

# On the Convergence of IRLS and Its Variants in Outlier-Robust Estimation: Main Paper + Supplementary Materials

Liangzu Peng  
Johns Hopkins University  
lpeng25@jhu.edu

Christian Kümmerle  
UNC Charlotte  
kummerle@uncc.edu

René Vidal  
University of Pennsylvania  
vidalr@seas.upenn.edu

## Abstract

*Outlier-robust estimation involves estimating some parameters (e.g., 3D rotations) from data samples in the presence of outliers, and is typically formulated as a non-convex and non-smooth problem. For this problem, the classical method called iteratively reweighted least-squares (IRLS) and its variants have shown impressive performance. This paper makes several contributions towards understanding why these algorithms work so well. First, we incorporate majorization and graduated non-convexity (GNC) into the IRLS framework and prove that the resulting IRLS variant is a convergent method for outlier-robust estimation. Moreover, in the robust regression context with a constant fraction of outliers, we prove this IRLS variant converges to the ground truth at a global linear and local quadratic rate for a random Gaussian feature matrix with high probability. Experiments corroborate our theory and show that the proposed IRLS variant converges within 5-10 iterations for typical problem instances of outlier-robust estimation, while state-of-the-art methods need at least 30 iterations. A basic implementation of our method is provided: <https://github.com/liangzu/IRLS-CVPR2023>*

*... attempts to analyze this difficulty [caused by infinite weights of IRLS for the  $\ell_p$ -loss] have a long history of proofs and counterexamples to incorrect claims.*

Khurram Aftab & Richard Hartley [1]

## 1. Introduction

### 1.1. The Outlier-Robust Estimation Problem

Many parameter estimation problems can be stated in the following general form. We are given some function  $r : \mathcal{C} \times \mathcal{D} \rightarrow [0, \infty)$ , called the *residual* function. Here  $\mathcal{D}$  is the domain of data samples  $d_1, \dots, d_m$ , and  $\mathcal{C} \subset \mathbb{R}^n$  is the *constraint set* where our (ground truth) *variable*  $v^*$  lies;  $\mathcal{C}$  can be convex such as an affine subspace, or non-convex such as a special orthogonal group  $\text{SO}(3)$ . We aim to recover  $v^*$  from data  $d_i$ 's. A simple example is *linear*

*regression*, where a sample  $d_i = (\mathbf{a}_i, y_i)$  consists of a feature vector  $\mathbf{a}_i \in \mathbb{R}^n$  and a scalar response  $y_i \in \mathbb{R}$ , and the residual function is  $r(v, d_i) := |\mathbf{a}_i^\top v - y_i|$ .

The sample  $d_i$  is called an *inlier*, if  $r(v^*, d_i) \approx 0$ . It is called an *outlier*, if the residual  $r(v^*, d_i)$  is *large* (vaguely speaking). If all samples are inliers, one usually prefers solving the following problem as a means to estimate  $v^*$ :

$$\min_{v \in \mathcal{C}} \sum_{i=1}^m r(v, d_i)^2 \quad (1)$$

Problem (1) is called *least-squares*, and is known since Legendre [47] and Gauss [30] in the linear regression context. Even before that, Boscovich [13] suggested to minimize (1) without the square. This unsquared version is called *least absolute deviation*, and is more robust to outliers than (1).

Consider the following formulation for *outlier-robust estimation* (i.e., a specific type of *M-estimators* [39, 61]):

$$\min_{v \in \mathcal{C}} \sum_{i=1}^m \rho(r(v, d_i)) \quad (2)$$

Here  $\rho : \mathbb{R} \rightarrow \mathbb{R}$  is some outlier-robust loss (the unsquared version of (1) corresponds to  $\rho(r) = |r|$  in (2)). Among many possible losses  $\rho$  [20, 24], we discuss two particular choices. The first is the  $\ell_p$ -loss  $\rho(r) = |r|^p/p$ ,  $p \in (0, 1]$ ; it has been used in several research fields, e.g., *geometric vision* [1, 20], *compressed sensing* [18, 23, 41], *matrix recovery* [42, 50, 51], and *subspace clustering* [27]. The other loss is due to Huber [39]:  $\rho(r) = \min\{r^2, c^2\}$ , with  $c > 0$  a hyper-parameter; it has later been named as *Talwar* [22, 54], *Huber-type skipped mean* [32], *truncated quadratic* [6, 11], and *truncated least-squares* (TLS) [4, 78, 86]. Both losses are highly robust to outliers but make solving (2) difficult, e.g., the objective of (2) becomes non-smooth or non-convex. This motivates the need to develop efficient and provably correct solvers for (2) with either of the two losses.

### 1.2. IRLS and Its Variants in Vision & Optimization

**The General Principle of IRLS.** As its name suggests, *iteratively reweighted least-squares* (IRLS) is a general algorithmic paradigm that alternates between defining a weight

for each sample and solving a weighted least squares problem. Specifically, IRLS initializes a variable  $v^{(0)} \in \mathcal{C}$ , and, for  $t = 0, 1, \dots$ , alternates between the following two steps:

Update weights  $w_i^{(t+1)}$  based on  $v^{(t)}$ ,  $\forall i = 1, \dots, m$  (3)

$$\text{Solve: } v^{(t+1)} \leftarrow \underset{v \in \mathcal{C}}{\operatorname{argmin}} \sum_{i=1}^m w_i^{(t+1)} r(v, d_i)^2 \quad (4)$$

This basic idea dates back to the seminal work of Weiszfeld [76]; see [9, 34] for some historical accounts. A well-known and general rule for the weight update is (cf. [1, 22, 53])

$$w_i^{(t+1)} \leftarrow \rho'(r_i^{(t)})/r_i^{(t)}, \quad r_i^{(t)} := r(v^{(t)}, d_i). \quad (5)$$

In a nutshell, the rationale behind rule (5) is to “connect” weighted least-squares (4) to outlier-robust estimation (2), allowing IRLS to optimize the latter (2). Indeed, [1] shows that IRLS with the weight update in (5) results in a non-increasing objective (2). Moreover, [1] gives conditions under which IRLS with (5) converges to a stationary point of (2). This confirms that one can apply IRLS to problem (2), as long as one can solve weighted least-squares<sup>1</sup> (4).

However, the conditions of the theorem of [1] are hard to verify, e.g., one condition requires the minimizer of (4) to be a continuous function of weights  $w_i^{(t+1)}$ . Moreover, as [1] commented, directly applying (5) to non-smooth or non-convex losses (e.g.,  $\ell_p$  or TLS) might create *significant theoretical and practical difficulties*, e.g., (5) is undefined at non-differentiable points. This suggests that rule (5) needs to be improved if the  $\ell_p$  or TLS loss is to be minimized.

**IRLS in A Tale of Two Losses.** For the non-smooth  $\ell_p$ -loss, (5) results in<sup>2</sup>  $w_i^{(t+1)} \leftarrow (r_i^{(t)})^{p-2}$ , which tends to infinity as  $r_i^{(t)} \rightarrow 0$ . A workaround is to truncate the residual by some positive number  $\epsilon$ , i.e.,  $w_i^{(t+1)} \leftarrow \max\{r_i^{(t)}, \epsilon\}^{p-2}$  [20, 25–27, 48, 49, 71]. While [1, 66] considered this to be “an *ad-hoc procedure*”, in the optimization literature, there do exist theoretical guarantees for IRLS with this revised weight update to converge, at least for some specific residual functions  $r$ , see, e.g., [8, 17, 48, 49].

For the non-smooth and non-convex TLS loss  $\rho(r) = \min\{r^2, c^2\}$ , (5) results in<sup>2</sup> a *hard thresholding* scheme: set  $w_i^{(t+1)}$  to 1 if  $r_i^{(t)} \leq c$ , or set it to 0 otherwise. IRLS fails with such a weight update if the outlier rate exceeds 10% for *category-level perception* as reported in [64]. This could be remedied in two ways, discussed next.

The first is to adopt a different hard thresholding method [10] from the optimization literature, which sets  $w_i^{(t+1)}$  to 1 if  $r_i^{(t)}$  is among the  $s$ -smallest of all residuals ( $s$  is a hyperparameter), or set  $w_i^{(t+1)}$  to 0 otherwise; this method is

<sup>1</sup>While solving weighted least-squares (4) can be hard, many solvers for geometric vision exist, see, e.g., [2, 5, 15, 37, 38, 55, 62, 64, 78, 84].

<sup>2</sup>Pretending that the  $\ell_p$  or TLS losses are differentiable everywhere.

robust up to 50% outliers for *robust regression*, and converges globally linearly under some conditions. Note that this IRLS variant is not meant to minimize the TLS loss.

The second remedy manifests itself if one applies rule (5) to some *smoothing* approximation  $\rho_\mu(r)$  of the TLS loss  $\rho(r) = \min\{r^2, c^2\}$ . The approximation of [12] is

$$\rho_\mu(r) = \begin{cases} r^2, & \text{if } r^2 \leq \frac{\mu c^2}{\mu+1}, \\ c^2, & \text{if } r^2 \geq \frac{\mu+1}{\mu} c^2, \\ 2c|r|\sqrt{\mu(\mu+1)} - \mu(c^2 + r^2), & \text{o/w.} \end{cases} \quad (6)$$

Since  $\rho_\mu \rightarrow \rho$  as  $\mu \rightarrow \infty$ , a natural strategy, called *graduated non-convexity* (GNC) [12], is to alternate between optimizing  $\rho_\mu$  and increasing  $\mu$  at each iteration  $t$ . The method used for increasing  $\mu$  is called a GNC *schedule* and the default schedule has been a *linear* one, i.e.,  $\mu^{(t+1)} \leftarrow \gamma \mu^{(t)}$  with some hyper-parameter  $\gamma > 1$  [45, 52, 64, 67, 78, 86]. For example, the GNC-TLS method [12, 78] incorporates this linear schedule within the IRLS framework (3)-(5) to approximate the TLS loss via  $\rho_\mu$ .

However, the great engineering intuition of [12] and its follow-up works [4, 45, 69, 78, 87] on GNC comes with the lack of theoretical guarantees, thus [80, 82] refer to GNC as a “*fast heuristic*” strategy. On the other hand, in the optimization literature, similar GNC twists for the  $\ell_p$ -loss have been empirically investigated [19, 73, 77] for compressed sensing and related problems, and empowered with global linear or local superlinear convergence rates [23, 41, 52, 58].

For outlier-robust estimation (2), either with general [4, 45, 78] or specific residual functions [27, 58], either with the  $\ell_p$  [27, 52, 58], TLS [4, 45, 64, 78], or even other losses [66, 84], combining IRLS and GNC has pushed the empirical performance to a certain limit, which other types of methods (e.g., RANSAC [29]) can hardly attain given the same time budget. On the other hand, theoretical guarantees for IRLS offered in the optimization literature are limited to specific problems (e.g., compressed sensing [23]), and, though related, cannot be applied directly to outlier-robust estimation (2). An intriguing but under-explored theoretical question is why IRLS, GNC, and the like work so well for outlier-robust estimation (2)—can we extend, not just apply, the insights from optimization to answer this question?

### 1.3. Our Contribution

We present an IRLS variant called GNC-IRLS<sub>p</sub> (Algorithm 1) for the outlier-robust estimation problem (2) and establish general convergence properties for general constraint sets  $\mathcal{C}$ , providing a well-founded framework for empirically successful GNC methods. We further elucidate how appropriately chosen update rules for the smoothing parameter  $\epsilon^{(t)}$  (Line 7) of GNC-IRLS<sub>p</sub> lead to a global and fast local convergence for outlier-robust estimation problems. More specifically, our contributions are as follows:

---

**Algorithm 1:** GNC-IRLS<sub>p</sub>

---

- 1 Input: data  $d_1, \dots, d_m$ ,  $p \in (0, 1]$ ;
  - 2 Let  $v^{(0)} \in \mathcal{C}$  with  $\|v^{(0)}\|_2 < \infty$  and  $\epsilon^{(0)} \in (0, \infty)$ ;
  - 3 For  $t \leftarrow 0, 1, 2, \dots$ :
  - 4   Compute the residual  $r_i^{(t)} \leftarrow r(v^{(t)}, d_i), \forall i$ ;
  - 5    $w_i^{(t+1)} \leftarrow \max\{r_i^{(t)}, \epsilon^{(t)}\}^{p-2}$ ;                      // Sec. 2
  - 6   Solve problem (4) and get  $v^{(t+1)}$ ;                      // Sec. 2
  - 7   Calculate  $\epsilon^{(t+1)}$  based on a GNC schedule; // Sec. 3
- 

- In Section 2, we consider outlier-robust estimation (2) for a general class of residual functions and constraints, and we prove that GNC-IRLS<sub>p</sub> converges to stationary points of (some majorizer of) the  $\ell_p$ -loss under suitable assumptions (Theorem 1). Moreover, the assumptions are easy to verify and satisfied by many geometric vision problems (see the appendix). This challenges the viewpoint of [1, 66] that truncating the residual (Line 4, Algorithm 1) is “*ad-hoc*”. Our proof is enabled by a *majorization* interpretation of GNC-IRLS<sub>p</sub>, and is motivated by [23, 51, 60]. As we will discuss, our result is more general than those of [23, 51, 60].

- In Section 3, we propose a *superlinear* GNC schedule for GNC-IRLS<sub>p</sub>, as opposed to a *linear* one. We prove that GNC-IRLS<sub>p</sub> with such a schedule converges to the ground truth at a global linear and local superlinear rate, with high probability (Theorem 2). Moreover, GNC-IRLS<sub>p</sub> provably enjoys quadratic rates starting from the first iteration. A theoretical drawback of this powerful result is that it has a “burn-in” period and is limited to the *robust regression* setting; this is harmless though when it comes to practical use. Our proof is motivated by [52]. Their result holds only for  $p = 1$ , and our contribution lies not only in overcoming the non-convexity for the case of  $p < 1$ , but in leveraging the non-convexity to obtain a faster convergence rate.

- In Section 4 we compare the performance of GNC-TLS and GNC-IRLS<sub>p</sub> for *point cloud registration*. GNC-IRLS<sub>p</sub> terminates in 10 iterations while GNC-TLS takes 30. This is because GNC-IRLS<sub>p</sub> uses a superlinear GNC schedule, while GNC-TLS uses a linear schedule.

- In Section 5, we endow the TLS loss with a *majorization* strategy and a *superlinear* GNC schedule, leading to an IRLS method that we call MS-GNC-TLS. With majorization we prove MS-GNC-TLS converges, which challenges the viewpoint of [80, 82] that GNC is “*heuristic*”. With the superlinear schedule, MS-GNC-TLS converges, say, at iteration 6, whereas GNC-TLS does so only at iteration 30.

## 2. GNC-IRLS<sub>p</sub>: Interpretation & Convergence

In this section we show that GNC-IRLS<sub>p</sub> is a convergent method, each iteration making steady progress towards

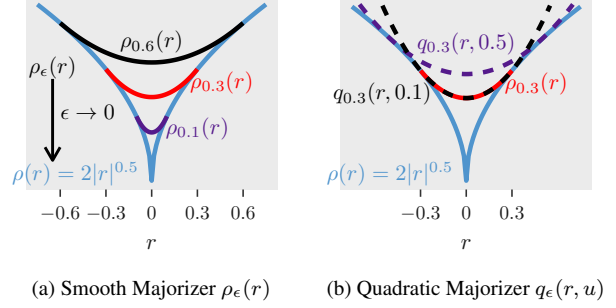


Figure 1. Two majorizers of  $\rho(r) = 2|r|^{0.5}$ ,  $\rho_\epsilon$  (7) and  $q_\epsilon$  (8).

minimizing (some majorizer of) the  $\ell_p$ -loss. We first show GNC-IRLS<sub>p</sub> involves *two-level majorization* (Section 2.1). Then we state our convergence result (Section 2.2).

### 2.1. Interpretation of GNC-IRLS<sub>p</sub>

For two functions  $f$  and  $g$  defined on  $\mathbb{R}$ , if  $f(r) \geq g(r)$  ( $\forall r \in \mathbb{R}$ ), we say  $f$  majorizes  $g$  or  $f$  is a majorizer of  $g$ . Behind the apparent alternating nature of GNC-IRLS<sub>p</sub>, it involves *two-level majorization*, as signified by the *smooth majorizer* and *quadratic majorizer*, introduced next.

**Smooth Majorizer.** As the main player in the first level of majorization, we define the *smooth majorizer*  $\rho_\epsilon : \mathbb{R} \rightarrow \mathbb{R}_{\geq 0}$  for each  $\epsilon > 0$  [58, 68] such that

$$\rho_\epsilon(r) = \begin{cases} \frac{1}{p}|r|^p, & |r| > \epsilon, \\ \frac{1}{2} \frac{r^2}{\epsilon^{2-p}} + (\frac{1}{p} - \frac{1}{2})\epsilon^p, & |r| \leq \epsilon. \end{cases} \quad (7)$$

The smooth majorizer  $\rho_\epsilon$  is a Huber-like loss [39] which coincides with the  $\ell_p$ -loss if  $|r| \geq \epsilon$  and is otherwise quadratic in  $r$ . Figure 1a shows that  $\rho_\epsilon$  majorizes the  $\ell_p$ -loss for  $p = 0.5$  and different values of  $\epsilon$ . More formally, we have:

**Lemma 1** ( $\rho_\epsilon(\cdot)$  is Smooth  $\ell_p$ -Majorizer). *For  $\rho(r) = \frac{1}{p}|r|^p$  and  $\rho_\epsilon(r)$  defined in (7), the following holds: (i)  $\rho_\epsilon(\cdot)$  is continuously differentiable, (ii)  $\rho(r) \leq \rho_\epsilon(r), \forall r \in \mathbb{R}$ , (iii)  $\epsilon' \leq \epsilon \Rightarrow \rho_{\epsilon'}(r) \leq \rho_\epsilon(r)$ , (iv)  $\rho(r) = \lim_{\epsilon \rightarrow 0} \rho_\epsilon(r)$ .*

*Remark 1 (Rethink Weight Update).* The weight update of Algorithm 1 coincides with rule (5) with  $\rho = \rho_{\epsilon^{(t)}}$ .

*Remark 2 (GNC for the  $\ell_p$ -Loss).* Lemma 1 prompts a GNC strategy of minimizing  $\rho_\epsilon$  (7) or even the  $\ell_p$ -loss: decrease  $\epsilon^{(t)}$  at each iteration  $t$  (Line 7, Algorithm 1).

**Quadratic Majorizer.** The smooth majorizer (7) is non-convex, and directly minimizing it can be hard. This is why the second level of majorization comes into play; the quadratic majorizer is the following quadratic function  $q_\epsilon$ :

$$q_\epsilon(r, u) = \rho_\epsilon(u) + \frac{1}{2} \cdot \frac{r^2 - u^2}{\max\{|u|, \epsilon\}^{2-p}}. \quad (8)$$

Note that  $q_\epsilon(r, u)$  is a shifted version of  $\rho_\epsilon(u)$  by a carefully chosen amount, which makes  $q_\epsilon(\cdot, u)$  into a majorizer of  $\rho_\epsilon(\cdot)$ . Indeed, Figure 1b shows that  $q_{0.3}(\cdot, u)$  majorizes  $\rho_{0.3}(\cdot)$  for  $u = 0.1$  and  $0.5$ . More formally, we have:

**Lemma 2** ( $q_\epsilon(\cdot, u)$  is Quadratic  $\ell_p$ -Majorizer). *With  $\rho(r) = \frac{1}{p}|r|^p$ ,  $\rho_\epsilon(r)$  and  $q_\epsilon(r, u)$  defined respectively in (7) and (8), we have  $\rho_\epsilon(u) = q_\epsilon(u, u)$  and  $\rho_\epsilon(r) \leq q_\epsilon(r, u)$ ,  $\forall r, u \in \mathbb{R}$ .*

**Remark 3** (Rethink Weighted Least-Squares). Recall  $r_i^{(t)} := r(v^{(t)}, d_i)$ . The WLS step (4) of Algorithm 1 minimizes the quadratic majorizer  $\sum_{i=1}^m q_{\epsilon^{(t)}}(r(\cdot, d_i), r_i^{(t)})$ :

$$\begin{aligned} v^{(t+1)} &\in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m \frac{r(v, d_i)^2}{\max\{|r_i^{(t)}|, \epsilon^{(t)}\}^{2-p}} \\ &= \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m q_{\epsilon^{(t)}}(r(v, d_i), r_i^{(t)}) \end{aligned}$$

GNC-IRLS<sub>p</sub> differs from the *majorization-minimization* paradigm [28, 60, 70] in that, at different iterations, GNC-IRLS<sub>p</sub> minimizes different quadratic majorizers, as controlled by the smoothing parameter  $\epsilon^{(t)}$ ; in so doing, it blends (quadratic) majorization-minimization with GNC.

## 2.2. Convergence of GNC-IRLS<sub>p</sub>

To obtain convergence results, we need appropriate assumptions on the constraint set  $\mathcal{C}$  and residual function  $r$ . The first assumption is standard (cf. [3, Section 4.2]):

**Assumption 1.**  $\mathcal{C}$  is non-empty and closed. The residual function  $r(v, d) : \mathcal{C} \times \mathcal{D} \rightarrow \mathbb{R}_{\geq 0}$  is *weakly coercive* in  $v$ :

$$\text{Either } \mathcal{C} \text{ is bounded or } \lim_{v \in \mathcal{C}, \|v\|_2 \rightarrow \infty} r(v, d) \rightarrow \infty. \quad (9)$$

Moreover, if  $\|v\|_2 \neq \infty$  then  $r(v, d) \neq \infty$ .

The next assumption is about *differentiability*:

**Assumption 2.** The residual function  $r(v, d)$  is continuous in  $v$  everywhere, and differentiable in  $v$  if  $r(v, d) \neq 0$ . Moreover,  $r(v, d)^2$  is continuously differentiable in  $v$ .

Assumptions 1 and 2 are mild and easy to verify. With these assumptions, we prove the following:

**Theorem 1** (Convergence of GNC-IRLS<sub>p</sub>). *Let  $\{v^{(t)}\}_t$  be the iterates of Algorithm 1 with  $\epsilon^{(t)}$  non-increasing and  $\epsilon := \lim_{t \rightarrow \infty} \epsilon^{(t)} > 0$ . Under Assumptions 1 and 2, every accumulation point of  $\{v^{(t)}\}_t$  is a stationary point<sup>3</sup> of*

$$\min_{v \in \mathcal{C}} \sum_{i=1}^m \rho_\epsilon(r(v, d_i)). \quad (10)$$

<sup>3</sup>Stationary points are in the sense of [3, Section 5.3]; they satisfy a certain geometric condition that every local minimizer of (10) fulfills.

