

# VARIABLE SELECTION IN QUANTILE REGRESSION

Yichao Wu and Yufeng Liu

*Princeton University and University of North Carolina*

*Abstract:* After its inception in Koenker and Bassett (1978), quantile regression has become an important and widely used technique to study the whole conditional distribution of a response variable and grown into an important tool of applied statistics over the last three decades. In this work, we focus on the variable selection aspect of penalized quantile regression. Under some mild conditions, we demonstrate the oracle properties of the SCAD and adaptive-LASSO penalized quantile regressions. For the SCAD penalty, despite its good asymptotic properties, the corresponding optimization problem is non-convex and as a result much harder to solve. In this work, we take advantage of the decomposition of the SCAD penalty function as the difference of two convex functions and propose to solve the corresponding optimization using the Difference Convex Algorithm (DCA).

*Key words and phrases:* DCA, LASSO, oracle, quantile regression, SCAD, variable selection.

## 1. Introduction

At the heart of statistics lies regression. Ordinary least squares regression (OLS) estimates the conditional mean function, *i.e.*, the mean response as a function of the regressors or predictors. Least absolute deviation regression (LADR) estimates the conditional median function, which has been shown to be more robust to outliers. In the seminal paper of Koenker and Bassett (1978), they generalized the idea of LADR and introduced quantile regression (QR) to estimate the conditional quantile function of the response. As a result, QR provides much more information about the conditional distribution of a response variable. It includes LADR as a special case. After its introduction, QR has attracted tremendous interest in the literature. It has been applied in many different areas: economics (Hendricks and Koenker, 1992; Koenker and Hallock, 2001), survival analysis (Yang, 1999; Koenker and Geling, 2001), Microarray study (Wang and He, 2007), growth chart (Wei, Pere, Koenker and He, 2006; Wei and He, 2006),

and so on. Li, Liu and Zhu (2007) considered quantile regression in reproducing kernel Hilbert spaces and proposed a very efficient algorithm to compute its entire solution path with respect to the tuning parameter.

Variable selection plays an important role in the model building process. In practice, it is very common that there are a large number of candidate predictor variables available and they are included in the initial stage of modeling for the consideration of removing potential modeling bias (Fan and Li, 2001). However, it is undesirable to keep irrelevant predictors in the final model since this makes it difficult to interpret the resultant model and may decrease its predictive ability. In the regularization framework, many different types of penalties have been introduced to achieve variable selection. The  $L_1$  penalty was used in the LASSO proposed by Tibshirani (1996) for variable selection. Fan and Li (2001) proposed a unified approach via nonconcave penalized least squares regression, which simultaneously performs variable selection and coefficient estimation. By choosing an appropriate nonconcave penalty function, this method keeps many merits of the best subset selection and ridge regression: it produces sparse solution; ensures the stability of model selection; provides unbiased estimates for large coefficients. These are the three desirable properties of a good penalty (Fan and Li, 2001). An example of such nonconcave penalties is the smoothly clipped absolute deviation (SCAD) function first introduced in Fan (1997) and studied further by Fan and Li (2001) to show its oracle properties in the penalized likelihood setting. Later on, a series of papers (Fan and Li, 2002, 2004; Fan and Peng, 2004; Hunter and Li, 2005) studied its further properties and new algorithms.

By using adaptive weights for penalizing different coefficients in the LASSO penalty, Zou (2006) introduced the adaptive LASSO and demonstrated its oracle properties. Similar results were also established in Yuan and Lin (2007) and Zhao and Yu (2006). Zhang and Lu (2007) studied the adaptive LASSO in proportional hazard models. Candès and Tao (2007) and Fan and Lv (2006) studied variable selection in the setting of dimensionality higher than the sample size.

Previously, Koenker (2004) applied the LASSO penalty to the mixed-effect quantile regression model for longitudinal data to encourage shrinkage in estimating the random effects. Li and Zhu (2005) developed the solution path of the  $L_1$  penalized quantile regression. Wang, Li and Jiang (2007) considered LADR

with the adaptive LASSO penalty. To our limited knowledge, there still lacks of study on variable selection in penalized quantile regression. In this work, we try to fill this void. Notice that the loss function used in quantile regression is not differentiable at the origin and, as a result, the general oracle properties for nonconcave penalized likelihood (Fan and Li, 2001) do not apply directly. Here, we extend the oracle properties of the SCAD and adaptive-LASSO penalties to the context of penalized quantile regression, including the LADR by Wang et al. (2007) as a special case.

The SCAD penalty is nonconvex, and consequently it is hard to solve the corresponding optimization problem. Motivated by the fact that the SCAD penalty function can be decomposed as the difference of two convex functions, we propose to use the Difference Convex algorithm (DCA) (see An and Tao, 1997) to solve the corresponding non-convex optimization problem. DCA minimizes a non-convex objective function by solving a sequence of convex minimization problems. At each iteration, it approximates the second convex function by a linear function. As a result, the objective function at each step is convex and it is much easier to optimize than the original non-convex optimization problem. In this sense, DCA turns out to be an instance of the MM algorithm since, at each step, DCA majorizes the nonconvex objective function and then performs minimization. One difference between DCA and Hunter and Li (2005)'s MM is that, at each iteration, DCA majorizes the nonconvex function using a linear approximation while Hunter and Li (2005)'s MM uses a quadratic approximation. We opt for DCA due to its clean formulation and simple implementation. In particular, for quantile regression, the resulting optimization at each iteration is a linear programming problem, thus more efficient. We recently learned that Zou and Li (2007) proposed a local linear approximation algorithm (LLA) to solve the SCAD optimization problem. Although both DCA and LLA perform iterative linear programming, unlike the LLA, DCA does not enforce symmetry in the approximation of the SCAD penalty.

The rest of the paper is organized as follows. Penalized quantile regressions with the SCAD and adaptive-LASSO penalties are introduced in Section 2. We present the asymptotic properties of the SCAD and adaptive-LASSO penalized quantile regressions in Section 3. Algorithms for handling their corresponding

optimization problems are proposed in Section 4. Sections 5 and 6 present numerical results on both simulated and real data respectively. We conclude the paper with Section 7.

## 2. Penalized linear quantile regression

Consider a sample  $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$  of size  $n$  from some unknown population, where  $\mathbf{x}_i \in \mathbb{R}^d$ . The conditional  $\tau$ -th quantile function  $f_\tau(\mathbf{x})$  is defined such that  $P(Y \leq f_\tau(\mathbf{X}) | \mathbf{X} = \mathbf{x}) = \tau$ , for  $0 < \tau < 1$ . By tilting the absolute loss function, Koenker and Bassett (1978) introduced the check function which is defined by  $\rho_\tau(r) = \tau r$  if  $r > 0$ , and  $-(1-\tau)r$  otherwise. In their seminal paper (Koenker and Bassett, 1978), they demonstrated that the  $\tau$ -th conditional quantile function can be estimated by solving the following minimization problem

$$\min_{f_\tau \in \mathcal{F}} \sum_{i=1}^n \rho_\tau(y_i - f_\tau(\mathbf{x}_i)). \quad (2.1)$$

To avoid over-fitting and improve generalization ability, as in Koenker, Ng and Portnoy (1994) and Koenker (2004), we consider the penalized version of (2.1) in the regularization framework

$$\min_{f_\tau \in \mathcal{F}} \sum_{i=1}^n \rho_\tau(y_i - f_\tau(\mathbf{x}_i)) + \lambda J(f_\tau), \quad (2.2)$$

where  $\lambda \geq 0$  is the regularization parameter and  $J(f_\tau)$  denotes the roughness penalty of the function  $f_\tau(\cdot)$ .

