

VARIABLE SELECTION IN QUANTILE REGRESSION

Yichao Wu and Yufeng Liu

North Carolina State University and University of North Carolina

Supplementary Material

This note contains Lemmas 2 and 3, and technical proofs.

Lemma 2 (Convexity Lemma). *Let $\{h_n(\mathbf{u}) : \mathbf{u} \in \mathbf{U}\}$ be a sequence of random convex functions defined on a convex, open subset \mathbf{U} of \mathbb{R}^d . Suppose $h(\mathbf{u})$ is a real-valued function on \mathbf{U} for which $h_n(\mathbf{u}) \rightarrow h(\mathbf{u})$ in probability, for each $\mathbf{u} \in \mathbf{U}$. Then for each compact subset K of \mathbf{U} ,*

$$\sup_{\mathbf{u} \in K} |h_n(\mathbf{u}) - h(\mathbf{u})| \rightarrow 0 \text{ in probability.}$$

The function $h(\cdot)$ is necessarily convex on \mathbf{U} .

Proof of Lemma 2. There are many versions of the proof for this well known Convexity Lemma. To save space, we skip its proof. Interested readers are referred to Pollard (1991). \square

Denote a linear approximation to $\rho_\tau(\epsilon_i - t)$ by $D_i = (1 - \tau)\{\epsilon_i < 0\} - \tau\{\epsilon_i \geq 0\}$. One intuitive interpretation of D_i is that D_i can be thought of as the first derivative of $\rho_\tau(\epsilon_i - t)$ at $t = 0$ (cf Pollard, 1991). Moreover, the condition that ϵ_i has the τ -th quantile zero implies $E(D_i) = 0$. Define $R_{i,n}(\mathbf{u}) = \rho_\tau(\epsilon_i - \mathbf{x}_i^T \mathbf{u} / \sqrt{n}) - \rho_\tau(\epsilon_i) - D_i \mathbf{x}_i^T \mathbf{u} / \sqrt{n}$, $W_n = \sum_{i=1}^n D_i \mathbf{x}_i / \sqrt{n}$, and $W_{n,11} = \sum_{i=1}^n D_i \mathbf{x}_i / \sqrt{n}$. Then $W_n \xrightarrow{L} N(\mathbf{0}, \tau(1 - \tau)\Sigma)$ and $W_{n,11} \xrightarrow{L} N(\mathbf{0}, \tau(1 - \tau)\Sigma_{11})$.

Lemma 3. *For model (3.1) with true parameter β_0 , denote $G_n(\mathbf{u}) = \sum_{i=1}^n [\rho_\tau(\epsilon_i - \mathbf{x}_i^T \mathbf{u} / \sqrt{n}) - \rho_\tau(\epsilon_i)]$, where $\epsilon_i = y_i - \mathbf{x}_i^T \beta_0$. Under Conditions (i) and (ii), we have, for any fixed \mathbf{u} ,*

$$G_n(\mathbf{u}) = \frac{f(0)}{2} \mathbf{u}^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \mathbf{u} + W_n^T \mathbf{u} + o_p(1). \tag{1}$$

Proof of Lemma 3. Note first that Condition (i) ensures that the function $M(t) = E(\rho_\tau(\epsilon_i - t) - \rho_\tau(\epsilon_i))$ has a unique minimizer at zero, and its Taylor expansion at origin has the following form $M(t) = \frac{f(0)}{2} t^2 + o(t^2)$. Hence, for large n , we have

$$\begin{aligned} E(G_n(\mathbf{u})) &= \sum_{i=1}^n M\left(\frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}}\right) = \sum_{i=1}^n \left[\frac{f(0)}{2} \left(\frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}}\right)^2 + o\left(\left(\frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}}\right)^2\right) \right] \\ &= \frac{f(0)}{2n} \mathbf{u}^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{u} + o\left(\frac{1}{2n} \mathbf{u}^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{u}\right). \end{aligned}$$

So, under Condition (ii), we have $E(G_n(\mathbf{u})) = \frac{f(0)}{2n} \mathbf{u}^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{u} + o(1)$.

