

Robust Estimation of High-Dimensional Mean Regression*

Jianqing Fan, Quefeng Li, Yuyan Wang

Department of Operations Research and Financial Engineering,
Princeton University, Princeton, NJ 08544

October 6, 2014

Abstract

Data subject to heavy-tailed errors are commonly encountered in various scientific fields, especially in the modern era with explosion of massive data. To address this problem, procedures based on quantile regression and Least Absolute Deviation (LAD) regression have been developed in recent years. These methods essentially estimate the conditional median (or quantile) function. They can be very different from the conditional mean functions when distributions are asymmetric and heteroscedastic. How can we efficiently estimate the mean regression functions in ultra-high dimensional setting with existence of only the second moment? To solve this problem, we propose a penalized Huber loss with diverging parameter to reduce biases created by the traditional Huber loss. Such a penalized robust approximate quadratic (RA-quadratic) loss will be called RA-Lasso. In the ultra-high dimensional setting, where the dimensionality can grow exponentially with the sample size, our results reveal that the RA-lasso estimator produces a consistent estimator at the same rate as the optimal rate under the light-tail situation. We further study the computational convergence of RA-Lasso and show that the composite gradient descent algorithm indeed produces a solution that admits the same optimal rate after sufficient iterations. As a byproduct, we also establish the concentration inequality for estimating population mean when there exists only the second moment. We compare RA-Lasso with

*Supported in part by NSF Grants DMS-1206464 and DMS-1406266 and NIH grants R01-GM072611-9 and NIH R01GM100474-4.

other regularized robust estimators based on quantile regression and LAD regression. Extensive simulation studies demonstrate the satisfactory finite-sample performance of RA-Lasso.

Key words: High dimension, Huber loss, M-estimator, Optimal rate, Robust regularization.

1 Introduction

Our era has witnessed the massive explosion of data and a dramatic improvement of technology in collecting and processing large data sets. We often encounter huge data sets that the number of features greatly surpasses the number of observations. It makes many traditional statistical analysis tools infeasible and poses great challenge on developing new tools. Regularization methods have been widely used for the analysis of high-dimensional data. These methods penalize the least squares or the likelihood function with the L_1 -penalty on the unknown parameters (Lasso, Tibshirani (1996)), or a folded concave penalty function such as the SCAD (Fan and Li, 2001) and the MCP (Zhang, 2010). However, these penalized least-squares methods are sensitive to the tails of the error distributions, particularly for ultrahigh dimensional covariates, as the maximum spurious correlation between the covariates and the realized noises can be large in those cases. As a result, theoretical properties are often obtained under light-tailed error distributions (Bickel, Ritov and Tsybakov, 2009; Fan and Lv, 2011).

To tackle the problem of heavy-tailed errors, robust regularization methods have been extensively studied. Li and Zhu (2008), Wu and Liu (2009) and Zou and Yuan (2008) developed robust regularized estimators based on quantile regression for the case of fixed dimensionality. Belloni and Chernozhukov (2011) studied the L_1 -penalized quantile regression in high dimensional sparse models. Fan, Fan, and Barut (2014) further considered an adaptively weighted L_1 penalty to alleviate the bias problem and showed the oracle property and asymptotic normality of the corresponding estimator. Other robust estimators were developed based on Least Absolute Deviation (LAD) regression. Wang (2013) studied the L_1 -penalized LAD regression and showed that the estimator achieves near oracle risk performance under the high dimensional setting.

The above methods essentially estimate the conditional *median (or quantile)* regression, instead of the conditional *mean* regression function. In the applications where the mean regression is of interest, these methods are not feasible unless a strong assumption is made that the distribution of errors is symmetric around zero. A simple example is the heteroscedastic linear model with asymmetric noise distribution. Another example is to estimate the conditional variance function

such as ARCH model (Engle, 1982). In these cases, the conditional mean and conditional median are very different. Another important example is to estimate large covariance matrix without assuming light-tails. We will explain this more in details in Section 5. In addition, LAD-based methods tend to penalize strongly on small errors. If only a small proportion of samples are outliers, they are expected to be less efficient than the least squares based method.

A natural question is then how to conduct ultrahigh dimensional mean regression when the tails of errors are not light? How to estimate the sample mean with very fast concentration when the distribution has only bounded second moment? These simple questions have not been carefully studied. LAD-based methods do not intend to answer these questions as they alter the problems of the study. This leads us to consider Huber loss as another way of robustification. The Huber loss (Huber, 1964) is a hybrid of squared loss for relatively small errors and absolute loss for relatively large errors, where the degree of hybridization is controlled by one tuning parameter. Unlike the traditional Huber loss, we allow the regularization parameter to diverge (or converge if its reciprocal is used) in order to reduce the bias induced by the Huber loss for estimating conditional mean regression function. In this paper, we consider the regularized approximate quadratic (RA-Lasso) estimator with an L_1 penalty and show that it admits the same L_2 error rate as the optimal error rate in the light-tail situation. In particular, if the distribution of errors is indeed symmetric around 0 (where the median and mean agree), this rate is the same as the regularized LAD estimator obtained in Wang (2013). Therefore, the RA-Lasso estimator does not lose efficiency in this special case. In practice, since the distribution of errors is unknown, RA-Lasso is more flexible than the existing methods in terms of estimating the conditional mean regression function.

A by-product of our method is that the RA-Lasso estimator of the population mean has the exponential type of concentration even in presence of the finite second moment. Catoni (2012) studied this type of problem and proposed a class of losses to result in a robust M -estimator of mean with exponential type of concentration. We further extend his idea to the sparse linear regression setting.

As done in many other papers, estimators with nice sampling properties are typically defined

through the optimization of a target function such as the penalized least-squares. The properties that are established are not the same as the ones that are computed. Following the framework of Agarwal, Negahban, and Wainwright (2012), we propose the composite gradient descent algorithm for solving the RA-Lasso estimator and develop the sampling properties by taking computational error into consideration. We show that the algorithm indeed produces a solution that admits the same optimal L_2 error rate as the theoretical estimator after sufficient number of iterations.

This paper is organized as follows. First, in Section 2, we introduce the RA-Lasso estimator and show that it has the same L_2 error rate as the optimal rate under light-tails. In Section 3, we study the property of the composite gradient descent algorithm for solving our problem and show that the algorithm produces a solution that performs as well as the solution in theory. In Section 4, we show the connection between Huber loss and Catoni loss and establish an concentration inequality for robust estimation of mean. The estimation of the error's variance is investigated in Section 5. Numerical studies are given in Section 6 and 7 to compare our method with two competitors. All technical proofs are presented in Section 8.

2 RA-Lasso estimator

We consider the linear regression model

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}^* + \epsilon_i, \quad (2.1)$$

where $\{\mathbf{x}_i\}_{i=1}^n$ are independent and identically distributed (i.i.d) p -dimensional covariate vectors, $\{\epsilon_i\}_{i=1}^n$ are i.i.d errors, and $\boldsymbol{\beta}^*$ is a p -dimensional regression coefficient vector. We consider the high-dimensional setting, where $\log(p) = o(n^b)$ for some constant $0 < b < 1$. We assume the distributions of \mathbf{x} and ϵ are independent and both have mean 0. Under this assumption, $\boldsymbol{\beta}^*$ is the mean effect of y conditioning on \mathbf{x} , which is assumed to be of interest.

To adapt for different magnitude of errors and robustify the estimation, we propose to use the

Huber loss (Huber, 1964):

$$\ell_\alpha(x) = \begin{cases} 2\alpha^{-1}|x| - \alpha^{-2} & \text{if } |x| > \alpha^{-1}; \\ x^2 & \text{if } |x| \leq \alpha^{-1}. \end{cases} \quad (2.2)$$

The Huber loss is quadratic for small values of x and linear for large values of x . The parameter α controls the blending of quadratic and linear penalization. The least squares and the LAD can be regarded as two extremes of the Huber loss for $\alpha = 0$ and $\alpha = \infty$, respectively. Deviated from the traditional Huber's estimator, the parameter α converges to zero in order to reduce the biases of estimating the mean regression function when the conditional distribution of ε_i is not symmetric. On the other hand, α can not shrink too fast in order to maintain the robustness. In this paper, we regard α as a tuning parameter, whose optimal value will be discussed later in this section. In practice, α needs to be tuned by some data-driven method. By letting α vary, we call $\ell_\alpha(x)$ the robust approximate quadratic (RA-quadratic) loss.

To estimate β^* , we propose to solve the following convex optimization problem:

$$\hat{\beta} = \operatorname{argmin}_{\beta} \frac{1}{n} \sum_{i=1}^n \ell_\alpha(y_i - \mathbf{x}_i^T \beta) + \lambda_n \sum_{j=1}^p |\beta_j|. \quad (2.3)$$

To assess the performance of $\hat{\beta}$, we study the property of $\|\hat{\beta} - \beta^*\|_2$, where $\|\cdot\|_2$ is the Euclidean norm of a vector. When λ_n converges to zero sufficiently fast, $\hat{\beta}$ is a natural M -estimator of $\beta_\alpha^* = \operatorname{argmin}_{\beta} \mathbb{E} \ell_\alpha(y - \mathbf{x}^T \beta)$, which is the population minimizer under the RA-quadratic loss and varies by α . In general, β_α^* differs from β^* . But, since the RA-quadratic loss approximates the quadratic loss as α tends to 0, β_α^* is expected to converge to β^* . This property will be established in Theorem 1. Therefore, we decompose the statistical error $\hat{\beta} - \beta^*$ into the approximation error $\beta_\alpha^* - \beta^*$ and the estimation error $\hat{\beta} - \beta_\alpha^*$. The statistical error $\|\hat{\beta} - \beta^*\|_2$ is then bounded by

$$\|\hat{\beta} - \beta^*\|_2 \leq \underbrace{\|\beta_\alpha^* - \beta^*\|_2}_{\text{approximation error}} + \underbrace{\|\hat{\beta} - \beta_\alpha^*\|_2}_{\text{estimation error}}.$$

In the following, we give the rate of the approximation and estimation error, respectively. We show that $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2$ admits the same rate as the optimal rate under light tails, as long as the two tuning parameters α and λ_n are properly chosen. We first give the rate of the approximation error under some moment conditions on \mathbf{x} and ϵ . We assume both $\boldsymbol{\beta}^*$ and $\boldsymbol{\beta}_\alpha^*$ are interior points of an L_2 ball with sufficiently large radius.

Theorem 1. *It holds that $\|\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*\|_2 = O(\alpha^{k-1})$, under the following conditions:*

(C1) $E|\epsilon|^k \leq M_k < \infty$, for some $k \geq 2$.

(C2) $0 < \kappa_l \leq \lambda_{\min}(E[\mathbf{xx}^T]) \leq \lambda_{\max}(E[\mathbf{xx}^T]) \leq \kappa_u < \infty$,

(C3) For any $\boldsymbol{\nu} \in \mathbb{R}^p$, $\mathbf{x}^T \boldsymbol{\nu}$ is sub-Gaussian with parameter at most $\kappa_0^2 \|\boldsymbol{\nu}\|_2^2$, i.e. $E \exp(t\mathbf{x}^T \boldsymbol{\nu}) \leq \exp(t^2 \kappa_0^2 \|\boldsymbol{\nu}\|_2^2 / 2)$, for any $t \in \mathbb{R}$.

Theorem 1 reveals that the approximation error vanishes faster if higher moments of error distribution exists. We next give the rate of the estimation error $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2$. This part differs from the existing work regarding the estimation error of high dimensional regularized M -estimator (Negahban, *et al.*, 2012; Agarwal, Negahban, and Wainwright, 2012) as the population minimizer $\boldsymbol{\beta}_\alpha^*$ now varies with α . However, we will show that the estimation error rate does not depend on α , given a uniform sparsity condition.

In order to be solvable in the high-dimensional setting, $\boldsymbol{\beta}^*$ is usually assumed to be sparse or weakly sparse, i.e. many elements of $\boldsymbol{\beta}^*$ are zero or small. By Theorem 1, $\boldsymbol{\beta}_\alpha^*$ converges to $\boldsymbol{\beta}^*$ as α goes to 0. In view of this fact, we assume that $\boldsymbol{\beta}_\alpha^*$ is uniformly weakly sparse when α is sufficiently small. In particular, we assume that there exists a small constant $r > 0$, such that $\boldsymbol{\beta}_\alpha^*$ belongs to an L_q -ball with a uniform radius R_q that

$$\sum_{j=1}^p |\beta_{\alpha,j}^*|^q \leq R_q, \quad (2.4)$$

for all $\alpha \in (0, r]$, and some $q \in [0, 1]$. When the conditional distribution of ϵ_i is symmetric, $\beta_{\alpha,j}^* = \beta_j^*$ for all α and j . Therefore the condition reduces to that $\boldsymbol{\beta}^*$ is in the L_q ball. In a special case where $q = 1$, it follows from Theorem 1 that if $\boldsymbol{\beta}^*$ belongs to the L_1 -ball with radius $R_1/2$ and

$r \leq [R_1/(2c_0\sqrt{p})]^{1/k-1}$, where c_0 is a generic constant, then β_α^* belongs to the L_1 -ball with radius R_1 for all $\alpha \in (0, r)$. For a general $q \in [0, 1)$, we assume a uniform upper bound R_q as in (2.4), which is allowed to diverge to infinity.

Since the RA-quadratic loss is convex, we show that with high probability the estimation error $\widehat{\Delta} = \widehat{\beta} - \beta_\alpha^*$ belongs to a star-shaped set, which depends on α and the threshold level η of signals.

Lemma 1. *Under Conditions (C1) and (C3), by choosing $\lambda_n = \kappa_\lambda \sqrt{\frac{\log p}{n}}$ and $\alpha \geq \frac{L\lambda_n}{4v}$, where κ_λ , v and L are some constants, with probability greater than $1 - 2p^{-c_0}$,*

$$\widehat{\Delta} = \widehat{\beta} - \beta_\alpha^* \in \mathbb{C}_{\alpha\eta} = \{\Delta \in \mathbb{R}^p : \|\Delta_{S_{\alpha\eta}^c}\|_1 \leq 3\|\Delta_{S_{\alpha\eta}}\|_1 + 4\|\beta_{S_{\alpha\eta}^c}^*\|_1\},$$

where $c_0 = \kappa_\lambda^2/(32v) - 1$, η is a positive constant, $S_{\alpha\eta} = \{j : |\beta_{\alpha,j}^*| > \eta\}$ and $\Delta_{S_{\alpha\eta}}$ denotes the subvector of Δ with indices in set $S_{\alpha\eta}$.

