

# Penalized expectile regression: an alternative to penalized quantile regression

Lina Liao<sup>1</sup> · Cheolwoo Park<sup>1</sup> · Hosik Choi<sup>2</sup>

Received: 18 March 2017 / Revised: 1 January 2018 / Published online: 19 February 2018  
© The Institute of Statistical Mathematics, Tokyo 2018

**Abstract** This paper concerns the study of the entire conditional distribution of a response given predictors in a heterogeneous regression setting. A common approach to address heterogeneous data is quantile regression, which utilizes the minimization of the  $L_1$  norm. As an alternative to quantile regression, we consider expectile regression, which relies on the minimization of the asymmetric  $L_2$  norm and detects heteroscedasticity effectively. We assume that only a small set of predictors is relevant to the response and develop penalized expectile regression with SCAD and adaptive LASSO penalties. With properly chosen tuning parameters, we show that the proposed estimators display oracle properties. A numerical study using simulated and real examples demonstrates the competitive performance of the proposed penalized expectile regression, and its combined use with penalized quantile regression would be helpful and recommended for practitioners.

**Keywords** Asymptotics · Expectile regression · Heteroscedasticity · Penalized regression · Variable selection

---

✉ Hosik Choi  
choi.hosik@gmail.com

Lina Liao  
linaliaostat@gmail.com

Cheolwoo Park  
cpark@uga.edu

<sup>1</sup> Department of Statistics, University of Georgia, Athens, GA 30602, USA

<sup>2</sup> Department of Applied Statistics, Kyonggi University, Suwon, Kyonggi-do 16227, Korea

## 1 Introduction

Quantiles and expectiles, which contain information about the full distribution for a random variable, are extensions of median and mean, respectively. Quantiles are percentiles of the cumulative distribution function of a random variable. For instance, if  $x$  is the  $\alpha$ th quantile of  $X$ ,  $P(X \leq x) = \alpha$ . Unlike quantiles, which is the minimal value of the tail, expectiles incorporate information about the expectation of  $X$ , conditional on  $X$  being in a tail of its distribution (Newey and Powell 1987).

In financial time series, expectiles emerge as an alternative to popular risk measures such as value at risk (VaR) and expected shortfall (ES), as they have desirable properties (Ziegel 2014). A risk measure is an estimated amount of capital to be reserved at a given risk level to prevent substantial losses. VaR at  $\alpha$ , the  $100\alpha$ th quantile of the return distribution, can be interpreted as the minimum potential loss at the  $100\alpha\%$  level. While VaR is the most widely used risk measure, it does not provide information regarding the potential magnitude of losses because VaR only depends on the tail probability. ES is an alternative risk measure that considers the magnitude of the potential losses in the lower tail. However, it is also known that ES can be too conservative, which could be a major disadvantage to commercial and investment banks. In contrast to VaR and ES, which only concern the lower tail, the expectile relies on both tails of the distribution to measure risk. The squared loss function makes the expectile more sensitive to extreme values and to the shape of the distribution than VaR, which can be beneficial when measuring potential losses because one wants a risk measure to be sensitive to extreme tail losses.

When dealing with heterogeneous data in regression, we frequently see that targeting only a mean function is not sufficient to capture a complete picture of the relationship between the response variable and predictors. In such a case, quantile regression (Koenker and Bassett 1978) based on an asymmetric  $L_1$  norm could be a more appropriate tool as it allows one to study the quantile structure of the conditional distribution. Quantile regression has been applied in various fields such as economics, survival analysis, and microarray studies. While quantile regression has a strong intuitive appeal, Newey and Powell (1987) point out three drawbacks: non-differentiability, inefficiency for Gaussian-like error distributions, and difficulty of calculation of a covariance matrix. They propose expectile regression based on an asymmetric  $L_2$  norm as an alternative way to analyze the complete conditional distribution of the response. Expectile regression generalizes ordinary mean regression, which is known to be efficient when typical assumptions, including homogeneity of errors, are met. It is also closely related to quantile regression, which is robust to outliers. Expectile regression has gained attention in several fields. Aigner et al. (1976) construct expectiles to estimate production frontiers. Sobotka et al. (2013a) investigate the relationship between women's education and fertility in Botswana via semiparametric expectile regression. Sobotka et al. (2013b) study statistical inference of semiparametric expectile regression. Schnabel and Eilers (2009) also demonstrate efficiency of expectiles over quantiles. Although expectile regression has found applications in various fields, to our knowledge, there has been little work with a penalized version of expectile regression. In this work, we fill this gap by investigating expectile regression with SCAD (Fan and Li 2001) and adaptive LASSO (Zou 2006) penalties.

A regularization approach has been extended to quantile regression because sets of important predictors might differ from quantile to quantile. Penalized quantile regression is capable of investigating the complete conditional distribution of the response variable and the sparsity pattern. [Li and Zhu \(2008\)](#) propose  $L_1$  norm quantile regression, which selects variables automatically and controls the variance of the fitted coefficients simultaneously. [Wu and Liu \(2009\)](#) demonstrate the oracle properties of quantile regression with SCAD and adaptive LASSO penalties when the number of variables,  $p$ , is fixed. [Wang et al. \(2012\)](#) study the behavior of quantile regression in ultra-high-dimensional data when  $p$  is much larger than  $n$ . [Belloni and Chernozhukov \(2011\)](#) study theory of quantile regression with lasso penalty in an ultra-high-dimensional setting, and [Belloni et al. \(2015\)](#) develop post-selection inference methods for quantile regression with lasso penalty in an ultra-high-dimensional setting.

The main objective of this paper is to develop penalized expectile regression when  $p$  is fixed as in [Wu and Liu \(2009\)](#). We assume that only a small number of predictors influence the conditional distribution of the response variable. We note that these sets of relevant predictors may vary for different segments of the conditional distribution. Therefore, consideration of different expectiles would enable us to explore the entire conditional distribution of the response variable and its sparsity pattern. The main contribution of the proposed work is threefold: (i) computation of expectile regression is straightforward and simple, (ii) theoretical development of expectile regression is more manageable with  $L_2$  norm and the estimation is more efficient because it uses the entire distribution information, and (iii) expectile regression is more sensitive to extreme values than quantile regression, which results in better detection of heteroscedasticity in the data when it is present. In the case of heteroscedasticity, penalized expectile regression yields superior performance in estimation and variable selection, as demonstrated in our simulation study. On the other hand, quantile regression is more interpretable and robust to outliers in a homogeneous setting than expectile regression. Therefore, we propose penalized expectile regression to complement penalized quantile regression; together they provide a more complete picture of the entire conditional distribution of a response given predictors for heterogeneous data.

Recently, [Gu and Zou \(2016\)](#) develop penalized expectile regression in a high-dimensional setting when  $p$  is larger than  $n$  and propose a way of detecting heteroscedasticity using a two-step procedure. While their theoretical work considers a more generalized setting than ours by letting  $p$  grow to infinity, the proposed work provides the following two additions to the regularized regression literature that is not covered in [Gu and Zou \(2016\)](#): (i) we discuss pros and cons of penalized quantile and expectile regression and conduct a thorough simulation study to compare the finite sample performance of both approaches, and (ii) we provide asymptotic distributions of penalized expectile regression with SCAD and adaptive LASSO penalties for both *i.i.d.* and non-*i.i.d.* random errors.

The remainder of the paper is organized as follows: In Sect. 2, we review expectiles and describe the relationship with quantiles. Section 3 describes the proposed penalized expectile regression using SCAD and adaptive LASSO penalties, our tuning selection method, and presents its theoretical properties. In Sect. 4, we compare penalized expectile and quantile regression via simulated examples. Section 5 analyzes a real

example and compares prediction and variable selection results with penalized quantile regression. We provide concluding remarks and future directions in Sect. 6. Appendix details the proofs of theorems presented in Sect. 3.

## 2 Expectile

The  $\tau$ th,  $\tau \in (0, 1)$ , expectile is defined as the value of  $m$  which minimizes

$$E[|\tau - I(Y < m)| (Y - m)^2]$$

for a random variable  $Y$ . Jones (1994) shows that expectiles are quantiles of a distribution  $G$  with an explicit relation to the original distribution  $F$ :

$$G(y) = \frac{P(y) - yF(y)}{2(P(y) - yF(y)) + (y - \mu)},$$

where  $P(y) = \int_{-\infty}^y xf(x)dx$  and  $\mu = \int_{-\infty}^{\infty} xf(x)dx$ . Table 1 (Kuan et al. 2009) contains the corresponding quantile  $\alpha$  to the expectile  $\tau = 0.01, 0.03, 0.05, 0.1, 0.25$ , respectively, for the Uniform( $-a, a$ ),  $N(0, 1)$ ,  $t(30)$ ,  $t(10)$ ,  $t(5)$  and  $t(3)$  distributions where  $t(d)$  indicates the distribution with  $d$  degrees of freedom, illustrating the relationship above. Note that for a given expectile, the corresponding quantile depends on the distribution.

Quantile regression and expectile regression, which generalize median regression and mean regression, respectively, are of great interest mainly for two reasons. First, data analysts often want to obtain a more complete relationship among the variables in a regression setting rather than simple mean or median information. Second, a mean regression function could be an incomplete summary for the relationship when the assumption of a common normal distribution fails to hold, or when a heterogeneous variation is present.

Suppose, from some unknown distribution, we draw a random sample  $\{(\mathbf{x}_i, y_i), i = 1, \dots, n\}$ , where  $\mathbf{x}_i$  and  $y_i$  denote the  $p$ -dimensional predictor and the response variable, respectively. Denote the vector of parameters as  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)^T$ . The quantile

**Table 1**  $\tau$  and corresponding  $\alpha$  for different distributions

$\tau$	$U(-a, a)$	$N(0, 1)$	$t(30)$	$t(10)$	$t(5)$
0.01	0.092	0.043	0.040	0.035	0.030
0.05	0.186	0.126	0.123	0.115	0.100
0.10	0.250	0.195	0.190	0.183	0.166
0.25	0.366	0.332	0.328	0.322	0.319
0.75	0.634	0.669	0.671	0.678	0.689
0.90	0.750	0.806	0.810	0.819	0.835
0.95	0.813	0.873	0.877	0.886	0.901
0.99	0.909	0.957	0.960	0.965	0.973

regression estimators, proposed by [Koenker and Bassett \(1978\)](#), are defined as those vectors which minimize the function

$$Q_n(\boldsymbol{\beta}; \alpha) = \sum_{i=1}^n r_\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta})$$

over  $\boldsymbol{\beta}$  in  $\mathbb{R}^p$  for fixed values of  $\alpha \in (0, 1)$ , where  $r_\alpha(\cdot)$  is a convex loss function of the form

$$r_\alpha(\epsilon) = |\alpha - I(\epsilon < 0)| \cdot |\epsilon|. \tag{1}$$

Similarly, the estimator of expectile regression can be obtained in the form of a vector that minimizes the following asymmetric least squares loss function:

$$R_n(\boldsymbol{\beta}; \tau) = \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\beta}),$$

over  $\boldsymbol{\beta}$  in  $\mathbb{R}^p$  for fixed values of  $\tau \in (0, 1)$ , where  $\rho_\tau(\cdot)$  is a convex loss function of the form

$$\rho_\tau(\epsilon) = |\tau - I(\epsilon < 0)| \cdot \epsilon^2. \tag{2}$$

We note that computation is no longer an issue for quantile regression due to recent advanced optimization techniques, but it is still true that computation in the expectile regression is more straightforward because it is based on the squared loss function in the form of (2). Thus, the optimization problem of expectile regression can be easily solved by many iteratively updated type algorithms, such as iteratively reweighted least squares (IRLS). Another attractive feature of expectile regression is that the expectile regression estimator depends on the shape of the entire distribution, while quantile regression estimator only relies on the percentiles of the estimated tail distribution. Hence, the expectile regression estimator contains additional information about the magnitude of the tail distribution and reflects the real value more accurately, especially for heavy-tailed distributions mentioned in Sect. 1. In terms of robustness, quantile regression is more resistant to outliers than expectile regression because quantile regression utilizes the  $L_1$  norm. However, sensitivity to extreme values can be beneficial if detecting heteroscedasticity in the data is of the main interest. We demonstrate the utility of expectile regression in detecting heteroscedasticity of the distribution in our simulation study in Sect. 4.

### 3 Penalized expectile regression

With large-sized data, selecting relevant variables and obtaining an interpretable model are important in regression analysis. To achieve these two goals for expectile regression, we apply regularization approaches using two different types of penalty functions.

We consider the same random sample in Sect. 2. Denote the true vector of parameters as  $\boldsymbol{\beta}_0 = (\beta_{10}, \dots, \beta_{p0})^T$ . Without loss of generality, we assume the first  $q$  elements of  $\boldsymbol{\beta}_0$  are nonzero and the last  $p - q$  elements are zero. That is,  $\boldsymbol{\beta}_0$  can be written

as  $\beta_0 = (\beta_{10}^T, \beta_{20}^T)^T$ , where  $\beta_{10}$  is a  $q$ -dimensional vector of nonzero elements and  $\beta_{20} = \mathbf{0}$ , a  $(p - q)$ -dimensional vector of zero. We decompose  $\beta$  and  $\mathbf{x}_i$  accordingly and write  $\beta = (\beta_1^T, \beta_2^T)^T$  and  $\mathbf{x}_i = (\mathbf{x}_{i1}^T, \mathbf{x}_{i2}^T)^T, i = 1, \dots, n$ . We focus on penalized linear expectile regression and assume that  $y$  depends on  $\mathbf{x}$  in a linear fashion. Namely, in our sample, we have the following linear model:

$$y_i = \mathbf{x}_i^T \beta + \epsilon_i = \mathbf{x}_{i1}^T \beta_1 + \mathbf{x}_{i2}^T \beta_2 + \epsilon_i, \quad i = 1, \dots, n. \tag{3}$$

For some predetermined  $\tau \in (0, 1)$ , the  $\tau$ th expectile of the random error  $\epsilon_i$  is zero.