Hence $G_n(\mathbf{u}) = E(G_n(\mathbf{u})) + W_n^T \mathbf{u} + \sum_{i=1}^n (R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u}))$. By routine calculation, we get $|R_{i,n}(\mathbf{u})| \leq |\mathbf{x}_i^T \mathbf{u} / \sqrt{n}| \cdot \{|\epsilon_i| \leq |\mathbf{x}_i^T \mathbf{u} / \sqrt{n}|\}$. For fixed \mathbf{u} , due to the cancelation of cross-product terms, we get

$$\begin{aligned}
E\left(\sum_{i=1}^n [R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u})]^2\right) &= \sum_{i=1}^n E(R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u}))^2 \\
&\leq \sum_{i=1}^n E(R_{i,n}(\mathbf{u}))^2 \\
&\leq \sum_{i=1}^n \left[\left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right|^2 E\{|\epsilon_i| \leq \left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right|\} \right] \\
&\leq \left(\sum_{i=1}^n \left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right|^2 \right) E\{|\epsilon_i| \leq \frac{\|\mathbf{u}\|}{\sqrt{n}} \max_{j=1,2,\dots,n} \|\mathbf{x}_j\|\} \\
&\rightarrow 0
\end{aligned} \tag{2}$$

as in Pollard (1991), where $\|\cdot\|$ denotes the Euclidean norm operator. Here the last step converging to zero holds because

$$\begin{aligned}
\sum_{i=1}^n \left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right|^2 &= \mathbf{u}^T \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T / n \right) \mathbf{u} \rightarrow \mathbf{u}^T \Sigma \mathbf{u} \\
\max_{j=1,2,\dots,n} \|\mathbf{x}_j\| / \sqrt{n} &\rightarrow 0 \text{ due to } \frac{\sum_{i=1}^n \|\mathbf{x}_i\|^2}{n} \rightarrow \text{trace}(\Sigma).
\end{aligned}$$

Equation (2) implies that $\sum_{i=1}^n (R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u})) = o_p(1)$. This completes the proof. \square

Before we start the proof of Theorem 1, we want to point that $W_n^T \mathbf{u} = E(W_n^T \mathbf{u}) + O_p(\sqrt{\text{Var}(W_n^T \mathbf{u})})$, together with $\text{Var}(W_n^T \mathbf{u}) = \sum_{i=1}^n E(D_i \mathbf{x}_i^T \mathbf{u} / \sqrt{n})^2 = \tau(1-\tau) \mathbf{u}^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \mathbf{u}$, implies that $W_n^T \mathbf{u} = O_p\left(\sqrt{\tau(1-\tau) \mathbf{u}^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \mathbf{u}}\right)$.

Proof of Theorem 1. We use the same strategy as in Fan and Li (2001). To prove Theorem 1, it is enough to show that for any given $\delta > 0$, there exists a large constant C such that

$$P\left\{ \inf_{\|\mathbf{u}\|=C} Q(\beta_0 + \mathbf{u}/\sqrt{n}) > Q(\beta_0) \right\} \geq 1 - \delta \tag{3}$$

which implies that with probability at least $1 - \delta$ there exists a local minimum in the ball $\{\beta_0 + \mathbf{u}/\sqrt{n} : \|\mathbf{u}\| \leq C\}$. This in turn implies that there exists a local minimizer such that

$\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\| = O_p(1/\sqrt{n})$, which is exactly what we want to show. Note that

$$\begin{aligned} & Q(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n}) - Q(\boldsymbol{\beta}_0) \\ &= \sum_{i=1}^n [\rho_\tau(y_i - \mathbf{x}_i^T(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n})) - \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0)] + n \sum_{j=1}^d [p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|)] \\ &\geq \sum_{i=1}^n [\rho_\tau(y_i - \mathbf{x}_i^T(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n})) - \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0)] + n \sum_{j=1}^s [p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|)], \end{aligned}$$

where s is the number of components in $\boldsymbol{\beta}_{10}$ and β_{j0} denotes the j -th component of $\boldsymbol{\beta}_{10}$. Due to Lemma 3, the first term on the right hand side is exactly $G_n(\mathbf{u}) = \frac{f(0)}{2} \mathbf{u}^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \mathbf{u} + W_n^T \mathbf{u} + o_p(1)$ for any fixed \mathbf{u} .

By applying the Convexity Lemma (Lemma 2) to $h_n(\mathbf{u}) = G_n(\mathbf{u}) - W_n^T \mathbf{u}$, we can strengthen this pointwise convergence to uniform convergence on any compact subset of \mathbb{R}^d .

Note that, for large n ,

$$n \sum_{j=1}^s [p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|)] = 0 \quad (4)$$

uniformly in any compact set of \mathbb{R}^d due to the facts that $\beta_{j0} > 0$ for $j = 1, 2, \dots, s$, SCAD penalty is flat for coefficient of magnitude larger than $a\lambda_n$, and $\lambda_n \rightarrow 0$.

