# Robustness and Tractability for Non-convex M-estimators

Ruizhi Zhang, Yajun Mei, Jianjun Shi, and Huan Xu

H. Milton Stewart School of Industrial and Systems Engineering,

Georgia Institute of Technology

June 7, 2019

## Abstract

We investigate two important properties of M-estimator, namely, robustness and tractability, in linear regression setting, when the observations are contaminated by some arbitrary outliers. Specifically, robustness means the statistical property that the estimator should always be close to the underlying true parameters *regardless of the distribution of the outliers*, and tractability indicates the computational property that the estimator can be computed efficiently, even if the objective function of the M-estimator is *non-convex*. In this article, by learning the landscape of the empirical risk, we show that under mild conditions, many M-estimators enjoy nice robustness and tractability properties simultaneously, when the percentage of outliers is small. We further extend our analysis to the high-dimensional setting, where the number of parameters is greater than the number of samples, $p \gg n$, and prove that when the proportion of outliers is small, the penalized M-estimators with $L_1$ penalty will enjoy robustness and tractability simultaneously. Our research provides an analytic approach to see the effects of outliers and tuning parameters on the robustness and tractability for some families of M-estimators. Simulation and case study are presented to illustrate the usefulness of our theoretical results for M-estimators under Welsch's exponential squared loss.

*Keywords:* computational tractability, gross error, high-dimensionality, non-convexity,robust regression, sparsity

1

# 1 Introduction

M-estimation plays an important role in linear regression due to its robustness and flexibility. From the statistical viewpoint, it has been shown that many M-estimators enjoy desirable robustness properties in the presence of outliers, as well as asymptotic normality when the data are normally distributed without outliers. Some general theoretical properties and review of robust M-estimators can be found in Bai et al. (1992); Huber and Ronchetti (2009); Cheng et al. (2010); Hampel et al. (2011); El Karoui et al. (2013). In the high-dimensional setting, where the dimensionality is greater than the number of samples, penalized M-estimators have been widely used to tackle the challenges of outliers and have been used for sparse recovery and variable selection, see Lambert-Lacroix and Zwald (2011); Li et al. (2011); Wang et al. (2013); Loh (2017). However, from the computational tractability perspective, it is often not easy to compute the M-estimators, since optimization problems over non-convex loss functions are usually involved. Moreover, the tractability issue may become more challenging when the data are contaminated by some arbitrary outliers, which is essentially the situation where robust M-estimator is designed to tackle.

This paper aims to investigate two important properties of M-estimators, *robustness* and *tractability*, simultaneously under *the gross error model*. Specifically, we assume the data generation model is $y_i = \langle \theta_0, x_i \rangle + \epsilon_i$, where $y_i \in \mathbb{R}, x_i \in \mathbb{R}^p$, , for $i = 1, \cdots, n$, and the noise term $\epsilon_i$'s are from Huber's gross error model (Huber, 1964): $\epsilon_i \sim (1 - \delta)f_0 + \delta g$, for $i = 1, \cdots, n$. Here, $f_0$ denotes the probability density function (pdf) of the noise of the normal samples, which has the desirable properties, such as zero mean and finite variance; $g$ denotes the pdf of the outliers (contaminations), which may also depend on the explanatory variable $x_i$, for $i = 1, \cdots, n$. One thing to notice is that we do not require the mean of $g$ to be 0. The parameter $\delta \in [0, 1]$, denotes the percentage of the contaminations, which is

also known as the contamination ratio in robust statistics literature. The gross error model indicates that for the $i^{th}$ sample, the residual term $\epsilon_i$ is generated from the pdf $f_0$ with probability $1 - \delta$, and from the pdf $g$ with probability $\delta$. It is important to point out that the residual $\epsilon_i$ is independent of $x_i$ and other $x_j$'s when it is from the pdf $f_0$, but can be dependent with the variable $x_i$ when it is from the pdf $g$.

In the first part of this paper, we start with the low-dimensional case when the dimension $p$ is fixed. We consider the robust M-estimation with a constraint on the $\ell_2$ norm of $\theta$. Mathematically, we study the following optimization problem:

$$\text{Minimize:} \quad \hat{R}_n(\theta) := \frac{1}{n} \sum_{i=1}^{n} \rho(y_i - \langle \theta, x_i \rangle), \tag{1}$$
$$\text{subject to:} \quad \|\theta\|_2 \leq r.$$

Here, $\rho : \mathbb{R} \to \mathbb{R}$ is the loss function, and is often *non-convex*. We consider the problem with the $\ell_2$ constraint due to three reasons: first, it is well know the constrainted optimization problem in (1) is equivalent to the unconstrained optimization problem with a $\ell_2$ regularizer. Therefore, it is related to the Ridge regression, which can alleviate multicollinearity amongst regression predictors. Second, by considering the problem of (1) in a compact ball with radius $r$, it guarantees the existence of the global optimal, which is necessary for establishing the tractability properties of the M-estimator. Finally, by working on the constrained optimization problem, we can avoid technical complications and establish the uniform convergence theorems of the empirical risk and population risk. Besides, the constrained M-estimators are widely used and studied in the literature, see Geyer et al. (1994); Mei et al. (2018); Loh (2017) for more details. To be consistent with the assumptions used in the literature, in the current work, we assume $r$ is a constant and the true parameter $\theta_0$ is inside of the ball.

In the second part, we extend our research to the high-dimensional case, where $p \gg n$

and the true parameter $\theta_0$ is sparse. In order to achieve the sparsity in the resulting estimator, we consider the penalized M-estimator with the $\ell_1$ regularizer:

$$\text{Minimize:} \quad \hat{L}_n(\theta) := \frac{1}{n} \sum_{i=1}^{n} \rho(y_i - \langle \theta, x_i \rangle) + \lambda_n ||\theta||_1, \tag{2}$$
$$\text{subject to:} \quad ||\theta||_2 \leq r.$$

Note the corresponding penalized M-estimator with the $\ell_2$ constraint is related to the Elastic net, which overcomes the limitations of the LASSO type regularization (Zou and Hastie, 2005).

In both parts, we will show that (in the finite sample setting,) the M-estimator obtained from (1) or (2) is robust in the sense that all stationary points of empirical risk function $\hat{R}_n(\theta)$ or $\hat{L}_n(\theta)$ are bounded in the neighborhood of the true parameter $\theta_0$ when the proportion of outliers is small. In addition, we will show that with a high probability, there is a unique stationary point of the empirical risk function, which is the global minimizer of (1) or (2) for some general (possibly nonconvex) loss functions $\rho$. This implies that the M-estimator can be computed efficiently. To illustrate our general theoretical results, we study some specific M-estimators with Huber's loss (Huber, 1964) and Welsch's exponential squared loss (Dennis Jr and Welsch, 1978), and explicitly discuss how the tuning parameter and percentage of outliers affect the robustness and tractability of the corresponding M-estimators.

Our research makes several fundamental contributions on the field of robust statistics and non-convex optimization. First, we demonstrate the uniform convergence results for the gradient and Hessian of the empirical risk to the population risk under the gross error model. Second, we provide nonasymptotic upper bound of the estimation error for the general M-estimators, which nearly achieve the minimax error bound in Chen et al. (2016). Third, we investigate the computational tractability of the general non-convex M-estimators

under the gross error model and show when the contamination ratio $\delta$ is small, there is only one unique stationary point of the empirical risk function. Therefore, efficient algorithms such as gradient descent or proximal gradient decent can be guaranteed to converge to a unique global minimizer irrespective of the initialization. Our general results also imply the following interesting and to some extent surprising statement: the percentage of outliers has an impact on the *tractability* of non-convex M-estimators. In a nutshell, the estimation and the corresponding optimization problem become more difficult both in terms of solution quality and computational efficiency when more outliers appear. While the former is well expected, we find the latter – that more outliers make M-estimators more difficult to numerically compute – an interesting and somewhat surprising discovery. Our simulation results and case study also verify this phenomenon.

**Related works**

Since Huber's pioneer work on robust M-estimators (Huber, 1964), many M-estimators with different choices of loss functions have been proposed, e.g., Huber's loss (Huber, 1964), Andrews sine loss (Andrews et al., 1972), Tukey's Bisquare loss (Beaton and Tukey, 1974), Welsch's exponential squared loss (Dennis Jr and Welsch, 1978), to name a few. From the statistical perspective, much research has been done to investigate the robustness of M-estimators such as large breakdown point (Donoho and Huber, 1983; Mizera and Müller, 1999; Alfons et al., 2013), finite influent function (Hampel et al., 2011) and asymptotic normality (Maronna and Yohai, 1981; Lehmann and Casella, 2006; El Karoui et al., 2013). Recently, in the high-dimensional context, regularized M-estimators have received a lot of attentions. Lambert-Lacroix and Zwald (2011) proposed a robust variable selection method by combing Huber's loss and adaptive lasso penalty. Li et al. (2011) show the nonconcave penalized M-estimation method can perform parameter estimation and variable selection

simultaneously. Welsch's exponential squared loss combined with adaptive lasso penalty is used by Wang et al. (2013) to construct a robust estimator for sparse estimation and variable selection. Chang et al. (2018) proposed a robust estimator by combining the Tukey's biweight loss with adaptive lasso penalty. Loh and Wainwright (2015) proved that under mild conditions, any stationary point of the non-convex objective function will close to the underlying true parameters. However, those statistical works did not discuss the computational tractability of the M-estimators even though many of these loss functions are non-convex.

During the last several years, non-convex optimization has attracted fast growing interests due to its ubiquitous applications in machine learning and in particular deep learning, such as dictionary learning (Mairal et al., 2009), phase retrieval (Candes et al., 2015), orthogonal tensor decomposition (Anandkumar et al., 2014) and training deep neural networks (Bengio, 2009). It is well known that there is no efficient algorithm that can guarantee to find the global optimal solution for general non-convex optimization.

Fortunately, in the context of estimating non-convex M-estimators for high-dimensional linear regression (*without outliers*), under some mild statistical assumptions, Loh (2017) establishes the uniqueness of the stationary point of the non-convex M-estimator when using some non-convex bounded regularizers instead of $\ell_1$ regularizer. By investigating the uniform convergence of gradient and Hessian of the empirical risk, Mei et al. (2018) prove that with a high probability, there exists one unique stationary point of the regularized empirical risk function with $\ell_1$ regularizer. Thus regardless of the initial points, many computational efficient algorithm such as gradient descent or proximal gradient descent algorithm could be applied and are guaranteed to converge to the global optimizer, which implies the high tractability of the M-estimator. However, their analysis is restricted to the standard linear regression setting without outliers. In particular, they assume the

distribution of the noise terms in the linear regression model should have some desirable properties such as zero mean, sub-gaussian and independent of feature vector $x$, which might not hold when the data are contaminated with outliers. To the best of our knowledge, no research has been done on analyzing the computational tractability properties of the non-convex M-estimators when data are contaminated by arbitrary outliers, although the very reason why M-estimators are proposed is to handle outliers in linear regression in the robust statistics literature. Our research is the first to fill the significant gap on the tractability of non-convex M-estimators. We prove that under mild assumptions, many M-estimators can tolerate a small amount of arbitrary outliers in the sense of keeping the tractability, even if the loss functions are non-convex.

**Notations.** Given $\mu, \nu \in \mathbb{R}^p$, their standard inner product is defined by $\langle \mu, \nu \rangle = \sum_{i=1}^{p} \mu_i \nu_i$. The $\ell_p$ norm of a vector $x$ is denoted by $||x||_p$. The $p$ by $p$ identity matrix is denoted by $I_{p \times p}$. Given a matrix $M \in \mathbb{R}^{m \times m}$, let $\lambda_{\max}(M), \lambda_{\min}(M)$ denote the largest and the smallest eigenvalue of $M$, respectively. The operator norm of $M$ is denoted by $||M||_{op}$, which is equal to $\max(\lambda_{\max}(M), -\lambda_{\min}(M))$ when $M \in \mathbb{R}^{m \times m}$. Let $B_q^p(a, r) = \{x \in \mathbb{R}^p : ||x - a||_q \leq r\}$, be the $\ell_q$ ball in the $\mathbb{R}^p$ space with center $a$ and radius $r$. Given a random variable $X$ with probability density function $f$, we denote the corresponding expectation by $\mathbf{E}_f$. We will often omit the density function subscript $f$ when it is clear from the context, the expectation is taken for all variables.

**Organization.** The rest of this article is organized as follows. In Section 2, we present the theorems about the robustness and tractability of general M-estimators under the low-dimensional setup when dimension $p$ is fixed and less than $n$. Then in Section 3, we consider the penalized M-estimator with $\ell_1$ regularizer in the high-dimensional regression when $p \gg n$. The $\ell_2$ error bounds of the estimation and the scenario when the M-estimator has nice tractability are provided. In Section 4, we discuss two special families of robust

estimator constructed by Huber's and Welsch's exponential loss as examples to illustrate our general theorems of robustness and tractability of M-estimators. Simulation results are presented in Section 5 and a case study is shown in Section 6 to illustrate the robustness and tractability properties when the data are contaminated by outliers. Concluding remarks are given in Section 7. We relegate all proofs to the Appendix due to space limits.