With a GNC schedule that creates a non-increasing sequence  $\{\epsilon^{(t)}\}_t$  convergent to  $\epsilon$ , GNC-IRLS<sub>p</sub> finds a stationary point of  $\rho_\epsilon$  (Theorem 1), and  $\rho_\epsilon$  approximates the  $\ell_p$ -loss very well if  $\epsilon$  is small (Lemma 1, Figure 1a). The convergence statement of “*accumulation points are stationary points*” in Theorem 1 is standard, and similar results can be found in optimization papers on IRLS or majorization-minimization, e.g., [23, Thm 5.3 (ii)], [51, Thm 3.2], [60, Thm 1], [68, Thm 11], [53, Proposition 5], [48, Thm 1]. However, to our knowledge, Theorem 1 is the only result that holds for a *general* constraint set  $\mathcal{C}$  and for minimizing *a sequence of majorizers* within the GNC framework.

Theorem 1 is proved by combing ideas of [23, 51] and [60], while generalizing their results. Unlike in Theorem 1,  $\mathcal{C}$  is assumed to be convex and  $\epsilon^{(t)} = \epsilon$  for all  $t$  in [60]. In [23, 51],  $\mathcal{C}$  is defined by linear equality constraints and the residual function  $r$  is very specific, unlike in Theorem 1. Finally, as reviewed in Section 1.2, the result of [1] requires a condition that is hard to verify and their result does not apply to IRLS with the GNC strategy.

While stationary points are not necessarily local minimizers<sup>3</sup>, convergence to them is perhaps the best one could guarantee in the setting where the objective (7) and constraint set  $\mathcal{C}$  can both be non-convex. That said, a stronger convergence theory is possible given more assumptions on the problem and data. We will explore this in Section 3.

## 3. Convergence Rates for Robust Regression

While Theorem 1 is general, it does not reveal any convergence *speed*. Here, we compromise on generality and prove that GNC-IRLS<sub>p</sub> converges rapidly for *robust regression* [54]. Consider the following problem setup:

**Problem 1** (Robust Regression). For a *feature matrix*  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_m]^\top \in \mathbb{R}^{m \times n}$  and a *response vector*  $\mathbf{y} = [y_1, \dots, y_m]^\top \in \mathbb{R}^m$ , assume there is a ground truth vector  $v^* = \mathbf{x}^* \in \mathbb{R}^n$  such that the residual vector  $\mathbf{A}\mathbf{x}^* - \mathbf{y}$  has  $k$  non-zero entries; i.e., there are  $k$  outliers and  $m - k$  inliers among data  $\{d_i\}_{i=1}^m = \{(\mathbf{a}_i, y_i)\}_{i=1}^m$ . The goal of robust regression is to recover  $\mathbf{x}^*$  from data  $\mathbf{A}$  and  $\mathbf{y}$ .

In Problem 1 we assume all inliers  $(\mathbf{a}_i, y_i)$  are noiseless, i.e.,  $r(v^*, d_i) = |\mathbf{a}_i^\top \mathbf{x}^* - y_i| = 0$ . The extension to the noisy case is not hard (cf. [52, Thm 2], [41, Thm A.1]).

The GNC schedule is closely related to the convergence rates of IRLS. Informally, [52] suggests that the *linear* GNC schedule (as is commonly seen) leads to a *linear* rate. However, it is possible for IRLS to attain *superlinear* rates. In particular, defining the *superlinear GNC schedule*

$$\epsilon^{(t+1)} \leftarrow \beta(\epsilon^{(t)})^{2-p}, \quad \beta > 0, \quad (11)$$

we prove the following result:

**Theorem 2.** Assume  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0, 1)$  entries. Initialize Algorithm 1 at  $\mathbf{x}^{(0)}$  and  $\epsilon^{(0)} > 0$  such that  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$ . Denote by  $r_{\min+}^*$  the smallest non-zero number among the set of residuals  $\{\|\mathbf{a}_i^\top \mathbf{x}^* - y_i\}\}_{i=1}^m$ . Define

$$\alpha := \frac{\sqrt{5} \cdot 2^{2-p}}{0.99 \cdot 0.516} \cdot \frac{1}{(r_{\min+}^*)^{1-p}} \cdot \frac{\sqrt{k} \cdot (1.01\sqrt{k} + \sqrt{n})}{(m-k)}. \quad (12)$$

Then the iterates  $\{\mathbf{x}^{(t)}\}_{t \geq 0}$  produced by GNC-IRLS<sub>p</sub> with  $p \in [0, 1]^4$  and the GNC schedule (11) with  $\beta \geq \alpha$  satisfy

$$\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \begin{cases} \beta^t \cdot \epsilon^{(0)} & p = 1 \\ \beta^{\frac{(2-p)t-1}{1-p}} \cdot (\epsilon^{(0)})^{(2-p)^t} & p \in [0, 1) \end{cases} \quad (13)$$

with probability at least  $1 - (P_0 + P_1 + P_2)$ , where

$$\begin{aligned} P_0 &:= \exp(-\tilde{\Omega}(n)), & P_1 &:= \exp(-\tilde{\Omega}(k-n)), \\ P_2 &:= \exp(-\tilde{\Omega}(m-k-n \log n)). \end{aligned} \quad (14)$$

We discuss several aspects of Theorem 2: (i) the probabilities (14), (ii) the condition  $\beta \geq \alpha$  (12), (iii) its relation to prior works, and (iv) the interaction of the GNC schedule (11), error bound (13) and condition  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$ .

(i) In the probability terms of (14),  $\tilde{\Omega}$  stands for the standard big- $\Omega$  notation, with the difference that  $\tilde{\Omega}$  also suppresses logarithmic terms. We wish  $P_0, P_1, P_2$  of (14) to be small so that (13) holds with high probability. This is true whenever  $n$  is large ( $P_0$ ),  $k \gg n$  ( $P_1$ ), and  $m-k \gg n \log n$  ( $P_2$ ). It seems counterintuitive to ask for the number  $k$  of outliers to be far larger than  $n$ , but the challenging case of Problem 1 occurs exactly when  $k$  is large. If  $k$  were small, then an alternative proof would give  $P_1 = \exp(-\tilde{\Omega}(n))$ . Such proof is much simpler, which is why we omit it.

(ii) We wish  $\alpha$  to be as small as possible as this would make it easier to set the factor  $\beta$  in the GNC schedule (11). Ignoring the values of the constants in (12),  $\alpha$  mainly involves two terms,  $r_{\min+}^*$  and  $O((k + \sqrt{kn})/(m-k))$ . Since  $r_{\min+}^*$  measures the minimum residual of outliers at the ground truth  $\mathbf{x}^*$ , we expect it to be a large constant. Since  $p \in [0, 1]$ , a large  $(r_{\min+}^*)^{1-p}$  would make  $\alpha$  small; on the other hand, for  $p = 1$ ,  $\alpha$  does not depend on  $r_{\min+}^*$  at all. Then note that we require a large  $k$  in (i), but this might make the second term  $O((k + \sqrt{kn})/(m-k))$  and therefore  $\alpha$  very large. The rescue is in the denominator:  $\alpha$  is small if the number  $m-k$  of inliers (or the inlier rate) is large.

(iii) Theorem 2 is motivated by [52, Thm 1], over which we make some improvements. First, the GNC schedule of [52] sets  $\epsilon^{(t+1)} \leftarrow \beta \epsilon^{(t)}$  if  $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2 \leq 2\beta \epsilon^{(t)}$ , or otherwise  $\epsilon^{(t+1)} \leftarrow \epsilon^{(t)}$ . We simplify and generalize it into (11). Also, [52] is limited to the case  $p = 1$ , but Theorem 2 holds for any  $p \in [0, 1]$ ; we derive some technical lemmas that overcome the challenges of the non-convex case  $p < 1$ .

The final point (iv) has more delicate interpretations and ramifications, and we discuss it in Sections 3.1-3.4 next.

<sup>4</sup>It is valid to run GNC-IRLS<sub>p</sub> with  $p = 0$ , as we justified in [58].

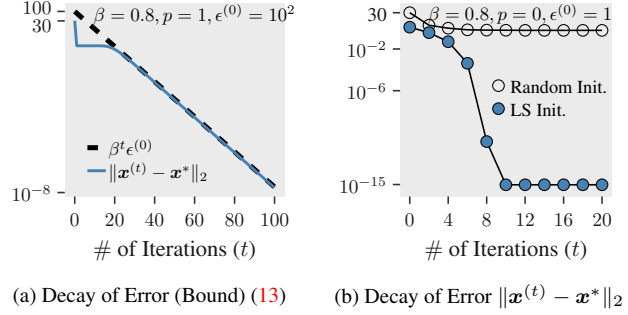


Figure 2. (2a, Section 3.1): Error bound  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \beta^t \epsilon^{(0)}$  (13) with initialization  $\mathbf{x}^{(0)} \sim \mathcal{N}(0, 100\mathbf{I}_n)$ . (2b, Section 3.2): Errors of GNC-IRLS<sub>0</sub> at each iteration with least-squares versus random initialization. 100 trials,  $k = 400$ ,  $m = 1000$ ,  $n = 10$ .

### 3.1. Global Linear Convergence at $p = 1$

For the error bound  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \beta^t \cdot \epsilon^{(0)}$  of (13) to make sense, one needs to set  $\beta < 1$ , then the condition  $\alpha \leq \beta$  in Theorem 2 implies  $\alpha < 1$ . As discussed, we have  $\alpha < 1$  if the inlier rate is large. Indeed, assuming  $k, m \gg n$  and bringing now the constant of (12) into the picture, we see that  $\alpha < 1$  amounts to  $m-k > 2.02\sqrt{5}/(0.99 \times 0.516)k$ . This defines an outlier rate below which Theorem 2 holds. This also implies Theorem 2 is optimal in an information-theoretical sense (e.g., it only requires  $m$  to be linear in  $k$ ).

For  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \beta^t \cdot \epsilon^{(0)}$  to be true, Theorem 2 requires  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$  (among other assumptions). Given any initialization  $\mathbf{x}^{(0)}$ , one can choose a large  $\epsilon^{(0)}$  such that  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$ , so [52] claimed this is a global linear convergence. But this claim is imprecise, e.g., if  $\mathbf{x}^{(0)}$  is the least-squares initialization and  $\epsilon^{(0)}$  is larger than all residuals  $|\mathbf{a}_i^\top \mathbf{x}^{(0)} - y_i|$ , then all weights  $w_i^{(1)}$  are equal to  $\epsilon^{(0)}$ , and we would get  $\mathbf{x}^{(1)} = \mathbf{x}^{(0)}$ . As such, the error would not decrease until  $\epsilon^{(t)}$  becomes smaller: Figure 2a shows that  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2$  “waits” for almost 20 iterations to decay together with the bound  $\beta^t \epsilon^{(0)}$  ( $\epsilon^{(0)} = 100$ ). This overlooked phenomenon caused by large  $\epsilon^{(0)}$  is what we call a burn-in period. Interestingly, the burn-in period does not mean that our bound (13) is incorrect, but just that it might be loose for large  $\epsilon^{(0)}$  in early iterations.

Figure 2a shows that GNC-IRLS<sub>1</sub> needs more than 100 iterations to reach machine accuracy. We improve this next, by considering  $p \in [0, 1)$  (Sections 3.2-3.4).

### 3.2. Local Quadratic Convergence at $p = 0$

Theorem 2 with  $p \in [0, 1)$  is better elaborated in the case  $p = 0$ , for which (13) gives  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq (\beta \epsilon^{(0)})^{2^t} / \beta$ . This corresponds to a quadratic convergence rate. Again, the error bound  $(\beta \epsilon^{(0)})^{2^t} / \beta$  only makes sense if  $\beta \epsilon^{(0)} < 1$ , or if we set  $\epsilon^{(0)}$  small (note that this time we do not require  $\beta < 1$ ). In turn, Theorem 2 would demand an initialization

$\mathbf{x}^{(0)}$  such that  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$ . As corroborated by Figure 2b, GNC-IRLS<sub>0</sub> with random (“bad”) initialization and small  $\epsilon^{(0)}$  fails, but the least-squares initialization seems to be good enough, allowing GNC-IRLS<sub>0</sub> to converge at a quadratic rate, within 10 iterations, where “the number of correct digits doubles at each iteration” [14, Section 9.5.3]. Powerful as it might seem, quadratic (and superlinear) convergence is doomed to be local and in general cannot hold for all initializations (*cf.* Newton’s method); we refer the reader to our prior work [58] for different insights into the quadratic rates of IRLS for robust regression.

We believe the *next best* convergence guarantees are these two: (i) prove that some IRLS variant has two-phase convergence, first global linear and then local quadratic, (ii) derive a suitable choice of  $\epsilon^{(0)}$ ,  $\beta$ , and  $\mathbf{x}^{(0)}$  for which quadratic convergence happens starting from the first iteration. We discuss these next in Section 3.3 and 3.4.

### 3.3. Graduated Rates From Linear to Quadratic?

Consider the following slight twist over Algorithm 1:

- With some initialization  $\mathbf{x}^{(0)}$  and a sufficiently large  $\epsilon^{(0)}$  such that  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$ , run Algorithm 1 with  $p = 1$  and GNC schedule (11), until  $\beta^t \epsilon^{(0)} < 1$ . Theorem 2 suggests that  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \beta^t \epsilon^{(0)}$ .
- Re-run Algorithm 1 with  $\mathbf{x}^{(0)} := \mathbf{x}^{(t)}$ ,  $\epsilon^{(0)} := \beta^t \epsilon^{(0)}$ ,  $p = 0$ , and schedule (11). Quadratic convergence (13) of Theorem 2 is now meaningful, since  $\beta \epsilon^{(0)} < 1$ .

Simply put, the above twist switches from  $p = 1$  to  $p = 0$  if  $\beta^t \epsilon^{(0)} < 1$ , resulting in a graduated rate from global linear to local quadratic. Such a graduated rate guarantee seems rare; we can only find it in [21]. Figure 3a shows that when we switch to  $p = 0$ , the convergence ensues in the next 10 iterations. A deficiency is that this twist also comes with a burn-in period (*cf.* Section 3.1), after which it is possible that the linear convergence phase is skipped and the quadratic convergence takes place directly (Figure 3a).

### 3.4. Quadratic Rates From The First Iteration?

The IRLS twist of Section 3.3 can take 30 iterations to converge if  $\epsilon^{(0)}$  is large (Figure 3a). But we also saw that, with the least-squares initialization and  $\beta \epsilon^{(0)} = 0.8 < 1$ , GNC-IRLS<sub>0</sub> converges within 10 iterations, at a quadratic rate (Figure 2b). We now argue that it is theoretically possible for GNC-IRLS<sub>0</sub> to have quadratic rates *starting from as early as the first iteration*. For this, we first prove:

**Proposition 1.** *Assume  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0, 1)$  entries with  $m \geq n$ . Let  $\mathbf{x}^\dagger := (\mathbf{A}^\top \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{y}$ . With probability at least  $1 - \exp(-\Omega(k)) - \exp(-\Omega(m))$ , we have*

$$\|\mathbf{x}^\dagger - \mathbf{x}^*\|_2 \leq \frac{(1.01\sqrt{k} + \sqrt{n}) \cdot \|\mathbf{A}\mathbf{x}^* - \mathbf{y}\|_2}{(0.99\sqrt{m} - \sqrt{n})^2}. \quad (15)$$

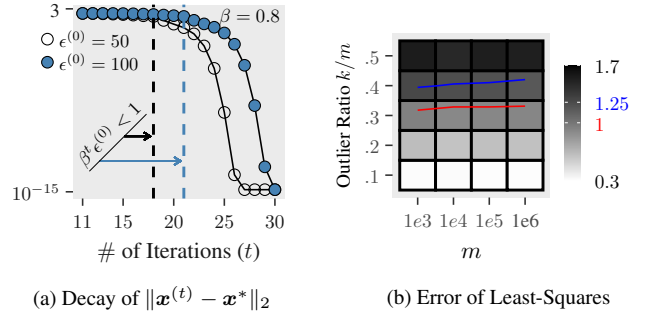


Figure 3. (3a, Section 3.3): From linear to quadratic rates, vertical lines indicating the transition takes place;  $k = 400$ ,  $m = 1000$ ,  $\mathbf{x}^{(0)} \sim \mathcal{N}(0, \mathbf{I}_n)$ . (3b, Section 3.4): Error  $\|\mathbf{x}^\dagger - \mathbf{x}^*\|_2$  of the least-squares estimator  $\mathbf{x}^\dagger$ . We set 100 trials,  $n = 10$ .

We wish  $\|\mathbf{x}^\dagger - \mathbf{x}^*\|_2 \leq 1$ ; if so we can set  $\epsilon^{(0)} = 1$  and  $\beta < 1$ , achieving quadratic rates with initialization  $\mathbf{x}^\dagger$  (Theorem 2). This is possible if  $k/m$  is small and  $m, k \gg n$ ; see (15). This is also empirically confirmed in Figure 3b, where  $\|\mathbf{x}^\dagger - \mathbf{x}^*\|_2 \leq 1$  for fewer than 30% outliers.

**Implementation Details.** The discussions so far suggest the following implementation of GNC-IRLS<sub>p</sub>. Set  $\epsilon^{(0)} = 1$ ,  $p = 0$ . Initialize it via least-squares. Set  $\beta$  smaller than 1; we always use  $\beta = 0.8$ . Theorem 1 suggests to let  $\{\epsilon^{(t)}\}_t$  converge to some  $\epsilon > 0$ . In the noiseless case, we set  $\epsilon = 10^{-16}$ . Otherwise, if we are given an *inlier threshold*  $c$  such that  $r(v^*, d_i) \leq c$  for all inliers  $d_i$ , then we set  $\epsilon = c$ .

## 4. Experiments: Lp Versus TLS

Here we compare GNC-TLS [12, 78] and GNC-IRLS<sub>p</sub>. For more extensive experiments of IRLS and its variants, see, *e.g.*, [1, 20, 25, 27, 46, 64, 65, 69, 78].

**Experimental Setup.** We contextualize our experiment in the application of *point cloud registration*. In this application, each sample  $d_i$  is a 3D point pair  $(\mathbf{y}_i, \mathbf{x}_i)$ , the variable  $v$  consists of a 3D rotation  $\mathbf{R}$  and translation  $\mathbf{t}$ , and the residual function is  $r(v; d_i) = \|\mathbf{y}_i - \mathbf{R}\mathbf{x}_i - \mathbf{t}\|_2$ . The corresponding weighted least-squares problem (4) is solved by eliminating the translation first and then applying SVD [38].

**Data.** We randomly sample  $k$  outlier point pairs  $(\mathbf{y}_j, \mathbf{x}_j)$  so that  $\mathbf{y}_j \sim \mathcal{N}(0, \mathbf{I}_3)$  and  $\mathbf{x}_j \sim \mathcal{N}(0, \mathbf{I}_3)$ ; here  $\mathbf{I}_3$  denotes the  $3 \times 3$  identity matrix. To get  $m - k$  inlier pairs  $(\mathbf{y}_i, \mathbf{x}_i)$ , we randomly sample  $\mathbf{x}_i$  from  $\mathcal{N}(0, \mathbf{I}_3)$  and compute  $\mathbf{y}_i = \mathbf{R}^* \mathbf{x}_i + \mathbf{t}^* + \epsilon_i$ . Here,  $\mathbf{R}^*$  and  $\mathbf{t}^*$  are randomly generated ground truth rotation and translation respectively, and  $\epsilon_i \sim \mathcal{N}(0, 0.01^2 \mathbf{I}_3)$  is some Gaussian noise. We set  $c^2 = 0.01^2 \times 5.54^2$ , so each inlier  $(\mathbf{y}_i, \mathbf{x}_i)$  satisfies  $\|\mathbf{y}_i - \mathbf{R}^* \mathbf{x}_i - \mathbf{t}^*\|_2 \leq c^2$  with probability  $\geq 1 - 10^{-6}$ .

**Metric.** Given a rotation  $\mathbf{R}$ , translation  $\mathbf{t}$ , and ground truth inlier index set  $\mathcal{I}^*$ , we can calculate the *average inlier resid-*

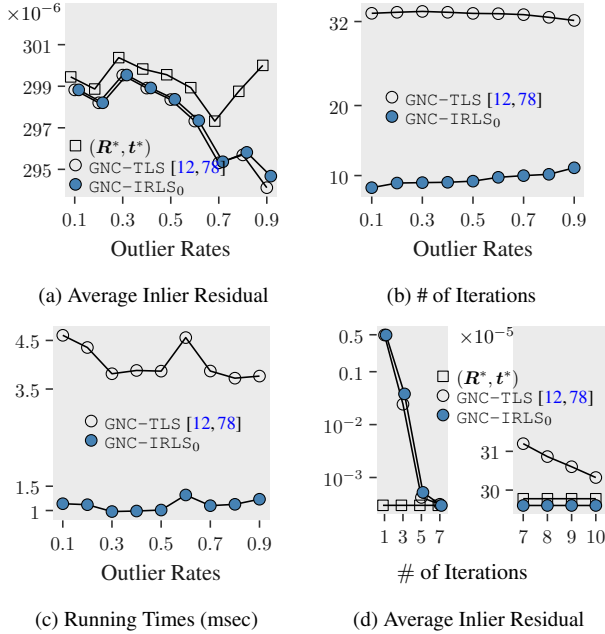


Figure 4. Comparison of GNC-IRLS<sub>0</sub> and GNC-TLS for point cloud registration. 100 trials,  $m = 1000$ . 4d: 900/1000 outliers.

ual  $\sum_{i \in \mathcal{I}^*} \|\mathbf{y}_i - \mathbf{R}\mathbf{x}_i - \mathbf{t}\|_2 / (m - k)$ . This is used to measure the errors made by the algorithms to evaluate.

**Results.** As the outlier rate varies from 10% to 90%, GNC-IRLS<sub>0</sub> and GNC-TLS entail almost the same average inlier residual (Figure 4a). Their errors are even smaller than those at the ground truth  $(\mathbf{R}^*, \mathbf{t}^*)$ , which suggests that the performance of both algorithms cannot be further improved for such experiments. But note that they could fail for more than 900/1000 outliers (which was reported in prior works, so we did not provide a plot here) and that the breakdown points will change for different data distributions and different geometric problems.

What can actually be improved is the convergence rate: GNC-IRLS<sub>0</sub> terminates in 10 iterations, while GNC-TLS takes 32 (Figure 4b), indicating that GNC-IRLS<sub>0</sub> is 3 times faster (Figure 4c). For fair<sup>5</sup> comparison, both methods are implemented to terminate under the same condition, that is whenever the difference of the minimum values of (4) between two consecutive iterations is smaller than  $10^{-10}$  (other thresholds, e.g.,  $10^{-6}$ ,  $10^{-16}$ , lead to similar results).

The errors of GNC-TLS decrease as fast as GNC-IRLS<sub>0</sub> (Figure 4d, left). At first glance, this seems counterintuitive because GNC-TLS comes with a linear GNC schedule (cf. Section 1.2) and is thus expected to converge linearly (cf. [52]), as opposed to the quadratic rate of GNC-IRLS<sub>0</sub> (cf. Theorem 2). With hindsight, this might be a natural consequence of the weighting strategy of GNC-TLS (cf. [78, Eq.

