delta, \Delta)] - z(A, r_A^*) \\ &\geq \frac{f}{2K} \bar{\zeta} (r_A^*)^2 - 4Mr_A^* \sqrt{\frac{|A|+1}{n}} - t. \end{aligned}$$

Now take $t = 8r_A^* \sqrt{\bar{\zeta}(|A|+1)\log(p)/n}$. It follows that for each $A \in \mathcal{G}_n^+$, with probability at least $b_n^A = 1 - \exp\{-2(|A|+1)\log(p)\}$, we have

$$\inf_{(\delta, \Delta) \in \mathcal{S}_A(r_A^*)} F(\delta, \Delta) \geq \frac{f\bar{\zeta}(r_A^*)^2}{2K} - r_A^* \sqrt{\frac{|A|+1}{n}} \left(8\sqrt{\bar{\zeta}\log(p)} + 4M\right) > 0,$$

which immediately implies that

$$\|\hat{\alpha}^A - \alpha\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \leq (r_A^*)^2.$$

Now by the Bonferroni inequality, we have $\|\hat{\alpha}^A - \alpha\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \leq (r_A^*)^2$ for all $A \in \mathcal{G}_n^+$ simultaneously with probability at least

$$\begin{aligned} b_n &= 1 - \sum_{|A|=s+1}^{J_n} \binom{p-s}{|A|-s} (1-b_n^A) \\ &= 1 - \sum_{k=1}^{J_n-s} \binom{p-s}{k} \exp\{-2(k+s+1)\log(p)\} \\ &\geq 1 - p^{-2(s+1)} \sum_{k=1}^{p-s} \binom{p-s}{k} \left(\frac{1}{p^2}\right)^k \\ &= 1 - p^{-2(s+1)} \left[\left(1 + \frac{1}{p^2}\right)^{p-s} - 1\right] \rightarrow 1 \quad \text{as } p \rightarrow \infty. \end{aligned} \tag{41}$$

Now we derive the upper bound for $\sup_{A \in \mathcal{G}_n^+} |\hat{Q}_n^A - Q_n^*|$. Let $\hat{\delta}^A = \hat{\alpha}^A - \alpha^*$ and $\hat{\Delta}^A = \hat{\beta}^A - \beta^*$. Observe that

$$\begin{aligned} |\hat{Q}_n^A - Q_n^*| &= \left| \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \{\rho_{\tau_k}(y_i - \hat{\alpha}_k^A - \mathbf{x}_i^T \hat{\beta}^A) - \rho_{\tau_k}(y_i - \alpha_k^* - \mathbf{x}_i^T \beta^*)\} \right| \\ &\leq |\mathbb{E}\{F(\hat{\delta}^A, \hat{\beta}^A)\}| + z(A, r_A^*). \end{aligned}$$

Similarly, we have

$$\begin{aligned} |\mathbb{E}\{F(\hat{\delta}^A, \hat{\beta}^A)\}| &= \left| \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{\hat{\delta}_k^A + \mathbf{x}_i^T \hat{\Delta}^A} [F(\alpha_k^* + t) - F(\alpha_k^*)] dt \right| \\ &\leq \frac{1}{nK} \sum_{i=1}^n \sum_{k=1}^K \int_0^{|\hat{\delta}_k^A + \mathbf{x}_i^T \hat{\Delta}^A|} [tf(\alpha_k^* + \bar{u}_{ik,t})] dt \\ &\leq \frac{\bar{f}}{2nK} \sum_{i=1}^n \sum_{k=1}^K (\hat{\delta}_k^A + \mathbf{x}_i^T \hat{\Delta}^A)^2 \\ &\leq \frac{1}{2} \bar{f} \bar{\zeta} (\|\hat{\delta}^A\|_2^2 + \|\hat{\Delta}^A\|_2^2) \leq \frac{1}{2} \bar{f} \bar{\zeta} (r_A^*)^2. \end{aligned}$$