In this work, we focus on linear quantile regression by setting  $f_\tau(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}_\tau$  where  $\boldsymbol{\beta}_\tau = (\beta_{\tau,1}, \beta_{\tau,2}, \dots, \beta_{\tau,d})^T$ , namely, the conditional quantile function is a linear function of the regressor  $\mathbf{x}$ . This form can be easily generalized to handle nonlinear quantile regression via basis expansion. For functions of linear form, there are many different types of penalty functions available, for example, the  $L_0$  penalty (also known as the entropy penalty) used in the best subset selection (Breiman, 1996), the  $L_1$  penalty (LASSO) (Tibshirani, 1996), the  $L_2$  penalty used in ridge regression (Hoerl and Kennard, 1988), the combination of the  $L_0$  and  $L_1$  penalties (Liu and Wu, 2007), and the  $L_q$  ( $q \geq 0$ ) penalties in bridge regression (Frank and Friedman, 1993). Fan and Li (2001) argued that a good penalty should possess the following three properties in its estimator:

unbiasedness, sparsity, and continuity. But, unfortunately, none of the  $L_q$  penalty family satisfies these three properties simultaneously. To remedy this problem, Fan and Li (2001) studied the SCAD penalty in the penalized likelihood setting, which achieves these three desirable properties simultaneously. Another penalty falling into this category is the adaptive-LASSO penalty studied by Zou (2006).

## 2.1 SCAD

Fan and Li (2001) demonstrated the oracle properties for the SCAD in the variable selection aspect and conjectured that the LASSO penalty does not possess the oracle properties. This conjecture was later confirmed by Zou (2006), who further proposed the adaptive LASSO and showed its oracle properties in penalized least squares regression.

The SCAD penalty is mathematically defined in terms of its first order derivative and is symmetric around the origin. For  $\theta > 0$ , its first order derivative is given by

$$p'_\lambda(\theta) = \lambda \left\{ I(\theta \leq \lambda) + \frac{(a\lambda - \theta)_+}{(a-1)\lambda} I(\theta > \lambda) \right\}, \quad (2.3)$$

where  $a > 2$  and  $\lambda > 0$  are tuning parameters. Note that the SCAD penalty function is symmetric, non-convex on  $[0, \infty)$ , and singular at the origin. One instance of the SCAD penalty function is plotted in the right panel of Figure 4.1. We can see that, around the origin, it takes the same form as the LASSO penalty and this leads to its sparsity property. But different from the LASSO penalty, the SCAD penalizes large coefficients equally while the LASSO penalty increases linearly as the magnitude of the coefficient increases. In this way, the SCAD results in unbiased penalized estimators for large coefficients. After putting the SCAD penalty in (2.2) with linear function  $f(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta}_\tau$ , the SCAD penalized quantile regression solves the following minimization problem

$$\min_{\boldsymbol{\beta}_\tau} \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\tau) + \sum_{j=1}^d p_\lambda(\beta_{\tau,j}).$$

## 2.2 Adaptive-LASSO

The adaptive-LASSO can be viewed as a generalization of the LASSO penalty. Basically the idea is to penalize the coefficients of different covariates at a different level by using adaptive weights. In the case of least squares regression, Zou

(2006) proposed to use as weights the reciprocal of the ordinal least squares estimates raised to some power. The straightforward generalization, for our case of quantile regression, is to use that of the solution of non-penalized quantile regression as weights. More explicitly, denote the solution of linear quantile regression by  $\tilde{\boldsymbol{\beta}}_\tau$ , namely,

$$\tilde{\boldsymbol{\beta}}_\tau = \underset{\boldsymbol{\beta}_\tau}{\operatorname{argmin}} \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\tau). \quad (2.4)$$

It can be shown that  $\tilde{\boldsymbol{\beta}}_\tau$  is a root- $n$  consistent estimator of  $\boldsymbol{\beta}_\tau$ . Then the adaptive-LASSO penalized quantile regression minimizes

$$\sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\tau) + \lambda \sum_{j=1}^d \tilde{w}_j |\beta_{\tau,j}|$$

with respect to  $\boldsymbol{\beta}_\tau$ , where the weights are set to be  $\tilde{w}_j = 1/|\tilde{\beta}_{\tau,j}|^\gamma$ ,  $j = 1, 2, \dots, d$ ; for some appropriately chosen  $\gamma > 0$ .

### 3. Asymptotic properties

In this section, we establish the oracle properties of the SCAD or adaptive-LASSO penalized quantile regression. To prove the asymptotic properties of penalized quantile regression, we need to lay out some basic assumptions on our data. We assume that our data set  $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$  consists of  $n$  observations from the following linear model:

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i = \mathbf{x}_{i1}^T \boldsymbol{\beta}_1 + \mathbf{x}_{i2}^T \boldsymbol{\beta}_2 + \epsilon_i, \quad i = 1, 2, \dots, n, \quad (3.1)$$

with  $P(\epsilon_i < 0) = \tau$  as in Condition (i). Here  $\mathbf{x}_i = (\mathbf{x}_{i1}^T, \mathbf{x}_{i2}^T)^T$ ,  $\boldsymbol{\beta} = (\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T$ ,  $\mathbf{x}_{i1} \in \mathbb{R}^s$ ,  $\mathbf{x}_{i2} \in \mathbb{R}^{d-s}$ , and the true regression coefficients are  $\boldsymbol{\beta}_1 = \boldsymbol{\beta}_{10}$  with each component being nonzero and  $\boldsymbol{\beta}_2 = \boldsymbol{\beta}_{20} = \mathbf{0}$  (as a result  $\boldsymbol{\beta}_0 = (\boldsymbol{\beta}_{10}^T, \boldsymbol{\beta}_{20}^T)^T$ ). This means that the first  $s$  regressors are important while the remaining  $p - s$  are noise variables.

For our theoretical results, we enforce the following technical conditions:

- (i) Error assumption (cf Pollard, 1991): The regression errors  $\{\epsilon_i\}$  are independent. They are identically distributed, with the  $\tau$ -th quantile zero and a continuous, positive density  $f(\cdot)$  in a neighborhood of zero.

- (ii) The design  $\mathbf{x}_i; i = 1, 2, \dots, n$ , is a deterministic sequence satisfying that the limit of  $(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T)/n$  as  $n \rightarrow \infty$  is a positive definite matrix, *i.e.*, there exists a positive definite matrix  $\Sigma$  such that  $\lim_{n \rightarrow \infty} (\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T)/n = \Sigma$ . Denote the top-left  $s$ -by- $s$  submatrix of  $\Sigma$  by  $\Sigma_{11}$  and the right-bottom  $(d - s)$ -by- $(d - s)$  submatrix of  $\Sigma$  by  $\Sigma_{22}$ .