<sup>5</sup>It is slightly unfair to GNC-IRLS<sub>0</sub> as its weights are typically larger.

(14)], (5), (6)): Weight 0 is set if the residual is *particularly large*, and this could completely rule out some *obvious* outliers at early iterations (and similarly for *particularly small* residuals), resulting in a fast decrease of errors. But this weighting scheme brings diminishing gains in later iterations, where the errors of GNC-TLS decrease only linearly (Figure 4d, right). The final observation is that GNC-IRLS<sub>0</sub> reaches an error smaller than that of  $(\mathbf{R}^*, \mathbf{t}^*)$  at iteration 7, but it requires a few more iterations to terminate (similarly for GNC-TLS). This implies the termination criterion is sub-optimal (it is hard to design a provably better one).

Finally, Figures 4b and 4d show that GNC-TLS has an error of  $\leq 10^{-3}$  already at iteration 10, but, unnecessarily, it terminates at iteration 32. We improve this in Section 5, without even changing the termination criterion.

## 5. MS-GNC-TLS: Improving GNC-TLS

... it indicates that GNC can fail, and that there is therefore no point in looking for a general proof of correctness.

Andrew Blake & Andrew Zisserman [12]

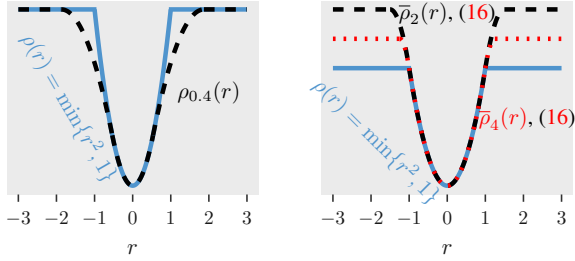
In this section, we improve GNC-TLS [12, 78] from two aspects, as respectively motivated by two ideas that we have developed for the  $\ell_p$ -loss, namely *majorization* (Section 2) and *superlinear GNC schedule* (Section 3). Majorization guarantees a monotonic decrease of the objective and the eventual convergence (cf. Theorem 1), and the superlinear GNC schedule speeds up convergence (cf. Theorem 2).

**Majorization.** To motivate the need for majorizing the TLS loss  $\rho(r) = \min\{r^2, c^2\}$ , recall GNC-TLS uses  $\rho_\mu$  (6) to approximate  $\rho(\cdot)$ . The issue is that  $\rho_\mu(\cdot)$  relaxes  $\rho(\cdot)$  and approximates it *from below*, and hence  $\rho_\mu(r) \leq \rho(r), \forall \mu > 0$  (Figure 5a). This makes a convergence analysis difficult.

We propose the following smooth function

$$\bar{\rho}_\mu(r) = \begin{cases} r^2, & \text{if } |r| \leq c, \\ \frac{\mu+1}{\mu} c^2, & \text{if } |r| \geq \frac{\mu+1}{\mu} c, \\ -\mu r^2 + 2(1+\mu)c|r| - (1+\mu)c^2, & \text{o/w,} \end{cases} \quad (16)$$

to majorize the TLS loss  $\rho(r)$ ; see Figure 5b. Since both  $\rho_\mu(r)$  (6) and  $\bar{\rho}_\mu(r)$  approach  $\rho(r)$  as  $\mu \rightarrow \infty$ , one might expect comparable performance. However, a crucial difference is that, with the majorizer  $\bar{\rho}_\mu(r)$ , convergence guarantees easily ensue. Indeed,  $\bar{\rho}_\mu(r)$  is akin to the smooth majorizer (7) of the  $\ell_p$ -loss, and one could construct a quadratic majorizer for  $\bar{\rho}_\mu(r)$ , which enables an IRLS + GNC scheme (cf. Remarks 1-3, Section 1.2). In particular, this IRLS vari-



(a) Approximation  $\rho_\mu(r)$  of [12]      (b) Smooth Majorizer (16)

Figure 5. The TLS loss  $\rho(r)$  and its surrogates.

ant involves (i) weight update using (5) with  $\rho = \bar{\rho}_{\mu^{(t)}}$ , i.e.,

$$w_i^{(t+1)} = \begin{cases} 1, & \text{if } r_i^{(t)} \leq c, \\ 0, & \text{if } r_i^{(t)} \geq \frac{\mu^{(t)}+1}{\mu^{(t)}}c, \\ \frac{c(1+\mu^{(t)})}{r_i^{(t)}} - \mu^{(t)}, & \text{o/w,} \end{cases} \quad (17)$$

and (ii) updating  $\mu^{(t+1)}$  based on some GNC schedule.

We prove the following result to accompany Theorem 1.

**Theorem 3** (Convergence of majorized GNC-TLS). *Let  $\{v^{(t)}\}_t$  be the iterates of IRLS with weight update (17) and a GNC schedule  $\{\mu^{(t)}\}_t$ . Assume  $\{v^{(t)}\}_t$  is bounded, i.e.,  $\|v^{(t)}\|_2 < \infty$  ( $\forall t$ ). Suppose  $\{\mu^{(t)}\}_t$  is non-decreasing and converges to  $\mu < \infty$ . Under Assumptions 1 and 2, every accumulation point of  $\{v^{(t)}\}_t$  is a stationary point of*

$$\min_{v \in \mathcal{C}} \sum_{i=1}^m \bar{\rho}_\mu(r(v, d_i)). \quad (18)$$

**Superlinear Schedule.** Motivated by (11) (with  $p = 0$ ) and discussions in Sections 3.2-3.3, we propose the update rule

$$\mu^{(t+1)} \leftarrow \begin{cases} \gamma \sqrt{\mu^{(t)}} & \mu^{(t)} \leq 1 \\ \gamma \mu^{(t)} & \mu^{(t)} > 1 \end{cases}, \quad \gamma > 1, \quad (19)$$

as our GNC schedule. Denote by MS-GNC-TLS the resulting IRLS method that optimizes (16) with schedule (19).

The intuition behind the superlinear schedule (19) is as follows. With (19), the interval  $(c, c + c/\mu^{(t)})$  of (17) that produces non-binary weights shrinks faster than the linear schedule  $\mu^{(t+1)} \leftarrow \gamma \mu^{(t)}$  (Figure 6a), thus the superlinear schedule makes it happen earlier that all weights become binary, which is a good indicator for convergence. Note though that this argument does not prove (MS-)GNC-TLS converges, as it does not exclude the case that (MS-)GNC-TLS could produce different binary weights at consecutive iterations (cf. [10] and [4, Thm 15]).

Under the setting of Figure 4c, MS-GNC-TLS takes 6 iterations to converge (Figure 6b). It is even faster than GNC-IRLS<sub>0</sub> as it benefits from combining soft and hard

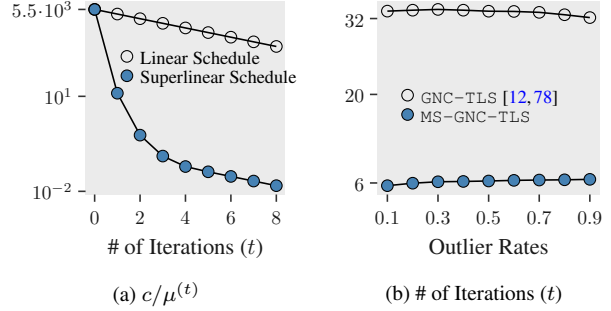


Figure 6. **6a:** Length  $c/\mu^{(t)}$  of the interval that corresponds to non-binary weights (17) with  $c = 0.0554$ ,  $\mu^{(0)} = 10^{-5}$ ,  $\gamma = 1.4$ . **6b:** Number of iterations at which the algorithms terminate.

thresholding (17). In this experiment, MS-GNC-TLS and GNC-TLS result in basically the same error upon convergence; it is just that GNC-TLS does not monotonically decrease the objective, and that its linear GNC schedule is more conservative than the proposed superlinear one.

**Implementation Details.** With the superlinear schedule (19),  $\mu^{(t)}$  increases very fast, so one could set  $\mu^{(0)} \leftarrow 10^{-15}$  such that MS-GNC-TLS can still terminate within 10 iterations. However, schedule (19) is *double-edged*: If  $\mu^{(t)}$  increases so fast that all residuals are larger than  $\frac{\mu^{(t)}+1}{\mu^{(t)}}c$ , then all weights would be zero as per (17) and MS-GNC-TLS might fail. Fortunately, this situation can be prevented if we *slow down*: replace  $\mu^{(t+1)} \leftarrow \gamma \sqrt{\mu^{(t)}}$  of (19) with  $\mu^{(t+1)} \leftarrow \gamma (\mu^{(t)})^{1/(2-p)}$  for a larger  $p \in (0, 1]$ .

## 6. Conclusion, Limitations, and Future Work

**Conclusion.** While IRLS and GNC have often been viewed as different techniques [44,45,64,78,83], we reconcile them with an emphasis on a theoretical understanding of convergence properties and their relation with GNC schedules. Two messages are (i) that a majorization strategy should be constructed for guaranteeing *convergence* (Theorems 1 and 3), and (ii) that a superlinear GNC schedule should be considered for guaranteeing *convergence rates* (Theorem 2).

**Limitations & Future Work.** IRLS and its variants would break down if the number  $m - k$  of inliers is close to the number  $n$  of variables, say if  $m - k < 3n$ . For geometric vision problems,  $n$  is small (e.g.,  $n = 6$  for point cloud registration), so IRLS might fail if, for example,  $m - k < 18$ . In fact, for small  $m$ , other methods (e.g., RANSAC [7, 29], outlier removal [56], or semidefinite relaxations [33, 57, 81]) are efficient, accurate, and are thus recommended.

A limitation of the TLS loss  $\rho(r) = \min\{r^2, c^2\}$  is the need to choose a threshold parameter  $c$ . Ideally, it should be chosen as small as possible but larger than every inlier residual; see [59] for a related discussion. Prior works [4, 69]



tried to dispense with  $c^2$ , but it was at the expense of introducing other parameters. This issue might be solved by changing  $c$  in a GNC style at each iteration, which implies future work of designing a GNC schedule for  $c$  and studying its interplay with another GNC parameter  $\mu^{(t)}$ . On the theory side, we note that extending the analysis of Theorem 2 beyond  $\ell_p$ -losses remains to be studied in future work.

**Acknowledgements.** This work was supported by grants NSF 1704458, NSF 1934979, ONR MURI 503405-78051, and the Northrop Grumman Mission Systems Research in Applications for Learning Machines (REALM) initiative.

## References

- [1] Khurram Aftab and Richard Hartley. Convergence of iteratively re-weighted least squares to robust M-estimators. In *IEEE Winter Conference on Applications of Computer Vision*, 2015. 1, 2, 3, 4, 6, 22
- [2] Chris Aholt, Sameer Agarwal, and Rekha Thomas. A QCQP approach to triangulation. In *European Conference on Computer Vision*, 2012. 2
- [3] Niclas Andréasson, Anton Evgrafov, and Michael Patriksson. *An Introduction to Continuous Optimization: Foundations and Fundamental Algorithms*. Courier Dover Publications, 2020. 4, 16, 23
- [4] Pasquale Antonante, Vasileios Tzoumas, Heng Yang, and Luca Carlone. Outlier-robust estimation: Hardness, minimally tuned algorithms, and applications. *IEEE Transactions on Robotics*, 2021. 1, 2, 8
- [5] K Somani Arun, Thomas S Huang, and Steven D Blostein. Least-squares fitting of two 3D point sets. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, (5):698–700, 1987. 2
- [6] Erik Ask, Olof Enqvist, and Fredrik Kahl. Optimal geometric fitting under the truncated  $L_2$ -norm. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2013. 1
- [7] Daniel Barath, Jana Noskova, Maksym Ivashechkin, and Jiri Matas. MAGSAC++, a fast, reliable and accurate robust estimator. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 8
- [8] Amir Beck. On the convergence of alternating minimization for convex programming with applications to iteratively reweighted least squares and decomposition schemes. *SIAM Journal on Optimization*, 25(1):185–209, 2015. 2
- [9] Amir Beck and Shoham Sabach. Weiszfeld’s method: Old and new results. *Journal of Optimization Theory and Applications*, 164(1):1–40, 2015. 2
- [10] Kush Bhatia, Prateek Jain, and Purushottam Kar. Robust regression via hard thresholding. *Advances in Neural Information Processing Systems*, 2015. 2, 8
- [11] Michael J Black and Anand Rangarajan. On the unification of line processes, outlier rejection, and robust statistics with applications in early vision. *International Journal of Computer Vision*, 19(1):57–91, 1996. 1
- [12] Andrew Blake and Andrew Zisserman. *Visual Reconstruction*. MIT Press, 1987. 2, 6, 7, 8
- [13] Roger Joseph Boscovich. De litteraria expeditione per pontificiam ditionem, et synopsis amplioris operis, ac habentur plura ejus ex exemplaria etiam sensorum impessa. *Bononiensi Scientiarum et Artum Instituto Atque Academia Commentarii*, 4:353–396, 1757. 1
- [14] Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge University Press, 2004. 6
- [15] Jesus Briales and Javier Gonzalez-Jimenez. Convex global 3D registration with lagrangian duality. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2
- [16] Lucas Brynte, Viktor Larsson, José Pedro Iglesias, Carl Olsson, and Fredrik Kahl. On the tightness of semidefinite relaxations for rotation estimation. *Journal of Mathematical Imaging and Vision*, 64(1):57–67, 2022. 14

- [17] Tony F Chan and Pep Mulet. On the convergence of the lagged diffusivity fixed point method in total variation image restoration. *SIAM Journal on Numerical Analysis*, 36(2):354–367, 1999. 2
- [18] Rick Chartrand. Exact reconstruction of sparse signals via nonconvex minimization. *IEEE Signal Processing Letters*, 14(10):707–710, 2007. 1
- [19] Rick Chartrand and Wotao Yin. Iteratively reweighted algorithms for compressive sensing. In *IEEE International Conference on Acoustics, Speech and Signal Processing*, 2008. 2
- [20] Avishek Chatterjee and Venu Madhav Govindu. Robust relative rotation averaging. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):958–972, 2017. 1, 2, 6, 22
- [21] Bintong Chen and Naihua Xiu. A global linear and local quadratic noninterior continuation method for nonlinear complementarity problems based on Chen–Mangasarian smoothing functions. *SIAM Journal on Optimization*, 9(3):605–623, 1999. 6
- [22] David Coleman, Paul Holland, Neil Kaden, Virginia Klema, and Stephen C Peters. A system of subroutines for iteratively reweighted least squares computations. *ACM Transactions on Mathematical Software*, 6(3):327–336, 1980. 1, 2
- [23] Ingrid Daubechies, Ronald DeVore, Massimo Fornasier, and C Sinan Güntürk. Iteratively reweighted least squares minimization for sparse recovery. *Communications on Pure and Applied Mathematics*, 63(1):1–38, 2010. 1, 2, 3, 4
- [24] DQF De Menezes, Diego Martinez Prata, Argimiro R Secchi, and José Carlos Pinto. A review on robust M-estimators for regression analysis. *Computers & Chemical Engineering*, 147:107254, 2021. 1
- [25] Tianjiao Ding, Yunchen Yang, Zhihui Zhu, Daniel P Robinson, René Vidal, Laurent Kneip, and Manolis C Tsakiris. Robust homography estimation via dual principal component pursuit. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 2, 6
- [26] Tianyu Ding, Zhihui Zhu, Tianjiao Ding, Yunchen Yang, René Vidal, Manolis C. Tsakiris, and Daniel Robinson. Noisy dual principal component pursuit. In *International Conference on Machine Learning*, 2019. 2, 14
- [27] Wenhua Dong, Xiao-jun Wu, and Josef Kittler. Sparse subspace clustering via smoothed  $\ell_p$  minimization. *Pattern Recognition Letters*, 125:206–211, 2019. 1, 2, 6
- [28] Taosha Fan and Todd Murphey. Majorization minimization methods for distributed pose graph optimization with convergence guarantees. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, 2020. 4
- [29] Martin A Fischler and Robert C Bolles. Random sample consensus: A paradigm for model fitting with applications to image analysis and automated cartography. *Communications of the ACM*, 24(6):381–395, 1981. 2, 8
- [30] Carl Friedrich Gauss. *Theoria motus corporum coelestium in sectionibus conicis solem ambientium auctore Carolo Fridrico Gauss*. sumtibus Frid. Perthes et IH Besser, 1809. 1
- [31] Stuart Geman and Donald E. McClure. Bayesian image analysis: An application to single photon emission tomography. In *Proceedings of the American Statistical Association*, 1985. 22
- [32] Frank R Hampel. The breakdown points of the mean combined with some rejection rules. *Technometrics*, 27(2):95–107, 1985. 1
- [33] Linus Härenstam-Nielsen, Niclas Zeller, and Daniel Cremers. Semidefinite relaxations for robust multiview triangulation. Technical report, arXiv:2301.11431 [cs.CV], 2023. 8
- [34] Richard Hartley. Tutorial notes on IRLS. <http://users.cecs.anu.edu.au/~hartley/Papers/PDF/Hartley:IRLS14.pdf>, 2014. Accessed: September 2022. 2
- [35] Richard Hartley, Jochen Trumpf, Yuchao Dai, and Hongdong Li. Rotation averaging. *International journal of computer vision*, 103(3):267–305, 2013. 14
- [36] Richard Hartley and Andrew Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, 2004. 14, 22
- [37] Joel A Hesch and Stergios I Roumeliotis. A direct least-squares (DLS) method for PnP. In *International Conference on Computer Vision*, 2011. 2, 13
- [38] Berthold KP Horn, Hugh M Hilden, and Shahriar Negahdaripour. Closed-form solution of absolute orientation using orthonormal matrices. *Journal of the Optical Society of America A*, 5(7):1127–1135, 1988. 2, 6, 13
- [39] Peter J Huber. Robust estimation of a location parameter. *The Annals of Mathematical Statistics*, 35(1):73–101, 1964. 1, 3
- [40] Laurent Kneip, Hongdong Li, and Yongduek Seo. UPnP: An optimal  $O(n)$  solution to the absolute pose problem with universal applicability. In *European Conference on Computer Vision*, 2014. 13
- [41] Christian Kümmerle, Claudio Mayrink Verdun, and Dominik Stöger. Iteratively reweighted least squares for basis pursuit with global linear convergence rate. *Advances in Neural Information Processing Systems*, 2021. 1, 2, 4
- [42] Christian Kümmerle and Claudio M Verdun. A scalable second order method for ill-conditioned matrix completion from few samples. In *International Conference on Machine Learning*, 2021. 1
- [43] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *Annals of Statistics*, pages 1302–1338, 2000. 33
- [44] Huu Le and Christopher Zach. A graduated filter method for large scale robust estimation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 8
- [45] Huu Le and Christopher Zach. Robust fitting with truncated least squares: A bilevel optimization approach. In *International Conference on 3D Vision*, 2021. 2, 8, 22
- [46] Seong Hun Lee and Javier Civera. HARA: A hierarchical approach for robust rotation averaging. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 6
- [47] Adrien Marie Legendre. *Nouvelles méthodes pour la détermination des orbites des comètes: avec un supplément contenant divers perfectionnemens de ces méthodes et leur application aux deux comètes de 1805*. Courcier, 1806. 1

- [48] Gilad Lerman and Tyler Maunu. Fast, robust and non-convex subspace recovery. *Information and Inference: A Journal of the IMA*, 7(2):277–336, 2018. 2, 4
- [49] Gilad Lerman, Michael B. McCoy, Joel A. Tropp, and Teng Zhang. Robust computation of linear models by convex relaxation. *Foundations of Computational Mathematics*, 15(2):363–410, 2015. 2
- [50] Goran Marjanovic and Victor Solo. On  $\ell_q$  optimization and matrix completion. *IEEE Transactions on Signal Processing*, 60(11):5714–5724, 2012. 1
- [51] Karthik Mohan and Maryam Fazel. Iterative reweighted algorithms for matrix rank minimization. *The Journal of Machine Learning Research*, 13(1):3441–3473, 2012. 1, 3, 4
- [52] Bhaskar Mukhoty, Govind Gopakumar, Prateek Jain, and Purushottam Kar. Globally-convergent iteratively reweighted least squares for robust regression problems. In *International Conference on Artificial Intelligence and Statistics*, 2019. 2, 3, 4, 5, 7
- [53] Peter Ochs, Alexey Dosovitskiy, Thomas Brox, and Thomas Pock. On iteratively reweighted algorithms for nonsmooth nonconvex optimization in computer vision. *SIAM Journal on Imaging Sciences*, 8(1):331–372, 2015. 2, 4
- [54] Dianne P O’Leary. Robust regression computation using iteratively reweighted least squares. *SIAM Journal on Matrix Analysis and Applications*, 11(3):466–480, 1990. 1, 4
- [55] Frank C Park and Bryan J Martin. Robot sensor calibration: Solving  $AX = XB$  on the Euclidean group. *IEEE Transactions on Robotics and Automation*, 10(5):717–721, 1994. 2
- [56] Álvaro Parra Bustos and Tat-Jun Chin. Guaranteed outlier removal for point cloud registration with correspondences. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(12):2868–2882, 2018. 8
- [57] Liangzu Peng, Mahyar Fazlyab, and René Vidal. Semidefinite relaxations of truncated least-squares in robust rotation search: Tight or not. In *European Conference on Computer Vision*, 2022. 8
- [58] Liangzu Peng, Christian Kümmerle, and René Vidal. Global linear and local superlinear convergence of IRLS for nonsmooth robust regression. In *Advances in Neural Information Processing Systems*, 2022. 2, 3, 5, 6
- [59] Liangzu Peng, Manolis C. Tsakiris, and René Vidal. ARCS: Accurate rotation and correspondences search. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022. 8
- [60] Meisam Razaviyayn, Mingyi Hong, and Zhi-Quan Luo. A unified convergence analysis of block successive minimization methods for nonsmooth optimization. *SIAM Journal on Optimization*, 23(2):1126–1153, 2013. 3, 4
- [61] William J.J. Rey. *Introduction to Robust and Quasi-Robust Statistical Methods*. Springer Science & Business Media, 1983. 1
- [62] David M Rosen, Luca Carlone, Afonso S Bandeira, and John J Leonard. SE-Sync: A certifiably correct algorithm for synchronization over the special Euclidean group. *The International Journal of Robotics Research*, 38(2-3):95–125, 2019. 2, 14
- [63] Gerald Schweighofer and Axel Pinz. Globally optimal  $O(n)$  solution to the PnP problem for general camera models. In *British Machine Vision Conference*, 2008. 13
- [64] Jingnan Shi, Heng Yang, and Luca Carlone. Optimal and robust category-level perception: Object pose and shape estimation from 2D and 3D semantic keypoints. Technical report, arXiv:2206.12498 [cs.CV], 2022. 2, 6, 8, 13
- [65] Yunpeng Shi and Gilad Lerman. Message passing least squares framework and its application to rotation synchronization. In *International Conference on Machine Learning*, 2020. 6
- [66] Chitturi Sidhartha and Venu Madhav Govindu. It is all in the weights: Robust rotation averaging revisited. In *International Conference on 3D Vision*, 2021. 2, 3
- [67] Torbjorn Smith and Olav Egeland. Dynamical pose estimation with graduated non-convexity for outlier robustness. *Modeling, Identification and Control*, 43(2):79–89, 2022. 2
- [68] Junxiao Song, Prabhu Babu, and Daniel P Palomar. Sparse generalized eigenvalue problem via smooth optimization. *IEEE Transactions on Signal Processing*, 63(7):1627–1642, 2015. 3, 4
- [69] Lei Sun. IMOT: General-purpose, fast and robust estimation for spatial perception problems with outliers. Technical report, arXiv:2204.01324v1 [cs.CV], 2022. 2, 6, 8
- [70] Ying Sun, Prabhu Babu, and Daniel P Palomar. Majorization-minimization algorithms in signal processing, communications, and machine learning. *IEEE Transactions on Signal Processing*, 65(3):794–816, 2016. 4
- [71] Manolis C. Tsakiris and René Vidal. Dual principal component pursuit. *Journal of Machine Learning Research*, 19(18):1–50, 2018. 2
- [72] Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018. 33
- [73] Sergey Voronin and Ingrid Daubechies. An iteratively reweighted least squares algorithm for sparse regularization. Technical report, arXiv:1511.08970v3 [math.NA], 2015. 2
- [74] Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019. 33
- [75] Jianqiao Wangni, Dahua Lin, Ji Liu, Kostas Daniilidis, and Jianbo Shi. Towards statistically provable geometric 3D human pose recovery. *SIAM Journal on Imaging Sciences*, 14(1):246–270, 2021. 13
- [76] Endre Weiszfeld. Sur le point pour lequel la somme des distances de n points donnés est minimum. *Tohoku Mathematical Journal*, 43:355–386, 1937. 2
- [77] David Wipf and Srikantan Nagarajan. Iterative reweighted  $\ell_1$  and  $\ell_2$  methods for finding sparse solutions. *IEEE Journal of Selected Topics in Signal Processing*, 4(2):317–329, 2010. 2
- [78] Heng Yang, Pasquale Antonante, Vasileios Tzoumas, and Luca Carlone. Graduated non-convexity for robust spatial perception: From non-minimal solvers to global outlier rejection. *IEEE Robotics and Automation Letters*, 5(2):1127–1134, 2020. 1, 2, 6, 7, 8
- [79] Heng Yang and Luca Carlone. In perfect shape: Certifiably optimal 3D shape reconstruction from 2D landmarks.

