

1 **HIGH-DIMENSIONAL GENERALIZATIONS OF ASYMMETRIC**
2 **LEAST SQUARES REGRESSION AND THEIR APPLICATIONS**

3 BY YUWEN GU AND HUI ZOU*

4 *University of Minnesota*

5 Asymmetric least squares regression is an important method that has
6 wide applications in statistics, econometrics and finance. The existing work
7 on asymmetric least squares only considers the traditional low dimension
8 and large sample setting. In this paper, we systematically study the Sparse
9 Asymmetric LEast Squares (SALES) regression under high dimensions where
10 the penalty functions include the Lasso and nonconvex penalties. We develop a
11 unified efficient algorithm for fitting SALES and establish its theoretical prop-
12 erties. As an important application, SALES is used to detect heteroscedasticity
13 in high-dimensional data. Another method for detecting heteroscedasticity is
14 the sparse quantile regression. However, both SALES and the sparse quantile
15 regression may fail to tell which variables are important for the conditional
16 mean and which variables are important for the conditional scale/variance,
17 especially when there are variables that are important for both the mean and
18 the scale. To that end, we further propose a COupled Sparse Asymmetric
19 LEast Squares (COSALES) regression which can be efficiently solved by an
20 algorithm similar to that for solving SALES. We establish theoretical proper-
21 ties of COSALES. In particular, COSALES using the SCAD penalty or MCP
22 is shown to consistently identify the two important subsets for the mean and
23 scale simultaneously, even when the two subsets overlap. We demonstrate the
24 empirical performance of SALES and COSALES by simulated and real data.

25 **1. Introduction.** High-dimensional data have received tremendous attention in
26 the last decade due to the advance of data collection technology. Sparse estimation,
27 which uses penalization or regularization techniques to perform variable selection
28 and estimation simultaneously, has become a mainstream approach for analyzing
29 high-dimensional data. Popular penalized estimators include the L_1 -type selectors
30 such as the Lasso (Tibshirani, 1996) and Dantzig (Candes and Tao, 2007) selectors
31 and the nonconvex penalized estimators such as the SCAD (Fan and Li, 2001)
32 and MCP (Zhang, 2010) estimators. Some embrace the L_1 -regularization for its
33 computational efficiency, while others prefer to use the nonconvex penalization due
34 to its oracle (Fan and Li, 2001) property.

35 The current literature on sparse estimation often assumes homoscedasticity. For
36 example, the existing theory for the sparse linear regression model is based on the
37 classical linear model assumption in which the mean function is linear and the errors
38 are i.i.d. with zero mean and constant variance. The heteroscedasticity issue is often
39 overlooked for theoretical convenience. However, heteroscedasticity often exists

*Supported a NSF grant DMS1505111

MSC 2010 subject classifications: Primary 62J07

Keywords and phrases: Asymmetric least squares, COSALES, High Dimensions, SALES

1 due to heterogeneity in measurement units or accumulation of outlying observations
2 from numerous sources of inputs. This is particularly relevant with high-dimensional
3 data. For example, in genomics experiments, tens of thousands of genes are often
4 analyzed simultaneously by microarrays and occasional outlying measurements
5 appearing in numerous experimental and data-preprocessing steps can accumulate
6 to form heteroscedasticity in the data obtained therein. These data sets are often of
7 high dimension since only a small number of subjects are available for the study.
8 Several studies on expression quantitative trait loci (eQTLs) (Wang, Wu and Li,
9 2012; Daye, Chen and Li, 2012) confirmed the presence of heteroscedasticity in
10 these high-dimensional data and it was shown that genetic variants have effects on
11 both the mean and the scale (i.e., standard deviation) of gene expression levels. In
12 such scenarios, it is important to incorporate heteroscedasticity to make inference
13 from the limited amount of data. To our knowledge, most existing work on high-
14 dimensional data analysis fails to address the heteroscedasticity issue.

15 The sparse quantile regression was proposed in Wang, Wu and Li (2012) to detect
16 heteroscedasticity in high-dimensional data. Quantile regression (Koenker and
17 Bassett, 1978) is appropriate under heteroscedasticity, because it uses an asymmetric
18 absolute value loss. The key word is “asymmetric,” not the absolute value loss. The
19 absolute value loss is computationally more challenging than the squared error
20 loss. Computational efficiency is always one of the primary considerations in high-
21 dimensional data analysis. This motivates us to study the asymmetric least squares
22 (ALS) regression under high dimensionality. The ALS regression has been studied
23 in Efron (1991). It is also known as the expectile regression in econometrics and
24 finance. See Newey and Powell (1987); Taylor (2008); Kuan, Yeh and Hsu (2009);
25 Xie, Zhou and Wan (2014). The key idea in ALS is to assign different squared error
26 loss to the positive and negative residuals, respectively. By doing so, one can infer a
27 more complete description of the conditional distribution than ordinary least squares
28 (OLS). Thus, ALS and quantile regression share a common virtue although they
29 differ technically. The most notable advantage of ALS over quantile regression is
30 that the former employs a smooth differentiable loss, which considerably alleviates
31 the computational effort involved and also makes the theoretical analysis more
32 amenable. These two are desirable properties under high dimensionality.

33 In this paper, we develop the methodology and theory for the Sparse Asym-
34 metric LEast Squares (SALES) regression and show its applications in detecting
35 heteroscedasticity in a general class of sparse models in which the set of relevant
36 covariates may vary from segment to segment on the conditional distribution. For
37 the nonconvex penalized SALES regression, we prove its strong oracle property.
38 We then discuss an important issue overlooked by existing methods dealing with
39 heteroscedasticity in high dimensional data, that is, how to exactly differentiate
40 the sets of relevant covariates for the mean and scale when they have overlaps. To

1 resolve this issue, we propose a novel COupled Sparse Asymmetric LEast Squares
 2 (COSALES) regression method to select important variables for the mean and scale
 3 of the conditional distribution simultaneously. The strong oracle property is also
 4 shown for the nonconvex penalized COSALES estimator. We develop novel efficient
 5 algorithms for computing both SALES and COSALES.

6 The remainder of the article is organized as follows. We study SALES in Section 2
 7 and demonstrate its application in detecting heteroscedasticity in Section 3. In
 8 Section 4, we introduce and study COSALES. The performance of COSALES
 9 is illustrated by two simulation examples. In Section 5, we apply SALES and
 10 COSALES to analyze a real microarray dataset. The proofs of all main theoretical
 11 results are relegated to Section 6.

12 2. High-Dimensional SALES Regression.

13 2.1. *Background and setup.* We start by defining the τ -mean of a random
 14 variable $Z \in \mathbb{R}$,

$$15 \quad (2.1) \quad \mathcal{E}^\tau(Z) \equiv \arg \min_{a \in \mathbb{R}} \mathbb{E}\{\Psi_\tau(Z - a)\}, \quad \tau \in (0, 1),$$

16 where $\Psi_\tau(u) = |\tau - I(u < 0)|u^2$ is the asymmetric squared error loss (see e.g.
 17 Newey and Powell, 1987; Efron, 1991) and $I(\cdot)$ represents the indicator function.
 18 Similar definition can be found in Efron (1991). As a matter of fact, our τ -mean
 19 corresponds to Efron's w -mean, where $w = \tau/(1 - \tau)$. Hereafter, we call \mathcal{E}^τ the
 20 asymmetric expectation operator (with asymmetry coefficient τ). Note that $\mathcal{E}^{0.5}$
 21 coincides with the usual expectation operator \mathbb{E} . The τ -mean is also called the
 22 τ -expectile in the econometrics literature (Newey and Powell, 1987). By varying
 23 τ , the τ -mean quantifies different "locations" of a distribution, and thus it can be
 24 viewed as a generalization of the mean and an alternative measure of "location" of a
 25 distribution.

26 The asymmetric squared error loss $\Psi_\tau(\cdot)$ gives rise to the ALS regression, in
 27 which the squared error loss is given different weights depending on whether
 28 the residual is positive or negative. Let $\mathbf{X} = (X_1, \dots, X_p)$ be the $n \times p$ design
 29 matrix with $X_j = (x_{1j}, \dots, x_{nj})^\top$, $j = 1, \dots, p$, and $\mathbf{y} = (y_1, \dots, y_n)^\top$ be the
 30 n -dimensional response vector. The design matrix may also be written as $\mathbf{X} =$
 31 $(\mathbf{x}_1, \dots, \mathbf{x}_n)^\top$, where $\mathbf{x}_i = (x_{i1}, \dots, x_{ip})^\top$, $i = 1, \dots, n$. The ALS regression is
 32 done via

$$33 \quad \hat{\boldsymbol{\beta}}_\tau^{\text{ALS}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \sum_{i=1}^n \Psi_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}).$$

34 When $\tau = 0.5$, the ALS regression reduces to the OLS regression. When $\tau \neq$
 35 0.5 , due to the asymmetric nature and relative smoothness of $\Psi_\tau(\cdot)$, the ALS

1 regression provides a convenient and computationally efficient way of summarizing
 2 the conditional distribution of a response variable given the covariates (Newey and
 3 Powell, 1987; Efron, 1991). Applications of the ALS regression include estimation
 4 of the value at risk and expected shortfall (Taylor, 2008; Kuan, Yeh and Hsu,
 5 2009), medical baseline correction (Eilers and Boelens, 2005), and small area
 6 estimation (Chambers and Tzavidis, 2006; Salvati et al., 2012) among others.

7 In the literature, the underlying model considered for studying the theoretical
 8 property of the ALS regression is

$$9 \quad (2.2) \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta}^\tau + \boldsymbol{\varepsilon}^\tau,$$

10 where $\boldsymbol{\beta}^\tau$ is a p -dimensional vector of unknown parameters and $\boldsymbol{\varepsilon}^\tau$ is the vector of n
 11 independent errors, which satisfy $\mathcal{E}^\tau(\varepsilon_i^\tau | \mathbf{x}_i) = 0$, $i = 1, \dots, n$ for some $\tau \in (0, 1)$.
 12 It follows that $\mathcal{E}^\tau(y_i | \mathbf{x}_i) = \mathbf{x}_i^\top \boldsymbol{\beta}^\tau$, which means that the conditional τ -mean of y_i is
 13 a linear combination of \mathbf{x}_i , $i = 1, \dots, n$. A similar model to (2.2) was considered
 14 in Wang, Wu and Li (2012) for quantile regression, where the conditional quantile
 15 of the response variable was modeled as a linear combination of the covariates. In
 16 model (2.2), it is important to realize that the coefficient vector $\boldsymbol{\beta}^\tau$ is allowed to
 17 change with τ , which makes modeling for different “locations” of the conditional
 18 distribution possible, and as a result heteroscedasticity in the data, when it exists,
 19 can be inspected by this model. For convenience, we will drop the superscript for
 20 $\boldsymbol{\beta}^\tau$ and $\boldsymbol{\varepsilon}^\tau$ when no confusion arises.

21 To accommodate high-dimensional data in model (2.2), we allow the number
 22 of covariates p to increase with the sample size n , and moreover, we are primarily
 23 interested in cases where p exceeds n ($p > n$). We adopt the sparsity assumption
 24 that only a small number of covariates contribute to the response. Suppose $\boldsymbol{\beta}^* =$
 25 $(\beta_1^*, \dots, \beta_p^*)^\top$ is the parameter vector of the true underlying model that generates the
 26 data and assume $\boldsymbol{\beta}^*$ is s -sparse, where $s = |A|$ with $A \equiv \text{supp}(\boldsymbol{\beta}^*) = \{j: \beta_j^* \neq$
 27 $0\}$.

28 *2.2. Methodology.* To select important variables and estimate $\boldsymbol{\beta}$ in model (2.2)
 29 when the dimension is high, let us consider the following penalized SALES regres-
 30 sion:

$$31 \quad (2.3) \quad \min_{\boldsymbol{\beta} \in \mathbb{R}^p} n^{-1} \sum_{i=1}^n \Psi_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}) + \sum_{j=1}^p p_\lambda(\beta_j),$$

32 where $\Psi_\tau(\cdot)$ is the asymmetric squared error loss and $p_\lambda(\cdot)$ is a nonnegative penalty
 33 function with regularization parameter $\lambda \in (0, \infty)$. In the remainder of this article,
 34 we mainly focus on the Lasso and nonconvex penalties.

2.2.1. *L₁-penalized SALES regression.* For ease of notation, let $\mathcal{L}_n(\boldsymbol{\beta}) = n^{-1} \sum_{i=1}^n \Psi_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta})$. The *L₁-penalized SALES estimator* or SALES Lasso estimator $\hat{\boldsymbol{\beta}}^{\text{lasso}}$ is defined as the solution to the minimization problem

$$(2.4) \quad \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \mathcal{L}_n(\boldsymbol{\beta}) + \lambda_{\text{lasso}} \sum_{j=1}^p |\beta_j|, \quad \lambda_{\text{lasso}} \in (0, \infty).$$

This is to take $p_\lambda(u) = \lambda|u|$ in (2.3). The Lasso is computationally attractive and can be solved by efficient algorithms such as the LARS (Efron et al., 2004), the coordinate descent method (Friedman, Hastie and Tibshirani, 2010) and the generalized coordinate descent algorithm (Yang and Zou, 2013).

For efficient computation of $\hat{\boldsymbol{\beta}}^{\text{lasso}}$ in (2.4), we propose an algorithm called SALES which combines the cyclic coordinate descent (Tseng, 2001) and proximal gradient algorithms (Parikh and Boyd, 2013). Our algorithm solves the following more general “weighted” *L₁-minimization problem*:

$$(2.5) \quad \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \mathcal{L}_n(\boldsymbol{\beta}) + \sum_{j=1}^p w_j |\beta_j|$$

with constants $w_j \geq 0$ for all j . Our consideration of formulation (2.5) is twofold. First, it not only can be directly applied to the SALES Lasso problem (2.4) by setting $w_j = \lambda_{\text{lasso}}$ for all j , but also can be used to solve the convex approximations to the nonconvex penalized SALES estimation (see step (a) of Algorithm 2). Second, leaving some coefficients unpenalized is simply a matter of setting their corresponding weights to zero. Doing so gives us the flexibility to decide which covariates should always be kept in the model. The algorithm is described as follows.

For $\mathbf{v} = (v_1, \dots, v_d)^\top \in \mathbb{R}^d$, denote $\mathbf{v}_{-k} = (v_1, \dots, v_{k-1}, v_{k+1}, \dots, v_d)^\top$ the subvector of \mathbf{v} with its k th component removed. Recover \mathbf{v} from \mathbf{v}_{-k} by writing $\mathbf{v} = [v_k, \mathbf{v}_{-k}]$. Let $\boldsymbol{\beta}^r = (\beta_1^r, \dots, \beta_p^r)^\top$ be the update of $\boldsymbol{\beta}$ after the r th ($r \geq 0$) cycle of the coordinate descent algorithm. For ease of notation, denote

$$\mathbf{b}_{-k}^{r+1} = (\beta_1^{r+1}, \dots, \beta_{k-1}^{r+1}, \beta_{k+1}^r, \dots, \beta_p^r)^\top, \quad 1 \leq k \leq p, \quad r \geq 0.$$

Applying the coordinate descent method, to update β_k in the $(r+1)$ th cycle, we solve the following minimization problem:

$$(2.6) \quad \min_{\beta_k \in \mathbb{R}} \ell_n(\beta_k; \mathbf{b}_{-k}^{r+1}) + w_k |\beta_k|,$$

where $\ell_n(\beta_k; \mathbf{b}_{-k}^{r+1}) = \mathcal{L}_n([\beta_k, \mathbf{b}_{-k}^{r+1}]) = n^{-1} \sum_{i=1}^n \Psi_\tau(y_i - \mathbf{x}_{i,-k}^\top \mathbf{b}_{-k}^{r+1} - x_{ik} \beta_k)$.

One can show that $\ell'_n(\beta_k; \mathbf{b}_{-k}^{r+1})$ is Lipschitz continuous with constant $L_k =$

1 $2\bar{c}n^{-1}\|X_k\|_2^2$, where $\|\cdot\|_2$ is the Euclidean norm. Thus, the proximal gradient
2 method can be employed to solve problem (2.6)

$$3 \quad (2.7) \quad \beta_k^{r,0} := \beta_k^r, \quad \beta_k^{r,s+1} := \mathbb{S}_{L_k^{-1}w_k}(\beta_k^{r,s} - L_k^{-1}\ell'_n(\beta_k^{r,s}; \mathbf{b}_{-k}^{r+1})), \quad s \geq 0,$$

4 where $\mathbb{S}_v(u) = \text{sgn}(u)(|u| - v)^+$ denotes the soft thresholding operator with
5 $u^+ = uI(u > 0)$. We let (2.7) run for s_k^r iterations and set $\beta_k^{r+1} := \beta_k^{r,s_k^r}$. Our
6 algorithm is summarized in Algorithm 1. We prove in Gu and Zou (2015) that
7 Algorithm 1 converges at least linearly.

Algorithm 1: SALES — The cyclic coordinate descent plus proximal gradient algorithm for solving the weighted L_1 -minimization problem (2.5)

1. Initialize the algorithm with $\beta^0 = (\beta_1^0, \dots, \beta_p^0)^\top$.
 2. For $r = 0, 1, 2, \dots, m - 1$,
 - (2.1) For $k = 1, \dots, p$,
 - (2.1.1) Initialize $\beta_k^{r,0} := \beta_k^r$.
 - (2.1.2) For $s = 0, 1, 2, \dots, s_k^r - 1$,
 - (2.1.2.1) Calculate $\beta_k^{r,s+1} := \mathbb{S}_{L_k^{-1}w_k}(\beta_k^{r,s} - L_k^{-1}\ell'_n(\beta_k^{r,s}; \mathbf{b}_{-k}^{r+1}))$.
 - (2.1.3) Set $\beta_k^{r+1} := \beta_k^{r,s_k^r}$.
 - (2.2) Set $\beta^{r+1} := (\beta_1^{r+1}, \dots, \beta_p^{r+1})^\top$.
 3. Output $\hat{\beta} := \beta^m$.
-

8 **2.2.2. Nonconvex penalized SALES regression.** Nonconvex penalties have been
9 used in a broad type of sparse regression models (Fan and Lv, 2011; Wang et al.,
10 2013; Fan, Xue and Zou, 2014). The most popular nonconvex penalties include
11 the smoothly clipped absolute deviation (SCAD) penalty (Fan and Li, 2001) and
12 the minimax concave penalty (MCP, Zhang, 2010). For some constant $\gamma > 2$, the
13 SCAD penalty is given by

$$14 \quad (2.8) \quad p_\lambda(u) = \lambda|u|I(|u| \leq \lambda) + \left\{ \lambda|u| - \frac{(\lambda - |u|)^2}{2(\gamma - 1)} \right\} I(\lambda < |u| \leq \gamma\lambda) \\ + \frac{(\gamma + 1)\lambda^2}{2} I(|u| > \gamma\lambda).$$

15 The use of $\gamma = 3.7$ for the SCAD penalty is recommended in Fan and Li (2001)
16 from a Bayesian perspective. The MCP is characterized by

$$17 \quad (2.9) \quad p_\lambda(u) = \lambda \left(|u| - \frac{u^2}{2\gamma\lambda} \right) I(|u| \leq \gamma\lambda) + \frac{\gamma\lambda^2}{2} I(|u| > \gamma\lambda)$$

1 for some $\gamma > 1$. The use of $\gamma = 2$ is suggested in [Zhang \(2010\)](#). In this article, we
 2 consider both SCAD and MCP penalized SALES regression.