It follows that with probability at least b_n ,

$$\begin{aligned} |\hat{Q}_n^A - Q_n^*| &\leq \frac{1}{2} \bar{f} \bar{\zeta} (r_A^*)^2 + 4Mr_A^* \sqrt{\frac{|A|+1}{n}} + 8r_A^* \sqrt{\frac{\bar{\zeta}(|A|+1)\log(p)}{n}} \\ &\leq \frac{128K^2(M + \bar{\zeta}^{1/2})^2 (|A|+1)\log(p)}{\bar{f}\bar{\zeta}} \frac{1}{n} \end{aligned}$$

holds for all $A \in \mathcal{G}_n^+$. Now going back to BIC, we have

$$\begin{aligned} &\Pr\left(\inf_{\lambda \in \Xi_n^+} \left[\log(\hat{Q}_n^{\hat{\lambda}} / \hat{Q}_n^{\mathcal{A}}) + (|\hat{\lambda}| - s) \frac{C_n \log(p)}{n}\right] > 0\right) \\ &\geq \Pr\left(\frac{C_n \log(p)}{n} - \sup_{A \in \mathcal{G}_n^+} \frac{\hat{Q}_n^{\mathcal{A}} - \hat{Q}_n^A}{(|A| - s) \hat{Q}_n^A} > 0\right). \end{aligned}$$

Since $\mathbb{E}(|\varepsilon|) < \infty$, it follows that $\mathbb{E}(Q_n^*) \leq \mathbb{E}(|\varepsilon|) + \sum_{k=1}^K |\alpha_k^*|/K < \infty$. Thus, we have

$$\widehat{Q}_n^A = Q_n^* - (Q_n^* - \widehat{Q}_n^A) = \mathcal{O}_P(1)$$

by noting that $J_n \log(p) = \mathcal{o}(n)$. Therefore, with probability at least b_n , we have

$$\sup_{A \in \mathcal{G}_n^+} \frac{\widehat{Q}_n^{\mathcal{A}} - \widehat{Q}_n^A}{(|A| - s)\widehat{Q}_n^A} \leq \sup_{A \in \mathcal{G}_n^+} \frac{\widehat{Q}_n^* - \widehat{Q}_n^A}{(|A| - s)\widehat{Q}_n^A} = \mathcal{O}_P\left(\frac{s \log(p)}{n}\right).$$

Since $s = \mathcal{O}(1)$ and C_n diverges with n , we have $s \log(p)/n = \mathcal{O}(C_n \log(p)/n)$. It follows that

$$\Pr\left(\frac{C_n \log(p)}{n} - \sup_{A \in \mathcal{G}_n^+} \frac{\widehat{Q}_n^{\mathcal{A}} - \widehat{Q}_n^A}{(|A| - s)\widehat{Q}_n^A} > 0\right) \rightarrow 1 \quad \text{as } n \rightarrow \infty,$$

which implies that $\Pr(\inf_{\lambda \in \Xi_n^+} [\text{BIC}^{\text{HL}}(\lambda) - \text{BIC}^{\text{HL}}(\lambda_n)] > 0) \rightarrow 1$ as $n \rightarrow \infty$.

Case II: underfitted models. For any $\lambda \in \Xi_n^-$, similar to Case I, we have

$$\begin{aligned} & \Pr\left(\inf_{\lambda \in \Xi_n^-} [\text{BIC}^{\text{HL}}(\lambda) - \text{BIC}^{\text{HL}}(\lambda_n)] > 0\right) \\ & \geq \Pr\left(\inf_{\lambda \in \Xi_n^-} \left[\log(\widehat{Q}_n^{\hat{A}_\lambda} / \widehat{Q}_n^{\mathcal{A}}) + (|\hat{A}_\lambda| - s) \frac{C_n \log(p)}{n}\right] > 0\right) + \mathcal{o}(1). \end{aligned}$$