### 3.1 SCAD penalty

The SCAD penalized quantile regression solves  $\min_{\boldsymbol{\beta}} Q(\boldsymbol{\beta})$ , where  $Q(\boldsymbol{\beta}) = \sum_{i=1}^n \rho_{\tau}(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) + n \sum_{j=1}^d p_{\lambda_n}(|\beta_j|)$ . As in Fan and Li (2001), we establish the root- $n$  consistency of our SCAD penalized estimator as in Theorem 1 when the tuning parameter  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ .

**Theorem 1 (consistency).** *Consider a sample  $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$  from model (3.1) satisfying Conditions (i) and (ii) with i.i.d.  $\epsilon_i$ 's. If  $\lambda_n \rightarrow 0$ , there exists a local minimizer  $\hat{\boldsymbol{\beta}}$  such that  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_0\| = O_p(n^{-1/2})$ .*

Under some further conditions, the sparsity property  $\hat{\boldsymbol{\beta}}_2 = \mathbf{0}$  of the SCAD penalized estimator can be obtained as in Lemma 1.

**Lemma 1 (Sparsity).** *For a sample  $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$  from model (3.1) satisfying Conditions (i) and (ii) with i.i.d.  $\epsilon_i$ 's, if  $\lambda_n \rightarrow 0$  and  $\sqrt{n}\lambda_n \rightarrow \infty$  as  $n \rightarrow \infty$ , then with probability tending to one, for any given  $\boldsymbol{\beta}_1$  satisfying  $\|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}\| = O_p(n^{-1/2})$  and any constant  $C$ ,*

$$Q((\boldsymbol{\beta}_1^T, \mathbf{0}^T)^T) = \min_{\|\boldsymbol{\beta}_2\| \leq Cn^{-1/2}} Q((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T).$$

Our next theorem shows that the penalized estimator also enjoys the popular asymptotic oracle property.

**Theorem 2 (Oracle).** *For a sample  $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$  from model (3.1) satisfying Conditions (i) and (ii) with i.i.d.  $\epsilon_i$ 's, if  $\lambda_n \rightarrow 0$  and  $\sqrt{n}\lambda_n \rightarrow \infty$  as  $n \rightarrow \infty$ , then with probability tending to one the root- $n$  consistent local minimizer  $\hat{\boldsymbol{\beta}} = (\hat{\boldsymbol{\beta}}_1^T, \hat{\boldsymbol{\beta}}_2^T)^T$  in Theorem 1 must satisfy:*

- (a) *Sparsity:  $\hat{\boldsymbol{\beta}}_2 = \mathbf{0}$ .*
- (b) *Asymptotic normality:  $\sqrt{n}(\hat{\boldsymbol{\beta}}_1 - \boldsymbol{\beta}_{10}) \xrightarrow{\mathcal{L}} N(\mathbf{0}, \tau(1 - \tau)\Sigma_{11}^{-1}/f(0)^2)$ , where  $\Sigma_{11}$  is defined in Condition (ii).*

**Remark 1.** Notice that the main difference between penalized quantile regression and more general penalized likelihood as considered in Fan and Li (2001) is that the check function in penalized quantile regression is non-differentiable at the origin. To handle the difficulty caused by this non-differentiability, we use the convexity lemma as previously used in Pollard (1991).

### 3.2 Adaptive-LASSO

The adaptive-LASSO penalized quantile regression solves  $\min_{\boldsymbol{\beta}} Q_1(\boldsymbol{\beta})$  where  $Q_1(\boldsymbol{\beta}) = \sum_{i=1}^n \rho_{\tau}(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) + n\lambda_n \sum_{j=1}^d \tilde{w}_j |\beta_j|$ . Denote  $\hat{\boldsymbol{\beta}}^{(AL)}$  as its solution.

**Theorem 3 (Oracle).** *Consider a sample  $\{(\mathbf{x}_i, y_i), i = 1, 2, \dots, n\}$  from model (3.1) satisfying Conditions (i) and (ii) with i.i.d.  $\epsilon_i$ 's. If  $\sqrt{n}\lambda_n \rightarrow 0$ , and  $n^{(\gamma+1)/2}\lambda_n \rightarrow \infty$ , then we have*

1. Sparsity:  $\hat{\boldsymbol{\beta}}_2^{(AL)} = \mathbf{0}$ .
2. Asymptotic normality:  $\sqrt{n}(\hat{\boldsymbol{\beta}}_1^{(AL)} - \boldsymbol{\beta}_{10}) \xrightarrow{\mathcal{L}} N(\mathbf{0}, \tau(1-\tau)\Sigma_{11}^{-1}/f(0)^2)$ .

### 3.3 Non i.i.d. random errors

The conclusions in Theorems 2 and 3 are based on the assumption of i.i.d. random errors. We can further extend the aforementioned oracle results to the case of non i.i.d. random errors. In the light of the work by Knight (1999), we make the following assumptions.

(N1) As  $n \rightarrow \infty$ ,  $\max_{1 \leq i \leq n} \mathbf{x}_i^T \mathbf{x}_i / n \rightarrow 0$ .

(N2) The random errors  $\epsilon_i$ 's are independent, but not identically distributed. Let  $F_i(t) = P(\epsilon_i \leq t)$  be the distribution function of  $\epsilon_i$ . We assume that each  $F_i(\cdot)$  is locally linear near zero (with a positive slope) and  $F_i(0) = \tau$  (i.e. the  $\tau$ -th quantile of  $\epsilon_i$  is zero).

Define  $\psi_{ni}(t) = \int_0^t \sqrt{n}(F_i(s/\sqrt{n}) - F_i(0))ds$ , which is a convex function for each  $n$  and  $i$ .

(N3) Assume that, for each  $\mathbf{u}$ ,  $\frac{1}{n} \sum_{i=1}^n \psi_{ni}(\mathbf{u}^T \mathbf{x}_i) \rightarrow \varsigma(\mathbf{u})$ , where  $\varsigma(\cdot)$  is a strictly convex function taking values in  $[0, \infty)$ .

**Corollary 1.** *Under Conditions (ii) and (N1), the results of Theorems 2 and 3 still hold in the case with non i.i.d. random errors satisfying (N2) and (N3).*



**Remark 2.** The assumption (N2) covers a class of general models with non *i.i.d.* random errors. For example, it includes the common location-scale shift model (Koenker, 2005). The proof follows directly using the results of Knight (1999). Further details of all proofs are provided in the on-line supplement materials at <http://www.stat.sinica.edu.tw/statistica>.

## 4. Algorithms

### 4.1 SCAD

Despite the excellent statistical properties of the SCAD penalized estimator, the corresponding optimization is a non-convex minimization problem and is much harder to solve than the LASSO penalized counterpart. In Fan and Li (2001), a unified least quadratic approximation (LQA) algorithm was proposed to solve the SCAD penalized likelihood optimization problem. Hunter and Li (2005) studied LQA under a more general MM-algorithm framework, where MM stands for minorize-maximize or majorize-minimize. A typical example of the MM algorithm is the well known EM algorithm.