We consider the following objective function of the penalized expectile regression model:

$$R_n(\beta; \tau) = \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^T \beta) + n \sum_{j=1}^p p_{\lambda_n}(|\beta_j|), \tag{4}$$

where  $\rho_\tau(\cdot)$  is defined in (2) and  $p_{\lambda_n}(\cdot)$  is a penalty function with tuning parameter  $\lambda_n$ . The estimator of the penalized expectile regression minimizes (4). Among various penalty functions, we consider SCAD and adaptive LASSO because of their properties of unbiasedness, sparsity, and continuity. For SCAD,

$$p_{\lambda_n}(|\beta_j|) = \begin{cases} \lambda_n |\beta_j| & \text{if } |\beta_j| \leq \lambda_n; \\ -\frac{|\beta_j|^2 - 2a\lambda_n |\beta_j| + \lambda_n^2}{2(a-1)} & \text{if } \lambda_n < |\beta_j| \leq a\lambda_n; \\ \frac{(a+1)\lambda_n^2}{2} & \text{if } |\beta_j| > a\lambda_n \end{cases} \tag{5}$$

for some  $a > 2$ . For adaptive LASSO,

$$p_{\lambda_n}(|\beta_j|) = \lambda_n w_j |\beta_j|,$$

where  $w_j$  is a prespecified weight. In the classical linear regression setting, Zou (2006) suggests to construct the adaptive weights using the least squares estimates  $\hat{\beta}(\text{ols})$ , i.e.,  $w_j = 1/|\hat{\beta}_j(\text{ols})|^\gamma$  for some chosen  $\gamma > 0$ . We use the expectile regression estimates without any penalty to construct our weights in the same fashion.

### 3.1 Asymptotic properties

We study the theoretical properties of the expectile regression with SCAD and adaptive LASSO penalties in a similar setting in Wu and Liu (2009). The following conditions are needed to facilitate the theoretical proofs.

**Condition 1:** For any given  $\tau \in (0, 1)$ , the errors  $\{\epsilon_i, i = 1, 2, \dots, n\}$  are independent and identically distributed, with  $\tau$ th expectile zero and a continuous, positive density  $f(\cdot)$  in a neighborhood of zero. We further assume  $E(\epsilon_i^4) < \infty$ .

**Condition 2:** The row vectors of the design matrix  $\mathbf{X}$ ,  $\{\mathbf{x}_i, i = 1, 2, \dots, n\}$ , are a deterministic sequence. We assume that there exists a positive definite matrix  $\Sigma$

such that  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T = \Sigma$ . In addition, we assume  $\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \mathbf{x}_i)^2 = O(1)$ . Denote the top-left  $q$ -by- $q$  submatrix of  $\Sigma$  by  $\Sigma_{11}$  and the right-bottom  $(p - q)$ -by- $(p - q)$  submatrix of  $\Sigma$  by  $\Sigma_{22}$ .

For SCAD expectile regression, we show consistency and oracle properties.

**Theorem 1** (Consistency). *Assume that Conditions 1 and 2 are satisfied. If  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ , there exists a local minimizer of  $R_n(\boldsymbol{\beta}; \tau)$  in (4),  $\hat{\boldsymbol{\beta}}^{(\text{SCAD})}$ , such that*

$$\|\hat{\boldsymbol{\beta}}^{(\text{SCAD})} - \boldsymbol{\beta}_0\| = O_p(n^{-\frac{1}{2}}).$$

Denote  $g_\tau(\epsilon_i) = \rho'_\tau(\epsilon_i - t) |_{t=0} = -2\tau v I(v \geq 0) - 2(1 - \tau)v I(v < 0)$ ,  $h_\tau(\epsilon_i) = \rho''_\tau(\epsilon_i - t) |_{t=0} = 2\tau I(v \geq 0) + 2(1 - \tau)I(v < 0)$ ,  $i = 1, \dots, n$ , and  $\sigma_{g_\tau}^2 = \text{Var}(g_\tau(\epsilon_i))$ ,  $\mu_{h_\tau} = E(h_\tau(\epsilon_i))$ .

**Theorem 2** (Oracle property) *Assume that Conditions 1 and 2 are satisfied. 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  $\hat{\boldsymbol{\beta}}^{(\text{SCAD})} = \begin{pmatrix} \hat{\boldsymbol{\beta}}_1^{(\text{SCAD})} \\ \hat{\boldsymbol{\beta}}_2^{(\text{SCAD})} \end{pmatrix}$  in Theorem 1 satisfies*

- (a) Sparsity:  $\hat{\boldsymbol{\beta}}_2^{(\text{SCAD})} = \mathbf{0}$ ;
- (b) Asymptotic normality:  $\sqrt{n}(\hat{\boldsymbol{\beta}}_1^{(\text{SCAD})} - \boldsymbol{\beta}_{10}) \xrightarrow{\mathcal{L}} N(\mathbf{0}, \sigma_{g_\tau}^2 / \mu_{h_\tau}^2 \Sigma_{11}^{-1})$ . Here  $\xrightarrow{\mathcal{L}}$  denotes convergence in distribution and  $\Sigma_{11}$  is defined in Condition 2.

The oracle property also holds for adaptive LASSO expectile regression.

**Theorem 3** (Oracle property). *Assume that Conditions 1 and 2 are satisfied. If  $\sqrt{n}\lambda_n \rightarrow 0$  and  $n^{(\nu+1)/2}\lambda_n \rightarrow \infty$  as  $n \rightarrow \infty$ , then the adaptive LASSO expectile regression estimator  $\hat{\boldsymbol{\beta}}^{(\text{AL})} = \begin{pmatrix} \hat{\boldsymbol{\beta}}_1^{(\text{AL})} \\ \hat{\boldsymbol{\beta}}_2^{(\text{AL})} \end{pmatrix}$ , which minimizes (4), satisfies*

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

*Remark 1* Notice that the asymptotic covariance matrix of the penalized expectile regression estimator requires the fourth moment of the errors while penalized quantile estimator in Wu and Liu (2009) requires the second moment. However, the asymptotic covariance matrix of the penalized quantile regression estimator relies on the density function of the errors at the origin, which is usually unknown and difficult to estimate.

The following corollary extends Theorems 2 and 3 to non-*i.i.d.* errors under the following assumptions:

**Condition 3:** As  $n \rightarrow \infty$ ,  $\frac{1}{n} \max_i \mathbf{x}_i^T \mathbf{x}_i \text{Var}(g_\tau(\epsilon_i)) \rightarrow 0$ .

**Condition 4:** We assume that there exist positive definite matrices  $\Sigma^{g_\tau}$  and  $\Sigma^{h_\tau}$  such that  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \text{Var}(g_\tau(\epsilon_i)) = \Sigma^{g_\tau}$  and  $\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top \text{E}(h_\tau(\epsilon_i)) = \Sigma^{h_\tau}$ . Denote the top-left  $q$ -by- $q$  submatrix of  $\Sigma^{g_\tau}$  by  $\Sigma_{11}^{g_\tau}$  and  $\Sigma^{h_\tau}$  by  $\Sigma_{11}^{h_\tau}$ , respectively.

**Corollary 1** For non-i.i.d. random errors that satisfy Conditions 3 and 4, Theorems 2 and 3 hold with the limiting distribution

$$N(\mathbf{0}, (\Sigma_{11}^{h_\tau})^{-1} \Sigma_{11}^{g_\tau} (\Sigma_{11}^{h_\tau})^{-1}).$$

The proofs of Theorems 1–3 and Corollary 1 are provided in Appendix.

*Remark 2* Gu and Zou (2016) consider heteroscedasticity in the sense that they allow a different model depending on a different expectile of the conditional distribution. However, they assume i.i.d. for the regression error  $\epsilon_i$  in their theory, while we consider both i.i.d. and non-i.i.d. cases in this paper. Our theory for non-i.i.d. can be applied to error distributions whose fourth moment exists as in Condition 1.

### 3.2 Computational algorithms

This subsection concerns the computational aspect of the proposed penalized expectile regression. Unlike quantile regression, which minimizes an asymmetrically weighted sum of absolute deviations of residuals, expectile regression solves a minimization problem of an asymmetrically weighted sum of squared residuals.

#### 3.2.1 Adaptive LASSO penalty

In the case of the adaptive LASSO, the optimization problem becomes

$$\arg \min_{\beta} \left\{ \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^\top \beta) + n\lambda_n \sum_{j=1}^p w_j |\beta_j| \right\}. \tag{6}$$

Problem (6) can be formulated as the standard form of quadratic programming. We take positive slack variables  $\xi = (\xi_i)_{i=1, \dots, n}$ ,  $\zeta = (\zeta_i)_{i=1, \dots, n}$  and  $\mathbf{v} = (v_j)_{j=1, \dots, p}$ , and then, the equivalent problem of (6) is given as

$$\min_{\xi, \zeta, \mathbf{v}, \beta} \left\{ \tau \sum_{i=1}^n \xi_i^2 + (1 - \tau) \sum_{i=1}^n \zeta_i^2 + n\lambda_n \sum_{j=1}^p w_j v_j \right\},$$

subject to  $\xi_i - \zeta_i = y_i - \mathbf{x}_i^\top \beta$ ,  $i = 1, \dots, n$ , and  $-v_j \leq \beta_j \leq v_j$ ,  $j = 1, \dots, p$ . To solve this, we use the Rmosek (Friebert 2014) interface to MOSEK (MOSEK ApS 2011), which is known to provide a flexible and reliable platform for convex programming (Koenker and Mizera 2014).

### 3.2.2 SCAD penalty

For the SCAD penalty with (5), because the objective function  $R_n(\boldsymbol{\beta}; \tau)$  in (4) is not convex, we apply the concave–convex procedure (CCCP) (Yuille and Rangarajan 2003; Zou and Li 2008; Kim et al. 2008). The CCCP searches for a local minimizer by successively minimizing the locally upper-tight convex function of the objective function. When the objective function can be decomposed into a sum of convex and concave functions, the locally upper-tight convex function can be obtained by locally linearizing the concave function. We name it as local linear approximation (LLA) algorithm (Zou and Li 2008).

For  $R_n(\boldsymbol{\beta}; \tau)$  in (4), note that  $R_n(\boldsymbol{\beta}; \tau) = R_n^{(\text{conv})}(\boldsymbol{\beta}; \tau) + R_n^{(\text{conc})}(\boldsymbol{\beta}; \tau)$ , where  $R_n^{(\text{conv})}(\boldsymbol{\beta}; \tau)$  is a convex function and  $R_n^{(\text{conc})}(\boldsymbol{\beta}; \tau)$  is a differentiable concave function given by

$$R_n^{(\text{conv})}(\boldsymbol{\beta}; \tau) = \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) + n\lambda_n \sum_{j=1}^p |\beta_j|,$$

$$R_n^{(\text{conc})}(\boldsymbol{\beta}; \tau) = \sum_{j=1}^p (np\lambda_n (|\beta_j|) - n\lambda_n |\beta_j|).$$

Given an initial solution  $\hat{\boldsymbol{\beta}}^{(1)}$ , the locally upper-tight convex function of  $R_n(\boldsymbol{\beta}; \tau)$  becomes

$$R_n^{(\text{conv})}(\boldsymbol{\beta}; \tau) + \{\partial R_n^{(\text{conc})}(\hat{\boldsymbol{\beta}}^{(1)}; \tau) / \partial \boldsymbol{\beta}\}^\top (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(1)}) + R_n^{(\text{conc})}(\hat{\boldsymbol{\beta}}^{(1)}; \tau).$$

Then it iteratively updates  $\hat{\boldsymbol{\beta}}^{(1)}$  with the minimizer of the locally upper-tight convex function, which leads  $\hat{\boldsymbol{\beta}}^{(1)}$  to converge eventually to a local minimizer, and the LLA algorithm is summarized as follows:

- Initialize  $\hat{\boldsymbol{\beta}}^{(1)}$  and  $\lambda_n > 0$ .
- For  $m = 1, 2, \dots$ , update the following equation until convergence:
 
$$\hat{\boldsymbol{\beta}}^{(m+1)} = \operatorname{argmin}_{\boldsymbol{\beta}} R_n^{(\text{conv})}(\boldsymbol{\beta}; \tau) + \{\partial R_n^{(\text{conc})}(\hat{\boldsymbol{\beta}}^{(m)}; \tau) / \partial \boldsymbol{\beta}\}^\top (\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}^{(m)}) + R_n^{(\text{conc})}(\hat{\boldsymbol{\beta}}^{(m)}; \tau).$$

Let  $\mathcal{V}_1^{(m)} = \{j : |\hat{\beta}_j^{(m)}| \leq \lambda_n\}$ ,  $\mathcal{V}_2^{(m)} = \{j : \lambda_n < |\hat{\beta}_j^{(m)}| \leq a\lambda_n\}$ , and  $\mathcal{V}_3^{(m)} = \{j : |\hat{\beta}_j^{(m)}| > a\lambda_n\}$ , where  $\hat{\boldsymbol{\beta}}^{(m)}$  is the solution of the  $m$ th iteration step in the LLA. The  $m$ th problem excluding irrelevant parameters is equivalent to minimizing

$$\sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) + n\lambda_n \sum_{j=1}^p w_j^{(m)} |\beta_j|, \tag{7}$$

where the weights  $w_j^{(m)}$  are defined as follows:  $w_j^{(m)} = 1, j \in \mathcal{V}_1^{(m)}$ ,  $w_j^{(m)} = 1 - (|\hat{\beta}_j^{(m)}| - \lambda_n) / (a - 1)\lambda_n, j \in \mathcal{V}_2^{(m)}$ , and  $w_j^{(m)} = 0, j \in \mathcal{V}_3^{(m)}$ . Then, problem (7) has

the same form as (6). In our numerical study, we set  $\hat{\beta}^{(1)} = 0$  in the LLA algorithm, which leads to solving the LASSO problem at the initial stage.