3 The main motivation for using the nonconvex penalties is to achieve the oracle
 4 property. For the SALES regression, the oracle estimator is

$$5 \quad (2.10) \quad \widehat{\boldsymbol{\beta}}^{\text{oracle}} = \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p: \boldsymbol{\beta}_{\text{Ac}} = \mathbf{0}} \mathcal{L}_n(\boldsymbol{\beta}).$$

6 In practice, the oracle estimator is infeasible, but it sets a benchmark for evaluation
 7 of other estimators. Many papers have shown that the nonconvex penalized least
 8 squares can find the oracle estimator with high probability ([Wang et al., 2013](#); [Fan,
 9 Xue and Zou, 2014](#)). In particular, [Fan, Xue and Zou \(2014\)](#) showed that the local
 10 linear approximation (LLA) algorithm ([Zou and Li, 2008](#)) converges to the oracle es-
 11 timator under regularity conditions. The LLA algorithm fits a sequence of weighted
 12 L_1 -regularization problems. Since we already have Algorithm 1 for computing any
 13 weighted L_1 -penalized SALES regression, we adopt the LLA algorithm for solving
 14 the nonconvex penalized SALES estimation problem (2.3). The details of the LLA
 15 algorithm are shown in Algorithm 2. Note that step (a) can be readily solved by
 16 Algorithm 1.

Algorithm 2: The local linear approximation (LLA) algorithm for solving the
 nonconvex penalized SALES estimation problem (2.3)

1. Initialize $\widehat{\boldsymbol{\beta}}^0 := \widehat{\boldsymbol{\beta}}^{\text{initial}}$. Compute weights $\widehat{w}_j^0 = p'_\lambda(|\widehat{\beta}_j^0|)$, $j = 1, \dots, p$.
2. For $m = 1, 2, \dots$, repeat the LLA iteration in (a) and (b) until convergence
 - (a) Solve the following convex optimization problem for $\widehat{\boldsymbol{\beta}}^m$

$$\widehat{\boldsymbol{\beta}}^m := \arg \min_{\boldsymbol{\beta} \in \mathbb{R}^p} \mathcal{L}_n(\boldsymbol{\beta}) + \sum_{j=1}^p \widehat{w}_j^{m-1} |\beta_j|.$$

- (b) Update the weights $\widehat{w}_j^m = p'_\lambda(|\widehat{\beta}_j^m|)$, $j = 1, \dots, p$.
-

17 In our numerical examples, we tried using both the SALES Lasso estimator
 18 and zero as the initial values of the LLA algorithm for computing the noncon-
 19 vex penalized SALES estimator. Our practice is based on theoretical results in
 20 Section 2.3.

21 **2.3. Theory.** In this section, we theoretically analyze the SALES regression.
 22 We consider the case where the covariates are from a fixed design.

23 The following notation will be used. For any vector $\mathbf{v} = (v_1, \dots, v_p)^T \in \mathbb{R}^p$
 24 and an arbitrary index set $I \subset \{1, \dots, p\}$, we write $\mathbf{v}_I = (v_j, j \in I)^T$ and denote

1 by $\mathbf{X}_I = (\mathbf{x}_j, j \in I)$ the submatrix consisting of the columns of \mathbf{X} with indices
 2 in I . The complement of I is denoted by $I^c = \{1, \dots, p\} \setminus I$. For $q \in [1, \infty]$, the
 3 L_q -norm of \mathbf{v} is denoted by $\|\mathbf{v}\|_q$. Sub-Gaussian norm (Rudelson and Vershynin,
 4 2013) of a random variable Z is denoted by $\|Z\|_{\text{SG}} = \sup_{k \geq 1} k^{-1/2} (\mathbb{E}|Z|^k)^{1/k}$.
 5 Let $a \vee b = \max(a, b)$ and $a \wedge b = \min(a, b)$ for real numbers a and b . For a
 6 differentiable function $f : \mathbb{R}^p \rightarrow \mathbb{R}$, we write $\nabla f(\mathbf{v}) = \partial f(\mathbf{v}) / \partial \mathbf{v}$ and $\nabla_I f(\mathbf{v}) =$
 7 $(\partial f(\mathbf{v}) / \partial v_j, j \in I)^T$. We use $\lambda_{\min}(\cdot)$ and $\lambda_{\max}(\cdot)$ to represent respectively the
 8 smallest and largest eigenvalues of a symmetric matrix. We also let $\underline{c} = \tau \wedge (1 - \tau)$
 9 and $\bar{c} = \tau \vee (1 - \tau)$.

10 2.3.1. L_1 -penalized SALES regression. The estimation accuracy of the Lasso
 11 has been extensively studied in the literature; see, for example, Negahban et al.
 12 (2012) and Ye and Zhang (2010). Let $\mathcal{C} = \{\boldsymbol{\delta} \in \mathbb{R}^p : \|\boldsymbol{\delta}_{A^c}\|_1 \leq 3\|\boldsymbol{\delta}_A\|_1 \neq 0\}$ be
 13 a cone in \mathbb{R}^p . Denote $\rho_{\min} = \lambda_{\min}(n^{-1}\mathbf{X}_A^T\mathbf{X}_A)$ and $\rho_{\max} = \lambda_{\max}(n^{-1}\mathbf{X}_A^T\mathbf{X}_A)$.
 14 We assume $\rho_{\min} > 0$ so that the important variables are not linearly dependent.
 15 To study the estimation accuracy of the SALES Lasso, we impose the following
 16 conditions on the design matrix \mathbf{X} and the random errors $\boldsymbol{\varepsilon}$.

17 (C1) The columns of \mathbf{X} are normalizable, that is, $M_0 = \max_{1 \leq j \leq p} \frac{\|X_j\|_2}{\sqrt{n}} \in$
 18 $(0, \infty)$.

19 (C2) The random errors ε_i are i.i.d. sub-Gaussian random variables satisfying
 20 $\mathcal{E}^\tau(\varepsilon_i) = 0, i = 1, \dots, n$.

21 (C3) $\kappa = \inf_{\boldsymbol{\delta} \in \mathcal{C}} \frac{\|\mathbf{X}\boldsymbol{\delta}\|_2^2}{n\|\boldsymbol{\delta}\|_2^2} \in (0, \infty)$.

22 (C4) $\varrho = \inf_{\boldsymbol{\delta} \in \mathcal{C}} \frac{\|\mathbf{X}\boldsymbol{\delta}\|_2^2}{n\|\boldsymbol{\delta}_A\|_1\|\boldsymbol{\delta}\|_\infty} \in (0, \infty)$.

23 Condition (C3) is called the restricted eigenvalue condition and has been fre-
 24 quently assumed in the literature to study the Lasso and Dantzig selectors. See Bickel,
 25 Ritov and Tsybakov (2009), Meier et al. (2009), and Negahban et al. (2012). Con-
 26 dition (C4), the generalized invertability factor (GIF) condition, is closely related
 27 to condition (C3) and has also been often adopted to study the Lasso and Dantzig
 28 selectors. See discussion of these conditions in Ye and Zhang (2010) and Huang and
 29 Zhang (2012). Both conditions (C3) and (C4) are crucial assumptions to establish
 30 estimation consistency of the Lasso for high-dimensional data.

31 THEOREM 1. Suppose in model (2.2) the true coefficients $\boldsymbol{\beta}^*$ are s -sparse and
 32 assume conditions (C1-C2) hold. Let $\widehat{\boldsymbol{\beta}}^{\text{lasso}}$ be any optimal solution to the SALES
 33 Lasso problem (2.4). Then with probability at least $1 - p_1^{\text{ALS}}$, $\|\widehat{\boldsymbol{\beta}}^{\text{lasso}} - \boldsymbol{\beta}^*\|_2 \leq$
 34 $3s^{1/2}\lambda_{\text{lasso}}(4\kappa\underline{c})^{-1}$ if condition (C3) holds, and $\|\widehat{\boldsymbol{\beta}}^{\text{lasso}} - \boldsymbol{\beta}^*\|_\infty \leq 3\lambda_{\text{lasso}}(4\varrho\underline{c})^{-1}$

1 if condition (C4) holds, where

$$2 \quad p_1^{ALS} = 2p \exp\left(-\frac{Cn\lambda_{lasso}^2}{4K_0^2 M_0^2}\right),$$

3 $K_0 = \|\Psi'_\tau(\varepsilon_i)\|_{SG}$ with $\Psi'_\tau(\cdot)$ being the derivative of $\Psi_\tau(\cdot)$, and $C > 0$ is an
4 absolute constant.

5 **REMARK 1.** In some applications, it is natural to leave a given subset of the
6 parameters unpenalized in the penalized framework (2.3). Let \mathcal{R} denote the index
7 set of such parameters. For example, when X_1 is a vector consisting of all ones,
8 $\mathcal{R} = \{1\}$ reflects the common practice of leaving the intercept term not penalized.
9 In this case, it is natural to modify the penalized SALES estimation problem (2.3)
10 to be

$$11 \quad \min_{\beta \in \mathbb{R}^p} \mathcal{L}_n(\beta) + \sum_{j \in \mathcal{R}^c} p_\lambda(\beta_j).$$

12 With Lasso penalty, the SALES algorithm can be readily used to solve the above case.
13 Moreover, similar theoretical analysis can be carried out with slight modifications.
14 For instance, in the SALES Lasso problem (2.4) we can define $A' \equiv \text{supp}(\beta_{\mathcal{R}^c}^*)$
15 and $\mathcal{C}' = \{\delta \in \mathbb{R}^p : \|\delta_{(A' \cup \mathcal{R})^c}\|_1 \leq 3\|\delta_{A' \cup \mathcal{R}}\|_1 \neq 0\}$. Conditions (C3) and (C4)
16 can be then modified respectively as

$$17 \quad \kappa' = \inf_{\delta \in \mathcal{C}'} \frac{\|\mathbf{X}\delta\|_2^2}{n\|\delta\|_2^2} \in (0, \infty) \quad \text{and} \quad \varrho' = \inf_{\delta \in \mathcal{C}'} \frac{\|\mathbf{X}\delta\|_2^2}{n\|\delta_{A' \cup \mathcal{R}}\|_1 \|\delta\|_\infty} \in (0, \infty).$$

18 To establish the selection consistency of the Lasso, it is almost necessary to
19 impose the irrepresentable condition; see [Zou \(2006\)](#) and [Zhao and Yu \(2006\)](#).
20 When the focus is on identifying the underlying sparsity pattern, the nonconvex
21 penalized regression is a competitive alternative as it requires weaker conditions to
22 achieve selection consistency.

23 **2.3.2. Nonconvex penalized SALES regression.** To offer a unified treatment of
24 the SCAD and MCP penalized SALES regression, our theoretical analysis handles
25 the following class of nonconvex penalties:

- 26 (P1) $p_\lambda(u) = p_\lambda(-u)$;
- 27 (P2) $p_\lambda(u)$ is nondecreasing and concave in $u \in [0, \infty)$ and $p_\lambda(0) = 0$;
- 28 (P3) $p_\lambda(u)$ is differentiable in $u \in (0, \infty)$;
- 29 (P4) $p'_\lambda(u) \geq a_1\lambda$ for $u \in (0, a_2\lambda]$ and $p'_\lambda(0) := p'_\lambda(0+) \geq a_1\lambda$;
- 30 (P5) $p'_\lambda(u) = 0$ for $u \in [a\lambda, \infty)$ with some prespecified constant $a > a_2$,

31 where a_1 and a_2 are fixed constants characteristic of the penalty functions. It is easy
32 to verify that both the SCAD penalty and MCP are in the above class.

1 We show that the sparse solutions obtained by the LLA algorithm in Section 2.2.2
 2 possess the oracle property. Assume sufficient signal strength in the nonzero com-
 3 ponents of β^*

$$4 \quad (A1) \quad \min_{j \in A} |\beta_j^*| > (a + 1)\lambda.$$

5 **THEOREM 2.** *Suppose in model (2.2) the true coefficients β^* are s -sparse and*
 6 *satisfy assumption (A1). Assume conditions (C1-C2) hold and take $\widehat{\beta}^{\text{lasso}}$ as the*
 7 *initial value. Let $a_0 = 1 \wedge a_2$. Take $\lambda \geq 3s^{1/2}\lambda_{\text{lasso}}(4a_0\kappa\underline{c})^{-1}$ when (C3) holds, or*
 8 *take $\lambda \geq 3\lambda_{\text{lasso}}(4a_0\underline{\rho}\underline{c})^{-1}$ when (C4) holds, or take $\lambda \geq [3s^{1/2}\lambda_{\text{lasso}}(4a_0\kappa\underline{c})^{-1}] \wedge$*
 9 *$[3\lambda_{\text{lasso}}(4a_0\underline{\rho}\underline{c})^{-1}]$ when both (C3) and (C4) hold. The LLA algorithm (Algorithm 2)*
 10 *converges to $\widehat{\beta}^{\text{oracle}}$ after two iterations with probability at least $1 - p_1^{\text{ALS}} - p_2^{\text{ALS}} -$*
 11 *p_3^{ALS} , where p_1^{ALS} is given in Theorem 1,*

$$12 \quad p_2^{\text{ALS}} = 2(p - s) \exp\left(-\frac{Ca_1^2 n \lambda^2}{4K_0^2 M_0^2}\right) + \Gamma(Q_1 \lambda; n, s, K_0, M_0, \rho_{\max}, \nu_0)$$

13 *and*

$$14 \quad p_3^{\text{ALS}} = \Gamma(2\underline{c}\rho_{\min}R; n, s, K_0, M_0, \rho_{\max}, \nu_0),$$

15 *where $Q_1 = a_1\underline{c}\rho_{\min}(2\underline{c}\rho_{\max}^{1/2}M_0)^{-1}$, $\nu_0 = \text{var}(\Psi'_\tau(\varepsilon_i))$, $R = \min_{j \in A} |\beta_j^*| - a\lambda$,*
 16 *K_0 is defined in Theorem 1 and $\Gamma(\cdot)$ is a function defined by*

$$17 \quad \Gamma(x; n, s, K, M, \rho, \nu) = 2s \exp\left(-\frac{Cnx^2}{K^2M^2s}\right) \\ \wedge 2 \exp\left(-\frac{C\nu^2[(n^{1/2}x - \nu\rho^{1/2}s^{1/2})^+]^2}{K^4\rho}\right),$$

18 *and $C > 0$ is an absolute constant.*

19 It is interesting to note that with the SCAD penalty or MCP, a three-step LLA
 20 algorithm starting from the zero vector may also work. Indeed, for these two
 21 penalties we have $p'_\lambda(0) = \lambda$, so if we can take $\lambda = \lambda_{\text{lasso}}$, this would give us the
 22 SALES Lasso estimator in the second step.

23 **COROLLARY 1.** *Assume the same framework of Theorem 2 and suppose the*
 24 *SCAD penalty (2.8) or MCP (2.9) is used. If condition (C3) holds and $4a_0\kappa\underline{c} \geq$*
 25 *$3s^{1/2}$, or if condition (C4) holds and $4a_0\underline{\rho}\underline{c} \geq 3$, or if both (C3) and (C4) hold and*
 26 *$[3s^{1/2}(\kappa)^{-1}] \wedge [3(\underline{\rho})^{-1}] \leq 4a_0\underline{c}$, the LLA algorithm (Algorithm 2) initialized by*
 27 *zero converges to the oracle estimator after three iterations with probability at least*
 28 *$1 - 2p \exp\{-Cn\lambda^2(4K_0^2M_0^2)^{-1}\} - p_2^{\text{ALS}} - p_3^{\text{ALS}}$, where p_2^{ALS} and p_3^{ALS} are given in*
 29 *Theorem 2.*

1 **3. Application of SALES: detecting heteroscedasticity.** Due to asymmetry
 2 of the squared error loss, the SALES regression (2.3) can be employed to detect
 3 heteroscedasticity in high-dimensional data. In the following, we use a simulation
 4 example to illustrate this application. For the nonconvex penalty functions used
 5 in the simulation, we fix $\gamma = 3.7$ for the SCAD penalty (2.8) and $\gamma = 2$ for the
 6 MCP (2.9).

7 **EXAMPLE 1.** We adopt a model from Wang, Wu and Li (2012). In the model,
 8 the covariates are generated in two steps. First, we generate copies of $(z_1, \dots, z_p)^T$
 9 from the multivariate normal distribution $N(\mathbf{0}, \Sigma)$ with $\Sigma = (0.5^{|i-j|})_{p \times p}$. In the
 10 second step, for each copy of $(z_1, \dots, z_p)^T$, we set $x_1 = \Phi(z_1)$ and $x_j = z_j$ for
 11 $j = 2, 3, \dots, p$, where $\Phi(\cdot)$ is the standard normal CDF. The response is then
 12 simulated from the following normal linear heteroscedastic model:

$$13 \quad (3.1) \quad y = x_6 + x_{12} + x_{15} + x_{20} + (0.7x_1)\varepsilon,$$

14 where $\varepsilon \sim N(0, 1)$ is independent of the covariates. This model was considered
 15 in Wang, Wu and Li (2012) for the sparse quantile regression, where a sample
 16 size $n = 300$ and covariate dimensions $p = 400$ and 600 were considered. We
 17 apply the SALES regression (2.3) instead to select active variables and estimate the
 18 coefficients for this model. For the purpose of demonstration, we choose $n = 300$
 19 and $p = 600$. A validation set of size $n = 300$ is generated independently to tune the
 20 regularization parameter by minimizing the validation error $\sum_{i \in \text{validation}} \Psi_\tau(y_i -$
 21 $\mathbf{x}_i^T \hat{\beta})$ for the computed estimate $\hat{\beta}$, where $\tau = 0.5$ and 0.85 are considered.

22 For comparison purpose, we included in this simulation the SALES Lasso (2.4)
 23 and two variations of the LLA algorithm for each nonconvex penalized SALES
 24 regression: the two-step LLA algorithm initialized by the Lasso estimator (SCAD*,
 25 MCP*), and the three-step LLA algorithm initialized by zero (SCAD⁰, MCP⁰).