Define $\text{BIC}^{\text{HL}}(A) = \log(n\widehat{Q}_n^A) + |A|C_n \log(p)/n$ and let $\mathcal{G}_n^- = \{A : |A| \leq J_n, \mathcal{A} \not\subset A\}$. We can see that $\{\hat{A}_\lambda : \lambda \in \Xi_n^-\} \subset \mathcal{G}_n^-$. It suffices to show $\inf_{A \in \mathcal{G}_n^-} \text{BIC}^{\text{HL}}(A) > \text{BIC}^{\text{HL}}(\mathcal{A})$ with probability tending to one as $n \rightarrow \infty$. For any $A \in \mathcal{G}_n^-$, let $\bar{A} = A \cup \mathcal{A}$. Let $\theta = \min_{j \in \mathcal{A}} |\beta_j^*|$. Since $A \not\subset \mathcal{A}$, we must have $\|\hat{\alpha}^A - \alpha^*\|_2^2 + \|\hat{\beta}^A - \beta^*\|_2^2 \geq \theta^2$. However, since $\bar{A} \supset \mathcal{A}$ and $|\bar{A}| \leq 2J_n$, using Lemma 9, we can similarly show as in Case I that $\|\hat{\alpha}^{\bar{A}} - \alpha^*\|_2^2 + \|\hat{\beta}^{\bar{A}} - \beta^*\|_2^2 \leq \theta^2$ with probability at least $b_n^{\bar{A}} = 1 - \exp\{-2(|\bar{A}| + 1) \log(p)\}$ as long as

$$\theta > 8K(f\zeta)^{-1} (2\sqrt{\zeta \log(p)} + M) \sqrt{\frac{2J_n + 1}{n}},$$

which is implied by the assumption $\sqrt{J_n \log(p)/n} = \mathcal{O}(\theta)$. It then follows that $\|\hat{\alpha}^{\bar{A}} - \alpha^*\|_2^2 + \|\hat{\beta}^{\bar{A}} - \beta^*\|_2^2 \leq \theta^2$ holds for all $A \in \mathcal{G}_n^-$ with probability at least $\tilde{b}_n \rightarrow 1$ as $n \rightarrow \infty$, where

$$\tilde{b}_n = 1 - \sum_{|\bar{A}|=s+1}^{2J_n} \binom{p-s}{|\bar{A}|-s} (1 - b_n^{\bar{A}}) \geq 1 - p^{-2(s+1)} \left[\left(1 + \frac{1}{p^2}\right)^{p-s} - 1 \right].$$

Therefore, there exists $a \in [0, 1]$, $\bar{\alpha}^{\bar{A}} = a\hat{\alpha}^{\bar{A}} + (1-a)\hat{\alpha}^{\mathcal{A}}$ and $\bar{\beta}^{\bar{A}} = a\hat{\beta}^{\bar{A}} + (1-a)\hat{\beta}^{\mathcal{A}}$ such that $\|\bar{\alpha}^{\bar{A}} - \alpha^*\|_2^2 + \|\bar{\beta}^{\bar{A}} - \beta^*\|_2^2 = \theta^2$. By convexity of ρ_τ and the fact that $\widehat{Q}_n^{\bar{A}} \geq \widehat{Q}_n^{\mathcal{A}}$, we have $\bar{Q}_n^{\bar{A}} = (nK)^{-1} \sum_{i=1}^n \sum_{k=1}^K \rho_{\tau_k}(y_i - \bar{\alpha}_k^{\bar{A}} - \mathbf{x}_i^T \bar{\beta}^{\bar{A}}) \leq \widehat{Q}_n^{\bar{A}}$. Note that $\widehat{Q}_n^{\bar{A}} \leq \widehat{Q}_n^{\mathcal{A}} \leq Q_n^*$. It follows that $\bar{Q}_n^{\bar{A}} - \widehat{Q}_n^{\bar{A}} \geq \bar{Q}_n^{\bar{A}} - Q_n^*$. For ease of notation, let $\bar{\delta}^{\bar{A}} = \bar{\alpha}^{\bar{A}} - \alpha^*$ and $\bar{\Delta}^{\bar{A}} = \bar{\beta}^{\bar{A}} - \beta^*$. It can be seen that

$$\bar{Q}_n^{\bar{A}} - Q_n^* \geq \mathbb{E}[F(\bar{\delta}^{\bar{A}}, \bar{\Delta}^{\bar{A}})] - z(\bar{A}, \theta).$$

Following similar arguments from Case I and noting that the support of $\bar{\beta}^{\bar{A}}$ is a subset of \bar{A} , we can show that with probability at least \tilde{b}_n , for all $A \in \mathcal{G}_n^-$, we have

$$\bar{Q}_n^{\bar{A}} - Q_n^* \geq \frac{f}{2K} \zeta \theta^2 - 4\theta \sqrt{\frac{|\bar{A}| + 1}{n}} (2\sqrt{\zeta \log(p)} + M).$$