We further verify a restricted strong convexity (RSC) condition, which has been shown to be critical in the study of high dimensional regularized M -estimator (Negahban, *et al.*, 2012; Agarwal, Negahban, and Wainwright, 2012). Let

$$\delta\mathcal{L}_n(\Delta, \beta) = \mathcal{L}_n(\beta + \Delta) - \mathcal{L}_n(\beta) - [\nabla\mathcal{L}_n(\beta)]^T \Delta, \quad (2.5)$$

where $\mathcal{L}_n(\beta) = \frac{1}{n} \sum_{i=1}^n \ell_\alpha(y_i - \mathbf{x}_i^T \beta)$, Δ is a p -dimensional vector and $\nabla\mathcal{L}_n(\beta)$ is the gradient of \mathcal{L}_n at the point of β .

Definition 1. *The loss function \mathcal{L}_n satisfies RSC condition on a set S with curvature $\kappa_{\mathcal{L}} > 0$ and tolerance $\tau_{\mathcal{L}}$ if*

$$\delta\mathcal{L}_n(\Delta, \beta) \geq \kappa_{\mathcal{L}} \|\Delta\|_2^2 - \tau_{\mathcal{L}}^2, \text{ for all } \Delta \in S.$$

Next, we show that with high probability the RA-quadratic loss (2.2) satisfies RSC for all $\Delta \in \mathbb{C}_{\alpha\eta} \cap \{\Delta : \|\Delta\|_2 \leq 1\}$ with uniform constants $\kappa_{\mathcal{L}}$ and $\tau_{\mathcal{L}}$ that do not depend on α . Lemma 2 is a preliminary result, based on which RSC is checked in Lemma 3.

Lemma 2. Under conditions (C1)-(C3), for all $\|\Delta\|_2 \leq 1$, there exist uniform positive constants κ_1 and κ_2 , such that

$$\delta\mathcal{L}_n(\Delta, \beta_\alpha^*) \geq \kappa_1 \|\Delta\|_2 (\|\Delta\|_2 - \kappa_2 \sqrt{(\log p)/n}) \|\Delta\|_1, \quad (2.6)$$

with probability at least $1 - c_1 \exp(-c_2 n)$ for some positive constants c_1 and c_2 .

Lemma 3. Suppose conditions (C1)-(C3) hold and assume that

$$8\kappa_2 \kappa_\lambda^{-q/2} \sqrt{R_q} \left(\frac{\log p}{n} \right)^{(1-q)/2} \leq 1, \quad (2.7)$$

with the choice $\eta = \lambda_n$, with probability at least $1 - c_1 \exp(-c_2 n)$, the RSC condition holds for $\delta\mathcal{L}_n(\Delta, \beta_\alpha^*)$ and $\Delta \in \mathbb{C}_{\alpha\eta} \cap \{\Delta : \|\Delta\|_2 \leq 1\}$ with $\kappa_{\mathcal{L}} = \kappa_1/2$ and $\tau_{\mathcal{L}}^2 = 4R_q \kappa_2 \kappa_\lambda^{1-q} \left(\frac{\log p}{n} \right)^{1-(q/2)}$.

Lemma 3 shows that, even though β_α^* is unknown and the set $\mathbb{C}_{\alpha\eta}$ depends on α , RSC holds with uniform constants that do not depend on α . This further gives the following upper bound of the estimation error $\|\widehat{\beta} - \beta_\alpha^*\|_2$, which also does not depend on α .

Theorem 2. Under conditions of Lemma 1 and (2.7), with probability at least $1 - 2p^{-c_0} - c_1 \exp(-c_2 n)$,

$$\|\widehat{\beta} - \beta_\alpha^*\|_2 = O(\sqrt{R_q}[(\log p)/n]^{1/2-q/4}).$$

Finally, Theorem 1 and 2 together lead to the following main result, which gives the rate of the statistical error $\|\widehat{\beta} - \beta^*\|_2$.

Theorem 3. Under conditions of Lemma 1 and (2.7), with probability at least $1 - 2p^{-c_0} - c_1 \exp(-c_2 n)$,

$$\|\widehat{\beta} - \beta^*\|_2 = O(\alpha^{k-1}) + O(\sqrt{R_q}[(\log p)/n]^{1/2-q/4}).$$

Next, we compare our result with the existing results regarding the robust estimation of high dimensional linear regression model.

1. When the distribution of ϵ is symmetric around 0, then $\beta_\alpha^* = \beta^*$ for any α , which has no approximation error. If ϵ has heavy tails in addition to being symmetric, we would like to choose α sufficiently large to robustify the estimation. It then follows from Theorem 2 that $\|\widehat{\beta} - \beta^*\|_2 = O_P(\sqrt{R_q}[(\log p)/n]^{1/2-q/4})$, where $R_q = \sum_{j=1}^p |\beta_j^*|^q$. The rate is the same as the minimax rate (Raskutti, Wainwright, and Yu, 2011) for weakly sparse model under the light tails. In a special case that $q = 0$, it gives $\|\widehat{\beta} - \beta^*\|_2 = O_P(\sqrt{s(\log p)/n})$, where s is the number of nonzero elements in β^* . This is the same rate as the regularized LAD estimator in Wang (2013) and the regularized quantile estimator in Belloni and Chernozhukov (2011). It suggests that our method does not lose efficiency for symmetric heavy-tailed errors.
2. If the distribution of ϵ is asymmetric around 0, the quantile and LAD based methods are inconsistent, since they estimate the median instead of the mean. Theorem 3 shows that our estimator still achieves the optimal rate given that $\alpha = o(\{R_q[(\log p)/n]^{1-\frac{q}{2}}\}^{\frac{1}{2(k-1)}})$ even though the k -th moment of ϵ is assumed. Recall from conditions in Lemma 1 that we also need to choose $\alpha > c\sqrt{(\log p)/n}$ for some constant c . Given the sparsity condition (2.7), α can be chosen to meet the above two requirements. In terms of estimating the conditional mean effect, errors with heavy but asymmetric tails give the case where the RA-Lasso has the biggest advantage over the other estimators.

In practice, the distribution of errors is unknown. However, we proved that our method is no worse than the existing methods for any type of errors, as long as the tuning parameters are chosen properly. Hence, our method is more flexible.

3 Geometric convergence of computational error

The gradient descent algorithm (Nesterov, 2007; Agarwal, Negahban, and Wainwright, 2012) is usually applied to solve the convex problem (2.3). For example, we can replace the RA-quadratic loss with its local quadratic approximation (LQA) and iteratively solve the following optimization

problem:

$$\widehat{\boldsymbol{\beta}}^{t+1} = \operatorname{argmin}_{\|\boldsymbol{\beta}\|_1 \leq \rho} \left\{ \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^t) + [\nabla \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^t)]^T (\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^t) + \frac{\gamma_u}{2} \|\boldsymbol{\beta} - \widehat{\boldsymbol{\beta}}^t\|_2^2 + \lambda_n \|\boldsymbol{\beta}\|_1 \right\}, \quad (3.1)$$

where γ_u is a fixed constant at each iteration, and the side constraint “ $\|\boldsymbol{\beta}\|_1 \leq \rho$ ” is introduced to guarantee good performance in the first few iterations and ρ is allowed to be sufficiently large. To solve (3.1), the update can be computed by a two-step procedure. We first solve (3.1) without the norm constraint by soft-thresholding the vector $\widehat{\boldsymbol{\beta}}^t - \frac{1}{\gamma_u} \nabla \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^t)$ at level λ_n and call the solution $\check{\boldsymbol{\beta}}$. If $\|\check{\boldsymbol{\beta}}\|_1 \leq \rho$, set $\widehat{\boldsymbol{\beta}}^{t+1} = \check{\boldsymbol{\beta}}$. Otherwise, $\widehat{\boldsymbol{\beta}}^{t+1}$ is obtained by further project $\check{\boldsymbol{\beta}}$ onto the L_1 -ball $\{\boldsymbol{\beta} : \|\boldsymbol{\beta}\|_1 \leq \rho\}$. The projection can be done (Duchi, *et al.*, 2008) by soft-thresholding $\check{\boldsymbol{\beta}}$ at level π_n , where π_n is given by the following procedure: (1) sort $\{|\check{\beta}_j|\}_{j=1}^p$ into $b_1 \geq b_2 \geq \dots \geq b_p$; (2) find $J = \max\{1 \leq j \leq p : b_j - (\sum_{r=1}^j b_r - \rho)/j > 0\}$ and let $\pi_n = (\sum_{r=1}^J b_r - \rho)/J$.

Agarwal, Negahban, and Wainwright (2012) considered the computational error of such first-order gradient descent method. They showed that, for a convex and differentiable loss functions $\ell(x)$ and decomposable penalty function $p(\boldsymbol{\beta})$, the error $\|\widehat{\boldsymbol{\beta}}^t - \boldsymbol{\beta}^*\|_2$ has the same rate as $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2$ for all sufficiently large t , where $\boldsymbol{\beta}^* = \operatorname{argmin}_{\boldsymbol{\beta}} \mathbb{E} \mathcal{L}(\mathbf{x}, y; \boldsymbol{\beta})$, and $\widehat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}_i, y_i; \boldsymbol{\beta}) + p(\boldsymbol{\beta})$. Different from their setup, our population minimizer $\boldsymbol{\beta}_\alpha^*$ varies by α . Nevertheless, as $\boldsymbol{\beta}_\alpha^*$ converges to the true effect $\boldsymbol{\beta}^*$, by a careful control of α , we can still show that $\|\widehat{\boldsymbol{\beta}}^t - \boldsymbol{\beta}^*\|_2$ has the same rate as $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2$, where $\widehat{\boldsymbol{\beta}}$ is the theoretical solution of (2.3) and $\widehat{\boldsymbol{\beta}}^t$ is as defined in (3.1).

The key is that the RA-quadratic loss function \mathcal{L}_n satisfies the restricted strong convexity (RSC) condition and the restricted smoothness condition (RSM) with some uniform constants, namely $\delta \mathcal{L}_n(\boldsymbol{\Delta}, \boldsymbol{\beta})$ as defined in (2.5) satisfies the following conditions:

$$\text{RSC} : \delta \mathcal{L}_n(\boldsymbol{\Delta}, \boldsymbol{\beta}) \geq \frac{\gamma_l}{2} \|\boldsymbol{\Delta}\|_2^2 - \tau_l \|\boldsymbol{\Delta}\|_1^2, \quad (3.2)$$

$$\text{RSM} : \delta \mathcal{L}_n(\boldsymbol{\Delta}, \boldsymbol{\beta}) \leq \frac{\gamma_u}{2} \|\boldsymbol{\Delta}\|_2^2 + \tau_u \|\boldsymbol{\Delta}\|_1^2, \quad (3.3)$$

for all $\boldsymbol{\beta}$ and $\boldsymbol{\Delta}$ in some set of interest, with parameters γ_l , τ_l , γ_u and τ_u that do not depend on α .

We show that such conditions hold with high probability.

Lemma 4. Under conditions (C1)-(C3), for all $\beta \in \mathbb{R}^p$ and $\Delta \in \{\Delta : \|\Delta\|_2 \leq 1\}$, with probability greater than $1 - c_1 \exp(-c_2 n)$, (3.2) and (3.3) hold with $\gamma_l = \kappa_1$, $\tau_l = \kappa_1 \kappa_2^2 (\log p) / (2n)$, $\gamma_u = 3\kappa_u$, $\tau_u = \kappa_u (\log p) / n$.

We further show in Theorem 4 that, whenever $R_q \left(\frac{\log p}{n}\right)^{1-(q/2)} = o(1)$, which is required for consistency of *any method* over the weak sparse L_q ball by the known minimax results (Raskutti, Wainwright, and Yu, 2011), it holds that $\|\widehat{\beta}^t - \widehat{\beta}\|_2 = o(\|\widehat{\beta} - \beta_\alpha^*\|_2)$ for sufficiently many iterations with lower bound specified in Theorem 4. Hence,

$$\begin{aligned} \|\widehat{\beta}^t - \beta^*\|_2 &\leq \|\widehat{\beta}^t - \widehat{\beta}\|_2 + \|\widehat{\beta} - \beta_\alpha^*\|_2 + \|\beta_\alpha^* - \beta^*\|_2 \\ &= o(\|\widehat{\beta} - \beta_\alpha^*\|_2) + \|\widehat{\beta} - \beta_\alpha^*\|_2 + \|\beta_\alpha^* - \beta^*\|_2 \\ &= O(\alpha^{k-1}) + O(\sqrt{R_q}[(\log p)/n]^{1/2-q/4}), \end{aligned}$$

which has the same rate as $\|\widehat{\beta} - \beta^*\|_2$. Hence, from a statistical point of view, there is no need to iterate beyond t steps.

Theorem 4. Under conditions of Theorem 3, suppose we choose λ_n as in Lemma 1 and also satisfying

$$\lambda_n \geq \frac{32\rho}{1-\kappa} \left(1 - \frac{64\kappa_u |S_{\alpha\eta}| \frac{\log p}{n}}{\bar{\gamma}_l}\right)^{-1} \left[1 + \kappa_1 \kappa_2^2 \left(\frac{\bar{\gamma}_l}{12\kappa_u} + \frac{128\kappa_u |S_{\alpha\eta}| \frac{\log p}{n}}{\bar{\gamma}_l}\right) + 8\kappa_u\right] \frac{\log p}{n},$$

where $|S_{\alpha\eta}|$ denotes the cardinality of set $S_{\alpha\eta}$ and $\bar{\gamma}_l = \gamma_l - 64\tau_l |S_{\alpha\eta}|$, then with probability at least $1 - p^{-c_0} - c_1 \exp(-c_2 n)$, we have

$$\|\widehat{\beta}^t - \widehat{\beta}\|_2^2 = O\left(R_q \left(\frac{\log p}{n}\right)^{1-(q/2)} \left[\|\widehat{\beta} - \beta_\alpha^*\|_2^2 + R_q \left(\frac{\log p}{n}\right)^{1-(q/2)}\right]\right),$$

for all iterations

$$t \geq \frac{2 \log((\phi_n(\widehat{\beta}^0) - \phi_n(\widehat{\beta})) / \delta^2)}{\log(1/\kappa)} + \log_2 \log_2 \left(\frac{\rho \lambda_n}{\delta^2}\right) \left(1 + \frac{\log 2}{\log(1/\kappa)}\right),$$

where $\phi_n(\boldsymbol{\beta}) = \mathcal{L}_n(\boldsymbol{\beta}) + \lambda_n \|\boldsymbol{\beta}\|_1$ and $\widehat{\boldsymbol{\beta}}^0$ is the initial value, $\delta = \varepsilon^2/(1 - \kappa)$ is the tolerance level, κ and ε are some constants as will be defined in (8.21) and (8.22), respectively.