- In *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020. 13
- [80] Heng Yang and Luca Carlone. One ring to rule them all: Certifiably robust geometric perception with outliers. In *Advances in Neural Information Processing Systems*, 2020. 2, 3
- [81] Heng Yang and Luca Carlone. Certifiable outlier-robust geometric perception: Exact semidefinite relaxations and scalable global optimization. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022. 8
- [82] Heng Yang, Jingnan Shi, and Luca Carlone. TEASER: Fast and certifiable point cloud registration. *IEEE Transactions on Robotics*, 37(2):314–333, 2021. 2, 3
- [83] Christopher Zach and Guillaume Bourmaud. Descending, lifting or smoothing: Secrets of robust cost optimization. In *European Conference on Computer Vision*, 2018. 8
- [84] Ji Zhao. An efficient solution to non-minimal case essential matrix estimation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020. 2
- [85] Yinqiang Zheng, Yubin Kuang, Shigeki Sugimoto, Kalle Astrom, and Masatoshi Okutomi. Revisiting the PnP problem: A fast, general and optimal solution. In *IEEE International Conference on Computer Vision*, 2013. 13
- [86] Pengwei Zhou, Xuexun Guo, Xiaofei Pei, and Ci Chen. T-TOAM: Truncated least squares Lidar-only odometry and mapping in real time. *IEEE Transactions on Geoscience and Remote Sensing*, 60:1–13, 2021. 1, 2
- [87] Qian-Yi Zhou, Jaesik Park, and Vladlen Koltun. Fast global registration. In *European Conference on Computer Vision*, 2016. 2
- [88] Xiaowei Zhou, Spyridon Leonardos, Xiaoyan Hu, and Kostas Daniilidis. 3D shape reconstruction from 2D landmarks: A convex formulation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2015. 13
- [89] Michael Zollhöfer, Matthias Nießner, Shahram Izadi, Christoph Rehmann, Christopher Zach, Matthew Fisher, Chenglei Wu, Andrew Fitzgibbon, Charles Loop, Christian Theobalt, et al. Real-time non-rigid reconstruction using an RGB-D camera. *ACM Transactions on Graphics*, 33(4):1–12, 2014. 22

## A. Structure of The Appendix

We structure the appendix as follows:

- In Section B, we present some geometric vision problems whose residuals and constraints satisfy Assumptions 1 and 2.
- In Section C, we prove Lemma 1, Lemma 2, and Theorem 1 of the main paper.
- In Section D, we discuss the design of smooth and quadratic majorizers for the TLS loss, and prove Theorem 3.
- In Section E, we extend our convergence theory for the  $\ell_p$  and TLS losses to a general class of loss functions. The result is Theorem 4, which can be applied to general classes of loss functions, residual functions, and constraints. We also list some losses in Example 1. As a side product, this section sheds light on how the weight update rule (5) shows up.
- In Section F, we prove Theorem 2 of the main paper. For this, we invoke Theorem 5 of Section G multiple times.
- In Section H, we prove Proposition 1 of the main paper.
- In Section I, we reference several results from high-dimensional statistics.

## B. Geometric Vision Problems Satisfy Assumptions 1 and 2

Here, we list several geometric vision problems that satisfy Assumptions 1 and 2, as one can verify very easily. Therefore, Theorems 1 and 3 of the main paper (and the more general Theorem 4 in Section E) can be applied to these problems, yielding convergence guarantees for the proposed IRLS variants (GNC-IRLS<sub>p</sub> and MS-GNC-TLS).

**Point Cloud Registration.** In the point cloud registration problem [38] that we discussed in the main paper, the variable  $v$  lies in the special Euclidean group  $\text{SE}(3) = \{(\mathbf{R}, \mathbf{t}) : \mathbf{R} \in \text{SO}(3), \mathbf{t} \in \mathbb{R}^3\}$ , each sample  $d_i$  is a point pair  $(\mathbf{y}_i, \mathbf{x}_i) \in \mathbb{R}^3 \times \mathbb{R}^3$ , and the residual function  $r(v, d_i)$  is

$$r(v, d_i) = \|\mathbf{y}_i - \mathbf{R}\mathbf{x}_i - \mathbf{t}\|_2.$$

One can easily check that Assumptions 1 and 2 are satisfied.

In the more general and more challenging problem of *category-level registration*, point  $\mathbf{x}_i$  is not given directly. We are instead given  $\mathbf{b}_{i1}, \dots, \mathbf{b}_{is} \in \mathbb{R}^3$  such that  $\mathbf{x}_i$  can be represented as a convex combination of  $\mathbf{b}_{ij}$ 's, i.e.,  $\mathbf{x}_i = \sum_{j=1}^s c_j \mathbf{b}_{ij}$  for some unknown non-negative coefficients  $c_j$ 's with  $\sum_{j=1}^s c_j = 1$ . This setting is under the assumption that the 3D point cloud object  $\{\mathbf{x}_i\}_{i=1}^m$  (e.g., a car) can be represented as a convex combination of  $s$  objects  $\{\{\mathbf{b}_{ij}\}_{i=1}^m\}_{j=1}^s$  of the same category (e.g.,  $s$  different cars). As such, our variable consists of a 3D rotation  $\mathbf{R}$ , translation  $\mathbf{t}$ , and coefficients  $c_j$ 's. Therefore, the constraint set  $\mathcal{C}$  becomes  $\text{SE}(3) \times \{(c_1, \dots, c_s) \in \mathbb{R}^s : c_j \geq 0, \forall j, c_1 + \dots + c_s = 1\}$ , and the residual function  $r(v, d_i)$  is

$$r(v, d_i) = \|\mathbf{y}_i - \mathbf{R}(\sum_{j=1}^s c_j \mathbf{b}_{ij}) - \mathbf{t}\|_2.$$

Both for this problem and a similar one where  $\mathbf{R}(\sum_{j=1}^s c_j \mathbf{b}_{ij}) + \mathbf{t}$  of the residual function  $r(v, d_i)$  is replaced by its 2D projection (cf. [64, 75, 79, 88]), Assumptions 1 and 2 are satisfied.

**Absolute Pose Estimation.** The *absolute pose estimation* problem has a few different formulations [37, 40, 63, 85]. In one of the formulations, the variable  $v$  lies in  $\text{SE}(3)$  and consists of a 3D rotation  $\mathbf{R} \in \text{SO}(3)$  and translation  $\mathbf{t} \in \mathbb{R}^3$ , each data sample  $d_i = (\mathbf{u}_i, \mathbf{x}_i) \in \mathbb{S}^2 \times \mathbb{R}^3$  consists of a unit vector  $\mathbf{u}_i$  and a 3D point  $\mathbf{x}_i$ , and the residual function is

$$r(v, d_i) = \|[\mathbf{u}_i]_{\times}(\mathbf{R}\mathbf{x}_i + \mathbf{t})\|_2,$$

where  $[\cdot]_{\times}$  denotes the cross product matrix of some 3D vector. One verifies that Assumptions 1 and 2 are satisfied.

**Essential Matrix Estimation.** In the problem of *essential matrix estimation*, our variable  $v$  lies in  $\overline{\text{SE}}(3) := \{(\mathbf{R}, \mathbf{t}/\|\mathbf{t}\|_2) : (\mathbf{R}, \mathbf{t}) \in \text{SE}(3)\}$ . Therefore, the translation variable is assumed to be a unit vector. The data sample  $d_i$  is a 2D point pair  $(\mathbf{y}_i, \mathbf{x}_i) \in \mathbb{R}^2 \times \mathbb{R}^2$  representing pixel locations in two different images respectively. The residual function in this case is

$$r(v, d_i) = \left| [\mathbf{y}_i^{\top} \ 1] ([\mathbf{t}]_{\times} \mathbf{R}) \begin{bmatrix} \mathbf{x}_i \\ 1 \end{bmatrix} \right|.$$

The matrix  $[t]_{\times} \mathbf{R}$  is called an essential matrix—thus the name of the problem. One easily verifies that Assumptions 1 and 2 are satisfied. And one further verifies the assumptions are true for related geometric vision problems (e.g., *fundamental matrix estimation, homography estimation, trifocal tensor estimation*). See, e.g., [36] for a summary of these problems.

**Hand-Eye Calibration.** In this *hand-eye calibration* problem (cf. [16, Example 3]), one aims to find the relative transformation  $v = \mathbf{R} \in \text{SO}(3)$  between a robot hand and a camera. Each data sample  $d_i$  consists of a pair of rotations  $(\mathbf{U}_i, \mathbf{V}_i) \in \text{SO}(3) \times \text{SO}(3)$ , each of which represents the location of the robot hand or the camera respectively in some global coordinate system. The residual function for this problem is

$$r(v, d_i) = \|\mathbf{U}_i \mathbf{R} - \mathbf{R} \mathbf{V}_i\|_{\text{F}}.$$

One can then verify that Assumptions 1 and 2 are satisfied.

**Rotation Averaging.** In this problem [35], our variable  $v$  consists of  $s$  3D rotations  $\mathbf{R}_1, \dots, \mathbf{R}_s$ . And we are given  $s(s+1)/2$  samples  $\{\tilde{\mathbf{R}}_{ij}\}_{1 \leq i < j \leq s} \subset \text{SO}(3)$ , and each sample  $d_{ij} = \tilde{\mathbf{R}}_{ij}$  is assumed to satisfy  $\tilde{\mathbf{R}}_{ij} \approx \mathbf{R}_i^{\top} \mathbf{R}_j$  if it is an inlier. Therefore, a residual function in this problem is

$$r(v, d_{ij}) = \|\mathbf{R}_i \tilde{\mathbf{R}}_{ij} - \mathbf{R}_j\|_{\text{F}}.$$

One verifies Assumptions 1 and 2 are satisfied. One further verifies that the assumptions also hold for a more general problem called SE(3) synchronization [62, Eq. (12)].

**Affine Hyperplane Recovery.** In this problem [26, Section 5], each data sample  $d$  is an  $n$ -dimensional point  $\mathbf{x} \in \mathbb{R}^n$ , the variable  $v = (\mathbf{u}, t) \in \mathbb{S}^{n-1} \times \mathbb{R}$  consists of some vector  $\mathbf{u} \in \mathbb{S}^{n-1}$  and a scale translation  $t \in \mathbb{R}$ . The variable  $(\mathbf{u}, t)$  defines an affine hyperplane in  $\mathbb{R}^n$ , and in the 3D case, it could represent some road in a real point cloud. The residual function is

$$r(v; d) = |\mathbf{x}_i^{\top} \mathbf{u} - t|.$$

It is easy to verify that Assumptions 1 and 2 are satisfied.

## C. Proof of Theorem 1

*Proof of Lemma 1.* To show statement (i) of Lemma 1, recall the definition of (7) and we note that the  $\rho_{\epsilon}(\cdot)$  is differentiable for all  $|\rho| \neq \epsilon$  and its derivative satisfies

$$\rho'_{\epsilon}(r) = \begin{cases} \frac{|r|^p}{r}, & \text{if } |r| > \epsilon, \\ \frac{r}{\epsilon^{2-p}}, & \text{if } |r| \leq \epsilon. \end{cases}$$

To check continuous differentiability of  $\rho_{\epsilon}(\cdot)$  at  $r = \epsilon$ , we compute

$$\lim_{h>0, h \rightarrow 0} \frac{\rho_{\epsilon}(\epsilon + h) - \rho_{\epsilon}(\epsilon)}{h} = \lim_{h>0, h \rightarrow 0} \frac{\frac{1}{p}(\epsilon + h)^p - \frac{1}{p}\epsilon^p}{h} = \epsilon^{p-1}$$

and

$$\lim_{h<0, h \rightarrow 0} \frac{\rho_{\epsilon}(\epsilon + h) - \rho_{\epsilon}(\epsilon)}{h} = \lim_{h>0, h \rightarrow 0} \frac{\frac{1}{2} \frac{(\epsilon+h)^2}{\epsilon^{2-p}} + \left(\frac{1}{p} - \frac{1}{2}\right)\epsilon^p - \frac{1}{p}\epsilon^p}{h} = \epsilon^{p-2} \lim_{h>0, h \rightarrow 0} \frac{\frac{1}{2}(\epsilon + h)^2 - \frac{1}{2}\epsilon^2}{h} = \epsilon^{p-1},$$

which shows that

$$\rho'_{\epsilon}(\epsilon) = \epsilon^{p-1} = \lim_{r \rightarrow \epsilon} \rho'_{\epsilon}(r), \tag{20}$$

and analogously, it can be shown that  $\rho'_{\epsilon}$  exists and is continuous at  $r = -\epsilon$ , which shows (i).

For (ii), due to  $0 < p \leq 1$ , we know that  $z \rightarrow z^{p/2}$  is concave, and therefore

$$\rho(r) = \frac{1}{p}|r|^p \leq \frac{1}{p}(r^2)^{p/2} \leq \frac{1}{p} \left[ (\epsilon^2)^{\frac{p}{2}} + \frac{p}{2}(\epsilon^2)^{\frac{p}{2}-1}(r^2 - \epsilon^2) \right] = \frac{1}{2} \frac{r^2}{\epsilon^{2-p}} + \left( \frac{1}{p} - \frac{1}{2} \right) \epsilon^p = \rho_{\epsilon}(r)$$

for  $|r| \leq \epsilon$ , and  $\rho(r) = \rho_{\epsilon}(r)$  for  $|r| > \epsilon$ , which shows (ii).

For (iii), we can define the function  $g_r(\epsilon) := \rho_\epsilon(r)$ , which is differentiable for  $\epsilon \geq |r|$ : Computing its derivative, we obtain

$$g_r'(\epsilon) = \left(\frac{1}{p} - \frac{1}{2}\right) p\epsilon^{p-1} + (p-2)\frac{1}{2}\frac{r^2}{\epsilon^{3-p}} = \epsilon^{p-1} \left[ \left(1 - \frac{r^2}{\epsilon^2}\right) + \frac{p}{2}\frac{r^2}{\epsilon^2} \right],$$

which is  $\geq 0$  for all  $\epsilon \geq |r|$ , so we conclude that  $g_r(\cdot)$  is monotonously increasing in this interval and  $\rho_{\epsilon'}(r) \leq \rho_\epsilon(r)$  if  $\epsilon' \leq \epsilon$  for all  $\epsilon \geq |r|$ , and  $\rho_{\epsilon'}(r) = \rho_\epsilon(r)$  if  $|r| \leq \epsilon' \leq \epsilon$ . If  $\epsilon' < |r| \leq \epsilon$ ,  $\rho_{\epsilon'}(r) = \rho(r) \leq \rho_\epsilon(r)$ , which finishes the proof of (iii).

(iv) finally follows as for each  $r \neq 0$ ,  $\rho_\epsilon(r)$  coincides with  $\rho(r)$  for  $\epsilon$  small enough.  $\square$

*Proof of Lemma 2.* Lemma 2 follows from the fact that  $\rho_\epsilon(\sqrt{r})$  is concave in  $r$ , and therefore, it is majorized by its tangent line at any given point  $u \in \mathbb{R}$ . More formally, it follows from (34) below with  $\theta = u$  and the fact that  $\frac{\rho_\epsilon'(u)}{|u|} = \frac{1}{\max(|u|, \epsilon)^{2-p}}$  for the smoothed  $\ell_p$ -loss  $\rho_\epsilon(\cdot)$ .  $\square$

*Proof of Theorem 1.* Recall the definition  $r_i^{(t)} := r(v^{(t)}, d_i), \forall t$ . We have

$$\begin{aligned} \sum_{i=1}^m \rho_{\epsilon^{(t+1)}}(r_i^{(t+1)}) &\stackrel{(i)}{\leq} \sum_{i=1}^m \rho_{\epsilon^{(t)}}(r_i^{(t+1)}) \\ &\stackrel{(ii)}{\leq} \sum_{i=1}^m q_{\epsilon^{(t)}}(r_i^{(t+1)}, r_i^{(t)}) \\ &\stackrel{(iii)}{\leq} \sum_{i=1}^m q_{\epsilon^{(t)}}(r_i^{(t)}, r_i^{(t)}) \\ &= \sum_{i=1}^m \rho_{\epsilon^{(t)}}(r_i^{(t)}). \end{aligned} \tag{21}$$

The three inequalities in (21) are due to (i) Lemma 1, (ii) Lemma 2, and (iii) Remark 3, respectively. Hence  $\{\sum_{i=1}^m \rho_{\epsilon^{(t)}}(r_i^{(t)})\}_t$  is a non-increasing sequence of non-negative numbers, thus convergent. We now show that  $\{v^{(t)}\}_t$  is bounded. This is true if  $\mathcal{C}$  is bounded, so we assume that  $\mathcal{C}$  is unbounded. Suppose for the sake of contradiction that  $\|v^{(t)}\|_2 \rightarrow \infty$  as  $t \rightarrow \infty$ . By Assumption 1, for every  $i = 1, \dots, m$  and as  $t \rightarrow \infty$ , we have  $r_i^{(t)} = r(v^{(t)}, d_i) \rightarrow \infty$  and thus  $\rho_{\epsilon^{(t)}}(r_i^{(t)}) \rightarrow \infty$ . Therefore

$$\sum_{i=1}^m \rho_{\epsilon^{(0)}}(r_i^{(0)}) \geq \lim_{t \rightarrow \infty} \sum_{i=1}^m \rho_{\epsilon^{(t)}}(r_i^{(t)}) = \infty. \tag{22}$$

But Algorithm 1 ensures  $\|v^{(0)}\|_2 < \infty$ , thus  $r_i^{(0)}$  is finite (by Assumption 1), and so is  $\rho_{\epsilon^{(0)}}(r_i^{(0)})$ , contradicting (22). We have thus proved that  $\{v^{(t)}\}_t$  is a bounded sequence. It now follows from the famous *Bolzano–Weierstrass* theorem that this bounded sequence  $\{v^{(t)}\}_t$  has a convergent subsequence. Let  $\{v^{(t_j)}\}_j$  be such subsequence that converges to some accumulation point, say  $\hat{v}$ . We will next show that  $\hat{v}$  is a stationary point of (10). In fact, we have

$$\begin{aligned} \sum_{i=1}^m q_{\epsilon^{(t_j+1)}}(r_i^{(t_j+1)}, r_i^{(t_j+1)}) &= \sum_{i=1}^m \rho_{\epsilon^{(t_j+1)}}(r_i^{(t_j+1)}) \\ &\stackrel{(i)}{\leq} \sum_{i=1}^m \rho_{\epsilon^{(t_j+1)}}(r_i^{(t_j+1)}) \\ &\stackrel{(ii)}{\leq} \sum_{i=1}^m \rho_{\epsilon^{(t_j)}}(r_i^{(t_j+1)}) \\ &\stackrel{(iii)}{\leq} \sum_{i=1}^m q_{\epsilon^{(t_j)}}(r_i^{(t_j+1)}, r_i^{(t_j)}) \\ &\stackrel{(iv)}{\leq} \sum_{i=1}^m q_{\epsilon^{(t_j)}}(r(v, d_i), r_i^{(t_j)}), \quad \forall v \in \mathcal{C}. \end{aligned} \tag{23}$$

The above inequalities are due respectively to (i)  $t_{j+1} \geq t_j + 1$  and (21), (ii) Lemma 1, (iii) Lemma 2, (iv) Remark 3. Note that  $\rho_{\epsilon(t)}$ ,  $q_{\epsilon(t)}$  are continuous for any  $t$ , and  $r(v, d_i)$  is continuous in  $v$  by Assumption 2; with  $t_j \rightarrow \infty$  in (23) we get

$$\sum_{i=1}^m q_{\epsilon} \left( r(\hat{v}, d_i), r(\hat{v}, d_i) \right) \leq \sum_{i=1}^m q_{\epsilon} \left( r(v, d_i), r(\hat{v}, d_i) \right), \quad \forall v \in \mathcal{C}.$$

With  $\hat{w}_i := \max\{r(\hat{v}, d_i), \epsilon\}^{2-p}$  and Remark 3, the above implies

$$\hat{v} \in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m q_{\epsilon} \left( r(v, d_i), r(\hat{v}, d_i) \right) = \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m \hat{w}_i \cdot r(v, d_i)^2.$$

Thus,  $\hat{v}$  must satisfy the geometric optimality condition [3, Section 5.3]

$$\left( \sum_{i=1}^m \hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2) \right)^{\top} \mathbf{b} \geq 0, \quad \forall \mathbf{b} \in T_{\hat{v}}\mathcal{C}, \quad (24)$$

where  $T_{\hat{v}}\mathcal{C}$  denotes the tangent cone of  $\mathcal{C}$  at  $\hat{v}$ .