Notice that in (2.3), the first order derivative of the SCAD penalty function on  $(0, \infty)$  is the sum of two components: the first one is a constant and the second one is a decreasing function on the range  $(0, \infty)$ . As a result, the SCAD penalty function can be decomposed as the difference of two convex functions. More explicitly, we have  $p_\lambda(\theta) = p_{\lambda,1}(\theta) - p_{\lambda,2}(\theta)$  where both  $p_{\lambda,1}(\cdot)$  and  $p_{\lambda,2}(\cdot)$  are convex functions and their derivatives for  $\theta > 0$  are given by

$$\begin{cases} p'_{\lambda,1}(\theta) &= \lambda \\ p'_{\lambda,2}(\theta) &= \lambda(1 - \frac{(a\lambda - \theta)_+}{(a-1)\lambda})I(\theta > \lambda). \end{cases} \quad (4.1)$$

For a particular set of parameters  $a = 3.7$  and  $\lambda = 2$ , this decomposition is graphically illustrated in Figure 4.1, where the left panel plots  $p_{\lambda,1}(\theta)$ , the central panel corresponds to  $p_{\lambda,2}(\theta)$ , and  $p_\lambda(\theta) = p_{\lambda,1}(\theta) - p_{\lambda,2}(\theta)$  are given in the right panel. The above decomposition of the SCAD penalty allows us to use the well studied DC algorithm. DCA was proposed by An and Tao (1997) to handle non-convex optimization. Later on, it was applied in machine learning (Liu, Shen and Doss, 2005b; Wu and Liu, 2007). DCA is a local algorithm and it decreases the objective value at each iteration. Due to its decomposition and approximation, DCA converges in finite steps. More details on DCA can be found in (Liu et al.,

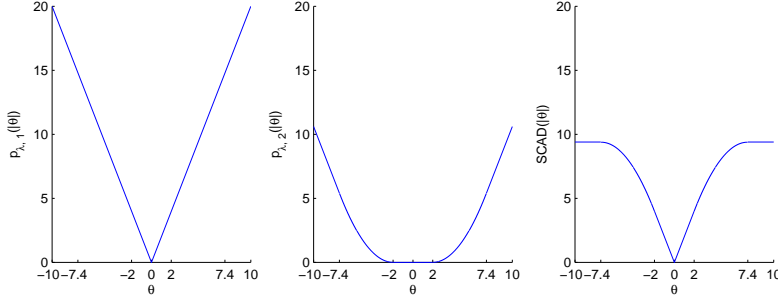


Figure 4.1: Decomposition of the SCAD penalty as  $p_\lambda(\theta) = p_{\lambda,1}(\theta) - p_{\lambda,2}(\theta)$ , with parameters  $\lambda = 2$  and  $a = 3.7$

2005b; Liu, Shen and Wong, 2005a).

Due to the above decomposition of the SCAD penalty, the objective function of the SCAD penalized quantile regression can be decomposed as  $Q_{vex}(\boldsymbol{\beta}) + Q_{cav}(\boldsymbol{\beta})$ , where  $Q_{vex}(\boldsymbol{\beta}) = \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) + n \sum_{j=1}^d p_{\lambda_n,1}(|\beta_j|)$  and  $Q_{cav}(\boldsymbol{\beta}) = -n \sum_{j=1}^d p_{\lambda_n,2}(|\beta_j|)$ .

**Algorithm 1:** Difference Convex Algorithm for minimizing  $Q(\boldsymbol{\beta}) = Q_{vex}(\boldsymbol{\beta}) + Q_{cav}(\boldsymbol{\beta})$

1. Initialize  $\boldsymbol{\beta}^{(0)}$ .
2. Repeat  $\boldsymbol{\beta}^{(t+1)} = \underset{\boldsymbol{\beta}}{\operatorname{argmin}}(Q_{vex}(\boldsymbol{\beta}) + \langle Q'_{cav}(\boldsymbol{\beta}^{(t)}), \boldsymbol{\beta} - \boldsymbol{\beta}^{(t)} \rangle)$  until convergence.

The difference convex algorithm solves the non-convex minimization problem via a sequence of convex subproblems (see Algorithm 1). Denote the solution at step  $t$  by  $\boldsymbol{\beta}^{(t)} = (\beta_1^{(t)}, \beta_2^{(t)}, \dots, \beta_p^{(t)})^T$ . Hence the derivative of the concave part at  $\boldsymbol{\beta}^{(t)}$  is given by

$$Q'_{cav}(\boldsymbol{\beta}^{(t)}) = -n(p'_{\lambda_n,2}(|\beta_1^{(t)}|) \operatorname{sign}(\beta_1^{(t)}), p'_{\lambda_n,2}(|\beta_p^{(t)}|) \operatorname{sign}(\beta_p^{(t)}), \dots, p'_{\lambda_n,2}(|\beta_p^{(t)}|) \operatorname{sign}(\beta_p^{(t)}))^T,$$

where  $p'_{\lambda_n,2}(\cdot)$  is defined in (4.1) and  $\operatorname{sign}(\cdot)$  is the sign function. In the  $(t+1)$ -th iteration, DCA approximates the second function by a linear function and solves the following optimization problem:

$$\min_{\boldsymbol{\beta}} \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}) + n \sum_{j=1}^d p_{\lambda_n,1}(|\beta_j|) - n \sum_{j=1}^d p'_{\lambda_n,2}(|\beta_j^{(t)}|) \operatorname{sign}(\beta_j^{(t)}) (\beta_j - \beta_j^{(t)}) \quad (4.2)$$

Here for the initialization in Step 1 of Algorithm 1, we use the solution of non-penalized linear quantile regression, namely,  $\beta^{(0)} = \tilde{\beta}_\tau$  where  $\tilde{\beta}_\tau$  is given by (2.4).

By introducing some slack variables, we can recast the above minimization problem (4.2) into the following linear programming problem:

$$\begin{aligned} \min \quad & \sum_{i=1}^n (\tau \xi_i + (1 - \tau) \zeta_i) + n \lambda_n \sum_{j=1}^d \nu_j - n \sum_{j=1}^d p'_{\lambda_n, 2}(|\beta_j^{(t)}|) \text{sign}(\beta_j^{(t)}) (\beta_j - \beta_j^{(t)}) \\ \text{subject to} \quad & \xi_i \geq 0, \zeta_i \geq 0, \xi_i - \zeta_i = y_i - \mathbf{x}_i^T \boldsymbol{\beta}, \quad i = 1, 2, \dots, n \\ & \nu_j \geq \beta_j, \nu_j \geq -\beta_j, \quad j = 1, 2, \dots, d, \end{aligned}$$

which can be easily solved by many optimization softwares. In contrast, at each iteration, the LQA (Fan and Li, 2001; Hunter and Li, 2005) needs to solve a quadratic programming problem and as a result they are less efficient. Our numerical studies in Section 5 show that the DCA is indeed more efficient.

## 4.2 Adaptive-LASSO

With the aid of slack variables, the adaptive-LASSO penalized quantile regression can also be casted into a linear programming problem as follows,

$$\begin{aligned} \min \quad & \sum_{i=1}^n (\tau \xi_i + (1 - \tau) \zeta_i) + n \lambda_n \sum_{j=1}^d \tilde{w}_j \eta_j \\ \text{subject to} \quad & \xi_i \geq 0, \zeta_i \geq 0, \xi_i - \zeta_i = y_i - \mathbf{x}_i^T \boldsymbol{\beta}, \quad i = 1, 2, \dots, n \quad (4.3) \\ & \eta_j \geq \beta_j, \eta_j \geq -\beta_j, \quad j = 1, 2, \dots, d. \end{aligned}$$

Here the weights  $\tilde{w}_j$ 's are appropriately chosen as discussed in Section 2.2. Note that the minimization problem (4.3) includes the LASSO penalized quantile regression as a special case by setting  $\tilde{w}_j = 1$  for  $j = 1, 2, \dots, d$ .