## 2   M-estimators in the low-dimensional regime

In this section, we investigate two key properties of M-estimators, namely *robustness* and *tractability*, in the setting of linear regression with arbitrary outliers in the low-dimensional regime where the dimension $p$ is fixed and smaller than the number of samples $n$. In terms of robustness, we show that under some mild conditions, any stationary point of the objective function in (1) will be well bounded in a neighborhood of the true parameter $\theta_0$. Moreover, the neighborhood shrinks when the proportion of outliers decreases. In terms of tractability, we show that when the proportion of outliers is small and the sample size is large, with a high probability, there is a *unique stationary point* of the empirical risk function, which is the global optimum (and hence the corresponding M-estimator). Consequently, many first order methods are guaranteed to converge to the global optimum, irrespective of initialization.

Before presenting our main theorems, we make the following mild assumptions on the loss function $\rho$, the explanatory or feature vectors $x_i$, and the idealized noise distribution $f_0$. We define the score function $\psi(z) := \rho'(z)$.

**Assumption 1. (a)** *The score function $\psi(z)$ is twice differentiable and odd in $z$ with $\psi(z) \geq 0$ for all $z \geq 0$. Moreover, we assume $\max\{||\psi(z)||_\infty, ||\psi'(z)||_\infty, ||\psi''(z)||_\infty\} \leq L_\psi$.*

**(b)** *The feature vector $x_i$ are i.i.d with zero mean and $\tau^2$-sub-Gaussain, that is $\mathbf{E}[e^{\langle \lambda, x_i \rangle}] \leq \exp(\frac{1}{2}\tau^2 ||\lambda||_2^2)$, for all $\lambda \in \mathbb{R}^p$.*

**(c)** *The feature vector $x_i$ spans all possible directions in $\mathbb{R}^p$, that is $\mathbf{E}[x_i x_i^T] \succeq \gamma \tau^2 I_{p \times p}$, for some $0 < \gamma \leq 1$.*

**(d)** *The idealized noise distribution $f_0(\epsilon)$ is symmetric. Define $h(z) := \int_{-\infty}^{\infty} f_0(\epsilon)\psi(z+\epsilon)d\epsilon$ and $h(z)$ satisfies $h(z) > 0$, for all $z > 0$ and $h'(0) > 0$.*

Assumption (a) requires the smoothness of the loss function in the objective function, which is crucial to study the tractability of the estimation problem; Assumption (b) assumes the sub-Gaussian design of the observed feature matrix; Assumption (c) assumes that the covariance matrix of the feature vector is positive semidefinite. We remark that the condition on $h(z)$ is mild. It is not difficult to show that it is satisfied if the idealized noise distribution $f_0(\epsilon)$ is strictly positive for all $\epsilon$ and decreasing for $\epsilon > 0$, e.g., if $f_0 =$ pdf of $N(0, \sigma^2)$.

Before presenting our main results in this section, we first define the population risk as follows:

$$R(\theta) = \mathbf{E}\hat{R}_n(\theta) = \mathbf{E}[\rho(Y - \langle \theta, X \rangle)]. \tag{3}$$

The high level idea is to analyze the population risk first, and then we build a link between the population risk and the empirical risk, which solves the original estimation problem. Theorem 1 below summarizes the results for the population risk function $R(\theta)$ in (3).

**Theorem 1.** *Assume that Assumption 1 holds and the true parameter $\theta_0$ satisfies $||\theta_0||_2 \leq r/3$.*

9

**(a)** *There exists a constant $\eta_0 = \frac{\delta}{1-\delta}C_1$ such that any stationary point $\theta^*$ of $R(\theta)$ satisfies $||\theta^* - \theta_0||_2 \leq \eta_0$, where $\delta$ is the contamination ratio, and $C_1$ is a positive constant that only depends on $\gamma, r, \tau, \psi(z)$ and the pdf $f_0$, but does not depend on the outlier pdf $g$.*

**(b)** *When $\delta$ is small, there exist a constant $\eta_1 = C_2 - C_3\delta > 0$, where $C_2, C_3$ are two positive constants that only depend on $\gamma, r, \tau, \psi(z)$ and the pdf $f_0$ but not depend on the outlier pdf $g$, such that*

$$\lambda_{\min}(\nabla^2 R(\theta)) > 0 \tag{4}$$

*for every $\theta$ with $||\theta_0 - \theta||_2 < \eta_1$.*

**(c)** *There is a unique stationary point of $R(\theta)$ in the ball $B_2^p(0, r)$ as long as $\eta_0 < \eta_1$ for a given contamination ratio $\delta$.*

It is useful to add some remarks for better understanding Theorem 1. First, recall that the noise term $\epsilon_i$ follows the gross error model: $\epsilon_i \sim (1 - \delta)f_0 + \delta g$, where the outlier pdf $g$ may also depend on $x_i$. While the true parameter $\theta_0$ may no longer be the stationary point of the population risk function $R(\theta)$, Theorem 1 implies that the stationary points of $R(\theta)$ will always bounded in a neighborhood of the true parameter $\theta_0$ when the percentage of contamination $\delta$ is small. This indicates the robustness of M-estimators in the population case.

Second, Theorem 1 asserts that when there are no outliers, i.e., $\delta = 0$, the stationary point is indeed the true parameter $\theta_0$. In addition, since the constant $\eta_0$ in (a) is an increasing function of $\delta$ whereas the constant $\eta_1$ in (b) is a decreasing function of $\delta$, stationary points of $R(\theta)$ may disperse from the true parameter $\theta_0$ and the strongly convex region

10

around $\theta_0$ will be decreasing, as the contamination ratio $\delta$ is increasing. This indicates the difficulty of optimization for large contamination ratio cases.

Third, part (c) is a direct result from part (a) and (b). Note that $\eta_0(\delta = 0) = 0 < \eta_1(\delta = 0) = C_2$, thus there exists a positive $\delta^*$, such that $\eta_0 < \eta_1$ for any $\delta < \delta^*$. A simple lower bound on $\delta^*$ is $C_3/(C_1 + C_2 + C_3)$, since $C_1\delta < (1 - \delta)(C_2 - C_3\delta)$ whenever $0 \leq \delta \leq C_3/(C_1 + C_2 + C_3)$.

Our next step is to link the empirical risk function (and the corresponding M-estimator) with the population version. To this end, we need the following lemma, which shows the global uniform convergence theorem of the sample gradient and Hessian.

**Lemma 1.** *Under Assumption 1, for any $\pi > 0$, there exists a constant $C_\pi$ depending on $\pi, \gamma, r, \tau, \psi(z), h(z)$ but independent of $p, n, \delta$ and $g$, such that for any $\delta \geq 0$, the following hold:*

**(a)** *The sample gradient converges uniformly to the population gradient in Euclidean norm, i.e., if $n \geq C_\pi p \log n$, we have*

$$\mathbf{P}\left(\sup_{\theta \in B_2^p(0,r)} ||\nabla \hat{R}_n(\theta) - \nabla R(\theta)||_2 \leq \tau\sqrt{\frac{C_\pi p \log n}{n}}\right) \geq 1 - \pi. \tag{5}$$

**(b)** *The sample Hessian converges uniformly to the population Hessian in operator norm, i.e., if $n \geq C_\pi p \log n$, we have*

$$\mathbf{P}\left(\sup_{\theta \in B_2^p(0,r)} ||\nabla^2 \hat{R}_n(\theta) - \nabla^2 R(\theta)||_{op} \leq \tau^2\sqrt{\frac{C_\pi p \log n}{n}}\right) \geq 1 - \pi. \tag{6}$$

We are now ready to present our main result about M-estimators by investigating the empirical risk function $\hat{R}_n(\theta)$.

11

**Theorem 2.** *Assume Assumption 1 holds and $||\theta_0||_2 \leq r/3$. Let us use the same notation $\eta_0$ and $\eta_1$ as in Theorem 1. Then for any $\pi > 0$, there exist constant $C_\pi$ depends on $\pi, \gamma, r, \tau, \psi, f_0$ but independent of $n, p, \delta$ and $g$, such that as $n \geq C_\pi p \log n$, the following statements hold with probability at least $1 - \pi$ :*

**(a)** *for all $||\theta - \theta_0||_2 > 2\eta_0$,*

$$\langle \theta - \theta_0, \nabla \widehat{R}_n(\theta) \rangle > 0. \tag{7}$$

**(b)** *for all $||\theta - \theta_0||_2 \leq \eta_1$,*

$$\lambda_{\min}(\nabla^2 \widehat{R}_n(\theta)) > 0. \tag{8}$$

*Thus, as long as $2\eta_0 < \eta_1$, $\widehat{R}_n(\theta)$ has a unique stationary point, which lies in the ball $B^p(0, r)$. This is the unique global optimal solution of (1), and denote this unique stationary point by $\widehat{\theta}_n$.*

**(c)** *There exists a positive constant $\kappa$ that depends on $\pi, \gamma, r, \psi, \delta, f_0$ but independent of $n, p$ and $g$, such that*

$$||\widehat{\theta}_n - \theta_0||_2 \leq \eta_0 + \frac{4\tau}{\kappa} \sqrt{\frac{C_\pi p \log n}{n}}. \tag{9}$$

A few remarks are in order. First, since $\eta_0$ is independent of $n, p$ and $g$, Theorem 2(a) asserts that the M-estimator which minimizes $\widehat{R}_n(\theta)$ is always bounded in the ball $B_2^p(\theta_0, 2\eta_0)$, regardless of $g$ (and hence the outliers observed). This indicates the robustness of the M-estimator, i.e., the estimates are not severely skewed by a small amount of "bad" outliers. Next, when the contamination ratio $\delta$ is small such that $2\eta_0 < \eta_1$, there is a unique stationary point of $\widehat{R}_n(\theta)$. Therefore, although the original optimization problem

12

(1) is non-convex and the sample contains some arbitrary outliers, the optimal solution of $\widehat{R}_n(\theta)$ can be computed efficiently via most off-the-shelf first-order algorithms such as gradient descent or stochastic gradient descent. This indicates the tractability of the M-estimator. Interestingly, as in the population risk case, the tractability is closely related to the amount of outliers – the problem is easier to optimize when the data contains fewer outliers. Finally, when the number of samples $n \gg p \log n$, the estimation error bound $\eta_0$ is as the order of $O(\delta + \sqrt{\frac{p \log n}{n}})$, which nearly achieves the minimax lower bound of $O(\delta + \sqrt{\frac{p}{n}})$ in Chen et al. (2016).

# 3    Penalized M-estimator in the high-dimensional regime

In this section, we investigate the tractability and the robustness of the penalized M-estimator in the high-dimension region where the dimension of parameter $p$ is much greater than the number of samples $n$. Specifically, we consider the same data generation model $y_i = \langle \theta_0, x_i \rangle + \epsilon_i$, where $y_i \in \mathbb{R}, x_i \in \mathbb{R}^p$, and the noise term $\epsilon_i$ are from Huber's gross error model (Huber, 1964): $\epsilon_i \sim (1 - \delta)f_0 + \delta g$. Moreover, we assume $p \gg n$ and the true parameter $\theta_0$ is sparse.

We consider the $\ell_1$-regularized M-estimation under a $\ell_2$-constraint on $\theta$:

$$\text{Minimize:} \quad \hat{L}_n(\theta) := \frac{1}{n} \sum_{i=1}^{n} \rho(y_i - \langle \theta, x_i \rangle) + \lambda_n ||\theta||_1, \tag{10}$$
$$\text{subject to:} \quad ||\theta||_2 \leq r.$$

Before presenting our main theorem, we need additional assumptions on the feature vector $x$.

**Assumption 2.** *The feature vector $x$ has a probability density function in $\mathbb{R}^p$. In addition,*

13

there exists constant $M > 1$ that is independent of dimension $p$ such that $||x||_\infty \leq M\tau$ almost sure.

The following lemma shows the uniform convergence of gradient and Hessian under the Huber's contamination model in the high-dimensional setting where $p >> n$.

**Lemma 2.** *Under assumption 1 and 2, there exist constants $C_1, C_2, T_0, L_0$ that depend on $r, \tau, \pi, \delta, L_\psi$, but independent of $n, p$, and $g$, such that the following hold:*

**a** *The sample directional gradient converges uniformly to the population directional gradient, along the direction $(\theta - \theta_0)$.*

$$\mathbf{P}\left(\sup_{\theta \in B_2^p(r)\backslash\{0\}} \frac{|\langle \nabla R_n(\theta) - \nabla R(\theta), \theta - \theta_0\rangle|}{||\theta - \theta_0||_1} \leq (T_0 + L_0\tau)\sqrt{\frac{C_1 \log(np)}{n}}\right) \geq 1 - \pi. \ (11)$$

**b** *As $n \geq C_2 s_0 \log(np)$, we have*

$$\mathbf{P}\left(\sup_{\theta \in B_2^p(r) \cap B_2^p(s_0), \nu \in B_2^p(1) \cap B_0^p(s_0)} |\langle \nu, \left(\nabla^2 R_n(\theta) - \nabla^2 R(\theta)\right)\nu\rangle| \leq \tau^2\sqrt{\frac{C_2 s_0 \log(np)}{n}}\right) \geq 1 - \pi.$$

Now we are ready for our main theorem.