4 Connection with Catoni loss

Catoni (2012) considered the estimation of the mean of heavy-tailed distribution with fast concentration. He proposed an M -estimator by solving

$$\sum_{i=1}^n \psi_c[\alpha(y_i - \theta)] = 0,$$

where the influence function $\psi_c(x)$ is chosen such that $-\log(1-x+x^2/2) \leq \psi_c(x) \leq \log(1+x+x^2/2)$. He showed that this M -estimator has the exponential type of concentration by only requiring the existence of the variance. It performed as well as the sample mean under the light-tail case. In Section 2, we essentially showed the same type of concentration for the RA-quadratic loss under the linear regression setting.

The estimation of mean can be regarded as a univariate linear regression where the covariate equals to 1. In that special case, we have a more explicit concentration result for the RA-mean estimator, which is the estimator that minimizes the RA-quadratic loss. Let $\{y_i\}_{i=1}^n$ be an i.i.d sample from some unknown distribution with $E(y_i) = \mu$ and $\text{var}(y_i) = \sigma^2$. The RA-mean estimator $\widehat{\mu}_\alpha$ of μ is the solution of

$$\sum_{i=1}^n \psi[\alpha(y_i - \mu)] = 0,$$

for parameter $\alpha \rightarrow 0$, where the influence function $\psi(x) = x$ if $|x| \leq 1$, $\psi(x) = 1$, if $x > 1$ and $\psi(x) = -1$ if $x < -1$. The following theorem gives the exponential type of concentration of $\widehat{\mu}_\alpha$ around μ .

Theorem 5. Assume $\frac{\log(1/\delta)}{n} \leq 1/8$ and let $\alpha = \sqrt{\frac{\log(1/\delta)}{nv^2}}$ where $v \geq \sigma$. Then,

$$P\left(|\widehat{\mu}_\alpha - \mu| \geq 4v\sqrt{\frac{\log(1/\delta)}{n}}\right) \leq 2\delta.$$

The above result provides fast concentration of the mean estimation with only two moments assumption. This is very useful for last scale hypothesis testing (Efron, 2010; Fan, Han, and Gu, 2012) and covariance matrix estimation (Bickel and Levina, 2008; Fan, Liao and Mincheva, 2013), where uniform convergence is required. Taking the estimation of large covariance matrix as an example, in order for the elements of the sample covariance matrix to converge uniformly, the aforementioned authors require the underlying multivariate distribution be sub-Gaussian. This restrictive assumptions can be removed if we apply the robust estimation with concentration bound. Regarding $\sigma_{ij} = E X_i X_j$ as the expected value of the random variable $X_i X_j$ (it is typically not the same as the median of $X_i X_j$), it can be estimated with accuracy

$$P \left(|\hat{\sigma}_{ij} - \sigma_{ij}| \geq 4v \sqrt{\frac{\log(1/\delta)}{n}} \right) \leq 2\delta,$$

where $v \geq \max_{i,j \leq p} \sqrt{\text{var}(X_i X_j)}$ and $\hat{\sigma}_{ij}$ is RA-mean estimator using data $\{X_{ik} X_{jk}\}_{k=1}^n$. Since there are only $O(p^2)$ elements, by taking $\delta = p^{-3}$ and the union bound, we have

$$\max_{i,j \leq p} \sqrt{n/\log p} |\hat{\sigma}_{ij} - \sigma_{ij}| \rightarrow 0,$$

when $E X_i^4 < \infty$. This robustified covariance estimator requires much weaker condition than the sample covariance and has far wide applicability than the sample covariance. It can be regularized further in the same way as the sample covariance matrix.

On the other hand, Catoni's idea could also be extended to the linear regression setting. Suppose we replace the RA-quadratic loss $\ell_\alpha(x)$ in (2.3) with Catoni loss

$$\ell_\alpha^c(x) = \frac{2}{\alpha} \int_0^x \psi_c(\alpha t) dt,$$

where the influence function $\psi_c(t)$ is given by

$$\psi_c(t) = \text{sgn}(t) \{ -\log(1 - |t| + t^2/2) I(|t| < 1) + \log(2) I(|t| \geq 1) \}.$$

Let $\widehat{\boldsymbol{\beta}}^c$ be the corresponding solution. Then, $\widehat{\boldsymbol{\beta}}^c$ has the same convergence rate as the RA-Lasso, when the second or the third moment of errors exists.

Theorem 6. *Suppose condition (C1) holds for $k = 2$ or 3 , (C2), (C3) and (2.7) hold, then with probability at least $1 - 2p^{-c_0} - c_1 \exp(-c_2 n)$,*

$$\|\widehat{\boldsymbol{\beta}}^c - \boldsymbol{\beta}^*\|_2 = O(\alpha^{k-1}) + O(\sqrt{R_q}[(\log p)/n]^{1/2-q/4}).$$

Unlike the RA-lasso, the order of bias of $\widehat{\boldsymbol{\beta}}^c$ cannot be further improved, even when higher moments of errors exist beyond the third order. The reason is that the Catoni loss is not exactly the quadratic loss over any finite intervals. Similar results regarding the computational error of $\widehat{\boldsymbol{\beta}}^c$ could also be established as in Theorem 4, since the RSC/RSM conditions also hold for Catoni loss with uniform constants.

5 Variance Estimation

We estimate $\sigma^2 = E\epsilon^2$ based on the RA-Lasso estimator and a cross-validation scheme. To ease the presentation, we assume the data set can be evenly divided into K folds with m observations in each fold. Then, we estimate σ^2 by

$$\widehat{\sigma}^2 = \frac{1}{K} \sum_{k=1}^K \frac{1}{m} \sum_{i \in \text{fold } k} (y_i - \mathbf{x}_i^T \widehat{\boldsymbol{\beta}}^{(-k)})^2,$$

where $\widehat{\boldsymbol{\beta}}^{(-k)}$ is the RA-Lasso estimator obtained by using data points outside the k -th fold. We show that $\widehat{\sigma}^2$ is asymptotically efficient.

Theorem 7. *Under conditions of Theorem 3, if $R_q(\log p)^{1-q/2}/n^{(1-q)/2} \rightarrow 0$ for $q \in [0, 1)$, and $\alpha = o\left(\{R_q[(\log p)/n]^{1-\frac{q}{2}}\}^{\frac{1}{2(k-1)}}\right)$, then*

$$\sqrt{n}(\widehat{\sigma}^2 - \sigma^2) \xrightarrow{\mathcal{D}} N(0, E\epsilon^4 - \sigma^4).$$

6 Simulation Studies

In this section, we assess the finite sample performance of the RA-Lasso and compare it with other methods through various models. We simulated data from the following high dimensional model

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}^* + \epsilon_i, \quad \mathbf{x}_i \sim N(0, I_p), \quad (6.1)$$

where we generated $n = 100$ observations and the number of parameters was chosen to be $p = 400$.

We chose the true regression coefficient vector as

$$\boldsymbol{\beta}^* = (3, \dots, 3, 0, \dots, 0)^T,$$

where the first 20 elements are all equal to 3 and the rest are all equal to 0. To involve various shapes of error distributions, we considered the following five scenarios:

1. Normal with mean 0 and variance 4 ($N(0,4)$);
2. Two times the t-distribution with degrees of freedom 3 ($2t_3$);
3. Mixture of Normal distribution(MixN): $0.5N(-1, 4) + 0.5N(8, 1)$;
4. Log-normal distribution (LogNormal): $\epsilon = e^{1+1.2Z}$, where Z is standard normal.
5. Weibull distribution with shape parameter = 0.3 and scale parameter = 0.5 (Weibull).

In order to meet the model assumption, the errors were standardized to have mean 0. Table 1 categorizes the five scenarios according to the shapes and tails of the error distributions.

	Light Tail	Heavy Tail
Symmetric	$N(0, 4)$	$2t_3$
Asymmetric	MixN	LogNormal, Weibull

Table 1: Summary of the shapes and tails of five error distributions

To obtain our estimator, we iteratively applied the gradient descent algorithm. We compared RA-Lasso with another two methods in high-dimensional setting: (a) Lasso: the penalized least-squares estimator with L_1 -penalty as in Tibshirani (1996); and (b) R-Lasso: the R-Lasso estimator in Fan, Fan, and Barut (2014), which is the same as the regularized LAD estimator with L_1 -penalty as in Wang (2013). Their performance under the five scenarios was evaluated by the following four measurements:

- (1) L_2 error, which is defined as $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2$.
- (2) L_1 error, which is defined as $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_1$.
- (3) Number of false positives (FP), which is number of noise covariates that are selected.
- (4) Number of false negatives (FN), which is number of signal covariates that are not selected.

We also measured the relative gain of RA-Lasso with respect to R-Lasso and Lasso, in terms of the difference to the oracle estimator. The oracle estimator $\widehat{\boldsymbol{\beta}}_{\text{oracle}}$ is defined to be the least square estimator by using the first 20 covariates only. Then, the relative gain of RA-Lasso with respect to Lasso ($\text{RG}_{\text{A,L}}$) in L_2 and L_1 norm are defined as

$$\frac{\|\widehat{\boldsymbol{\beta}}_{\text{Lasso}} - \boldsymbol{\beta}^*\|_2 - \|\widehat{\boldsymbol{\beta}}_{\text{oracle}} - \boldsymbol{\beta}^*\|_2}{\|\widehat{\boldsymbol{\beta}}_{\text{RA-Lasso}} - \boldsymbol{\beta}^*\|_2 - \|\widehat{\boldsymbol{\beta}}_{\text{oracle}} - \boldsymbol{\beta}^*\|_2} \quad \text{and} \quad \frac{\|\widehat{\boldsymbol{\beta}}_{\text{Lasso}} - \boldsymbol{\beta}^*\|_1 - \|\widehat{\boldsymbol{\beta}}_{\text{oracle}} - \boldsymbol{\beta}^*\|_1}{\|\widehat{\boldsymbol{\beta}}_{\text{RA-Lasso}} - \boldsymbol{\beta}^*\|_1 - \|\widehat{\boldsymbol{\beta}}_{\text{oracle}} - \boldsymbol{\beta}^*\|_1}.$$

The relative gain of RA-Lasso with respect to R-Lasso ($\text{RG}_{\text{A,R}}$) is defined similarly.

For RA-Lasso, the tuning parameters λ_n and α were chosen optimally based on 100 independent validation datasets. We ran a 2-dimensional grid search to find the best (λ_n, α) pair that minimizes the mean L_2 -loss of the 100 validation datasets. Such an optimal pair was then used in the simulations. Similar method was applied in choosing the tuning parameters in Lasso and R-Lasso.

The above simulation model is based on the additive model (6.1), in which error distribution is independent of covariates. However, this homoscedastic model makes the conditional mean and the conditional median differ only by a constant. To further examine the deviations between the

mean regression and median regression, we also simulated the data from the heteroscedastic model

$$y_i = \mathbf{x}_i^T \boldsymbol{\beta}^* + c^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}^*)^2 \epsilon_i, \quad \mathbf{x}_i \sim N(0, I_p), \quad (6.2)$$

where the constant $c = \sqrt{3}\|\boldsymbol{\beta}^*\|^2$ makes $E[c^{-1}(\mathbf{x}_i^T \boldsymbol{\beta}^*)^2] = 1$. Note that $\mathbf{x}_i^T \boldsymbol{\beta}^* \sim N(0, \|\boldsymbol{\beta}^*\|^2)$ and therefore c is chosen so that the average noise levels is the same as that of ϵ_i . For both the homoscedastic and the heteroscedastic models, we ran 100 simulations for each scenario. The mean of each performance measurement is reported in Table 2 and 3, respectively.

		Lasso	R-Lasso	RA-Lasso	RG_{A,L}	RG_{A,R}
N(0, 4)	L_2 loss	4.54	4.40	4.53	1.00	0.96
	L_1 loss	27.21	29.11	27.21	1.00	1.08
	FP, FN	52.10, 0.09	66.36, 0.17	52.10, 0.09		
2t₃	L_2 loss	6.04	5.10	5.47	1.14	0.91
	L_1 loss	35.22	33.07	30.42	1.19	1.10
	FP, FN	47.13, 0.34	65.84, 0.22	41.34, 0.28		
MixN	L_2 loss	6.14	6.44	6.13	1.00	1.06
	L_1 loss	40.46	46.18	38.48	1.06	1.23
	FP, FN	65.99, 0.34	80.31, 0.33	58.05, 0.39		
LogNormal	L_2 loss	11.08	12.16	10.10	1.14	1.30
	L_1 loss	53.17	57.18	51.58	1.04	1.14
	FP, FN	26.5, 15.00	27.20, 6.90	37.20, 3.90		
Weibull	L_2 loss	7.77	7.11	6.62	1.23	1.10
	L_1 loss	55.65	50.49	42.93	1.34	1.20
	FP, FN	78.70, 0.71	77.13, 0.56	62.27, 0.52		

Table 2: Simulation results of Lasso, R-Lasso and RA-Lasso under homoscedastic model (6.1).

Tables 2 and 3 indicate that our method had the biggest advantage when the errors were asymmetric and heavy-tailed (LogNormal and Weibull). In this case, R-Lasso had larger L_1 and L_2 errors due to the bias for estimating the conditional median instead of the mean. Even though Lasso did not have bias in the loss component, it did not perform well due to its sensitivity to outliers. The advantage of our method is more pronounced in the heteroscedastic model than in the homoscedastic model. Both of them clearly indicate that if the errors come from asymmetric and heavy-tailed distributions, our method is better than both Lasso and R-Lasso. When the

		Lasso	R-Lasso	RA-Lasso	RG_{A,L}	RG_{A,R}
N(0, 4)	L_2 loss	4.60	4.34	4.60	1.00	0.93
	L_1 loss	27.16	27.14	27.15	1.00	1.00
	FP, FN	48.78, 0.10	58.25, 0.27	48.78, 0.10		
$2t_3$	L_2 loss	8.08	6.71	6.70	1.26	1.01
	L_1 loss	41.16	42.76	38.52	1.08	1.12
	FP, FN	55.33, 0.67	71.67, 0.33	45.33, 0.33		
MixN	L_2 loss	6.26	6.54	6.25	1.00	1.06
	L_1 loss	41.26	46.95	39.25	1.06	1.23
	FP, FN	65.98, 0.34	80.30, 0.32	58.80, 0.34		
LogNormal	L_2 loss	10.86	9.19	8.48	1.43	1.13
	L_1 loss	57.52	57.18	53.20	1.10	1.09
	FP, FN	29.70, 5.70	54.10, 2.00	54.30, 1.50		
Weibull	L_2 loss	7.40	8.81	5.53	1.53	1.92
	L_1 loss	40.95	47.82	34.65	1.23	1.48
	FP, FN	38.87, 0.96	35.31, 2.90	58.15, 0.39		

Table 3: Simulation results of Lasso, R-Lasso and RA-Lasso under heteroscedastic model (6.2).

errors were symmetric and heavy-tailed ($2t_3$), our estimator performed closely as R-Lasso, both of which outperformed Lasso. The above two cases evidently showed that RA-Lasso was robust to the outliers and did not lose efficiency when the errors were indeed symmetric. Under the light-tailed scenario, if the errors were asymmetric (MixN), our method performed similarly as Lasso. R-Lasso performed worse, since it had bias. For the regular setting ($N(0, 4)$), where the errors were light-tailed and symmetric, the three methods were comparable with each other.