We now verify  $\nabla \rho_{\epsilon}(r(\hat{v}, d_i)) = \frac{1}{2} \hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2)$ . Indeed, if  $r(\hat{v}, d_i) = 0$ , then by the definition of  $\rho_{\epsilon}$  (7) we have

$$\nabla \rho_{\epsilon}(r(\hat{v}, d_i)) = \nabla \left( \frac{1}{2} \frac{r(\hat{v}, d_i)^2}{\epsilon^{2-p}} + \left( \frac{1}{p} - \frac{1}{2} \right) \epsilon^p \right) = \frac{1}{2} \nabla \left( \frac{r(\hat{v}, d_i)^2}{\epsilon^{2-p}} \right) = \frac{1}{2} \hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2).$$

On the other hand,  $r(\hat{v}, d_i) \neq 0$ , then  $r(v, d_i)$  is differentiable at  $v = \hat{v}$  by Assumption 2. Therefore, applying the chain rule and the equality  $\rho'_{\epsilon}(r(\hat{v}, d_i)) = \hat{w}_i \cdot r(\hat{v}, d_i)$ , we get

$$\nabla \rho_{\epsilon}(r(\hat{v}, d_i)) = \rho'_{\epsilon}(r(\hat{v}, d_i)) \cdot \nabla r(\hat{v}, d_i) = \hat{w}_i \cdot r(\hat{v}, d_i) \cdot \nabla r(\hat{v}, d_i) = \frac{1}{2} \hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2).$$

Substitute this into (24), and we obtain

$$\left( \sum_{i=1}^m \nabla \rho_{\epsilon}(r(\hat{v}, d_i)) \right)^{\top} \mathbf{b} \geq 0, \quad \forall \mathbf{b} \in T_{\hat{v}}\mathcal{C},$$

which means  $\hat{v}$  is a stationary point of (10). The proof is complete.  $\square$

## D. Majorizing Truncated Least-Squares and Proof of Theorem 3

In this section, we consider the TLS loss with some hyper-parameter  $c^2$  ( $c > 0$ ),

$$\rho(r) = \min\{r^2, c^2\}. \quad (25)$$

As mentioned in the main paper, here we present a smooth majorizer and a quadratic majorizer for the TLS loss (25), in Section D.1 and Section D.2 respectively. After that, we provide a proof of Theorem 3. The development of this section follows much from Section 2, but this time we discuss how the majorizers are designed in more detail.

### D.1. A Smooth Majorizer for The TLS Loss

In designing a majorizer  $\bar{\rho}_{\mu}$  for the TLS loss  $\rho$  (25), we expect the following properties that it should yield:

- (S1)  $\bar{\rho}_{\mu}$  is continuously differentiable.
- (S2)  $\bar{\rho}(r) \leq \bar{\rho}_{\mu}(r)$  for every  $r \in \mathbb{R}$ .
- (S3)  $\mu' \leq \mu \Rightarrow \bar{\rho}_{\mu}(r) \leq \bar{\rho}_{\mu'}(r)$ .
- (S4)  $\bar{\rho}_{\mu} \rightarrow \rho$  as  $\mu \rightarrow \infty$ .



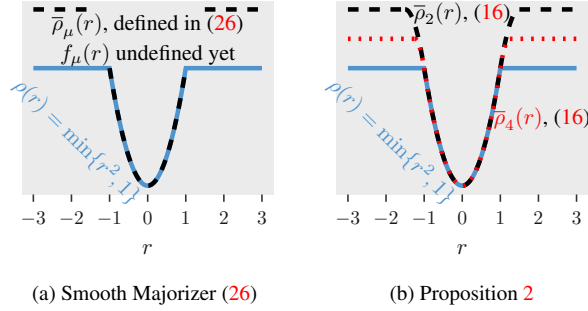


Figure 7. The TLS loss  $\rho(r)$  (25) and its smooth majorizer  $\bar{\rho}_\mu(r)$  (cf. (26), (28), and (16)).

The above four properties correspond to those of Lemma 1. (S1) ensures the weighting strategy (5) makes sense. (S2) ensures majorization. (S3) is in different direction than that of Lemma 1; this is because we are increasing  $\mu$  to approach  $\rho$ , whereas  $\epsilon$  is decreased in Lemma 1 to achieve similar effects. (S4) ensures the approximation quality, and as a side effect, it guarantees binary weights as  $\mu \rightarrow \infty$ . Such majorizers do exist, as we will show next.

We choose majorizers based on the following idea. First we make sure  $\bar{\rho}_\mu(r) = r^2$  whenever  $r^2 \in [0, c^2]$ . Then we set  $\bar{\rho}_\mu(r) = \frac{(\mu+1)}{\mu}c^2$  if  $r^2 \geq \nu_\mu^2 c^2$ , where  $\nu_\mu$  is some number to be determined and we require  $\nu_\mu > 1$  and  $\nu_\mu \rightarrow 1$  as  $\mu \rightarrow \infty$ . Specifically, we parameterize  $\bar{\rho}_\mu(r)$  as

$$\bar{\rho}_\mu(r) = \begin{cases} r^2, & \text{if } r^2 \in [0, c^2], \\ f_\mu(|r|), & \text{if } r^2 \in (c^2, \nu_\mu^2 c^2), \\ \frac{\mu+1}{\mu}c^2, & \text{if } r^2 \in [\nu_\mu^2 c^2, \infty), \end{cases} \quad (26)$$

where the function  $f : \mathbb{R} \rightarrow \mathbb{R}$  is to be determined. This is illustrated in Figure 7a (dashed, black). This construction ensures (S4). It remains to find  $f_\mu$  and  $\nu_\mu$  to fulfill properties (S1)-(S3). In view of the requirement of continuous differentiability (S1),  $f_\mu$  itself must be continuously differentiable, which means that it should satisfy the following constraints with  $\nu_\mu$ :

$$\begin{aligned} f_\mu(c) &= c^2 \\ f'_\mu(c) &= 2c \\ f'_\mu(\nu_\mu c) &= 0 \end{aligned} \quad (27)$$

Interestingly, if  $f_\mu$  belongs to some specific class of functions (e.g., quadratic functions), then it is determined as the *unique* solution of (27). And if so, this construction (26) automatically satisfies the remaining properties, namely (S2) and (S3), which makes  $\rho_\mu$  (26) into a smooth majorizer, as desired. More formally, we have the following proposition:

**Proposition 2.** *Let  $\rho$  be the truncated least-squares loss (25). Let  $\rho_\mu$  be defined in (26), where  $f_\mu : \mathbb{R} \rightarrow \mathbb{R}$  is a quadratic function and  $\nu_\mu > 1$  with  $\nu_\mu \rightarrow 1$  as  $\mu \rightarrow \infty$ . If  $\bar{\rho}_\mu$  satisfies properties (S1)-(S4), then it is uniquely determined by the following expressions:*

$$\bar{\rho}_\mu(r) = \begin{cases} r^2, & \text{if } r^2 \in [0, c^2], \\ -\mu r^2 + 2(1 + \mu)c|r| - (1 + \mu)c^2, & \text{if } r^2 \in (c^2, \frac{(\mu+1)^2}{\mu^2}c^2), \\ \frac{\mu+1}{\mu}c^2, & \text{if } r^2 \in [\frac{(\mu+1)^2}{\mu^2}c^2, \infty). \end{cases} \quad (28)$$

*Pictorial Proof of Proposition 2.* We need to “attach” a quadratic function  $f_\mu$  into Figure 7a so that it “connects” the black dashed curves of Figure 7a and makes the overall curve continuously differentiable. We attach such  $f_\mu(|r|) = -\mu r^2 + 2(1 + \mu)c|r| - (1 + \mu)c^2$  (28) in Figure 7b with different values of  $\mu$ , where, visually, all properties (S1)-(S4) are fulfilled.  $\square$

*Proof of Proposition 2.* Note that  $\bar{\rho}_\mu$  is symmetric with  $\bar{\rho}_\mu(r) = \bar{\rho}_\mu(-r)$ , so we assume  $r \geq 0$  without loss of generality. Since  $f_\mu$  is a quadratic function, we write  $f_\mu(r) = ar^2 + br + d$ , where  $a, b$  and  $d$  are some real-valued coefficients to be

determined. These coefficients depend on  $\mu$ , but we drop the indices for simplicity. Note that  $\nu_\mu c \neq 0$  and  $f$  is quadratic, for  $f'_\mu(\nu_\mu c) = 0$  of (27) to hold, we must have  $\nu_\mu c = -\frac{b}{2a}$  and  $f(-\frac{b}{2a}) = \frac{\mu+1}{\mu}c^2$ . Hence, (27) gives the following equations:

$$\begin{aligned} ac^2 + bc + d &= c^2 \\ 2ac + b &= 2c \\ -\frac{b}{2a} &= \nu_\mu c \\ -\frac{b^2}{4a} + d &= \frac{\mu+1}{\mu}c^2 \end{aligned}$$

Note that  $\mu$  and  $c$  are given, the above are four equations with four unknowns, which admit a unique solution:

$$\begin{aligned} a &= -\mu \\ b &= 2c(1 + \mu) \\ d &= -(1 + \mu)c^2 \\ \nu_\mu &= \frac{1 + \mu}{\mu} \end{aligned}$$

This gives a unique expression of  $\bar{\rho}_\mu$  (28). By construction,  $\bar{\rho}_\mu$  is continuously differentiable (S1) and approximates the TLS loss  $\rho$  as  $\mu \rightarrow \infty$  (S4). It is a simple routine to verify that  $\bar{\rho}_\mu$  also satisfies (S2) and (S3). The proof is complete.  $\square$

*Remark 4 (Other classes of functions).* The majorizer of the TLS loss exists also for other classes of functions. For example, if  $f_\mu : \mathbb{R} \rightarrow \mathbb{R}$  is a trigonometric function parameterized as  $f_\mu(r) = a \sin(b(r - c)) + c^2$ , then one can solve (27) for  $a$  and  $b$ . Pictorially, this amounts to attaching a trigonometric function to Figure 7a that ensures continuous differentiability. However, we did not find more benefits of doing so than what (28) gives, so we omitted this construction.

*Remark 5 (Other quadratic functions  $f_\mu$ ).* While Proposition 2 shows  $f_\mu$  is a unique quadratic function, it is based on the fact that  $\rho_\mu$  is parameterized as (26). Different parameterizations lead to different quadratic functions, and finding an optimal one that accelerates convergence is left to future work.

## D.2. A Quadratic Majorizer for The TLS Loss

We will construct a quadratic majorizer  $q_\mu(r, \cdot) : \mathbb{R} \rightarrow \mathbb{R}$  of  $\bar{\rho}_\mu(r)$  (28). We expect  $q_\mu(r, \cdot)$  to satisfy three properties:

- (Q1)  $q_\mu(r, \theta)$  is a quadratic *symmetric* function in  $r$ ; by *symmetric* we mean that  $q_\mu(r, \theta) = q_\mu(-r, \theta)$  for any  $r, \theta \in \mathbb{R}$ .
- (Q2)  $\bar{\rho}_\mu(r) \leq q_\mu(r, \theta)$  for every  $r, \theta \in \mathbb{R}$ .
- (Q3)  $\bar{\rho}_\mu(\theta) = q_\mu(\theta, \theta)$  for every  $\theta \in \mathbb{R}$ .

(Q1) and (Q2) are to make sure that  $q_\mu(r, \theta)$  is indeed a quadratic majorizer of  $\bar{\rho}_\mu(r)$ ; requiring  $q_\mu(r, \theta)$  to be symmetric in  $r$  is because  $\bar{\rho}_\mu(r)$  is symmetric in  $r$ . (Q3) makes sure that  $q_\mu(r, \theta)$  is tight and, as such, minimizing  $q_\mu(r, \theta)$  is expected to lead to minimization of  $\bar{\rho}_\mu(\theta)$ .

Constructing such a quadratic function to majorize  $\bar{\rho}_\mu(\theta)$  at a point  $\theta$  is not hard. In particular, note that  $\bar{\rho}_\mu(r)$  is locally quadratic, or otherwise concave (or in particular linear) in  $r$ . Its concave part (Figure 8, red) can be majorized at any given point  $\theta$  by any tangent line passing through  $\theta$  (Figure 8, blue), and furthermore by a quadratic convex function (Figure 8, black); the latter is what we need. Next, we make this argument more formal.

To start with, let us assume  $\theta > 0$ . If  $\theta \in [0, c]$  we could simply set  $q_\mu(r, \theta) = r^2$ , which fulfills all desired properties (Q1)-(Q3). So let us next focus on the case  $\theta > c$ . To satisfy (Q1), write  $q(r, \theta) = ar^2 + d$  for some unknown coefficients  $a$  and  $d$ , which are to be determined and implicitly depend on  $\theta$ . (Q3) gives us

$$a\theta^2 + d = \bar{\rho}_\mu(\theta),$$

and we would expect that (Q2) gives us one more equation so that we can solve them for  $a$  and  $d$ . Since  $\bar{\rho}_\mu(\theta)$  is concave in  $(c, \infty)$ , it is majorized by its tangent line passing through  $\theta$  with derivative  $\bar{\rho}'_\mu(\theta)$ , and thus by a quadratic (convex) function passing through  $\theta$  with derivative  $\bar{\rho}'_\mu(\theta)$ . Therefore, for  $q(r, \theta)$  to majorize  $\bar{\rho}_\mu(r)$  when  $\theta > c$ , it suffices to have the equality

$$\frac{\partial q(r, \theta)}{\partial r} = \bar{\rho}'_\mu(\theta) \Leftrightarrow 2a = \bar{\rho}'_\mu(\theta). \quad (29)$$

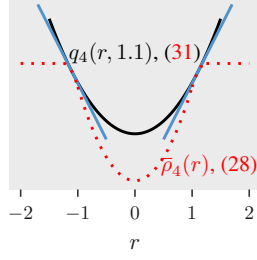


Figure 8. The quadratic majorizer (31) for the TLS loss.

Now, solve the above two equations for  $a$  and  $d$  and generalize to the case  $\theta^2 \geq c^2$ . Then we get the following expression of the quadratic majorizer  $q(\cdot, \theta)$  for  $\theta^2 \geq c^2$ :

$$q_\mu(r, \theta) = \bar{\rho}_\mu(\theta) + \frac{\bar{\rho}'_\mu(\theta)}{2|\theta|}(r^2 - \theta^2) = \begin{cases} \bar{\rho}_\mu(\theta) + \left(\frac{c(1+\mu)}{|\theta|} - \mu\right) \cdot (r^2 - \theta^2), & \text{if } \theta^2 \in \left(c^2, \frac{(\mu+1)^2}{\mu^2}c^2\right), \\ \bar{\rho}_\mu(\theta), & \text{if } \theta^2 \geq \frac{(\mu+1)^2}{\mu^2}c^2. \end{cases} \quad (30)$$

We now summarize our discovery about  $q_\mu(r, \theta)$ :

**Proposition 3.** Recall the definition of the smooth majorizer  $\bar{\rho}_\mu(\cdot)$  (28) for the TLS loss. The quadratic function

$$q_\mu(r, \theta) = \begin{cases} r^2, & \text{if } \theta^2 \in [0, c^2], \\ \bar{\rho}_\mu(\theta) + \left(\frac{c(1+\mu)}{|\theta|} - \mu\right) \cdot (r^2 - \theta^2), & \text{if } \theta^2 \in \left(c^2, \frac{(\mu+1)^2}{\mu^2}c^2\right), \\ \bar{\rho}_\mu(\theta), & \text{if } \theta^2 \geq \frac{(\mu+1)^2}{\mu^2}c^2, \end{cases} \quad (31)$$

is a quadratic majorizer of  $\bar{\rho}_\mu(r)$  that fulfills (Q1), (Q2), and (Q3).

*Proof.* Note that, by design,  $q_\mu(r, \theta)$  (31) fulfills (Q1) and (Q3). However, it should be noted that, by design, the majorizer (30) only makes sure  $q_\mu(r, \theta) \geq \bar{\rho}_\mu(r)$  for any  $r$  and  $\theta$  that fulfills  $r^2 > c^2$  and  $\theta^2 > c^2$  (and similarly for the case  $\theta^2 \leq c^2$ ); therefore, slightly more work is needed to verify (Q2).

If  $\theta^2 \leq c^2$ , we have

$$q_\mu(r, \theta) = r^2, \geq \begin{cases} r^2, & \text{if } r^2 \in [0, c^2], \\ -\mu r^2 + 2(1+\mu)c|\theta| - (1+\mu)c^2, & \text{if } r^2 \in \left(c^2, \frac{(\mu+1)^2}{\mu^2}c^2\right), = \bar{\rho}_\mu(r). \\ \frac{\mu+1}{\mu}c^2, & \text{if } r^2 \in \left[\frac{(\mu+1)^2}{\mu^2}c^2, \infty\right), \end{cases}$$

On the other hand, consider the case  $\theta^2 \geq c^2$ . We need to prove  $q_\mu(r, \theta) \geq \bar{\rho}_\mu(r)$  whenever  $r^2 \leq c^2$  and  $\theta^2 > c^2$ , in which case  $\bar{\rho}_\mu(r)$  is quadratic and the inequality to prove is  $q_\mu(r, \theta) \geq r^2$ . This is true if  $\theta^2 \geq (\mu+1)^2 c^2 / \mu^2$ , so we consider the case  $\theta^2 \in \left(c^2, \frac{(\mu+1)^2}{\mu^2}c^2\right)$ . In this case, we have

$$\begin{aligned} q_\mu(r, \theta) &= \bar{\rho}_\mu(\theta) + \left(\frac{c(1+\mu)}{|\theta|} - \mu\right) \cdot (r^2 - \theta^2) \\ &= -\mu\theta^2 + 2(1+\mu)c|\theta| - (1+\mu)c^2 + \left(\frac{c(1+\mu)}{|\theta|} - \mu\right) \cdot (r^2 - \theta^2) \\ &= (1+\mu)c|\theta| - (1+\mu)c^2 + \left(\frac{c(1+\mu)}{|\theta|} - \mu\right) \cdot r^2 \\ &= (1+\mu) \cdot \left(c|\theta| - c^2 + \frac{c}{|\theta|}r^2 - r^2\right) + r^2 \\ &= (1+\mu) \cdot (|\theta| - c) \cdot \left(c - \frac{r^2}{|\theta|}\right) + \bar{\rho}_\mu(r) \\ &\geq \bar{\rho}_\mu(r), \end{aligned}$$

by which the proof is complete.  $\square$

*Remark 6.* Similar statements in Remarks 1 and 3 can also be made for MS-GNC-TLS and its smooth and quadratic majorizers. In particular, first note that the coefficient of the quadratic term of  $q_\mu(r, \theta)$  is exactly  $\bar{\rho}_\mu(\theta)/(2|\theta|)$ , which coincides with the weight update rule (17) of MS-GNC-TLS if  $\theta$  is the residual  $r_i^{(t)}$  at the  $t$ -th iteration. Therefore, the weighted least-squares problem that MS-GNC-TLS solves at each iteration  $t$  is minimizing the quadratic majorizer

$$\begin{aligned} v^{(t+1)} &\in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m w_i^{(t+1)} r(v, d_i)^2 \\ &= \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m \frac{\bar{\rho}_\mu(r_i^{(t)})}{2r_i^{(t)}} r(v, d_i)^2 \\ &= \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m q_{\mu^{(t)}}(r(v, d_i), r_i^{(t)}). \end{aligned}$$

### D.3. Proof of Theorem 3

The proof follows very similarly from that of Theorem 1, therefore we provide a high-level proof here and refer to the proof of Theorem 1 for omitted steps. Recall  $r_i^{(t)} := r(v^{(t)}, d_i), \forall t$ . Similarly to (21), based on majorization properties of smooth and quadratic majorizers, we can show that the objective, despite keeping changing, has non-increasing values:

$$\sum_{i=1}^m \bar{\rho}_{\mu^{(t+1)}}(r_i^{(t+1)}) \leq \sum_{i=1}^m \bar{\rho}_{\mu^{(t)}}(r_i^{(t)})$$

Since  $\{v^{(t)}\}_t$  is a bounded sequence, it has a convergent subsequence, say,  $v_j^{(t_j)}$ , which converges to  $\hat{v}$ . It remains to show that  $\hat{v}$  is a stationary point of (18). Towards this end, one can follow the reasoning of (23) and show that

$$\sum_{i=1}^m q_{\mu^{(t_j+1)}}(r_i^{(t_j+1)}, r_i^{(t_j+1)}) \leq \sum_{i=1}^m q_{\mu^{(t_j)}}(r(v, d_i), r_i^{(t_j)}), \forall v \in \mathcal{C}.$$

As  $t_j \rightarrow \infty$ , the above equality becomes

$$\begin{aligned} \sum_{i=1}^m q_\mu(r(\hat{v}, d_i), r(\hat{v}, d_i)) &\leq \sum_{i=1}^m q_\mu(r(v, d_i), r(\hat{v}, d_i)), \forall v \in \mathcal{C} \\ \Leftrightarrow \hat{v} \in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m q_\mu(r(v, d_i), r(\hat{v}, d_i)) &= \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m \hat{w}_i \cdot r(v, d_i)^2, \end{aligned}$$

where  $\hat{w}_i$  is defined to be the weight (17) as  $t_j$  goes to infinity, that is

$$\hat{w}_i = \begin{cases} 1, & \text{if } r(\hat{v}, d_i) \leq c, \\ 0, & \text{if } r(\hat{v}, d_i) \geq \frac{\mu+1}{\mu}c, \\ \frac{c(1+\mu)}{r(\hat{v}, d_i)} - \mu, & \text{if } r(\hat{v}, d_i) \in (c, \frac{\mu+1}{\mu}c). \end{cases}$$

Therefore,  $\hat{v}$  must satisfy the geometric condition (cf. (24))

$$\left( \sum_{i=1}^m \hat{w}_i \cdot \nabla(r(\hat{v}, d_i)^2) \right)^\top \mathbf{b} \geq 0, \quad \forall \mathbf{b} \in T_{\hat{v}}\mathcal{C},$$

which is equivalent to (see the proof of Theorem 1 for a similar derivation)

$$\left( \sum_{i=1}^m \nabla \bar{\rho}_\mu(r(\hat{v}, d_i)) \right)^\top \mathbf{b} \geq 0, \quad \forall \mathbf{b} \in T_{\hat{v}}\mathcal{C},$$

meaning that  $\hat{v}$  is a stationary point of (18).