26 Let $\hat{\beta}$ be the coefficient estimates from a given method. Based on 100 replicates,
 27 the following measurements are calculated to evaluate the sparsity recovery and
 28 estimation performance of that method:

- 29 $|\hat{A}|$: the average size of the active set $\hat{A} = \{j : \hat{\beta}_j \neq 0\}$ of $\hat{\beta}$.
- 30 p_a : proportion of the event $A \subset \hat{A}$, where A is the active set of β^* . When
 31 $\tau = 0.5$, $A = \{6, 12, 15, 20\}$ and when $\tau \neq 0.5$, $A = \{1, 6, 12, 15, 20\}$.
- 32 p_1 : proportion of the event that $\{1\} \subset \hat{A}$.
- 33 R_1 : the average L_1 risk $\|\hat{\beta} - \beta^*\|_1$.
- 34 R_2 : the average L_2 risk $\|\hat{\beta} - \beta^*\|_2$.

35 The simulation results are shown in Table 1. The following conclusions can be
 36 made:

- 37 (1) The variable x_1 in the scale function is often not recovered by penalized
 38 least-squares ($\tau = 0.5$). However, when several τ -means (e.g., $\tau = 0.85$)

- 1 are inspected together, it is possible to detect this variable with high proba-
 2 bility. This shows that indeed the SALES regression can be used to detect
 3 heteroscedasticity.
- 4 (2) Compared to the SALES Lasso, the nonconvex penalized SALES regression
 5 selects much fewer irrelevant covariates and has better estimation accuracy.
- 6 (3) The three-step LLA algorithm starting from zero produces similar results to
 7 the two-step LLA algorithm starting from the Lasso solution.

TABLE 1

Numerical summary of simulation results from the Lasso, SCAD and MCP penalized SALES regression for model (3.1): $y = x_6 + x_{12} + x_{15} + x_{20} + (0.7x_1)\varepsilon$. The sparsity recovery performance is measured by the selected active set size $|\hat{A}|$, the proportion p_α of covering the true active set and the proportion p_1 of selecting the signature variable X_1 . The estimation accuracy is measured by the L_1 risk R_1 and the L_2 risk R_2 . The results are shown as averages over 100 replicates with standard errors listed in the parentheses when available

	Method	$ \hat{A} $	p_α	p_1	R_1	R_2
$\tau = 0.5$	SALES-Lasso	25.82 (1.15)	100%	0%	0.399 (0.015)	0.120 (0.003)
	SALES-SCAD*	7.75 (0.68)	100%	0%	0.103 (0.006)	0.049 (0.002)
	SALES-SCAD ⁰	6.65 (0.68)	100%	0%	0.100 (0.006)	0.050 (0.002)
	SALES-MCP*	6.39 (0.48)	100%	0%	0.099 (0.005)	0.049 (0.002)
	SALES-MCP ⁰	5.75 (0.29)	100%	0%	0.093 (0.004)	0.049 (0.002)
$\tau = 0.85$	SALES-Lasso	34.17 (1.26)	100%	100%	0.714 (0.016)	0.249 (0.005)
	SALES-SCAD*	7.52 (0.51)	100%	100%	0.160 (0.009)	0.083 (0.005)
	SALES-SCAD ⁰	8.19 (0.59)	100%	100%	0.166 (0.007)	0.084 (0.003)
	SALES-MCP*	6.30 (0.25)	100%	100%	0.148 (0.005)	0.079 (0.003)
	SALES-MCP ⁰	6.35 (0.23)	100%	100%	0.147 (0.005)	0.078 (0.003)

8 **4. High-Dimensional COSALES Regression.** In Section 3, we showed that
 9 the SALES regression provides a means of detecting heteroscedasticity in high-
 10 dimensional data. Indeed, in the linear heteroscedastic model (3.1), the signature
 11 variable x_1 , which appears in the scale function, was detected through comparison
 12 of different τ -means. However, in high-dimensional heteroscedastic models, often
 13 of more interest are the sparsity patterns in both the mean and the scale functions
 14 of the conditional distribution. The SALES regression and methods proposed by
 15 other authors, for example, Wang, Wu and Li (2012), are not sufficient to fulfill this
 16 task. To see it, consider a linear heteroscedastic model in which the active set for
 17 the mean is $\{1, 2\}$ and the active set for the scale is $\{1, 3\}$. Suppose the SALES
 18 regression can exactly recover the active variables. Then the method picks x_1 and
 19 x_2 when $\tau = 0.5$ and hopefully x_1, x_2 , and x_3 when $\tau \neq 0.5$. A natural question is
 20 whether the scale function depends on x_1 . With the SALES regression, we cannot

1 answer this question. This motivates us to consider the COSALES regression for a
 2 general class of models and gain some insight into analyzing heteroscedasticity in
 3 high-dimensional data.

4 4.1. *Formulation and computation.* Consider the following model of systematic
 5 heteroscedasticity,

$$6 \quad (4.1) \quad y_i = \mathbf{x}_i^T \boldsymbol{\gamma} + (\mathbf{x}_i^T \boldsymbol{\omega}) \varepsilon_i, \quad i = 1, \dots, n,$$

7 where ε_i are i.i.d. random errors that are independent of the covariates and that have
 8 distribution F_0 with $\mathbb{E}(\varepsilon_i) = \int_{\mathbb{R}} x dF_0(x) = 0$; $\boldsymbol{\gamma}$ and $\boldsymbol{\omega}$ are unknown p -dimensional
 9 parameter vectors controlling the conditional mean and scale; and $\boldsymbol{\omega}$ is assumed
 10 to satisfy $\mathbf{x}_i^T \boldsymbol{\omega} > 0$ for all i . The intercept can be included by letting $x_{i1} = 1$. The
 11 linear scale model of heteroscedasticity (4.1) is an important model considered by
 12 many authors (Koenker and Bassett, 1982; Efron, 1991; Koenker and Zhao, 1994)
 13 for analyzing heteroscedasticity.

14 Let $A_1 \equiv \text{supp}(\boldsymbol{\gamma}^*) = \{j: \gamma_j^* \neq 0\}$ and $A_2 \equiv \text{supp}(\boldsymbol{\omega}^*) = \{j: \omega_j^* \neq 0\}$ be
 15 the active sets of $\boldsymbol{\gamma}^*$ and of $\boldsymbol{\omega}^*$, respectively. Suppose $|A_1| = s_1$ and $|A_2| = s_2$. Let
 16 $e_\tau = \mathcal{E}^\tau(\varepsilon_1)$ be the τ -mean of the random error for $\tau \in (0, 1)$. It follows that the
 17 τ -mean of y_i given \mathbf{x}_i is $\mathcal{E}^\tau(y_i | \mathbf{x}_i) = \mathbf{x}_i^T (\boldsymbol{\gamma} + \boldsymbol{\omega} e_\tau)$. To select significant variables
 18 in both the mean and the scale functions, we now propose the COSALES regression.
 19 Write $\boldsymbol{\varphi} = \boldsymbol{\omega} e_\tau$. Note that we omit the dependency of $\boldsymbol{\varphi}$ on τ to ease exposition.
 20 In the COSALES regression, we will deal with $\boldsymbol{\varphi}$ instead of $\boldsymbol{\omega}$. However, when
 21 $e_\tau \neq 0$, it should be noted that since $\text{supp}(\boldsymbol{\varphi}) = \text{supp}(\boldsymbol{\omega})$, the selection result on $\boldsymbol{\varphi}$
 22 applies to $\boldsymbol{\omega}$. Moreover, $\boldsymbol{\omega}$ can be estimated up to a scale from the estimate of $\boldsymbol{\varphi}$.
 23 Ideally, if the distribution F_0 of ε_i is known, exact estimation of $\boldsymbol{\omega}$ is possible.

For some $\tau \in (0, 1)$ and $\tau \neq 0.5$, let

$$S_n(\boldsymbol{\gamma}, \boldsymbol{\varphi}) = n^{-1} \sum_{i=1}^n \{\Psi_{0.5}(y_i - \mathbf{x}_i^T \boldsymbol{\gamma}) + \Psi_\tau(y_i - \mathbf{x}_i^T \boldsymbol{\gamma} - \mathbf{x}_i^T \boldsymbol{\varphi})\}.$$

24 The COSALES regression tries to minimize

$$25 \quad (4.2) \quad Q_n(\boldsymbol{\gamma}, \boldsymbol{\varphi}) = S_n(\boldsymbol{\gamma}, \boldsymbol{\varphi}) + \sum_{j=1}^p p_{\lambda_1}(\gamma_j) + \sum_{j=1}^p p_{\lambda_2}(\varphi_j),$$

26 over $\boldsymbol{\gamma}, \boldsymbol{\varphi} \in \mathbb{R}^p$, where $p_{\lambda_1}(\cdot)$ and $p_{\lambda_2}(\cdot)$ are penalty functions with regularization
 27 parameters $\lambda_1, \lambda_2 \in (0, \infty)$, respectively. Let $\hat{\boldsymbol{\gamma}}^{\text{oracle}}$ and $\hat{\boldsymbol{\varphi}}^{\text{oracle}}$ be the oracle
 28 estimators of $\boldsymbol{\gamma}$ and $\boldsymbol{\varphi} = \boldsymbol{\omega} e_\tau$, respectively, in model (4.1),

$$29 \quad (4.3) \quad (\hat{\boldsymbol{\gamma}}^{\text{oracle}}, \hat{\boldsymbol{\varphi}}^{\text{oracle}}) = \arg \min_{\boldsymbol{\gamma}, \boldsymbol{\varphi} \in \mathbb{R}^p: \boldsymbol{\gamma}_{A_1^c} = \mathbf{0}, \boldsymbol{\varphi}_{A_2^c} = \mathbf{0}} S_n(\boldsymbol{\gamma}, \boldsymbol{\varphi}).$$

30 In what follows, let us focus on the Lasso and nonconvex penalties.

1 4.1.1. *L₁-penalized COSALES regression.* For $\lambda_1^{\text{lasso}}, \lambda_2^{\text{lasso}} \in (0, \infty)$, the L_1 -
 2 penalized COSALES estimators or the COSALES Lasso estimators of γ and φ can
 3 be achieved simultaneously by

$$4 \quad (4.4) \quad (\hat{\gamma}^{\text{lasso}}, \hat{\varphi}^{\text{lasso}}) = \arg \min_{\gamma, \varphi \in \mathbb{R}^p} S_n(\gamma, \varphi) + \lambda_1^{\text{lasso}} \|\gamma\|_1 + \lambda_2^{\text{lasso}} \|\varphi\|_1.$$

5 We note that problem (4.4) is a special case of the minimization problem in step (a)
 6 of Algorithm 4 (Section 4.1.2) and efficient computation of the solutions can be
 7 carried out by an algorithm similar to Algorithm 1. The algorithm applies the cyclic
 8 coordinate descent and proximal gradient descent methods to γ and φ alternately.
 9 We call this algorithm COSALES and display it in Algorithm 3. Note that COSALES
 10 solves the general coupled weighted L_1 -minimization problem

$$11 \quad (4.5) \quad \min_{\gamma, \varphi \in \mathbb{R}^p} S_n(\gamma, \varphi) + \sum_{j=1}^p w_j |\gamma_j| + \sum_{j=1}^p v_j |\varphi_j|.$$

12 To facilitate the presentation, in Algorithm 3, we let γ^r and φ^r be the updates of
 13 γ and φ respectively after the r th cycle of the coordinate descent algorithm and
 14 denote

$$15 \quad \mathbf{g}_{-k}^{r+1} = (\gamma_1^{r+1}, \dots, \gamma_{k-1}^{r+1}, \gamma_{k+1}^r, \dots, \gamma_p^r), \quad 1 \leq k \leq p, \quad r \geq 0,$$

16 and

$$17 \quad \mathbf{p}_{-k}^{r+1} = (\varphi_1^{r+1}, \dots, \varphi_{k-1}^{r+1}, \varphi_{k+1}^r, \dots, \varphi_p^r), \quad 1 \leq k \leq p, \quad r \geq 0.$$

18 Theoretical justification of the estimation accuracy of the COSALES Lasso will be
 19 deferred to the next section.

20 4.1.2. *Nonconvex penalized COSALES regression.* In (4.2), let $p_{\lambda_1}(\cdot)$ and $p_{\lambda_2}(\cdot)$
 21 be nonconvex penalties having properties (P1-P5). This nonconvex penalized COS-
 22 ALES estimation problem can be solved by the LLA algorithm shown in Algo-
 23 rithm 4. Note that the minimization problem in step (a) was solved in Algorithm 3.
 24 Oracle properties of the sparse solutions will be established in the following section.
 25

26 4.2. *Theory.* In this section, we show the selection and estimation accuracy of
 27 the COSALES regression for both Lasso and nonconvex penalties.

28 4.2.1. *L₁-penalized COSALES regression.* For the Lasso problem (4.4), let
 29 $\check{M} = (\lambda_1^{\text{lasso}}/\lambda_2^{\text{lasso}}) \vee (\lambda_2^{\text{lasso}}/\lambda_1^{\text{lasso}})$ and define set $A_0 = (A_1, A'_2)$, where $A'_2 =$
 30 $\{j+p: \omega_j^* \neq 0\}$. For $M \geq 1$, define $\mathcal{C}_M = \{\delta \in \mathbb{R}^{2p}: \|\delta_{A_0}\|_1 \leq M \|\delta_{A_0}\|_1 \neq 0\}$.
 31 For $k = 1, 2$, let $\rho_{k\bullet\min} = \lambda_{\min}(n^{-1} \mathbf{X}_{A_k}^T \mathbf{X}_{A_k})$ and $\rho_{k\bullet\max} = \lambda_{\max}(n^{-1} \mathbf{X}_{A_k}^T \mathbf{X}_{A_k})$.

Algorithm 3: COSALES — The coordinate descent plus proximal gradient algorithm for solving the coupled weighted L_1 -minimization problem (4.5)

1. Initialize the algorithm with $\boldsymbol{\gamma}^0 = (\gamma_1^0, \dots, \gamma_p^0)^\top$ and $\boldsymbol{\varphi}^0 = (\varphi_1^0, \dots, \varphi_p^0)^\top$.
 2. For $r = 1, \dots, m - 1$,
 - (2.1) For $k = 1, \dots, p$,
 - (2.1.1) Initialize $\gamma_k^{r,0} := \gamma_k^r$.
 - (2.1.2) For $s = 0, 1, \dots, s_{1k}^r - 1$,
 - (2.1.2.1) Compute $\gamma_k^{r,s+1} := \mathbb{S}_{L_{1k}^{-1}w_k}(\gamma_k^{r,s} - L_{1k}^{-1}h'_n(\gamma_k^{r,s}; \mathbf{g}_{-k}^{r+1}, \boldsymbol{\varphi}^r))$, where $L_{1k} = (2\bar{c} + 1)n^{-1}\|X_k\|_2^2$; $h_n(\gamma_k; \mathbf{g}_{-k}^{r+1}, \boldsymbol{\varphi}^r) = S_n([\gamma_k, \mathbf{g}_{-k}^{r+1}], \boldsymbol{\varphi}^r)$.
 - (2.1.3) Set $\gamma_k^{r+1} := \gamma_k^{r,s_{1k}^r}$.
 - (2.2) Set $\boldsymbol{\gamma}^{r+1} := (\gamma_1^{r+1}, \dots, \gamma_p^{r+1})^\top$.
 - (2.3) For $k = 1, \dots, p$,
 - (2.3.1) Initialize $\varphi_k^{r,0} := \varphi_k^r$.
 - (2.3.2) For $s = 0, 1, \dots, s_{2k}^r - 1$,
 - (2.3.2.1) Compute $\varphi_k^{r,s+1} := \mathbb{S}_{L_{2k}^{-1}v_k}(\varphi_k^{r,s} - L_{2k}^{-1}h'_n(\varphi_k^{r,s}; \boldsymbol{\gamma}^{r+1}, \mathbf{p}_{-k}^{r+1}))$, where $L_{2k} = 2\bar{c}n^{-1}\|X_k\|_2^2$; $h_n(\varphi_k; \boldsymbol{\gamma}^{r+1}, \mathbf{p}_{-k}^{r+1}) = S_n(\boldsymbol{\gamma}^{r+1}, [\varphi_k, \mathbf{p}_{-k}^{r+1}])$.
 - (2.3.3) Set $\varphi_k^{r+1} := \varphi_k^{r,s_{2k}^r}$.
 - (2.4) Set $\boldsymbol{\varphi}^{r+1} := (\varphi_1^{r+1}, \dots, \varphi_p^{r+1})^\top$.
 3. Output $\hat{\boldsymbol{\gamma}} := \boldsymbol{\gamma}^m$ and $\hat{\boldsymbol{\varphi}} := \boldsymbol{\varphi}^m$.
-

Algorithm 4: The local linear approximation (LLA) algorithm for solving the nonconvex penalized COSALES estimation problem (4.2)

1. Initialize $\hat{\boldsymbol{\gamma}}^0 = \hat{\boldsymbol{\gamma}}^{\text{initial}}$ and $\hat{\boldsymbol{\varphi}}^0 = \hat{\boldsymbol{\varphi}}^{\text{initial}}$. Compute weights

$$\hat{w}_j^0 = p'_{\lambda_1}(|\hat{\gamma}_j^0|), \bar{w}_j^0 = p'_{\lambda_2}(|\hat{\varphi}_j^0|), j = 1, \dots, p.$$
 2. For $m = 1, 2, \dots$, repeat the LLA iteration in (a) and (b) until convergence.
 - (a) Solve the following convex optimization problem for $\hat{\boldsymbol{\gamma}}^m$ and $\hat{\boldsymbol{\varphi}}^m$

$$\min_{\boldsymbol{\gamma}, \boldsymbol{\varphi} \in \mathbb{R}^p} S_n(\boldsymbol{\gamma}, \boldsymbol{\varphi}) + \sum_{j=1}^p \hat{w}_j^{m-1} |\gamma_j| + \sum_{j=1}^p \bar{w}_j^{m-1} |\varphi_j|.$$
 - (b) Update the weights

$$\hat{w}_j^m = p'_{\lambda_1}(|\hat{\gamma}_j^m|), \bar{w}_j^m = p'_{\lambda_2}(|\hat{\varphi}_j^m|), j = 1, \dots, p.$$
-

1 Denote $\phi_{\min} = \rho_{1\bullet\min} \wedge \rho_{2\bullet\min}$ and $\phi_{\max} = \rho_{1\bullet\max} \vee \rho_{2\bullet\max}$. Assume $\phi_{\min} > 0$.
 2 Let \mathbf{I}_2 be a 2×2 identity matrix and let \otimes denote the Kronecker product. To establish
 3 an error bound on the COSALES Lasso estimators, the following conditions on the
 4 design matrix \mathbf{X} and the random errors ε are imposed:

- 5 (C1') The columns of \mathbf{X} is normalizable, that is, $M_0 = \max_{1 \leq j \leq p} \frac{\|X_j\|_2}{\sqrt{n}} \in (0, \infty)$.
 6 (C2') $M_1 = \|\mathbf{X}^T \boldsymbol{\omega}^*\|_\infty \in (0, \infty)$.
 7 (C3') The random errors ε_i are i.i.d. mean zero sub-Gaussian random variables.
 8 (C4') $\bar{\kappa} = \kappa(3\check{M}) \in (0, \infty)$, where $\kappa(M) = \inf_{\boldsymbol{\delta} \in \mathcal{C}_M} \boldsymbol{\delta}^T [\mathbf{I}_2 \otimes (n^{-1} \mathbf{X}^T \mathbf{X})] \boldsymbol{\delta} / \|\boldsymbol{\delta}\|_2^2$.
 9 (C5') $\bar{\varrho} = \varrho(3\check{M}) \in (0, \infty)$, where $\varrho(M) = \inf_{\boldsymbol{\delta} \in \mathcal{C}_M} \frac{\boldsymbol{\delta}^T [\mathbf{I}_2 \otimes (n^{-1} \mathbf{X}^T \mathbf{X})] \boldsymbol{\delta}}{\|\boldsymbol{\delta}_{A_0}\|_1 \|\boldsymbol{\delta}\|_\infty}$.