In conclusion, RA-Lasso is more flexible than Lasso and R-Lasso. The tuning parameter α automatically adapts to errors with different shapes and tails. It enables RA-Lasso to render consistently satisfactory results under all scenarios.

7 Real Data Example

In this section, we use a microarray data to illustrate the performance of Lasso, R-Lasso and RA-Lasso. Huang, *et al.* (2011) studied the role of innate immune system on the development of atherosclerosis by examining gene profiles from peripheral blood of 119 patients. The data

were collected using Illumina HumanRef8 V2.0 Bead Chip and are available on Gene Expression Omnibus. The original study showed that the toll-like receptors (TLR) signaling pathway plays an important role on triggering the innate immune system in face of atherosclerosis. Under this pathway, the “TLR8” gene was found to be a key atherosclerosis-associated gene. To further study the relationship between this key gene and the other genes, we regressed it on another 464 genes from 12 different pathways (TLR, CCC, CIR, IFNG, MAPK, RAPO, EXAPO, INAPO, DRS, NOD, EPO, CTR) that are related to the TLR pathway. We applied Lasso, R-Lasso and RA-Lasso to this data. The tuning parameters for all methods were chosen by using five-fold cross validation. Figure 1 shows our choice of the penalization parameter based on the cross validation results. For RA-Lasso, the choice of α was insensitive to the results and was fixed at 5. We then applied the three methods with the above choice of tuning parameters to select significant genes. The QQ-plots of the residuals from the three methods are shown in Figure 2. The selected genes by the three methods are reported in Table 4. After the selection, we regressed the expression of TLR8 gene on the selected genes, the t -values from the refittings are also reported in Table 4.

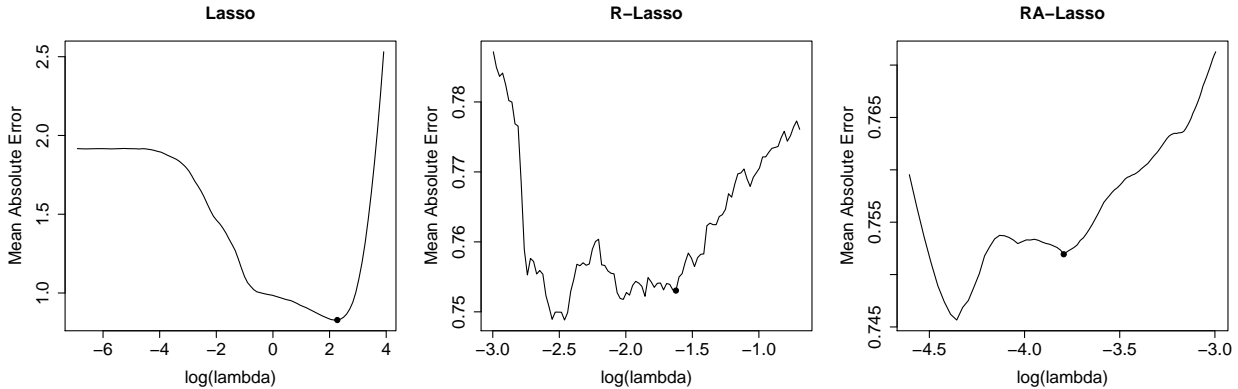


Figure 1: Five-fold cross validation results: black dot marks the choice of the penalization parameter.

Table 4 shows that Lasso only selected one gene. R-Lasso selected 17 genes. Our proposed RA-Lasso selected 34 genes. Eight genes (CSF3, IL10, AKT1, TOLLIP, TLR1, SHC1, EPOR, and TJP1) found by R-Lasso were also selected by RA-Lasso. Compared with Lasso and R-Lasso, our method selected more genes, which could be useful for a second-stage confirmatory study. It is

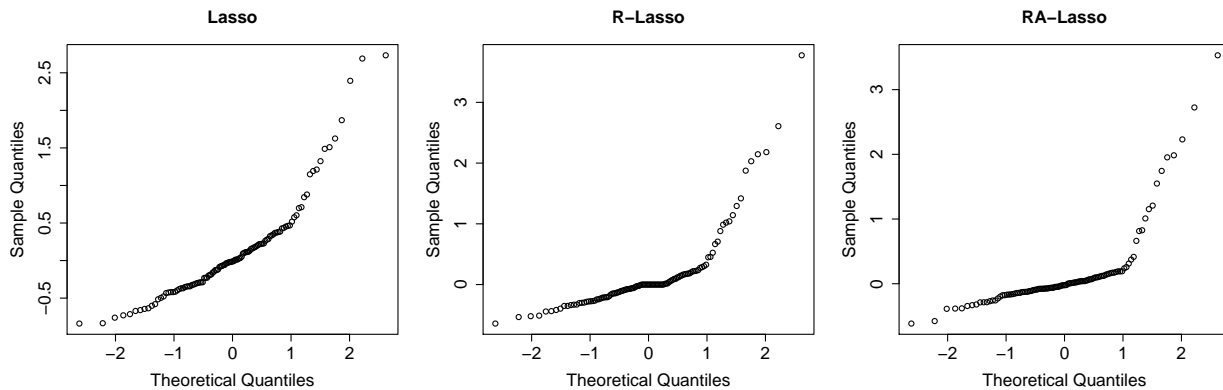


Figure 2: QQ plots of the residuals from three methods.

clearly seen from Figure 2 that the residuals from the fitted regressions had heavy right tail and skewed distribution. We know from the simulation studies in Section 6 that RA-Lasso tends to perform better than Lasso and R-Lasso in this situation. For further investigation, we randomly chose 24 patients as the test set; applied three methods to the rest patients to obtain the estimated coefficients, which in return were used to predict the responses of 24 patients. We repeated the random splitting 100 times, the boxplots of the Mean Absolute/Squared Error of predictions are shown in Figure 3. RA-Lasso has better predictions than Lasso and R-Lasso.

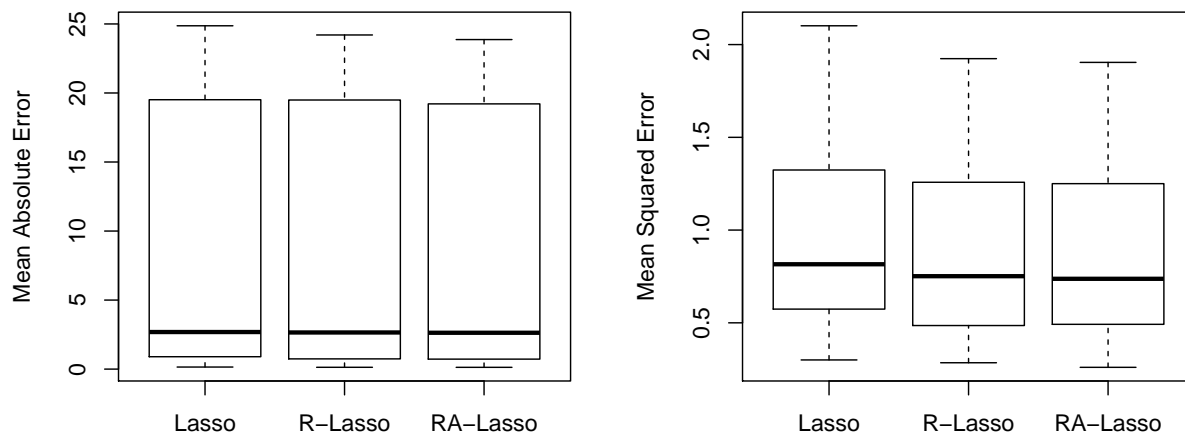


Figure 3: Boxplot of the Mean Absolute/Squared Error of predictions.

Lasso	CRK 0.23						
R-Lasso	CSF3	IL10	AKT1	KPNB1	TLR2	GRB2	MAPK1
	-2.46	2.24	1.68	1.49	1.41	-1.06	0.98
	DAPK2	TOLLIP	TLR1	TLR3	SHC1	PSMD1	F12
	0.7	-0.68	0.52	0.33	-0.28	0.27	0.24
RA-Lasso	EPOR	TJP1	GAB2				
	-0.17	-0.12	-0.01				
	CSF3	CD3E	BTK	CLSPN	RELA	AKT1	IRS2
	-2.95	2.67	2.37	1.93	1.88	1.61	1.55
	IL10	MAP2K4	PMAIP1	BCL2L11	AKT3	DUSP10	IRF4
	1.52	1.17	-1.14	-1.13	-1.01	0.97	-0.95
	IFI6	TLR1	PSMB8	KPNB1	IFNG	FADD	TJP1
	0.86	0.82	0.79	0.77	-0.74	0.65	-0.57
CR2	IL2	PSMC2	HSPA8	SHC1	SPI1	IFNA6	
0.57	-0.47	0.38	-0.35	-0.33	-0.28	0.28	
FYN	EPOR	MASP1	PRKCZ	TOLLIP	BAK1		
-0.24	0.24	-0.24	0.24	-0.19	0.14		

Table 4: Selected genes by Lasso, R-Lasso and RA-Lasso.

8 Proofs

Proof of Theorem 1. Let $\ell(x) = x^2$. Since β^* minimizes $E\ell(y - \mathbf{x}^T\beta)$, it follows from condition(C2) that

$$E[\ell(y - \mathbf{x}^T\beta_\alpha^*) - \ell(y - \mathbf{x}^T\beta^*)] = (\beta_\alpha^* - \beta^*)^T E(\mathbf{xx}^T)(\beta_\alpha^* - \beta^*) \geq \kappa_l \|\beta_\alpha^* - \beta^*\|_2^2. \quad (8.1)$$

Let $g_\alpha(x) = \ell(x) - \ell_\alpha(x) = (|x| - \alpha^{-1})^2 I(|x| > \alpha^{-1})$. Then, since β_α^* is the minimizer of $E\ell_\alpha(y - \mathbf{x}^T\beta)$, we have

$$\begin{aligned} & E[\ell(y - \mathbf{x}^T\beta_\alpha^*) - \ell(y - \mathbf{x}^T\beta^*)] \\ = & E[\ell(y - \mathbf{x}^T\beta_\alpha^*) - \ell_\alpha(y - \mathbf{x}^T\beta_\alpha^*)] + E[\ell_\alpha(y - \mathbf{x}^T\beta_\alpha^*) - \ell_\alpha(y - \mathbf{x}^T\beta^*)] \\ & + E[\ell_\alpha(y - \mathbf{x}^T\beta^*) - \ell(y - \mathbf{x}^T\beta^*)] \\ \leq & E[g_\alpha(y - \mathbf{x}^T\beta_\alpha^*)] - E[g_\alpha(y - \mathbf{x}^T\beta^*)]. \end{aligned}$$

By Taylor's expansion, we have

$$\mathbb{E}[\ell(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*) - \ell_\alpha(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*)] \leq 2 \mathbb{E}[(z - \alpha^{-1})I(z > \alpha^{-1})|\mathbf{x}^T(\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*)|], \quad (8.2)$$

where $\tilde{\boldsymbol{\beta}}$ is a vector lying between $\boldsymbol{\beta}^*$ and $\boldsymbol{\beta}_\alpha^*$ and $z = |y - \mathbf{x}^T \tilde{\boldsymbol{\beta}}|$. With \mathbb{E}_ϵ denoting the conditional expectation with respect to ϵ given \mathbf{x} , we have

$$\mathbb{E}_\epsilon[(z - \alpha^{-1})I(z > \alpha^{-1})] \leq \mathbb{E}_\epsilon[zI(z > \alpha^{-1})] \leq \alpha^{k-1} \mathbb{E}_\epsilon z^k.$$

Therefore, $\mathbb{E}[\ell(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*) - \ell(y - \mathbf{x}^T \boldsymbol{\beta}^*)]$ is further bounded by

$$2\alpha^{k-1} \mathbb{E}[|y - \mathbf{x}^T \tilde{\boldsymbol{\beta}}|^k |\mathbf{x}^T(\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*)|] \leq 2(2\alpha)^{k-1} \mathbb{E}[(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k) |\mathbf{x}^T(\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*)|], \quad (8.3)$$

where the constant M_k is defined in Condition (C1). Next, we show that $\lambda_{\max}(\mathbb{E}[(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k)^2 \mathbf{x}\mathbf{x}^T]) = O(1)$. Let $\boldsymbol{\nu}$ be a p -dimensional vector with $\|\boldsymbol{\nu}\|_2 = 1$. By the Cauchy-Schwartz inequality,

$$\mathbb{E}[(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k)^2 (\mathbf{x}^T \boldsymbol{\nu})^2] \leq [\mathbb{E}(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k)^4]^{1/2} [\mathbb{E}(\mathbf{x}^T \boldsymbol{\nu})^4]^{1/2}.$$

By (C3), $\mathbf{x}^T \boldsymbol{\nu}$ is sub-Gaussian with parameter κ_0^2 . Under the assumption that $\boldsymbol{\beta}^*$ and $\boldsymbol{\beta}_\alpha^*$ are interior points of an L_2 -ball with sufficiently large radius, $\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)$ is sub-Gaussian with parameter $\kappa_0^2 \|\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2^2$, which is no larger than $\kappa_0^2 \|\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*\|_2^2 = O(1)$. Using the moment results of sub-Gaussian random variables (Rivasplata, 2012), $\mathbb{E}(\mathbf{x}^T \boldsymbol{\nu})^4 \leq 16\kappa_0^4 = O(1)$. Similarly, $\mathbb{E}|\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^{4k} \leq \mathbb{E}|\mathbf{x}^T(\boldsymbol{\beta}^* - \boldsymbol{\beta}_\alpha^*)|^{4k} = O(1)$. Therefore, $\mathbb{E}[(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k)^4] = O(1)$. Hence, by definition, $\lambda_{\max}(\mathbb{E}[(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k)^2 \mathbf{x}\mathbf{x}^T]) = O(1)$. Using this result and (8.3),

$$\begin{aligned} \mathbb{E}[\ell(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*) - \ell(y - \mathbf{x}^T \boldsymbol{\beta}^*)] &\leq 2(2\alpha)^{k-1} [\lambda_{\max}(\mathbb{E}[(M_k + |\mathbf{x}^T(\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^k)^2 \mathbf{x}\mathbf{x}^T])]^{1/2} \|\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*\|_2 \\ &= O(\alpha^{k-1} \|\boldsymbol{\beta}_\alpha^* - \boldsymbol{\beta}^*\|_2). \end{aligned}$$

This together with (8.1) completes the proof. \square

Proof of Lemma 1. First of all, it follows from Lemma 1 of Negahban, *et al.* (2012) that $\widehat{\Delta} = \widehat{\beta} - \beta_\alpha^* \in \mathbb{C}_{\alpha\eta}$ on the event $\{\lambda_n \geq 2 \|\nabla \mathcal{L}_n(\beta_\alpha^*)\|_\infty\}$. Hence, we need to show that the event $\{\lambda_n \geq 2 \|\nabla \mathcal{L}_n(\beta_\alpha^*)\|_\infty\}$ holds with high probability. The latter will be established by using Bernstein's inequality along with the union bound.