We choose the tuning parameter  $\lambda$  to minimize the validation error

$$\sum_{i=1}^n \rho_{\tau}(y_i - \mathbf{x}_i^T \hat{\beta})$$

for the computed estimate  $\hat{\beta}$  as in Gu and Zou (2016) to make a fair comparison in our numerical study.

## 4 Simulation

In this section, we compare our proposed penalized expectile regression with the penalized quantile regression (Wu and Liu 2009) and the SALES (Gu and Zou 2016), using the SCAD and adaptive LASSO penalties. The SALES solve the optimization problem in (4) using the cyclic coordinate descent and proximal gradient algorithms. We consider two simulation settings.

**Setting 1:** We generate our data from the linear model with non-*i.i.d.* random errors studied by Kocherginsky et al. (2005) and Wu and Liu (2009),

$$Y = 1 + X_1 + X_2 + X_3 + (1 + X_3)\epsilon, \quad (8)$$

where  $X_2 = X_1 + X_3 + Z$ , with both  $X_1$  and  $Z$  generated from the standard normal distribution and  $X_3$  generated from the uniform distribution on  $[0, 1]$ . The variables  $X_1$ ,  $X_3$ ,  $Z$  and  $\epsilon$  are mutually independent.

**Setting 2:** We use a similar example in Wang et al. (2012) and Gu and Zou (2016):

$$Y = X_6 + X_{12} + X_{15} + X_{20} + (1 + 0.7X_1)\epsilon. \quad (9)$$

We first generate  $(Z_1, \dots, Z_p)^T$  from the multivariate normal distribution with zero mean and the  $(i, j)$ th element of the covariance matrix  $0.5^{|i-j|}$ . Then, we set  $X_1 = \Phi(Z_1)$  and  $X_j = Z_j$  for  $j = 2, 3, \dots, p$ , where  $\Phi$  is the standard normal cumulative distribution function. Note that  $X_1$  does not have an impact on the mean but only on the variance.

For each case, we generate independent training and validation data sets of size  $n = 100$  and  $200$ , and testing data set of size  $n = 10,000$ . The training data set is used to obtain the penalized expectile regression estimate given a fixed expectile value  $\tau$  and a tuning parameter  $\lambda_n$ . Let  $\lambda_n$  range from  $0.01$  to  $10$  with the gap,  $\log(\lambda_{i+1}) - \log(\lambda_i) = 0.1$ , and it is chosen by the validation error described in Sect. 3.2. Then, a test error whose definition is given below is computed on the testing data set. We repeat the whole procedure  $100$  times and evaluate the performance of the penalized expectile and quantile regression. For penalized expectile regression, we consider five different expectiles  $\tau = 0.1, 0.25, 0.5, 0.75, 0.9$ . Notice that expectiles and quantiles have one-to-one correspondence given a specific distribution (Table 1). Therefore, it is

fair to compare the test error of penalized  $\tau$ th expectile regression estimate and the corresponding penalized  $\alpha$ th quantile regression estimate. We apply the LLA algorithm to both penalized expectile and quantile regressions.

For a comprehensive comparison among penalized expectile/quantile regressions and SALES, we consider the standard normal and a heavy-tailed distribution  $t(10)$  for  $\epsilon$  in Setting 1, and the standard normal distribution only in Setting 2. For each distribution of  $\epsilon$ , we study two different dimensions of the predictors,  $p = 10$  and  $50$ , by adding independent noise variables following the standard normal distribution,  $N(0, 1)$  to the original  $X_1, X_2$ , and  $X_3$  in Setting 1 and  $X_1, X_6, X_{12}, X_{15}$ , and  $X_{20}$  in Setting 2. For fair comparison, we calculate the corresponding quantile for each expectile for a given distribution. Since test errors use different loss functions, we report the absolute errors (AE) defined as the absolute distance between the estimated and the true parameters, i.e.,  $\sum_{j=1}^p |\hat{\beta}_j - \beta_j|$ . In the results, ET denotes the proposed expectile method, ES the SALES, and QT the quantile regression.

Table 2 reports the absolute errors for model (8) with  $\epsilon$  generated from the standard normal and  $t(10)$  distribution with  $p = 10, 50$  and  $n = 100, 200$ . As expected, all methods produce smaller errors for lower dimension and larger sample size. For both distributions and both small and large predictors, it can be seen that there is no single estimator that dominates the performance in terms of the absolute errors. In comparison between ALASSO ET and SCAD ET, ALASSO ET tends to yield smaller errors for high  $\tau = 0.75, 0.9$  while SCAD ET for low  $\tau = 0.1, 0.25, 0.5$ . They both produce smaller absolute errors than the SALES methods (ALASSO ES and SCAD ES) in many cases. We see that the performance of penalized expectile and quantile regression is competitive. Penalized expectile regression tends to perform better for  $p=10$ , and penalized quantile regression tends to perform better for  $p = 50$  for both sample sizes.

Tables 3 and 4 report the variable selection results for Setting 1. Again, all methods show better selection for lower dimension and larger sample size. It can be seen that every method successfully selects  $X_2$  with high probability ( $P_2$ ) for all of the expectiles or quantiles. For  $X_1$ , ALASSO ET chooses it with the highest probability ( $P_1$ ) for each level of expectile or quantile across different dimensional cases (0.91–1.00) followed by SCAD ET (0.86–1.00). Both ES and QT produce similar but slightly lower probabilities (0.74–1.00 and 0.79–1.00, respectively) compared to ET for both penalties. On the other hand, selecting  $X_3$  is a challenging task for the six methods in the sense that  $X_3$  can be easily dominated by  $(1 + X_3)\epsilon$  in (8), because  $X_3 \sim U(0, 1)$  and  $\epsilon \sim N(0, 1)$ . The proportions including  $X_3$  ( $P_3$ ) for ET are higher than those of ES in most cases and are similar to those of QT for both penalties. In conclusion, the proposed penalized expectile regression (ET) tends to produce smaller absolute errors, performs better variable selection for different distributions and dimensions across various  $\tau$  values compared to the SALES (ES), and shows competitive performance compared to the penalized quantile regression in terms of absolute errors and variable selection.

Table 5 reports the absolute errors for model (9) with  $\epsilon$  generated from the standard normal distribution with  $p = 20, 50$  and  $n = 100, 200$ . Again, there is no dominant winner in terms of the absolute errors. It can be seen that ET and ES perform similarly and tend to yield smaller errors than QT in most cases.

**Table 2** Absolute errors for Setting 1

<i>n</i>	<i>p</i>	$\tau$	ALASSO			SCAD			
			ET	ES	QT	ET	ES	QT	
100	10	0.10	Normal	1.48(0.04)	1.48(0.04)	1.57(0.05)	1.37(0.06)	1.52(0.03)	1.48(0.06)
			<i>t</i> (10)	1.75(0.06)	1.76(0.06)	1.75(0.05)	1.69(0.07)	1.80(0.06)	1.74(0.08)
		0.25	Normal	1.31(0.04)	1.30(0.04)	1.39(0.05)	1.16(0.05)	1.43(0.03)	1.22(0.05)
			<i>t</i> (10)	1.46(0.05)	1.45(0.05)	1.46(0.04)	1.32(0.06)	1.61(0.05)	1.34(0.06)
			Normal	1.11(0.05)	1.11(0.05)	1.28(0.05)	1.09(0.05)	1.39(0.04)	1.15(0.06)
	0.50	0.10	<i>t</i> (10)	1.21(0.05)	1.22(0.05)	1.28(0.05)	1.18(0.05)	1.56(0.04)	1.23(0.06)
			Normal	1.13(0.05)	1.14(0.05)	1.23(0.06)	1.13(0.05)	1.50(0.04)	1.25(0.07)
		0.75	<i>t</i> (10)	1.27(0.06)	1.28(0.06)	1.33(0.06)	1.34(0.06)	1.74(0.05)	1.46(0.08)
			Normal	1.40(0.06)	1.39(0.06)	1.43(0.06)	1.48(0.07)	1.79(0.05)	1.55(0.08)
			<i>t</i> (10)	1.59(0.07)	1.61(0.08)	1.75(0.08)	1.86(0.09)	2.12(0.07)	1.93(0.09)
50	10	0.10	Normal	1.80(0.06)	1.78(0.06)	1.93(0.07)	1.67(0.08)	1.77(0.06)	1.65(0.07)
			<i>t</i> (10)	1.93(0.07)	1.93(0.07)	2.08(0.07)	1.78(0.08)	1.99(0.06)	1.88(0.08)
		0.25	Normal	1.67(0.06)	1.68(0.06)	1.81(0.06)	1.51(0.06)	1.78(0.05)	1.40(0.07)
			<i>t</i> (10)	1.72(0.06)	1.73(0.06)	1.92(0.07)	1.66(0.06)	1.92(0.05)	1.48(0.07)
			Normal	1.65(0.06)	1.62(0.06)	1.82(0.07)	1.47(0.06)	1.92(0.05)	1.34(0.08)
	0.50	0.10	<i>t</i> (10)	1.62(0.06)	1.61(0.06)	1.72(0.07)	1.61(0.06)	2.01(0.05)	1.28(0.06)
			Normal	1.58(0.06)	1.60(0.06)	1.81(0.07)	1.65(0.07)	2.04(0.05)	1.44(0.07)
		0.75	<i>t</i> (10)	1.63(0.06)	1.62(0.06)	1.66(0.06)	1.71(0.07)	2.10(0.05)	1.40(0.07)
			Normal	1.82(0.07)	1.82(0.07)	1.97(0.07)	1.79(0.08)	2.09(0.06)	1.85(0.09)
			<i>t</i> (10)	1.81(0.08)	1.80(0.08)	2.11(0.08)	2.03(0.09)	2.20(0.05)	2.11(0.11)

Table 2 continued

<i>n</i>	<i>p</i>	$\tau$	ALASSO			SCAD			
			ET	ES	QT	ET	ES	QT	
200	10	0.10	Normal	1.33(0.03)	1.34(0.03)	1.37(0.03)	1.33(0.04)	1.37(0.03)	1.26(0.04)
		$t(10)$	1.51(0.04)	1.52(0.05)	1.53(0.04)	1.49(0.05)	1.51(0.04)	1.54(0.05)	
		0.25	Normal	1.12(0.03)	1.13(0.03)	1.18(0.03)	1.01(0.04)	1.24(0.03)	0.97(0.04)
		$t(10)$	1.22(0.04)	1.24(0.04)	1.26(0.04)	1.08(0.04)	1.35(0.03)	1.06(0.05)	
		0.50	Normal	0.76(0.03)	0.77(0.03)	0.96(0.04)	0.76(0.04)	0.99(0.04)	0.76(0.04)
	$t(10)$	0.90(0.04)	0.89(0.04)	0.92(0.04)	0.83(0.04)	1.18(0.04)	0.81(0.04)		
	0.75	Normal	0.68(0.03)	0.68(0.03)	0.82(0.04)	0.82(0.04)	1.04(0.04)	0.92(0.05)	
	$t(10)$	0.86(0.04)	0.87(0.04)	0.88(0.05)	0.96(0.05)	1.27(0.04)	0.96(0.05)		
	0.90	Normal	0.93(0.05)	0.94(0.05)	1.08(0.05)	1.21(0.05)	1.39(0.05)	1.35(0.06)	
	$t(10)$	1.25(0.05)	1.26(0.06)	1.21(0.06)	1.46(0.07)	1.74(0.05)	1.45(0.05)		
50	10	0.10	Normal	1.44(0.03)	1.44(0.03)	1.54(0.04)	1.25(0.04)	1.38(0.03)	1.30(0.04)
		$t(10)$	1.65(0.05)	1.64(0.05)	1.63(0.05)	1.50(0.06)	1.62(0.05)	1.47(0.06)	
		0.25	Normal	1.32(0.04)	1.32(0.03)	1.38(0.03)	1.09(0.05)	1.41(0.03)	1.01(0.04)
		$t(10)$	1.43(0.05)	1.43(0.05)	1.42(0.04)	1.21(0.04)	1.55(0.04)	1.10(0.05)	
		0.50	Normal	1.08(0.04)	1.09(0.04)	1.18(0.05)	0.84(0.05)	1.49(0.03)	0.78(0.04)
	$t(10)$	1.21(0.06)	1.20(0.05)	1.23(0.05)	1.06(0.05)	1.62(0.04)	0.87(0.05)		
	0.75	Normal	0.93(0.05)	0.92(0.05)	1.04(0.05)	0.92(0.05)	1.62(0.04)	0.87(0.05)	
	$t(10)$	1.09(0.05)	1.10(0.05)	1.16(0.06)	1.10(0.05)	1.79(0.05)	1.00(0.05)		
	0.90	Normal	1.08(0.05)	1.07(0.05)	1.27(0.06)	1.23(0.05)	1.73(0.05)	1.32(0.06)	
	$t(10)$	1.36(0.06)	1.36(0.06)	1.44(0.07)	1.48(0.06)	2.04(0.06)	1.46(0.07)		