10 **THEOREM 3.** *In model (4.1), suppose the true parameter vectors $\boldsymbol{\gamma}^*$ and $\boldsymbol{\omega}^*$*
 11 *are respectively s_1 -sparse and s_2 -sparse and assume conditions (C1'-C3') hold. Let*
 12 *$\hat{\boldsymbol{\gamma}}^{\text{lasso}}$ and $\hat{\boldsymbol{\varphi}}^{\text{lasso}}$ be any optimal solutions to the L_1 -penalized COSALES estimation*
 13 *problem (4.4). Then with probability at least $1 - \pi_1^{\text{ALS}}$,*

$$14 \quad \left\| \begin{pmatrix} \hat{\boldsymbol{\gamma}}^{\text{lasso}} \\ \hat{\boldsymbol{\varphi}}^{\text{lasso}} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\gamma}^* \\ \boldsymbol{\varphi}^* \end{pmatrix} \right\|_2 \leq 3(s_1 + s_2)^{1/2} (\lambda_1^{\text{lasso}} \vee \lambda_2^{\text{lasso}}) (2\bar{\kappa}c_0)^{-1}$$

15 *if condition (C4') holds and*

$$16 \quad \left\| \begin{pmatrix} \hat{\boldsymbol{\gamma}}^{\text{lasso}} \\ \hat{\boldsymbol{\varphi}}^{\text{lasso}} \end{pmatrix} - \begin{pmatrix} \boldsymbol{\gamma}^* \\ \boldsymbol{\varphi}^* \end{pmatrix} \right\|_\infty \leq 3(\lambda_1^{\text{lasso}} \vee \lambda_2^{\text{lasso}}) (2\bar{\varrho}c_0)^{-1}$$

17 *if condition (C5') holds, where*

$$18 \quad \pi_1^{\text{ALS}} = 2p \exp\left(-\frac{Cn(\lambda_1^{\text{lasso}})^2}{4M_0^2 M_1^2 (K_1 + K_2)^2}\right) + 2p \exp\left(-\frac{Cn(\lambda_2^{\text{lasso}})^2}{4M_0^2 M_1^2 K_2^2}\right),$$

19 $c_0 = 2^{-1}[(1 + 4\bar{c}) - (1 + 16\bar{c}^2)^{1/2}]$, $K_1 = \|\varepsilon_i\|_{SG}$, $K_2 = \|\Psi'_\tau(\varepsilon_i - e_\tau)\|_{SG}$, and
 20 $C > 0$ is an absolute constant.

21 **4.2.2. Nonconvex penalized COSALES regression.** We show that the oracle
 22 estimators $\hat{\boldsymbol{\gamma}}^{\text{oracle}}$ and $\hat{\boldsymbol{\varphi}}^{\text{oracle}}$ can be achieved with overwhelming probability by
 23 Algorithm 4 under rather general conditions. Indeed, suppose the minimal signal
 24 strength of $\boldsymbol{\gamma}^*$ and $\boldsymbol{\omega}^*$ satisfies

$$25 \quad (\text{A0}') \quad \min_{j \in A_1} |\gamma_j^*| > (a + 1)\lambda_1 \quad \text{and} \quad \min_{j \in A_2} |\omega_j^*| > (a + 1)|e_\tau|^{-1}\lambda_2.$$

26 **THEOREM 4.** *Suppose in model (4.1) $\boldsymbol{\gamma}^*$ and $\boldsymbol{\omega}^*$ are respectively s_1 -sparse and*
 27 *s_2 -sparse and satisfy assumption (A0'). Take $\hat{\boldsymbol{\gamma}}^{\text{lasso}}$ and $\hat{\boldsymbol{\varphi}}^{\text{lasso}}$ as the initial values*
 28 *and assume conditions (C1'-C3') hold. Take $\lambda \geq 3s^{1/2}(\lambda_1^{\text{lasso}} \vee \lambda_2^{\text{lasso}})(2a_0c_0\bar{\kappa})^{-1}$*

1 when (C4') holds, or take $\lambda \geq 3(\lambda_1^{\text{lasso}} \vee \lambda_2^{\text{lasso}})(2a_0c_0\bar{\varrho})^{-1}$ when (C5') holds, or
 2 take $\lambda \geq 3(\lambda_1^{\text{lasso}} \vee \lambda_2^{\text{lasso}})(2a_0c_0)^{-1}[(s^{1/2}\bar{\kappa}^{-1}) \wedge \bar{\varrho}^{-1}]$ when both (C4') and (C5')
 3 hold. The LLA algorithm (Algorithm 4) converges to the oracle estimators $\hat{\gamma}^{\text{oracle}}$
 4 and $\hat{\varphi}^{\text{oracle}}$ in two iterations with probability at least $1 - \pi_1^{\text{ALS}} - \pi_2^{\text{ALS}} - \pi_3^{\text{ALS}}$, where
 5 π_1^{ALS} is given in Theorem 3,

$$\begin{aligned} & \pi_2^{\text{ALS}} = \Gamma(2^{-1}Q_2\lambda; n, s_1, K_1 + K_2, M_0M_1, M_1^2\rho_{1\bullet\max}, \nu_1) \\ & \quad + \Gamma(2^{-1}Q_2\lambda; n, s_2, K_2, M_0M_1, M_1^2\rho_{2\bullet\max}, \nu_2) \\ & \quad + 2(p - s_1) \exp\left(-\frac{Ca_1^2n\lambda^2}{4M_0^2M_1^2(K_1 + K_2)^2}\right) + 2(p - s_2) \exp\left(-\frac{Ca_1^2n\lambda^2}{4M_0^2M_1^2K_2^2}\right), \end{aligned}$$

7 and

$$\begin{aligned} & \pi_3^{\text{ALS}} = \Gamma(2^{-1}c_0\phi_{\min}\bar{R}; n, s_1, K_1 + K_2, M_0M_1, M_1^2\rho_{1\bullet\max}, \nu_1) \\ & \quad + \Gamma(2^{-1}c_0\phi_{\min}\bar{R}; n, s_2, K_2, M_0M_1, M_1^2\rho_{2\bullet\max}, \nu_2), \end{aligned}$$

9 where $s = s_1 + s_2$, $\lambda = \lambda_1 \wedge \lambda_2$, $Q_2 = a_1c_0\phi_{\min}[2(1 + 2\bar{c})M_0\phi_{\max}^{1/2}]^{-1}$, $\nu_1 =$
 10 $\text{var}(\varepsilon_i + \Psi'_\tau(\varepsilon_i - e_\tau))$, $\nu_2 = \text{var}(\Psi'_\tau(\varepsilon_i - e_\tau))$, $\bar{R} = (\min_{j \in A_1} |\gamma_j^*| - a\lambda_1) \wedge$
 11 $(\min_{j \in A_2} |\varphi_j^*| - a\lambda_2)$, $C > 0$ is an absolute constant, c_0, K_1, K_2 are given in
 12 Theorem 3, and $\Gamma(\cdot)$ is given in Theorem 2.

13 For SCAD and MCP penalized COSALES regressions, the LLA algorithm
 14 (Algorithm 4) starting from the zero vector can also be used as long as we can take
 15 $\lambda_k = \lambda_k^{\text{lasso}}$, $k = 1, 2$.

16 **COROLLARY 2.** Assume the same framework of Theorem 4 and suppose the
 17 SCAD penalty (2.8) or MCP (2.9) is used. If condition (C4') holds and $2a_0c_0\bar{\kappa} \geq$
 18 $3\check{M}s^{1/2}$, or if condition (C5') holds and $2a_0c_0\bar{\varrho} \geq 3\check{M}$, or if both (C4') and (C5')
 19 hold and $3\check{M}[(s^{1/2}\bar{\kappa}^{-1}) \wedge \bar{\varrho}^{-1}] \leq 2a_0c_0$, then the LLA algorithm (Algorithm 4)
 20 initialized by zero converges to the oracle estimators $\hat{\gamma}^{\text{oracle}}$ and $\hat{\varphi}^{\text{oracle}}$ after three
 21 iterations with probability at least $1 - \check{\pi}_1^{\text{ALS}} - \pi_2^{\text{ALS}} - \pi_3^{\text{ALS}}$, where

$$\check{\pi}_1^{\text{ALS}} = 2p \exp\left(-\frac{Cn\lambda_1^2}{4M_0^2M_1^2(K_1 + K_2)^2}\right) + 2p \exp\left(-\frac{Cn\lambda_2^2}{4M_0^2M_1^2K_2^2}\right),$$

23 π_2^{ALS} and π_3^{ALS} are given in Theorem 4, and $s = s_1 + s_2$.

24 **REMARK 2.** We can easily modify (4.2) to allow certain subsets of coefficients
 25 not to be penalized. Let \mathcal{R}_1 and \mathcal{R}_2 be the index sets of unpenalized components of
 26 γ and φ , respectively. Then (4.2) can be modified as

$$\min_{\gamma, \varphi \in \mathbb{R}^p} S_n(\gamma, \varphi) + \sum_{j \in \mathcal{R}_1^c} p_{\lambda_1}(\gamma_j) + \sum_{j \in \mathcal{R}_2^c} p_{\lambda_2}(\varphi_j).$$

1 The COSALES algorithm can be readily used to solve the above problem. Moreover,
2 similar theoretical results can be established with slight modifications.

3 *4.3. Simulation examples.* We demonstrate the selection and estimation ac-
4 curacy of the COSALES regression through two numerical simulations. For the
5 nonconvex penalties used in both simulations, we fix $\gamma = 3.7$ for the SCAD penalty
6 and $\gamma = 2$ for the MCP.

7 **EXAMPLE 2.** We consider the same model (3.1) that was used in Example 1,
8 but different from the approach used there, we estimate the coefficients through the
9 nonconvex penalized COSALES regression (4.2). Again we choose $p = 600$ and
10 independently simulate a training set of size $n = 300$ for fitting and a validation
11 set of size $n = 300$ for tuning. The tuning parameter is selected by minimizing
12 the validation error $\sum_{i \in \text{validation}} \{\Psi_{0.5}(y_i - \mathbf{x}_i^T \hat{\boldsymbol{\gamma}}) + \Psi_{\tau}(y_i - \mathbf{x}_i^T \hat{\boldsymbol{\gamma}} - \mathbf{x}_i^T \hat{\boldsymbol{\varphi}})\}$ for the
13 computed estimates $\hat{\boldsymbol{\gamma}}$ and $\hat{\boldsymbol{\varphi}}$. We pick a fairly extreme τ -value ($\tau = 0.95$) for easy
14 separation of the conditional mean and scale functions. Both the COSALES Lasso
15 and two variations of the LLA algorithm for each of the SCAD and MCP penalized
16 COSALES regressions are implemented.

17 Based on 100 independent runs, the following measurements are calculated
18 to evaluate the sparsity recovery and estimation performance of the COSALES
19 estimators:

20 $|\hat{A}_1|, |\hat{A}_2|$: the average size of the active sets for $\hat{\boldsymbol{\gamma}}$ and $\hat{\boldsymbol{\varphi}}$, respectively, $\hat{A}_1 =$
21 $\{j: \hat{\gamma}_j \neq 0\}$ and $\hat{A}_2 = \{j: \hat{\varphi}_j \neq 0\}$.

22 p_{a_1}, p_{a_2} : proportions of the events $A_1 \subset \hat{A}_1$ and $A_2 \subset \hat{A}_2$, respectively, where
23 $A_1 = \{6, 12, 15, 20\}$ denotes the active set of $\boldsymbol{\gamma}^*$ and $A_2 = \{1\}$ denotes the
24 active set of $\boldsymbol{\varphi}^*$.

25 $R_1^{\boldsymbol{\gamma}}, R_1^{\boldsymbol{\varphi}}$: the average L_1 risks, $R_1^{\boldsymbol{\gamma}} = \|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*\|_1$ and $R_1^{\boldsymbol{\varphi}} = \|\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi}^*\|_1$.

26 $R_2^{\boldsymbol{\gamma}}, R_2^{\boldsymbol{\varphi}}$: the average L_2 risks, $R_2^{\boldsymbol{\gamma}} = \|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}^*\|_2$ and $R_2^{\boldsymbol{\varphi}} = \|\hat{\boldsymbol{\varphi}} - \boldsymbol{\varphi}^*\|_2$.

27 The results are summarized in Table 2, from which we can draw the following
28 conclusions:

- 29 (1) The COSALES regression (with Lasso or nonconvex penalties) can recover
30 the sparse patterns in both the mean and scale functions with overwhelming
31 probabilities.
- 32 (2) The COSALES Lasso tends to select a lot more irrelevant covariates and
33 has much larger estimation errors than the nonconvex penalized COSALES
34 regression (with the SCAD penalty or MCP).
- 35 (3) The three-step LLA algorithm starting from zero produces similar results to
36 the two-step LLA algorithm starting from the Lasso solution.

37 **EXAMPLE 3.** In this example, we simulate data from the following normal

TABLE 2

Numerical summary of simulation results from the Lasso, SCAD and MCP penalized COSALES regression for model (3.1) $y = x_6 + x_{12} + x_{15} + x_{20} + (0.7x_1)\varepsilon$. The selection accuracy is measured by the number of selected variables $|\hat{A}_1|$ and $|\hat{A}_2|$, and the proportions p_{a_1} and p_{a_2} of covering the true active sets. The estimation accuracy is measured by the L_1 risks R_1^γ and R_1^φ , and the L_2 risks R_2^γ and R_2^φ . The results are shown as averages over 100 replicates with standard errors listed in the parentheses. A fairly extreme τ -value ($\tau = 0.95$) is used in the simulation for easy separation of the mean and scale

Method	$ \hat{A}_1 $	$ \hat{A}_2 $	p_{a_1}	p_{a_2}	R_1^γ	R_1^φ	R_2^γ	R_2^φ
COSALES-Lasso	26.88 (1.04)	13.36 (0.45)	100% (0)	100% (0)	0.407 (0.012)	0.378 (0.008)	0.124 (0.002)	0.294 (0.006)
COSALES-SCAD*	7.24 (0.10)	1.01 (0.01)	100% (0)	100% (0)	0.095 (0.004)	0.072 (0.005)	0.048 (0.002)	0.072 (0.005)
COSALES-SCAD ⁰	8.85 (0.57)	1.01 (0.01)	100% (0)	100% (0)	0.107 (0.005)	0.065 (0.005)	0.049 (0.002)	0.065 (0.005)
COSALES-MCP*	6.46 (0.38)	1.01 (0.01)	100% (0)	100% (0)	0.089 (0.004)	0.070 (0.005)	0.045 (0.002)	0.070 (0.005)
COSALES-MCP ⁰	7.08 (0.44)	1.01 (0.01)	100% (0)	100% (0)	0.102 (0.006)	0.067 (0.005)	0.052 (0.003)	0.067 (0.005)

1 linear heteroscedastic model

2 (4.6)
$$y = x_6 + x_{12} + x_{15} + x_{20} + (0.7x_1 + 0.7x_{12})\varepsilon,$$

3 where the covariates are simulated by setting $x_1 = \Phi(z_1)$, $x_{12} = \Phi(z_{12})$, and $x_j =$
4 z_j , $j \neq 1, 12$, where $(z_1, \dots, z_p)^T \sim N(\mathbf{0}, \Sigma)$ with $\Sigma = (0.5^{|i-j|})_{p \times p}$, and $\Phi(\cdot)$ is
5 the CDF of the standard normal distribution. The random error $\varepsilon \sim N(0, 1)$. Note
6 that in model (3.1), the active sets of the true parameter vectors do not overlap, so the
7 SALES regression can detect active variables in the scale. However, in model (4.6)
8 the active set for the mean, $A_1 = \{6, 12, 15, 20\}$, overlaps with the active set for the
9 scale, $A_2 = \{1, 12\}$. Thus, the SALES regression cannot recover the variable x_{12}
10 in the scale function. We show by this Monte Carlo simulation that the COSALES
11 regression can recover the sparse patterns in both the mean and scale functions. We
12 fix $p = 600$ and independently simulate a training set of size $n = 500$ for fitting and
13 a validation set of the same size for tuning. We select the regularization parameter by
14 minimizing the validation error $\sum_{i \in \text{validation}} \{ \Psi_{0.5}(y_i - \mathbf{x}_i^T \hat{\gamma}) + \Psi_\tau(y_i - \mathbf{x}_i^T \hat{\gamma} - \mathbf{x}_i^T \hat{\varphi}) \}$
15 for the computed estimate $\hat{\gamma}$ and $\hat{\varphi}$. In order to separate the mean and scale easily,
16 we again pick $\tau = 0.95$. We implement the COSALES Lasso and two variations of
17 the LLA algorithm as were done in Examples 2 for each of the SCAD and MCP
18 penalized COSALES regressions.

19 Based on 100 independent runs, the same measurements of performance as in
20 Example 2 are calculated to evaluate the sparsity recovery and estimation accuracy of

1 the COSALES estimation. The results are summarized in Table 3. Same conclusions
2 in Example 2 can be drawn here.