The gradient of \mathcal{L}_n ,

$$\nabla \mathcal{L}_n(\beta_\alpha^*) = \frac{1}{n} \sum_{i=1}^n \frac{2}{\alpha} \psi[\alpha(y_i - \mathbf{x}_i^T \beta_\alpha^*)] \mathbf{x}_i, \quad (8.4)$$

where $\psi(x) = x$, for $|x| \leq 1$; $\psi(x) = 1$, for $x > 1$; and $\psi(x) = -1$, for $x < -1$. Using $\alpha^{-1} |\psi(\alpha x)| \leq |x|$ and assumption (C3), we have

$$\begin{aligned} \mathbb{E}\{2\alpha^{-1} \psi[\alpha(y_i - \mathbf{x}_i^T \beta_\alpha^*)] x_{ij}\}^2 &\leq 4 \mathbb{E}\{(y_i - \mathbf{x}_i^T \beta_\alpha^*)^2 x_{ij}^2\} \\ &\leq 8 \mathbb{E}\{(\epsilon_i^2 + |\mathbf{x}_i^T (\beta_\alpha^* - \beta^*)|^2) x_{ij}^2\} \\ &\leq v, \end{aligned}$$

where v is a constant depending on κ_0 and M_2 and the last inequality follows from a similar argument as in the proof of Theorem 1. By (C3) and that $|\psi(x)| \leq 1$, $\psi[\alpha(y_i - \mathbf{x}_i^T \beta_\alpha^*)] x_{ij}$ is also sub-Gaussian. For any $k \geq 3$, using the relation between the k th moment and the second moment of sub-Gaussian random variables (Rivasplata, 2012),

$$\mathbb{E} |\psi[\alpha(y_i - \mathbf{x}_i^T \beta_\alpha^*)] x_{ij}|^k \leq \frac{k!}{2} L^{k-2} \mathbb{E} |\psi[\alpha(y_i - \mathbf{x}_i^T \beta_\alpha^*)] x_{ij}|^2,$$

where L is a constant depending on κ_0 only. Hence,

$$\mathbb{E} |2\alpha^{-1} \psi[\alpha(y_i - \mathbf{x}_i^T \beta_\alpha^*)] x_{ij}|^k \leq \frac{k!}{2} (2L/\alpha)^{k-2} v.$$

By Bernstein inequality (Proposition 2.9 of Massart and Picard (2007)) and note that $\mathbb{E}(\frac{2}{\alpha} \psi[\alpha(y_i -$

$\mathbf{x}_i^T \boldsymbol{\beta}_\alpha^* \mathbf{x}_i) = \mathbf{0}$, we have

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n \frac{2}{\alpha} \psi[\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*)] x_{ij} \right| \geq \sqrt{\frac{2vt}{n}} + \frac{Lt}{\alpha n} \right) \leq 2 \exp(-t).$$

Let $t = n\lambda_n^2/(32v)$ and observe that $\frac{2Lt}{\alpha n} \leq \sqrt{\frac{2vt}{n}}$ by the choice of λ_n and α . We have

$$P \left(\left| \frac{1}{n} \sum_{i=1}^n \frac{2}{\alpha} \psi[\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*)] x_{ij} \right| \geq \frac{\lambda_n}{2} \right) \leq 2 \exp \left(-\frac{n\lambda_n^2}{32v} \right).$$

It then follows from union inequality that

$$P \left(\left\| \frac{1}{n} \sum_{i=1}^n \frac{2}{\alpha} \psi[\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*)] \mathbf{x}_i \right\|_\infty > \frac{\lambda_n}{2} \right) \leq 2 \exp \left(-\frac{n\lambda_n^2}{32v} + \log p \right) = 2p^{-c_0},$$

where $c_0 = \kappa_\lambda^2/(32v) - 1$. This completes the proof. \square

Proof of Lemma 2. Denote $\mathcal{L}_n(\boldsymbol{\beta}) = \frac{1}{n} \sum_{i=1}^n \ell_\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta})$. Applying a second-order Taylor expansion to $\mathcal{L}_n(\boldsymbol{\beta})$ between $\boldsymbol{\beta}_\alpha^*$ and $\boldsymbol{\beta}_\alpha^* + \boldsymbol{\Delta}$, we conclude that for some $v \in [0, 1]$,

$$\delta \mathcal{L}_n(\boldsymbol{\Delta}, \boldsymbol{\beta}_\alpha^*) = \frac{1}{n} \sum_{i=1}^n \psi'[\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^* + v \mathbf{x}_i^T \boldsymbol{\Delta})] (\mathbf{x}_i^T \boldsymbol{\Delta})^2, \quad (8.5)$$

where $\psi'(x) = 1$ for $|x| \leq 1$, and $\psi'(x) = 0$ otherwise. Note that each term in (8.5) is nonnegative. However, the quadratic component in (8.5) is not Lipschitz continuous with a bounded Lipschitz coefficient. In order to apply the contraction theorem of Ledoux and Talagrand (1991), we introduce a truncation function that is Lipschitz and bound (8.5) from below. Let

$$\varphi_t(u) = u^2 I(|u| \leq t/2) + (t - u)^2 I(t/2 \leq |u| \leq t), \quad (8.6)$$

where $I(\cdot)$ is the indicator function. Clearly, $\varphi_t(u) \leq u^2$ and satisfies the Lipschitz condition with

Lipschitz coefficient bounded by $2t$. We first show

$$\delta \mathcal{L}_n(\mathbf{\Delta}, \boldsymbol{\beta}_\alpha^*) \geq \frac{1}{n} \sum_{i=1}^n \varphi_\tau(\mathbf{x}_i^T \mathbf{\Delta} I(|y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*| \leq T)), \quad (8.7)$$

for $0 < \alpha \leq 1/(T + \tau)$, where the thresholds T and τ will be chosen as in (8.10).

Let $A_i = \psi'[\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^* + v \mathbf{x}_i^T \mathbf{\Delta})](\mathbf{x}_i^T \mathbf{\Delta})^2$. We need only to show that

$$A_i \geq \varphi_\tau(\mathbf{x}_i^T \mathbf{\Delta} I(|y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*| \leq T)).$$

When $|y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*| > T$ or $|\mathbf{x}_i^T \mathbf{\Delta}| > \tau$, the right hand side is zero and the inequality holds trivially.

Thus, we need only to consider the case $|y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*| \leq T$ and $|\mathbf{x}_i^T \mathbf{\Delta}| \leq \tau$. In this case,

$$|\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^* + v \mathbf{x}_i^T \mathbf{\Delta})| \leq \alpha(T + \tau) \leq 1,$$

and hence $\psi'[\alpha(y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^* + v \mathbf{x}_i^T \mathbf{\Delta})] = 1$. Using this,

$$A_i = (\mathbf{x}_i^T \mathbf{\Delta})^2 \geq \varphi_\tau(\mathbf{x}_i^T \mathbf{\Delta}) = \varphi_\tau(\mathbf{x}_i^T \mathbf{\Delta} I(|y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*| \leq T)).$$

Using (8.7), to prove the Lemma, we need to show that, for any fixed $\delta \in (0, 1]$, with high probability

$$\mathbb{P}_n \varphi_\tau(\mathbf{x}^T \mathbf{\Delta} I(|y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*| \leq T)) \geq \kappa_1 \|\mathbf{\Delta}\|_2 \{ \|\mathbf{\Delta}\|_2 - \kappa_2 \sqrt{(\log p)/n} \|\mathbf{\Delta}\|_1 \}, \text{ for all } \|\mathbf{\Delta}\|_2 = \delta, \quad (8.8)$$

where constants κ_1 and κ_2 do not depend on α and

$$\mathbb{P}_n \varphi_\tau(\mathbf{x}^T \mathbf{\Delta} I(|y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*| \leq T)) = \frac{1}{n} \sum_{i=1}^n \varphi_\tau(\mathbf{x}_i^T \mathbf{\Delta} I(|y_i - \mathbf{x}_i^T \boldsymbol{\beta}_\alpha^*| \leq T)).$$

This is equivalent to proving (8.8) for $\delta = 1$. Indeed, from the definition (8.6), for any $d > 0$ and

$z \in \mathbb{R}$, we have $\varphi_d(dz) = d^2\varphi_1(z)$. Thus, the event

$$\mathbb{P}_n \varphi_{\tau_1}(\mathbf{x}^T \boldsymbol{\Delta} I(|y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*| \leq T)) \geq \kappa_1 \{1 - \kappa_2 \sqrt{(\log p)/n} \|\boldsymbol{\Delta}\|_1\}, \text{ for all } \|\boldsymbol{\Delta}\|_2 = 1 \quad (8.9)$$

is the same as the event

$$\mathbb{P}_n \varphi_{\tau_1}(\mathbf{x}^T (\boldsymbol{\Delta}/\|\boldsymbol{\Delta}\|_2) I(|y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*| \leq T)) \geq \kappa_1 \{1 - \kappa_2 \sqrt{(\log p)/n} \|\boldsymbol{\Delta}\|_1/\|\boldsymbol{\Delta}\|_2\},$$

which equals to the event (8.8) with $\tau = \delta\tau_1$.

To establish (8.9), let us consider its complementary event. Define

$$f_{\boldsymbol{\Delta}}(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\Delta} I(|y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*| \leq T) \quad \text{and} \quad g_{\boldsymbol{\Delta}}(\mathbf{x}) = \varphi_\tau(f_{\boldsymbol{\Delta}}(\mathbf{x})).$$

Let $\mathbb{S}_2(1)$ be the unit sphere with L_2 -radius one, and $\mathbb{S}_1(t)$ be the sphere of L_1 -radius t , which is to be chosen later. The complementary event of (8.9) is given by

$$\left\{ \mathbb{P}_n [g_{\boldsymbol{\Delta}}(\mathbf{x})] < \kappa_1 \{1 - \kappa_2 \sqrt{(\log p)/n} \|\boldsymbol{\Delta}\|_1\}, \text{ for some } \boldsymbol{\Delta} \in \mathbb{S}_2(1) \right\}.$$

Our goal is to show that the probability of this event is very small, which is demonstrated through the following three steps.

(a) First, we show that with the following choice of truncation

$$T^2 = 1024\kappa_0^4 \kappa_l^{-2} M_k^{2/k} \quad \text{and} \quad \tau^2 = \max\{32\kappa_0^2 \log(12\kappa_l^{-1}\kappa_0^2), 1\}, \quad (8.10)$$

for any fixed $\boldsymbol{\Delta} \in \mathbb{S}_2(1)$, we have

$$\mathbb{E}[g_{\boldsymbol{\Delta}}(\mathbf{x})] \geq \kappa_l/2. \quad (8.11)$$

(b) Second, with $Z(t) = \sup_{\boldsymbol{\Delta} \in \mathbb{S}_2(1) \cap \mathbb{S}_1(t)} |\mathbb{P}_n [g_{\boldsymbol{\Delta}}(\mathbf{x})] - \mathbb{E}[g_{\boldsymbol{\Delta}}(\mathbf{x})]|$, we prove the tail probability

bound for $Z(t)$ is bounded by

$$P(Z(t) \geq \kappa_l/4 + 40\tau^2\kappa_0t\sqrt{(\log p)/n}) \leq \exp(-c'_1n - c'_2t^2 \log p), \quad (8.12)$$

for each given t .

(c) Finally, we use a standard peeling argument (Alexander, 1987; Van de Geer, 2000) to establish

$$P\left\{\exists \Delta \in \mathbb{S}_2(1) : Z(\|\Delta\|_1) \geq \kappa_l/4 + 40\tau^2\kappa_0\|\Delta\|_1\sqrt{(\log p)/n}\right\} \leq \exp(-c_1n - c_2 \log p).$$

The result (c) together with (8.11) show that the probability of the complementary event of (8.9) with $\kappa_1 = \kappa_l/4$ and $\kappa_2 = 40\tau^2\kappa_0\kappa_1^{-1}$ is bounded by $\exp(-c_1n - c_2 \log p)$, which completes the proof.

We first prove (8.11). In fact, by condition (C2), for any $\Delta \in \mathbb{S}_2(1)$, $E[(\mathbf{x}^T \Delta)^2] \geq \kappa_l \|\Delta\|_2^2 = \kappa_l$. So, it suffices to show that $E[(\mathbf{x}^T \Delta)^2 - g_\Delta(\mathbf{x})] \leq \kappa_l/2$.

Note that, $g_\Delta(\mathbf{x}) = (\mathbf{x}^T \Delta)^2$ for all \mathbf{x} such that $|y - \mathbf{x}^T \beta_\alpha^*| \leq T$ and $|\mathbf{x}^T \Delta| \leq \tau/2$. Therefore, we have

$$E[(\mathbf{x}^T \Delta)^2 - g_\Delta(\mathbf{x})] \leq E[(\mathbf{x}^T \Delta)^2 I(|y - \mathbf{x}^T \beta_\alpha^*| \geq T)] + E[(\mathbf{x}^T \Delta)^2 I(|\mathbf{x}^T \Delta| \geq \tau/2)]. \quad (8.13)$$

To bound the first term on the right hand side of (8.13), it follows from the Cauchy-Schwartz inequality that

$$E[(\mathbf{x}^T \Delta)^2 I(|y - \mathbf{x}^T \beta_\alpha^*| \geq T)] \leq [E(\mathbf{x}^T \Delta)^4]^{1/2} [P(|y - \mathbf{x}^T \beta_\alpha^*| \geq T)]^{1/2}.$$

Since $\mathbf{x}^T \Delta$ is sub-Gaussian with parameter at most κ_0^2 by assumption (C3), we have $E(\mathbf{x}^T \Delta)^4 \leq 16\kappa_0^4$. Meanwhile, for any $0 < \alpha \leq 1/(T + \tau)$, it follows from the Chebyshev inequality and Theorem 1 that

$$T^2 P(|y - \mathbf{x}^T \beta_\alpha^*| \geq T) \leq E[(y - \mathbf{x}^T \beta_\alpha^*)^2]$$

$$\begin{aligned}
&\leq 2 \mathbb{E} \epsilon^2 + 2 \mathbb{E}[\mathbf{x}^T (\boldsymbol{\beta}^* - \boldsymbol{\beta}_\alpha^*)]^2 \\
&\leq 2M_k^{2/k} + O(\alpha^{2k-2}) \\
&\leq 4M_k^{2/k}.
\end{aligned}$$

To bound the second term on the right hand side of (8.13), by the concentration inequality of sub-Gaussian variables, we have

$$P(|\mathbf{x}^T \boldsymbol{\Delta}| \geq \tau/2) \leq 2 \exp\{-\tau^2/(8\kappa_0^2)\}.$$

Then, by the choice of T and τ in (8.10),

$$\mathbb{E}[(\mathbf{x}^T \boldsymbol{\Delta})^2 I(|y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*| \geq T)] \leq \frac{\kappa_l}{4} \quad \text{and} \quad \mathbb{E}[(\mathbf{x}^T \boldsymbol{\Delta})^2 I(|\mathbf{x}^T \boldsymbol{\Delta}| \geq \tau/2)] \leq \frac{\kappa_l}{4}.$$

Hence, (8.11) follows.