**Theorem 3.** *Assume that Assumption 1 and Assumption 2 hold and the true parameter $\theta_0$ satisfies $||\theta_0||_2 \leq r/3$ and $||\theta_0||_0 \leq s_0$. Then there exist constants $C, C_0, C_1, C_2$ that are dependent on $(\rho, L_\psi, \tau^2, r, \gamma, \pi)$ but independent on $(\delta, s_0, n, p, M)$ such that as $n \geq Cs_0 \log p$ and $\lambda_n = C_0 M \sqrt{\frac{\log p}{n}} + \frac{C_1}{\sqrt{s_0}}\delta$, the following hold with probability as least $1 - \pi$ :*

**(a)** *All stationary points of problem (10) are in $B_2^p(\theta_0, \eta_0 + \frac{\sqrt{s_0}}{1-\delta}\lambda_n C_2)$*

**(b)** *As long as $n$ is large enough such that $n \geq Cs_0 \log^2 p$ and the contamination ratio $\delta$ is small such that $(\eta_0 + \frac{1}{1-\delta}\sqrt{s_0}\lambda_n C_2) \leq \eta_1$, the problem (10) has a unique local stationary point which is also the global minimizer.*

14

The proof of Theorem 3 is based on several lemmas, which are postponed to the appendix. We believe that some of our lemmas are of interest in their own right. Theorem 3 implies the estimation error of the penalized M-estimator is bounded as the order of $O(\delta + \sqrt{\frac{s_0 \log p}{n}})$, which achieves the minimax estimation rate (Chen et al., 2016). Moreover, it implies that the penalized M-estimator has good tractability when the percentage of outliers $\delta$ is small.

# 4 Example

In this section, we use some examples to illustrate our general theoretical results about the robustness and tractability of M-estimators. In the first subsection, we consider the low-dimensional regime and study a family of M-estimators with a specific loss function known as Huber's loss (Huber, 1964). In the second subsection, we consider the high-dimensional regime and study the penalized M-estimator with Welsch's exponential squared loss (Dennis Jr and Welsch, 1978; Rey, 2012; Wang et al., 2013). In both subsections, we will derive the explicit expression of the two critical radius $\eta_0$, $\eta_1$ and discuss the robustness and tractability of the corresponding M-estimators.

## 4.1 M-estimator via Huber's loss

In this subsection, we illustrate the general results presented in Section 2 by studying the Huber's loss function (Huber, 1964)

$$\rho_\alpha(t) = \begin{cases} \frac{1}{2}t^2, & \text{if } |t| \leq \alpha \\ \alpha(|t| - \alpha/2), & \text{if } |t| > \alpha. \end{cases} \tag{12}$$

15

where $\alpha > 0$ is a tuning parameter. The corresponding M-estimator is obtained by solving the optimization problem

$$\min_{\theta} \quad \hat{R}_n(\theta) := \frac{1}{n} \sum_{i=1}^{n} \rho_\alpha(y_i - \langle \theta, x_i \rangle), \tag{13}$$

$$\text{subject to} \quad ||\theta||_2 \leq r.$$

First, note the loss function $\rho_\alpha(t)$ in (12) is convex. Thus, the corresponding M-estimator should be tractable even though there are some outliers. Second, when $\alpha$ goes to 0, $\rho_\alpha(t)$ will converges to $t^2/2$. Thus, the least square estimator is a special case of the M-estimator obtained from (13), which is not robust to outliers. Third, for fixed $\alpha > 0$, $\rho'_\alpha(t), \rho''_\alpha(t)$ are all bounded. Intuitively, this implies that the impact of outlier observations of $y_i$ will be controlled and thus the corresponding statistical procedure will be robust.

We now study the robustness and tractability of the M-estimator of (13) based on our framework in Theorem 2. In order to emphasize on the effects of the tuning parameter $\alpha$ and the contamination ratio $\delta$ on the robustness property and tractability property, we consider a simplified assumption on the feature vector $x_i$ and the pdf of idealized residual $f_0$.

**Assumption 3. (a)** *The feature vector $x_i$ are i.i.d multivariate Gaussian distribution $N(0, \tau^2 I_{p \times p})$.*

**(b)** *The idealized noise pdf $f_0(\epsilon)$ has Gaussian distribution $N(0, \sigma^2)$.*

**(c)** *Assume the true parameter $||\theta_0||_2 \leq r/3$.*

**Corollary 1.** *Under Assumption 3, for any $\delta, \alpha \geq 0$, there exist two constants $\eta_0(\delta, \alpha), \eta_1(\delta, \alpha)$ :*

$$\eta_0(\delta, \alpha) = \frac{\delta}{1 - \delta} \frac{4\sqrt{2\pi}\sigma^3}{(\alpha^2 + 3\sigma^2)\tau} e^{\frac{\alpha^2 + 22\tau^2 r^2}{2\sigma^2}} \tag{14}$$

$$\eta_1(\delta, \alpha) = +\infty, \tag{15}$$

16

*such that when the number of data points n is large, with high probability, any stationary*
*points of the empirical risk function $\hat{R}_n(\theta)$ in (13) belongs in the ball $B_2^p(\theta_0, 2\eta_0(\delta, \alpha))$. More-*
*over, the empirical risk function $\hat{R}_n(\theta)$ in (13) is strongly convex in the ball $B_2^p(\theta_0, \eta_1(\delta, \alpha))$.*
*Thus, there exists a unique stationary point of $\hat{R}_n(\theta)$, which is the corresponding M-estimator.*

Note $\eta_1(\delta, \alpha) = \infty$, which means the corresponding Huber's estimator will be tractable, no matter there are outliers or not. This is consistent with the fact that the Huber's loss function is convex. Moreover, it is interesting to see the special case of Corollary 1 with $\alpha = +\infty$, which reduces to the least square estimator. As we can see, with $\delta > 0$, we have $\eta_0(\delta, \alpha = +\infty) = +\infty$, which implies the solution of the optimization problem in (13) can be arbitrarily in the ball $B_2^p(0, r = 10)$, even when the proportion of outliers is small. Thus it is not robust to the outliers. This recovers the well-known fact: the least square estimator is easy to compute, but is very sensitive to outliers.

Additionally, for another special case with $\delta = 0$ and $\alpha > 0$, we have $\eta_0(\delta = 0, \alpha) = 0$, which means the true parameter $\theta_0$ is the unique stationary point of the risk function. This implies the Huber's estimator is consistent when there are no outliers.

## 4.2  Penalized M-estimator via Welsch's exponential squared loss

In this subsection, we illustrate the general results presented in Section 3 by considering a family of M-estimators with a specific nonconvex loss function known as Welsch's exponential squared loss (Dennis Jr and Welsch, 1978; Rey, 2012; Wang et al., 2013),

$$\rho_\alpha(t) = \frac{1 - \exp(-\alpha t^2/2)}{\alpha}, \tag{16}$$

17

where $\alpha \geq 0$ is a tuning parameter. The corresponding penalized M-estimator is obtained by solving the optimization problem

$$\min_{\theta} \quad \hat{L}_n(\theta) := \frac{1}{n} \sum_{i=1}^{n} \rho_\alpha(y_i - \langle \theta, x_i \rangle) + \lambda_n ||\theta||_1, \tag{17}$$

$$\text{subject to} \quad ||\theta||_2 \leq r.$$

The non-convex loss function $\rho_\alpha(t)$ in (16) has been used in other contexts such as robust estimation and robust hypothesis testing, see Ferrari and Yang (2010); Qin and Priebe (2017), as it has many nice properties. First, it is a smooth function of both $\alpha$ and $t$, and the gradient and Hessian are well-defined. Second, when $\alpha$ goes to 0, $\rho_\alpha(t)$ will converges to $t^2/2$. Thus, the LASSO estimator is a special case of the M-estimator obtained from (17). Third, for fixed $\alpha > 0$, $\rho_\alpha(t), \rho'_\alpha(t), \rho''_\alpha(t)$ are all bounded. Intuitively, this implies that the impact of outlier observations of $y_i$ will be controlled and thus the corresponding statistical procedure will be robust.

We now study the robustness and tractability of the penalized M-estimator of (17) based on our framework in Theorem 3. When $\alpha$ goes to 0, the M-estimator reduces to the LASSO estimator, which can be computed easily. However, it is also known to be very sensitive to the outliers (Alfons et al., 2013). On the other hand, when $\alpha$ increases, the estimator becomes more robust, but may lose tractability due to the highly non-convexity of the function $\rho_\alpha(t)$ as well as the presence of the outliers.

In order to emphasize on the relation between the tuning parameter $\alpha$ and the contamination ratio $\delta$, we consider a simplified assumption on the feature vector $x_i$ and the pdf of idealized residual $f_0$.

**Assumption 4. (a)** *The feature vector $x_i$ are i.i.d multivariate uniform distribution $[-\tau, \tau]^p$.*

**(b)** *The idealized noise pdf $f_0(\epsilon)$ has Gaussian distribution $N(0, \sigma^2)$.*

**(c)** *The true parameter $||\theta_0||_2 \le r/3$.*

With Assumption 4 and Theorem 3, we can get the following corollary, which characterizes the robustness and tractability of the penalized M-estimator with Welsch's exponential squared loss in (17):

**Corollary 2.** *Assume that Assumption 4 holds and the true parameter $\theta_0$ satisfies $||\theta_0||_2 \le r/3$, for any $\pi \in (0,1)$, there exist a constant $C_\pi$ such that if choose $\lambda_n = 2C_\pi \tau \sqrt{\frac{\log p}{n}} + \frac{\alpha \tau}{2} \frac{\delta}{\sqrt{s_0}}$, as $n \gg s_0 \log p$, the following hold with probability as least $1 - \pi$ :*

**(a)** *All stationary points of problem (17) are in $B_2^p(\theta_0, (1+2\tau)\eta_0)$*

**(b)** *The empirical risk function $\hat{L}_n(\theta)$ are strong convex in the ball $B_2^p(\theta_0, \eta_1)$*

**(c)** *As long as $n$ is large enough and the contamination ratio $\delta$ is small such that $(1+2\tau)\,\eta_0 \le \eta_1$, the problem (17) has a unique local stationary point which is also the global minimizer.*

*Here*

$$\eta_0(\delta, \alpha) \;=\; \frac{\delta}{1-\delta}\sqrt{\frac{e}{\alpha}} \frac{4(1+\alpha\sigma^2)^{3/2}}{\tau} e^{\frac{32\alpha r^2 \tau^2}{3(1+\alpha\sigma^2)}} \tag{18}$$

$$\eta_1(\delta, \alpha) \;=\; \frac{1}{3\sqrt{3\alpha}(1+\alpha\sigma^2)^{3/2}\tau}\left[\tau^2 - \delta(\tau^2 + (1+\alpha\sigma^2)^{3/2})\right]. \tag{19}$$

It is interesting to see the special case of Corollary 2 with $\alpha = 0$, which reduces to the LASSO estimator. On the one hand, with $\alpha = 0$, we have $\eta_1(\delta, \alpha = 0) = +\infty$ for any $\delta > 0$. This means that the corresponding risk function is strongly convex in the entire region of $B_2^p(0, r = 10)$, and hence it is always tractable. On the other hand, since $\eta_0(\delta, \alpha = 0) = +\infty$, the solution of the optimization problem in (17) can be arbitrarily in the ball $B_2^p(0, r = 10)$, even when the proportion of outliers is small. Thus it is not

19

robust to the outliers. This recovers the well-known fact: the LASSO estimator is easy to compute, but is very sensitive to outliers.

Additionally, for another special case with $\delta = 0$ and $\alpha > 0$, we have $\eta_0(\delta = 0, \alpha) = 0$, which means the true parameter $\theta_0$ is the unique stationary point of the risk function. This implies the Welsch's estimator has nice tractability when there is no outliers. However, when the percentage of outlier $\delta$ is increasing, $\eta_1(\delta, \alpha)$ will decrease, which implies more outliers will reduce the tractability of the M-estimator.

# 5   Simulation results

In this section, we report the simulation results by using Welsch's exponential loss (Dennis Jr and Welsch, 1978) when the data are contaminated, using synthetic data setting. We first generate covariates $x_i \sim N(0, I_{p \times p})$ and responses $y_i = \langle \theta_0, x_i \rangle + \epsilon_i$, where $||\theta_0||_2 = 1$. We consider the case when the residual term $\epsilon_i$ have gross error model with contamination ratio $\delta$, i.e., $\epsilon_i \sim (1 - \delta)N(0, 1) + \delta N(\mu_i, 3^2)$ where $\mu_i = ||x_i||_2^2 + 1$. The outlier distribution is chosen to highlight the effects of outliers when they are dependent on $x_i$ and has non-zero mean.