**Table 3** Variable selection results with  $n = 100$  in Setting 1

$n$	$p$	$\tau$		ALASSO			SCAD			
				ET	ES	QT	ET	ES	QT	
100	10	0.10	$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.98 <sub>(0.01)</sub>	0.98 <sub>(0.01)</sub>	
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.25 <sub>(0.04)</sub>	0.24 <sub>(0.04)</sub>	0.22 <sub>(0.04)</sub>	0.38 <sub>(0.05)</sub>	0.04 <sub>(0.02)</sub>	0.37 <sub>(0.05)</sub>	
		0.25	$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.98 <sub>(0.01)</sub>	
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.43 <sub>(0.05)</sub>	0.44 <sub>(0.05)</sub>	0.36 <sub>(0.05)</sub>	0.44 <sub>(0.05)</sub>	0.13 <sub>(0.03)</sub>	0.53 <sub>(0.05)</sub>	
		0.50	$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.64 <sub>(0.05)</sub>	0.64 <sub>(0.05)</sub>	0.56 <sub>(0.05)</sub>	0.68 <sub>(0.05)</sub>	0.29 <sub>(0.05)</sub>	0.73 <sub>(0.04)</sub>	
	0.75	$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>							
		$P_3$	0.82 <sub>(0.04)</sub>	0.82 <sub>(0.04)</sub>	0.75 <sub>(0.04)</sub>	0.89 <sub>(0.03)</sub>	0.49 <sub>(0.05)</sub>	0.88 <sub>(0.03)</sub>		
	0.90	$P_1$	0.97 <sub>(0.02)</sub>	0.98 <sub>(0.01)</sub>	0.94 <sub>(0.02)</sub>	0.97 <sub>(0.02)</sub>	0.91 <sub>(0.03)</sub>	0.93 <sub>(0.03)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>						
		$P_3$	0.82 <sub>(0.04)</sub>	0.83 <sub>(0.04)</sub>	0.79 <sub>(0.04)</sub>	0.94 <sub>(0.02)</sub>	0.52 <sub>(0.05)</sub>	0.92 <sub>(0.03)</sub>		
	50	0.10		$P_1$	0.93 <sub>(0.03)</sub>	0.93 <sub>(0.03)</sub>	0.93 <sub>(0.03)</sub>	0.89 <sub>(0.03)</sub>	0.90 <sub>(0.03)</sub>	0.85 <sub>(0.04)</sub>
				$P_2$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>
				$P_3$	0.15 <sub>(0.04)</sub>	0.14 <sub>(0.03)</sub>	0.08 <sub>(0.03)</sub>	0.24 <sub>(0.04)</sub>	0.00 <sub>(0.00)</sub>	0.26 <sub>(0.04)</sub>
0.25			$P_1$	0.95 <sub>(0.02)</sub>	0.95 <sub>(0.02)</sub>	0.89 <sub>(0.03)</sub>	0.91 <sub>(0.03)</sub>	0.92 <sub>(0.03)</sub>	0.88 <sub>(0.03)</sub>	
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.18 <sub>(0.04)</sub>	0.20 <sub>(0.04)</sub>	0.22 <sub>(0.04)</sub>	0.33 <sub>(0.05)</sub>	0.00 <sub>(0.00)</sub>	0.33 <sub>(0.05)</sub>	
0.50			$P_1$	0.95 <sub>(0.02)</sub>	0.94 <sub>(0.02)</sub>	0.93 <sub>(0.03)</sub>	0.94 <sub>(0.02)</sub>	0.92 <sub>(0.03)</sub>	0.91 <sub>(0.03)</sub>	
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.41 <sub>(0.05)</sub>	0.41 <sub>(0.05)</sub>	0.36 <sub>(0.05)</sub>	0.52 <sub>(0.05)</sub>	0.00 <sub>(0.00)</sub>	0.53 <sub>(0.05)</sub>	
0.75		$P_1$	0.95 <sub>(0.02)</sub>	0.95 <sub>(0.02)</sub>	0.93 <sub>(0.03)</sub>	0.91 <sub>(0.03)</sub>	0.87 <sub>(0.03)</sub>	0.89 <sub>(0.03)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>	0.98 <sub>(0.01)</sub>						
		$P_3$	0.53 <sub>(0.05)</sub>	0.53 <sub>(0.05)</sub>	0.39 <sub>(0.05)</sub>	0.62 <sub>(0.05)</sub>	0.02 <sub>(0.01)</sub>	0.66 <sub>(0.05)</sub>		
0.90		$P_1$	0.91 <sub>(0.03)</sub>	0.91 <sub>(0.03)</sub>	0.85 <sub>(0.04)</sub>	0.86 <sub>(0.03)</sub>	0.72 <sub>(0.05)</sub>	0.79 <sub>(0.04)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.98 <sub>(0.01)</sub>		
		$P_3$	0.53 <sub>(0.05)</sub>	0.54 <sub>(0.05)</sub>	0.46 <sub>(0.05)</sub>	0.68 <sub>(0.05)</sub>	0.00 <sub>(0.00)</sub>	0.68 <sub>(0.05)</sub>		

Table 6 reports the variable selection results with  $n = 100$  for Setting 2. We do not report the results for  $n = 200$  because the results are similar. All methods successfully select the four variables in the mean ( $X_6, X_{12}, X_{15}$ , and  $X_{20}$ ) with high probabilities. For  $X_1$ , all methods show low probabilities ( $P_1$ ) because it is not part of the mean terms. Overall, SCAD ES has the lowest probabilities (0–0.03), and ALASSO ES performs similarly with ALASSO ET (0.05–0.26), followed by ALASSO QT (0.14–0.33). On the other hand, SCAD QT shows the highest probabilities (0.16–0.40). In conclusion, the proposed penalized expectile regression (ET) tends to produce similar

**Table 4** Variable selection results with  $n = 200$  for Setting 1

$n$	$p$	$\tau$		ALASSO			SCAD			
				ET	ES	QT	ET	ES	QT	
200	10	0.10	$P_1$	1.00 <sub>(0.00)</sub>						
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.20 <sub>(0.04)</sub>	0.22 <sub>(0.04)</sub>	0.14 <sub>(0.03)</sub>	0.29 <sub>(0.05)</sub>	0.08 <sub>(0.03)</sub>	0.30 <sub>(0.05)</sub>	
		0.25	$P_1$	1.00 <sub>(0.00)</sub>						
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.45 <sub>(0.05)</sub>	0.46 <sub>(0.05)</sub>	0.35 <sub>(0.05)</sub>	0.51 <sub>(0.05)</sub>	0.21 <sub>(0.04)</sub>	0.64 <sub>(0.05)</sub>	
		0.50	$P_1$	1.00 <sub>(0.00)</sub>						
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.87 <sub>(0.03)</sub>	0.87 <sub>(0.03)</sub>	0.84 <sub>(0.04)</sub>	0.85 <sub>(0.04)</sub>	0.62 <sub>(0.05)</sub>	0.92 <sub>(0.03)</sub>	
	0.75	$P_1$	1.00 <sub>(0.00)</sub>							
		$P_2$	1.00 <sub>(0.00)</sub>							
		$P_3$	0.97 <sub>(0.02)</sub>	0.97 <sub>(0.02)</sub>	0.93 <sub>(0.03)</sub>	0.96 <sub>(0.02)</sub>	0.82 <sub>(0.04)</sub>	0.99 <sub>(0.01)</sub>		
	0.90	$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>							
		$P_3$	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.86 <sub>(0.03)</sub>	1.00 <sub>(0.00)</sub>		
	50	0.10	0.10	$P_1$	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.98 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.97 <sub>(0.02)</sub>
				$P_2$	1.00 <sub>(0.00)</sub>					
				$P_3$	0.06 <sub>(0.02)</sub>	0.05 <sub>(0.02)</sub>	0.06 <sub>(0.02)</sub>	0.18 <sub>(0.04)</sub>	0.00 <sub>(0.00)</sub>	0.16 <sub>(0.04)</sub>
0.25			$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.98 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.22 <sub>(0.04)</sub>	0.22 <sub>(0.04)</sub>	0.17 <sub>(0.04)</sub>	0.42 <sub>(0.05)</sub>	0.00 <sub>(0.00)</sub>	0.44 <sub>(0.05)</sub>	
0.50			$P_1$	1.00 <sub>(0.00)</sub>						
			$P_2$	1.00 <sub>(0.00)</sub>						
			$P_3$	0.63 <sub>(0.05)</sub>	0.62 <sub>(0.05)</sub>	.51 <sub>(0.05)</sub>	0.79 <sub>(0.04)</sub>	0.05 <sub>(0.02)</sub>	0.80 <sub>(0.04)</sub>	
0.75		$P_1$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>							
		$P_3$	0.87 <sub>(0.03)</sub>	0.87 <sub>(0.03)</sub>	0.84 <sub>(0.04)</sub>	0.94 <sub>(0.02)</sub>	0.21 <sub>(0.04)</sub>	0.96 <sub>(0.02)</sub>		
0.90		$P_1$	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.97 <sub>(0.02)</sub>	0.99 <sub>(0.01)</sub>		
		$P_2$	1.00 <sub>(0.00)</sub>							
		$P_3$	0.89 <sub>(0.03)</sub>	0.89 <sub>(0.03)</sub>	0.93 <sub>(0.03)</sub>	1.00 <sub>(0.00)</sub>	0.17 <sub>(0.04)</sub>	0.99 <sub>(0.01)</sub>		

absolute errors to those of ES and smaller errors than those of QT in this simulation setting. In terms of variable selection, ES selects  $X_1$  the fewest times, followed by ET. QT often selects  $X_1$ , which might yield higher absolute errors as a result.

### 5 Real data analysis

The data set comes from a study that investigates different methodological problems associated with clean air using housing market data by [Harrison and Rubinfeld \(1978\)](#).

**Table 5** Absolute errors for Setting 2

<i>n</i>	<i>p</i>	$\tau$	ALASSO			SCAD		
			ET	ES	QT	ET	ES	QT
100	20	0.10	1.83 <sub>(0.06)</sub>	1.84 <sub>(0.06)</sub>	2.14 <sub>(0.10)</sub>	1.84 <sub>(0.08)</sub>	1.56 <sub>(0.05)</sub>	2.01 <sub>(0.11)</sub>
		0.25	1.58 <sub>(0.05)</sub>	1.56 <sub>(0.05)</sub>	1.76 <sub>(0.06)</sub>	1.49 <sub>(0.05)</sub>	1.41 <sub>(0.04)</sub>	1.77 <sub>(0.11)</sub>
		0.50	1.51 <sub>(0.04)</sub>	1.51 <sub>(0.04)</sub>	1.74 <sub>(0.07)</sub>	1.37 <sub>(0.05)</sub>	1.38 <sub>(0.04)</sub>	1.54 <sub>(0.07)</sub>
		0.75	1.54 <sub>(0.05)</sub>	1.53 <sub>(0.05)</sub>	1.73 <sub>(0.06)</sub>	1.30 <sub>(0.04)</sub>	1.40 <sub>(0.04)</sub>	1.41 <sub>(0.06)</sub>
		0.90	1.67 <sub>(0.06)</sub>	1.68 <sub>(0.06)</sub>	1.91 <sub>(0.08)</sub>	1.44 <sub>(0.06)</sub>	1.53 <sub>(0.05)</sub>	1.62 <sub>(0.08)</sub>
	50	0.10	2.20 <sub>(0.09)</sub>	2.19 <sub>(0.09)</sub>	2.60 <sub>(0.11)</sub>	1.80 <sub>(0.09)</sub>	1.72 <sub>(0.09)</sub>	2.14 <sub>(0.13)</sub>
		0.25	1.91 <sub>(0.07)</sub>	1.90 <sub>(0.06)</sub>	2.29 <sub>(0.09)</sub>	1.57 <sub>(0.06)</sub>	1.52 <sub>(0.04)</sub>	1.84 <sub>(0.11)</sub>
		0.50	1.80 <sub>(0.06)</sub>	1.80 <sub>(0.06)</sub>	2.05 <sub>(0.07)</sub>	1.47 <sub>(0.04)</sub>	1.54 <sub>(0.04)</sub>	1.35 <sub>(0.04)</sub>
		0.75	1.85 <sub>(0.07)</sub>	1.84 <sub>(0.06)</sub>	2.04 <sub>(0.08)</sub>	1.44 <sub>(0.04)</sub>	1.55 <sub>(0.05)</sub>	1.38 <sub>(0.06)</sub>
		0.90	2.01 <sub>(0.07)</sub>	2.00 <sub>(0.07)</sub>	2.38 <sub>(0.10)</sub>	1.49 <sub>(0.05)</sub>	1.67 <sub>(0.06)</sub>	1.79 <sub>(0.09)</sub>
200	20	0.10	1.44 <sub>(0.04)</sub>	1.46 <sub>(0.04)</sub>	1.66 <sub>(0.06)</sub>	1.52 <sub>(0.05)</sub>	1.30 <sub>(0.05)</sub>	1.60 <sub>(0.06)</sub>
		0.25	1.24 <sub>(0.03)</sub>	1.24 <sub>(0.03)</sub>	1.39 <sub>(0.04)</sub>	1.26 <sub>(0.04)</sub>	1.15 <sub>(0.03)</sub>	1.39 <sub>(0.06)</sub>
		0.50	1.16 <sub>(0.03)</sub>	1.16 <sub>(0.03)</sub>	1.28 <sub>(0.03)</sub>	1.10 <sub>(0.03)</sub>	1.11 <sub>(0.02)</sub>	1.15 <sub>(0.03)</sub>
		0.75	1.14 <sub>(0.03)</sub>	1.14 <sub>(0.03)</sub>	1.25 <sub>(0.04)</sub>	1.01 <sub>(0.03)</sub>	1.15 <sub>(0.02)</sub>	1.04 <sub>(0.04)</sub>
		0.90	1.18 <sub>(0.04)</sub>	1.20 <sub>(0.04)</sub>	1.30 <sub>(0.04)</sub>	0.99 <sub>(0.04)</sub>	1.27 <sub>(0.03)</sub>	1.07 <sub>(0.04)</sub>
	50	0.10	1.42 <sub>(0.05)</sub>	1.43 <sub>(0.05)</sub>	1.53 <sub>(0.05)</sub>	1.40 <sub>(0.05)</sub>	1.20 <sub>(0.05)</sub>	1.48 <sub>(0.05)</sub>
		0.25	1.27 <sub>(0.03)</sub>	1.26 <sub>(0.03)</sub>	1.41 <sub>(0.04)</sub>	1.18 <sub>(0.03)</sub>	1.11 <sub>(0.02)</sub>	1.23 <sub>(0.04)</sub>
		0.50	1.22 <sub>(0.03)</sub>	1.21 <sub>(0.02)</sub>	1.37 <sub>(0.03)</sub>	1.11 <sub>(0.02)</sub>	1.11 <sub>(0.02)</sub>	1.16 <sub>(0.03)</sub>
		0.75	1.25 <sub>(0.03)</sub>	1.25 <sub>(0.03)</sub>	1.43 <sub>(0.05)</sub>	1.08 <sub>(0.03)</sub>	1.16 <sub>(0.03)</sub>	1.12 <sub>(0.05)</sub>
		0.90	1.31 <sub>(0.03)</sub>	1.31 <sub>(0.03)</sub>	1.43 <sub>(0.04)</sub>	1.06 <sub>(0.03)</sub>	1.21 <sub>(0.03)</sub>	1.08 <sub>(0.03)</sub>