Next, we give the tail bound as in (b). Indeed, for any $\boldsymbol{\Delta} \in \mathbb{S}_2(1)$, we have $\|g_{\boldsymbol{\Delta}}\|_\infty \leq \tau^2$. Therefore, by Massart concentration inequality (Theorem 14.2 of Bühlmann and Van De Geer (2011)), for any $z > 0$, we have $P(Z(t) \geq \mathbb{E} Z(t) + z) \leq \exp(-\frac{nz^2}{32\tau^4})$. By choosing $z = \kappa_l/4 + 16\tau^2\kappa_0 t \sqrt{(\log p)/n}$, we have

$$P(Z(t) \geq \mathbb{E} Z(t) + z) \leq \exp\left(-\frac{n\kappa_l^2}{512\tau^4} - 8\kappa_0^2 t^2 \log p\right). \quad (8.14)$$

Next, we bound $\mathbb{E} Z(t)$. Let $\{\omega_i\}_{i=1}^n$ be an i.i.d. sequence of Rademacher variables. A symmetrization theorem (Theorem 14.3 of Bühlmann and Van De Geer (2011)) yields

$$\mathbb{E}[Z(t)] \leq 2 \mathbb{E} \left[\sup_{\boldsymbol{\Delta} \in \mathbb{S}_2(1) \cap \mathbb{S}_1(t)} \left| \frac{1}{n} \sum_{i=1}^n \omega_i g_{\boldsymbol{\Delta}}(\mathbf{x}_i) \right| \right] = 2 \mathbb{E} \left[\sup_{\boldsymbol{\Delta} \in \mathbb{S}_2(1) \cap \mathbb{S}_1(t)} \left| \frac{1}{n} \sum_{i=1}^n \omega_i \varphi_\tau(f_{\boldsymbol{\Delta}}(\mathbf{x}_i)) \right| \right].$$

By definition, the function φ_τ is Lipschitz with parameter at most $2\tau \leq 2\tau^2$ and $\varphi_\tau(0) = 0$. Therefore, by the Ledoux-Talagrand contraction theorem (Ledoux and Talagrand (1991), p.112),

we have

$$\begin{aligned}
\mathbb{E}[Z(t)] &\leq 8\tau^2 \mathbb{E} \left[\sup_{\Delta \in \mathbb{S}_2(1) \cap \mathbb{S}_1(t)} \left| \frac{1}{n} \sum_{i=1}^n \omega_i f_{\Delta}(\mathbf{x}_i) \right| \right] \\
&= 8\tau^2 \mathbb{E} \left[\sup_{\Delta \in \mathbb{S}_2(1) \cap \mathbb{S}_1(t)} \left| \frac{1}{n} \sum_{i=1}^n \omega_i \mathbf{x}_i^T \Delta I(|y_i - \mathbf{x}_i^T \beta_{\alpha}^*| \leq T) \right| \right] \\
&\leq 8\tau^2 t \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \omega_i \mathbf{x}_i I(|y_i - \mathbf{x}_i^T \beta_{\alpha}^*| \leq T) \right\|_{\infty}.
\end{aligned}$$

Since the variables $\{x_{ij}\}_{i=1}^n$ are zero-mean i.i.d. sub-Gaussian with parameter at most κ_0^2 , so are $\{\omega_i x_{ij} I(|y_i - \mathbf{x}_i^T \beta_{\alpha}^*| \leq T)\}_{i=1}^n$. Since $\mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \omega_i \mathbf{x}_i I(|y_i - \mathbf{x}_i^T \beta_{\alpha}^*| \leq T) \right\|_{\infty}$ is the maxima of p such terms, known bounds on the expectation of sub-Gaussian maxima (e.g. see Ledoux and Talagrand (1991), p.79) yield

$$\mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \omega_i \mathbf{x}_i I(|y_i - \mathbf{x}_i^T \beta_{\alpha}^*| \leq T) \right\|_{\infty} \leq 3\kappa_0 \sqrt{(\log p)/n}.$$

Hence,

$$\mathbb{E}[Z(t)] \leq 24\tau^2 \kappa_0 t \sqrt{(\log p)/n}. \tag{8.15}$$

Combining (8.14) and (8.15), we have

$$P\left(Z(t) \geq \kappa_l/4 + 40\tau^2 \kappa_0 t \sqrt{(\log p)/n}\right) \leq \exp(-c'_1 n - c'_2 t^2 \log p),$$

where constants c'_1 and c'_2 depends on κ_l and κ_0 only. This result holds for each given t .

Next, we furnish the peeling argument in (c). Let $h(\|\Delta\|_1) = \kappa_l/8 + 20\tau^2 \kappa_0 \|\Delta\|_1 \sqrt{(\log p)/n}$ and $B = \{\exists \Delta \in \mathbb{S}_2(1) : Z(\|\Delta\|_1) \geq 2h(\|\Delta\|_1)\}$. Since $h(\|\Delta\|_1) \geq \kappa_l/8$, the set can be covered by partition $\{B_m\}_{m=1}^{\infty}$ with $B_m = \{\Delta \in \mathbb{S}_2(1) : 2^{m-4} \kappa_l \leq h(\|\Delta\|_1) \leq 2^{m-3} \kappa_l\}$. Thus, by union bound,

$$P(B) \leq \sum_{m=1}^{\infty} P(\Delta \in B_m \text{ such that } Z(\|\Delta\|_1) \geq 2h(\|\Delta\|_1))$$

$$\leq \sum_{m=1}^{\infty} P(Z(\|\Delta\|_1) \geq 2^{m-3}\kappa_l)$$

since $Z(\|\Delta\|_1) \geq 2^{m-3}\kappa_l$ for $\Delta \in B_m$. By letting $2^{m-3}\kappa_l = \kappa_l/4 + 40\tau^2\kappa_0 t\sqrt{(\log p)/n}$ as in (8.12) and solving for t , by (8.12), we obtain

$$\begin{aligned} P(B) &\leq \sum_{m=1}^{\infty} \exp\left(-c'_1 n - \frac{c'_2 \kappa_l^2 (2^{m-1} - 1)^2 n}{\tau^4 \kappa_0^2}\right) \\ &\leq \exp(-c'_1 n) + \sum_{m=2}^{\infty} \exp\left(-c'_1 n - \frac{c'_2 n \kappa_l^2 2^{2m-4}}{\tau^4 \kappa_0^2}\right) \\ &\leq c_1 \exp(-c_2 n), \end{aligned}$$

where the last inequality follows from sum of geometric series. \square

Proof of Lemma 3. Note that,

$$R_q \geq \sum_{j=1}^p |\beta_{\alpha,j}^*|^q \geq \sum_{j \in S_{\alpha\eta}} |\beta_{\alpha,j}^*|^q \geq \eta^q |S_{\alpha\eta}|. \quad (8.16)$$

Therefore, $|S_{\alpha\eta}| \leq \eta^{-q} R_q$. Let $S_{\alpha\eta}^c = \{1, 2, \dots, p\} \setminus S_{\alpha\eta}$, we have

$$\|\beta_{S_{\alpha\eta}^c}^*\|_1 = \sum_{j \in S_{\alpha\eta}^c} |\beta_{\alpha,j}^*| = \sum_{j \in S_{\alpha\eta}^c} |\beta_{\alpha,j}^*|^q |\beta_{\alpha,j}^*|^{1-q} \leq R_q \eta^{1-q}. \quad (8.17)$$

Hence, for any $\Delta \in \mathbb{C}_{\alpha\eta}$, we have

$$\|\Delta\|_1 = \|\Delta_{S_{\alpha\eta}}\|_1 + \|\Delta_{S_{\alpha\eta}^c}\|_1 \leq 4\|\Delta_{S_{\alpha\eta}}\|_1 + 4\|\beta_{S_{\alpha\eta}^c}^*\|_1.$$

By the Cauchy-Schwartz inequality and (8.17), we can bound further that

$$\|\Delta\|_1 \leq 4\sqrt{|S_{\alpha\eta}|} \|\Delta\|_2 + 4R_q \eta^{1-q} \leq 4R_q^{1/2} \eta^{-q/2} \|\Delta\|_2 + 4R_q \eta^{1-q}.$$

It then follows from Lemma 2 that

$$\begin{aligned}\delta\mathcal{L}_n(\mathbf{\Delta}, \beta_\alpha^*) &\geq \kappa_1 \|\mathbf{\Delta}\|_2 \{ \|\mathbf{\Delta}\|_2 - \kappa_2 \sqrt{(\log p)/n} [4R_q^{1/2} \eta^{-q/2} \|\mathbf{\Delta}\|_2 + 4R_q \eta^{1-q}] \} \\ &= \left(\kappa_1 - 4\kappa_1 \kappa_2 R_q^{1/2} \eta^{-q/2} \sqrt{(\log p)/n} \right) \|\mathbf{\Delta}\|_2^2 - 4\kappa_2 R_q \eta^{1-q} \sqrt{(\log p)/n}.\end{aligned}$$

With $\lambda_n = \kappa_\lambda \sqrt{(\log p)/n}$ and $\eta = \lambda_n$, it holds that

$$4\kappa_1 \kappa_2 R_q^{1/2} \eta^{-q/2} \sqrt{\frac{\log p}{n}} = 4\kappa_1 \kappa_2 R_q^{1/2} \kappa_\lambda^{-q/2} \left(\frac{\log p}{n} \right)^{(1-q)/2},$$

which is no larger than $\kappa_1/2$ under assumption (2.7). On the other hand,

$$4R_q \kappa_2 \eta^{1-q} \sqrt{\frac{\log p}{n}} = 4R_q \kappa_2 \kappa_\lambda^{1-q} \left(\frac{\log p}{n} \right)^{1-(q/2)}.$$

Therefore, RSC holds with $\kappa_{\mathcal{L}} = \frac{\kappa_1}{2}$ and $\tau_{\mathcal{L}}^2 = 4R_q \kappa_2 \kappa_\lambda^{1-q} \left(\frac{\log p}{n} \right)^{1-(q/2)}$. \square

Proof of Theorem 2. Let A_1 and A_2 denote the events that Lemma 1 and Lemma 3 hold, respectively. By Theorem 1 of Negahban, *et al.* (2012), within $A_1 \cap A_2$, it holds that

$$\begin{aligned}\|\mathbf{\Delta}\|_2^2 &\leq 9 \frac{\lambda_n^2}{\kappa_{\mathcal{L}}^2} |S_{\alpha\eta}| + \frac{\lambda_n}{\kappa_{\mathcal{L}}^2} \{ 2\tau_{\mathcal{L}}^2 + 4\|\beta_{S_{\alpha\eta}^c}^*\|_1 \} \\ &\leq \frac{36\lambda_n^2 R_q}{\kappa_1^2 \eta^q} + \frac{4\lambda_n}{\kappa_1^2} \left\{ 8R_q \kappa_2 \kappa_\lambda^{1-q} \left(\frac{\log p}{n} \right)^{1-(q/2)} + 4R_q \eta^{1-q} \right\} \\ &\stackrel{(i)}{=} \frac{36}{\kappa_1^2} R_q \lambda_n^{2-q} + \frac{16}{\kappa_1^2} R_q \lambda_n^{2-q} \left\{ 2\kappa_2 \left(\frac{\log p}{n} \right)^{\frac{1-q}{2}} + 1 \right\} \\ &= O(R_q \lambda_n^{2-q}) = O\left(R_q [(\log p)/n]^{1-(q/2)} \right),\end{aligned}$$

where (i) follows from the choice of $\eta = \lambda_n$. On the other hand, by Lemma 1 and 3, $P(A_1 \cap A_2) \geq 1 - 2p^{-c_0} - c_1 \exp(-c_2 n)$. \square

Proof of Lemma 4. From the proof of Lemma 2, we can see that (2.6) indeed holds for all β and

$\Delta \in \{\Delta : \|\Delta\|_2 \leq 1\}$ that

$$\delta\mathcal{L}_n(\Delta, \beta) \geq \kappa_1 \|\Delta\|_2^2 - \kappa_1 \kappa_2 \|\Delta\|_2 \|\Delta\|_1 \sqrt{(\log p)/n}.$$

Using the fact that $ab \leq (a^2 + b^2)/2$, we conclude that

$$\delta\mathcal{L}_n(\Delta, \beta) \geq \kappa_1 \|\Delta\|_2^2 - \left(\frac{1}{2} \kappa_1 \|\Delta\|_2^2 + \frac{1}{2} \kappa_1 \kappa_2^2 \|\Delta\|_1^2 \left(\frac{\log p}{n} \right) \right).$$

Therefore, (3.2) holds with $\gamma_l = \kappa_1$ and $\tau_l = \kappa_1 \kappa_2^2 (\log p)/(2n)$. Meanwhile, since $|\psi'(\cdot)| \leq 1$, it follows from (8.5) that

$$\delta\mathcal{L}_n(\Delta, \beta) \leq \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \Delta)^2.$$

Under the sub-Gaussianity assumption (C3), it follows from some existing work (e.g. page 18 of Loh and Wainwright (2013)) that, with probability great than $1 - c_1 \exp(-c_2 n)$, it holds that

$$\frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \Delta)^2 \leq \kappa_u \left(\frac{3}{2} \|\Delta\|_2^2 + \frac{\log p}{n} \|\Delta\|_1^2 \right),$$

where c_1 and c_2 are some generic constants. Hence, (3.3) holds with $\gamma_u = 3\kappa_u$ and $\tau_u = \kappa_u (\log p)/n$. □

Proof of Theorem 4. We prove the theorem by the following two steps:

(a) We first show that, for any $\delta^2 \geq \varepsilon^2/(1 - \kappa)$, $\phi(\widehat{\beta}^t) - \phi(\widehat{\beta}) \leq \delta^2$, for all t greater than the right hand side of (8.20), where $\kappa \in [0, 1)$ is a contraction constant and ε is a tolerance parameter, which will be given in (8.21) and (8.22), respectively.