Now for all $A \in \mathcal{G}_n^-$, applying the inequality $\log(1+x) \geq \min\{\log(2), x/2\}, \forall x \geq 0$, we have

$$\begin{aligned} \text{BIC}^{\text{HL}}(A) - \text{BIC}^{\text{HL}}(\bar{A}) &= \log\left(1 + \frac{\widehat{Q}_n^A - \widehat{Q}_n^{\bar{A}}}{\widehat{Q}_n^{\bar{A}}}\right) + (|A| - |\bar{A}|) \frac{C_n \log(p)}{n} \\ &\geq \min\left\{\log(2), \frac{\widehat{Q}_n^{\bar{A}} - Q_n^*}{\widehat{Q}_n^{\bar{A}}}\right\} - \frac{C_n s \log(p)}{n}. \end{aligned}$$

Since $\sqrt{C_n s \log(p)/n} = \sigma(\theta)$ and $\widehat{Q}_n^{\bar{A}} = \mathcal{O}_P(1)$, it can be seen that with probability tending to one, we have $\inf_{A \in \mathcal{G}_n^-} \text{BIC}^{\text{HL}}(A) - \text{BIC}^{\text{HL}}(\bar{A}) > 0$. Following similar arguments as in Case I, we can show $\text{BIC}^{\text{HL}}(\bar{A}) \geq \inf_{S \supset \mathcal{A}, |S| \leq 2J_n} \text{BIC}^{\text{HL}}(S) \geq \text{BIC}^{\text{HL}}(\mathcal{A})$ with probability tending to one. Case II then follows by noting that

$$\begin{aligned} &\inf_{A \in \mathcal{G}_n^-} [\text{BIC}^{\text{HL}}(A) - \text{BIC}^{\text{HL}}(\mathcal{A})] \\ &= \inf_{A \in \mathcal{G}_n^-} [\text{BIC}^{\text{HL}}(A) - \text{BIC}^{\text{HL}}(\bar{A}) + \text{BIC}^{\text{HL}}(\bar{A}) - \text{BIC}^{\text{HL}}(\mathcal{A})] \\ &\geq \inf_{A \in \mathcal{G}_n^-} [\text{BIC}^{\text{HL}}(A) - \text{BIC}^{\text{HL}}(\bar{A})]. \end{aligned}$$

□

APPENDIX B

NUMERICAL PROPERTIES OF THE CQR ORACLE SOLUTION

Recall that the CQR oracle estimator is obtained through regression on the true set of variables

$$(\widehat{\alpha}^o, \widehat{\beta}^o) := \arg \min_{(\alpha, \beta): \beta_{\mathcal{A}^c} = \mathbf{0}} \sum_{k=1}^K w_k \sum_{i=1}^n \rho_{\tau_k}(y_i - \alpha_k - \mathbf{x}_i^T \beta).$$

For ease of exposition, we will restrict the scope of variables under consideration to those in \mathcal{A} . Specifically, let $\mathbf{a} = \alpha$, $\mathbf{b} = \beta_{\mathcal{A}} \in \mathbb{R}^s$ and $\mathbf{z}_i = \mathbf{x}_{i\mathcal{A}}$, $i = 1, \dots, n$. The oracle solution can be equivalently obtained through the following minimization problem

$$(\widehat{\mathbf{a}}, \widehat{\mathbf{b}}) := \arg \min_{\mathbf{a}, \mathbf{b}} \sum_{k=1}^K w_k \sum_{i=1}^n \rho_{\tau_k}(y_i - a_k - \mathbf{z}_i^T \mathbf{b}).$$