We apply its corrected version, which is available at [http://lib.stat.cmu.edu/datasets/boston\\_corrected.txt](http://lib.stat.cmu.edu/datasets/boston_corrected.txt) as in Wu and Liu (2009). In total, there are 506 observations and 16 variables, among which *CMEDV* (corrected median value of owner-occupied homes) is the response and the other 14 non-constant predictors include *LON* (longitude), *LAT* (latitude), *CRIM* (crime rate by town), *ZN* (proportion of residential land zoned for large lots by town), *INDUS* (proportion non-retail business acres per town), *CHAS* (Charles River dummy: 1 if tract bounds the river; 0 if not), *NOX* (nitrogen oxide concentration), *RM* (average number of rooms), *AGE* (proportion of owner-occupied homes built prior to 1940), *DIS* (weighted distances to five employment centers in Boston), *TAX* (property tax rate), *PTRATIO* (pupil-teacher ratio by town), *B* (black population proportion), and *LSTAT* (proportion of lower socioeconomic status population). As in Wu and Liu (2009), we standardize the response *CMEDV* and the 14 continuous predictors. In the application of the penalized expectile regression, we consider *CMEDV* as the response, *CHAS* and all the other standardized continuous predictors and their corresponding squares as predictors, for a total of 27 predictors.

We follow a similar procedure as in Sect. 4. First, we randomly split the data set into training, validation, and testing data sets with size 200, 200, and 106, respectively. Second, we select the tuning parameter  $\lambda$  with the validation error introduced in

**Table 6** Variable selection results with  $n = 100$  for Setting 2

$n$	$p$	$\tau$		ALASSO			SCAD		
				ET	ES	QT	ET	ES	QT
100	20	0.10	$P_1$	0.25 <sub>(0.04)</sub>	0.26 <sub>(0.04)</sub>	0.25 <sub>(0.04)</sub>	0.38 <sub>(0.05)</sub>	0.00 <sub>(0.00)</sub>	0.40 <sub>(0.05)</sub>
			$P_6$	1.00 <sub>(0.00)</sub>					
			$P_{12}$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>				
			$P_{15}$	1.00 <sub>(0.00)</sub>					
			$P_{20}$	1.00 <sub>(0.00)</sub>					
		0.25	$P_1$	0.13 <sub>(0.03)</sub>	0.14 <sub>(0.03)</sub>	0.20 <sub>(0.04)</sub>	0.17 <sub>(0.04)</sub>	0.01 <sub>(0.01)</sub>	0.30 <sub>(0.05)</sub>
			$P_6$	1.00 <sub>(0.00)</sub>					
			$P_{12}$	1.00 <sub>(0.00)</sub>					
			$P_{15}$	1.00 <sub>(0.00)</sub>					
			$P_{20}$	1.00 <sub>(0.00)</sub>					
	0.50	$P_1$	0.17 <sub>(0.04)</sub>	0.18 <sub>(0.04)</sub>	0.18 <sub>(0.04)</sub>	0.15 <sub>(0.04)</sub>	0.01 <sub>(0.01)</sub>	0.26 <sub>(0.04)</sub>	
		$P_6$	1.00 <sub>(0.00)</sub>						
		$P_{12}$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>					
		$P_{15}$	1.00 <sub>(0.00)</sub>						
		$P_{20}$	1.00 <sub>(0.00)</sub>						
	0.75	$P_1$	0.22 <sub>(0.04)</sub>	0.20 <sub>(0.04)</sub>	0.25 <sub>(0.04)</sub>	0.21 <sub>(0.04)</sub>	0.03 <sub>(0.02)</sub>	0.32 <sub>(0.05)</sub>	
		$P_6$	1.00 <sub>(0.00)</sub>						
		$P_{12}$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	
		$P_{15}$	1.00 <sub>(0.00)</sub>						
		$P_{20}$	1.00 <sub>(0.00)</sub>						
0.90	$P_1$	0.25 <sub>(0.04)</sub>	0.21 <sub>(0.04)</sub>	0.33 <sub>(0.05)</sub>	0.31 <sub>(0.05)</sub>	0.02 <sub>(0.01)</sub>	0.40 <sub>(0.05)</sub>		
	$P_6$	1.00 <sub>(0.00)</sub>							
	$P_{12}$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>		
	$P_{15}$	1.00 <sub>(0.00)</sub>							
	$P_{20}$	1.00 <sub>(0.00)</sub>							
50	0.10	$P_1$	0.13 <sub>(0.03)</sub>	0.12 <sub>(0.03)</sub>	0.22 <sub>(0.04)</sub>	0.18 <sub>(0.04)</sub>	0.00 <sub>(0.00)</sub>	0.27 <sub>(0.04)</sub>	
		$P_6$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.97 <sub>(0.02)</sub>	
		$P_{12}$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.98 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.97 <sub>(0.02)</sub>	
		$P_{15}$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.98 <sub>(0.01)</sub>	
		$P_{20}$	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.95 <sub>(0.02)</sub>	
	0.25	$P_1$	0.10 <sub>(0.03)</sub>	0.06 <sub>(0.02)</sub>	0.14 <sub>(0.03)</sub>	0.10 <sub>(0.03)</sub>	0.00 <sub>(0.00)</sub>	0.16 <sub>(0.04)</sub>	
		$P_6$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	
		$P_{12}$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>					
		$P_{15}$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>					
		$P_{20}$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.98 <sub>(0.01)</sub>	
	0.50	$P_1$	0.05 <sub>(0.02)</sub>	0.05 <sub>(0.02)</sub>	0.10 <sub>(0.03)</sub>	0.06 <sub>(0.02)</sub>	0.00 <sub>(0.00)</sub>	0.16 <sub>(0.04)</sub>	
		$P_6$	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	
		$P_{12}$	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>					
		$P_{15}$	1.00 <sub>(0.00)</sub>						

**Table 6** continued

<i>n</i>	<i>p</i>	$\tau$	ALASSO			SCAD		
			ET	ES	QT	ET	ES	QT
		<i>P</i> <sub>20</sub>	1.00 <sub>(0.00)</sub>					
	0.75	<i>P</i> <sub>1</sub>	0.09 <sub>(0.03)</sub>	0.10 <sub>(0.03)</sub>	0.15 <sub>(0.04)</sub>	0.09 <sub>(0.03)</sub>	0.00 <sub>(0.00)</sub>	0.16 <sub>(0.04)</sub>
		<i>P</i> <sub>6</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>
		<i>P</i> <sub>12</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>				
		<i>P</i> <sub>15</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>
		<i>P</i> <sub>20</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>
	0.90	<i>P</i> <sub>1</sub>	0.13 <sub>(0.03)</sub>	0.14 <sub>(0.03)</sub>	0.20 <sub>(0.04)</sub>	0.22 <sub>(0.04)</sub>	0.00 <sub>(0.00)</sub>	0.26 <sub>(0.04)</sub>
		<i>P</i> <sub>6</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.96 <sub>(0.02)</sub>
		<i>P</i> <sub>12</sub>	0.99 <sub>(0.01)</sub>	0.99 <sub>(0.01)</sub>	0.97 <sub>(0.02)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	0.96 <sub>(0.02)</sub>
		<i>P</i> <sub>15</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.97 <sub>(0.02)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.94 <sub>(0.02)</sub>
		<i>P</i> <sub>20</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.99 <sub>(0.01)</sub>	1.00 <sub>(0.00)</sub>	1.00 <sub>(0.00)</sub>	0.97 <sub>(0.02)</sub>

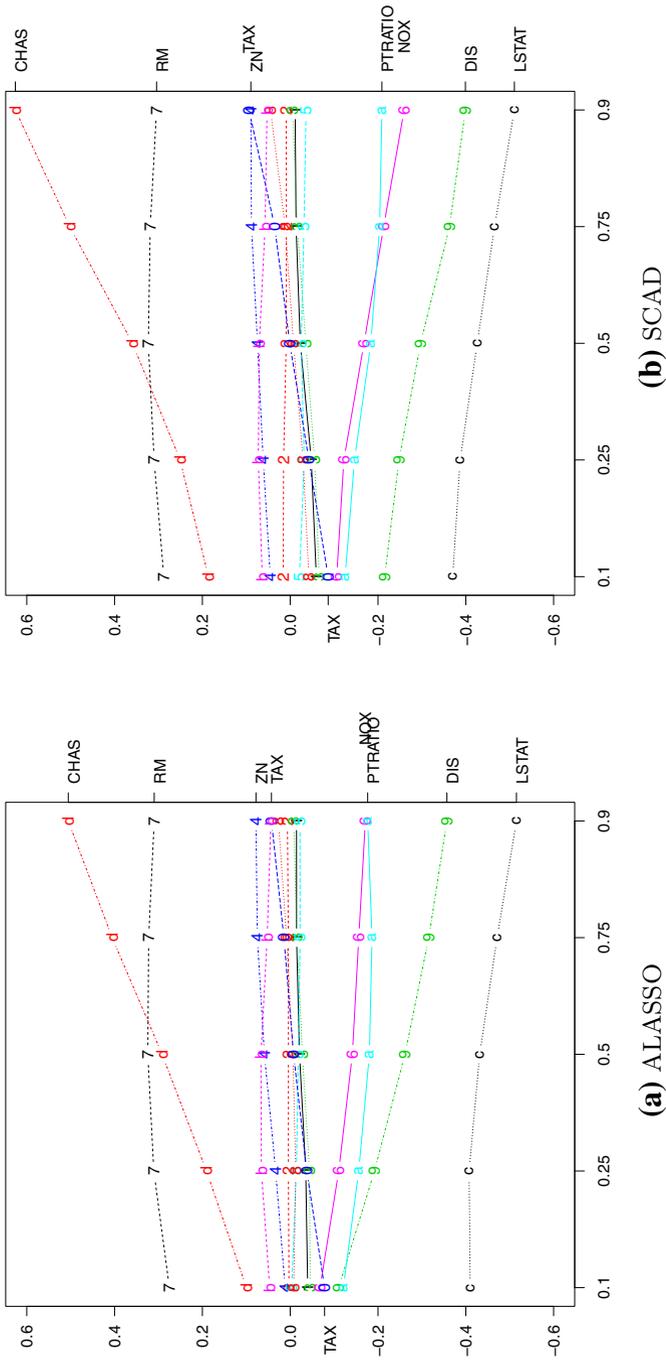
Sect. 3.2. Then, we evaluate the test error by the average expectile check loss on the testing data. We replicate 100 times and report its average values. We still consider the five different expectiles  $\tau = 0.1, 0.25, 0.5, 0.75$  and  $0.9$ . Although the underlying distribution is unknown, by assuming the normal distribution, we apply the SALES and penalized quantile regression and compare their performance. For the quantile regression, we use the quantiles that correspond to the five expectiles above, i.e.,  $\alpha = 0.19, 0.33, 0.5, 0.66$  and  $0.8$ .

Table 7 displays the test error and variable selection results for cases of  $p = 14$  and  $p = 27$ . In the table, TE is the average expectile check loss on the testing data. Size refers to the average number of predictors selected in the final model. For the penalized expectile regression ET and ES, both SCAD and adaptive LASSO produce similar test errors. Although they are not directly comparable, both penalized expectile regression methods yield smaller test errors than penalized quantile regression across different segments of the distribution. In terms of variable selection, for both SCAD and adaptive LASSO penalties, ET tends to produce sparser models than ES and QT for  $p = 14$ . When  $p = 27$ , ALASSO QT produces sparser models than the others except at  $\tau = 0.1$ , and SCAD ET tends to produce larger models.

Figure 1 shows the estimated coefficients by the proposed penalized expectile regression with both ALASSO and SCAD when  $p = 14$ . It can be seen that some regression coefficients vary as  $\tau$  changes, for example CHAS, DIS, and TAX. In particular, the estimate for TAX shows an interesting pattern because the estimate is increasing as  $\tau$  increases and is close to zero when  $\tau = 0.5$ . Hence, we can interpret that the property tax rate is not an important variable in the mean function, but is important in other expectiles.