## 5. Monte Carlo study

In this section, we first use one example to compare three different algorithms (LQA, MM, and DCA) for the SCAD penalized quantile regression and show the advantage of our new DC algorithm for the SCAD. Hence, we choose the DCA for the SCAD in the remaining numerical studies. In the remaining examples, we study the finite-sample variable selection performance of different penalized quantile regressions. Here we want to point out that the intercept

term is included in penalized quantile regression for all data analysis in this paper. For the SCAD penalty, we do not tune the parameter  $a$ . Following Fan and Li (2001)'s suggestion, we set  $a = 3.7$  to reduce the computation burden. The number of zero coefficients is evaluated as follows: an estimate is treated as zero if its absolute value is smaller than  $10^{-6}$ .

The data for Examples 5.1 and 5.2 are generated from the following linear model

$$y = \mathbf{x}^T \boldsymbol{\beta} + \sigma \epsilon, \quad (5.1)$$

where  $\boldsymbol{\beta} = (3, 1.5, 0, 0, 2, 0, 0, 0)^T$ . The components of  $\mathbf{x}$  and  $\epsilon$  are standard normal. The correlation between any two components  $x_i$  and  $x_j$  is set to be  $\rho^{|i-j|}$  with  $\rho = 0.5$ . This model has been considered by many authors (Tibshirani, 1996; Fan and Li, 2001; Zou, 2006, to name a few).

Denote the sample size of training data sets by  $n$ . Throughout this section, an independent tuning data set and testing data set of size  $n$  and  $100n$  respectively are generated exactly in the same way as the training data set. The tuning parameter  $\lambda$  is selected via a grid search based on the tuning error in terms of the check loss function evaluated on the tuning data. Similarly defined testing errors on the testing data set will be reported. More explicitly, a test error refers to the average check loss on the independent testing data set.

### **Example 5.1 Comparison of LQA, MM, and DCA for the SCAD**

In this example, we generate data from model (5.1) with size  $n = 60$  and different algorithms for the SCAD penalized quantile regression are compared. Table 5.1 summarizes the results of 100 repetitions for two cases:  $\sigma = 1$  and  $\sigma = 3$ . Average test errors, numbers of correct and wrong zero coefficients, and CPU times with standard deviations in their corresponding parentheses are reported. It suggests that, while giving very similar test errors, these three algorithms produce different numbers of zero coefficients. On average, DCA gives significantly more zeros. Remarkably, we notice that on average DCA takes much less CPU-time than LQA and MM as we expected. The reason is that in each iteration, DCA solves a linear programming while LQA and MM require a quadratic programming, as discussed at the end of Section 4.1. For the MM algorithm, we set Hunter and Li (2005)'s parameter  $\tau$  to be  $10^{-6}$  in their Equation (3.12).

Table 5.1: Simulation results for Example 5.1 with  $n = 60$ ,  $\tau = 0.5$

	Method	Test Error	no. of zeros		CPU-time in seconds
			Correct	Wrong	
$\sigma = 1$	LQA	0.4189 (0.0158)	2.76 (1.40)	0.00 (0.00)	14.44 ( 20.63)
	MM	0.4189 (0.0161)	4.08 (1.20)	0.00 (0.00)	26.29 (155.87)
	DCA	0.4193 (0.0163)	4.40 (1.02)	0.00 (0.00)	0.26 ( 0.15)
$\sigma = 3$	LQA	1.2856 (0.0692)	2.30 (1.48)	0.06 (0.24)	9.56 ( 13.06)
	MM	1.2807 (0.0647)	3.68 (1.63)	0.14 (0.38)	12.39 ( 32.98)
	DCA	1.2822 (0.0642)	3.84 (1.57)	0.15 (0.39)	0.21 ( 0.11)

Table 5.2: Simulation results for Example 5.2

$\tau$	$n = 100, \sigma = 1$				$n = 100, \sigma = 3$			
	Method	Test Error	no. of zeros		Test Error	no. of zeros		
			Correct	Wrong		Correct	Wrong	
.25	$L_1$	0.3378 (0.0111)	3.16 (1.47)	0.00 (0.00)	0.9976 (0.0347)	1.87 (1.49)	0.00 (0.00)	
	SCAD	0.3296 (0.0091)	3.98 (1.66)	0.00 (0.00)	0.9968 (0.0364)	3.94 (1.60)	0.01 (0.10)	
	adapt- $L_1$	0.3288 (0.0090)	4.21 (1.25)	0.00 (0.00)	0.9944 (0.0318)	3.00 (1.52)	0.00 (0.00)	
	Oracle	0.3282 (0.0087)	5.00 (0.00)	0.00 (0.00)	0.9873 (0.0284)	5.00 (0.00)	0.00 (0.00)	
.5	$L_1$	0.4143 (0.0108)	2.92 (1.48)	0.00 (0.00)	1.2379 (0.0392)	2.00 (1.48)	0.00 (0.00)	
	SCAD	0.4101 (0.0119)	4.00 (1.50)	0.00 (0.00)	1.2336 (0.0385)	4.01 (1.64)	0.02 (0.14)	
	adapt- $L_1$	0.4081 (0.0101)	4.31 (1.00)	0.00 (0.00)	1.2339 (0.0385)	3.20 (1.48)	0.01 (0.10)	
	Oracle	0.4072 (0.0099)	5.00 (0.00)	0.00 (0.00)	1.2248 (0.0348)	5.00 (0.00)	0.00 (0.00)	
.75	$L_1$	0.3364 (0.0093)	3.24 (1.42)	0.00 (0.00)	0.9893 (0.0324)	2.09 (1.36)	0.00 (0.00)	
	SCAD	0.3286 (0.0100)	4.05 (1.42)	0.00 (0.00)	0.9866 (0.0336)	4.30 (1.35)	0.06 (0.24)	
	adapt- $L_1$	0.3307 (0.0083)	4.51 (1.05)	0.00 (0.00)	0.9827 (0.0325)	3.73 (1.30)	0.01 (0.10)	
	Oracle	0.3266 (0.0091)	5.00 (0.00)	0.00 (0.00)	0.9747 (0.0241)	5.00 (0.00)	0.00 (0.00)	

Because of its superior performance, the DCA algorithm is used for the remaining data analysis to solve the SCAD penalized quantile regression.

**Example 5.2 Comparison of finite-sample variable selection performance**

We generate data from model (5.1) to compare the finite-sample variable selection performance of the  $L_1$ , the SCAD, and the adaptive- $L_1$  with the oracle. Simulation results of different settings are reported in Table 5.2. We can see that the reported test errors are very similar. But on average, the SCAD and the adaptive- $L_1$  give more zero coefficients than the  $L_1$ . This confirms the superiority of the SCAD and the adaptive- $L_1$  as shown in our theoretical results.