In the first part, we consider the low-dimensional case when the dimension $p = 10$. Specifically, we generate $n = 200$ pairs of data $(y_i, x_i)_{i=1,..,n}$ with dimension $p = 10$ and with different choices of contamination ratios $\delta$. We use projected gradient descent to solve the optimization problem in (13) with $r = 10$. In order to make the iteration points be inside the ball, we will project the points back into $B_2^p(0, r = 10)$ if they fall out of the ball. The step size is fixed as 1. In order to test the tractability of the M-estimator, we run gradient descent algorithm with 20 random initial values in the ball $B_2^p(0, r = 10)$ to see whether the gradient descent algorithm can converge to the same stationary point or not. Denote $\hat{\theta}(k)$ as the $k^{th}$ iteration points, Figure 1 shows the convergence of the gradient descent

algorithm for the exponential loss with the choice of $\alpha = 0.1$ under the gross error model with different $\delta$. From Figure 1 we observe when the proportion of outliers is small (i.e., $\delta \leq 0.1$,) gradient descent could converge to the same stationary point fast. However, when the contamination ratio $\delta$ becomes larger, gradient descent may not converge to the same point for different initial points, indicating the loss of tractability *for the same objective function* with increasing proportion of outliers. Those observations are consistent to our Theorem 2, which asserts the M-estimator is tractable when the contamination ratio $\delta$ is small.

To illustrate the robustness of the M-estimator, we generate 100 realizations of $(Y, X)$ and run gradient descent algorithm with different initial values. The average estimation errors between the M-estimator and the true parameter $\theta_0$ are presented in Figure 2. As we can see, when $\delta = 0$, all estimators have small estimation errors, which are well expected as those M-estimators are consistent without outliers (Huber, 1964; Huber and Ronchetti, 2009). However, for the M-estimator with $\alpha = 0$, i.e., the least square estimator, the estimation error will increase dramatically as the proportion of outliers increases. This confirms that the least square estimator is not robust to the outliers.

Meanwhile, when $\alpha = 0.1$, the overall estimation error does not increase much even with 40% outliers, which clearly demonstrate the robustness of the M-estimator. Note that when $\alpha$ is further increased from 0.1 to 0.3, although the estimator error is still very small for $\delta \leq 0.2$, it will increase dramatically when $\delta$ is greater than 0.2. We believe that two reasons contribute to this phenomenon: robustness starts to decrease when $\alpha$ becomes too large; and more importantly, the algorithm fails to find the global optimum due to multiple stationary points when $\alpha$ is large. Thus for each $\alpha$, there exists a critical bound of $\delta$, such that the estimator will be robust and tractable efficiently when the proportion of outliers is smaller than that bound.

21

In the second part, we present our results in the high-dimensional region when $p = 400$. Data $(y_i, x_i)$ are generated from the same gross error model in the previous simulation study, with the true parameter $\theta_0$ a sparse vector with 10 nonzero entries. All nonzero entries are set to be $1/\sqrt{10}$. We use proximal gradient descent algorithm to solve problem (10). Similarly, we will project the points back into $B_2^p(0, r = 10)$ if they fall out of the ball. Figure 3 shows the convergence of the proximal gradient descent algorithm for the nonconvex exponential loss with the choice of $\alpha = 0.1$ and $L_1$ regularizer with the parameter $\lambda = 0.1$ under the gross error model with different $\delta$. From Figure 3 we observe when the percentage of outliers is small, the algorithm will converge to the same stationary point fast, which implies there is only one unique stationary point. When $\delta$ is larger, the converge rate become slower, which implies there may exist another stationary points. Those simulation results reflect our theoretical result for the tractability of the penalized M-estimator in high-dimensional regression.
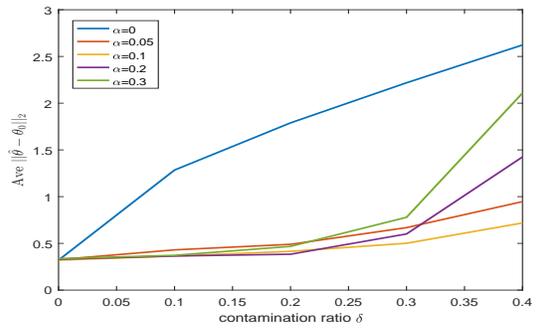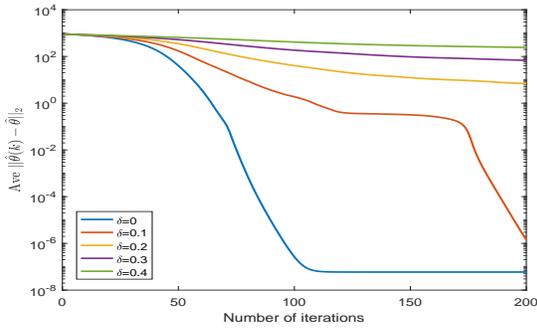


Figure 1: The convergence of gradient descent algorithm for different $\delta$. Y-axis is with log scale.

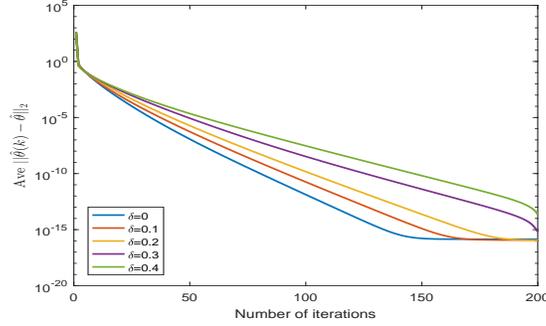Figure 2: The estimation error for different $\alpha$ and $\delta$

Figure 3: The convergence of gradient descent algorithm for different $\delta$. Y-axis is with log scale.

# 6 Case study

In this section, we present a case study of the robust regression problem for the Airfoil Self-Noise dataset (Brooks et al., 2014), which is available on UCI Machine Learning Repository. The dataset was processed by NASA and is commonly used for regression study to learn the relation between the airfoil self-noise and five explanatory variables. Specifically, the dataset contain the following 5 explanatory variables: Frequency (in Hertzs), Angle of attack (in degrees), Chord length,(in meters), Free-stream velocity (in meters per second), and Suction side displacement thickness (in meters). There are 1503 observations in the dataset. The response variable is Scaled sound pressure level (in decibels). In this section, the five explanatory variables are scaled to have zero mean and unit variance. Then, we corrupt the response by adding noise $\epsilon$ from the same gross error model as the previous section: $\epsilon_i \sim (1-\delta)N(0,1) + \delta N(\mu_i, 3^2)$ with $\mu_i = ||x_i||_2^2 + 1$.

We consider the M-estimator using Welsch's exponential loss (Dennis Jr and Welsch, 1978) on the dataset to validate the tractability and the robustness of the corresponding M-estimator. First, we run 100 Monte Carlo simulations. At each time, we split the dataset

23

which consists of 1503 pairs of data into a training dataset of size 1000 and a testing dataset of size 503. Then for the training dataset, we use gradient descent method with 20 different initial values to update the iteration points.

Figure 4 shows the average distance between each iteration point and the optimal point with the choice of $\alpha = 0.7$ and step size 0.5. Clearly, when $\delta$ is smaller than 0.3, gradient descent will converge to the same local minimizer, which implies the uniqueness of the stationary point. This result demonstrates the nice tractability of the M-estimator under the gross error model when the proportion of outliers is small. Then, using the optimal point as the M-estimator, we calculate the prediction error, which is the mean square error on the testing data. Figure 5 shows the average prediction error on the testing data. As we can see, the prediction error with the choice of $\alpha = 0$ will increase dramatically when the percentage of outliers increases. In contrast, the prediction errors of M-estimators with $\alpha = 0.4$ is stable even with a large percentage of outliers. This illustrates the robustness of M-estimators for some positive $\alpha$.
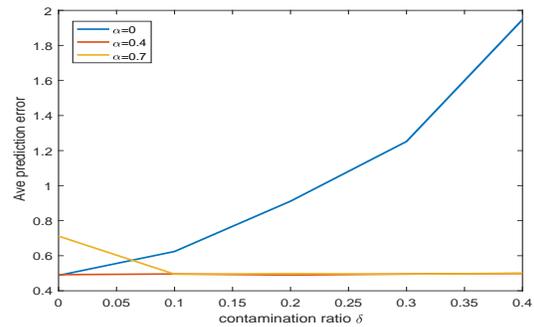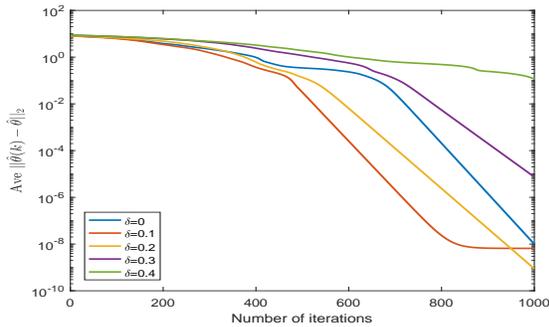


Figure 4: The convergence of gradient descent algorithm for different $\delta$. Y-axis is with log scale.

Figure 5: The prediction error for different $\alpha$ and $\delta$

# 7    Conclusions

In this paper, we investigate the robustness and computational tractability of general (non-convex) M-estimators in both classical low-dimensional regime and modern high-dimensional regime. In terms of *robustness*, in the low-dimensional regime, we show the estimation error of the M-estimator is as the order of $O(\delta + \sqrt{\frac{p \log n}{n}})$, which nearly achieves the minimax lower bound of $O(\delta + \sqrt{\frac{p}{n}})$ in Chen et al. (2016). In the high-dimensional regime, we show the estimation error of the penalized M-estimator has the estimation error as the order of $O(\delta + \sqrt{\frac{s_0 \log p}{n}})$, which achieves the minimax estimation rate (Chen et al., 2016).

In terms of *tractability*, our theoretical results imply under sufficient conditions, when the percentage of arbitrary outliers is small, the general M-estimator could have good computational tractability since it has only one unique stationary point, even if the loss function is non-convex. Therefore, M-estimators can tolerate certain level of outliers by keeping both estimation accuracy and computation efficiency. Both simulation and real data case study are conducted to validate our theoretical results about the robustness and tractability of M-estimation in the presence of outliers.

# 8  Appendix

**Proof of Lemma 1**: In order to prove the uniform convergency theorem, it is suffice to verify assumption 1, 2 and 3 in Mei et al. (2018). Specifically, first, we will verify that the directional gradient of the population risk is sub-Gaussian (Assumption 1 in Mei et al. (2018)). Note the directional gradient of the population risk is given by $\langle \nabla \rho(Y - \langle X, \theta \rangle), \nu \rangle = \psi(Y - \langle X, \theta \rangle) \langle X, \nu \rangle$. Since $|\psi(Y - \langle X, \theta \rangle)| \leq L_\psi$, and $\langle X, \nu \rangle$ is mean zero and $\tau^2$-sub-Gaussian by our assumption 1, due to Lemma 1 in Mei et al. (2018), there exists a universal constant $C_1$, such that $\langle \nabla \rho(Y - \langle X, \theta \rangle), \nu \rangle$ is $C_1 L_\psi \tau^2-$sub-Gaussian. Second, we will verify that the directional Hessian of the loss is sub-exponential (Assumption 2 in Mei et al. (2018)). The directional Hessian of the loss gives $\langle \nabla^2 \rho(Y - \langle X, \theta \rangle) \nu, \nu \rangle = \psi'(Y - \langle X, \theta \rangle) \langle X, \nu \rangle^2$. Since $|\psi'(Y - \langle X, \theta \rangle)| \leq L_\psi$, by Lemma 1 in Mei et al. (2018), $\langle \nabla^2 \rho(Y - \langle X, \theta \rangle) \nu, \nu \rangle$ is $C_2 \tau^2$-sub-exponential. Third, let $H = ||\nabla^2 R(\theta_0)||_{op}$ and $J^* = \mathbf{E}\left[\sup_{\theta_1 \neq \theta_2} \frac{||(\psi'(\theta_1) - \psi'(\theta_2)) x x^T||_{op}}{||\theta_1 - \theta_2||_2}\right]$. Then, we can show $H \leq L_\psi \tau^2$ and $J^* \leq L_\psi (p\tau^2)^{3/2}$. Therefore, there exists a constant $c_h$ such that $H \leq \tau^2 p^{c_h}$ and $J^* \leq \tau^3 p^{c_h}$, which verifies the assumption 3 in Mei et al. (2018). Therefore, the uniform convergence of gradient and Hessian in theorem 1 in Mei et al. (2018) holds for our gross error model. $\qquad\square$

**Proof of Theorem 1**: Part (a): It is suffice to show that $\langle \theta - \theta_0, \nabla R(\theta) \rangle > 0$ for all $||\theta - \theta_0||_2 > \eta_0$. Note by Assumption 1(d), we have $h(z) = \int_{-\infty}^{+\infty} \psi(z + \epsilon) f_0(\epsilon) d\epsilon > 0$ as $z > 0$ and $h'(0) > 0$. Define $H(s) := \inf_{0 \leq z \leq s} \frac{h(z)}{z}$, it is easy to see that $H(s) > 0$ for all $s > 0$.