Now let $\mathbf{u}_k = (\mathbf{y} - a_k \mathbf{1}_n - \mathbf{Z}\mathbf{b})_+$ and $\mathbf{v}_k = (\mathbf{y} - a_k \mathbf{1}_n - \mathbf{Z}\mathbf{b})_-$, $k = 1, \dots, K$, where $\mathbf{Z} = (\mathbf{z}_1, \dots, \mathbf{z}_n)^T$ and the positive and negative parts are taken componentwisely. Also, let $\mathbf{u} = (\mathbf{u}_1^T, \dots, \mathbf{u}_K^T)^T$ and $\mathbf{v} = (\mathbf{v}_1^T, \dots, \mathbf{v}_K^T)^T$. Then the above regression problem can be cast into the following linear program of standard form

$$\begin{aligned} &\text{minimize} && c^T x \\ &\text{subject to} && Ax = b \\ &&& x \succeq 0, \end{aligned}$$

where $b = \mathbf{1}_K \otimes \mathbf{y}$, (\otimes : Kronecker product) and

$$\begin{aligned} x &= (\mathbf{a}_+^T, \mathbf{a}_-^T, \mathbf{b}_+^T, \mathbf{b}_-^T, \mathbf{u}^T, \mathbf{v}^T)^T, \\ c &= (\mathbf{0}_K^T, \mathbf{0}_K^T, \mathbf{0}_p^T, \mathbf{0}_p^T, w_1 \tau_1 \mathbf{1}_n^T, \dots, w_K \tau_K \mathbf{1}_n^T, w_1 (1 - \tau_1) \mathbf{1}_n^T, \dots, w_K (1 - \tau_K) \mathbf{1}_n^T)^T, \\ A &= \begin{pmatrix} \mathbf{1}_n & \cdots & \mathbf{0} & -\mathbf{1}_n & \cdots & \mathbf{0} & \mathbf{Z} & -\mathbf{Z} & \mathbf{I}_n & -\mathbf{I}_n \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \mathbf{0} & \cdots & \mathbf{1}_n & \mathbf{0} & \cdots & -\mathbf{1}_n & \mathbf{Z} & -\mathbf{Z} & \mathbf{I}_n & -\mathbf{I}_n \end{pmatrix}_{(nK) \times (2K+2s+2n)}. \end{aligned}$$

Without loss of generality, assume that $\mathbf{1}_n \notin \text{Span}(\mathbf{Z})$, where $\text{Span}(\mathbf{Z})$ denotes the column span of \mathbf{Z} . Write

$$D = \begin{pmatrix} \mathbf{1}_n & \cdots & \mathbf{0} & \mathbf{Z} \\ \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \cdots & \mathbf{1}_n & \mathbf{Z} \end{pmatrix}_{(nK) \times (K+s)}.$$

The rows of D will be denoted by d_i^\top , $i = 1, \dots, nK$. Let \mathcal{H} be the collection of $(K+s)$ -element subsets of $\{1, \dots, nK\}$. For $h \in \mathcal{H}$, let $D(h)$ denote the submatrix of D with rows $\{d_i^\top, i \in h\}$ and $b(h)$ be the $(K+s)$ -vector with coordinates $\{b_i, i \in h\}$. We also let $\bar{h} = \{1, \dots, nK\} \setminus h$ for $h \in \mathcal{H}$. Let $H = \{h \in \mathcal{H} : |D(h)| \neq 0\}$. By similar arguments as in Section 6.2 of [3], one can verify that the vertices of the polyhedron $\{x : Ax = b, x \geq 0\}$ are given by

$$\begin{aligned} (\mathbf{a}^\top, \mathbf{b}^\top)^\top &= [D(h)]^{-1} b(h) \\ \mathbf{u}(h) &= \mathbf{v}(h) = \mathbf{0} \\ \mathbf{u}(\bar{h}) &= \left[b(\bar{h}) - D(\bar{h}) \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \right]_+ \\ \mathbf{v}(\bar{h}) &= \left[b(\bar{h}) - D(\bar{h}) \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} \right]_- \end{aligned}$$

for all $h \in H$. According to the simplex algorithm (see, e.g., [29], Chapter 3), the optimal solution to this linear program is among the above set of vertices. Recall that y has a density with respect to the Lebesgue measure. It can be seen that with probability one, there are at most $K(K+s)$ zero residuals for which $y_i - \hat{a}_k - \mathbf{z}_i^\top \hat{\mathbf{b}} = 0$, $1 \leq i \leq n$, $1 \leq k \leq K$, given each optimal solution $(\hat{\mathbf{a}}, \hat{\mathbf{b}})$. Otherwise, suppose that there exist $h \in H$ and $i \in \bar{h}$ such that $\mathbf{u}(i) = \mathbf{v}(i) = 0$. Then since $D(h)$ is non-singular, it follows that