**Example 5.3 Dimensionality larger than the sample size**

For the adaptive LASSO penalty, an initial consistent estimator is required to derive the adaptive weights. Due to the work by He and Shao (2000), the solution of the linear quantile regression as defined by (2.4) is still consistent even for the case that the dimensionality increases as the sample size but at a speed slower than some root of the sample size. However, it is not clear how to find an consistent initial solution for deriving the adaptive weights in the case of

Table 5.3: Simulation results for Example 5.3 with sample size  $n = 100$ . Here the superscript  $*$  in  $\text{adapt-}L_1^*$  indicates that the adaptive weights of the adaptive- $L_1$  penalty are based on the solution of the  $L_2$  penalized quantile regression.

$\tau$	Method	Test Error	no. of zeros	
			Correct	Wrong
.25	$L_1$	0.1744 (0.0073)	113.65 (4.30)	0.00 (0.00)
	SCAD-DCA	0.1673 (0.0049)	116.72 (1.09)	0.00 (0.00)
	$\text{adapt-}L_1^*$	0.1684 (0.0045)	115.47 (3.16)	0.00 (0.00)
	Oracle	0.1668 (0.0038)	117.00 (0.00)	0.00 (0.00)
.5	$L_1$	0.2150 (0.0073)	112.37 (5.67)	0.00 (0.00)
	SCAD-DCA	0.2094 (0.0043)	116.38 (1.33)	0.00 (0.00)
	$\text{adapt-}L_1^*$	0.2101 (0.0057)	114.54 (4.58)	0.00 (0.00)
	Oracle	0.2089 (0.0038)	117.00 (0.00)	0.00 (0.00)
.75	$L_1$	0.1749 (0.0068)	112.47 (6.50)	0.00 (0.00)
	SCAD-DCA	0.1692 (0.0063)	116.59 (1.18)	0.00 (0.00)
	$\text{adapt-}L_1^*$	0.1705 (0.0055)	115.02 (4.79)	0.00 (0.00)
	Oracle	0.1680 (0.0048)	117.00 (0.00)	0.00 (0.00)

dimensionality larger than the sample size. For the aforementioned case with  $p > n$ , we propose to perform the  $L_2$  penalized quantile regression first and use the obtained solution to derive the adaptive weight for the adaptive- $L_1$  penalty. In this example, we compare the performance of these different penalties in the case with more predictor variables than the sample size.

Our datasets in this example are generated from model (5.1), augmented with 102 more independent noise variables  $x_9, x_{10}, \dots, x_{110}$ . Adding more independent noise variables makes the estimation much harder. In order to make the estimation possible, the variance of random error  $\epsilon$  is set to be  $\sigma^2 = 0.5^2$ ; each of these additional noise variable is identically distributed as  $N(0, 0.5^2)$  and independent of each other. The results based on 100 repetitions with sample size 100 are reported in Table 5.3. It is evident from Table 5.3 that both the SCAD and adaptive- $L_1$  penalties improve over the  $L_1$  penalty in terms of prediction accuracy as well as variable selection capability, even in the more difficult case of  $p > n$ . This results also validate our proposed procedure of using the  $L_2$  penalized solution to derive the adaptive weights for the adaptive- $L_1$  penalty.

#### Example 5.4 Non *i.i.d* random errors

In this example, we consider the case with non *i.i.d.* random errors to check the robustness of our methods. Our data is generated from model 2 of Kocherginsky, He and Mu (2005), defined as follows

$$y = 1 + x_1 + x_2 + x_3 + (1 + x_3)\epsilon, \quad (5.2)$$

Table 5.4: Simulation results for Example 5.4 with sample size  $n = 100$ .

$\tau$	Method	Test Error	no. of zeros	
			Correct	Wrong
.25	$L_1$	0.4944 (0.0110)	2.27 (1.72)	0.72 (0.45)
	SCAD-DCA	0.4919 (0.0136)	4.73 (0.62)	0.51 (0.50)
	adapt- $L_1$	0.4926 (0.0119)	3.64 (1.27)	0.65 (0.48)
	Oracle	0.4925 (0.0133)	5.00 (0.00)	0.00 (0.00)
.5	$L_1$	0.6272 (0.0145)	1.37 (1.57)	0.36 (0.48)
	SCAD-DCA	0.6157 (0.0118)	4.52 (0.69)	0.19 (0.42)
	adapt- $L_1$	0.6196 (0.0107)	3.31 (1.32)	0.29 (0.46)
	Oracle	0.6157 (0.0117)	5.00 (0.00)	0.00 (0.00)
.75	$L_1$	0.5081 (0.0146)	1.44 (1.59)	0.25 (0.44)
	SCAD-DCA	0.4942 (0.0140)	4.72 (0.55)	0.06 (0.24)
	adapt- $L_1$	0.5008 (0.0147)	3.59 (1.08)	0.16 (0.37)
	Oracle	0.4935 (0.0132)	5.00 (0.00)	0.00 (0.00)

where  $x_1$  and  $x_3$  are generated from the standard normal distribution and the uniform distribution on  $[0, 1]$ ,  $x_2 = x_1 + x_3 + z$  with  $z$  being standard normal, and  $\epsilon \sim N(0, 1)$ . The variables  $x_1$ ,  $x_3$ ,  $z$ , and  $\epsilon$  are mutually independent. To study the effect of variable selection, we include five more independent standard normal noise variables,  $x_4, x_5, \dots, x_8$ , which are independent of each other.

The results based on 100 repetitions with sample size  $n = 100$  are reported in Table 5.4, in the same format as in Example 5.2. Similarly, we can see the improvement in terms of test errors. Moreover, both the SCAD and adaptive- $L_1$  penalties can identify more correction zero coefficients than the  $L_1$  does. In this case, all these three penalties tend to produce more wrong zero coefficients in the final model comparing to Example 5.2. The potential reason is that  $x_2$  is highly correlated with  $x_1$  and  $x_3$  with  $x_2 = x_1 + x_3 + z$ . Nevertheless, comparing the  $L_1$  penalty, the SCAD and adaptive- $L_1$  penalties on average lead to less wrong zero coefficients.

## 6. Real data

In Harrison and Rubinfeld (1978), they studied various methodological issues related to the use of housing data to estimate the demand for clean air. In particular, the Boston House Price Dataset was used. This dataset is available online at [http://lib.stat.cmu.edu/datasets/boston\\_corrected.txt](http://lib.stat.cmu.edu/datasets/boston_corrected.txt), with some corrections and augmentation with the latitude and longitude of each observation which is called Corrected Boston House Price Data. There are 506 observations, 15 non-constant predictor variables, and one response variable corrected median value of owner-occupied homes (CMEDV). They include longitude (LON),

Table 6.5: Results of the Corrected Boston House Price Data

Method	$\tau = 0.25$		$\tau = 0.5$		$\tau = 0.75$	
	Test Error	no. of zeros	Test Error	no. of zeros	Test Error	no. of zeros
$L_1$	0.1339 (0.0107)	11.10 (3.14)	0.1832 (0.0215)	9.30 (4.16)	0.1813 (0.0419)	7.10 (4.72)
SCAD	0.1367 (0.0164)	14.20 (2.78)	0.1862 (0.0257)	12.40 (4.40)	0.1920 (0.0799)	12.40 (3.86)
adapt- $L_1$	0.1346 (0.0130)	13.60 (3.20)	0.1840 (0.0216)	11.10 (5.67)	0.1776 (0.0403)	12.10 (3.98)

Note: In this table, the DCA is chosen for the SCAD.

latitude (LAT), crime rate (CRIM), proportion of area zoned with large lots (ZN), proportion of non-retail business acres per town (INDUS), Charles River as a dummy variable (= 1 if tract bounds river; 0 otherwise) (CHAS), nitric oxides concentration (NOX), average number of rooms per dwelling (RM), proportion of owner-occupied units built prior to 1940 (AGE), weighted distances to five Boston employment centres (DIS), index of accessibility to radial highways (RAD), property tax rate (TAX), pupil-teacher ratio by town (PTRATIO), black population proportion town (B), and lower status population proportion (LSTAT). For simplicity, we exclude the categorical variable RAD. We also standardize the response variable CMEDV and 13 of the remaining predictor variables except the indicator variable CHAS. Penalized quantile regression is applied with the standardized CMEDV as the response. We use 27 predictor variables in the penalized quantile regression, including the variable CHAS, these 13 standardized predictor variables and their squares.