Then, we have

$$
\begin{aligned}
\langle \theta - \theta_0, \nabla R(\theta) \rangle &= \mathbf{E}\left[\mathbf{E}[\psi(z+\epsilon)z | z = \langle \theta_0 - \theta, X \rangle]\right] \\
&= (1-\delta)\mathbf{E}[h(\langle \theta - \theta_0, X \rangle)\langle \theta - \theta_0, X \rangle] + \delta \mathbf{E}\left[\mathbf{E}_g(\psi(z+\epsilon)z | z = \langle \theta_0 - \theta, X \rangle)\right] \\
&\geq (1-\delta)H(s)\mathbf{E}[\langle \theta - \theta_0, X \rangle^2 I_{(|\langle \theta - \theta_0, X \rangle| \leq s)}] - \delta L_\psi \mathbf{E}|\langle \theta_0 - \theta, X \rangle| \\
&= (1-\delta)H(s)\mathbf{E}[\langle \theta - \theta_0, X \rangle^2 - \langle \theta - \theta_0, X \rangle^2 I_{(|\langle \theta - \theta_0, X \rangle| > s)}] - \delta L_\psi \mathbf{E}|\langle \theta - \theta_0, X \rangle| \\
&\geq (1-\delta)H(s)\left[\mathbf{E}[\langle \theta - \theta_0, X \rangle^2] - \left(\mathbf{E}[\langle \theta - \theta_0, X \rangle^4] \cdot \mathbf{P}(|\langle \theta - \theta_0, X \rangle| > s)\right)^{1/2}\right] \\
&\quad - \delta L_\psi (\mathbf{E}|\langle \theta - \theta_0, X \rangle|^2)^{1/2} \\
&\overset{(i)}{\geq} (1-\delta)H(s)||\theta - \theta_0||_2^2 \tau^2 \left(\gamma - \sqrt{c_2 \mathbf{P}(|\langle \theta - \theta_0, X \rangle| > s)}\right) - \delta L_\psi ||\theta - \theta_0||_2 \tau \\
&\overset{(ii)}{\geq} (1-\delta)H(s)||\theta - \theta_0||_2^2 \tau^2 \left(\gamma - \sqrt{\frac{c_2 \mathbf{E}(|\langle \theta - \theta_0, X \rangle|^4)}{s^4}}\right) - \delta L_\psi ||\theta - \theta_0||_2 \tau \\
&\geq (1-\delta)H(s)||\theta - \theta_0||_2^2 \tau^2 \left(\gamma - \sqrt{\frac{c_2 \cdot c_2 \tau^4 ||\theta - \theta_0||_2^4}{s^4}}\right) - \delta L_\psi ||\theta - \theta_0||_2 \tau \\
&\geq (1-\delta)H(s)||\theta - \theta_0||_2^2 \tau^2 \left(\gamma - \frac{c_2 \tau^2 ||\theta - \theta_0||_2^2}{s^2}\right) - \delta L_\psi ||\theta - \theta_0||_2 \tau \\
&\geq (1-\delta)H(s)||\theta - \theta_0||_2^2 \tau^2 \left(\gamma - \frac{16 c_2 \tau^2 r^2}{9 s^2}\right) - \delta L_\psi ||\theta - \theta_0||_2 \tau.
\end{aligned}
$$

Here (i) holds from the fact that if $X$ has mean zero and is $\tau^2$-sub-Gaussian, then for all $u \in \mathbb{R}^p$,

$$
\begin{aligned}
\mathbf{E}|\langle u, X \rangle|^2 &\leq ||u||_2^2 \tau^2, \\
\mathbf{E}|\langle u, X \rangle|^4 &\leq c_2 ||u||_2^4 \tau^4,
\end{aligned}
$$

where $c_2$ is a constant (Boucheron et al., 2013). (ii) holds from Chebyshev's inequality. Thus, a choice of $\tilde{s} = \frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}}$ will ensure that

$$
\langle \theta - \theta_0, \nabla R(\theta) \rangle \geq (1-\delta)\frac{3}{4}H\left(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}}\right)||\theta - \theta_0||_2^2 \tau^2 \gamma - \delta L_\psi ||\theta - \theta_0||_2 \tau, \qquad (20)
$$

27

which is greater than 0 when

$$\|\theta - \theta_0\|_2 > \frac{\delta L_\psi}{(1-\delta)\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma} := \eta_0. \tag{21}$$

Therefore, there are no stationary point outside of the ball $B_2^p(\theta_0, \eta_0)$.

Part(b): We first look at the minimum eigenvalue of the Hessian $\nabla^2 R(\theta)$ at $\theta = \theta_0$. For any $u \in \mathbb{R}^p, \|u\|_2 = 1$,

$$\begin{aligned}
\langle u, \nabla^2 R(\theta_0) u \rangle &= (1-\delta)\mathbf{E}_{f_0}[\psi'(\epsilon)\langle X, u \rangle^2] + \delta\mathbf{E}_g[\psi'(\epsilon)\langle X, u \rangle^2] \\
&= (1-\delta)\mathbf{E}_{f_0}[\psi'(\epsilon)]\mathbf{E}[\langle X, u \rangle^2] + \delta\mathbf{E}_g[\psi'(\epsilon)\langle X, u \rangle^2] \\
&\geq (1-\delta)h'(0)\gamma\tau^2 - \delta L_\psi\tau^2.
\end{aligned}$$

Therefore, we have the minimum eigenvalue of $\nabla^2 R(\theta_0)$ is greater than 0 as long as $\delta < \frac{h'(0)\gamma}{h'(0)\gamma + L_\psi}$.

Then we look at the operator norm of $\nabla^2 R(\theta) - \nabla^2 R(\theta_0)$. For any $u \in \mathbb{R}^p, \|u\|_2 = 1$,

$$\begin{aligned}
|\langle u, (\nabla^2 R(\theta) - \nabla^2 R(\theta_0))u \rangle| &= |\mathbf{E}[(\psi'(\langle X, \theta_0 - \theta \rangle + \epsilon) - \psi'(\epsilon))\langle X, u \rangle^2]| \\
&= |\mathbf{E}[\psi''(\xi)\langle X, \theta_0 - \theta \rangle\langle X, u \rangle^2]| \\
&\leq \mathbf{E}|\psi''(\xi)|\mathbf{E}|\langle X, \theta_0 - \theta \rangle\langle X, u \rangle^2| \\
&\leq L_\psi\{\mathbf{E}[\langle X, \theta_0 - \theta \rangle^2]\mathbf{E}[\langle X, u \rangle^4]\}^{1/2} \\
&\leq L_\psi(\|\theta_0 - \theta\|_2^2\tau^2 c_2\tau^4)^{1/2} \\
&= L_\psi\sqrt{c_2}\|\theta_0 - \theta\|_2\tau^3.
\end{aligned}$$

Hence, taking

$$\|\theta - \theta_0\|_2 \leq \eta_1 := \frac{(1-\delta)h'(0)\gamma - \delta L_\psi}{2\sqrt{c_2}\tau L_\psi} \tag{22}$$

28

guarantees that $(\nabla^2 R(\theta) - \nabla^2 R(\theta_0))_{op} \leq \frac{(1-\delta)h'(0)\gamma\tau^2 - \delta L_\psi \tau^2}{2}$. Therefore, for all $\theta \in B_2^p(\theta_0, \eta_1)$, we have

$$\lambda_{\min}(\nabla^2 R(\theta)) \geq \kappa := \frac{(1-\delta)h'(0)\gamma - \delta L_\psi}{2}\tau^2, \tag{23}$$

which yields there is at most one minimizer of $R(\theta)$ in the ball $B_2^p(\theta_0, \eta_1)$, as long as $\delta < \frac{h'(0)\gamma}{h'(0)\gamma + L_\psi}$.

Part (c): Note $R(\theta)$ is a continuous function on $B_2^p(r)$. Thus there exists a global minimizer, denoted by $\theta^*$. Since we have shown that there is no stationary points outside the ball $B_2^p(\theta_0, \eta_0)$, $\theta^*$ should be in the ball $B_2^p(\theta_0, \eta_0)$. Therefore, as long as $\eta_1 > \eta_0$, i.e.,

$$\frac{(1-\delta)h'(0)\gamma - \delta L_\psi}{2\sqrt{c_2}\tau L_\psi} > \frac{\delta L_\psi}{(1-\delta)\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma}, \tag{24}$$

there exists and only exists a unique stationary point of $R(\theta)$, which is also the global optimum $\theta^*$. $\qquad\square$

**Proof of Theorem 2** Based on Lemma 1, there exists a constant $C$ such that when $n \geq Cp\log p$,

$$\mathbf{P}\left(\sup_{\theta \in B^p(0,r)} ||\nabla \hat{R}_n(\theta) - \nabla R(\theta)||_2 \leq \tau\delta L_\psi\right) \geq 1 - \pi \tag{25}$$

$$\mathbf{P}\left(\sup_{\theta \in B^p(0,r)} ||\nabla^2 \hat{R}_n(\theta) - \nabla^2 R(\theta)||_{op} \leq \kappa/2\right) \geq 1 - \pi. \tag{26}$$

Part (a): Note

$$\langle \theta - \theta_0, \nabla \widehat{R}_n(\theta) \rangle \geq \langle \theta - \theta_0, \nabla R(\theta) \rangle - ||\nabla \hat{R}_n(\theta) - \nabla R(\theta)||_2 ||\theta - \theta_0||_2 \tag{27}$$

$$\geq (1-\delta)\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})||\theta - \theta_0||_2^2 \tau^2 \gamma - 2\tau\delta L_\psi ||\theta - \theta_0||_2 \tag{28}$$

which is greater than 0 when

$$||\theta - \theta_0||_2 > \frac{2\delta L_\psi}{(1-\delta)\frac{3}{4}L(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma} = 2\eta_0. \tag{29}$$

29

Therefore, there are no stationary points outside of the ball $B_2^p(\theta_0, 2\eta_0)$.

Part (b): For the least eigenvalue of the empirical Hessian in $B_2^p(\theta_0, \eta_1)$, we have

$$
\begin{aligned}
\inf_{||\theta-\theta_0||_2 \leq \eta_1} \lambda_{\min}(\nabla^2 \widehat{R}_n(\theta)) &\geq \inf_{||\theta-\theta_0||_2 \leq \eta_1} \lambda_{\min}(\nabla^2 R(\theta)) - \sup_{\theta \in B^p(0,\eta_1)} ||\nabla^2 \widehat{R}_n(\theta) - \nabla^2 R(\theta)||_{op} \\
&\geq \kappa - \kappa/2 = \kappa/2 > 0.
\end{aligned} \tag{30}
$$

This lead to the conclusion that, $\widehat{R}_n(\theta)$ is strong convex inside the ball $B_2^p(\theta_0, \eta_1)$.

Part(c): When $2\eta_0 < \eta_1$, by strong convexity of $\widehat{R}_n(\theta)$ in $B_2^p(\theta_0, \eta_1)$, there exists a unique local minimizer, which is in $B_2^p(\theta_0, 2\eta_0)$. We denote the unique local minimizer as $\widehat{\theta}_n$.

By Theorem 1, there is a unique stationary point of the population risk function $R(\theta)$ in the ball $B_2^p(\theta_0, \eta_0)$. Suppose $\theta^*$ is the unique stationary point of $R(\theta)$. By Taylor expansion of $\widehat{R}_n(\theta)$ at the point $\theta^*$, there exists a $\tilde{\theta}$ in $B^p(\theta_0, 2\eta_0)$, such that

$$
\widehat{R}_n(\widehat{\theta}_n) = \widehat{R}_n(\theta^*) + \langle \widehat{\theta}_n - \theta^*, \nabla \widehat{R}_n(\theta^*) \rangle + \frac{1}{2}(\widehat{\theta}_n - \theta^*)' \nabla^2 \widehat{R}_n(\tilde{\theta})(\widehat{\theta}_n - \theta^*) \leq \widehat{R}_n(\theta^*). \tag{31}
$$

Since by equation (30), the least eigenvalue of $\nabla^2 \widehat{R}_n(\tilde{\theta})$ is greater than $\kappa/2$, which lead to

$$
\frac{\kappa}{4}||\widehat{\theta}_n - \theta^*||_2^2 \leq \langle \theta^* - \widehat{\theta}_n, \nabla \widehat{R}_n(\theta^*) \rangle \leq ||\theta^* - \widehat{\theta}_n||_2 ||\nabla \widehat{R}_n(\theta^*)||_2, \tag{32}
$$

which yield

$$
||\widehat{\theta}_n - \theta^*||_2 \leq \frac{4}{\kappa}||\nabla \widehat{R}_n(\theta^*)||_2. \tag{33}
$$

By Theorem 1, $||\theta_0 - \theta^*||_2 < \eta_0$, combined with equation (33) and the uniform convergency theorem in Lemma 1 yield

$$
||\widehat{\theta}_n - \theta_0||_2 \leq \eta_0 + \frac{4\tau}{\kappa}\sqrt{\frac{C * p \log n}{n}}. \tag{34}
$$

□

**Proof of lemma 2:** From the Theorem 3 in Mei et al. (2018), the uniform convergency theorem of our Lemma 2 holds if Assumption 4, 5 in Mei et al. (2018) hold under the contaminated model with outliers. Here we will show under our assumption 1 and 2, there exist constants $T_0$ and $L_0$ such that

**a** For all $\theta \in B_2^p(r)$, $Y \in \mathbb{R}$, $X \in \mathbb{R}^p$, $||\nabla_\theta \rho(Y - \langle X, \theta \rangle)||_\infty \leq T_0 M$

**b** There exist functions $h_1 : \mathbb{R} \times \mathbb{R}^{p+1} \to \mathbb{R}$, and $h_2 : \mathbb{R}^{p+1} \to \mathbb{R}^p$, such that

$$\langle \nabla_\theta \rho(Y - \langle X, \theta \rangle), \theta - \theta_0 \rangle = h_1(\langle \theta - \theta_0, h_2(Y, X) \rangle), Y, X). \tag{35}$$

In addition, $h_1(t, Y, X)$ is $L_0 M$- Lipschitz to its first argument $t$, $h_1(0, Y, X) = 0$, and $h_2(Y, X)$ is mean-zero and $\tau^2$-sub-Gaussian.