(b) We use RSC condition (3.2) to transform the upper bound of $\phi(\widehat{\beta}^t) - \phi(\widehat{\beta})$ into the upper bound of $\|\widehat{\beta}^t - \widehat{\beta}\|_2$.

For step (a), since our loss function is convex, we apply Theorem 2 of Agarwal, Negahban, and Wainwright (2012). In order for our proof to be self-contained, we cite their theorem as the follows:

[Theorem 2 of Agarwal, Negahban, and Wainwright (2012)] Suppose for any data set Z_1^n , the loss function $\mathcal{L}_n(\cdot; Z_1^n)$ is convex and differentiable and the regularizer \mathcal{R} is a norm. Consider

the optimization problem of $\hat{\theta} = \operatorname{argmin}_{\mathcal{R}(\theta) \leq \rho} \{\mathcal{L}_n(\theta; Z_1^n) + \lambda_n \mathcal{R}(\theta)\}$ for a radius ρ such that θ^* is feasible, where $\theta^* = \operatorname{argmin} \mathbb{E} \mathcal{L}_n(\theta; Z_1^n)$, and a regularization parameter λ_n satisfying bound

$$\lambda_n \geq 2\mathcal{R}^*(\nabla \mathcal{L}_n(\theta^*)), \quad (8.18)$$

where \mathcal{R}^* is the dual norm of the regularizer. In addition, suppose that the loss function \mathcal{L}_n satisfies the RSC/RSM condition with parameters (γ_l, τ_l) and (γ_u, τ_u) , respectively. Let $(\mathcal{M}, \bar{\mathcal{M}}^\perp)$ be any \mathcal{R} -decomposable pair of subspaces such that

$$\kappa = \left\{ 1 - \frac{\bar{\gamma}_l}{4\gamma_u} + \frac{64\Psi^2(\bar{\mathcal{M}})\tau_u}{\bar{\gamma}_l} \right\} \xi \in [0, 1) \quad \text{and} \quad \frac{32\rho}{1-\kappa} \xi \chi \leq \lambda_n, \quad (8.19)$$

where $\Psi(\bar{\mathcal{M}}) = \sup_{\theta \in \bar{\mathcal{M}} \setminus \{0\}} \mathcal{R}(\theta) / \|\theta\|_2$, $\bar{\gamma}_l = \gamma_l - 64\tau_l\Psi^2(\bar{\mathcal{M}})$, $\xi = (1 - 64\tau_u\bar{\gamma}_l^{-1}\Psi^2(\bar{\mathcal{M}}))^{-1}$, and $\chi = 2(\bar{\gamma}_l/(4\gamma_u) + 128\tau_u\bar{\gamma}_l^{-1}\Psi^2(\bar{\mathcal{M}}))\tau_l + 8\tau_u + 2\tau_l$. Denote $\varepsilon^2 = 8\xi\chi \left(6\Psi(\bar{\mathcal{M}})\|\hat{\theta} - \theta^*\|_2 + 8\mathcal{R}(\Pi_{\mathcal{M}^\perp}(\theta^*)) \right)^2$, where $\Pi_{\mathcal{M}^\perp}(\theta^*)$ is the projection of θ^* onto \mathcal{M}^\perp . Then for any $\delta^2 \geq \varepsilon^2/(1-\kappa)$, we have $\phi_n(\hat{\theta}^t) - \phi_n(\hat{\theta}) \leq \delta^2$ for all

$$t \geq \frac{2 \log((\phi_n(\theta^0) - \phi_n(\hat{\theta}))/\delta^2)}{\log(1/\kappa)} + \log_2 \log_2 \left(\frac{\rho\lambda_n}{\delta^2} \right) \left(1 + \frac{\log 2}{\log(1/\kappa)} \right), \quad (8.20)$$

where $\phi_n(\theta) = \mathcal{L}_n(\theta; Z_1^n) + \lambda_n \mathcal{R}(\theta)$, $\hat{\theta}^t$ is the solution by the gradient descent algorithm after t^{th} iteration, and θ^0 is the initial value of θ .

In fact, Theorem 2 of Agarwal, Negahban, and Wainwright (2012) is a deterministic statement for all choices of pairs $(\mathcal{M}, \bar{\mathcal{M}}^\perp)$. From Lemma 1 and Lemma 4, we have shown that with our choice of λ_n , the RA-quadratic loss function satisfy (8.18) and RSC/RSM with probability at least $1 - 2p^{-c_0}$ and $1 - c_1 \exp(-c_2 n)$, respectively. Hence, Theorem 2 of Agarwal, Negahban, and Wainwright (2012) applies to our problem with high probability. We further choose the pair $(\mathcal{M}, \bar{\mathcal{M}}^\perp) = (S_{\alpha\eta}, S_{\alpha\eta}^c)$ and give the explicit expression of constants for our problem as the follows:

$$\kappa = \left\{ 1 - \frac{\bar{\gamma}_l}{4\gamma_u} + \frac{64\kappa_u |S_{\alpha\eta}| \frac{\log p}{n}}{\bar{\gamma}_l} \right\} \left(1 - \frac{64\kappa_u |S_{\alpha\eta}| \frac{\log p}{n}}{\bar{\gamma}_l} \right)^{-1} \quad (8.21)$$

$$\varepsilon^2 = 8\xi\chi \left(6\sqrt{|S_{\alpha\eta}|} \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2 + 8\|\boldsymbol{\beta}_{S_{\alpha\eta}^c}^*\|_1 \right)^2, \quad (8.22)$$

where $\bar{\gamma}_l = \kappa_1 - 32\kappa_1\kappa_2^2|S_{\alpha\eta}|(\log p)/n$, $\xi = \{1 - 64\kappa_u|S_{\alpha\eta}|(\log p)/(n\bar{\gamma}_l)\}^{-1}$, and $\chi = 2\{\bar{\gamma}_l/(4\gamma_u) + 128\tau_u|S_{\alpha\eta}|/\bar{\gamma}_l + 1\}\tau_l + 8\tau_u$. It remains to check (8.19). By (8.21), $\kappa \in [0, 1)$ is equivalent to requiring

$$|S_{\alpha\eta}| \frac{\log p}{n} < \frac{\bar{\gamma}_l^2}{1536\kappa_u^2} \quad (8.23)$$

With $\eta = \lambda_n$, it follows from (8.16) that

$$|S_{\alpha\eta}| \frac{\log p}{n} \leq R_q \eta^{-q} \frac{\log p}{n} \leq \kappa_\lambda^{-q} R_q \left(\frac{\log p}{n} \right)^{1-(q/2)}.$$

Hence, (8.23) holds when n is sufficiently large. Moreover, from (8.19) we need

$$\lambda_n \geq \frac{32\rho}{1-\kappa} \left(1 - \frac{64\kappa_u|S_{\alpha\eta}| \frac{\log p}{n}}{\bar{\gamma}_l} \right)^{-1} \left[1 + \kappa_1\kappa_2^2 \left(\frac{\bar{\gamma}_l}{12\kappa_u} + \frac{128\kappa_u|S_{\alpha\eta}| \frac{\log p}{n}}{\bar{\gamma}_l} \right) + 8\kappa_u \right] \frac{\log p}{n},$$

which is satisfied under the stated assumption. It then follows from Theorem 2 of Agarwal, Negahban, and Wainwright (2012) that, for any $\delta^2 \geq \varepsilon^2/(1-\kappa)$, $\phi(\widehat{\boldsymbol{\beta}}^t) - \phi(\widehat{\boldsymbol{\beta}}) \leq \delta^2$, for all iterations t greater than the right hand side of (8.20).

For step (b), it follows from RSC condition that

$$\mathcal{L}_n(\widehat{\boldsymbol{\beta}}^t) - \mathcal{L}_n(\widehat{\boldsymbol{\beta}}) - [\nabla \mathcal{L}_n(\widehat{\boldsymbol{\beta}})]^T (\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}) \geq \frac{\gamma_l}{2} \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 - \tau_l \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1^2.$$

Then we have

$$\begin{aligned} \phi(\widehat{\boldsymbol{\beta}}^t) - \phi(\widehat{\boldsymbol{\beta}}) &= \mathcal{L}_n(\widehat{\boldsymbol{\beta}}^t) - \mathcal{L}_n(\widehat{\boldsymbol{\beta}}) + \lambda_n (\|\widehat{\boldsymbol{\beta}}^t\|_1 - \|\widehat{\boldsymbol{\beta}}\|_1) \\ &\geq [\nabla \mathcal{L}_n(\widehat{\boldsymbol{\beta}})]^T (\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}) + \lambda_n (\|\widehat{\boldsymbol{\beta}}^t\|_1 - \|\widehat{\boldsymbol{\beta}}\|_1) + \frac{\gamma_l}{2} \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 - \tau_l \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1^2. \end{aligned}$$

Since $\widehat{\boldsymbol{\beta}}$ is the minimizer of $\phi(\boldsymbol{\beta})$, by the first-order condition, $[\nabla \mathcal{L}_n(\widehat{\boldsymbol{\beta}}) + \lambda_n \nabla \|\widehat{\boldsymbol{\beta}}\|_1]^T (\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}) \geq 0$.

Therefore,

$$\phi(\widehat{\boldsymbol{\beta}}^t) - \phi(\widehat{\boldsymbol{\beta}}) \geq -\lambda_n [\nabla \|\widehat{\boldsymbol{\beta}}\|_1]^T (\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}) + \lambda_n (\|\widehat{\boldsymbol{\beta}}^t\|_1 - \|\widehat{\boldsymbol{\beta}}\|_1) + \frac{\gamma_l}{2} \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 - \tau_l \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1^2.$$

By the convexity of the L_1 -norm, $\|\widehat{\boldsymbol{\beta}}^t\|_1 - \|\widehat{\boldsymbol{\beta}}\|_1 - [\nabla \|\widehat{\boldsymbol{\beta}}\|_1]^T (\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}) \geq 0$. Hence,

$$\phi(\widehat{\boldsymbol{\beta}}^t) - \phi(\widehat{\boldsymbol{\beta}}) \geq \frac{\gamma_l}{2} \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 - \tau_l \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1^2. \quad (8.24)$$

Next, we bound $\|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1$. It follows from Lemma 3 of Agarwal, Negahban, and Wainwright (2012) that

$$\|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1 \leq 2 \left(2\sqrt{|S_{\alpha\eta}|} \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2 + 4\sqrt{|S_{\alpha\eta}|} \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2 + 4\|\boldsymbol{\beta}_{S_{\alpha\eta}^c}^*\|_1 + \delta^2/\lambda_n \right),$$

where δ is defined as in (a). Then, by the Cauchy-Schwartz inequality,

$$\|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_1^2 \leq 16 \left(4|S_{\alpha\eta}| \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 + 16|S_{\alpha\eta}| \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + 16\|\boldsymbol{\beta}_{S_{\alpha\eta}^c}^*\|_1^2 + \delta^4/\lambda_n^2 \right). \quad (8.25)$$

Equations (8.24) and (8.25) together with results in (a) imply that,

$$\delta^2 \geq \frac{\gamma_l}{2} \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 - 16\tau_l \left(4|S_{\alpha\eta}| \|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 + 16|S_{\alpha\eta}| \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + 16\|\boldsymbol{\beta}_{S_{\alpha\eta}^c}^*\|_1^2 + \delta^4/\lambda_n^2 \right).$$

Letting $\tilde{\gamma}_l = \gamma_l/2 - 64\tau_l|S_{\alpha\eta}|$, we have

$$\|\widehat{\boldsymbol{\beta}}^t - \widehat{\boldsymbol{\beta}}\|_2^2 \leq \frac{1}{\tilde{\gamma}_l} \left(\delta^2 + \frac{16\tau_l\delta^4}{\lambda_n^2} \right) + \frac{256\tau_l}{\tilde{\gamma}_l} (|S_{\alpha\eta}| \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + \|\boldsymbol{\beta}_{S_{\alpha\eta}^c}^*\|_1^2). \quad (8.26)$$

We now bound the second term in (8.26). By (8.16) and (8.17), we have

$$\begin{aligned} |S_{\alpha\eta}| \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + \|\boldsymbol{\beta}_{S_{\alpha\eta}^c}^*\|_1^2 &\leq R_q \eta^{-q} \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + R_q^2 \eta^{2-2q} \\ &\leq R_q \kappa_\lambda^{-q} \left(\frac{\log p}{n} \right)^{-q/2} \|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + \kappa_\lambda^{-q} R_q^2 \left(\frac{\log p}{n} \right)^{1-q} \\ &\leq \kappa_\lambda^{-q} R_q \left(\frac{\log p}{n} \right)^{-q/2} \left[\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^*\|_2^2 + R_q \left(\frac{\log p}{n} \right)^{1-(q/2)} \right]. \end{aligned} \quad (8.27)$$

Meanwhile, from (a) we have

$$\begin{aligned}
\delta^2 &= \frac{\varepsilon^2}{1-\kappa} = \frac{8\xi\chi}{1-\kappa} \left(6\sqrt{|S_{\alpha\eta}|} \|\widehat{\beta} - \beta_\alpha^*\|_2 + 8\|\beta_{S_{\alpha\eta}^c}^*\|_1 \right)^2 \\
&\leq \frac{8\xi\chi}{1-\kappa} (72|S_{\alpha\eta}| \|\widehat{\beta} - \beta_\alpha^*\|_2^2 + 128\|\beta_{S_{\alpha\eta}^c}^*\|_1^2) \\
&\leq \frac{1024\xi\chi}{1-\kappa} (|S_{\alpha\eta}| \|\widehat{\beta} - \beta_\alpha^*\|_2^2 + \|\beta_{S_{\alpha\eta}^c}^*\|_1^2).
\end{aligned} \tag{8.28}$$

Since $\bar{\gamma}_l \asymp 1$, $\kappa \asymp 1$, $\xi \asymp 1$, $\chi \asymp \frac{\log p}{n}$, and $\tau_l \asymp \frac{\log p}{n}$, it follows from (8.26), (8.27) and (8.28) that

$$\|\widehat{\beta}^t - \widehat{\beta}\|_2^2 = O \left(R_q \left(\frac{\log p}{n} \right)^{1-(q/2)} \left[\|\widehat{\beta} - \beta_\alpha^*\|_2^2 + R_q \left(\frac{\log p}{n} \right)^{1-(q/2)} \right] \right).$$

□

Proof of Theorem 5. The proof follows the same spirit of the proof of Proposition 2.4 of Catoni (2012). The influence function $\psi(x)$ satisfies

$$-\log(1-x+x^2) \leq \psi(x) \leq \log(1+x+x^2).$$

Using this and independence, with $r(\theta) = \frac{1}{\alpha n} \sum_{i=1}^n \psi[\alpha(Y_i - \theta)]$, we have

$$\begin{aligned}
\mathbb{E} \{ \exp[\alpha n r(\theta)] \} &\leq \left(\mathbb{E} \{ \exp\{ \psi[\alpha(Y_i - \theta)] \} \} \right)^n \\
&\leq \{ 1 + \alpha(\mu - \theta) + \alpha^2[\sigma^2 + (\mu - \theta)^2] \}^n \\
&\leq \exp \{ n\alpha(\mu - \theta) + n\alpha^2[v^2 + (\mu - \theta)^2] \}.
\end{aligned}$$

Similarly, $\mathbb{E} \{ \exp[-\alpha n r(\theta)] \} \leq \exp \{ -n\alpha(\mu - \theta) + n\alpha^2[v^2 + (\mu - \theta)^2] \}$. Define

$$\begin{aligned}
B_+(\theta) &= \mu - \theta + \alpha[v^2 + (\mu - \theta)^2] + \frac{\log(1/\delta)}{n\alpha} \\
B_-(\theta) &= \mu - \theta - \alpha[v^2 + (\mu - \theta)^2] - \frac{\log(1/\delta)}{n\alpha}
\end{aligned}$$

By Chebyshev inequality,

$$P(r(\theta) > B_+(\theta)) \leq \frac{\mathbb{E} \{ \exp[\alpha n r(\theta)] \}}{\exp\{\alpha n(\mu - \theta) + n\alpha^2[v^2 + (\mu - \theta)^2] + \log(1/\delta)\}} \leq \delta$$

Similarly, $P(r(\theta) < B_-(\theta)) \leq \delta$.