In each repetition, we randomly split all the 506 observations into training, tuning and testing data sets of size 150, 150, and 206 respectively. The performance over 10 repetitions of the penalized quantile regression with different penalties and different quantiles is summarized in Table 6.5. The results indicate that different penalties give similar test errors but the SCAD and adaptive- $L_1$  use less variables than the  $L_1$  does.

## 7. Discussion

In this work, we study penalized quantile regression with the SCAD and the adaptive-LASSO penalties. We show that they enjoy the oracle properties as established by Fan and Li (2001) and Zou (2006) even though the check function is non-differentiable at the origin. To handle the non-convex optimization problem of the SCAD penalized quantile regression, we propose to use the Difference Convex algorithm. The new algorithm is very efficient as confirmed by the simulation results given in Example 5.1.



Notice that DCA is a very general algorithm. It can be easily extended to apply to a more general SCAD penalized likelihood setting as long as the likelihood part is convex. For example, in SCAD penalized least squares regression, each iteration involves a quadratic programming problem. Similarly, DCA can be applied to the SCAD SVM (Zhang, Ahn, Lin and Park, 2006).

### Acknowledgment

We would like to thank Professor Jianqing Fan for his helpful comments and suggestions. This research is partially supported by NSF grant DMS-06-06577 and NIH grant R01-GM07261.

### References

- AN, L. T. H. and TAO, P. D. (1997). Solving a class of linearly constrained indefinite quadratic problems by d.c. algorithms. *Journal of Global Optimization*, **11** 253–285.
- BREIMAN, L. (1996). Heuristics of instability and stabilization in model selection. *The Annals of Statistics*, **24** 2350–2383.
- CANDES, E. and TAO, T. (2007). The dantzig selector: statistical estimation when  $p$  is much larger than  $n$ . *The Annals of Statistics*. To appear.
- FAN, J. (1997). Comments on “wavelets in statistics: A review”, by a. antoniadis. *Journal of Italian Statistical Society*, **6** 131–138.
- FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, **96** 1348–1360.
- FAN, J. and LI, R. (2002). Variable selection for coxs proportional hazards model and frailty model. *The Annals of Statistics*, **30** 74–99.
- FAN, J. and LI, R. (2004). New estimation and model selection procedures for semiparametric modeling in longitudinal data analysis. *Journal of the American Statistical Association*, **99** 710–723.
- FAN, J. and LV, J. (2006). Sure independence screening for ultra-high dimensional feature space. Submitted.

- FAN, J. and PENG, H. (2004). Nonconcave penalized likelihood with a diverging number of parameters. *The Annals of Statistics*, **32** 928–961.
- FRANK, I. and FRIEDMAN, J. (1993). A statistical view of some chemometrics regression tools. *Technometrics*, **35** 109–148.
- GEYER, C. J. (1994). On the asymptotics of constrained m-estimation. *The Annals of Statistics*, **22** 1993–2010.
- HARRISON, D. and RUBINFELD, D. L. (1978). Hedonic housing prices and the demand for clean air. *Journal of Environmental Economics and Management* 81–102.
- HE, X. and SHAO, Q.-M. (2000). On parameters of increasing dimensions. *Journal of Multivariate Analysis*, **73** 120–135.
- HENDRICKS, W. and KOENKER, R. (1992). Hierarchical spline models for conditional quantiles and the demand for electricity. *Journal of the American Statistical Association*, **87** 58–68.
- HOERL, A. and KENNARD, R. (1988). Ridge regression. In *Encyclopedia of Statistical Sciences*, vol. 8. Wiley, New York, 129–136.
- HUNTER, D. R. and LI, R. (2005). Variable selection using mm algorithm. *The Annals of Statistics*, **33** 1617–1642.
- KNIGHT, K. (1999). Asymptotics for  $L_1$ -estimators of regression parameters under heteroscedasticity. *The Canadian Journal of Statistics*, **27** 497–507.
- KOCHERGINSKY, M., HE, X. and MU, Y. (2005). Practical confidence intervals for regression quantiles. *Journal of Computational and Graphical Statistics*, **14** 41–55.
- KOENKER, R. (2004). Quantile regression for longitudinal data. *Journal of Multivariate Analysis*, **2004** 74–89.
- KOENKER, R. (2005). *Quantile regression*, Cambridge University Press.
- KOENKER, R. and BASSETT, G. (1978). Regression quantiles. *Econometrica*, **46** 33–50.

- KOENKER, R. and GELING, R. (2001). Reappraising medfly longevity: a quantile regression survival analysis. *Journal of the American Statistical Association*, **96** 458–468.
- KOENKER, R. and HALLOCK, K. (2001). Quantile regression. *Journal of Economic Perspectives*, **15** 143–156.
- KOENKER, R., NG, P. and PORTNOY, S. (1994). Quantile smoothing splines. *Biometrika*, **81** 673–680.
- LI, Y., LIU, Y. and ZHU, J. (2007). Quantile regression in reproducing kernel hilbert spaces. *Journal of the American Statistical Association*, **102** 255–268.
- LI, Y. and ZHU, J. (2005).  $l_1$ -norm quantile regressions. *Journal of Computational and Graphical Statistics*, To appear.
- LIU, S., SHEN, X. and WONG, W. (2005a). Computational development of  $\psi$ -learning. In *The SIAM 2005 International Data Mining Conf.* 1–12.
- LIU, Y., SHEN, X. and DOSS, H. (2005b). Multicategory  $\psi$ -learning and support vector machine: computational tools. *Journal of Computational and Graphical Statistics*, **14** 219–236.
- POLLARD, D. (1991). Asymptotics for least absolute deviation regression estimators. *Econometric Theory*, **7** 186–199.
- TIBSHIRANI, R. J. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, Series B*, **58** 267–288.
- WANG, H. and HE, X. (2007). Detecting differential expressions in genechip microarray studies: A quantile approach. *Journal of the American Statistical Association* 104–112.
- WANG, H., LI, G. and JIANG, G. (2007). Robust regression shrinkage and consistent variable selection through the lad-lasso. *Journal of Business & Economic Statistics*, **25** 347–355.
- WEI, Y. and HE, X. (2006). Conditional growth charts (with discussions). *Annals of Statistics*, **34** 2069–2031.