## E. Convergence Theory of IRLS +GNC for A General Class of Losses

In this section we develop a general convergence theory for IRLS and GNC applied to the outlier-robust estimation problem (2). We repeat the optimization problem here for convenience:

$$\min_{v \in \mathcal{C}} \sum_{i=1}^m \rho_{\zeta}(r(v, d_i)) \quad (32)$$

This time we added a subscript  $\zeta \in \mathbb{R} \cup \{\infty\}$  to the loss  $\rho_{\zeta}$  for notational convenience. Since we are now working with a general loss function  $\rho_{\zeta}$ , we need some assumptions on it for deriving meaningful convergence results. Our previous experience suggests that  $\rho_{\zeta}$  should admit some smooth majorizer and quadratic majorizer. We start with assuming the first:

**Assumption 3** (Existence of A Smooth Majorizer). There exists a function  $\rho_{\epsilon} : \mathbb{R} \rightarrow \mathbb{R}$  parametrized by  $\epsilon > 0$ , which satisfies the following properties: (i)  $\rho_{\epsilon}$  is continuously differentiable, (ii)  $\rho_{\epsilon} \geq \rho_{\zeta}$  for every  $\epsilon > 0$ , (iii)  $\rho_{\epsilon}(r) \geq \rho_{\epsilon}'(r)$  as long as  $\epsilon \geq \epsilon'$ , (iv)  $\rho_{\epsilon}(r) \rightarrow \rho_{\zeta}(r)$  whenever  $\epsilon \rightarrow \zeta$ , (v)  $\rho_{\epsilon}$  is symmetric, i.e.,  $\rho_{\epsilon}(r) = \rho_{\epsilon}(-r)$ .

A simple case is when  $\rho_{\zeta}$  is already continuously differentiable and symmetric, which makes itself a smooth majorizer. We knew from Lemma 1 that the  $\ell_p$ -loss admits a smooth majorizer, and in this case we have  $\zeta = 0$ . We knew from Section D.1 that the TLS loss admits a smooth majorizer (if one of the inequalities in the property (iii) is reversed), and in this case we have  $\zeta = \infty$ . Indeed, smooth majorizers of  $\rho_{\zeta}$  are not hard to find. A simple way to do so is to plot  $\rho_{\zeta}$  and then specify a ‘‘higher’’ continuously differentiable curve (recall Figures 1a and 5b, Section D.1).

From a loss function  $\rho_{\zeta}$  and its smooth majorizer  $\rho_{\epsilon}$ , a quadratic majorizer  $q_{\epsilon}(r, \theta)$  ensues if  $\rho_{\epsilon}(\sqrt{r})$  is concave and differentiable in  $r$  over  $[0, \infty)$ . In particular, the concave function  $\rho_{\epsilon}(\sqrt{r})$  ( $r \geq 0$ ) is majorized by its tangent line at a given point  $\theta \geq 0$ , meaning that

$$\rho_{\epsilon}(\sqrt{r}) \leq (\rho_{\epsilon}(\sqrt{\theta}))' \cdot (r - \theta) + \rho_{\epsilon}(\sqrt{\theta}) = \begin{cases} \frac{\rho_{\epsilon}'(\sqrt{\theta})}{2\sqrt{\theta}} \cdot (r - \theta) + \rho_{\epsilon}(\sqrt{\theta}), & \forall r \geq 0, \forall \theta > 0, \\ (\rho_{\epsilon}(\sqrt{\theta}))' \cdot r + \rho_{\epsilon}(0), & \forall r \geq 0, \theta = 0. \end{cases} \quad (33)$$

Make a substitution of variables (say  $\sqrt{r} \leftarrow |r|$  and  $\sqrt{\theta} \leftarrow |\theta|$ ) in (33), and we arrive at

$$\rho_{\epsilon}(r) \leq \frac{\rho_{\epsilon}'(\theta)}{2\theta} \cdot (r^2 - \theta^2) + \rho_{\epsilon}(\theta), \quad \forall r \geq 0, \forall \theta \neq 0,$$

and  $\rho_{\epsilon}(r) \leq (\rho_{\epsilon}(\sqrt{\theta}))' \cdot r^2 + \rho_{\epsilon}(0)$  if  $\theta = 0$ . In other words, we now have the quadratic majorizer

$$q_{\epsilon}(r, \theta) = \begin{cases} \frac{\rho_{\epsilon}'(\theta)}{2|\theta|} \cdot (r^2 - \theta^2) + \rho_{\epsilon}(\theta), & \text{if } \theta \neq 0, \\ (\rho_{\epsilon}(\sqrt{\theta}))' \cdot r^2 + \rho_{\epsilon}(0), & \text{if } \theta = 0. \end{cases} \quad (34)$$

One verifies that the quadratic majorizers for the  $\ell_p$  and TLS losses are special cases of (34). Note that in the above development, we have distinguished the case  $\theta \neq 0$  and  $\theta = 0$ . This is because  $\sqrt{\theta}$  is not differentiable at  $\theta = 0$ , and the chain rule cannot be applied directly. That said, for the case  $\theta = 0$ , the derivative  $(\rho_{\epsilon}(\sqrt{\theta}))'$  can be easily evaluated once the exact form of  $\rho_{\epsilon}$  is given. Finally, since  $\rho_{\epsilon}(\sqrt{\theta})$  is differentiable in  $\theta$  and  $\rho_{\epsilon}(\cdot)$  is continuously differentiable, we can write

$$\begin{aligned} (\rho_{\epsilon}(\sqrt{\theta}))' \Big|_{\theta=0} &= \lim_{\theta=0} (\rho_{\epsilon}(\sqrt{\theta}))' \\ &= \lim_{\theta=0} \frac{\rho_{\epsilon}'(\sqrt{\theta})}{2\sqrt{\theta}} \\ &= \lim_{\theta=0} \frac{\rho_{\epsilon}'(\theta)}{2\theta}. \end{aligned} \quad (35)$$

This equality will be useful in the sequel.

We have therefore justified the following assumption:

**Assumption 4** (Existence of A Quadratic Majorizer). For a loss function  $\rho_{\zeta}$ , assume its smooth majorizer  $\rho_{\epsilon}$  composed with the square root function is concave and differentiable on  $[0, \infty)$ , that is  $\rho_{\epsilon}(\sqrt{r})$  is concave and differentiable in  $r$  over  $[0, \infty)$ .

**Example 1.** Here are several examples of continuously differentiable losses  $\rho(\cdot)$  for which  $\rho_\epsilon(\sqrt{r})$  is differentiable and concave in  $r$ . See [36, Appendix 6.8], [1], and [20, Table 1] for more examples.

- The smooth majorizers for the  $\ell_p$ -loss and TLS loss, as we have discussed.
- The Geman-McClure loss [31], which is defined as

$$\rho_\epsilon(r) = \frac{\epsilon c^2 r^2}{\epsilon c^2 + r^2}.$$

- A surrogate function of the TLS loss used in [45, 89]:

$$\rho_\epsilon(r) = \begin{cases} r^2(2 - \frac{r^2}{c^2}), & \text{if } r^2 \leq c^2, \\ c^2, & \text{if } r^2 > c^2. \end{cases}$$

Note that this loss function does not depend on  $\epsilon$ . It is a smooth majorizer of itself.

- The Blake-Zisserman loss (as [36, Appendix 6.8] called):

$$\rho_\epsilon(r) = -\log(\exp(-r^2) + c)$$

With the smooth and quadratic majorizers of  $\rho_\zeta$ , we can now use the IRLS framework with the weight update rule

$$w_i^{(t+1)} \leftarrow \begin{cases} \frac{\rho'_{\epsilon^{(t)}}(r_i^{(t)})}{2r_i^{(t)}}, & \text{if } r_i^{(t)} \neq 0, \\ \left( \rho_{\epsilon^{(t)}}(\sqrt{r_i^{(t)}}) \right)', & \text{if } r_i^{(t)} = 0, \end{cases} \quad (36)$$

and some GNC schedule  $\{\epsilon^{(t)}\}_t$  to optimize  $\rho_\zeta$ . Compared to rule (5), our rule (36) explicitly takes the case  $r_i^{(t)} = 0$  into account and makes it mathematically correct (computationally, this makes no difference if all the differentiability assumptions are satisfied). We are now ready to state the main result of this section:

**Theorem 4.** Let  $\{v^{(t)}\}_t$  be the iterates generated by IRLS with weight update rule (36) and some GNC schedule  $\{\epsilon^{(t)}\}_t$ . Assume  $\|v^{(t)}\|_2 < \infty$  for every  $t$ . Suppose that  $\{\epsilon^{(t)}\}_t$  is non-increasing and converges to  $\zeta$  with  $\zeta > 0$ . Under Assumptions 1-4, every accumulation point of  $\{v^{(t)}\}_t$  is a stationary point of the outlier-robust estimation problem (32).

*Proof of Theorem 4.* The proof follows closely from that of Theorem 1. Recall the definition  $r_i^{(t)} = r(v^{(t)}, d_i)$ . Note that, by the weighting strategy (36),  $v^{(t+1)}$  is a global minimizer of the quadratic majorizer (34) (also recall Remark 3):

$$\begin{aligned} v^{(t+1)} &\in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m w_i^{(t+1)} r(v, d_i)^2 \\ &= \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m q_{\epsilon^{(t)}}(r(v, d_i), r_i^{(t)}) \end{aligned}$$

Since  $\{v^{(t)}\}_t$  is bounded, it has a subsequence  $\{v^{(t_j)}\}_j$  convergent to some point say  $\hat{v} \in \mathcal{C}$ . For this subsequence, it holds that (by identical arguments as in (21) and (23))

$$\sum_{i=1}^m q_{\epsilon^{(t_{j+1})}}(r_i^{(t_{j+1})}, r_i^{(t_{j+1})}) \leq \sum_{i=1}^m q_{\epsilon^{(t_j)}}(r(v, d_i), r_i^{(t_j)}), \quad \forall v \in \mathcal{C}.$$

Since  $\rho'_{\epsilon^{(t_j)}}$  is continuously differentiable, we see that  $q_{\epsilon^{(t_j)}}(34)$  is continuous in both of its parameters. Therefore, as  $t_j \rightarrow \infty$ , the above inequality becomes

$$\begin{aligned} \sum_{i=1}^m q_\zeta(r(\hat{v}, d_i), r(\hat{v}, d_i)) &\leq \sum_{i=1}^m q_\zeta(r(v, d_i), r(\hat{v}, d_i)), \quad \forall v \in \mathcal{C} \\ \Rightarrow \hat{v} \in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m q_\zeta(r(v, d_i), r(\hat{v}, d_i)) &= \hat{v} \in \operatorname{argmin}_{v \in \mathcal{C}} \sum_{i=1}^m \hat{w}_i \cdot r(v, d_i)^2 \end{aligned}$$

where  $\hat{w}_i$  is the weight of the  $i$ -th sample as  $t_j \rightarrow \infty$ , that is

$$\begin{aligned}\hat{w}_i &= \frac{\rho'_\zeta(r(\hat{v}, d_i))}{2r(\hat{v}, d_i)} && \text{if } r(\hat{v}, d_i) \neq 0 \\ \hat{w}_i &= \left( \rho_\zeta(\sqrt{r(\hat{v}, d_i)}) \right)' && \text{if } r(\hat{v}, d_i) = 0\end{aligned}$$

Global optimality of  $\hat{v}$  for the weighted least-squares problem in the limit implies the optimality condition [3, Section 5.3]

$$\left( \sum_{i=1}^m \hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2) \right)^\top \mathbf{b} \geq 0, \quad \forall \mathbf{b} \in T_{\hat{v}}\mathcal{C}.$$

As before, we finish the proof by showing that  $\hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2)$  is equal to  $\nabla \rho_\zeta(r(\hat{v}, d_i))$ . Indeed, if  $r(\hat{v}, d_i) \neq 0$  then

$$\begin{aligned}\hat{w}_i \cdot \nabla (r(\hat{v}, d_i)^2) &= \hat{w}_i \cdot 2r(\hat{v}, d_i) \cdot \nabla (r(\hat{v}, d_i)) \\ &= \rho'_\zeta(r(\hat{v}, d_i)) \cdot \nabla (r(\hat{v}, d_i)) \\ &= \nabla \rho_\zeta(r(\hat{v}, d_i)).\end{aligned}$$

If  $r(\hat{v}, d_i) = 0$ , using (35) and basic limit operations we obtain

$$\begin{aligned}\hat{w}_i \cdot \nabla_{r(\hat{v}, d_i)=0} (r(\hat{v}, d_i)^2) &= \left( \rho_\zeta(\sqrt{r(\hat{v}, d_i)}) \right)' \cdot 2r(\hat{v}, d_i) \cdot \nabla_{r(\hat{v}, d_i)=0} (r(\hat{v}, d_i)) \\ &= \nabla_{r(\hat{v}, d_i)=0} (r(\hat{v}, d_i)) \cdot \lim_{r(\hat{v}, d_i) \rightarrow 0} 2r(\hat{v}, d_i) \cdot \lim_{r(\hat{v}, d_i) \rightarrow 0} \frac{\rho'_\zeta(r(\hat{v}, d_i))}{2r(\hat{v}, d_i)} \\ &= \nabla_{r(\hat{v}, d_i)=0} (r(\hat{v}, d_i)) \cdot \lim_{r(\hat{v}, d_i) \rightarrow 0} \rho'_\zeta(r(\hat{v}, d_i)) \\ &= \nabla_{r(\hat{v}, d_i)=0} (r(\hat{v}, d_i)) \cdot \rho'_\zeta(r(\hat{v}, d_i)) \Big|_{r(\hat{v}, d_i)=0} \\ &= \nabla_{r(\hat{v}, d_i)=0} \left( \rho_\zeta(r(\hat{v}, d_i)) \right)\end{aligned}$$

The proof is now complete.  $\square$

## F. Proof of Theorem 2

Here we collect notations that are used throughout our proofs. Let  $p \in [0, 1]$ . Let the weight  $w_i^{(t+1)} \in \mathbb{R}$ , vector  $\mathbf{x}^{(t)} \in \mathbb{R}^n$ , and smoothing parameter  $\epsilon^{(t)} \in \mathbb{R}$  be defined as per Algorithm 1. Let  $\mathbf{W}^{(t+1)} \in \mathbb{R}^{m \times m}$  be the diagonal matrix with  $w_i^{(t+1)}$  its diagonal, *i.e.*,  $\mathbf{W}^{(t+1)} = \text{diag}(w_1^{(t+1)}, \dots, w_m^{(t+1)})$ . The ground truth signal is denoted by  $\mathbf{x}^*$  and  $\mathbf{r}^* = [r_1^*, \dots, r_m^*]^\top \in \mathbb{R}^m$  with  $r_i^* := \mathbf{a}_i^\top \mathbf{x}^* - y_i$  is the associated residual vector, which is assumed to be  $k$ -sparse. Denote by  $S^*$  the support of  $\mathbf{r}^*$ , that is  $S^* := \{i : r_i^* \neq 0\} \subset \{1, \dots, m\}$ .

By our assumption, we have  $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \leq \epsilon^{(0)}$ . Therefore, under other assumptions of Theorem 2, we can invoke Theorem 5 to obtain

$$\|\mathbf{x}^{(1)} - \mathbf{x}^*\|_2 \leq \beta(\epsilon^{(0)})^{2-p}$$

with probability at least  $1 - P'_0 - P'_1 - P'_2$ , where each  $P'_i$  is defined in (38). In fact, with a union bound, we can invoke Theorem 5 for  $t$  times and get

$$\|\mathbf{x}^{(t_0+1)} - \mathbf{x}^*\|_2 \leq \beta(\epsilon^{(t_0)})^{2-p}, \quad \forall t_0 = 0, 1, \dots, t-1, \quad (37)$$

with probability at least  $1 - tP'_0 - tP'_1 - tP'_2$ . Via simple algebraic manipulation, we see that (37) implies the desired bound (13) of Theorem 2. It is easy to show that the probability term  $1 - tP'_0 - tP'_1 - tP'_2$  matches the probability term  $1 - P_0 - P_1 - P_2$  of Theorem 2. For example, we have

$$tP'_0 = \exp(-\Omega(n - \log m - \log t)) = \exp(-\tilde{\Omega}(n)) = P_0.$$

Note that the logarithmic term  $\log t$  is suppressed by  $\tilde{\Omega}$  as  $t$  is a constant. Similarly we have  $tP'_1 = P_1$  and  $tP'_2 = P_2$ . The proof is complete, provided we can prove Theorem 5. This is the most technical part and is done in the next section.

## G. Theorem 5 and Its Proof

**Theorem 5.** Suppose  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \epsilon^{(t)}$  and  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0, 1)$  entries. Denote by  $r_{\min+}^*$  the smallest non-zero residual among  $\{|\mathbf{a}_i^\top \mathbf{x}^* - y_i|\}_{i=1}^m$ . Define

$$\alpha := \frac{\sqrt{5} \cdot 2^{2-p}}{0.99 \cdot 0.516} \cdot \frac{1}{(r_{\min+}^*)^{1-p}} \cdot \frac{\sqrt{k} \cdot (1.01\sqrt{k} + \sqrt{n})}{(m-k)}.$$

If  $\alpha \leq \beta$  ( $\beta$  is defined in (11)), then we have

$$\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2 \leq \beta(\epsilon^{(t)})^{2-p}$$

with probability at least  $1 - P'_0 - P'_1 - P'_2$ , where

$$\begin{aligned} P'_0 &:= \exp(-\Omega(n - \log m)) \\ P'_1 &:= \exp(-\Omega(k - n)) \\ P'_2 &:= \exp(-\Omega(m - k - n \log n)). \end{aligned} \quad (38)$$

*Proof of Theorem 5.* To begin with, consider the event  $\mathcal{E}$

$$\mathcal{E} := \{\sqrt{0.5n} \leq M \leq \sqrt{5n}\}, \quad M := \max_{i=1, \dots, m} \|\mathbf{a}_i\|_2 \quad (39)$$

By Lemma 10 and a union bound we have with probability at least  $1 - m \cdot \exp(-n) - \exp(-n/16) = 1 - P'_0$  that

$$\mathbb{P}(\mathcal{E}) \geq 1 - m \cdot \exp(-n) - \exp(-n/16) = 1 - \exp(-\Omega(n - \log m)). \quad (40)$$

As indicated by (40), we will heavily use the big- $\Omega$  notation to suppress lower-order and constant terms.

Our proof will condition on the event  $\mathcal{E}$ . In other words, We assume  $\mathcal{E}$  happens, and makes derivations. After all of the derivations say in Proposition 4 and Proposition 5, we will get some high probability bounds. We can apply a union bound and take the probability bound  $1 - \exp(-\Omega(n - \log m))$  for  $\mathcal{E}$  into account.

Note that  $\mathbf{r}^* = \mathbf{A}\mathbf{x}^* - \mathbf{y}$ , so we have

$$\begin{aligned} \mathbf{x}^{(t+1)} &= (\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{y} \\ &= (\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{W}^{(t+1)} (\mathbf{A}\mathbf{x}^* - \mathbf{r}^*) \\ &= \mathbf{x}^* - (\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{A})^{-1} \mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*. \end{aligned}$$

This gives us an upper bound on the distance between  $\mathbf{x}^{(t+1)}$  and  $\mathbf{x}^*$ :

$$\|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2 \leq \frac{\|\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*\|_2}{\lambda_{\min}(\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{A})} \quad (41)$$

We can finish the proof by invoking Proposition 4 and Proposition 5, which provides high probability bounds for  $\|\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*\|_2$  and  $\lambda_{\min}(\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{A})$  respectively. Specifically, conditioned on  $\mathcal{E}$  (39), Proposition 4 gives

$$\|\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*\|_2 \leq \sqrt{5} \cdot 2^{2-p} \cdot \left(\frac{1}{r_{\min+}^*}\right)^{1-p} \cdot \sqrt{k} \cdot (1.01\sqrt{k} + \sqrt{n})$$

with probability at least  $1 - \exp(-\Omega(k - n)) = 1 - P'_1$ . On the other hand, by assumption we have  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \epsilon^{(t)}$ . Therefore, conditioned on  $\mathcal{E}$  (39), we can invoke Proposition 5 and obtain

$$\lambda_{\min}(\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{A}) \geq 0.99 \cdot 0.516 \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}$$

with probability at least  $1 - \exp(-\Omega(m - k - n \log n)) = 1 - P'_2$ . Combining the above with a union bound, we have

$$\begin{aligned} \|\mathbf{x}^{(t+1)} - \mathbf{x}^*\|_2 &\leq \frac{\sqrt{5} \cdot 2^{2-p}}{0.99 \cdot 0.516} \cdot \frac{1}{(r_{\min+}^*)^{1-p}} \cdot \frac{\sqrt{k} \cdot (1.01\sqrt{k} + \sqrt{n})}{(m-k)} \cdot (\epsilon^{(t)})^{2-p} \\ &= \alpha \cdot (\epsilon^{(t)})^{2-p} \\ &\leq \beta \cdot (\epsilon^{(t)})^{2-p}. \end{aligned}$$

Here we used the definition of  $\alpha$  and the condition  $\alpha \leq \beta$ . The proof is complete.  $\square$



### G.1. Proposition 4 and Its Proof

**Proposition 4.** Assume  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \epsilon^{(t)}$  and  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0, 1)$  entries. Assume event  $\mathcal{E}$  (39) happens. Then

$$\|\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*\|_2 \leq \sqrt{5} \cdot 2^{2-p} \cdot \left(\frac{1}{r_{\min+}^*}\right)^{1-p} \cdot \sqrt{k} \cdot (1.01\sqrt{k} + \sqrt{n})$$

with probability at least  $1 - \exp(-\Omega(k - n))$ .

*Proof.* Define  $\mathbf{z}^{(t)} := \mathbf{x}^{(t)} - \mathbf{x}^*$  and  $\bar{\mathbf{z}}^{(t)} := \mathbf{z}^{(t)} / \|\mathbf{z}^{(t)}\|_2$ . Define  $r_i^* := \mathbf{a}_i^\top \mathbf{x}^* - y_i$ . We first bound  $\|\mathbf{W}^{(t+1)} \mathbf{r}^*\|_2^2$ . Recall that  $S^*$  is the support of  $\mathbf{r}^*$  and  $0 < r_{\min+}^* \leq |r_i^*|$  for any  $i \in S^*$ . Then we have

$$\|\mathbf{W}^{(t+1)} \mathbf{r}^*\|_2^2 = \sum_{i \in S^*} (w_i^{(t+1)} r_i^*)^2 \leq \left(\frac{1}{r_{\min+}^*}\right)^{2-2p} \sum_{i \in S^*} (w_i^{(t+1)})^2 \cdot (r_i^*)^2 \cdot (r_i^*)^{2-2p}.$$

We need to bound  $(w_i^{(t+1)})^2 \cdot (r_i^*)^{4-2p}$  for every  $i \in S^*$ . To do so, we consider three cases.