**Table 7** Test errors and variable selection results for real data

$p$	$\tau$	ALASSO			SCAD			
		ET	ES	QT	ET	ES	QT	
14	0.10	TE	0.058(0.001)	0.056(0.001)	0.099(0.002)	0.056(0.001)	0.056(0.001)	0.096(0.002)
		Size	8.83(0.17)	11.81(0.25)	9.96(0.13)	12.26(0.24)	11.26(0.23)	12.48(0.17)
	0.25	TE	0.105(0.002)	0.104(0.002)	0.135(0.004)	0.104(0.002)	0.105(0.002)	0.132(0.004)
		Size	9.90(0.14)	11.92(0.16)	10.31(0.14)	12.15(0.22)	11.70(0.23)	12.15(0.17)
	0.50	TE	0.154(0.004)	0.153(0.004)	0.163(0.005)	0.154(0.004)	0.154(0.004)	0.160(0.005)
		Size	9.64(0.15)	10.95(0.19)	9.91(0.15)	10.83(0.23)	11.15(0.24)	11.69(0.20)
0.75	TE	0.164(0.004)	0.164(0.004)	0.175(0.006)	0.164(0.004)	0.164(0.004)	0.174(0.006)	
	Size	9.17(0.19)	10.34(0.23)	9.06(0.17)	9.98(0.26)	10.55(0.29)	10.40(0.25)	
27	0.90	TE	0.133(0.004)	0.133(0.004)	0.173(0.007)	0.134(0.004)	0.132(0.004)	0.172(0.006)
		Size	8.94(0.21)	10.24(0.27)	8.31(0.20)	9.20(0.29)	9.50(0.36)	9.80(0.31)
	0.10	TE	0.066(0.003)	0.065(0.003)	0.080(0.003)	0.069(0.003)	0.063(0.003)	0.084(0.003)
		Size	10.99(0.62)	12.58(0.71)	13.90(0.31)	14.16(0.77)	13.70(0.62)	17.54(0.57)
	0.25	TE	0.090(0.003)	0.089(0.003)	0.100(0.004)	0.092(0.004)	0.089(0.003)	0.104(0.005)
		Size	16.35(0.53)	17.24(0.49)	15.25(0.26)	20.67(0.55)	17.13(0.54)	18.99(0.55)
0.50	TE	0.108(0.004)	0.108(0.004)	0.115(0.005)	0.110(0.004)	0.111(0.004)	0.114(0.005)	
	Size	17.33(0.32)	19.17(0.39)	15.74(0.32)	21.11(0.54)	19.15(0.55)	20.35(0.59)	
0.75	TE	0.104(0.004)	0.102(0.004)	0.115(0.005)	0.103(0.004)	0.107(0.004)	0.112(0.005)	
	Size	18.16(0.30)	20.19(0.41)	16.32(0.34)	21.73(0.49)	20.74(0.57)	21.57(0.54)	
0.90	TE	0.080(0.003)	0.078(0.003)	0.104(0.005)	0.079(0.003)	0.080(0.003)	0.103(0.004)	
	Size	17.68(0.43)	20.11(0.45)	15.88(0.33)	21.39(0.50)	20.74(0.57)	20.14(0.48)	



**Fig. 1** Regression coefficient estimates at different  $\tau$  values with  $p = 14$  for the proposed penalized exponential regression

## 6 Discussion

In this work, we propose the penalized expectile regression using SCAD and adaptive LASSO penalties and show their oracle theoretical properties introduced by [Fan and Li \(2001\)](#), [Zou \(2006\)](#), and [Wu and Liu \(2009\)](#). Also, we compare the proposed method with the penalized quantile regression in [Wu and Liu \(2009\)](#) through a simulation study and real data analysis, and demonstrate its superior performance in terms of estimation and variable selection when heteroscedasticity is present. However, these advantages come at a cost: (i) although there is one-to-one match between the expectiles and quantiles, the expectiles do not have strong interpretability, and (ii) when data are homogeneous with some outliers rather than heterogeneous, penalized expectile regression is more affected by these outliers, while penalized quantile regression is more robust. The proposed work shows some attractive aspects of penalized expectile regression, and we recommend combining its use with penalized quantile regression.

An important problem in the field of penalized regression is construction of confidence intervals for the model coefficients. As [Javanmard and Montanari \(2014\)](#) point out, it is challenging to derive an exact sampling distribution of the parameter estimators in penalized regression due to the use of the optimization procedure, which is a main obstacle for conducting statistical inference using confidence intervals or hypothesis testing. [Javanmard and Montanari \(2014\)](#) propose an efficient algorithm for obtaining confidence intervals and  $p$  values by constructing a de-biased version of regularized  $M$ -estimators in a high-dimensional setting. [Zhang and Zhang \(2014\)](#) propose a hypothesis testing procedure for high-dimensional data using a low-dimensional projection approach. [Lockhart et al. \(2014\)](#) develop the covariance test statistic to determine the significance of regression coefficients in the sequence of models visited along the LASSO solution path. Also, resampling methods for hypothesis testing have been studied; [Chatterjee and Lahiri \(2010\)](#) apply the residual bootstrap approach to the LASSO estimator, and [Minnier et al. \(2011\)](#) propose a perturbation-based procedure to approximate the distribution of a general class of penalized parameter estimates, which leads to the estimation of the covariance matrix and confidence regions. We suggest statistical inference for penalized quantile and expectile regression as our future work.

**Acknowledgements** This work is part of the first author's dissertation. The third author's research was supported by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (NRF-2013R1A1A2007611).

## Appendix: Proofs of theorems

### Proof of Theorem 1

Following [Wu and Liu \(2009\)](#), it is sufficient to show that for any given  $\delta > 0$ , there exists a large constant  $C$  such that

$$P \left\{ \inf_{\|\mathbf{u}\|=C} R_n(\boldsymbol{\beta}_0 + \mathbf{u}/\sqrt{n}) > R_n(\boldsymbol{\beta}_0) \right\} \geq 1 - \delta. \quad (10)$$

It implies that there exists a local minimizer satisfying  $\|\hat{\beta} - \beta_0\| = O_p(n^{-\frac{1}{2}})$ . Now consider

$$\begin{aligned} &R_n(\beta_0 + \mathbf{u}/\sqrt{n}) - R_n(\beta_0) \\ &= \sum_{i=1}^n \left( \rho_\tau(y_i - \mathbf{x}_i^\top \beta_0 - \mathbf{x}_i^\top \mathbf{u}/\sqrt{n}) - \rho_\tau(y_i - \mathbf{x}_i^\top \beta_0) \right) \\ &\quad + n \sum_{j=1}^p \left( p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|) \right). \end{aligned}$$

Because  $p'_{\lambda_n}(\theta) = \lambda_n \{I(\theta \leq \lambda_n) + \frac{(a\lambda_n - \theta)_+}{(a-1)\lambda_n} I(\theta > \lambda_n)\} \geq 0$  for some  $a > 2$  and  $\theta > 0$ , and  $p_{\lambda_n}(0) = 0$ ,

$$n(p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|)) = n(p_{\lambda_n}(|u_j/\sqrt{n}|) - p_{\lambda_n}(0)) \geq 0$$

for  $j = q + 1, \dots, p$ . Hence,

$$\begin{aligned} &R_n(\beta_0 + \mathbf{u}/\sqrt{n}) - R_n(\beta_0) \tag{11} \\ &\geq \sum_{i=1}^n \left( \rho_\tau(y_i - \mathbf{x}_i^\top \beta_0 - \mathbf{x}_i^\top \mathbf{u}/\sqrt{n}) - \rho_\tau(y_i - \mathbf{x}_i^\top \beta_0) \right) \\ &\quad + n \sum_{j=1}^q \left( p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|) \right). \end{aligned}$$

We first consider the second term on the right-hand side of (11). For  $j = 1, \dots, q$ ,

$$\begin{aligned} &n(p_{\lambda_n}(|\beta_{j0} + u_j/\sqrt{n}|) - p_{\lambda_n}(|\beta_{j0}|)) \\ &= n \left( p'_{\lambda_n}(|\beta_{j0}|) \operatorname{sgn}(\beta_{j0}) \frac{u_j}{\sqrt{n}} + \frac{p''_{\lambda_n}(|\beta_{j0}|)}{2} \left( \frac{u_j}{\sqrt{n}} \right)^2 + o\left( \frac{p''_{\lambda_n}(|\beta_{j0}|)}{n} \right) \right) \\ &= O\left( \sqrt{n} \max_{1 \leq j \leq q} p'_{\lambda_n}(|\beta_{j0}|) + \max_{1 \leq j \leq q} p''_{\lambda_n}(|\beta_{j0}|) \right). \end{aligned}$$

For large  $n$ ,

$$\begin{aligned} p'_{\lambda_n}(|\beta_{j0}|) &= \lambda_n \left( I(|\beta_{j0}| \leq \lambda_n) + \frac{(a\lambda_n - |\beta_{j0}|)_+}{(a-1)\lambda_n} I(|\beta_{j0}| > \lambda_n) \right) \\ &= \frac{(a\lambda_n - |\beta_{j0}|)_+}{a-1} \rightarrow 0 \text{ as } \lambda_n \rightarrow 0, \\ p''_{\lambda_n}(|\beta_{j0}|) &= -\frac{1}{a-1} I(\lambda_n < |\beta_{j0}| < a\lambda_n) \rightarrow 0 \text{ as } \lambda_n \rightarrow 0. \end{aligned}$$

Denote the first and second derivatives of  $\rho_\tau(\epsilon_i - t)$  at  $t = 0$  as follows:

$$g_\tau(\epsilon_i) = \rho'_\tau(\epsilon_i - t) |_{t=0} = -2\tau\epsilon_i I(\epsilon_i \geq 0) - 2(1 - \tau)\epsilon_i I(\epsilon_i < 0),$$

$$h_\tau(\epsilon_i) = \rho_\tau''(\epsilon_i - t) |_{t=0} = 2\tau I(\epsilon_i \geq 0) + 2(1 - \tau)I(\epsilon_i < 0).$$

Then  $E(g_\tau(\epsilon_i)) = 0$ . Denote  $\text{Var}(g_\tau(\epsilon_i)) = \sigma_{g_\tau}^2$ ,  $E(h_\tau(\epsilon_i)) = \mu_{h_\tau} > 0$  and  $\text{Var}(h_\tau(\epsilon_i)) = \sigma_{h_\tau}^2$ ,  $i = 1, \dots, n$ . According to model (3.1),  $\epsilon_i = y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_0$ ,  $i = 1, \dots, n$ . Now we consider the first term on the right-hand side of (11):

$$\begin{aligned} & \sum_{i=1}^n \left( \rho_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_0 - \mathbf{x}_i^\top \mathbf{u} / \sqrt{n}) - \rho_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_0) \right) \\ &= \sum_{i=1}^n \left( g_\tau(\epsilon_i) \frac{\mathbf{x}_i^\top \mathbf{u}}{\sqrt{n}} + \frac{h_\tau(\epsilon_i)}{2} \left( \frac{\mathbf{x}_i^\top \mathbf{u}}{\sqrt{n}} \right)^2 + o\left(\frac{1}{n}\right) \right). \end{aligned}$$

We note that

$$\begin{aligned} \sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_i^\top \mathbf{u}}{\sqrt{n}} &= E \left( \sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_i^\top \mathbf{u}}{\sqrt{n}} \right) + O_p \left( \sqrt{\text{Var} \left( \sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_i^\top \mathbf{u}}{\sqrt{n}} \right)} \right) \\ &= O_p \left( \sqrt{\mathbf{u}^\top \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top}{n} \mathbf{u} \sigma_{g_\tau}^2} \right), \end{aligned}$$

and

$$\begin{aligned} \sum_{i=1}^n \frac{h_\tau(\epsilon_i)}{2} \left( \frac{\mathbf{x}_i^\top \mathbf{u}}{\sqrt{n}} \right)^2 &= \frac{\mu_{h_\tau}}{2} \mathbf{u}^\top \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top}{n} \mathbf{u} + O_p \left( \sqrt{\frac{1}{4} \sum_{i=1}^n \left( \mathbf{u}^\top \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top}{n} \mathbf{u} \right)^2 \sigma_{h_\tau}^2} \right) \\ &= \frac{\mu_{h_\tau}}{2} \mathbf{u}^\top \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top}{n} \mathbf{u} + o_p(1). \end{aligned}$$

Therefore,  $R_n(\boldsymbol{\beta}_0 + \mathbf{u} / \sqrt{n}) - R_n(\boldsymbol{\beta}_0)$  is dominated by  $\frac{\mu_{h_\tau}}{2} \mathbf{u}^\top \frac{\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top}{n} \mathbf{u}$ , for  $\|\mathbf{u}\| = C$ , where  $C$  is sufficiently large. In conclusion, there exists a local minimizer of  $R_n(\boldsymbol{\beta})$ ,  $\hat{\boldsymbol{\beta}}^{(\text{SCAD})}$ , such that  $\|\hat{\boldsymbol{\beta}}^{(\text{SCAD})} - \boldsymbol{\beta}_0\| = O_p(n^{-\frac{1}{2}})$ , if  $\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ .  $\square$

**Proof of Theorem 2**

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

$$\begin{aligned} & R_n((\boldsymbol{\beta}_1^\top, \mathbf{0}^\top)^\top) - R_n((\boldsymbol{\beta}_1^\top, \boldsymbol{\beta}_2^\top)^\top) \\ &= \sum_{i=1}^n \left( \rho_\tau(y_i - \mathbf{x}_{i1}^\top \boldsymbol{\beta}_1) - \rho_\tau(y_i - \mathbf{x}_{i1}^\top \boldsymbol{\beta}_1 - \mathbf{x}_{i2}^\top \boldsymbol{\beta}_2) \right) - n \sum_{j=q+1}^p p_{\lambda_n}(|\boldsymbol{\beta}_j|) \\ &= \sum_{i=1}^n \left( g_\tau(\epsilon_i) \mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) + \frac{h_\tau(\epsilon_i)}{2} (\mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}))^2 + o((\mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}))^2) \right) \end{aligned}$$

$$\begin{aligned}
 & - \sum_{i=1}^n \left( g_{\tau}(\epsilon_i) \mathbf{x}_i^T ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T + \frac{h_{\tau}(\epsilon_i)}{2} \left( \mathbf{x}_{i1}^T ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T \right)^2 \right. \\
 & \left. + o \left( \mathbf{x}_{i1}^T ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T \right)^2 \right) - n \sum_{j=q+1}^p p_{\lambda_n}(|\beta_j|). \tag{12}
 \end{aligned}$$