TABLE 3

Numerical summary of simulation results from the the Lasso, SCAD and MCP penalized COSALES regression for model (4.6): $y = x_6 + x_{12} + x_{15} + x_{20} + (0.7x_1 + 0.7x_{12})\varepsilon$. The selection accuracy is measured by the number of selected variables $|\hat{A}_1|$ and $|\hat{A}_2|$, and the proportions p_{α_1} and p_{α_2} of covering the true active sets. The estimation accuracy is measured by the L_1 risks R_1^γ and R_1^φ , and the L_2 risks R_2^γ and R_2^φ . The results are shown as averages over 100 replicates with standard errors listed in the parentheses. A fairly extreme τ -value ($\tau = 0.95$) is used in the simulation for easy separation of the mean and scale

Method	$ \hat{A}_1 $	$ \hat{A}_2 $	p_{α_1}	p_{α_2}	R_1^γ	R_1^φ	R_2^γ	R_2^φ
COSALES-Lasso	27.92 (0.98)	12.67 (0.49)	100% (0)	100% (0)	0.719 (0.018)	0.450 (0.011)	0.249 (0.006)	0.282 (0.008)
COSALES-SCAD*	6.80 (0.52)	2.06 (0.04)	100% (0)	100% (0)	0.167 (0.008)	0.210 (0.014)	0.089 (0.004)	0.161 (0.010)
COSALES-SCAD ⁰	5.70 (0.25)	2.02 (0.01)	100% (0)	100% (0)	0.157 (0.006)	0.199 (0.013)	0.090 (0.003)	0.148 (0.009)
COSALES-MCP*	5.95 (0.35)	2.06 (0.03)	100% (0)	100% (0)	0.153 (0.006)	0.221 (0.015)	0.086 (0.003)	0.165 (0.010)
COSALES-MCP ⁰	6.00 (0.36)	2.04 (0.02)	100% (0)	100% (0)	0.180 (0.009)	0.205 (0.014)	0.098 (0.004)	0.154 (0.010)

3 **5. Real Data Example.** We apply the SALES and COSALES regressions to a
4 real data set reported in Scheetz et al. (2006). The data set consists of gene expres-
5 sion levels of more than 31,000 probes obtained from 120 rats. The expressions are
6 analyzed on a logarithmic scale (base 2). As was done in Scheetz et al. (2006), we
7 exclude the probes that were not expressed in the eye or that lacked sufficient varia-
8 tion. Among those 18,976 probes left, we study how the expressions of other genes
9 are associated with the gene *TRIM32* (probe 1389163.at). This gene was found to be
10 associated with Bardet–Biedl syndrome, which is a disorder that affects many parts
11 of the body including the retina. For all the other genes, we first standardize them
12 and select the 3,000 probes with the largest variances. These 3,000 probes are then
13 ranked according to the magnitude of the correlations between their expressions and
14 that of probe 1389163.at. We choose the top 300 probes with the largest correlations
15 in magnitude for the analysis.

16 The third column of Table 4 lists the number of active variables selected by the
17 SALES regressions with Lasso, SCAD and MCP penalties, fitted on the whole data
18 set of 120 subjects. For both SCAD and MCP penalized SALES regressions, the two
19 variations of the LLA algorithm were used. The tuning parameter for each method
20 is selected by five-fold cross-validation. The last two columns of Table 4 summarize

1 the results from 50 random partitions. Each partition randomly splits the data into a
 2 training set with 80 observations and a validation set with 40 observations. We fit the
 3 model with the training set using five-fold cross-validation for tuning and calculate
 4 the predicted loss $(1/40) \sum_{i \in \text{validation}} \Psi_\tau(y_i - \hat{\beta}_0 - \mathbf{x}_i^\top \hat{\beta})$ based on the validation
 5 set. The average number of active variables selected and the average predicted loss
 6 are calculated from the 50 partitions with their respective standard errors listed
 7 in the parentheses. Table 4 reveals two interesting findings. First, the nonconvex
 8 penalized SALES regression selects less variables than the SALES Lasso, but there
 9 is no obvious improvement of the nonconvex penalized SALES regression over the
 10 SALES Lasso in terms of predicted loss. Second, for all SALES regressions, the
 11 number of variables selected is different at different values of τ (0.3, 0.5 and 0.7).
 12 This is an indication of heteroscedasticity in the data.

13 To further explore the heterogeneous scale, we also apply the COSALES re-
 14 gression to the data. The results are summarized in Table 5. Columns 2 and 3
 15 display the number of variables selected for the mean ($|\hat{A}_1|$) and scale ($|\hat{A}_2|$), and
 16 the number of variables that overlap ($|\hat{A}_1 \cap \hat{A}_2|$) for each method. For all penal-
 17 ties, τ is set to be 0.7 in the COSALES regression. Random partitions are done
 18 in the same way as the SALES regression and the predicted loss is calculated via
 19 $(1/40) \sum_{i \in \text{validation}} \Psi_{0.5}(y_i - \hat{\gamma}_0 - \mathbf{x}_i^\top \hat{\gamma}) + \Psi_\tau(y_i - \hat{\gamma}_0 - \mathbf{x}_i^\top \hat{\gamma} - \hat{\varphi}_0 - \mathbf{x}_i^\top \hat{\varphi})$. The
 20 results for the random partitions are shown in columns 4 to 6. It can be seen that
 21 the COSALES regression reveals more information about the heterogeneous scale
 22 which cannot be otherwise detected in the SALES regression or the sparse quantile
 23 regression (Wang, Wu and Li, 2012) due to overlaps.

24 **6. Proofs.** In this section, we give the proofs of the main theoretical results
 25 stated in previous sections. First of all, let us state two lemmas on the properties
 26 of the asymmetric squared error loss $\Psi_\tau(\cdot)$ given in (2.1). These properties play an
 27 important role in the proofs of many results to be presented below. Let $w_\tau(u) =$
 28 $|\tau - I(u < 0)|$ and recall that $\underline{c} = \tau \wedge (1 - \tau)$ and $\bar{c} = \tau \vee (1 - \tau)$.

29 **LEMMA 1.** *The asymmetric squared error loss $\Psi_\tau(\cdot)$ is continuously differ-*
 30 *entiable, but is not twice differentiable at zero when $\tau \neq 0.5$. Moreover, for any*
 31 *$u, u_0 \in \mathbb{R}$ and $\tau \in (0, 1)$, we have*

$$32 \quad \underline{c}(u - u_0)^2 \leq \Psi_\tau(u) - \Psi_\tau(u_0) - \Psi'_\tau(u_0)(u - u_0) \leq \bar{c}(u - u_0)^2.$$

33 *It follows that $\Psi_\tau(\cdot)$ is strongly convex.*

34 **LEMMA 2.** *For any $u, u_0 \in \mathbb{R}$ and $\tau \in (0, 1)$, we have*

$$35 \quad 2\underline{c}|u - u_0| \leq |\Psi'_\tau(u) - \Psi'_\tau(u_0)| \leq 2\bar{c}|u - u_0|.$$

36 *It follows immediately that $\Psi'_\tau(\cdot)$ is Lipschitz continuous.*

TABLE 4

Analysis of microarray data using SALES regressions with Lasso, SCAD and MCP penalties. Three different values of τ (0.3, 0.5 and 0.7) are used for each method. The number of active variables selected using the whole data set is given in column 3. The average number of active variables selected and average predicted loss $(1/40) \sum_{i \in \text{validation}} \Psi_{\tau}(y_i - \hat{\beta}_0 - \mathbf{x}_i^T \hat{\beta})$ listed in columns 4 and 5 are calculated from 50 random partitions of the original data with standard errors listed in parentheses

Method	τ	All data		Random partition	
		$ \hat{A} $	$ \hat{A} $	$ \hat{A} $	Predicted loss
SALES-Lasso	0.3	22	22.00 (1.51)	0.007 (0.00055)	
	0.5	25	25.38 (1.94)	0.005 (0.00036)	
	0.7	20	21.90 (1.66)	0.005 (0.00022)	
SALES-SCAD*	0.3	19	16.02 (2.09)	0.006 (0.00048)	
	0.5	13	15.52 (1.80)	0.006 (0.00043)	
	0.7	11	13.54 (1.98)	0.005 (0.00037)	
SALES-SCAD ⁰	0.3	16	16.60 (2.03)	0.006 (0.00054)	
	0.5	17	17.22 (2.36)	0.007 (0.00048)	
	0.7	14	14.82 (2.18)	0.005 (0.00030)	
SALES-MCP*	0.3	14	15.82 (2.56)	0.006 (0.00053)	
	0.5	12	12.66 (2.58)	0.008 (0.00054)	
	0.7	10	9.66 (1.78)	0.006 (0.00035)	
SALES-MCP ⁰	0.3	11	11.74 (1.47)	0.006 (0.00057)	
	0.5	13	13.24 (2.75)	0.007 (0.00058)	
	0.7	13	14.18 (3.36)	0.006 (0.00034)	

TABLE 5

Analysis of microarray data using COSALES regressions with Lasso, SCAD and MCP penalties. In this analysis, $\tau = 0.7$ is used. The number of active variables selected for the mean and scale using the whole data set is given in columns 2 and 3. The average number of active variables selected for the mean and scale and average predicted loss $(1/40) \sum_{i \in \text{validation}} \Psi_{0.5}(y_i - \hat{\gamma}_0 - \mathbf{x}_i^T \hat{\gamma}) + \Psi_{\tau}(y_i - \hat{\gamma}_0 - \mathbf{x}_i^T \hat{\gamma} - \hat{\varphi}_0 - \mathbf{x}_i^T \hat{\varphi})$ listed in columns 4 to 6 are calculated from 50 random partitions of the original data with standard errors listed in parentheses

Method	All data			Random partition			Predicted loss
	$ \hat{A}_1 $	$ \hat{A}_2 $	$ \hat{A}_1 \cap \hat{A}_2 $	$ \hat{A}_1 $	$ \hat{A}_2 $	$ \hat{A}_1 \cap \hat{A}_2 $	
COSALES-Lasso	22	10	9	22.62 (1.21)	9.80 (1.10)	7.86 (0.93)	0.010 (0.00056)
COSALES-SCAD*	19	7	6	18.92 (0.79)	5.58 (0.50)	3.90 (0.31)	0.011 (0.00067)
COSALES-SCAD ⁰	20	5	4	20.22 (0.98)	5.82 (0.64)	3.92 (0.44)	0.011 (0.00072)
COSALES-MCP*	10	3	1	10.96 (2.32)	3.08 (1.40)	1.38 (0.74)	0.014 (0.00096)
COSALES-MCP ⁰	10	4	3	12.94 (1.83)	4.56 (1.04)	1.46 (0.42)	0.012 (0.00083)

1 **PROOF OF LEMMA 1.** It is easy to see that $\underline{c} \leq w_\tau(u) \leq \bar{c}$ for any $u \in \mathbb{R}$. Note
 2 that $\Psi'_\tau(u) = 2w_\tau(u)u$, which is continuous and which is not differentiable at
 3 $u = 0$ when $\tau \neq 0.5$. To show the inequalities, consider the following situations. If
 4 $w_\tau(u) \geq w_\tau(u_0)$, it follows that

$$\begin{aligned} & \Psi_\tau(u) - \Psi_\tau(u_0) - \Psi'_\tau(u_0)(u - u_0) \\ &= w_\tau(u)u^2 - w_\tau(u_0)u_0^2 - 2w_\tau(u_0)u_0(u - u_0) \\ 5 \quad &= w_\tau(u_0)(u - u_0)^2 + \{w_\tau(u) - w_\tau(u_0)\}u^2 \\ &\geq w_\tau(u_0)(u - u_0)^2 \geq \underline{c}(u - u_0)^2. \end{aligned}$$

6 Otherwise, if $w_\tau(u) < w_\tau(u_0)$, then we know that $\underline{c} = w_\tau(u)$, $\bar{c} = w_\tau(u_0)$ and
 7 $u_0u \leq 0$. It follows that

$$\begin{aligned} & \Psi_\tau(u) - \Psi_\tau(u_0) - \Psi'_\tau(u_0)(u - u_0) \\ 8 \quad &= \underline{c}u^2 - \bar{c}u_0^2 - 2\bar{c}u_0(u - u_0) \\ &\geq \underline{c}u^2 - 2\underline{c}u_0u + \underline{c}u_0^2 = \underline{c}(u - u_0)^2. \end{aligned}$$

9 Therefore, the first inequality holds. Similarly, we can show the second inequality.
 10 □

11 **PROOF OF LEMMA 2.** If $u = 0$ or $u_0 = 0$, then the inequalities hold trivially. If
 12 $uu_0 > 0$, we know that $w_\tau(u) = w_\tau(u_0)$. It follows that

$$13 \quad 2\underline{c}|u - u_0| \leq |\Psi'_\tau(u) - \Psi'_\tau(u_0)| = 2w_\tau(u)|u - u_0| \leq 2\bar{c}|u - u_0|.$$

14 If instead, $uu_0 < 0$, there are two cases: $u > 0, u_0 < 0$ or $u < 0, u_0 > 0$. For the
 15 first case, we have

$$16 \quad 2\underline{c}|u - u_0| \leq |\Psi'_\tau(u) - \Psi'_\tau(u_0)| = 2\tau u - 2(1 - \tau)u_0 \leq 2\bar{c}|u - u_0|.$$

17 For the second case, we have

$$18 \quad 2\underline{c}|u - u_0| \leq |\Psi'_\tau(u) - \Psi'_\tau(u_0)| = -2(1 - \tau)u + 2\tau u_0 \leq 2\bar{c}|u - u_0|.$$

19 This completes the proof. □

20 The following lemma deals with sub-Gaussian random variables.

21 **LEMMA 3.** *Suppose that $Z, Z_1, \dots, Z_n \in \mathbb{R}$ are i.i.d. sub-Gaussian random*
 22 *variables. Let $\mathbf{Z} = (Z_1, \dots, Z_n)^T$, $K = \|Z\|_{SG}$, $Z^+ = \max(Z, 0)$ and $Z^- =$
 23 $\max(-Z, 0)$.*

- 1 (1) If $\mathbb{E}(Z) = 0$, then there exists an absolute constant $C > 0$ such that for any
2 $\mathbf{a} = (a_1, \dots, a_n)^T \in \mathbb{R}^n$ and any $t \geq 0$,

$$3 \quad \mathbb{P}(|\mathbf{a}^T \mathbf{Z}| \geq t) \leq 2 \exp\left(-\frac{Ct^2}{K^2 \|\mathbf{a}\|_2^2}\right).$$

- 4 (2) Let \mathbf{A} be a fixed $m \times n$ matrix. If $\mathbb{E}(Z) = 0$ and $\text{var}(Z) = 1$, then there
5 exists an absolute constant $C > 0$ such that for any $t \geq 0$,

$$6 \quad \mathbb{P}(|\|\mathbf{AZ}\|_2 - \|\mathbf{A}\|_F| \geq t) \leq 2 \exp\left(-\frac{Ct^2}{K^4 \|\mathbf{A}\|_2^2}\right),$$

7 where $\|\mathbf{A}\|_F$ and $\|\mathbf{A}\|_2$ represent the Frobenius and L_2 norms of matrix \mathbf{A}
8 respectively.

- 9 (3) Let \mathbf{A} be a fixed $m \times n$ matrix. Let $\mathbf{e}_j \in \mathbb{R}^m$ be the unit vector with its j th
10 component one, $j = 1, \dots, m$. Suppose $M \equiv \max_{1 \leq j \leq m} n^{-1/2} \|\mathbf{A}^T \mathbf{e}_j\|_2 \in$
11 $(0, \infty)$ and $\rho \equiv \lambda_{\max}(n^{-1} \mathbf{A} \mathbf{A}^T) \in (0, \infty)$. If $\mathbb{E}(Z) = 0$ and $\nu = \text{var}(Z) \in$
12 $(0, \infty)$, then there exists an absolute constant $C > 0$ such that for any $t \geq 0$,

$$13 \quad \begin{aligned} & \mathbb{P}(\|n^{-1} \mathbf{AZ}\|_2 \geq t) \leq \Gamma(t; n, m, K, M, \rho, \nu) \\ & = 2m \exp\left(-\frac{Cnt^2}{K^2 M^2 m}\right) \wedge 2 \exp\left(-\frac{C\nu^2[(n^{1/2}t - \nu m^{1/2} \rho^{1/2})^+]^2}{K^4 \rho}\right). \end{aligned}$$

- 14 (4) The random variables Z^+ and Z^- are also sub-Gaussian. Moreover, for any
15 $c_1, c_2 \in \mathbb{R}$, $c_1 Z^+ + c_2 Z^-$ is sub-Gaussian.

16 **PROOF OF LEMMA 3.** (1) This part follows directly from Proposition 5.10
17 of [Vershynin \(2010\)](#).

18 (2) This part follows from Theorem 2.1 of [Rudelson and Vershynin \(2013\)](#).

19 (3) On one hand, we have by part (1) that

$$20 \quad \mathbb{P}\left(\left\|\frac{\mathbf{A}}{n} \mathbf{Z}\right\|_2 \geq t\right) \leq \mathbb{P}\left(\left\|\frac{\mathbf{A}}{\sqrt{n}} \mathbf{Z}\right\|_\infty \geq \frac{t\sqrt{n}}{\sqrt{m}}\right) \leq 2m \exp\left(-\frac{Cnt^2}{K^2 M^2 m}\right).$$

21 One the other hand, note that $\|n^{-1/2} \mathbf{A}\|_F = \sqrt{\text{Tr}(\mathbf{A} \mathbf{A}^T / n)} \leq \sqrt{m\rho}$ and
22 $\|n^{-1/2} \mathbf{A}\|_2^2 = \lambda_{\max}(\mathbf{A}^T \mathbf{A} / n) = \lambda_{\max}(\mathbf{A} \mathbf{A}^T / n) = \rho$. We have by part (2)

$$23 \quad \begin{aligned} & \mathbb{P}\left(\left\|\frac{\mathbf{A}}{n} \mathbf{Z}\right\|_2 \geq t\right) \leq \mathbb{P}\left(\left\|\frac{\mathbf{A}}{\sqrt{n}} \frac{\mathbf{Z}}{\nu}\right\|_2 - \left\|\frac{\mathbf{A}}{\sqrt{n}}\right\|_F \geq \frac{t\sqrt{n}}{\nu} - \sqrt{m\rho}\right) \\ & \leq \mathbb{P}\left(\left|\left\|\frac{\mathbf{A}}{\sqrt{n}} \frac{\mathbf{Z}}{\nu}\right\|_2 - \left\|\frac{\mathbf{A}}{\sqrt{n}}\right\|_F\right| \geq \left(\frac{t\sqrt{n}}{\nu} - \sqrt{m\rho}\right)^+\right) \\ & \leq 2 \exp\left(-\frac{C\nu^2[(n^{1/2}t - \nu m^{1/2} \rho^{1/2})^+]^2}{K^4 \rho}\right). \end{aligned}$$

1 (4) Note that by definition, we have $K \in (0, \infty)$ and $(\mathbb{E}|Z|^p)^{1/p} \leq K\sqrt{p}$,
 2 $\forall p \geq 1$. It follows immediately that $(\mathbb{E}|Z^+|^p)^{1/p} \leq (\mathbb{E}|Z|^p)^{1/p} \leq K\sqrt{p}$
 3 and $(\mathbb{E}|Z^-|^p)^{1/p} \leq (\mathbb{E}|Z|^p)^{1/p} \leq K\sqrt{p}$, $\forall p \geq 1$. Now by Lemma 5.5
 4 of Vershynin (2010), we conclude that Z^+ and Z^- are both sub-Gaussian.
 5 For any $c_1, c_2 \in \mathbb{R}$, by Minkowski inequality,

$$6 \quad (\mathbb{E}|c_1 Z^+ + c_2 Z^-|^p)^{1/p} \leq |c_1|(\mathbb{E}|Z^+|^p)^{1/p} + |c_2|(\mathbb{E}|Z^-|^p)^{1/p} \\ \leq (|c_1| + |c_2|)K\sqrt{p}, \quad \forall p \geq 1.$$

7 By Lemma 5.5 of Vershynin (2010) again, we can see that $c_1 Z^+ + c_2 Z^-$ is
 8 also sub-Gaussian. This completes the proof. \square

10 Now we are ready to prove Theorems 1 and 2. Lemmas 4 and 5 are presented to
 11 facilitate the proofs.

12 LEMMA 4. Let $\zeta = (\zeta_i, 1 \leq i \leq n)^T$ with $\zeta_i = \Psi'_\tau(\varepsilon_i) = 2|\tau - I(\varepsilon_i < 0)|\varepsilon_i$.