Based on all the above, $Q(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n}) - Q(\boldsymbol{\beta}_0)$ is dominated by the quadratic term $f(0)\mathbf{u}^T (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T) \mathbf{u} / (2n)$ for $\|\mathbf{u}\|$ equal to sufficiently large C . Hence Condition (ii) implies that (3) holds as we have desired and this completes the proof. \square

Proof of Lemma 1. For any $\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10} = O_p(n^{-1/2})$, $0 < \|\boldsymbol{\beta}_2\| \leq Cn^{-1/2}$,

$$\begin{aligned} & Q((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T) - Q((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T) \\ &= [Q((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T) - Q((\boldsymbol{\beta}_{10}^T, \mathbf{0}^T)^T)] - [Q((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T) - Q((\boldsymbol{\beta}_{10}^T, \mathbf{0}^T)^T)] \\ &= G_n(\sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T)^T) - G_n(\sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T) - n \sum_{j=s+1}^d p_{\lambda_n}(|\beta_j|) \quad (5) \\ &= \frac{f(0)}{2} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T) \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T)^T + \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T) W_n \\ &\quad - \frac{f(0)}{2} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T) \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T - \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T) W_n \\ &\quad + o(1) + o_p(1) - n \sum_{j=s+1}^d p_{\lambda_n}(|\beta_j|) \end{aligned}$$

The conditions $\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10} = O_p(n^{-1/2})$ and $0 < \|\boldsymbol{\beta}_2\| \leq Cn^{-1/2}$ imply that

$$\frac{f(0)}{2} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T) \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T)^T = O_p(1)$$

$$\frac{f(0)}{2} \sqrt{n} ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T) \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n} ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T = O_p(1)$$

and

$$\begin{aligned} & \sqrt{n} ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T) W_n - \sqrt{n} ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T) W_n = -\sqrt{n} (\mathbf{0}^T, \boldsymbol{\beta}_2^T) W_n \\ &= \sqrt{n} \sqrt{\tau(1-\tau) \boldsymbol{\beta}_2^T \Sigma_{22} \boldsymbol{\beta}_2} (1 + o_p(1)). \end{aligned}$$

Note that

$$\begin{aligned} n \sum_{j=s+1}^d p_{\lambda_n}(|\beta_j|) &\geq n \lambda_n \left(\liminf_{\lambda \rightarrow 0} \liminf_{\theta \rightarrow 0^+} \frac{p'_\lambda(\theta)}{\lambda} \right) \left(\sum_{j=s+1}^d |\beta_j| \right) (1 + o(1)) \\ &= n \lambda_n \left(\sum_{j=s+1}^d |\beta_j| \right) (1 + o(1)), \end{aligned}$$

where the last step follows based on the fact that $\liminf_{\lambda \rightarrow 0} \liminf_{\theta \rightarrow 0^+} \frac{p'_\lambda(\theta)}{\lambda} = 1$.

Then $\sqrt{n} \lambda_n \rightarrow \infty$ implies that $n \lambda_n = \sqrt{n}(\sqrt{n} \lambda_n)$ is of higher order than \sqrt{n} . This implies that, in (5), the last term dominates in magnitude and, as a result, $Q((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T) - Q((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T) < 0$ for large n . The completes the proof. \square

Proof of Theorem 2. Similarly as in Fan and Li (2001), Part (a) holds simply due to Lemma 1. Next we prove part (b). By Theorem 1, we can show that there exists a root- n consistent minimizer $\hat{\boldsymbol{\beta}}_1$ of $Q((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T)$ as a function of $\boldsymbol{\beta}_1$.

From the proof of Theorem 1, we see that $\sqrt{n}(\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_{10})$ minimizes $G_n((\boldsymbol{\theta}^T, \mathbf{0}^T)^T) + n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0} + \frac{\theta_j}{\sqrt{n}}|)$, where $\boldsymbol{\theta} = (\theta_1, \theta_2, \dots, \theta_s)^T \in \mathbb{R}^s$. Notice that, as in the proof of Theorem 1, Lemma 3 and the convexity lemma imply that

$$\begin{aligned} G_n((\boldsymbol{\theta}^T, \mathbf{0}^T)^T) &= \frac{f(0)}{2} (\boldsymbol{\theta}^T, \mathbf{0}^T) \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} (\boldsymbol{\theta}^T, \mathbf{0}^T)^T + (\boldsymbol{\theta}^T, \mathbf{0}^T) W_n + o_p(1) \\ &= \frac{f(0)}{2} \boldsymbol{\theta}^T \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n} \boldsymbol{\theta} + \boldsymbol{\theta}^T \sum_{i=1}^n D_i \mathbf{x}_{i1} / \sqrt{n} + o_p(1) \end{aligned}$$