Part (a). The gradient of the loss is

$$\nabla_\theta \rho(Y - \langle X, \theta \rangle) = -\psi(Y - \langle X, \theta \rangle) X. \tag{36}$$

By assumption 1, we have $|-\psi(Y - \langle X, \theta \rangle)| \leq L_\psi$. By assumption 2, we have $||X||_\infty \leq M\tau$. Therefore, (a) is satisfied with parameter $T_0 = L_\psi \tau$.

Part (b). Note

$$\langle \nabla_\theta \rho(Y - \langle X, \theta \rangle), \theta - \theta_0 \rangle = -\psi(Y - \langle X, \theta \rangle) \langle X, \theta - \theta_0 \rangle. \tag{37}$$

We take $h_2(Y, X) = X$, $t = \langle X, \theta - \theta_0 \rangle$ and $h_1(t, Y, X) = -\psi(Y - t - \langle X, \theta_0 \rangle) t$. Clearly, we have $h_1(0, Y, X) = 0$ and $h_2(Y, X)$ is mean 0 and $\tau^2$-sub-Gaussian. Furthermore, note $|t| \leq 2rM\tau$, we have

$$|\frac{\partial}{\partial t} h_1(t, Y, X)| = |\psi'(Y - t - \langle X, \theta_0 \rangle) t - \psi(Y - t - \langle X, \theta_0 \rangle)| \tag{38}$$

$$\leq 2ML_\psi r\tau + L_\psi \tag{39}$$

$$\leq (2L_\psi r\tau + L_\psi) M. \tag{40}$$

31

Therefore, $h_1(t, X, Y)$ is at most $(2L_\psi r\tau + L_\psi)M$-Lipschitz in its first argument $t$. By part (a) and part (b), we can see assumption 4, 5 are satisfied under the gross error model, which prove the uniform convergency theorem in our Lemma 2. $\qquad\square$

**Proof of theorem 3:** We decompose the proof into four technical lemmas. First, in Lemma 3, we prove there cannot be any stationary points of the regularized empirical risk $\hat{L}_n$ in (10) outside the region $\mathbb{A}$, which is a cone with $\mathbb{A} = \{\theta_0 + \Delta : ||\Delta_{S_0^c}||_1 \leq 3||\Delta_{S_0}||_1\}$. Then in Lemma 4, we show there cannot be any stationary points outside the region $B_2^p(\theta_0, r_s)$ where $r_s$ is the statistical radius which is not less than $\eta_0$ in Theorem 1. In Lemma 5, we argue that all stationary points should have support size less or equal to $cs_0 \log p$. Finally, in Lemma 6, we show there cannot be two stationary points in $B_2^p(\theta_0, \eta_1) \cap \mathbb{A}$. Note $\hat{L}_n(\theta)$ is a continuous function, which indicates the existence of the global minimizer. Therefore, we can conclude there is and only is one unique stationary point of the regularized empirical risk $\hat{L}_n$ as long as $r_s < \eta_1$.

To start with those lemmas, we define the subgradient of $\hat{L}_n$ at $\theta$ as:

$$\partial\hat{L}_n(\theta) = \{\nabla R_n(\theta) + \lambda_n \nu : \nu \in \partial||\theta||_1\}. \tag{41}$$

Therefore, the optimality condition implies that $\theta$ is a stationary point of $\hat{L}_n$ if and only if $\mathbf{0} \in \partial\hat{L}_n(\theta)$. To simplify notations, all constants in the following lemmas are dependent on $(\rho, L_\psi, \tau^2, r, \gamma, \pi)$ but independent on $\delta, s_0, n, p, M$.

**Lemma 3.** *Let* $S_0 = supp(\theta_0)$ *and* $s_0 = |S_0|$. *Define a cone* $\mathbb{A} = \{\theta_0 + \Delta : ||\Delta_{S_0^c}||_1 \leq 3||\Delta_{S_0}||_1\} \subseteq \mathbb{R}^p$. *For any* $\pi > 0$, *there exist constants* $C_0, C_1$ *such that letting* $\lambda_n \geq C_0 M \sqrt{\frac{\log p}{n}} + \delta\frac{C_1}{\sqrt{s_0}}$, *with probability at least* $1-\pi$, $\hat{L}_n(\theta)$ *has no stationary points in* $B_2^p(0, r) \cap \mathbb{A}^c$:

$$\langle z(\theta), \theta - \theta_0 \rangle > 0, \quad \forall \theta \in B_2^p(0, r) \cap \mathbb{A}^c, z(\theta) \in \partial\hat{L}_n(\theta) \tag{42}$$

*Proof.* For any $z(\theta) \in \partial \hat{L}_n(\theta)$, it can be written as $z(\theta) = \nabla \hat{R}_n(\theta) + \lambda_n \nu(\theta)$, where $\nu(\theta) \in \partial ||\theta||_1$. Therefore, we have

$$\langle z(\theta), \theta - \theta_0 \rangle = \langle \nabla R(\theta), \theta - \theta_0 \rangle + \langle \nabla \hat{R}_n(\theta) - \nabla R(\theta), \theta - \theta_0 \rangle + \lambda_n \langle \nu(\theta), \theta - \theta_0 \rangle \quad (43)$$

Note by (20) we have

$$\langle \theta - \theta_0, \nabla R(\theta) \rangle \geq (1 - \delta) \frac{3}{4} H(\frac{8\tau r}{3} \sqrt{\frac{c_2}{\gamma}}) ||\theta - \theta_0||_2^2 \tau^2 \gamma - \delta L_\psi ||\theta - \theta_0||_2 \tau. \quad (44)$$

By lemma 2, for any $\pi > 0$, there exists a constant $C_\pi$ such that

$$\mathbf{P}(\sup_{0 < ||\theta||_2 < r} \frac{|\langle \nabla \hat{R}_n(\theta) - \nabla R(\theta), \theta - \theta_0 \rangle|}{||\theta - \theta_0||_1} \leq C_\pi M \sqrt{\frac{\log p}{n}}) > 1 - \pi. \quad (45)$$

Letting $\Delta = \theta - \theta_0$, we have

$$\langle \nu(\theta), \theta - \theta_0 \rangle = \langle \nu(\theta)_{S_0^c}, \Delta_{S_0^c} \rangle + \langle \nu(\theta)_{S_0}, \Delta_{S_0} \rangle \geq ||\Delta_{S_0^c}||_1 - ||\Delta_{S_0}||_1 \quad (46)$$

Plugging (44),(45),(46) into (43) yields

$$\langle z(\theta), \theta - \theta_0 \rangle \geq (1 - \delta) \frac{3}{4} H(\frac{8\tau r}{3} \sqrt{\frac{c_2}{\gamma}}) ||\theta - \theta_0||_2^2 \tau^2 \gamma - \delta L_\psi ||\theta - \theta_0||_2 \tau \quad (47)$$

$$- C_\pi M \sqrt{\frac{\log p}{n}} (||\Delta_{S_0^c}||_1 + ||\Delta_{S_0}||_1) + \lambda_n (||\Delta_{S_0^c}||_1 - ||\Delta_{S_0}||_1). \quad (48)$$

Let $\lambda_n \geq 2 C_\pi M \sqrt{\frac{\log p}{n}} + C_2$, we have

$$\langle z(\theta), \theta - \theta_0 \rangle \geq (1 - \delta) \frac{3}{4} H(\frac{8\tau r}{3} \sqrt{\frac{c_2}{\gamma}}) ||\theta - \theta_0||_2^2 \tau^2 \gamma - \delta L_\psi ||\theta - \theta_0||_2 \tau$$

$$+ C_\pi M \sqrt{\frac{\log p}{n}} (||\Delta_{S_0^c}||_1 - 3||\Delta_{S_0}||_1) + C_2 (||\Delta_{S_0^c}||_1 - ||\Delta_{S_0}||_1). \quad (49)$$

Next, we will find the lower bound of $||\Delta_{S_0^c}||_1 - ||\Delta_{S_0}||_1$ under the constraint of $||\Delta_{S_0^c}||_1 - 3||\Delta_{S_0}||_1 \geq 0$. Note by Cauchy inequality, we have

$$||\Delta||_2^2 \geq \frac{||\Delta_{S_0^c}||_1^2}{p - s_0} + \frac{||\Delta_{S_0}||_1^2}{s_0} \quad (50)$$

33

Therefore, under the constraint of $||\Delta_{S_0^c}||_1 - 3||\Delta_{S_0}||_1 \geq 0$, the minimal value of $||\Delta_{S_0^c}||_1 - ||\Delta_{S_0}||_1$ is obtained when $||\Delta_{S_0^c}||_1 - 3||\Delta_{S_0}||_1 = 0$ and $||\Delta||_2^2 = \frac{||\Delta_{S_0^c}||_1^2}{p-s_0} + \frac{||\Delta_{S_0}||_1^2}{s_0}$. By solving the two equations yield

$$||\Delta_{S_0^c}||_1 = 3\sqrt{\frac{(p-s_0)s_0}{8s_0+p}}||\Delta||_2 \tag{51}$$

$$||\Delta_{S_0}||_1 = \sqrt{\frac{(p-s_0)s_0}{8s_0+p}}||\Delta||_2 \tag{52}$$

and $||\Delta_{S_0^c}||_1 - ||\Delta_{S_0}||_1 \geq 2\sqrt{\frac{(p-s_0)s_0}{8s_0+p}}||\Delta||_2$. Combined with (49), setting $C_1 = \frac{L_\psi \tau}{2}$ and $C_2 = C_1\frac{\delta}{\sqrt{s_0}}$ yield $2\sqrt{\frac{(p-s_0)s_0}{8s_0+p}}C_2 \geq \delta L_\psi \tau$, which implies $\langle z(\theta), \theta - \theta_0 \rangle > 0$, as long as $\theta \in \mathbb{A}^c$, i.e., $||\Delta_{S_0^c}||_1 - 3||\Delta_{S_0}||_1 > 0$. $\qquad\square$

**Lemma 4.** *For any $\pi > 0$, $\theta \in \mathbb{A}$, $z(\theta) \in \partial\hat{L}_n(\theta)$, there exist constants $C_0$, $C_1$ such that with probability at least $1 - \pi$,*

$$\langle z(\theta), \theta - \theta_0 \rangle > 0 \tag{53}$$

*as long as $||\theta - \theta_0||_2 > r_s$, where*

$$r_s = \frac{\delta}{1-\delta}C_0 + \frac{4\sqrt{s_0}}{1-\delta}(M\sqrt{\frac{\log p}{n}} + \lambda_n)C_1. \tag{54}$$

*Proof.* Since for any $\theta \in \mathbb{A}$, we have $||\theta - \theta_0||_1 \leq 4\sqrt{s_0}||\theta - \theta_0||_2$. Combining with (43) yields

$$\langle z(\theta), \theta - \theta_0 \rangle \geq \langle \nabla R(\theta), \theta - \theta_0 \rangle - C_\pi M\sqrt{\frac{\log p}{n}}||\theta - \theta_0||_1 - \lambda_n||\theta - \theta_1||_1 \tag{55}$$

$$\geq (1-\delta)\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})||\theta - \theta_0||_2^2\tau^2\gamma - \delta L_\psi||\theta - \theta_0||_2\tau \tag{56}$$

$$-(C_\pi M\sqrt{\frac{\log p}{n}} + \lambda_n)4\sqrt{s_0}||\theta - \theta_0||_2, \tag{57}$$

34

which is greater than 0 as long as

$$||\theta - \theta_0||_2 \geq \frac{\delta L_\psi + (C_\pi M\sqrt{\frac{\log p}{n}} + \lambda_n)4\sqrt{s_0}}{(1-\delta)\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma} := r_s. \tag{58}$$

Taking $C_0 = \frac{L_\psi}{\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma}$ and $C_1 = \frac{\max(1,C_\pi)}{\frac{3}{4}H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma}$ give the result of $r_s$ in equation (54). $\square$

**Lemma 5.** *If $\delta \leq 1/2$, for any $\pi$, there exist constants $C_0, C_1, C$ such that letting $\lambda_n \geq C_0 M\sqrt{(\log p)/n} + \delta C_1/\sqrt{s_0}$, with probability at least $(1-\pi)$, any stationary points of $\hat{L}_n(\theta)$ in $B_2^p(\theta_0, r_s) \cap \mathbb{A}$ has support size $|S(\hat{\theta})| \leq Cs_0 \log p$.*