By Condition 2 and following the proof of Theorem 1,

$$\begin{aligned}
 \sum_{i=1}^n g_{\tau}(\epsilon_i) \mathbf{x}_{i1}^T (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) &= O_p \left( \sqrt{\sqrt{n}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n} \sqrt{n}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) \sigma_{g_{\tau}}^2}} \right) \\
 &= O_p(1), \\
 \sum_{i=1}^n g_{\tau}(\epsilon_i) \mathbf{x}_i^T ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T & \\
 &= O_p \left( \sqrt{\sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T) \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T \sigma_{g_{\tau}}^2}} \right) \\
 &= O_p(1), \\
 \frac{h_{\tau}(\epsilon_i)}{2} \left( \mathbf{x}_{i1}^T (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) \right)^2 &= \frac{\mu_{h_{\tau}}}{2} \sqrt{n}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^T}{n} \sqrt{n}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) + o_p(1) \\
 &= O_p(1), \\
 \frac{h_{\tau}(\epsilon_i)}{2} \left( \mathbf{x}_i^T ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T \right)^2 & \\
 &= \frac{\mu_{h_{\tau}}}{2} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T) \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^T}{n} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^T, \boldsymbol{\beta}_2^T)^T + o_p(1) = O_p(1).
 \end{aligned}$$

Now we consider the last term on the right-hand side of (12). For  $j = q + 1, \dots, p$ ,

$$\begin{aligned}
 p_{\lambda_n}(|\beta_j|) &= \lim_{\theta \rightarrow 0^+} p_{\lambda_n}(\theta) + \lim_{\theta \rightarrow 0^+} p'_{\lambda_n}(\theta) |\beta_j| + o(|\beta_j|) \\
 &= \lambda_n |\beta_j| + o(|\beta_j|).
 \end{aligned}$$

Therefore,  $n \sum_{j=q+1}^p p_{\lambda_n}(|\beta_j|) = n \lambda_n \left( \sum_{j=q+1}^p (|\beta_j| + o(|\beta_j|/\lambda_n)) \right)$ . Because  $\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10} = O_p(n^{-\frac{1}{2}})$ ,  $o(|\beta_j|/\lambda_n) = o\left(\frac{1}{\sqrt{n} \lambda_n}\right)$ . We note that  $\sqrt{n} \lambda_n \rightarrow \infty, n \lambda_n \rightarrow \infty$  as  $n \rightarrow \infty$  and  $R_n((\boldsymbol{\beta}_1^T, \boldsymbol{\theta}^T)^T) - R_n((\boldsymbol{\beta}_1^T, \boldsymbol{\beta}_2^T)^T)$  is dominated by

$$-n \sum_{j=q+1}^p p_{\lambda_n}(|\beta_j|).$$

Consequently,

$$R_n((\beta_1^T, \mathbf{0}^T)^T) - R_n((\beta_1^T, \beta_2^T)^T) \rightarrow -\infty \text{ as } n \rightarrow \infty.$$

This completes the proof of part(a) of the theorem. □

(b) From Theorem 1 and part(a), we know  $\hat{\beta}_1$  is a root- $n$  consistent local minimizer of  $R_n((\beta_1^T, \mathbf{0}^T)^T)$ . Let  $\theta_1 = \sqrt{n}(\beta_1 - \beta_{10})$ , i.e.,  $\beta_1 = \beta_{10} + \theta_1/\sqrt{n}$ . Then

$$\begin{aligned} R_n((\beta_1^T, \mathbf{0}^T)^T) &= \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_{i1}^T \beta_1) + n \sum_{j=1}^q p_{\lambda_n}(|\beta_j|) \\ &= \sum_{i=1}^n \rho_\tau(y_i - \mathbf{x}_{i1}^T \beta_{10} - \mathbf{x}_{i1}^T \theta_1/\sqrt{n}) + n \sum_{j=1}^q p_{\lambda_n}(|\beta_{j0} + \theta_j/\sqrt{n}|) \\ &\triangleq Q_n(\theta_1). \end{aligned}$$

Because  $\hat{\theta}_1 = \sqrt{n}(\hat{\beta}_1^{(SCAD)} - \beta_{10})$  is a local minimizer of  $Q_n(\theta_1)$ ,

$$\frac{\partial Q_n(\theta_1)}{\partial \theta_j} \Big|_{\theta_1 = \hat{\theta}_1} = 0,$$

for  $j = 1, \dots, q$ . Now we decompose the derivative of  $Q_n(\theta_1)$  by parts:

$$\begin{aligned} \rho_\tau(y_i - \mathbf{x}_{i1}^T \beta_{10} - \mathbf{x}_{i1}^T \theta_1/\sqrt{n}) &= \rho_\tau(\epsilon_i) + g_\tau(\epsilon_i) \left( -\frac{\mathbf{x}_{i1}^T \theta_1}{\sqrt{n}} \right) \\ &\quad + \frac{h_\tau(\epsilon_i)}{2} \left( -\frac{\mathbf{x}_{i1}^T \theta_1}{\sqrt{n}} \right)^2 + o(1), \\ \frac{\partial}{\partial \theta_j} \rho_\tau(y_i - \mathbf{x}_{i1}^T \beta_{10} - \mathbf{x}_{i1}^T \theta_1/\sqrt{n}) &= -g_\tau(\epsilon_i) \frac{x_{ij}}{\sqrt{n}} + h_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}^T \theta_1}{n} x_{ij}, \\ p_{\lambda_n}(|\beta_{j0} + \theta_j/\sqrt{n}|) &= p_{\lambda_n}(|\beta_{j0}|) + p'_{\lambda_n}(|\beta_{j0}|) \text{sgn}(\beta_{j0}) \frac{\theta_j}{\sqrt{n}} \\ &\quad + \frac{p''_{\lambda_n}(|\beta_{j0}|)}{2} \left( \frac{\theta_j}{\sqrt{n}} \right)^2 + o\left(\frac{1}{n}\right). \end{aligned}$$

Therefore, as  $n \rightarrow \infty$ ,

$$n \frac{\partial}{\partial \theta_j} p_{\lambda_n}(|\beta_{j0} + \theta_j/\sqrt{n}|) = n \left( p'_{\lambda_n}(|\beta_{j0}|) \text{sgn}(\beta_{j0}) \frac{1}{\sqrt{n}} + p''_{\lambda_n}(|\beta_{j0}|) \frac{\theta_j}{n} \right) \rightarrow 0. \tag{13}$$

From the proof of Theorem 1, (13) holds. Plugging them in  $\frac{\partial Q_n(\boldsymbol{\theta}_1)}{\partial \theta_j} \Big|_{\boldsymbol{\theta}_1 = \hat{\boldsymbol{\theta}}_1} = 0$ , for  $j = 1, \dots, q$ , we have

$$\begin{aligned} 0 &= \sum_{i=1}^n \left( -g_\tau(\epsilon_i) \frac{x_{ij}}{\sqrt{n}} + h_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}^\top \hat{\boldsymbol{\theta}}_1}{n} x_{ij} \right) \\ &\quad + \sum_{j=1}^q n \left( p'_{\lambda_n}(|\beta_{j0}|) \text{sgn}(\beta_{j0}) \frac{1}{\sqrt{n}} + p''_{\lambda_n}(|\beta_{j0}|) \frac{\hat{\theta}_j}{n} \right), \\ \mu_{h_\tau} \sum_{i=1}^n \frac{\mathbf{x}_{i1}^\top \hat{\boldsymbol{\theta}}_1}{n} x_{ij} &= \sum_{i=1}^n \left( g_\tau(\epsilon_i) \frac{x_{ij}}{\sqrt{n}} - \frac{(h_\tau(\epsilon_i) - \mu_{h_\tau}) \mathbf{x}_{i1}^\top \hat{\boldsymbol{\theta}}_1}{n} x_{ij} \right) \\ &\quad - \sum_{j=1}^q n \left( p'_{\lambda_n}(|\beta_{j0}|) \text{sgn}(\beta_{j0}) \frac{1}{\sqrt{n}} + p''_{\lambda_n}(|\beta_{j0}|) \frac{\hat{\theta}_j}{n} \right), \\ \mu_{h_\tau} \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \hat{\boldsymbol{\theta}}_1 &= \sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}} - \sum_{i=1}^n \frac{(h_\tau(\epsilon_i) - \mu_{h_\tau}) \mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \hat{\boldsymbol{\theta}}_1 - \sum_{j=1}^q \mathbf{m}_{\lambda_n}(\hat{\boldsymbol{\theta}}_1, \boldsymbol{\beta}_{10}), \\ \hat{\boldsymbol{\theta}}_1 &= \left( \mu_{h_\tau} \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \right)^{-1} \left( \sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}} - \sum_{i=1}^n \frac{(h_\tau(\epsilon_i) - \mu_{h_\tau}) \mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \hat{\boldsymbol{\theta}}_1 \right. \\ &\quad \left. - \sum_{j=1}^q \mathbf{m}_{\lambda_n}(\hat{\boldsymbol{\theta}}_1, \boldsymbol{\beta}_{10}) \right), \end{aligned}$$

where  $\mathbf{m}_{\lambda_n}(\hat{\boldsymbol{\theta}}_1, \boldsymbol{\beta}_{10})$  is defined as a  $q$ -dimensional vector with the  $j$ th element  $n(p'_{\lambda_n}(|\beta_{j0}|) \text{sgn}(\beta_{j0}) \frac{1}{\sqrt{n}} + p''_{\lambda_n}(|\beta_{j0}|) \frac{\hat{\theta}_j}{n})$ . According to (13) and Condition 2, as  $n \rightarrow \infty$ ,  $\mu_{h_\tau} \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \rightarrow \mu_{h_\tau} \Sigma_{11}$ ,  $\sum_{i=1}^n \frac{(h_\tau(\epsilon_i) - \mu_{h_\tau}) \mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \rightarrow 0$ , and  $\mathbf{m}_{\lambda_n}(\hat{\boldsymbol{\theta}}_1, \boldsymbol{\beta}_{10}) \rightarrow 0$ . In addition,  $E\left(g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}}\right) = \mathbf{0}$ ,  $i = 1, \dots, n$ ,

$$\begin{aligned} \text{Var}\left(\sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}}\right) &= \sigma_{g_\tau}^2 \frac{\sum_{i=1}^n \mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \rightarrow \sigma_{g_\tau}^2 \Sigma_{11}, \\ \sum_{i=1}^n E\left(\|g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}}\|^2 I\left(\|g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}}\| > \xi\right)\right) &\leq \sum_{i=1}^n \frac{E\|g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}}\|^4}{\xi^2} \\ &= \frac{1}{\xi^2} E\left(g_\tau^4(\epsilon_i)\right) \sum_{i=1}^n \left(\frac{\mathbf{x}_{i1}^\top \mathbf{x}_{i1}}{n}\right)^2 \rightarrow 0, \end{aligned}$$

for any  $\xi > 0$ . Applying Lindeberg–Feller CLT, we have

$$\sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}} \xrightarrow{\mathcal{L}} \mathbf{w}_1 \sim N(\mathbf{0}, \sigma_{g_\tau}^2 \Sigma_{11}).$$

By Slutsky’s theorem,  $\hat{\boldsymbol{\theta}}_1 \xrightarrow{\mathcal{L}} \left(\mu_{h_\tau} \Sigma_{11}\right)^{-1} \mathbf{w}_1$ . Then, we can conclude,

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_1^{(\text{SCAD})} - \boldsymbol{\beta}_{10}) \xrightarrow{\mathcal{L}} N(\mathbf{0}, \sigma_{g_\tau}^2 / \mu_{h_\tau}^2 \Sigma_{11}^{-1}).$$

This completes the proof. □

**Proof of Theorem 3**

We first prove the asymptotic normality in part (b). Let  $\boldsymbol{\theta} = \sqrt{n}(\boldsymbol{\beta} - \boldsymbol{\beta}_0)$ . Then, we have

$$\begin{aligned} V_n(\boldsymbol{\theta}) &\triangleq R_n(\boldsymbol{\beta}_0 + \boldsymbol{\theta}/\sqrt{n}) - R_n(\boldsymbol{\beta}_0) \\ &= \sum_{i=1}^n \left( g_\tau(\epsilon_i) \left( -\frac{\mathbf{x}_i^\top \boldsymbol{\theta}}{\sqrt{n}} \right) + \frac{h_\tau(\epsilon_i)}{2} \left( -\frac{\mathbf{x}_i^\top \boldsymbol{\theta}}{\sqrt{n}} \right)^2 + o\left(\frac{1}{n}\right) \right) \\ &\quad + n\lambda_n \sum_{j=1}^p w_j (|\beta_{j0} + \theta_j/\sqrt{n}| - |\beta_{j0}|) \\ &= \sum_{i=1}^n g_\tau(\epsilon_i) \left( -\frac{\mathbf{x}_i^\top \boldsymbol{\theta}}{\sqrt{n}} \right) + \frac{\mu_{h_\tau}}{2} \boldsymbol{\theta}^\top \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{n} \boldsymbol{\theta} + \frac{1}{2} \boldsymbol{\theta}^\top \sum_{i=1}^n \left( \frac{(h_\tau(\epsilon_i) - \mu_{h_\tau}) \mathbf{x}_i \mathbf{x}_i^\top}{n} \right) \boldsymbol{\theta} \\ &\quad + o_p(1) + n\lambda_n \sum_{j=1}^p w_j (|\beta_{j0} + \theta_j/\sqrt{n}| - |\beta_{j0}|). \end{aligned} \tag{14}$$