uniformly in any compact subset of \mathbb{R}^s . Notice that, for large n , $n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0} + \theta_j / \sqrt{n}|) = n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0}|)$ uniformly in any compact set of \mathbb{R}^s , due to (4). Hence we have

$$\begin{aligned} & G_n((\boldsymbol{\theta}^T, \mathbf{0}^T)^T) + n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0} + \frac{\theta_j}{\sqrt{n}}|) \\ &= \frac{1}{2} \boldsymbol{\theta}^T (f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n}) \boldsymbol{\theta} + (\sum_{i=1}^n D_i \mathbf{x}_{i1} / \sqrt{n})^T \boldsymbol{\theta} + n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0}|) + r_n(\boldsymbol{\theta}) \\ &= \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\zeta}_n)^T (f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n}) (\boldsymbol{\theta} - \boldsymbol{\zeta}_n) - \frac{1}{2} \boldsymbol{\zeta}_n^T (f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n}) \boldsymbol{\zeta}_n + n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0}|) + r_n(\boldsymbol{\theta}) \end{aligned}$$

where $\zeta_n = -(f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n})^{-1} W_{n,11}$ and the residual $r_n(\boldsymbol{\theta}) \rightarrow 0$ in probability uniformly in any compact subset of \mathbb{R}^s . Notice further that the term $n \sum_{j=1}^s p_{\lambda_n}(|\beta_{j0}|)$ does not depend on $\boldsymbol{\theta}$. So this implies that, for large n , the local minimizer $\hat{\boldsymbol{\theta}}$ is very close to ζ_n and satisfies $\hat{\boldsymbol{\theta}} - \zeta_n = o_p(1)$.

That is, the minimizer $\hat{\boldsymbol{\theta}}$ satisfies $\hat{\boldsymbol{\theta}} = -(f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n})^{-1} (\sum_{i=1}^n D_i \mathbf{x}_{i1} / \sqrt{n}) + o_p(1)$. Hence $\sqrt{n}(\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_{10}) = -(f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n})^{-1} (\sum_{i=1}^n D_i \mathbf{x}_{i1} / \sqrt{n}) + o_p(1)$. Applying *Slutsky's* theorem, we get $\sqrt{n}f(0)\Sigma_{11}(\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_{10}) \xrightarrow{\mathcal{L}} N(\mathbf{0}, \tau(1 - \tau)\Sigma_{11})$. This completes the proof. \square

Proof of Theorem 3. Note that

$$\begin{aligned} & Q_1(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n}) - Q_1(\boldsymbol{\beta}_0) \\ &= \sum_{i=1}^n [\rho_\tau(y_i - \mathbf{x}_i^T(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n})) - \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_0)] + n\lambda_n \sum_{j=1}^d [\tilde{w}_j | \beta_{j0} + u_j/\sqrt{n} | - \tilde{w}_j | \beta_{j0} |]. \end{aligned}$$

We consider the second term first, for $j = 1, 2, \dots, s$, we have $\beta_{j0} \neq 0$; as a result, $\tilde{w}_j \xrightarrow{P} |\beta_{j0}|^{-\gamma}$; hence $n\lambda_n[\tilde{w}_j | \beta_{j0} + u_j/\sqrt{n} | - \tilde{w}_j | \beta_{j0} |] \xrightarrow{P} 0$ as $\sqrt{n}(|\beta_{j0} + u_j/\sqrt{n}| - |\beta_{j0}|) \rightarrow u_j \text{sign}(\beta_{j0})$ and $\sqrt{n}\lambda_n \rightarrow 0$. On the other hand, for $j = s+1, s+2, \dots, d$, the true coefficient $\beta_{j0} = 0$; so $\sqrt{n}\lambda_n \tilde{w}_j = n^{(1+\gamma)/2} \lambda_n (\sqrt{n}|\tilde{\beta}_j|)^{-\gamma}$ with $\sqrt{n}\tilde{\beta}_j = O_p(1)$; so it follows that $n\lambda_n[\tilde{w}_j | \beta_{j0} + u_j/\sqrt{n} | - \tilde{w}_j | \beta_{j0} |] \xrightarrow{P} \infty$ when $u_j \neq 0$ and $= 0$ otherwise due to $\sqrt{n} | u_j/\sqrt{n} | = |u_j|$ for large n . These facts and the result of Lemma 3 imply that