13 (1) For any $\beta, \delta \in \mathbb{R}^p$, $\langle \nabla \mathcal{L}_n(\beta + \delta) - \nabla \mathcal{L}_n(\beta), \delta \rangle \geq 2c\|\mathbf{X}\delta\|_2^2/n$.

14 (2) For any $d > 0$, $P(\|\hat{\beta}^{\text{oracle}} - \beta^*\|_2 \geq d) \leq P(\|n^{-1}\mathbf{X}_A^T \zeta\|_2 \geq 2c\rho_{\min}d)$.

15 PROOF OF LEMMA 4. The first part follows from the strong convexity of $\Psi_\tau(\cdot)$.
 16 Specifically, by Lemma 1, we have

$$17 \quad \mathcal{L}_n(\beta + \delta) - \mathcal{L}_n(\beta) - \langle \nabla \mathcal{L}_n(\beta), \delta \rangle \geq c\|\mathbf{X}\delta\|_2^2/n, \\ \mathcal{L}_n(\beta) - \mathcal{L}_n(\beta + \delta) - \langle \nabla \mathcal{L}_n(\beta + \delta), -\delta \rangle \geq c\|\mathbf{X}\delta\|_2^2/n.$$

18 Summing up the above two inequalities yields the desired result in part (1).

19 For the second part, let $\hat{\delta} = \hat{\beta}^{\text{oracle}} - \beta^*$. By definition of $\hat{\beta}^{\text{oracle}}$, we have
 20 $\hat{\delta}_{Ac} = \mathbf{0}$ and $\nabla_A \mathcal{L}_n(\hat{\beta}^{\text{oracle}}) = \mathbf{0}$. Now by part (1) we have

$$21 \quad 2c\rho_{\min}\|\hat{\delta}\|_2^2 = 2c\rho_{\min}\|\hat{\delta}_A\|_2^2 \leq 2c\hat{\delta}_A^T (\mathbf{X}_A^T \mathbf{X}_A/n) \hat{\delta}_A = 2c\|\mathbf{X}\hat{\delta}\|_2^2/n \\ \leq \langle \nabla \mathcal{L}_n(\hat{\beta}^{\text{oracle}}) - \nabla \mathcal{L}_n(\beta^*), \hat{\delta} \rangle = \langle -\nabla_A \mathcal{L}_n(\beta^*), \hat{\delta}_A \rangle \\ \leq \|\nabla_A \mathcal{L}_n(\beta^*)\|_2 \|\hat{\delta}_A\|_2 = \|n^{-1}\mathbf{X}_A^T \zeta\|_2 \|\hat{\delta}\|_2,$$

22 which implies that $2c\rho_{\min}\|\hat{\beta}^{\text{oracle}} - \beta^*\|_2 \leq \|n^{-1}\mathbf{X}_A^T \zeta\|_2$. The result of part (2)
 23 then follows. \square

24 PROOF OF THEOREM 1. Let $\hat{\delta} = \hat{\beta}^{\text{lasso}} - \beta^*$ and $z_\infty^* = \|\nabla \mathcal{L}_n(\beta^*)\|_\infty$. Note
 25 that $\hat{\beta}^{\text{lasso}}$ satisfies the Karush–Kuhn–Tucker (KKT) condition

$$26 \quad \nabla \mathcal{L}_n(\hat{\beta}^{\text{lasso}}) + \mathbf{g} = \mathbf{0},$$

1 where $g_j = \lambda_{\text{lasso}} \text{sgn}(\hat{\beta}_j^{\text{lasso}})$ if $\hat{\beta}_j^{\text{lasso}} \neq 0$ and $g_j \in [-\lambda_{\text{lasso}}, \lambda_{\text{lasso}}]$ if $\hat{\beta}_j^{\text{lasso}} = 0$. It
 2 follows that $\hat{\beta}_j^{\text{lasso}} g_j = \lambda_{\text{lasso}} |\hat{\beta}_j^{\text{lasso}}|, \forall j$. Since $\beta_{Ac}^* = \mathbf{0}$, we have $\hat{\delta}_{Ac} = \hat{\beta}_{Ac}^{\text{lasso}}$. By
 3 Lemma 4 and Hölder's inequality, we get

$$\begin{aligned} & 0 \leq 2\underline{c} \|\mathbf{X}\hat{\delta}\|_2^2/n \leq \langle \nabla \mathcal{L}_n(\hat{\beta}^{\text{lasso}}) - \nabla \mathcal{L}_n(\beta^*), \hat{\delta} \rangle = \langle -\mathbf{g} - \nabla \mathcal{L}_n(\beta^*), \hat{\delta} \rangle \\ (6.1) \quad & = \langle \hat{\delta}_A, -\mathbf{g}_A - \nabla_A \mathcal{L}_n(\beta^*) \rangle + \langle \hat{\beta}_{Ac}^{\text{lasso}}, -\mathbf{g}_{Ac} - \nabla_{Ac} \mathcal{L}_n(\beta^*) \rangle \\ & \leq (z_\infty^* + \lambda_{\text{lasso}}) \|\hat{\delta}_A\|_1 + (z_\infty^* - \lambda_{\text{lasso}}) \|\hat{\delta}_{Ac}\|_1. \end{aligned}$$

5 Under the event $\mathcal{E} = \{z_\infty^* \leq 2^{-1} \lambda_{\text{lasso}}\}$, from (6.1) we get

$$6 \quad \|\hat{\delta}_{Ac}\|_1 \leq \frac{z_\infty^* + \lambda_{\text{lasso}}}{z_\infty^* - \lambda_{\text{lasso}}} \|\hat{\delta}_A\|_1 \leq 3 \|\hat{\delta}_A\|_1,$$

7 which implies that $\hat{\delta} \in \mathcal{C}$. Now under \mathcal{E} , by condition (C3), it follows from (6.1)
 8 that

$$9 \quad 2\underline{c}\kappa \|\hat{\delta}\|_2^2 \leq (3/2) \lambda_{\text{lasso}} \|\hat{\delta}_A\|_1 \leq (3/2) \lambda_{\text{lasso}} s^{1/2} \|\hat{\delta}_A\|_2 \leq (3/2) \lambda_{\text{lasso}} s^{1/2} \|\hat{\delta}\|_2,$$

10 and similarly by condition (C4) and (6.1), we get

$$11 \quad 2\underline{c}\varrho \|\hat{\delta}\|_\infty \leq 2\underline{c} \|\mathbf{X}\hat{\delta}\|_2^2 / (n \|\hat{\delta}_A\|_1) \leq (3/2) \lambda_{\text{lasso}}.$$

12 Thus, we have

$$\begin{aligned} & \text{P}(\|\hat{\delta}\|_2 \leq 3s^{1/2} \lambda_{\text{lasso}} (4\underline{c}\kappa)^{-1} \cap \|\hat{\delta}\|_\infty \leq 3\lambda_{\text{lasso}} (4\underline{c}\varrho)^{-1}) \\ 13 \quad & \geq \text{P}(z_\infty^* \leq 2^{-1} \lambda_{\text{lasso}}) \geq 1 - \text{P}(\|n^{-1} \mathbf{X}^T \zeta\|_\infty \geq 2^{-1} \lambda_{\text{lasso}}). \end{aligned}$$

14 Note that $\zeta_i = \Psi_\tau'(\varepsilon_i) = 2\tau\varepsilon_i^+ - 2(1-\tau)\varepsilon_i^-$. It follows from Lemma 3 and
 15 $\mathcal{E}^\tau(\varepsilon_i) = 0$ that ζ_i are i.i.d. mean zero sub-Gaussian random variables. Now by the
 16 union bound argument and Lemma 3 again

$$17 \quad \text{P}(\|n^{-1} \mathbf{X}^T \zeta\|_\infty \geq 2^{-1} \lambda_{\text{lasso}}) \leq 2p \exp\left(-\frac{Cn\lambda_{\text{lasso}}^2}{4K_0^2 M_0^2}\right) = 1 - p_1^{\text{ALS}}.$$

18 This completes the proof. \square

19 **LEMMA 5.** *Under the assumptions of Theorem 2, the probability that the LLA*
 20 *algorithm (Algorithm 2) initialized by $\hat{\beta}^{\text{lasso}}$ converges to $\hat{\beta}^{\text{oracle}}$ after two iterations*
 21 *is at least $1 - p_1 - p_2 - p_3$, where*

$$\begin{aligned} & p_1 = \text{P}(\|\hat{\beta}^{\text{lasso}} - \beta^*\|_\infty > a_0\lambda), \\ 22 \quad & p_2 = \text{P}(\|\nabla_{Ac} \mathcal{L}_n(\hat{\beta}^{\text{oracle}})\|_\infty \geq a_1\lambda), \\ & p_3 = \text{P}(\min_{j \in A} |\hat{\beta}_j^{\text{oracle}}| < a\lambda). \end{aligned}$$

1 **PROOF OF LEMMA 5.** The convexity of $\mathcal{L}_n(\boldsymbol{\beta})$ follows from Lemma 1. Let
2 $\mathcal{S} = \{\boldsymbol{\beta} \in \mathbb{R}^p : \boldsymbol{\beta}_{A^c} = \mathbf{0}\}$. Note that $\widehat{\boldsymbol{\beta}}^{\text{oracle}} \in \mathcal{S}$. For any $\boldsymbol{\beta} \in \mathcal{S}$, let $\bar{\mathcal{L}}_n(\boldsymbol{\beta}_A) \equiv$
3 $n^{-1} \sum_{i=1}^n \Psi_\tau(y_i - \mathbf{x}_{iA}^\top \boldsymbol{\beta}_A) = \mathcal{L}_n(\boldsymbol{\beta})$. Then $\nabla \bar{\mathcal{L}}_n(\boldsymbol{\beta}_A) = -n^{-1} \sum_{i=1}^n \mathbf{x}_{iA} \Psi'_\tau(y_i -$
4 $\mathbf{x}_{iA}^\top \boldsymbol{\beta}_A)$. Now for any $\boldsymbol{\beta}$ and $\boldsymbol{\beta}' \in \mathcal{S}$, by Lemma 1 again, we get

$$5 \quad \bar{\mathcal{L}}_n(\boldsymbol{\beta}_A) \geq \bar{\mathcal{L}}_n(\boldsymbol{\beta}'_A) + \langle \nabla \bar{\mathcal{L}}_n(\boldsymbol{\beta}'_A), \boldsymbol{\beta}_A - \boldsymbol{\beta}'_A \rangle + \underline{c}(\boldsymbol{\beta}_A - \boldsymbol{\beta}'_A)^\top \frac{\mathbf{X}_A^\top \mathbf{X}_A}{n} (\boldsymbol{\beta}_A - \boldsymbol{\beta}'_A).$$

6 Since \mathbf{X}_A is of full column rank by assumption, we can see that $\bar{\mathcal{L}}_n(\boldsymbol{\beta}_A)$ is strongly
7 convex with respect to $\boldsymbol{\beta}_A$ and, therefore, $\widehat{\boldsymbol{\beta}}^{\text{oracle}}$ is the unique solution of prob-
8 lem (2.10) with $\nabla \bar{\mathcal{L}}_n(\widehat{\boldsymbol{\beta}}_A^{\text{oracle}}) = \mathbf{0}$. The lemma then follows from Theorems 1 and
9 2 in [Fan, Xue and Zou \(2014\)](#). \square

10 **PROOF OF THEOREM 2.** Let $\widehat{\boldsymbol{\delta}} = \widehat{\boldsymbol{\beta}}^{\text{lasso}} - \boldsymbol{\beta}^*$. Assume both (C3) and (C4) hold.
11 The other cases where either (C3) or (C4) holds are similar. From Lemma 5 and
12 Theorem 1, we immediately get

$$13 \quad \begin{aligned} p_1 &\leq \mathbb{P}(\|\widehat{\boldsymbol{\delta}}\|_\infty > [3s^{1/2}\lambda_{\text{lasso}}(4\kappa\underline{c})^{-1}] \wedge [3\lambda_{\text{lasso}}(4\underline{\rho}\underline{c})^{-1}]) \\ &\leq \mathbb{P}(\|\widehat{\boldsymbol{\delta}}\|_2 > 3s^{1/2}\lambda_{\text{lasso}}(4\kappa\underline{c})^{-1}) \vee \mathbb{P}(\|\widehat{\boldsymbol{\delta}}\|_\infty > 3\lambda_{\text{lasso}}(4\underline{\rho}\underline{c})^{-1}) \leq p_1^{\text{ALS}}. \end{aligned}$$

14 To derive the bound for p_2 , by the triangular inequality, it suffices to show bounds
15 for $\mathbb{P}(\|\nabla_{A^c} \mathcal{L}_n(\boldsymbol{\beta}^*)\|_\infty \geq 2^{-1}a_1\lambda)$ and $\mathbb{P}(\|\nabla_{A^c} \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^{\text{oracle}}) - \nabla_{A^c} \mathcal{L}_n(\boldsymbol{\beta}^*)\|_\infty \geq$
16 $2^{-1}a_1\lambda)$. By the union bound argument and Lemma 3,

$$17 \quad \begin{aligned} \mathbb{P}(\|\nabla_{A^c} \mathcal{L}_n(\boldsymbol{\beta}^*)\|_\infty \geq 2^{-1}a_1\lambda) &= \mathbb{P}(\| -n^{-1} \mathbf{X}_{A^c}^\top \boldsymbol{\zeta} \|_\infty \geq 2^{-1}a_1\lambda) \\ &\leq 2(p-s) \exp\left(-\frac{Ca_1^2 n \lambda^2}{4M_0^2 K_0^2}\right). \end{aligned}$$

18 Let $\mathbf{d} = (d_i, i = 1, \dots, n)^\top$ with $d_i = \Psi'_\tau(y_i - \mathbf{x}_i^\top \widehat{\boldsymbol{\beta}}^{\text{oracle}}) - \Psi'_\tau(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}^*)$. By
19 Cauchy–Schwarz inequality and Lemma 2, we get

$$20 \quad \begin{aligned} &\|\nabla_{A^c} \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^{\text{oracle}}) - \nabla_{A^c} \mathcal{L}_n(\boldsymbol{\beta}^*)\|_\infty \\ &= n^{-1} \max_{j \in A^c} |\sum_{i=1}^n d_i x_{ij}| \leq n^{-1} \max_{j \in A^c} (\|\mathbf{d}\|_2 \|X_j\|_2) \\ &\leq (2\bar{c}M_0) [(\widehat{\boldsymbol{\beta}}_A^{\text{oracle}} - \boldsymbol{\beta}_A^*)^\top (n^{-1} \mathbf{X}_A^\top \mathbf{X}_A) (\widehat{\boldsymbol{\beta}}_A^{\text{oracle}} - \boldsymbol{\beta}_A^*)]^{1/2} \\ &\leq (2\bar{c}\rho_{\max}^{1/2} M_0) \|\widehat{\boldsymbol{\beta}}^{\text{oracle}} - \boldsymbol{\beta}^*\|_2. \end{aligned}$$

21 It follows from Lemma 4 and Lemma 3 that

$$22 \quad \begin{aligned} &\mathbb{P}(\|\nabla_{A^c} \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^{\text{oracle}}) - \nabla_{A^c} \mathcal{L}_n(\boldsymbol{\beta}^*)\|_\infty \geq 2^{-1}a_1\lambda) \\ &\leq \mathbb{P}\left(\|\widehat{\boldsymbol{\beta}}^{\text{oracle}} - \boldsymbol{\beta}^*\|_2 \geq \frac{a_1\lambda}{4\bar{c}\rho_{\max}^{1/2} M_0}\right) \leq \mathbb{P}(\|n^{-1} \mathbf{X}_A^\top \boldsymbol{\zeta}\|_2 \geq Q_1\lambda) \\ &\leq \Gamma(Q_1\lambda; n, s, K_0, M_0, \rho_{\max}, \nu_0). \end{aligned}$$

1 This establishes the desired upper bound for p_2 . To show the upper bound for p_3 ,
 2 let $R = \min_{j \in A} |\beta_j^*| - a\lambda$ and observe that

$$\begin{aligned} p_3 &= \mathbb{P}(\min_{j \in A} |\hat{\beta}_j^{\text{oracle}}| < a\lambda) \leq \mathbb{P}(\|\hat{\boldsymbol{\beta}}^{\text{oracle}} - \boldsymbol{\beta}^*\|_\infty > R) \\ &\leq \mathbb{P}(\|\hat{\boldsymbol{\beta}}^{\text{oracle}} - \boldsymbol{\beta}^*\|_2 > R) \leq \mathbb{P}(\|n^{-1} \mathbf{X}_A^T \boldsymbol{\zeta}\|_2 \geq 2c\rho_{\min} R). \end{aligned}$$

3
 4 Similarly, by Lemma 3 we obtain

$$\mathbb{P}(\|n^{-1} \mathbf{X}_A^T \boldsymbol{\zeta}\|_2 \geq 2c\rho_{\min} R) \leq \Gamma(2c\rho_{\min} R; n, s, K_0, M_0, \rho_{\max}, \nu_0),$$

5
 6 which completes the proof. \square

7 Let us now prove the results for the COSALES estimation. To simplify notation,
 8 let $\boldsymbol{\varpi} = (\boldsymbol{\gamma}^T, \boldsymbol{\varphi}^T)^T$. It follows that $\text{supp}(\boldsymbol{\varpi}^*) = A_0$. Let $\lambda_{\text{lasso}} = \lambda_1^{\text{lasso}} \wedge \lambda_2^{\text{lasso}}$ and
 9 $\Lambda_{\text{lasso}} = \lambda_1^{\text{lasso}} \vee \lambda_2^{\text{lasso}}$. We first present a lemma to facilitate the proofs.