*Proof.* Let $\hat{\theta} \in B_2^p(\theta_0, r_s) \cap \mathbb{A}$ be a stationary point of $\hat{L}_n(\theta)$ in (10). Then we have

$$\nabla R_n(\hat{\theta}) + \lambda_n \nu(\hat{\theta}) = 0, \tag{59}$$

where $\nu(\hat{\theta}) \in ||\hat{\theta}||_1$. Thus, we have

$$\left(\nabla R_n(\hat{\theta})\right)_j = \pm\lambda_n, \quad \forall j \in S(\hat{\theta}) \tag{60}$$

Note $|\psi(y_i - \langle x_i, \theta_0\rangle)| \leq L_\psi$ and $\langle x_i, e_j\rangle$ is $\tau^2$-subgaussian with mean 0. Then there exists an absolute constant $c_0$ such that $\psi(y_i - \langle x_i, \theta_0\rangle)\langle x_i, e_j\rangle$ is $c_0 L_\psi^2 \tau^2$-subgaussian, see Lemma 1(d) in Mei et al. (2018). Thus we have $\frac{1}{n}\sum_{i=1}^n \psi(y_i - \langle x_i, \theta_0\rangle)\langle x_i, e_j\rangle$ is $c_0 L_\psi^2 \tau^2/n$-subgaussian with mean $\langle \nabla R(\theta_0), e_j\rangle$. Moreover, note $|\langle \nabla R(\theta_0), e_j\rangle| = |\delta \mathbf{E}_g \psi(y_i - \langle x_i, \theta_0\rangle)\langle x_i, e_j\rangle| \leq \delta L_\psi \mathbf{E}|\langle x_i, e_j\rangle| \leq \delta L_\psi \tau$, we have for any $t > 0$,

$$\mathbf{P}(|\frac{1}{n}\sum_{i=1}^n \psi(y_i - \langle x_i, \theta_0\rangle)\langle x_i, e_j\rangle| \geq t + \delta L_\psi \tau)$$

$$\leq \mathbf{P}(|\frac{1}{n}\sum_{i=1}^n \psi(y_i - \langle x_i, \theta_0\rangle)\langle x_i, e_j\rangle - \langle \nabla R(\theta_0), e_j\rangle| \geq t) \leq 2\exp(-\frac{t^2 n}{2c_0 L_\psi^2 \tau^2}). \tag{61}$$

Thus, we can get

$$\mathbf{P}\left(||\nabla R_n(\theta_0)||_\infty > t + \delta L_\psi \tau\right) \quad \leq \quad p \max_{1 \leq j \leq p} \mathbf{P}\left(|\frac{1}{n}\sum_{i=1}^n \psi(y_i - \langle x_i, \theta_0 \rangle)\langle x_i, e_j \rangle| > t + \delta L_\psi \tau\right)$$

$$\leq \quad 2p \exp(-\frac{t^2 n}{2c_0 L_\psi^2 \tau^2}). \tag{62}$$

Thus, a choice of $t = L_\psi \tau \sqrt{\frac{2c_0(\log p + \log 6/\pi)}{n}}$ and $C = \sqrt{c_0 \log 6/\pi}$ will guarantee that

$$\mathbf{P}\left(||\nabla \hat{R}_n(\theta_0)||_\infty > L_\psi \tau (C\sqrt{\frac{\log p}{n}} + \delta)\right) \quad \leq \pi/3 \tag{63}$$

Let $\lambda_n \geq 2L_\psi \tau(C\sqrt{\frac{\log p}{n}} + \delta)$, we have the event $(||\nabla R_n(\theta_0)||_\infty < \lambda_n/2)$ happens with the probability at least $1 - \pi/3$. Under this event, combing with (60) yields

$$\lambda_n/2 \leq \left|\left(\nabla R_n(\theta_0) - \nabla R_n(\hat{\theta})\right)_j\right|, \quad \forall j \in S(\hat{\theta}). \tag{64}$$

Squaring and summing over $j \in S(\hat{\theta})$, we have

$$\lambda_n^2 |S(\hat{\theta})| \quad \leq \quad 4\left\|\left(\nabla \hat{R}_n(\theta_0) - \nabla \hat{R}_n(\hat{\theta})\right)_{S(\hat{\theta})}\right\|_2^2 \tag{65}$$

$$= \quad 4\left\|\left(\frac{1}{n}\sum_{i=1}^n \left(\psi(y_i - \langle \theta_0, x_i \rangle) - \psi(y_i - \langle \hat{\theta}, x_i \rangle)\right) x_i\right)_{S(\hat{\theta})}\right\|_2^2 \tag{66}$$

$$= \quad 4\left\|\left(\frac{1}{n}\sum_{i=1}^n \left(\psi'(y_i - \langle \beta_i, x_i \rangle)\right) \langle \theta_0 - \hat{\theta}, x_i \rangle x_i\right)_{S(\hat{\theta})}\right\|_2^2 \tag{67}$$

$$\leq \quad 4L_\psi^2 \left\|\left(\frac{1}{n}\sum_{i=1}^n \langle \theta_0 - \hat{\theta}, x_i \rangle x_i\right)_{S(\hat{\theta})}\right\|_2^2 \tag{68}$$

where $\beta_i$ are located on the line between $\theta_0$ and $\hat{\theta}$ obtained by intermediate value theorem.

36

Moreover, by Minkowski inequality and Cauchy-Schwarz inequality yield

$$\left\|\left(\frac{1}{n}\sum_{i=1}^{n}\langle\theta_0-\hat{\theta},x_i\rangle x_i\right)_{S(\hat{\theta})}\right\|_2 \leq \frac{1}{n}\sum_{i=1}^{n}|\langle\theta_0-\hat{\theta},x_i\rangle|\left\|(x_i)_{S(\hat{\theta})}\right\|_2$$

$$\leq \frac{1}{n}\left((\sum_{i=1}^{n}|\langle\theta_0-\hat{\theta},x_i\rangle|^2)(\sum_{i=1}^{n}\|(x_i)_{S(\hat{\theta})}\|_2^2)\right)^{1/2} \quad (69)$$

Due to the restricted smoothness property of the sub-Gaussian random variables Mei et al. (2018), there exists a constant $c_1$ depending on $\pi$ such that with probability at least $1-\pi/3$, as $n \geq c_1 s_0 \log p$, we have

$$\sup_{\theta\in\mathbb{A}}\frac{\frac{1}{n}(\sum_{i=1}^{n}|\langle\theta_0-\theta,x_i\rangle|^2)}{\|\theta-\theta_0\|_2^2} \leq 3\tau^2. \quad (70)$$

Therefore, with probability at least $1-\pi/3$, we have

$$\sup_{\theta\in\mathbb{A}\cap B^p(\theta_0,r_s)}\frac{1}{n}(\sum_{i=1}^{n}|\langle\theta_0-\hat{\theta},x_i\rangle|^2) \leq 3\tau^2\sup_{\theta\in\mathbb{A}\cap B^p(\theta_0,r_s)}\|\theta-\theta_0\|_2^2 \leq 3\tau^2 r_s^2. \quad (71)$$

Moreover, by Lemma 13 in Mei et al. (2018), for any $\pi$, there exists constant $c_2$ depending on $\pi$ such that

$$\mathbf{P}(\frac{1}{n}\sum_{i=1}^{n}\|(x_i)_{S(\hat{\theta})}\|_2^2 > c_2\tau^2\log p) \leq \pi/3. \quad (72)$$

By (63,71,72), as well as (69), at least $1-\pi$,

$$\lambda_n^2|S(\hat{\theta})| \leq 4L_\psi^2 3\tau^2 r_s^2 c_2\tau^2\log p \quad (73)$$

$$= Cr_s^2\log p \quad (74)$$

By equation (54) we have

$$r_s^2 \leq C_0(\frac{\delta}{1-\delta})^2 + \frac{s_0}{(1-\delta)^2}(M^2\frac{\log p}{n}+\lambda_n^2)C_1 \quad (75)$$

37

Taking $\lambda_n \geq C_2 M \sqrt{(\log p)/n} + C_3\delta/\sqrt{s_0}$ gives us

$$|S(\hat{\theta})| \leq (C_4\frac{s_0}{(1-\delta)^2} + s_0 C_5)\log p \tag{76}$$

$$= Cs_0\log p \tag{77}$$

$\square$

**Lemma 6.** *For any positive constants $C_0$ and $\pi$, letting $r_0 = C_0 s_0 \log p$, there exist constant $C_1$ such that when $n \geq C_1 s_0 \log^2 p$,*

$$\mathbf{P}(\sup_{\theta \in B_2^p(\theta_0,r)\cap B_0^p(0,r_0)}\sup_{\nu \in B_2^p(0,1)\cap B_0^p(0,r_0)}\langle\nu, (\nabla^2\hat{R}_n(\theta) - \nabla^2 R(\theta))\nu\rangle \leq \kappa/2) \geq 1 - \pi. \tag{78}$$

*Moreover, the regularized empirical risk $\hat{L}_n(\theta)$ in (10) cannot have two stationary points in the region $B_2^p(\theta_0, \eta_1) \cap B_0^p(0, r_0/2)$.*

*Proof.* According to (23), we have

$$\inf_{\theta \in B_2^p(\theta_0,\eta_1)} \lambda_{\min}(\nabla^2 R(\theta)) \geq \kappa. \tag{79}$$

By lemma 2, there exists constant $C$ such that when $n \geq Cs_0 \log^2 p$,

$$\mathbf{P}\left(\inf_{\theta \in B_2^p(\theta_0,\eta_1)\cap B_0^p(0,r_0)}\inf_{\nu \in B_2^p(0,1)\cap B_0^p(0,r_0)}\langle\nu, (\nabla^2\hat{R}_n(\theta))\nu\rangle \geq \kappa/2\right) \leq \pi. \tag{80}$$

Suppose $\theta_1, \theta_2$ are two distinct stationary points of $\hat{L}_n(\theta)$ in $B_2^p(\theta_0, \eta_1) \cap B_0^p(0, r_0/2)$. Define $u = \frac{\theta_2 - \theta_1}{||\theta_1 - \theta_2||_2}$. Since $\theta_1$ and $\theta_2$ are $r_0/2$-sparse, $u$ is $r_0$ sparse, as well as $\theta_1 + tu$ for any $t \in \mathbb{R}$. Therefore,

$$\langle\nabla\hat{R}_n(\theta_2), u\rangle = \langle\nabla\hat{R}_n(\theta_1), u\rangle + \int_0^{||\theta_1-\theta_2||_2}\langle u, \nabla^2\hat{R}_n(\theta_1 + tu)u\rangle dt$$

$$\geq \langle\nabla\hat{R}_n(\theta_1), u\rangle + \frac{\kappa}{2}||\theta_2 - \theta_1||_2. \tag{81}$$

38

Note the regularization term $\lambda_n||\theta||_1$ is convex, we have for any subgradients $\nu(\theta_1) \in \partial||\theta_1||_1$, $\nu(\theta_2) \in \partial||\theta_2||_1$,

$$\lambda_n\langle\nu(\theta_2), u\rangle \geq \lambda_n\langle\nu(\theta_1), u\rangle. \tag{82}$$

Adding (81) with (82) gives

$$\langle\nabla\hat{R}_n(\theta_2) + \lambda_n\nu(\theta_2), u\rangle \geq \langle\nabla\hat{R}_n(\theta_1) + \lambda_n\nu(\theta_1), u\rangle + \frac{\kappa}{2}||\theta_2 - \theta_1||_2, \tag{83}$$

which is contradict with the assumption that $\theta_1$ and $\theta_2$ are two distinct stationary points of $\hat{L}_n(\theta)$. $\qquad\square$

**Proof of Theorem 3**. Now we are ready to prove Theorem 3. By Lemma 3 and Lemma 4, as $n \geq Cs_0 \log p$, letting $\lambda_n \geq C_0 M\sqrt{\frac{\log p}{n}} + \delta\frac{C_1}{\sqrt{s_0}}$, all stationary points of $L_n(\theta)$ are in $B_2^p(\theta_0, r_s) \cap \mathbb{A} \cap B_0^p(C_1 s_0 \log p)$, where $r_s$ is defined in (54), $\mathbb{A}$ is the cone defined in Lemma 3. This proves Theorem 3(a). Moreover, by Lemma 5, Lemma 6, as $n \geq C_2 s_0 \log^2 p$, $\hat{L}_n(\theta)$ cannot have two distinct stationary points in $B_2^p(\theta_0, \eta_1) \cap \mathbb{A} \cap B_0^p(C_1 s_0 \log p)$. Thus, as long as $\eta_1 \geq r_s$, there is only one unique stationary point of the regularized empirical risk function $\hat{L}_n(\theta)$, which is the corresponding regularized M-estimator of (10). This proves Theorem 3 (b).