$$Q_1(\boldsymbol{\beta}_0 + \frac{\mathbf{u}}{\sqrt{n}}) - Q_1(\boldsymbol{\beta}_0) \xrightarrow{\mathcal{L}} V(\mathbf{u}) = \begin{cases} \frac{f(0)}{2} \mathbf{u}_1 \Sigma_{11} \mathbf{u}_1 + W_{n,11}^T \mathbf{u}_1 & \text{when } u_j = 0 \text{ for } j \geq s+1 \\ \infty & \text{otherwise,} \end{cases}$$

where $\mathbf{u}_1 = (u_1, u_2, \dots, u_s)^T$. Noticing that $Q_1(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n}) - Q_1(\boldsymbol{\beta}_0)$ is convex in \mathbf{u} and V has a unique minimizer, the epi-convergence results of Geyer (1994) imply that

$$\text{argmin } Q_1(\boldsymbol{\beta}_0 + \frac{\mathbf{u}}{\sqrt{n}}) = \sqrt{n}(\hat{\boldsymbol{\beta}}^{(AL)} - \boldsymbol{\beta}_0) \xrightarrow{\mathcal{L}} \text{argmin } V(\mathbf{u}),$$

which establishes the asymptotic normality part. Next we show the consistency property of model selection.

For any $\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10} = O_p(n^{-1/2})$, $0 < \|\boldsymbol{\beta}_2\| \leq Cn^{-1/2}$,

$$\begin{aligned} & Q_1((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T) - Q_1((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T) \\ &= [Q_1((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T) - Q_1((\boldsymbol{\beta}_{10}^T, \mathbf{0}^T)^T)] - [Q_1((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T) - Q_1((\boldsymbol{\beta}_{10}^T, \mathbf{0}^T)^T)] \\ &= G_n(\sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \mathbf{0}^T)^T) - G_n(\sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T) - n\lambda_n \sum_{j=s+1}^d \tilde{w}_j |\beta_j|. \end{aligned}$$

Note here the first two terms are exactly the same as in (5) and hence can be bounded similarly.

However the third term goes to $-\infty$ as $n \rightarrow \infty$ due to the following

$$n\lambda_n \sum_{j=s+1}^d \tilde{w}_j |\beta_j| = (n^{(1+\gamma)/2}\lambda_n)\sqrt{n} \sum_{j=s+1}^d \left| (\sqrt{n}|\tilde{\beta}_j|)^{-\gamma} \right| |\beta_j| \rightarrow \infty.$$

Hence the condition that $n^{(1+\gamma)/2}\lambda_n \rightarrow \infty$ implies that $n\lambda_n \sum_{j=s+1}^d \tilde{w}_j |\beta_j|$ is of higher order than any other terms and dominates as a result. This in turn implies that $Q_1((\beta_1^T, \mathbf{0}^T)^T) - Q_1((\beta_1^T, \beta_2^T)^T) < 0$ for large n . This proves the consistency of model selection of adaptive lasso penalized quantile regression. \square

Proof of Corollary 1. Notice from the proofs of Theorem 1, Lemma 1, Theorem 2, and Theorem 3, it is enough to establish an asymptotic approximation similar as (1).

Note that the check function $\rho_\tau(\cdot)$ can be rewritten as $\rho_\tau(r) = |r|/2 + (\tau - 1/2)r$. Hence,

$$\begin{aligned} G_n(\mathbf{u}) &= \sum_{i=1}^n [\rho_\tau(\epsilon_i - \mathbf{x}_i^T \mathbf{u} / \sqrt{n}) - \rho_\tau(\epsilon_i)] \\ &= \sum_{i=1}^n \frac{-\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \left(\frac{\text{sign}(\epsilon_i)}{2} + \left(\tau - \frac{1}{2}\right) \right) + \sum_{i=1}^n \int_0^{\frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}}} (I(\epsilon_i \leq s) - I(\epsilon_i \leq 0)) ds \end{aligned}$$

as in Knight (1999). By the same argument as in Knight (1999), we get $G_n(\mathbf{u}) \xrightarrow{\mathcal{L}} -\mathbf{u}^T \mathbf{V} + \zeta(\mathbf{u})$ for some multivariate normal random vector \mathbf{V} with mean zero. Then the result follows from the strictly convexity of $\zeta(\cdot)$. \square