Let θ_+ be the smallest solution of the quadratic equation $B_+(\theta_+) = 0$ and θ_- be the largest solution of the equation $B_-(\theta_-) = 0$. Under the assumption that $\frac{\log(1/\delta)}{n} \leq 1/8$ and the choice of $\alpha = \sqrt{\frac{\log(1/\delta)}{nv^2}}$, we have $\alpha^2 v^2 + \frac{\log(1/\delta)}{n} \leq 1/4$. Therefore,

$$\begin{aligned} \theta_+ &= \mu + 2 \left(\alpha v^2 + \frac{\log(1/\delta)}{\alpha n} \right) \left(1 + \sqrt{1 - 4 \left(\alpha^2 v^2 + \frac{\log(1/\delta)}{n} \right)} \right)^{-1} \\ &\leq \mu + 2 \left(\alpha v^2 + \frac{\log(1/\delta)}{\alpha n} \right). \end{aligned}$$

Similarly,

$$\begin{aligned} \theta_- &= \mu - 2 \left(\alpha v^2 + \frac{\log(1/\delta)}{\alpha n} \right) \left(1 + \sqrt{1 - 4 \left(\alpha^2 v^2 + \frac{\log(1/\delta)}{n} \right)} \right)^{-1} \\ &\geq \mu - 2 \left(\alpha v^2 + \frac{\log(1/\delta)}{\alpha n} \right). \end{aligned}$$

With $\alpha = \sqrt{\frac{\log(1/\delta)}{nv^2}}$, $\theta_+ \leq \mu + 4v\sqrt{\frac{\log(1/\delta)}{n}}$, $\theta_- \geq \mu - 4v\sqrt{\frac{\log(1/\delta)}{n}}$. Since the map $\theta \mapsto r(\theta)$ is non-increasing, under event $\{B_-(\theta) \leq r(\theta) \leq B_+(\theta)\}$

$$\mu - 4v\sqrt{\frac{\log(1/\delta)}{n}} \leq \theta_- \leq \hat{\mu}_\alpha \leq \theta_+ \leq \mu + 4v\sqrt{\frac{\log(1/\delta)}{n}},$$

i.e. $|\hat{\mu}_\alpha - \mu| \leq 4v\sqrt{\frac{\log(1/\delta)}{n}}$. Meanwhile, $P(B_-(\theta) \leq r(\theta) \leq B_+(\theta)) > 1 - 2\delta$. □

Proof of Theorem 6. First, we prove that the approximation error rate $\|\beta_\alpha^{c*} - \beta^*\|_2 = O(\alpha^{k-1})$, where $\beta_\alpha^{c*} = \operatorname{argmin}_\beta \mathbb{E} \ell_\alpha^c(y - \mathbf{x}^T \beta^*)$ is the population minimizer under the Catoni loss. Let

$g_\alpha(x) = \ell(x) - \ell_\alpha^c(x) = \int_0^x [2t - \frac{2}{\alpha}\psi_c(\alpha t)]dt$. It follows from (8.2) that

$$\mathbb{E}[\ell(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^*) - \ell(y - \mathbf{x}^T \boldsymbol{\beta}^*)] \leq \mathbb{E}[|g'_\alpha(y - \mathbf{x}^T \tilde{\boldsymbol{\beta}}) \mathbf{x}^T (\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*)|],$$

where $\tilde{\boldsymbol{\beta}}$ is a vector lying between $\boldsymbol{\beta}^*$ and $\boldsymbol{\beta}_\alpha^{c*}$. Since $|(\psi_c)'''| \leq 3$, by the second-order Taylor expansion with an integral remainder,

$$|g'_\alpha(x)| = |2x - \frac{2}{\alpha}\psi_c(\alpha x)| = \left| \frac{\alpha^2}{3} \int_0^x (\psi_c)'''(\alpha s)(x-s)^2 ds \right| \leq \alpha^2 |x|^3. \quad (8.29)$$

Hence, by the Cauchy-Schwartz inequality,

$$\begin{aligned} \mathbb{E}[\ell(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^{c*}) - \ell(y - \mathbf{x}^T \boldsymbol{\beta}^*)] &\leq \alpha^2 \mathbb{E}[|y - \mathbf{x}^T \tilde{\boldsymbol{\beta}}|^3 |\mathbf{x}^T (\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*)|] \\ &\leq 4\alpha^2 \mathbb{E}[(M_3 + |\mathbf{x}^T (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^3) |\mathbf{x}^T (\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*)|] \\ &\leq 4\alpha^2 [\lambda_{\max}(\mathbb{E}[(M_3 + |\mathbf{x}^T (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^3)^2 \mathbf{x} \mathbf{x}^T])]^{1/2} \|\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*\|_2, \end{aligned}$$

which is of order $O(\alpha^2 \|\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*\|_2)$, as $\lambda_{\max}(\mathbb{E}[(M_3 + |\mathbf{x}^T (\tilde{\boldsymbol{\beta}} - \boldsymbol{\beta}^*)|^3)^2 \mathbf{x} \mathbf{x}^T]) = O(1)$ by an analogous argument as in the proof of Theorem 1. Similarly as in (8.1),

$$\mathbb{E}[\ell(y - \mathbf{x}^T \boldsymbol{\beta}_\alpha^{c*}) - \ell(y - \mathbf{x}^T \boldsymbol{\beta}^*)] \geq \kappa_l \|\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*\|_2^2.$$

Hence, $\|\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*\|_2 = O(\alpha^2)$. If ϵ only has the second moment exist, by a first-order Taylor expansion of $g'_\alpha(x)$ similarly as in (8.29), we have $\|\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*\|_2 = O(\alpha)$. Next, since $(\psi_c)'(0) = 1$, by the same argument as in the proof of Lemma 3, RSC holds for Catoni's loss with probability no less than $1 - c_1 \exp(-c_2 n)$. Hence, similarly as in Theorem 2, $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}_\alpha^{c*}\|_2 = O(\sqrt{R_q}[(\log p)/n]^{1/2-q/4})$. This together with $\|\boldsymbol{\beta}_\alpha^{c*} - \boldsymbol{\beta}^*\|_2 = O(\alpha^{k-1})$ completes the proof. \square

Proof of Theorem 7. First of all, observe that

$$\hat{\sigma}^2 - \sigma^2 = \frac{1}{K} \sum_{k=1}^K \frac{1}{m} \sum_{i \in \text{fold } k} \left(\epsilon_i - (\mathbf{x}_i^T \hat{\boldsymbol{\beta}}^{(-k)} - \mathbf{x}_i^T \boldsymbol{\beta}^*) \right)^2 - \sigma^2$$

$$= \frac{1}{n} \sum_{i=1}^n \epsilon_i^2 - \sigma^2 - \frac{1}{K} \sum_{k=1}^K \frac{2}{m} \sum_{i \in \text{fold } k} \epsilon_i \mathbf{x}_i^T (\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*) + \frac{1}{K} \sum_{k=1}^K \frac{1}{m} \sum_{i \in \text{fold } k} (\mathbf{x}_i^T (\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*))^2.$$

Given that $E\epsilon^4$ exists, by Central Limit Theorem, $\sqrt{n}(\frac{1}{n} \sum_{i=1}^n \epsilon_i^2 - \sigma^2) \xrightarrow{D} \mathcal{N}(0, E\epsilon^4 - \sigma^4)$. Let $z_i = \mathbf{x}_i^T (\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*)$. We now need to prove that the last two terms are negligible. Conditioning on data outside the k th fold,

$$E \left\{ \frac{1}{m} \left(\sum_{i \in \text{fold } k} \epsilon_i z_i \right)^2 \right\} = E(\epsilon_i^2 z_i^2) \leq \sigma^2 \kappa_0^2 \|\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*\|_2^2.$$

Hence, $m^{-1/2} \sum_{i \in \text{fold } k} \epsilon_i \mathbf{x}_i^T (\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*) = O_P(\|\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*\|_2) = o_P(1)$, where the last equality follows from Theorem 3. By an analogous argument, we have

$$\frac{1}{m} \sum_{i \in \text{fold } k} (\mathbf{x}_i^T (\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*))^2 = O_p(\|\hat{\boldsymbol{\beta}}^{(-k)} - \boldsymbol{\beta}^*\|_2^2) = O_p(\max\{\alpha^{2(k-1)}, R_q[(\log p)/n]^{1-q/2}\}) = o(1/\sqrt{n}).$$

This completes the proof of the Theorem. □

References

- Agarwal, A., Negahban, S., and Wainwright, M. J. (2012). Fast global convergence of gradient methods for high-dimensional statistical recovery. *The Annals of Statistics* **40**, 2452–2482.
- Alexander, K. S. (1987). Rates of growth and sample moduli for weighted empirical processes indexed by sets. *Probability Theory and Related Fields* **75**, 379–423.
- Belloni, A. and Chernozhukov, V. (2011). L_1 -penalized quantile regression in high-dimensional sparse models. *The Annals of Statistics* **39**, 82–130.
- Bickel, P. J. and Levina, E. (2008). Covariance regularization by thresholding. *Ann. Statist.*, **36**, 2577–2604.

- Bickel, P.J., Ritov, Y. and Tsybakov, A. (2009). Simultaneous analysis of Lasso and Dantzig selector. *Ann. Statist.*, **37**, 1705-1732.
- Bühlmann, P. and Van De Geer, S. (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer.
- Catoni, O. (2012). Challenging the empirical mean and empirical variance: a deviation study. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques* **48**, 1148–1185.
- Duchi, J., Shalev-Shwartz, S., Singer, Y., and Chandra, T. (2008). Efficient projections onto the L_1 -ball for learning in high dimensions. *Proceedings of the 25th international conference on Machine learning*, 272–279.
- Efron, B. (2010). Correlated z-values and the accuracy of large-scale statistical estimates. *Jour Ameri. Statist. Assoc.*, **105**, 1042-1055.
- Engle, R.F. (1982). Autoregressive conditional heteroscedasticity with estimates of the variance of U.K. inflation. *Econometrica*, **50**, 987-1008.
- Fan, J., Fan, Y., and Barut, E. (2014). Adaptive robust variable selection. *The Annals of Statistics* **42**, 324–351.
- Fan, J., Han, X., and Gu, W.(2012). Estimating false discovery proportion under arbitrary covariance dependence (with discussion). *Journal of American Statistical Association*, **107**, 1019–1048.
- Fan, J. and Li, R. (2001). Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* **96**, 1348–1360.
- Fan, J., Liao, Y., and Mincheva, M. (2013). Large covariance estimation by thresholding principal orthogonal complements (with discussion). *Jour. Roy. Statist. Soc. B*, **75**, 603-680.
- Fan, J. and Lv, J. (2011). Non-concave penalized likelihood with NP-Dimensionality. *IEEE – Information Theory*, **57**, 5467-5484.

- Huang, C. C., Liu, K., Pope, R. M., Du, P., Lin, S., Rajamannan, N. M., et al. (2011). Activated TLR signaling in atherosclerosis among women with lower Framingham risk score: the multi-ethnic study of atherosclerosis. *PLoS ONE*, **6**, e21067.
- Huber, P. J. (1964). Robust estimation of a location parameter. *The Annals of Mathematical Statistics* **35**, 73–101.
- Ledoux, M. and Talagrand, M. (1991). *Probability in Banach Spaces: isoperimetry and processes*. Springer.
- Li, Y. and Zhu, J. (2008). L_1 -norm quantile regression. *Journal of Computational and Graphical Statistics* **17**, 163–185.
- Loh, P.-L. and Wainwright, M. J. (2013). Regularized M -estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Advances in Neural Information Processing Systems*, 476–484.
- Massart, P. and Picard, J. (2007). *Concentration inequalities and model selection*. Springer.
- Negahban, S. N., Ravikumar, P., Wainwright, M. J., and Yu, B. (2012). A unified framework for high-dimensional analysis of M -estimators with decomposable regularizers. *Statistical Science* **27**, 538–557.
- Nesterov, Y. (2007). Gradient methods for minimizing composite objective function. *Technical Report 76, Center for Operations Research and Econometrics (CORE), Catholic Univ. Lowain (UCL)*.
- Raskutti, G., Wainwright, M. J., and Yu, B. (2011). Minimax rates of estimation for high-dimensional linear regression over L_q -balls. *Information Theory, IEEE Transactions on* **57**, 6976–6994.
- Rivasplata, O. (2012). Subgaussian random variables: an expository note.

- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society, B*, **58**, 267–288.
- Van de Geer, S. (2000). *Empirical Processes in M-estimation*. Cambridge university press Cambridge.
- Wang, L. (2013). The L_1 penalized LAD estimator for high dimensional linear regression. *Journal of Multivariate Analysis* **120**, 135–151.
- Wu, Y. and Liu, Y. (2009). Variable selection in quantile regression. *Statistica Sinica* **19**, 801.
- Zhang, C.-H. (2010). Nearly unbiased variable selection under minimax concave penalty. *The Annals of Statistics* **38**, 894–942.
- Zou, H. and Yuan, M. (2008). Composite quantile regression and the oracle model selection theory. *The Annals of Statistics* **36**, 1108–1126.