$$b_i = d_i^\top \begin{pmatrix} \mathbf{a} \\ \mathbf{b} \end{pmatrix} = d_i^\top [D(h)]^{-1} b(h),$$

which implies that b_i is a linear combination of $b(h)$. By the assumption that y has a density and the structure of b , this occurs with probability zero unless $b_i = b_j$ for some $j \in h$. However, there are at most $(K-1)$ such i 's for each $j \in h$. This means with probability one, at each vertex, there are at most $K(K+s)$ indices i for which $\mathbf{u}(i) = \mathbf{v}(i) = 0$.

APPENDIX C
CONVERGENCE CRITERION FOR THE ADMM ALGORITHM

We adopt the following convergence criterion recommended by [30] for the ADMM algorithm (Algorithm 1) we propose to solve the weighted lasso penalized CQR problem:

$$\begin{aligned} & \left\| \begin{pmatrix} \mathbb{X}_1 \\ -\mathbb{X}_2 \end{pmatrix} \boldsymbol{\varphi}^r + \begin{pmatrix} \text{vec}(\mathbf{Z}^r) \\ \boldsymbol{\gamma}^r \end{pmatrix} - \begin{pmatrix} \mathbf{Y} \\ \mathbf{0} \end{pmatrix} \right\|_2 \leq \varepsilon_1 \sqrt{nK+p} \\ & + \varepsilon_2 \cdot \max \left\{ \left\| \begin{pmatrix} \mathbb{X}_1 \\ -\mathbb{X}_2 \end{pmatrix} \boldsymbol{\varphi}^r \right\|_2, \left\| \begin{pmatrix} \text{vec}(\mathbf{Z}^r) \\ \boldsymbol{\gamma}^r \end{pmatrix} \right\|_2, \|\mathbf{Y}\|_2 \right\}, \\ & \sigma \|\mathbb{X}_1^T \{\text{vec}(\mathbf{Z}^r) - \text{vec}(\mathbf{Z}^{r-1})\} - \mathbb{X}_2^T (\boldsymbol{\gamma}^r - \boldsymbol{\gamma}^{r-1})\|_2 \leq \varepsilon_1 \sqrt{p+K} \\ & + \varepsilon_2 \cdot \|\mathbb{X}_1^T \text{vec}(\mathbf{U}^r) - \mathbb{X}_2^T \mathbf{v}^r\|_2, \end{aligned}$$

where ε_1 and ε_2 are the tolerances taking small positive values.

REFERENCES

- [1] R. Koenker and G. Bassett, "Regression quantiles," *Econometrica: Journal of the Econometric Society*, vol. 46, no. 1, pp. 33–50, Jan. 1978.
- [2] K. Knight, "Limiting distributions for L_1 regression estimators under general conditions," *The Annals of Statistics*, vol. 26, no. 2, pp. 755–770, 1998.
- [3] R. Koenker, *Quantile Regression*. Cambridge, United Kingdom: Cambridge University Press, 2005.
- [4] H. Zou and M. Yuan, "Composite quantile regression and the oracle model selection theory," *The Annals of Statistics*, vol. 36, no. 3, pp. 1108–1126, 2008.
- [5] R. Koenker, "A note on L-estimates for linear models," *Statistics & probability letters*, vol. 2, no. 6, pp. 323–325, 1984.
- [6] Z. Zhao and Z. Xiao, "Efficient regressions via optimally combining quantile information," *Econometric theory*, vol. 30, no. 6, pp. 1272–1314, 2014.
- [7] H. Zou, "The adaptive lasso and its oracle properties," *Journal of the American Statistical Association*, vol. 101, no. 476, pp. 1418–1429, 2006.
- [8] A. Belloni and V. Chernozhukov, " ℓ_1 -penalized quantile regression in high-dimensional sparse models," *The Annals of Statistics*, vol. 39, no. 1, pp. 82–130, 02 2011. [Online]. Available: <http://dx.doi.org/10.1214/10-AOS827>
- [9] B. Kai, R. Li, and H. Zou, "Local composite quantile regression smoothing: an efficient and safe alternative to local polynomial regression," *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, vol. 72, no. 1, pp. 49–69, 2010.
- [10] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society, Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.
- [11] J. Fan and R. Li, "Variable selection via nonconcave penalized likelihood and its oracle properties," *Journal of the American Statistical Association*, vol. 96, no. 456, pp. 1348–1360, 2001. [Online]. Available: <http://pubs.amstat.org/doi/abs/10.1198/016214501753382273>
- [12] C.-H. Zhang, "Nearly unbiased variable selection under minimax concave penalty," *The Annals of Statistics*, vol. 38, no. 2, pp. 894–942, 2010.
- [13] L. Wang, Y. Wu, and R. Li, "Quantile regression for analyzing heterogeneity in ultra-high dimension," *Journal of the American Statistical Association*, vol. 107, no. 497, pp. 214–222, 2012.
- [14] E. Candes and T. Tao, "The Dantzig selector: Statistical estimation when p is much larger than n ," *The Annals of Statistics*, vol. 35, no. 6, pp. 2313–2351, Dec. 2007.
- [15] P. J. Bickel, Y. Ritov, and A. B. Tsybakov, "Simultaneous analysis of Lasso and Dantzig selector," *The Annals of Statistics*, vol. 37, no. 4, pp. 1705–1732, Aug. 2009.
- [16] A. Belloni, V. Chernozhukov, and L. Wang, "Square-root lasso: pivotal recovery of sparse signals via conic programming," *Biometrika*, vol. 98, no. 4, pp. 791–806, 2011.