**Case 1.** If  $|r_i^*| > 2 \cdot |\mathbf{a}_i^\top \mathbf{z}^{(t)}|$  then we have

$$\begin{aligned} |\mathbf{a}_i^\top \mathbf{x}^{(t)} - y_i| &= |\mathbf{a}_i^\top \mathbf{x}^{(t)} - \mathbf{a}_i^\top \mathbf{x}^* + r_i^*| \\ &= |\mathbf{a}_i^\top \mathbf{z}^{(t)} + r_i^*| \\ &\geq |r_i^*| - |\mathbf{a}_i^\top \mathbf{z}^{(t)}| \\ &\geq |r_i^*|/2, \end{aligned}$$

which by definition of the weights  $w_i^{(t+1)}$  of Algorithm 1 implies

$$(w_i^{(t+1)})^2 \cdot (r_i^*)^{4-2p} \leq \left(\frac{|r_i^*|}{|\mathbf{a}_i^\top \mathbf{x}^{(t)} - y_i|}\right)^{4-2p} \leq 2^{4-2p} = 4^{2-p}.$$

**Case 2.** If  $|\mathbf{a}_i^\top \mathbf{z}^{(t)}| \neq 0$  and  $|r_i^*| \leq 2 \cdot |\mathbf{a}_i^\top \mathbf{z}^{(t)}|$ , then, with  $w_i^{(t+1)} \leq (\epsilon^{(t)})^{p-2}$  (Algorithm 1), the assumption  $\|\mathbf{z}^{(t)}\|_2 \leq \epsilon^{(t)}$ , and the definition  $\bar{\mathbf{z}}^{(t)} := \mathbf{z}^{(t)} / \|\mathbf{z}^{(t)}\|_2$ , we obtain

$$\begin{aligned} (w_i^{(t+1)})^2 \cdot (r_i^*)^{4-2p} &\leq 2^{4-2p} \cdot \frac{(\mathbf{a}_i^\top \mathbf{z}^{(t)})^{4-2p}}{(\epsilon^{(t)})^{4-2p}} \\ &\leq 4^{2-p} \cdot (\mathbf{a}_i^\top \bar{\mathbf{z}}^{(t)})^{4-2p}. \end{aligned}$$

**Case 3.** If  $|\mathbf{a}_i^\top \mathbf{z}^{(t)}| = 0$ , then Case 1 implies that  $(w_i^{(t+1)})^2 \cdot (r_i^*)^{4-2p} \leq 1$ .

Combining the above, we obtain that

$$\begin{aligned} \|\mathbf{W}^{(t+1)} \mathbf{r}^*\|_2^2 &\leq \left(\frac{1}{r_{\min+}^*}\right)^{2-2p} \sum_{i \in S^*} \max \left\{ 1, 4^{2-p}, 4^{2-p} \cdot (\mathbf{a}_i^\top \bar{\mathbf{z}}^{(t)})^{4-2p} \right\} \\ &= \left(\frac{1}{r_{\min+}^*}\right)^{2-2p} \sum_{i \in S^*} \max \left\{ 4^{2-p}, 4^{2-p} \cdot (\mathbf{a}_i^\top \bar{\mathbf{z}}^{(t)})^{4-2p} \right\} \\ &\leq 4^{2-p} \cdot \left(\frac{1}{r_{\min+}^*}\right)^{2-2p} \cdot \left(k + \sum_{i \in S^*} (\mathbf{a}_i^\top \bar{\mathbf{z}}^{(t)})^{4-2p}\right) \\ &\stackrel{(i)}{\leq} 4^{2-p} \cdot \left(\frac{1}{r_{\min+}^*}\right)^{2-2p} \cdot (k + 4k) \\ &= 5 \cdot 4^{2-p} \cdot \left(\frac{1}{r_{\min+}^*}\right)^{2-2p} \cdot k \end{aligned}$$

where step (i) follows from Lemma 5 which assumes event  $\mathcal{E}$  happens and it holds with probability at least  $1 - \exp(-\Omega(k - n))$ . Therefore, we have

$$\|\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*\|_2 \leq \|\mathbf{W}^{(t+1)} \mathbf{r}^*\|_2 \cdot \sqrt{\lambda_{\max} \left( \sum_{i \in S^*} \mathbf{a}_i \mathbf{a}_i^\top \right)},$$

we can invoke Lemma 11 with  $\delta_1 = 0.01$  to upper bound the maximum eigenvalue of  $\sum_{i \in S^*} \mathbf{a}_i \mathbf{a}_i^\top$ . Therefore, invoking a union bound, it holds with probability at least  $1 - \exp(-\Omega(k - n)) - \exp(-\Omega(k)) = 1 - \exp(-\Omega(k - n))$  that

$$\|\mathbf{A}^\top \mathbf{W}^{(t+1)} \mathbf{r}^*\|_2 \leq \sqrt{5} \cdot 2^{2-p} \cdot \left( \frac{1}{r_{\min}^*} \right)^{1-p} \cdot \sqrt{k} \cdot (1.01\sqrt{k} + \sqrt{n}).$$

The proof is now complete.  $\square$

### G.1.1 Auxiliary Lemmas for Proposition 4

**Lemma 3** (Locally Lipschitz Continuity). *Let  $p \in [0, 1]$  and  $c > 0$ . The function  $\phi : [-d, d] \rightarrow \mathbb{R}$  with  $\phi(a) = a^{4-2p}$  is convex and differentiable, and thus for any  $a, b \in [-d, d]$  we have*

$$\phi(a) - \phi(b) \leq \phi'(a) \cdot (a - b) \leq (4 - 2p) \cdot d^{3-2p} \cdot |a - b|.$$

In other words,  $\phi$  is Lipschitz over  $[-d, d]$  with constant  $(4 - 2p) \cdot d^{3-2p}$ .

*Proof.* This follows directly from the definition of convexity for differentiable functions, and the fact that  $\phi'(a) = (4 - 2p) \cdot a^{3-2p}$  is bounded above by  $(4 - 2p) \cdot c^{3-2p}$ .  $\square$

**Lemma 4** (Moments of Gaussian Scalar). *For  $p \in [0, 1]$  and a standard Gaussian variable  $u \sim \mathcal{N}(0, 1)$ , we have*

$$\mathbb{E}[u^q] \leq 4.$$

*Proof.* we have  $\mathbb{E}[u^2] = 1$  and  $\mathbb{E}[u^4] = (4 - 1)!! = 3$ .<sup>6</sup> Therefore, for  $p \in [0, 1]$  it holds that

$$\begin{aligned} \mathbb{E}[u^{4-2p}] &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} u^{4-2p} \exp(-u^2/2) du \\ &= \frac{\sqrt{2}}{\sqrt{\pi}} \left( \int_0^1 u^{4-2p} \exp(-u^2/2) du + \int_1^{\infty} u^{4-2p} \exp(-u^2/2) du \right) \\ &\leq \frac{\sqrt{2}}{\sqrt{\pi}} \left( \int_0^1 u^2 \exp(-u^2/2) du + \int_1^{\infty} u^4 \exp(-u^2/2) du \right) \\ &\leq \frac{\sqrt{2}}{\sqrt{\pi}} \left( \int_0^{\infty} u^2 \exp(-u^2/2) du + \int_0^{\infty} u^4 \exp(-u^2/2) du \right) \\ &= \mathbb{E}[u^2] + \mathbb{E}[u^4] = 4. \end{aligned}$$

The proof is complete. Note that, mathematically, this bound 4 might be sub-optimal, as one would expect that  $\mathbb{E}[u^{4-2p}]$  takes values between  $\mathbb{E}[u^2] = 1$  and  $\mathbb{E}[u^4] = 3$ . However, it is harmless for our purpose as it is just a constant.  $\square$

**Lemma 5.** *Suppose  $\mathbf{a}_i \in \mathbb{R}^n$  has i.i.d.  $\mathcal{N}(0, 1)$  entries,  $\forall i = 1, \dots, k$ . Assume event  $\mathcal{E}$  (39) happens and  $p \in [0, 1]$ . We have*

$$\max_{\mathbf{z} \in \mathbb{S}^{n-1}} \sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z})^{4-2p} \leq 4k \quad (42)$$

with probability at least  $1 - \exp(-\Omega(k - n))$ .

<sup>6</sup>More generally, it holds that  $\mathbb{E}[u^{2q}] = (2q - 1)!!$  for every positive integer  $q$ . Here  $(\cdot)!!$  denotes the *double factorial* of some number.

*Proof.* Note that event  $\mathcal{E}$  (39) happens by assumption. In the proof we implicitly condition on  $\mathcal{E}$  (39). First, by (39), we know that  $(\mathbf{a}_1^\top \mathbf{z})^{4-2p}, \dots, (\mathbf{a}_k^\top \mathbf{z})^{4-2p}$  are independent random variables, and they are bounded with  $0 \leq |\mathbf{a}_i^\top \mathbf{z}|^{4-2p} \leq \max_{i=1, \dots, m} \|\mathbf{a}_i\|_2^{4-2p} \leq M^{4-2p}$ .

**Pointwise Bound.** We first consider a fixed  $\mathbf{z} \in \mathbb{S}^{n-1}$ , and derive a high probability upper bound for  $\sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z})^{4-2p}$ . Since  $\mathbf{a}_i^\top \mathbf{z} \sim \mathcal{N}(0, 1)$  for every  $i = 1, \dots, k$ , it follows from Lemma 4 that

$$\sum_{i=1}^k \mathbb{E}[(\mathbf{a}_i^\top \mathbf{z})^{4-2p}] \leq 4k. \quad (43)$$

By the standard Hoeffding bound on independent zero mean bounded random variables (i), for any  $\delta > 0$ , we have

$$\mathbb{P}\left(\sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z})^{4-2p} \geq (4 + \delta)k\right) \stackrel{(43)}{\leq} \mathbb{P}\left(\sum_{i=1}^k \left((\mathbf{a}_i^\top \mathbf{z})^{4-2p} - \mathbb{E}[(\mathbf{a}_i^\top \mathbf{z})^{4-2p}]\right) \geq \delta k\right) \quad (44)$$

$$\stackrel{(i)}{\leq} \exp\left(-\frac{2k\delta^2}{M^{8-4p}}\right) =: P_\delta \quad (45)$$

**Union Bound.** Similarly to the proof of Lemma 8 we consider a  $(1/4)$ -net  $\mathcal{N}_{0.25}$  of the sphere  $\mathbb{S}^{n-1}$ ;  $\mathcal{N}_{0.25}$  has at most  $9^n$  points. A direct application of the union bound yields

$$\sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z})^{4-2p} \geq (4 + \delta)k, \quad \forall \mathbf{z} \in \mathcal{N}_{0.25},$$

with probability at most  $9^n \cdot P_\delta$ .

**Approximation Bound.** With (43), we now bound  $\sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z})^{4-2p}$  for every  $\mathbf{z} \in \mathbb{S}^{n-1}$ . Indeed, for any  $\mathbf{z}_1 \in \mathbb{S}^n$ , there is some  $\mathbf{z}_2 \in \mathcal{N}_{0.25}$  such that  $\|\mathbf{z}_1 - \mathbf{z}_2\|_2 \leq 1/4$ . Lemma 3 implies that

$$\left| \sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z}_1)^{4-2p} - \sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z}_2)^{4-2p} \right| \leq \sum_{i=1}^k (4-2p) \cdot M^{3-2p} \cdot |\mathbf{a}_i^\top \mathbf{z}_1 - \mathbf{a}_i^\top \mathbf{z}_2| \leq (1-p/2) \cdot M^{4-2p} \cdot k. \quad (46)$$

Now, take (45) with  $\delta = (1-p/2) \cdot M^{4-2p}$ , and we have  $P_\delta = \exp(-k(2-p)^2/2)$ . Combining (45) with (46) gives

$$\mathbb{P}\left(\sum_{i=1}^k (\mathbf{a}_i^\top \mathbf{z}_1)^{4-2p} \geq 4k\right) \leq 9^n \exp(-k(2-p)^2/2), \quad \forall \mathbf{z}_1 \in \mathbb{S}^{n-1}.$$

In other words, (42) holds with probability at least  $1 - \exp(-\Omega(k-n))$ . The proof is complete.  $\square$

## G.2. Proposition 5 and Its Proof

**Proposition 5.** Let  $\mathbb{B}^{(t+1)}(\mathbf{x}^*; \epsilon^{(t)})$  be the  $\ell_2$ -ball centered at the ground truth  $\mathbf{x}^*$  with radius  $\epsilon^{(t)}$ , that is

$$\mathbb{B}^{(t+1)}(\mathbf{x}^*; \epsilon^{(t)}) := \left\{ \mathbf{x} \in \mathbb{R}^n : \|\mathbf{x} - \mathbf{x}^*\|_2 \leq \epsilon^{(t)} \right\}.$$

Associate each vector  $\mathbf{x}$  in the ball with an  $m \times m$  diagonal matrix of weights,  $\mathbf{W} = \text{diag}(w_1, \dots, w_m)$ , where  $w_i = \max\{|\mathbf{a}_i^\top \mathbf{x} - y_i|, \epsilon^{(t)}\}^{p-2}$ . This gives a product set

$$\mathcal{Q}_p^{(t+1)} := \left\{ (\mathbf{W}, \mathbf{x}) = (\text{diag}(w_1, \dots, w_m), \mathbf{x}) \in \mathbb{R}^{m \times m} \times \mathbb{B}^{(t+1)}(\mathbf{x}^*; \epsilon^{(t)}) : w_i = \max\{|\mathbf{a}_i^\top \mathbf{x} - y_i|, \epsilon^{(t)}\}^{p-2} \right\}. \quad (47)$$

Assume event  $\mathcal{E}$  (39) happens. Then we have

$$\lambda_{\min}(\mathbf{A}^\top \mathbf{W} \mathbf{A}) \geq 0.99 \cdot 0.516 \cdot (m-k) \cdot (\epsilon^{(t)})^{p-2}, \quad \forall (\mathbf{W}, \mathbf{x}) \in \mathcal{Q}_p^{(t+1)},$$

with probability at least  $1 - \exp(-\Omega(m-k-n \log n))$ .

Thus, the same conclusion holds for the iterate  $(\mathbf{W}^{(t+1)}, \mathbf{x}^{(t)})$  of our Algorithm 1 with at least the same probability, as long as  $(\mathbf{W}^{(t+1)}, \mathbf{x}^{(t)}) \in \mathcal{Q}_p^{(t+1)}$  or equivalently  $\|\mathbf{x}^{(t)} - \mathbf{x}^*\|_2 \leq \epsilon^{(t)}$ .

*Proof.* Note that we assume event  $\mathcal{E}$  (39) takes place. The derivation here implicitly conditions on this event.

Since  $\mathbf{A}^\top \mathbf{W} \mathbf{A} = \sum_{i=1}^m w_i \mathbf{a}_i \mathbf{a}_i^\top \succeq \sum_{i \in (S^*)^c} w_i \mathbf{a}_i \mathbf{a}_i^\top$ , we will bound the minimum eigenvalue of the latter. Recall the definition of  $M$  in (39), and consider a  $\tau$ -net  $\mathcal{N}_\tau$  over the ball  $\mathbb{B}^{(t+1)}(\mathbf{x}^*; \epsilon^{(t)})$ , with

$$\tau = \frac{0.516 \epsilon^{(t)}}{400 \cdot (2-p) \cdot M^3}.$$

Since event  $\mathcal{E}$  happens (39), we have  $M \in [\sqrt{0.5n}, 5n]$ , therefore

$$\frac{\tau M}{\epsilon^{(t)}} = \frac{0.516}{400 \cdot (2-p) \cdot M^2} \leq \frac{0.516}{400 \cdot (2-p) \cdot n/2} \leq 1, \quad (48)$$

$$\frac{\tau}{\epsilon^{(t)}} = \frac{0.516}{400 \cdot (2-p) \cdot M^3} \leq 1. \quad (49)$$

Note that Lemma 13 and (49) imply that this  $\tau$ -net is of size at most

$$T_\tau := \left( \frac{3\epsilon^{(t)}}{\tau} \right)^n = \left( \frac{1200 \cdot (2-p) \cdot M^3}{0.516} \right)^n.$$

Note that every  $\mathbf{x}_\tau \in \mathcal{N}_\tau$  induces a weight matrix  $\mathbf{W}_\tau = \text{diag}(w_{\tau 1}, \dots, w_{\tau m})$  such that  $(\mathbf{W}_\tau, \mathbf{x}_\tau) \in \mathcal{Q}_p^{(t+1)}$ . We can now invoke Lemma 8 with a union bound, to obtain

$$\lambda_{\min} \left( \sum_{i \in (S^*)^c} w_{\tau i} \mathbf{a}_i \mathbf{a}_i^\top \right) \geq 0.995 \cdot 0.516 \cdot (m-k) \cdot (\epsilon^{(t)})^{p-2}, \quad \forall \mathbf{x}_\tau \in \mathcal{N}_\tau,$$

which holds with probability at least  $1 - P_\tau$ , where  $P_\tau := T_\tau \cdot \exp(-\Omega(m-k-n))$  with  $\Omega(\cdot)$  hiding constants from Lemma 8. Since for any  $(\mathbf{W}, \mathbf{x}) \in \mathcal{Q}_p^{(t+1)}$ , there exists some  $\mathbf{x}_\tau \in \mathcal{N}_\tau$  with  $\|\mathbf{x}_\tau - \mathbf{x}\|_2 \leq \tau$ , combine this with (48) and we can now invoke Lemma 9, which implies

$$\begin{aligned} \left| \lambda_{\min} \left( \sum_{i \in (S^*)^c} w_i \mathbf{a}_i \mathbf{a}_i^\top \right) - \lambda_{\min} \left( \sum_{i \in (S^*)^c} w_{\tau i} \mathbf{a}_i \mathbf{a}_i^\top \right) \right| &\leq 2 \cdot (2-p) \cdot (m-k) \cdot \frac{\tau M^3}{(\epsilon^{(t)})^{3-p}} \\ &= 0.005 \cdot 0.516 \cdot (m-k) \cdot (\epsilon^{(t)})^{p-2}. \end{aligned}$$

Combining the above with a union bound, we now have

$$\lambda_{\min} \left( \sum_{i \in (S^*)^c} w_i \mathbf{a}_i \mathbf{a}_i^\top \right) \geq 0.99 \cdot 0.516 \cdot (m-k) \cdot (\epsilon^{(t)})^{p-2}, \quad \forall (\mathbf{W}, \mathbf{x}) \in \mathcal{Q}_p^{(t+1)}$$

with probability at least  $1 - P_\tau$ . To finish the proof, it remains to simplify the expression of  $P_\tau = T_\tau \cdot \exp(-\Omega(m-k-n))$ . Note that  $M \leq \sqrt{5n}$ , which implies  $T_\tau = O(M^{3n}) = O(n^{3n/2})$  and  $\log T_\tau = O(n \log n)$ . Therefore

$$\begin{aligned} P_\tau &= \exp(-\Omega(m-k-n) + \log T_\tau) \\ &= \exp(-\Omega(m-k-n - n \log n)) \\ &= \exp(-\Omega(m-k-n \log n)). \end{aligned}$$

The proof of Proposition 5 is now complete. □

### G.2.1 Auxiliary Lemmas for Proposition 5

**Lemma 6** (Binomial Approximation). *For any  $a \in [0, 1]$  and  $p \in [0, 1]$ , we have:*

$$(1+a)^{2-p} - 1 \leq 2 \cdot (2-p) \cdot a$$

*Proof.* Let  $p$  be fixed. Define a function  $f : [0, 1] \rightarrow \mathbb{R}$  such that

$$f(a) = (1 + a)^{2-p} - 1 - 2 \cdot (2 - p) \cdot a.$$

We see that  $f(0) = 0$ , and

$$\begin{aligned} f'(a) &= (2 - p) \cdot (1 + a)^{1-p} - 2 \cdot (2 - p) \\ &= (2 - p) \cdot ((1 + a)^{1-p} - 2). \end{aligned}$$

Since  $1 \leq 1 + a \leq 2$  for  $0 \leq a \leq 1$  and  $1 - p \in [0, 1]$ , we know that  $(1 + a)^{1-p} \leq 2$  and thus  $f'(a) \leq 0$ . Hence,  $f$  is a decreasing function in  $[0, 1]$  with  $f(a) \leq f(0) = 0$  for any  $a \in [0, 1]$ .  $\square$

**Lemma 7** (Pointwise Expectation). *Suppose for every  $i = 1, \dots, m$ ,  $\mathbf{a}_i \in \mathbb{R}^n$  has i.i.d.  $\mathcal{N}(0, 1)$  entries. Recall the definition of the product set  $\mathcal{Q}_p^{(t+1)}$  (47). Let  $\mathbf{v} \in \mathbb{S}^{n-1}$  and  $(\mathbf{W}, \mathbf{x}) = (\text{diag}(w_1, \dots, w_m), \mathbf{x}) \in \mathcal{Q}_p^{(t+1)}$  be fixed. Then it holds that*

$$(m - k) \cdot (\epsilon^{(t)})^{p-2} \geq \mathbb{E} \left[ \sum_{i \in (S^*)^c} w_i (\mathbf{a}_i^\top \mathbf{v})^2 \right] \geq 0.516 \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}. \quad (50)$$

*Proof.* The first inequality of (50) is obvious, and we now prove the second. For any  $i \in (S^*)^c$  we have  $y_i = \mathbf{a}_i^\top \mathbf{x}^*$ . Note that  $\|\mathbf{x} - \mathbf{x}^*\|_2 \leq \epsilon^{(t)}$ . Define  $\mathbf{z} = \mathbf{x} - \mathbf{x}^*$  and  $\bar{\mathbf{z}} = \mathbf{z} / \|\mathbf{z}\|_2$ . Then

$$\begin{aligned} w_i^{(t+1)} &= \max \{ |\mathbf{a}_i^\top \mathbf{x} - y_i|, \epsilon^{(t)} \}^{p-2} \\ &= (\epsilon^{(t)})^{p-2} \cdot \max \left\{ \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\epsilon^{(t)}}, 1 \right\}^{p-2} \\ &= (\epsilon^{(t)})^{p-2} \cdot \max \left\{ \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \cdot \frac{\|\mathbf{z}\|_2}{\epsilon^{(t)}}, 1 \right\}^{p-2} \\ &\geq (\epsilon^{(t)})^{p-2} \cdot \max \{ |\mathbf{a}_i^\top \bar{\mathbf{z}}|, 1 \}^{p-2}. \end{aligned}$$