From the proof of Theorem 2,

$$\begin{aligned} \sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_i}{\sqrt{n}} &\xrightarrow{\mathcal{L}} \mathbf{w} \sim N(\mathbf{0}, \sigma_{g_\tau}^2 \Sigma), \\ \frac{\mu_{h_\tau}}{2} \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{n} &\rightarrow \frac{\mu_{h_\tau}}{2} \Sigma, \\ \frac{1}{2} \sum_{i=1}^n \frac{(h_\tau(\epsilon_i) - \mu_{h_\tau}) \mathbf{x}_i \mathbf{x}_i^\top}{n} &\rightarrow 0. \end{aligned}$$

Now we consider the last term of (14). For  $1 \leq j \leq q$ ,

$$w_j \xrightarrow{\mathcal{P}} |\beta_{j0}|^{-\gamma}, \sqrt{n}(|\beta_{j0} + \theta_j/\sqrt{n}| - |\beta_{j0}|) \rightarrow \theta_j \operatorname{sgn}(\beta_{j0}).$$

By Slutsky’s theorem,  $n\lambda_n \sum_{j=1}^p w_j (|\beta_{j0} + \theta_j/\sqrt{n}| - |\beta_{j0}|) \xrightarrow{\mathcal{P}} 0$  because  $\sqrt{n}\lambda_n \rightarrow 0$  as  $n \rightarrow \infty$ . For  $q + 1 \leq j \leq p$ ,  $\beta_{j0} = 0$  and

$$\sqrt{n}(|\beta_{j0} + \theta_j/\sqrt{n}| - |\beta_{j0}|) = |\theta_j|, \sqrt{n}\lambda_n w_j = \lambda_n n^{(\gamma+1)/2} (|\sqrt{n}\tilde{\beta}_j|)^{-\gamma},$$

where  $\tilde{\beta}_j$  is the  $j$ th element of  $\tilde{\boldsymbol{\beta}}$  defined in (3.5) and  $\sqrt{n}\tilde{\beta}_j = O_p(1)$ . Therefore

$$n\lambda_n \sum_{j=1}^p w_j (|\beta_{j0} + \theta_j/\sqrt{n}| - |\beta_{j0}|) \begin{cases} \xrightarrow{\mathcal{P}} \infty & \text{if } \theta_j \neq 0, \\ = 0 & \text{if } \theta_j = 0. \end{cases}$$

Applying Slutsky’s theorem again, we have  $V_n(\boldsymbol{\theta}) \xrightarrow{\mathcal{L}} V(\boldsymbol{\theta})$  for every  $\boldsymbol{\theta}$ . Here,

$$V(\boldsymbol{\theta}) = \begin{cases} \frac{\mu_{h_\tau}}{2} \boldsymbol{\theta}_1^\top \Sigma_{11} \boldsymbol{\theta}_1 + \mathbf{w}_1^\top \boldsymbol{\theta}_1 & \text{if } \theta_j = 0, q + 1 \leq j \leq p, \\ \infty & \text{otherwise.} \end{cases}$$

where  $\mathbf{w}_1 = (w_1, w_2, \dots, w_q)^\top \sim N(\mathbf{0}, \sigma_{g_\tau}^2 \Sigma_{11})$  and  $\boldsymbol{\theta}_1 = (\theta_1, \theta_2, \dots, \theta_q)^\top$ . We note that  $V_n(\boldsymbol{\theta})$  is convex and the unique minimum of  $V(\boldsymbol{\theta})$  is

$$((-\mu_{h_\tau} \Sigma_{11})^{-1} \mathbf{w}_1)^\top, \mathbf{0}^\top)^\top.$$

With the epi-convergence results of Geyer (1994) and Knight and Fu (2000), we have

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_1^{(AL)} - \boldsymbol{\beta}_{10}) = \hat{\boldsymbol{\theta}}_1 \xrightarrow{\mathcal{L}} -(b\Sigma_{11})^{-1} \mathbf{w}_1 \sim N(\mathbf{0}, \sigma_{g_\tau}^2 / \mu_{h_\tau}^2 \Sigma_{11}^{-1})$$

and

$$\sqrt{n}(\hat{\boldsymbol{\beta}}_2^{(AL)} - \boldsymbol{\beta}_{20}) = \hat{\boldsymbol{\theta}}_2 \xrightarrow{\mathcal{L}} \mathbf{0}$$

where  $\hat{\boldsymbol{\theta}}_2 = (\hat{\theta}_{q+1}, \hat{\theta}_{q+2}, \dots, \hat{\theta}_p)^\top$ , which proves the asymptotic normality property.

Next, we show the sparsity property. For any  $\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10} = O_p(n^{-\frac{1}{2}})$ ,  $0 < \|\boldsymbol{\beta}_2\| \leq Cn^{-\frac{1}{2}}$ , following the proof of Theorem 2, we have

$$\begin{aligned} & R_n((\boldsymbol{\beta}_1^\top, \mathbf{0}^\top)^\top) - R_n((\boldsymbol{\beta}_1^\top, \boldsymbol{\beta}_2^\top)^\top) \\ &= \sum_{i=1}^n \left( g_\tau(\epsilon_i) \mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) + \frac{h_\tau(\epsilon_i)}{2} (\mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}))^2 + o((\mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}))^2) \right) \\ &\quad - \sum_{i=1}^n \left( g_\tau(\epsilon_i) \mathbf{x}_{i1}^\top ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top + \frac{h_\tau(\epsilon_i)}{2} (\mathbf{x}_{i1}^\top ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top)^2 \right) \\ &\quad + o_p\left( (\mathbf{x}_{i1}^\top ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top)^2 \right) - n\lambda_n \sum_{j=q+1}^p w_j (|\beta_j|). \end{aligned} \tag{15}$$

The first two terms are bounded in the same way as the proof of Theorem 2:

$$\sum_{i=1}^n g_\tau(\epsilon_i) \mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) = O_p \left( \sqrt{\sqrt{n}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \sqrt{n}(\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) \sigma_{g_\tau}^2} \right)$$

$$\begin{aligned}
 &= O_p(1), \\
 &\sum_{i=1}^n g_\tau(\epsilon_i) \mathbf{x}_i^\top ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top \\
 &= O_p \left( \sqrt{\sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{n} \sqrt{n}((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top \sigma_{g_\tau}^2}} \right) = O_p(1), \\
 &\frac{h_\tau(\epsilon_i)}{2} (\mathbf{x}_{i1}^\top (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}))^2 \\
 &= \frac{\mu_{h_\tau}}{2} \sqrt{n} (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top \sum_{i=1}^n \frac{\mathbf{x}_{i1} \mathbf{x}_{i1}^\top}{n} \sqrt{n} (\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10}) + o_p(1) = O_p(1), \\
 &\frac{h_\tau(\epsilon_i)}{2} (\mathbf{x}_i^\top ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top)^2 \\
 &= \frac{\mu_{h_\tau}}{2} \sqrt{n} ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top \sum_{i=1}^n \frac{\mathbf{x}_i \mathbf{x}_i^\top}{n} \sqrt{n} ((\boldsymbol{\beta}_1 - \boldsymbol{\beta}_{10})^\top, \boldsymbol{\beta}_2^\top)^\top + o_p(1) = O_p(1).
 \end{aligned}$$

For the third term on the right-hand side of (15),

$$n \lambda_n \sum_{j=q+1}^p w_j |\beta_j| = n^{(\gamma+1)/2} \lambda_n \sqrt{n} \sum_{j=q+1}^p (\sqrt{n} \tilde{\beta}_j)^{-\gamma} |\beta_j| \rightarrow \infty,$$

because  $\sqrt{n} \tilde{\beta}_j = O_p(1)$  and  $n^{(\gamma+1)/2} \lambda_n \rightarrow \infty$ . Therefore,

$$R_n((\boldsymbol{\beta}_1^\top, \mathbf{0}^\top)^\top) - R_n((\boldsymbol{\beta}_1^\top, \boldsymbol{\beta}_2^\top)^\top) \rightarrow -\infty \text{ as } n \rightarrow \infty.$$

This implies  $\hat{\boldsymbol{\beta}}_2^{(AL)} = \mathbf{0}$ . □

**Proof of Corollary 1**

From the proof of Theorem 2, it can be shown that

$$\sum_{i=1}^n g_\tau(\epsilon_i) \frac{\mathbf{x}_{i1}}{\sqrt{n}} \xrightarrow{\mathcal{L}} \mathbf{w}_1 \sim N(\mathbf{0}, \Sigma_{11}^g),$$

and  $\hat{\boldsymbol{\theta}}_1 \xrightarrow{\mathcal{L}} (\Sigma_{11}^h)^{-1} \mathbf{w}_1$ . □

**References**

Aigner, D., Amemiya, T., Poirier, D. (1976). On the estimation of production frontiers: Maximum likelihood estimation of the parameters of a discontinuous density function. *International Economic Review*, 17, 377–396.

- Belloni, A., Chernozhukov, V. (2011).  $\ell_1$ -penalized quantile regression in high-dimensional sparse models. *The Annals of Statistics*, 39, 82–130.
- Belloni, A., Chernozhukov, V., Kato, K. (2015). Uniform post-selection inference for least absolute deviation regression and other z-estimation problems. *Biometrika*, 102, 77–94.
- Chatterjee, A., Lahiri, S. N. (2010). Asymptotic properties of the residual bootstrap for Lasso estimators. *Proceedings of the American Mathematical Society*, 138, 4497–4509.
- Fan, J., Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association*, 96, 1348–1360.
- Friberg, H. A. (2014). *Rmosek: The r-to-mosek optimization interface*. r package version 1.2.5.1.
- Geyer, C. J. (1994). On the asymptotics of constrained M-estimation. *The Annals of Statistics*, 22, 1993–2010.
- Gu, Y., Zou, H. (2016). High-dimensional generalizations of asymmetric least squares regression and their applications. *The Annals of Statistics*, 44, 2661–2694.
- Harrison, D., Rubinfeld, D. L. (1978). Hedonic housing prices and the demand for clean air. *Environmental Economics and Management*, 5, 81–102.
- Javanmard, A., Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research*, 15, 2869–2909.
- Jones, M. C. (1994). Expectiles and m-quantiles are quantiles. *Statistics & Probability Letters*, 20, 149–153.
- Kim, Y., Choi, H., Oh, H. S. (2008). Smoothly clipped absolute deviation on high dimensions. *Journal of the American Statistical Association*, 103, 1665–1673.
- Knight, K., Fu, W. (2000). Asymptotics for lasso-type estimators. *The Annals of Statistics*, 28, 1356–1378.
- Kocherginsky, M., He, X., Mu, Y. (2005). Practical confidence intervals for regression quantiles. *Journal of Computational and Graphical Statistics*, 14(1), 41–55.
- Koenker, R., Bassett, G. (1978). Regression quantiles. *Econometrica*, 46, 33–50.
- Koenker, R., Mizera, I. (2014). Convex optimization in R. *Journal of Statistical Software*, 60, 1–23.
- Kuan, C. M., Yeh, J. H., Hsu, Y. C. (2009). Assessing value at risk with care, the conditional autoregressive expectile models. *Journal of Econometrics*, 150, 261–270.
- Li, Y., Zhu, J. (2008).  $l_1$ -norm quantile regression. *Journal of Computational and Graphical Statistics*, 17, 1–23.
- Lockhart, R., Taylor, J., Tibshirani, R. J., Tibshirani, R. (2014). A significance test for the Lasso. *The Annals of Statistics*, 42, 413–468.
- Minnier, J., Tian, T., Cai, T. (2011). A perturbation method for inference on regularized regression estimates. *Journal of the American Statistical Association*, 106, 1371–1382.
- MOSEK ApS D. (2011). *The mosek optimization tool manual. version 7.0*. <https://www.mosek.com>.
- Newey, W. K., Powell, J. L. (1987). Asymmetric least squares estimation and testing. *Econometrica*, 55, 819–847.
- Schnabel, S. K., Eilers, P. H. C. (2009). Optimal expectile smoothing. *Computational Statistics and Data Analysis*, 53, 4168–4177.
- Sobotka, F., Radice, R., Marra, G., Kneib, T. (2013a). Estimating the relationship between women's education and fertility in Botswana by using an instrumental variable approach to semiparametric expectile regression. *Journal of the Royal Statistical Society, Series B*, 62, 25–45.
- Sobotka, F., Radice, R., Marra, G., Kneib, T. (2013b). On confidence intervals for semiparametric expectile regression. *Statistics and Computing*, 23, 135–148.
- Wang, L., Wu, Y., Li, R. (2012). Quantile regression for analyzing heterogeneity in ultra-high dimension. *Journal of the American Statistical Association*, 107, 214–222.
- Wu, Y., Liu, Y. (2009). Variable selection in quantile regression. *Statistica Sinica*, 19, 801–817.
- Yuille, A. L., Rangarajan, A. (2003). The concave–convex procedure. *Neural Computation*, 15, 915–936.
- Zhang, C. H., Zhang, S. S. (2014). Confidence intervals for low-dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B*, 76, 217–242.
- Ziegel, J. (2014). Coherence and elicibility. *Mathematical Finance*, 26, 901–918.
- Zou, H. (2006). The adaptive lasso and its oracle properties. *Journal of the American Statistical Association*, 101, 1418–1429.
- Zou, H., Li, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models. *The Annals of Statistics*, 36, 1509–1533.