- [17] J. Fan and J. Lv, "Nonconcave penalized likelihood with NP-dimensionality," *IEEE Transactions on Information Theory*, vol. 57, no. 8, pp. 5467–5484, Aug. 2011.
- [18] J. Fan, L. Xue, and H. Zou, "Strong oracle optimality of folded concave penalized estimation," *The Annals of Statistics*, vol. 42, no. 3, pp. 819–849, 06 2014.
- [19] J. Lv and Y. Fan, "A unified approach to model selection and sparse recovery using regularized least squares," *The Annals of Statistics*, vol. 37, no. 6A, pp. 3498–3528, Dec. 2009.
- [20] H. Zou and R. Li, "One-step sparse estimates in nonconcave penalized likelihood models," *The Annals of Statistics*, vol. 36, no. 4, pp. 1509–1533, Aug. 2008.
- [21] Y. Kim, J.-J. Jeon *et al.*, "Consistent model selection criteria for quadratically supported risks," *The Annals of Statistics*, vol. 44, no. 6, pp. 2467–2496, 2016.
- [22] E. R. Lee, H. Noh, and B. U. Park, "Model selection via bayesian information criterion for quantile regression models," *Journal of the American Statistical Association*, vol. 109, no. 505, pp. 216–229, 2014.
- [23] Y. Gu, J. Fan, L. Kong, S. Ma, and H. Zou, "ADMM for high-dimensional sparse penalized quantile regression," *Technometrics*, vol. 60, no. 3, pp. 319–331, 2018. [Online]. Available: <https://doi.org/10.1080/00401706.2017.1345703>
- [24] D. Pollard, "Empirical processes: Theory and applications," *NSF-CBMS Regional Conference Series in Probability and Statistics*, vol. 2, pp. i–86, 1990. [Online]. Available: <http://www.jstor.org/stable/4153175>
- [25] P. Bühlmann and S. van de Geer, *Statistics for High-Dimensional Data: Methods, Theory and Applications*. Springer, 2011.
- [26] A. van der Vaart and J. Wellner, *Weak convergence and empirical processes*. Springer, New York, 1996.
- [27] M. Ledoux and M. Talagrand, *Probability in Banach Spaces, Isoperimetry and Processes*, 1st ed. Springer-Verlag Berlin Heidelberg, 1991.
- [28] U. Haagerup, "The best constants in the Khintchine inequality," *Studia Mathematica*, vol. 3, no. 70, pp. 231–283, 1981.
- [29] D. Bertsimas and J. N. Tsitsiklis, *Introduction to linear optimization*. Athena Scientific, 1997.
- [30] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein, "Distributed optimization and statistical learning via the alternating direction method of multipliers," *Foundations and Trends® in Machine Learning*, vol. 3, no. 1, pp. 1–122, 2011.