- WEI, Y., PERE, A., KOENKER, R. and HE, X. (2006). Quantile regression methods for reference growth curves. *Statistics in Medicine*, **25** 1369–1382.
- WU, Y. and LIU, Y. (2006). Robust truncated-hinge-loss support vector machines. *Journal of the American Statistical Association*. **102** 974–983.
- YANG, S. (1999). Censored median regression using weighted empirical survival and hazard functions. *Journal of the American Statistical Association*, **94** 137–145.
- YUAN, M. and LIN, Y. (2007). On the nonnegative garrote estimator. *Journal of the Royal Statistical Society, Series B*. **69** 143–161.
- ZHANG, H. H., AHN, J., LIN, X. and PARK, C. (2006). Gene selection using support vector machines with nonconvex penalty. *Bioinformatics*, **22** 88–95.
- ZHANG, H. H. and LU, W. (2007). Adaptive-lasso for cox’s proportional hazard model. *Biometrika*. **94** 691–703.
- ZHAO, P. and YU, B. (2006). On model selection consistency of lasso. *Journal of Machine Learning Research*, **7** 2541–2563.
- ZOU, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, **101** 1418–1429.
- ZOU, H. and LI, R. (2007). One-step Sparse Estimates in Nonconcave Penalized Likelihood Models. *The Annals of Statistics*, To appear.

Department of Operations Research and Financial Engineering, Princeton University, Princeton NJ 08544, U.S.A.

E-mail: yichaowu@princeton.edu

Department of Statistics and Operations Research, The University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, U.S.A.

E-mail: yfliu@email.unc.edu

## On-line Supplement

**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  $\epsilon_i$  has the  $\tau$ -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  $W_{n,11} = \sum_{i=1}^n D_i \mathbf{x}_i \mathbf{x}_i^T / \sqrt{n}$ . Then  $W_n \xrightarrow{\mathcal{L}} N(\mathbf{0}, \tau(1 - \tau)\Sigma)$  and  $W_{n,11} \xrightarrow{\mathcal{L}} N(\mathbf{0}, \tau(1 - \tau)\Sigma_{11})$ .

**Lemma 3.** *For model (3.1) with true parameter  $\beta_0$ , 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  $\mathbf{u}$ ,*

$$G_n(\mathbf{u}) = \frac{f(0)}{2} \mathbf{u}^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \mathbf{u} + W_n^T \mathbf{u} + o_p(1). \quad (7.1)$$

*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  $|R_{i,n}(\mathbf{u})| \leq |\mathbf{x}_i^T \mathbf{u} / \sqrt{n}| \{ |\epsilon_i| \leq |\mathbf{x}_i^T \mathbf{u} / \sqrt{n}| \}$ . For fixed  $\mathbf{u}$ , due to the cancelation of cross-product terms, we get

$$\begin{aligned}
E\left(\sum_{i=1}^n [R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u})]^2\right) &= \sum_{i=1}^n E(R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u}))^2 \\
&\leq \sum_{i=1}^n E(R_{i,n}(\mathbf{u}))^2 \\
&\leq \sum_{i=1}^n \left[ \left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right|^2 E\{ |\epsilon_i| \leq \left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right| \} \right] \\
&\leq \left( \sum_{i=1}^n \left| \frac{\mathbf{x}_i^T \mathbf{u}}{\sqrt{n}} \right|^2 \right) E\{ |\epsilon_i| \leq \frac{\|\mathbf{u}\|}{\sqrt{n}} \max_{j=1,2,\dots,n} \|\mathbf{x}_j\| \} \\
&\rightarrow 0
\end{aligned} \tag{7.2}$$

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 (7.2) implies that  $\sum_{i=1}^n (R_{i,n}(\mathbf{u}) - ER_{i,n}(\mathbf{u})) = o_p(1)$ . Hence 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

$$P \left\{ \inf_{\|\mathbf{u}\|=C} Q(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n}) > Q(\boldsymbol{\beta}_0) \right\} \geq 1 - \delta \quad (7.3)$$

which implies that with probability at least  $1 - \delta$  there exists a local minimum in the ball  $\{\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n} : \|\mathbf{u}\| \leq C\}$ . 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  $\beta_{j0}$  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  $n$ ,

$$n \sum_{j=1}^s [p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|)] = 0 \quad (7.4)$$

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 (7.3) holds as we have desired and this completes the proof.  $\square$

*Proof of Lemma 1.* For any  $\beta_1 - \beta_{10} = O_p(n^{-1/2})$ ,  $0 < \|\beta_2\| \leq Cn^{-1/2}$ ,

$$\begin{aligned}
& Q((\beta_1^T, \mathbf{0}^T)^T) - Q((\beta_1^T, \beta_2^T)^T) \\
&= [Q((\beta_1^T, \mathbf{0}^T)^T) - Q((\beta_{10}^T, \mathbf{0}^T)^T)] - [Q((\beta_1^T, \beta_2^T)^T) - Q((\beta_{10}^T, \mathbf{0}^T)^T)] \\
&= G_n(\sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T)^T) - G_n(\sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T)^T) - n \sum_{j=s+1}^d p_{\lambda_n}(|\beta_j|) \quad (7.5) \\
&= \frac{f(0)}{2} \sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T)^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T)^T + \sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T) W_n \\
&\quad - \frac{f(0)}{2} \sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T)^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T)^T - \sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T) W_n \\
&\quad + o(1) + o_p(1) - n \sum_{j=s+1}^d p_{\lambda_n}(|\beta_j|)
\end{aligned}$$

The conditions  $\beta_1 - \beta_{10} = O_p(n^{-1/2})$  and  $0 < \|\beta_2\| \leq Cn^{-1/2}$  imply that

$$\frac{f(0)}{2} \sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T)^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T)^T = O_p(1)$$

$$\frac{f(0)}{2} \sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T)^T \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T)^T = O_p(1)$$

and

$$\begin{aligned}
& \sqrt{n}((\beta_1 - \beta_{10})^T, \mathbf{0}^T) W_n - \sqrt{n}((\beta_1 - \beta_{10})^T, \beta_2^T) W_n = -\sqrt{n}(\mathbf{0}^T, \beta_2^T) W_n \\
&= \sqrt{n} \sqrt{\tau(1-\tau) \beta_2^T \Sigma_{22} \beta_2} (1 + o_p(1)).
\end{aligned}$$

Note that

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

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 (7.5), the last term dominates in magnitude and, as a result,  $Q((\beta_1^T, \mathbf{0}^T)^T) - Q((\beta_1^T, \beta_2^T)^T) < 0$  for large  $n$ . This 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  $\hat{\beta}_1$  of  $Q((\beta_1^T, \mathbf{0}^T)^T)$  as a function of  $\beta_1$ .

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

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

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 (7.4). Hence we have

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

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  $r_n(\theta) \rightarrow 0$  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  $\theta$ . So this implies that, for large  $n$ , the local minimizer  $\hat{\theta}$  is very close to  $\zeta_n$  and satisfies  $\hat{\theta} - \zeta_n = o_p(1)$ .

That is, the minimizer  $\hat{\theta}$  satisfies  $\hat{\theta} = -(f(0) \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n})^{-1} (\sum_{i=1}^n D_i \mathbf{x}_{i1} / \sqrt{n}) +$

$o_p(1)$ . 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 (7.5) 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 (7.1).

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  $\zeta(\cdot)$ .  $\square$