10 **LEMMA 6.** *Let $\boldsymbol{\varepsilon} = (\varepsilon_i, 1 \leq i \leq n)^T$ and $\boldsymbol{\eta} = (\eta_i, 1 \leq i \leq n)^T$, where
 11 $\eta_i = \Psi'_\tau(\varepsilon_i - e_\tau)$. Also, let $\mathbf{W} = \text{diag}\{\mathbf{x}_i^T \boldsymbol{\omega}^*, 1 \leq i \leq n\}$.*

12 (1) *For $\boldsymbol{\varpi}, \boldsymbol{\delta} \in \mathbb{R}^{2p}$, $\langle \nabla S_n(\boldsymbol{\varpi} + \boldsymbol{\delta}) - \nabla S_n(\boldsymbol{\varpi}), \boldsymbol{\delta} \rangle \geq n^{-1} c_0 \|(\mathbf{I}_2 \otimes \mathbf{X}) \boldsymbol{\delta}\|_2^2$, where
 13 \mathbf{I}_2 is a 2×2 identity matrix and $c_0 = 2^{-1}[(1 + 4\underline{c}) - (1 + 16\underline{c}^2)^{1/2}] > 0$.
 14 (2) For $d > 0$, $\mathbb{P}(\|\hat{\boldsymbol{\varpi}}^{\text{oracle}} - \boldsymbol{\varpi}^*\|_2 > d) \leq \mathbb{P}(\|\nabla_{A_0} S_n(\boldsymbol{\varpi}^*)\|_2 \geq c_0 \phi_{\min} d)$,
 15 where*

$$\nabla_{A_0} S_n(\boldsymbol{\varpi}^*) = -n^{-1} \begin{pmatrix} \mathbf{X}_{A_1}^T \mathbf{W}(\boldsymbol{\varepsilon} + \boldsymbol{\eta}) \\ \mathbf{X}_{A_2}^T \mathbf{W} \boldsymbol{\eta} \end{pmatrix}.$$

17 **PROOF OF LEMMA 6.** The first part follows directly from the strong convexity
 18 of the (asymmetric) squared error loss. Specifically, note that since c_0 is the smaller
 19 eigenvalue of the 2×2 matrix $\begin{pmatrix} 1 + 2\underline{c} & 2\underline{c} \\ 2\underline{c} & 2\underline{c} \end{pmatrix}$, we have

$$\begin{aligned} S_n(\boldsymbol{\varpi} + \boldsymbol{\delta}) - S_n(\boldsymbol{\varpi}) - \langle \nabla S_n(\boldsymbol{\varpi}), \boldsymbol{\delta} \rangle &\geq \frac{1}{2n} \boldsymbol{\delta}^T \left[\begin{pmatrix} 1 + 2\underline{c} & 2\underline{c} \\ 2\underline{c} & 2\underline{c} \end{pmatrix} \otimes (\mathbf{X}^T \mathbf{X}) \right] \boldsymbol{\delta} \\ &\geq (2n)^{-1} c_0 \|(\mathbf{I}_2 \otimes \mathbf{X}) \boldsymbol{\delta}\|_2^2. \end{aligned}$$

21 Similarly, $S_n(\boldsymbol{\varpi}) - S_n(\boldsymbol{\varpi} + \boldsymbol{\delta}) - \langle \nabla S_n(\boldsymbol{\varpi} + \boldsymbol{\delta}), -\boldsymbol{\delta} \rangle \geq (2n)^{-1} c_0 \|(\mathbf{I}_2 \otimes \mathbf{X}) \boldsymbol{\delta}\|_2^2$.
 22 Result (1) then follows by summing up the above two inequalities.

23 Let $\hat{\boldsymbol{\delta}} = \hat{\boldsymbol{\varpi}}^{\text{oracle}} - \boldsymbol{\varpi}^*$. Note that $\hat{\boldsymbol{\delta}}_{A_0^c} = \mathbf{0}$ and $\nabla_{A_0} S_n(\hat{\boldsymbol{\varpi}}^{\text{oracle}}) = \mathbf{0}$. From
 24 result (1) we have

$$\begin{aligned} c_0 \phi_{\min} \|\hat{\boldsymbol{\delta}}\|_2^2 &= c_0 \phi_{\min} \|\hat{\boldsymbol{\delta}}_{A_0}\|_2^2 \leq n^{-1} c_0 \|(\mathbf{I}_2 \otimes \mathbf{X}) \hat{\boldsymbol{\delta}}\|_2^2 \\ &\leq \langle \nabla S_n(\hat{\boldsymbol{\varpi}}^{\text{oracle}}) - \nabla S_n(\boldsymbol{\varpi}^*), \hat{\boldsymbol{\delta}} \rangle = \langle -\nabla_{A_0} S_n(\boldsymbol{\varpi}^*), \hat{\boldsymbol{\delta}}_{A_0} \rangle \\ &\leq \|\nabla_{A_0} S_n(\boldsymbol{\varpi}^*)\|_2 \|\hat{\boldsymbol{\delta}}\|_2^2. \end{aligned}$$

26 Result (2) follows immediately. \square

1 PROOF OF THEOREM 3. Let $\widehat{\boldsymbol{\delta}}_1 = \widehat{\boldsymbol{\gamma}}^{\text{lasso}} - \boldsymbol{\gamma}^*$, $\widehat{\boldsymbol{\delta}}_2 = \widehat{\boldsymbol{\varphi}}^{\text{lasso}} - \boldsymbol{\varphi}^*$, $\widehat{\boldsymbol{\delta}} = (\widehat{\boldsymbol{\delta}}_1^{\text{T}}, \widehat{\boldsymbol{\delta}}_2^{\text{T}})^{\text{T}}$,
 2 $z_{1\infty}^* = \|\partial S_n(\boldsymbol{\varpi}^*)/\partial \boldsymbol{\gamma}\|_{\infty}$, and $z_{2\infty}^* = \|\partial S_n(\boldsymbol{\varpi}^*)/\partial \boldsymbol{\varphi}\|_{\infty}$. By Lemma 6 and
 3 similar arguments in the proof of Theorem 1, it can shown that

$$\begin{aligned} 0 &\leq n^{-1}c_0\|(\mathbf{I}_2 \otimes \mathbf{X})\widehat{\boldsymbol{\delta}}\|_2^2 \leq \langle \nabla S_n(\widehat{\boldsymbol{\varpi}}^{\text{lasso}}) - \nabla S_n(\boldsymbol{\varpi}^*), \widehat{\boldsymbol{\delta}} \rangle \\ (6.2) \quad &\leq (z_{1\infty}^* + \lambda_1^{\text{lasso}})\|\widehat{\boldsymbol{\delta}}_{1A_1}\|_1 + (z_{1\infty}^* - \lambda_1^{\text{lasso}})\|\widehat{\boldsymbol{\delta}}_{1A_1^c}\|_1 \\ &\quad + (z_{2\infty}^* + \lambda_2^{\text{lasso}})\|\widehat{\boldsymbol{\delta}}_{2A_2}\|_1 + (z_{2\infty}^* - \lambda_2^{\text{lasso}})\|\widehat{\boldsymbol{\delta}}_{2A_2^c}\|_1. \end{aligned}$$

5 Under events $\mathcal{E}_1 = \{z_{1\infty}^* \leq 2^{-1}\lambda_1^{\text{lasso}}\}$ and $\mathcal{E}_2 = \{z_{2\infty}^* \leq 2^{-1}\lambda_2^{\text{lasso}}\}$, it follows
 6 from (6.2) that

$$\begin{aligned} 2^{-1}\lambda_{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{A_0^c}\|_1 &\leq 2^{-1}\lambda_1^{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{1A_1^c}\|_1 + 2^{-1}\lambda_2^{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{2A_2^c}\|_1 \\ &\leq (3/2)\lambda_1^{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{1A_1}\|_1 + (3/2)\lambda_2^{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{2A_2}\|_1 \leq (3/2)\Lambda_{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{A_0}\|_1, \end{aligned}$$

8 which implies that $\widehat{\boldsymbol{\delta}} \in \mathcal{C}_{3\bar{M}}$. Now under conditions (C4'-C5'), we have from (6.2)
 9 that

$$\begin{aligned} c_0\bar{\kappa}\|\widehat{\boldsymbol{\delta}}\|_2^2 &\leq n^{-1}c_0\|(\mathbf{I}_2 \otimes \mathbf{X})\widehat{\boldsymbol{\delta}}\|_2^2 \leq (3/2)\Lambda_{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{A_0}\|_1 \\ &\leq (3/2)\Lambda_{\text{lasso}}(s_1 + s_2)^{1/2}\|\widehat{\boldsymbol{\delta}}\|_2 \end{aligned}$$

11 and that

$$c_0\bar{\varrho}\|\widehat{\boldsymbol{\delta}}\|_{\infty}\|\widehat{\boldsymbol{\delta}}_{A_0}\|_1 \leq n^{-1}c_0\|(\mathbf{I}_2 \otimes \mathbf{X})\widehat{\boldsymbol{\delta}}\|_2^2 \leq (3/2)\Lambda_{\text{lasso}}\|\widehat{\boldsymbol{\delta}}_{A_0}\|_1.$$

13 It follows that under events \mathcal{E}_1 and \mathcal{E}_2 , we have $\|\widehat{\boldsymbol{\delta}}\|_2 \leq 3(s_1 + s_2)^{1/2}\Lambda_{\text{lasso}}(2\bar{\kappa}c_0)^{-1}$
 14 and $\|\widehat{\boldsymbol{\delta}}\|_{\infty} \leq 3\Lambda_{\text{lasso}}(2\bar{\varrho}c_0)^{-1}$. Recall that in Lemma 6 ε_i and $\eta_i = \Psi'_{\tau}(\varepsilon_i - e_{\tau})$
 15 are both mean zero sub-Gaussian random variables with $K_1 = \|\varepsilon_i\|_{\text{SG}}$ and $K_2 =$
 16 $\|\eta_i\|_{\text{SG}}$. It follows that $\varepsilon_i + \eta_i$ is also sub-Gaussian, and moreover, $\|\varepsilon_i + \eta_i\|_{\text{SG}} \leq$
 17 $K_1 + K_2$. Since $M_1 = \|\mathbf{X}\boldsymbol{\omega}^*\|_{\infty}$, we have

$$\begin{aligned} &\text{P}(\|\widehat{\boldsymbol{\delta}}\|_2 \leq 3(s_1 + s_2)^{1/2}\Lambda_{\text{lasso}}(2\bar{\kappa}c_0)^{-1} \cap \|\widehat{\boldsymbol{\delta}}\|_{\infty} \leq 3\Lambda_{\text{lasso}}(2\bar{\varrho}c_0)^{-1}) \\ &\geq \text{P}(\mathcal{E}_1 \cap \mathcal{E}_2) \geq 1 - \text{P}(\mathcal{E}_1^c) - \text{P}(\mathcal{E}_2^c) \\ &= 1 - \text{P}(\|n^{-1}\mathbf{X}^{\text{T}}\mathbf{W}(\boldsymbol{\varepsilon} + \boldsymbol{\eta})\|_{\infty} > 2^{-1}\lambda_1^{\text{lasso}}) - \text{P}(\|n^{-1}\mathbf{X}^{\text{T}}\mathbf{W}\boldsymbol{\eta}\|_{\infty} > 2^{-1}\lambda_2^{\text{lasso}}) \\ &\geq 1 - 2p \exp\left(-\frac{Cn(\lambda_1^{\text{lasso}})^2}{4M_0^2M_1^2(K_1 + K_2)^2}\right) - 2p \exp\left(-\frac{Cn(\lambda_2^{\text{lasso}})^2}{4M_0^2M_1^2K_2^2}\right). \end{aligned}$$

19 Theorem 3 then follows. \square

20 The proof of Theorem 4 relies on the following lemma.

1 LEMMA 7. Under assumptions of Theorem 4, the LLA algorithm (Algorithm 4)
 2 initialized by $\hat{\gamma}^{\text{lasso}}$ and $\hat{\varphi}^{\text{lasso}}$ converges to the oracle estimators $\hat{\gamma}^{\text{oracle}}$ and $\hat{\varphi}^{\text{oracle}}$
 3 in two iterations with probability at least $1 - \pi_1 - \pi_2 - \pi_3$, where

$$\pi_1 = \text{P}(\|\hat{\gamma}^{\text{lasso}} - \gamma^*\|_\infty > a_0\lambda_1, \|\hat{\varphi}^{\text{lasso}} - \varphi^*\|_\infty > a_0\lambda_2),$$

$$4 \quad \pi_2 = \text{P}(\|\partial S_n(\widehat{\varpi}^{\text{oracle}})/\partial \gamma_{A_1^c}\|_\infty \geq a_1\lambda_1, \|\partial S_n(\widehat{\varpi}^{\text{oracle}})/\partial \varphi_{A_2^c}\|_\infty \geq a_1\lambda_2),$$

$$\pi_3 = \text{P}(\min_{j \in A_1} |\hat{\gamma}_j^{\text{oracle}}| < a\lambda_1, \min_{j \in A_2} |\hat{\varphi}_j^{\text{oracle}}| < a\lambda_2).$$

5 PROOF OF LEMMA 7. The convexity of $S_n(\gamma, \varphi)$ follows immediately from
 6 Lemma 1,

$$7 \quad S_n(\gamma, \varphi) \geq S_n(\gamma', \varphi') + \langle \nabla_\gamma S_n(\gamma', \varphi'), \gamma - \gamma' \rangle + \langle \nabla_\varphi S_n(\gamma', \varphi'), \varphi - \varphi' \rangle \\
+ 2^{-1} \begin{pmatrix} \gamma - \gamma' \\ \varphi - \varphi' \end{pmatrix}^\top \left[\begin{pmatrix} 1 + 2\underline{c} & 2\underline{c} \\ 2\underline{c} & 2\underline{c} \end{pmatrix} \otimes (n^{-1} \mathbf{X}^\top \mathbf{X}) \right] \begin{pmatrix} \gamma - \gamma' \\ \varphi - \varphi' \end{pmatrix}.$$

8 Restrict $S_n(\gamma, \varphi)$ to the set $\mathcal{S} = \{\gamma, \varphi \in \mathbb{R}^p : \gamma_{A_1^c} = \mathbf{0}, \varphi_{A_2^c} = \mathbf{0}\}$ and define for
 9 any $(\gamma, \varphi) \in \mathcal{S}$

$$10 \quad \check{S}_n(\gamma_{A_1}, \varphi_{A_2}) = n^{-1} \sum_{i=1}^n \{\Psi_{0.5}(y_i - \mathbf{x}_{iA_1}^\top \gamma_{A_1}) + \Psi_\tau(y_i - \mathbf{x}_{iA_1}^\top \gamma_{A_1} - \mathbf{x}_{iA_2}^\top \varphi_{A_2})\}.$$

11 It follows immediately that for any $(\gamma, \varphi), (\gamma', \varphi') \in \mathcal{S}$,

$$\check{S}_n(\gamma_{A_1}, \varphi_{A_2}) \geq \check{S}_n(\gamma'_{A_1}, \varphi'_{A_2}) + \langle \nabla_{\gamma_{A_1}} \check{S}_n(\gamma'_{A_1}, \varphi'_{A_2}), \gamma_{A_1} - \gamma'_{A_1} \rangle \\
+ \langle \nabla_{\varphi_{A_2}} \check{S}_n(\gamma'_{A_1}, \varphi'_{A_2}), \varphi_{A_2} - \varphi'_{A_2} \rangle \\
12 \quad + 2^{-1} c_0 (\gamma_{A_1} - \gamma'_{A_1})^\top (n^{-1} \mathbf{X}_{A_1}^\top \mathbf{X}_{A_1}) (\gamma_{A_1} - \gamma'_{A_1}) \\
+ 2^{-1} c_0 (\varphi_{A_2} - \varphi'_{A_2})^\top (n^{-1} \mathbf{X}_{A_2}^\top \mathbf{X}_{A_2}) (\varphi_{A_2} - \varphi'_{A_2}),$$

13 where $c_0 = 2^{-1}[(1 + 4\underline{c}) - (1 + 16\underline{c}^2)^{1/2}]$. Since both \mathbf{X}_{A_1} and \mathbf{X}_{A_2} are of full
 14 column ranks by assumption, we can see that $\check{S}_n(\gamma_{A_1}, \varphi_{A_2})$ is strongly convex and
 15 thus the oracle estimators $\hat{\gamma}^{\text{oracle}}$ and $\hat{\varphi}^{\text{oracle}}$ are the unique solution of problem (4.3).

16 Let \mathcal{E}_1 be the event that $\|\hat{\gamma}^{\text{lasso}} - \gamma^*\|_\infty \leq a_0\lambda_1$ and $\|\hat{\varphi}^{\text{lasso}} - \varphi^*\|_\infty \leq$
 17 $a_0\lambda_2$. Under \mathcal{E}_1 and Assumption (A0'), on one hand we have $\min_{j \in A_1} |\hat{\gamma}_j^{\text{lasso}}| \geq$
 18 $\min_{j \in A_1} |\gamma_j^*| - \|\hat{\gamma}^{\text{lasso}} - \gamma^*\|_\infty > a\lambda_1$, implying that $p'_{\lambda_1}(|\hat{\gamma}_j^{\text{lasso}}|) = 0$ for $j \in A_1$.
 19 On the other hand, we have $\|\hat{\gamma}_{A_1^c}^{\text{lasso}}\|_\infty \leq \|\hat{\gamma}^{\text{lasso}} - \gamma^*\|_\infty \leq a_2\lambda_1$, indicating that
 20 $p'_{\lambda_1}(|\hat{\gamma}_j^{\text{lasso}}|) \geq a_1\lambda_1$ for $j \in A_1^c$. Similarly, we can show that $p'_{\lambda_2}(|\hat{\varphi}_j^{\text{lasso}}|) = 0$ for
 21 $j \in A_2$ and $p'_{\lambda_2}(|\hat{\varphi}_j^{\text{lasso}}|) \geq a_1\lambda_2$ for $j \in A_2^c$.

22 Let $\hat{\gamma}^1$ and $\hat{\varphi}^1$ be the update after the first iteration of the LLA algorithm. Then
 23 under \mathcal{E}_1 , $\hat{\gamma}^1$ and $\hat{\varphi}^1$ are minimizers of

$$24 \quad \mathcal{Q}_n(\gamma, \varphi) = S_n(\gamma, \varphi) + \sum_{j \in A_1^c} p'_{\lambda_1}(|\hat{\gamma}_j^{\text{lasso}}|) |\gamma_j| + \sum_{j \in A_2^c} p'_{\lambda_2}(|\hat{\varphi}_j^{\text{lasso}}|) |\varphi_j|.$$

1 By definition of the oracle estimators, $\partial S_n(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}})/\partial \gamma_j = 0$ for $j \in A_1$
 2 and $\partial S_n(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}})/\partial \varphi_j = 0$ for $j \in A_2^c$. Also, $\hat{\gamma}_{A_1^c}^{\text{oracle}} = \mathbf{0}$ and $\hat{\varphi}_{A_2^c}^{\text{oracle}} =$
 3 $\mathbf{0}$. Now let \mathcal{E}_2 be the event that $\max_{j \in A_1^c} |\partial \mathcal{L}(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}})/\partial \gamma_j| < a_1 \lambda_1$ and
 4 that $\max_{j \in A_2^c} |\partial \mathcal{L}(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}})/\partial \varphi_j| < a_1 \lambda_2$. It follows from the convexity of
 5 $S_n(\gamma, \varphi)$ that

$$\begin{aligned} & Q_n(\gamma, \varphi) - Q_n(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}}) \\ & \geq \sum_{j \in A_1^c} \frac{\partial}{\partial \gamma_j} S_n(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}}) \gamma_j + \sum_{j \in A_2^c} \frac{\partial}{\partial \varphi_j} S_n(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}}) \varphi_j \\ & \quad + \sum_{j \in A_1^c} p'_{\lambda_1}(|\hat{\gamma}_j^{\text{lasso}}|) |\gamma_j| + \sum_{j \in A_2^c} p'_{\lambda_2}(|\hat{\varphi}_j^{\text{lasso}}|) |\varphi_j|. \end{aligned}$$

7 Under \mathcal{E}_2 , this implies that $Q_n(\gamma, \varphi) \geq Q_n(\hat{\gamma}^{\text{oracle}}, \hat{\varphi}^{\text{oracle}})$ for any $\gamma \in \mathbb{R}^p$ and
 8 $\varphi \in \mathbb{R}^p$. The strict inequality holds unless $\gamma_j = 0$ for all $j \in A_1^c$ and $\varphi_j = 0$ for all
 9 $j \in A_2^c$. By the uniqueness of the oracle estimators, we must have $\hat{\gamma}^1 = \hat{\gamma}^{\text{oracle}}$ and
 10 $\hat{\varphi}^1 = \hat{\varphi}^{\text{oracle}}$.