**Proof of Corollary 1:** Huber's loss function is defined by

$$\rho_\alpha(t) = \begin{cases} \frac{1}{2}t^2, & \text{if } |t| \leq \alpha \\ \alpha(|t| - \alpha/2), & \text{if } |t| > \alpha. \end{cases} \tag{84}$$

the corresponding score function would be

$$\psi_\alpha(t) = \rho'_\alpha(t) = \begin{cases} t, & \text{if } |t| \leq \alpha \\ sign(t)\alpha, & \text{if } |t| > \alpha. \end{cases} \tag{85}$$

39

Note for any $\alpha > 0$, all of $\psi(t)$, $\psi'(t)$ and $\psi''(t)$ are bounded. Specifically, we have $|\psi_\alpha(t)| \le \alpha$, $|\psi'(t)| = |\psi''(t)| = 0$. Therefore, the assumptions in Theorem 1 and Theorem 2 are satisfied. It is suffice to find the explicit expression of $\eta_0$ and $\eta_1$ in equation (21) and (22). Since $|\psi'(t)| = |\psi''(t)| = 0$, it is easy to see $\eta_1 = +\infty$, which implies the Huber's estimator has nice computational tractability, regardless the choice of tuning parameter $\alpha$ and the percentage of outliers $\delta$. Moreover, to find the explicit expression of $\eta_0$, according to Assumption 3, we have $c_2 = 3, \gamma = 1$. Thus, we can calculate

$$
\begin{aligned}
h(z) &= \int_{-\infty}^{+\infty} \psi_\alpha(z + \epsilon) f_0(\epsilon) d\epsilon = \int_{-\infty}^{\infty} \psi_\alpha(t) f_0(t - z) dt \\
&= \int_0^\alpha t \left[ f_0(t - z) - f_0(t + z) \right] dt + \alpha \int_\alpha^{+\infty} \left[ f_0(t - z) - f_0(t + z) \right] dt \\
&\ge \int_0^\alpha t \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{t^2 + z^2}{2\sigma^2}} \left( \frac{tz}{\sigma^2} \right) dt + \alpha \int_\alpha^{+\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{t^2 + z^2}{2\sigma^2}} \left( \frac{tz}{\sigma^2} \right) dt \\
&\ge \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{\alpha^2 + z^2}{2\sigma^2}} \int_0^\alpha t \left( \frac{tz}{\sigma^2} \right) dt + \frac{z\alpha}{\sigma^2} e^{-\frac{z^2}{2\sigma^2}} \int_\alpha^{+\infty} t \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{t^2}{2\sigma^2}} dt \\
&= \frac{z\alpha^3}{3\sqrt{2\pi}\sigma^3} e^{-\frac{z^2 + \alpha^2}{2\sigma^2}} + \frac{z\alpha}{\sqrt{2\pi}\sigma} e^{-\frac{z^2 + \alpha^2}{2\sigma^2}}
\end{aligned}
$$

Therefore we have $H(s) = \left( \frac{\alpha^3}{3\sqrt{2\pi}\sigma^3} + \frac{\alpha}{\sqrt{2\pi}\sigma} \right) e^{-\frac{s^2 + \alpha^2}{2\sigma^2}}$. By equation (21) in the proof of Theorem 1 yields

$$
\begin{aligned}
\eta_0(\delta, \alpha) &= \frac{\delta L_\psi}{(1 - \delta) \frac{3}{4} H\left( \frac{8\tau r}{3} \sqrt{\frac{c_2}{\gamma}} \right) \tau \gamma} \tag{86} \\
&= \frac{\delta}{1 - \delta} \frac{4\sqrt{2\pi}\sigma^3}{(\alpha^2 + 3\sigma^2)\tau} e^{\frac{\alpha^2 + 22\tau^2 r^2}{2\sigma^2}}, \tag{87}
\end{aligned}
$$

which complete the proof. $\qquad\square$

**Proof of Corollary 2**: When the loss function is defined by $\rho_\alpha(t) = \frac{1 - e^{-\alpha t^2/2}}{\alpha}$, the corresponding score function would be $\psi_\alpha(t) = \rho'_\alpha(t) = te^{-\alpha t^2/2}$. Moreover, we can get

40

$\psi'_\alpha(t) = e^{-\alpha t^2/2}(1 - \alpha t^2)$ and $\psi''_\alpha(t) = e^{-\alpha t^2/2}\alpha(\alpha t^2 - 3)$. Note for any $\alpha > 0$, all of $\psi_\alpha(t)$, $\psi'_\alpha(t)$ and $\psi''_\alpha(t)$ are bounded.

$$
\begin{aligned}
|\psi_\alpha(t)| &\leq \sqrt{\frac{e}{\alpha}} \\
|\psi'_\alpha(t)| &\leq \max\{1, 2e^{-1.5}\} = 1 \\
|\psi''_\alpha(t)| &\leq \max\{e^{-(3+\sqrt{6})/2}\sqrt{(18 + 6\sqrt{6})\alpha}, e^{-(3-\sqrt{6})/2}\sqrt{(18 - 6\sqrt{6})\alpha}\} \leq 1.5\sqrt{\alpha}.
\end{aligned}
$$

Therefore, the Assumption 1 is satisfied. It is suffice to find the explicit expression of $\eta_0$ and $\eta_1$ in equation (21) and (22). In order to have an accurate expression, we will use the individual bound of $\psi_\alpha(t), \psi'_\alpha(t), \psi''_\alpha(t)$ instead of the universal bound $L_\psi$. Specifically, according to Assumption 4, $x_i$ is $\tau^2$-sub-Gaussian, $c_2 = 3, \gamma = 1/3$. Thus, we can calculate $h(z) = \int_{-\infty}^{+\infty} \psi_\alpha(z + \epsilon)f_0(\epsilon)d\epsilon = \frac{z}{(1+\alpha\sigma^2)^{3/2}}e^{-\frac{\alpha z^2}{2(1+\alpha\sigma^2)}}$ and $H(s) = \frac{1}{(1+\alpha\sigma^2)^{3/2}}e^{-\frac{\alpha s^2}{2(1+\alpha\sigma^2)}}$. By equation (21) in the proof of Theorem 1 yields

$$
\begin{aligned}
\eta_0(\delta, \alpha) &= \frac{\delta L_\psi}{(1 - \delta)\frac{3}{4}H\left(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}}\right)\tau\gamma} \qquad &(88) \\
&= \frac{\delta}{1 - \delta}\sqrt{\frac{e}{\alpha}}\frac{4(1 + \alpha\sigma^2)^{3/2}}{\tau}e^{\frac{32\alpha r^2\tau^2}{3(1+\alpha\sigma^2)}} \qquad &(89)
\end{aligned}
$$

Similarly, we can calculate $h'(0) = E_{f_0}\psi'_\alpha(\epsilon) = \frac{1}{(1+\alpha\sigma^2)^{3/2}}$. Note $|\psi'_\alpha(t)| \leq 1, |\psi''_\alpha(t)| \leq 1.5\sqrt{\alpha}$, by equation (22) in the proof of Theorem 1 yields

$$
\begin{aligned}
\eta_1(\delta, \alpha) &= \frac{(1 - \delta)h'(0)\tau^2 - \delta}{2\sqrt{3} \times 1.5\sqrt{\alpha}\tau} \qquad &(90) \\
&= \frac{1}{3\sqrt{3\alpha}(1 + \alpha\sigma^2)^{3/2}\tau}\left[\tau^2 - \delta(\tau^2 + (1 + \alpha\sigma^2)^{3/2})\right]. \qquad &(91)
\end{aligned}
$$

According to equation (58) in the proof of Theorem 3, we have with high probability, all stationary points of the empirical risk function $\hat{L}_n(\theta)$ in (17) are inside the ball $B_2^p(\theta_0, r_s)$,

41

where

$$r_s = \eta_0 + \frac{12 C_\pi \tau \sqrt{(s_0 \log p)/n} + 2\tau \delta L_\psi}{(1-\delta)\frac{3}{4} H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma} \tag{92}$$

$$= (1+2\tau)\eta_0 + \frac{16 C_\pi \tau \sqrt{(s_0 \log p)/n}}{(1-\delta) H(\frac{8\tau r}{3}\sqrt{\frac{c_2}{\gamma}})\tau\gamma}. \tag{93}$$

Therefore, as $n >> s_0 \log p$, we have $r_s \approx (1+2\tau)\eta_0$, which completes the proof. $\square$

# References

Alfons, A., Croux, C., and Gelper, S. (2013). Sparse least trimmed squares regression for analyzing high-dimensional large data sets. *The Annals of Applied Statistics*, 7(1):226–248.

Anandkumar, A., Ge, R., Hsu, D., Kakade, S. M., and Telgarsky, M. (2014). Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research*, 15(1):2773–2832.

Andrews, D. F., Bickel, P. J., Hampel, F. R., Huber, P. J., Rogers, W. H., and W.Tukey, J. (1972). *Robust Estimates of Location: Survey and Advances*. Princeton University Press.

Bai, Z., Rao, C. R., and Wu, Y. (1992). M-estimation of multivariate linear regression parameters under a convex discrepancy function. *Statistica Sinica*, 2(1):237–254.

Beaton, A. E. and Tukey, J. W. (1974). The fitting of power series, meaning polynomials, illustrated on band-spectroscopic data. *Technometrics*, 16(2):147–185.

Bengio, Y. (2009). Learning deep architectures for ai. *Foundations and trends® in Machine Learning*, 2(1):1–127.

Boucheron, S., Lugosi, G., and Massart, P. (2013). *Concentration inequalities: A nonasymptotic theory of independence.* Oxford university press.

Brooks, T., Pope, S., and Marcolini, M. (2014). Uci machine learning repository.

Candes, E. J., Li, X., and Soltanolkotabi, M. (2015). Phase retrieval from coded diffraction patterns. *Applied and Computational Harmonic Analysis*, 39(2):277–299.

Chang, L., Roberts, S., and Welsh, A. (2018). Robust lasso regression using tukey's biweight criterion. *Technometrics*, 60(1):36–47.

Chen, M., Gao, C., and Ren, Z. (2016). A general decision theory for hubers $\epsilon$-contamination model. *Electronic Journal of Statistics*, 10(2):3752–3774.

Cheng, G., Huang, J. Z., et al. (2010). Bootstrap consistency for general semiparametric m-estimation. *The Annals of Statistics*, 38(5):2884–2915.

Dennis Jr, J. E. and Welsch, R. E. (1978). Techniques for nonlinear least squares and robust regression. *Communications in Statistics-Simulation and Computation*, 7(4):345–359.

Donoho, D. L. and Huber, P. J. (1983). The notion of breakdown point. *A festschrift for Erich L. Lehmann*, pages 157–184.

El Karoui, N., Bean, D., Bickel, P. J., Lim, C., and Yu, B. (2013). On robust regression with high-dimensional predictors. *Proceedings of the National Academy of Sciences*, 110(36):14557–14562.

Ferrari, D. and Yang, Y. (2010). Maximum lq-likelihood estimation. *The Annals of Statistics*, 38(2):753–783.

Geyer, C. J. et al. (1994). On the asymptotics of constrained $m$-estimation. *The Annals of Statistics*, 22(4):1993–2010.

Hampel, F. R., Ronchetti, E. M., Rousseeuw, P. J., and Stahel, W. A. (2011). *Robust statistics: the approach based on influence functions.* John Wiley & Sons.

Huber, P. J. (1964). Robust estimation of a location parameter. *The annals of mathematical statistics*, 35(1):73–101.

Huber, P. J. and Ronchetti, E. (2009). *Robust statistics.* New York: Wiley.

Lambert-Lacroix, S. and Zwald, L. (2011). Robust regression through the hubers criterion and adaptive lasso penalty. *Electronic Journal of Statistics*, 5:1015–1053.

Lehmann, E. L. and Casella, G. (2006). *Theory of point estimation.* Springer Science & Business Media.

Li, G., Peng, H., and Zhu, L. (2011). Nonconcave penalized m-estimation with a diverging number of parameters. *Statistica Sinica*, 21:391–419.

Loh, P.-L. (2017). Statistical consistency and asymptotic normality for high-dimensional robust $m$-estimators. *The Annals of Statistics*, 45(2):866–896.

Loh, P.-L. and Wainwright, M. J. (2015). Regularized m-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16:559–616.

Mairal, J., Bach, F., Ponce, J., and Sapiro, G. (2009). Online dictionary learning for sparse coding. In *Proceedings of the 26th annual international conference on machine learning*, pages 689–696. ACM.

Maronna, R. A. and Yohai, V. J. (1981). Asymptotic behavior of general m-estimates for regression and scale with random carriers. *Zeitschrift für Wahrscheinlichkeitstheorie und verwandte Gebiete*, 58(1):7–20.

Mei, S., Bai, Y., Montanari, A., et al. (2018). The landscape of empirical risk for nonconvex losses. *The Annals of Statistics*, 46(6A):2747–2774.

Mizera, I. and Müller, C. H. (1999). Breakdown points and variation exponents of robust m-estimators in linear models. *The Annals of Statistics*, 27(4):1164–1177.

Qin, Y. and Priebe, C. E. (2017). Robust hypothesis testing via lq-likelihood. *Statistica Sinica*, 27(4):1793–1813.

Rey, W. J. (2012). *Introduction to robust and quasi-robust statistical methods*. Springer Science & Business Media.

Wang, X., Jiang, Y., Huang, M., and Zhang, H. (2013). Robust variable selection with exponential squared loss. *Journal of the American Statistical Association*, 108(502):632–643.

Zou, H. and Hastie, T. (2005). Regularization and variable selection via the elastic net. *Journal of the royal statistical society: series B (statistical methodology)*, 67(2):301–320.