This leads to

$$\begin{aligned} \mathbb{E} \left[ \sum_{i \in (S^*)^c} w_i (\mathbf{a}_i^\top \mathbf{v})^2 \right] &\geq (\epsilon^{(t)})^{p-2} \cdot \mathbb{E} \left[ \sum_{i \in (S^*)^c} \max \{ |\mathbf{a}_i^\top \bar{\mathbf{z}}|, 1 \}^{p-2} \cdot (\mathbf{a}_i^\top \mathbf{v})^2 \right] \\ &= (\epsilon^{(t)})^{p-2} \cdot (m - k) \cdot \mathbb{E}_{\mathbf{a} \sim \mathcal{N}(0, \mathbf{I}_n)} \left[ \frac{(\mathbf{a}^\top \mathbf{v})^2}{\max \{ |\mathbf{a}^\top \bar{\mathbf{z}}|, 1 \}^{2-p}} \right] \\ &\geq (\epsilon^{(t)})^{p-2} \cdot (m - k) \cdot c_p, \end{aligned}$$

where  $c_p$  is defined as

$$c_p := \inf_{\mathbf{u} \in \mathbb{S}^{n-1}, \mathbf{q} \in \mathbb{S}^{n-1}} \mathbb{E}_{\mathbf{a} \sim \mathcal{N}(0, \mathbf{I}_n)} \left[ \frac{(\mathbf{a}^\top \mathbf{q})^2}{\max \{ |\mathbf{a}^\top \mathbf{u}|, 1 \}^{2-p}} \right].$$

It then remains to show that  $c_p \geq 0.516$ . By rotation invariance, we can assume without loss of generality that  $\mathbf{u} =$

$[u_1, 0, \dots, 0]^\top$  and  $\mathbf{q} = [q_1, q_2, 0, \dots, 0]^\top$ . Then we have

$$\begin{aligned}
c_p &= \inf_{\substack{u_1: u_1^2=1 \\ q_1, q_2: q_1^2+q_2^2=1}} \mathbb{E}_{a_1 \sim \mathcal{N}(0,1), a_2 \sim \mathcal{N}(0,1)} \left[ \frac{a_1^2 q_1^2 + a_2^2 q_2^2 + 2a_1 a_2 q_1 q_2}{\max\{|a_1 u_1|, 1\}^{2-p}} \right] \\
&= \inf_{q_1, q_2: q_1^2+q_2^2=1} \mathbb{E}_{a_1 \sim \mathcal{N}(0,1)} \left[ \frac{a_1^2 q_1^2 + q_2^2}{\max\{|a_1|, 1\}^{2-p}} \right] \\
&= \frac{2}{\sqrt{2\pi}} \cdot \inf_{q_1, q_2: q_1^2+q_2^2=1} \int_0^1 (t^2 q_1^2 + q_2^2) e^{-t^2/2} dt + \int_1^\infty \frac{t^2 q_1^2 + q_2^2}{t^{2-p}} \cdot e^{-t^2/2} dt \\
&\geq \inf_{q_1, q_2: q_1^2+q_2^2=1} 0.516q_1^2 + 0.849q_2^2 \\
&\geq 0.516.
\end{aligned}$$

In the above, the (lower bound on the) integral was calculated by the MATLAB function `integral` with  $p = 0$ . This finishes the proof.  $\square$

**Lemma 8** (Pointwise Concentration). *Suppose  $\mathbf{A} \in \mathbb{R}^{m \times n}$  has i.i.d.  $\mathcal{N}(0, 1)$  entries. Let  $S^*$  be the support of  $\mathbf{A}\mathbf{x}^* - \mathbf{y}$ . Recall the definition of the product set  $\mathcal{Q}_p^{(t+1)}$  (47). With  $(\mathbf{W}, \mathbf{x}) = (\text{diag}(w_1, \dots, w_m), \mathbf{x}) \in \mathcal{Q}_p^{(t+1)}$  fixed, we have*

$$\lambda_{\min} \left( \sum_{i \in (S^*)^c} w_i \mathbf{a}_i \mathbf{a}_i^\top \right) \geq 0.995 \cdot 0.516 \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}$$

with probability at least  $1 - \exp(-\Omega(m - k - n))$ .

*Proof.* Define  $\mathbf{B}_i := w_i \mathbf{a}_i \mathbf{a}_i^\top$  and  $\mathbf{B} := \sum_{i \in (S^*)^c} \mathbf{B}_i$ . We now bound the minimum eigenvalue of  $\mathbf{B}$ . We first derive a high probability lower bound for  $\sum_{i \in (S^*)^c} g_i^2$ , where  $g_i := \sqrt{w_i} \mathbf{a}_i^\top \mathbf{v}$  with a fixed spherical vector  $\mathbf{v} \in \mathbb{S}^{n-1}$ . Then we apply a union bound.

**Bounding  $\sum_{i \in (S^*)^c} g_i^2$ .** Since  $w_i \leq (\epsilon^{(t)})^{p-2}$ , the sub-Gaussian norm  $\|g_i\|_{\psi_2}$  can be bounded above:

$$\begin{aligned}
\|g_i\|_{\psi_2} &\leq \|(\epsilon^{(t)})^{p/2-1} \cdot \mathbf{a}_i^\top \mathbf{v}\|_{\psi_2} \\
&= (\epsilon^{(t)})^{p/2-1} \cdot \|\mathbf{a}_i^\top \mathbf{v}\|_{\psi_2} \\
&= (\epsilon^{(t)})^{p/2-1} \cdot \|a\|_{\psi_2},
\end{aligned}$$

where  $a \sim \mathcal{N}(0, 1)$ , whose sub-Gaussian norm  $\|a\|_{\psi_2}$  is bounded above by some constant  $C_1$ , so we have  $\|g_i\|_{\psi_2} \leq (\epsilon^{(t)})^{p/2-1} \cdot C_1$ . Thus,  $g_i$  is a sub-Gaussian random variable with parameter  $(\epsilon^{(t)})^{p/2-1} \cdot C_1$ , which implies that  $g_i^2 = w_i \cdot (\mathbf{a}_i^\top \mathbf{v})^2$  is sub-exponential with parameter  $(\epsilon^{(t)})^{p-2} \cdot C_1^2$ . Furthermore, using the centering technique we know that  $g_i^2 - \mathbb{E}[g_i^2]$  is also sub-exponential, with parameter  $2(\epsilon^{(t)})^{p-2} \cdot C_1^2$ . Invoking Bernstein's inequality (Lemma 12), we have for any  $\delta > 0$  that

$$\left| \sum_{i \in (S^*)^c} (g_i^2 - \mathbb{E}[g_i^2]) \right| \geq \delta \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}$$

with probability at most

$$\begin{aligned}
P_\delta &:= 2 \cdot \exp \left[ -C_2 \cdot \min \left\{ \frac{\delta^2 \cdot (m - k)^2 \cdot (\epsilon^{(t)})^{2p-4}}{4(\epsilon^{(t)})^{2p-4} \cdot C_1^4 \cdot (m - k)}, \frac{\delta \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}}{2 \cdot (\epsilon^{(t)})^{p-2} \cdot C_1^2} \right\} \right] \\
&= 2 \cdot \exp \left[ -C_2 \cdot \min \left\{ \frac{\delta^2}{4 \cdot C_1^4}, \frac{\delta}{2 \cdot C_1^2} \right\} \cdot (m - k) \right],
\end{aligned}$$

where  $C_2$  is a universal constant.

**Union Bound.** Consider a  $(1/4)$ -net  $\mathcal{N}_{0.25}$  of the sphere  $\mathbb{S}^{n-1}$ , which has at most  $9^n$  points (Lemma 13). Lemma 14 implies

$$\begin{aligned}\|\mathbf{B} - \mathbb{E}[\mathbf{B}]\|_2 &\leq 2 \cdot \max_{\mathbf{v} \in \mathcal{N}_{0.25}} |\mathbf{v}^\top (\mathbf{B} - \mathbb{E}[\mathbf{B}]) \mathbf{v}| \\ &= 2 \cdot \max_{\mathbf{v} \in \mathcal{N}_{0.25}} \sum_{i \in (S^*)^c} |g_i^2 - \mathbb{E}[g_i^2]|.\end{aligned}$$

Thus, combine the above with a union bound and to obtain

$$\|\mathbf{B} - \mathbb{E}[\mathbf{B}]\|_2 \leq 2 \cdot \delta \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}$$

with probability at least  $1 - 9^n \cdot P_\delta$ . Thus, with at least this probability we have

$$\mathbf{v}^\top \mathbf{B} \mathbf{v} \geq \mathbf{v}^\top \mathbb{E}[\mathbf{B}] \mathbf{v} - 2 \cdot \delta \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}, \forall \mathbf{v} \in \mathbb{S}^{n-1}.$$

In particular, for an eigenvector  $\mathbf{v}_0$  of  $\mathbf{B}$  that corresponds to its minimum eigenvalue, we have

$$\begin{aligned}\lambda_{\min}(\mathbf{B}) &\geq \mathbf{v}_0^\top \mathbb{E}[\mathbf{B}] \mathbf{v}_0 - 2 \cdot \delta \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2} \\ &\geq \lambda_{\min}(\mathbb{E}[\mathbf{B}]) - 2 \cdot \delta \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2} \\ &\geq 0.516 \cdot (1 - 2 \cdot \delta) \cdot (m - k) \cdot (\epsilon^{(t)})^{p-2}\end{aligned}$$

with probability at least  $1 - 9^n \cdot P_\delta$ . In the last step, we used the lower bound of  $\lambda_{\min}(\mathbb{E}[\mathbf{B}])$  derived in Lemma 7. Now, with  $\delta = 0.0025$ , and noticing that  $1 - 9^n \cdot P_\delta$  is of order  $\exp(-\Omega(m - k - n))$ , we finished the proof.  $\square$

**Lemma 9** (Lipschitz Continuity of Minimum Eigenvalues). *Recall  $M := \max_{i=1, \dots, m} \|\mathbf{a}_i\|_2$  (39) and the definition of the product set  $\mathcal{Q}_p^{(t+1)}$  (47). Let  $(\mathbf{W}_1, \mathbf{x}_1)$  and  $(\mathbf{W}_2, \mathbf{x}_2)$  be two elements of  $\mathcal{Q}_p^{(t+1)}$  with  $\|\mathbf{x}_1 - \mathbf{x}_2\|_2 \leq \tau$ . Write  $\mathbf{W}_1 = \text{diag}(w_{1i}, \dots, w_{1m})$  and  $\mathbf{W}_2 = \text{diag}(w_{2i}, \dots, w_{2m})$ . If  $\tau M \leq \epsilon^{(t)}$ , then*

$$\left| \lambda_{\min} \left( \sum_{i \in (S^*)^c} w_{1i} \mathbf{a}_i \mathbf{a}_i^\top \right) - \lambda_{\min} \left( \sum_{i \in (S^*)^c} w_{2i} \mathbf{a}_i \mathbf{a}_i^\top \right) \right| \leq 2 \cdot (2 - p) \cdot (m - k) \cdot \frac{\tau M^3}{(\epsilon^{(t)})^{3-p}}.$$

*Proof.* The proof is divided into four cases, depending on whether the  $i$ -th weights,  $w_{1i}$  and  $w_{2i}$ , are truncated by  $\epsilon^{(t)}$  or not.

**Case 1.** If both  $|\mathbf{a}_i^\top \mathbf{x}_1 - y_i| \leq \epsilon^{(t)}$  and  $|\mathbf{a}_i^\top \mathbf{x}_2 - y_i| \leq \epsilon^{(t)}$ , then we have  $w_{1i} = w_{2i}$ .

**Case 2.** Assume  $|\mathbf{a}_i^\top \mathbf{x}_1 - y_i| \leq \epsilon^{(t)}$  and  $|\mathbf{a}_i^\top \mathbf{x}_2 - y_i| > \epsilon^{(t)}$ . First note that

$$|\mathbf{a}_i^\top \mathbf{x}_2 - y_i| = |\mathbf{a}_i^\top \mathbf{x}_2 - \mathbf{a}_i^\top \mathbf{x}_1 + \mathbf{a}_i^\top \mathbf{x}_1 - y_i| \leq \|\mathbf{a}_i\|_2 \cdot \tau + |\mathbf{a}_i^\top \mathbf{x}_1 - y_i|. \quad (51)$$

Then we have

$$\begin{aligned}|w_{i1} - w_{i2}| &= \left| (\epsilon^{(t)})^{p-2} - |\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{p-2} \right| \\ &= \frac{1}{(\epsilon^{(t)})^{2-p}} - \frac{1}{|\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{2-p}} \\ &\leq \frac{1}{(\epsilon^{(t)})^{2-p}} - \frac{1}{(\|\mathbf{a}_i\|_2 \cdot \tau + |\mathbf{a}_i^\top \mathbf{x}_1 - y_i|)^{2-p}} \\ &\leq \frac{1}{(\epsilon^{(t)})^{2-p}} - \frac{1}{(\|\mathbf{a}_i\|_2 \cdot \tau + \epsilon^{(t)})^{2-p}} \\ &= \frac{1}{(\epsilon^{(t)})^{2-p}} \left( 1 - \frac{1}{\left( \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{\epsilon^{(t)}} + 1 \right)^{2-p}} \right) \\ &\leq \frac{1}{(\epsilon^{(t)})^{2-p}} \left( 1 - \frac{1}{1 + 2 \cdot (2 - p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{\epsilon^{(t)}}} \right).\end{aligned}$$

In the last step we used the assumption  $\|\mathbf{a}_i\|_2 \cdot \tau \leq \epsilon^{(t)}$  and Lemma 6. Then we get

$$\begin{aligned} |w_{i1} - w_{i2}| &\leq \frac{1}{(\epsilon^{(t)})^{2-p}} \cdot \frac{2 \cdot (2-p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{\epsilon^{(t)}}}{1 + 2 \cdot (2-p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{\epsilon^{(t)}}} \\ &= 2 \cdot (2-p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{(\epsilon^{(t)})^{3-p}} \cdot \frac{1}{1 + 2 \cdot (2-p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{\epsilon^{(t)}}} \\ &\leq 2 \cdot (2-p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{(\epsilon^{(t)})^{3-p}}. \end{aligned}$$

**Case 3.** If  $|\mathbf{a}_i^\top \mathbf{x}_1 - y_i| > \epsilon^{(t)}$  and  $|\mathbf{a}_i^\top \mathbf{x}_2 - y_i| \leq \epsilon^{(t)}$ , then we can bound  $|w_{i1} - w_{i2}|$  similarly as in Case 2 (by symmetry), and obtain the same upper bound.

**Case 4.** Suppose  $|\mathbf{a}_i^\top \mathbf{x}_1 - y_i| > \epsilon^{(t)}$  and  $|\mathbf{a}_i^\top \mathbf{x}_2 - y_i| > \epsilon^{(t)}$ , then, assuming  $|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p}$  is larger than or equal to  $|\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{2-p}$  without loss of generality, we have

$$\begin{aligned} |w_{i1} - w_{i2}| &= \left| \frac{1}{|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p}} - \frac{1}{|\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{2-p}} \right| \\ &= \frac{|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p} - |\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{2-p}}{|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p} \cdot |\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{2-p}} \\ &= \frac{\frac{|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p}}{|\mathbf{a}_i^\top \mathbf{x}_2 - y_i|^{2-p}} - 1}{|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p}}. \end{aligned}$$

Similarly in (51), we have  $|\mathbf{a}_i^\top \mathbf{x}_1 - y_i| \leq |\mathbf{a}_i^\top \mathbf{x}_2 - y_i| + \|\mathbf{a}_i\|_2 \cdot \tau$ , and consequently

$$\begin{aligned} |w_{i1} - w_{i2}| &\leq \frac{\left| \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{|\mathbf{a}_i^\top \mathbf{x}_2 - y_i|} + 1 \right|^{2-p} - 1}{|\mathbf{a}_i^\top \mathbf{x}_1 - y_i|^{2-p}} \\ &\leq \frac{\left| \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{\epsilon^{(t)}} + 1 \right|^{2-p} - 1}{(\epsilon^{(t)})^{2-p}} \\ &\leq 2 \cdot (2-p) \cdot \frac{\|\mathbf{a}_i\|_2 \cdot \tau}{(\epsilon^{(t)})^{3-p}} \end{aligned}$$

where we used the assumption  $\|\mathbf{a}_i\|_2 \cdot \tau \leq \epsilon^{(t)}$  and Lemma 6. To summarize, we have shown

$$\begin{aligned} \sum_{i \in (S^*)^c} |w_{i1} - w_{i2}| &\leq 2 \cdot (2-p) \cdot (m-k) \cdot \frac{\tau}{(\epsilon^{(t)})^{3-p}} \cdot \max_{i \in (S^*)^c} \|\mathbf{a}_i\|_2 \\ &\leq 2 \cdot (2-p) \cdot (m-k) \cdot \frac{\tau M}{(\epsilon^{(t)})^{3-p}}. \end{aligned}$$

Next, suppose without loss of generality that the minimum eigenvalue of  $\sum_{i \in (S^*)^c} w_{1i} \mathbf{a}_i \mathbf{a}_i^\top$  is larger than or equal to that of  $\sum_{i \in (S^*)^c} w_{2i} \mathbf{a}_i \mathbf{a}_i^\top$ , and denote by  $\mathbf{v}_0 \in \mathbb{S}^{n-1}$  an eigenvector of  $\sum_{i \in (S^*)^c} w_{2i} \mathbf{a}_i \mathbf{a}_i^\top$  corresponding to its minimum eigenvalue. Then we have

$$\begin{aligned} \lambda_{\min} \left( \sum_{i \in (S^*)^c} w_{1i} \mathbf{a}_i \mathbf{a}_i^\top \right) - \lambda_{\min} \left( \sum_{i \in (S^*)^c} w_{2i} \mathbf{a}_i \mathbf{a}_i^\top \right) &\leq \sum_{i \in (S^*)^c} w_{1i} (\mathbf{a}_i^\top \mathbf{v}_0)^2 - \sum_{i \in (S^*)^c} w_{2i} (\mathbf{a}_i^\top \mathbf{v}_0)^2 \\ &\leq \left( \max_{i \in (S^*)^c} (\mathbf{a}_i^\top \mathbf{v}_0)^2 \right) \sum_{i \in (S^*)^c} |w_{i1} - w_{i2}| \\ &\leq 2 \cdot (2-p) \cdot (m-k) \cdot \frac{\tau M^3}{(\epsilon^{(t)})^{3-p}} \end{aligned}$$



and we finished the proof.  $\square$

## H. Proof of Proposition 1

Similarly to (41), for the least-squares estimator  $\mathbf{x}^\dagger$  we have

$$\|\mathbf{x}^\dagger - \mathbf{x}^*\|_2 \leq \frac{\|\mathbf{A}^\top(\mathbf{A}\mathbf{x}^* - \mathbf{y})\|_2}{\lambda_{\min}(\mathbf{A}^\top\mathbf{A})} \leq \frac{\|\mathbf{A}\|_2 \cdot \|(\mathbf{A}\mathbf{x}^* - \mathbf{y})\|_2}{\lambda_{\min}(\mathbf{A}^\top\mathbf{A})}.$$

Invoke Lemma 11 with  $\delta_1 = \delta_2 = 0.01$  to bound  $\|\mathbf{A}\|_2$  and  $\lambda_{\min}(\mathbf{A}^\top\mathbf{A})$ , then use a union bound, and we get that the desired bound (15) holds with probability at least  $1 - \exp(-\Omega(k)) - \exp(-\Omega(m))$ . The proof is complete.

## I. Auxiliary and Known Results in High Dimensional Statistics

**Lemma 10** (Lemma 1 of [43]). *Assume  $\mathbf{a} \in \mathbb{R}^n$  has i.i.d.  $\mathcal{N}(0, 1)$  entries, then for any  $\delta_1 > 0$  and  $\delta_2 > 0$ , it holds that*

$$\begin{aligned} \mathbb{P}\left(\|\mathbf{a}\|_2^2 > n + 2\sqrt{n\delta_1} + 2\delta_1\right) &\leq \exp(-\delta_1) \\ \mathbb{P}\left(\|\mathbf{a}\|_2^2 < n - 2\sqrt{n\delta_2}\right) &\leq \exp(-\delta_2). \end{aligned}$$

In particular, set  $\delta_1 = n$  and  $\delta_2 = n/16$ , and we get

$$\begin{aligned} \mathbb{P}\left(\|\mathbf{a}\|_2^2 > 5n\right) &\leq \exp(-n), \\ \mathbb{P}\left(\|\mathbf{a}\|_2^2 < n/2\right) &\leq \exp(-n/16). \end{aligned}$$

Finally, if  $\mathbf{a}_i \in \mathbb{R}^n$  has i.i.d.  $\mathcal{N}(0, 1)$  entries for every  $i = 1, \dots, s$ , then

$$\mathbb{P}\left(\max_{i=1, \dots, s} \|\mathbf{a}_i\|_2^2 > 5n\right) \leq s \cdot \exp(-n).$$

**Lemma 11** (Theorem 6.1 and Example 6.2 of [74]). *Suppose  $\mathbf{a}_i \in \mathbb{R}^n$  has i.i.d.  $\mathcal{N}(0, 1)$  entries for each  $i = 1, \dots, m$ . Then, for any  $\delta_1 > 0$ , it holds with probability at most  $\exp(-k\delta_1^2/2)$  that*

$$\lambda_{\max}\left(\sum_{i=1}^k \mathbf{a}_i \mathbf{a}_i^\top\right) \geq ((1 + \delta_1) \cdot \sqrt{k} + \sqrt{n})^2.$$

Moreover, if  $m \geq n$  then, for any  $\delta_2 \in (0, 1)$ , it holds with probability at most  $\exp(-m\delta_2^2/2)$  that

$$\lambda_{\min}\left(\sum_{i=1}^m \mathbf{a}_i \mathbf{a}_i^\top\right) \leq ((1 - \delta_2) \cdot \sqrt{m} - \sqrt{n})^2.$$

**Lemma 12** (Bernstein's inequality, cf. Theorem 2.8.1 of [72]). *Let  $X_1, \dots, X_N$  be independent, mean zero, sub-exponential random variables. Then, for every  $\delta \geq 0$ , we have*

$$\mathbb{P}\left[\left|\sum_{i=1}^N X_i\right| \geq \delta\right] \leq 2 \exp\left[-C \min\left(\frac{\delta^2}{\sum_{i=1}^N \|X_i\|_{\psi_1}^2}, \frac{\delta}{\max_i \|X_i\|_{\psi_1}}\right)\right].$$

Here  $C$  is some absolute constant and  $\|\cdot\|_{\psi_1}$  denotes the sub-exponential norm.

**Lemma 13** (Corollary 4.2.13 of [72]). *For any  $\varepsilon > 0$ , an  $\varepsilon$ -net of the unit Euclidean ball of  $\mathbb{R}^n$  has at most  $(2/\varepsilon + 1)^n$  points. This is also true for the unit Euclidean sphere  $\mathbb{S}^{n-1}$ .*

**Lemma 14** (Exercise 4.4.3 of [72]). *For a symmetric matrix  $\mathbf{B}$  and an  $\varepsilon$ -net  $\mathcal{N}$  with  $\varepsilon \in [0, 1/2)$ , we have*

$$\|\mathbf{B}\|_2 \leq \frac{1}{1 - 2\varepsilon} \cdot \max_{\mathbf{v} \in \mathcal{N}} |\mathbf{v}^\top \mathbf{B} \mathbf{v}|.$$