11 Let \mathcal{E}_3 be the event that $\min_{j \in A_1} |\hat{\gamma}_j^{\text{oracle}}| \geq a \lambda_1$ and $\min_{j \in A_2} |\hat{\varphi}_j^{\text{oracle}}| \geq a \lambda_2$.
 12 Once the oracle estimators are obtained after the first iteration, under \mathcal{E}_3 , we can
 13 see that $p'_{\lambda_1}(|\hat{\gamma}_j^{\text{oracle}}|) = 0$ for $j \in A_1$, $p'_{\lambda_1}(|\hat{\gamma}_j^{\text{oracle}}|) \geq a_1 \lambda_1$ for $j \in A_1^c$ and
 14 $p'_{\lambda_2}(|\hat{\varphi}_j^{\text{oracle}}|) = 0$ for $j \in A_2$, $p'_{\lambda_2}(|\hat{\varphi}_j^{\text{oracle}}|) \geq a_1 \lambda_2$ for $j \in A_2^c$. By similar
 15 arguments, it can be shown that the second iteration of the LLA algorithm will
 16 still yield the oracle estimators, which means the algorithm converges to the oracle
 17 estimators hereafter. This completes the proof. \square

18 **PROOF OF THEOREM 4.** Let $\hat{\delta} = \widehat{\varpi}^{\text{lasso}} - \varpi^*$. Assume both (C4') and (C5')
 19 hold. The other cases where either (C4') or (C5') holds are similar. It follows from
 20 Theorem 3 that

$$\begin{aligned} \pi_1 & \leq \mathbb{P}(\|\hat{\delta}\|_\infty > a_0 \lambda) \leq \mathbb{P}(\|\hat{\delta}\|_\infty > 3\Lambda_{\text{lasso}}(2c_0)^{-1} [(s^{1/2} \bar{\kappa}^{-1}) \wedge \bar{\varrho}^{-1}]) \\ & \leq \mathbb{P}(\|\hat{\delta}\|_2 > 3s^{1/2} \Lambda_{\text{lasso}}(2c_0 \bar{\kappa})^{-1}) \vee \mathbb{P}(\|\hat{\delta}\|_\infty > 3\Lambda_{\text{lasso}}(2c_0 \bar{\varrho})^{-1}) \leq \pi_1^{\text{ALS}}. \end{aligned}$$

22 Next, note that $\pi_2 \leq \mathbb{P}(\|\nabla_{A_0^c} S_n(\widehat{\varpi}^{\text{oracle}})\|_\infty \geq a_1 \lambda)$. By the triangular inequality,
 23 it suffices to show upper bounds for respectively $\mathbb{P}(\|\nabla_{A_0^c} S_n(\varpi^*)\|_\infty \geq 2^{-1} a_1 \lambda)$
 24 and $\mathbb{P}(\|\nabla_{A_0^c} S_n(\widehat{\varpi}^{\text{oracle}}) - \nabla_{A_0^c} S_n(\varpi^*)\|_\infty \geq 2^{-1} a_1 \lambda)$. First, by the union bound
 25 argument we have

$$\begin{aligned} & \mathbb{P}(\|\nabla_{A_0^c} S_n(\varpi^*)\|_\infty \geq 2^{-1} a_1 \lambda) \\ & \leq \mathbb{P}(\|n^{-1} \mathbf{X}_{A_1^c}^T \mathbf{W}(\varepsilon + \eta)\|_\infty \geq 2^{-1} a_1 \lambda) + \mathbb{P}(\|n^{-1} \mathbf{X}_{A_2^c}^T \mathbf{W} \eta\|_\infty \geq 2^{-1} a_1 \lambda) \\ & \leq 2(p - s_1) \exp\left(-\frac{Ca_1^2 n \lambda^2}{4M_0^2 M_1^2 (K_1 + K_2)^2}\right) + 2(p - s_2) \exp\left(-\frac{Ca_1^2 n \lambda^2}{4M_0^2 M_1^2 K_2^2}\right). \end{aligned}$$

1 Now let $\bar{d}_i = \Psi'_\tau(y_i - \mathbf{x}_i^\top \hat{\gamma}^{\text{oracle}} - \mathbf{x}_i^\top \hat{\varphi}^{\text{oracle}}) - \Psi'_\tau(y_i - \mathbf{x}_i^\top \gamma^* - \mathbf{x}_i^\top \varphi^*)$ and set
 2 $\bar{\mathbf{d}} = (\bar{d}_i, 1 \leq i \leq n)^\top$. It follows that

$$\begin{aligned} & \|\nabla_{A_0} S_n(\widehat{\boldsymbol{\omega}}^{\text{oracle}}) - \nabla_{A_0} S_n(\boldsymbol{\omega}^*)\|_\infty \leq M_0(\|\mathbf{X}(\hat{\gamma}^{\text{oracle}} - \gamma^*)\|_2 + \|\bar{\mathbf{d}}\|_2)/\sqrt{n} \\ & \leq M_0[(1 + 2\bar{c})\|\mathbf{X}_{A_1}(\hat{\gamma}_{A_1}^{\text{oracle}} - \gamma_{A_1}^*)\|_2 + (2\bar{c})\|\mathbf{X}_{A_2}(\hat{\varphi}_{A_2}^{\text{oracle}} - \varphi_{A_2}^*)\|_2]/\sqrt{n} \\ & \leq (1 + 2\bar{c})M_0\phi_{\max}^{1/2}\|\widehat{\boldsymbol{\omega}}^{\text{oracle}} - \boldsymbol{\omega}^*\|_2. \end{aligned}$$

4 By Lemma 6 and Lemma 3, we get

$$\begin{aligned} & \mathbb{P}(\|\nabla_{A_0} S_n(\widehat{\boldsymbol{\omega}}^{\text{oracle}}) - \nabla_{A_0} S_n(\boldsymbol{\omega}^*)\|_\infty \geq 2^{-1}a_1\lambda) \\ & \leq \mathbb{P}\left(\|\widehat{\boldsymbol{\omega}}^{\text{oracle}} - \boldsymbol{\omega}^*\|_2 \geq \frac{a_1\lambda}{2(1 + 2\bar{c})M_0\phi_{\max}^{1/2}}\right) \\ & \leq \mathbb{P}\left(\left\|\frac{1}{n}\begin{pmatrix} \mathbf{X}_{A_1}^\top \mathbf{W}(\boldsymbol{\varepsilon} + \boldsymbol{\eta}) \\ \mathbf{X}_{A_2}^\top \mathbf{W}\boldsymbol{\eta} \end{pmatrix}\right\|_2 \geq Q_2\lambda\right) \\ & \leq \mathbb{P}(\|n^{-1}\mathbf{X}_{A_1}^\top \mathbf{W}(\boldsymbol{\varepsilon} + \boldsymbol{\eta})\|_2 \geq 2^{-1}Q_2\lambda) + \mathbb{P}(\|n^{-1}\mathbf{X}_{A_2}^\top \mathbf{W}\boldsymbol{\eta}\|_2 \geq 2^{-1}Q_2\lambda) \\ & \leq \Gamma(2^{-1}Q_2\lambda; n, s_1, K_1 + K_2, M_0M_1, M_1^2\rho_{1\bullet\max}, \nu_1) \\ & \quad + \Gamma(2^{-1}Q_2\lambda; n, s_2, K_2, M_0M_1, M_1^2\rho_{2\bullet\max}, \nu_2). \end{aligned}$$

6 This completes the upper bound for π_2 . To derive the upper bound for π_3 , note
 7 that by Assumption (A0') we have $\min_{j \in A_1} |\gamma_j^*| \geq (a + 1)\lambda_1$ and $\min_{j \in A_2} |\varphi_j^*| \geq$
 8 $(a + 1)\lambda_2$. Observe that $\min_{j \in A_1} |\hat{\gamma}_j^{\text{oracle}}| \geq \min_{j \in A_1} |\gamma_j^*| - \|\hat{\gamma}^{\text{oracle}} - \gamma^*\|_\infty$ and
 9 $\min_{j \in A_2} |\hat{\varphi}_j^{\text{oracle}}| \geq \min_{j \in A_2} |\varphi_j^*| - \|\hat{\varphi}^{\text{oracle}} - \varphi^*\|_\infty$, and it follows that

$$\begin{aligned} \pi_3 & \leq \mathbb{P}\left(\|\widehat{\boldsymbol{\omega}}^{\text{oracle}} - \boldsymbol{\omega}^*\|_\infty > \bar{R}\right) \leq \mathbb{P}\left(\|\widehat{\boldsymbol{\omega}}^{\text{oracle}} - \boldsymbol{\omega}^*\|_2 > \bar{R}\right) \\ & \leq \mathbb{P}\left(\left\|\frac{1}{n}\begin{pmatrix} \mathbf{X}_{A_1}^\top \mathbf{W}(\boldsymbol{\varepsilon} + \boldsymbol{\eta}) \\ \mathbf{X}_{A_2}^\top \mathbf{W}\boldsymbol{\eta} \end{pmatrix}\right\|_2 \geq c_0\phi_{\min}\bar{R}\right) \\ & \leq \mathbb{P}(\|n^{-1}\mathbf{X}_{A_1}^\top \mathbf{W}(\boldsymbol{\varepsilon} + \boldsymbol{\eta})\|_2 \geq \frac{1}{2}c_0\phi_{\min}\bar{R}) + \mathbb{P}(\|n^{-1}\mathbf{X}_{A_2}^\top \mathbf{W}\boldsymbol{\eta}\|_2 \geq \frac{1}{2}c_0\phi_{\min}\bar{R}) \\ & \leq \Gamma(2^{-1}c_0\phi_{\min}\bar{R}; n, s_1, K_1 + K_2, M_0M_1, M_1^2\rho_{1\bullet\max}, \nu_1) \\ & \quad + \Gamma(2^{-1}c_0\phi_{\min}\bar{R}; n, s_2, K_2, M_0M_1, M_1^2\rho_{2\bullet\max}, \nu_2). \end{aligned}$$

11

□

12 **Acknowledgments.** The authors sincerely thank the editor, associate editor,
 13 and three referees for their helpful comments and suggestions that led to substantial
 14 improvement of the paper.

15

SUPPLEMENTARY MATERIAL

1 **Supplement to “High-dimensional generalizations of asymmetric least squares**
2 **regression and their applications”**. The supplementary material includes the iter-
3 ation complexity analysis of the SALES algorithm.

4 **References.**

- 5 BICKEL, P. J., RITOV, Y. and TSYBAKOV, A. B. (2009). Simultaneous analysis of Lasso and Dantzig
6 selector. *The Annals of Statistics* **37** 1705–1732.
- 7 CANDÈS, E. and TAO, T. (2007). The Dantzig selector: Statistical estimation when p is much larger
8 than n . *The Annals of Statistics* **35** 2313–2351.
- 9 CHAMBERS, R. and TZAVIDIS, N. (2006). M -quantile models for small area estimation. *Biometrika*
10 **93** 255–268.
- 11 DAYE, Z. J., CHEN, J. and LI, H. (2012). High-Dimensional Heteroscedastic Regression with an
12 Application to eQTL Data Analysis. *Biometrics* **68** 316–326.
- 13 EFRON, B. (1991). Regression percentiles using asymmetric squared error loss. *Statistica Sinica* **1**
14 93–125.
- 15 EFRON, B., HASTIE, T., JOHNSTONE, I., TIBSHIRANI, R. et al. (2004). Least angle regression. *The*
16 *Annals of Statistics* **32** 407–499.
- 17 EILERS, P. H. and BOELEN, H. F. (2005). Baseline correction with asymmetric least squares
18 smoothing. *Leiden University Medical Centre Report*.
- 19 FAN, J. and LI, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle
20 properties. *Journal of the American Statistical Association* **96** 1348–1360.
- 21 FAN, J. and LV, J. (2011). Nonconcave penalized likelihood with NP-dimensionality. *IEEE Transac-*
22 *tions on Information Theory* **57** 5467–5484.
- 23 FAN, J., XUE, L. and ZOU, H. (2014). Strong oracle optimality of folded concave penalized estimation.
24 *The Annals of Statistics* **42** 819–849.
- 25 FRIEDMAN, J., HASTIE, T. and TIBSHIRANI, R. (2010). Regularization paths for generalized linear
26 models via coordinate descent. *Journal of Statistical Software* **33** 1–22.
- 27 GU, Y. and ZOU, H. (2015). Supplement to “High-Dimensional Generalizations of Asymmetric Least
28 Squares Regression and Their Applications”.
- 29 HUANG, J. and ZHANG, C.-H. (2012). Estimation and selection via absolute penalized convex
30 minimization and its multistage adaptive applications. *The Journal of Machine Learning Research*
31 **13** 1839–1864.
- 32 KOENKER, R. and BASSETT, G. (1978). Regression quantiles. *Econometrica: journal of the Econo-*
33 *metric Society* **46** 33–50.
- 34 KOENKER, R. and BASSETT, G. (1982). Robust tests for heteroscedasticity based on regression
35 quantiles. *Econometrica: Journal of the Econometric Society* **50** 43–61.
- 36 KOENKER, R. and ZHAO, Q. (1994). L -estimation for linear heteroscedastic models. *Journal of*
37 *Nonparametric Statistics* **3** 223–235.
- 38 KUAN, C.-M., YEH, J.-H. and HSU, Y.-C. (2009). Assessing value at risk with CARE, the condi-
39 tional autoregressive expectile models. *Journal of Econometrics* **150** 261–270.
- 40 MEIER, L., VAN DE GEER, S., BÜHLMANN, P. et al. (2009). High-dimensional additive modeling.
41 *The Annals of Statistics* **37** 3779–3821.
- 42 NEGAHBAN, S. N., RAVIKUMAR, P., WAINWRIGHT, M. J. and YU, B. (2012). A Unified Framework
43 for High-Dimensional Analysis of M -Estimators with Decomposable Regularizers. *Statistical*
44 *Science* **27** 538–557.
- 45 NEWEY, W. K. and POWELL, J. L. (1987). Asymmetric least squares estimation and testing. *Econo-*
46 *metrica* **55** 819–847.
- 47 PARIKH, N. and BOYD, S. (2013). Proximal algorithms. *Foundations and Trends in optimization* **1**
48 123–231.

- 1 RUDELSON, M. and VERSHYNIN, R. (2013). Hanson-Wright inequality and sub-gaussian concentra-
 2 tion. *Electron. Commun. Probab.* **18** 1-9.
- 3 SALVATI, N., TZAVIDIS, N., PRATESI, M. and CHAMBERS, R. (2012). Small area estimation via
 4 M-quantile geographically weighted regression. *Test* **21** 1–28.
- 5 SCHEETZ, T. E., KIM, K.-Y. A., SWIDERSKI, R. E., PHILP, A. R., BRAUN, T. A., KNUDT-
 6 SON, K. L., DORRANCE, A. M., DiBONA, G. F., HUANG, J., CASAVANT, T. L. et al. (2006).
 7 Regulation of gene expression in the mammalian eye and its relevance to eye disease. *Proceedings*
 8 *of the National Academy of Sciences* **103** 14429–14434.
- 9 TAYLOR, J. W. (2008). Estimating value at risk and expected shortfall using expectiles. *Journal of*
 10 *Financial Econometrics* **6** 231–252.
- 11 TIBSHIRANI, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal*
 12 *Statistical Society. Series B (Methodological)* **58** 267–288.
- 13 TSENG, P. (2001). Convergence of a block coordinate descent method for nondifferentiable minimiza-
 14 tion. *Journal of optimization theory and applications* **109** 475–494.
- 15 VERSHYNIN, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv*
 16 *preprint arXiv:1011.3027v7*.
- 17 WANG, L., WU, Y. and LI, R. (2012). Quantile regression for analyzing heterogeneity in ultra-high
 18 dimension. *Journal of the American Statistical Association* **107** 214–222.
- 19 WANG, L., KIM, Y., LI, R. et al. (2013). Calibrating nonconvex penalized regression in ultra-high
 20 dimension. *The Annals of Statistics* **41** 2505–2536.
- 21 XIE, S., ZHOU, Y. and WAN, A. T. (2014). A varying-coefficient expectile model for estimating
 22 Value at Risk. *Journal of Business & Economic Statistics* **32** 576–592.
- 23 YANG, Y. and ZOU, H. (2013). An efficient algorithm for computing the HHSVM and its generaliza-
 24 tions. *Journal of Computational and Graphical Statistics* **22** 396–415.
- 25 YE, F. and ZHANG, C.-H. (2010). Rate Minimality of the Lasso and Dantzig Selector for the ℓ_q
 26 Loss in ℓ_r Balls. *The Journal of Machine Learning Research* **11** 3519–3540.
- 27 ZHANG, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The*
 28 *Annals of Statistics* **38** 894–942.
- 29 ZHAO, P. and YU, B. (2006). On model selection consistency of Lasso. *The Journal of Machine*
 30 *Learning Research* **7** 2541–2563.
- 31 ZOU, H. (2006). The Adaptive Lasso and Its Oracle Properties. *Journal of the American Statistical*
 32 *Association* **101** 1418-1429.
- 33 ZOU, H. and LI, R. (2008). One-step sparse estimates in nonconcave penalized likelihood models.
 34 *Annals of statistics* **36** 1509–1533.

35 YUWEN GU
 SCHOOL OF STATISTICS,
 UNIVERSITY OF MINNESOTA,
 MINNEAPOLIS, MN, USA
 E-MAIL: guxxx192@umn.edu

HUI ZOU
 SCHOOL OF STATISTICS,
 UNIVERSITY OF MINNESOTA,
 MINNEAPOLIS, MN, USA
 E-MAIL: zouxx019@umn.edu