

High Dimensional Robust Sparse Regression

Liu Liu
liuliu@utexas.edu

Yanyao Shen
shenyanyao@utexas.edu

Tianyang Li
lty@cs.utexas.edu

Constantine Caramanis
constantine@utexas.edu

The University of Texas at Austin

Abstract

We provide a novel – and to the best of our knowledge, the first – algorithm for high dimensional sparse regression with constant fraction of corruptions in explanatory and/or response variables. Our algorithm recovers the true sparse parameters with sub-linear sample complexity, in the presence of a constant fraction of arbitrary corruptions. Our main contribution is a robust variant of Iterative Hard Thresholding. Using this, we provide accurate estimators: when the covariance matrix in sparse regression is identity, our error guarantee is near information-theoretically optimal. We then deal with robust sparse regression with unknown structured covariance matrix. We propose a filtering algorithm which consists of a novel randomized outlier removal technique for robust sparse mean estimation that may be of interest in its own right: the filtering algorithm is flexible enough to deal with unknown covariance. Also, it is orderwise more efficient computationally than the ellipsoid algorithm. Using sub-linear sample complexity, our algorithm achieves the best known (and first) error guarantee. We demonstrate the effectiveness on large-scale sparse regression problems with arbitrary corruptions.

1 Introduction

Learning in the presence of arbitrarily (even adversarially) corrupted outliers in the training data has a long history in Robust Statistics [27, 23, 42], and has recently received much renewed attention. The high dimensional setting poses particular challenges as outlier removal via preprocessing is essentially impossible when the number of variables scales with the number of samples. We propose a computationally efficient estimator for outlier-robust sparse regression that has near-optimal sample complexity, and is the first algorithm resilient to a constant fraction of arbitrary outliers with corrupted covariates and/or response variables. Unless we specifically mention otherwise, all future mentions of outliers mean corruptions in covariates and/or response variables.

We assume that the authentic samples are independent and identically distributed (i.i.d.) drawn from an uncorrupted distribution P , where P represents the linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \xi_i$, where $\mathbf{x}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma})$, and $\boldsymbol{\beta}^* \in \mathbb{R}^d$ is the true parameter (see Section 1.3 for complete details and definitions). To model the corruptions, the adversary can choose an arbitrary ϵ -fraction of the authentic samples, and replace them with arbitrary values. We refer to the observations after corruption as ϵ -corrupted samples (Definition 1.1). This corruption model allows the adversary to select an ϵ -fraction of authentic samples to delete and corrupt, hence it is stronger than Huber’s ϵ -contamination model [26], where the adversary independently corrupts each sample with probability ϵ .

Outlier-robust regression is a classical problem within robust statistics (e.g., [39]), yet even in the low-dimensional setting, efficient algorithms robust to corruption in the covariates have proved elusive, until recent breakthroughs in [38, 15] and [31], which built on important results in Robust Mean Estimation [13, 33] and Sums of Squares [4], respectively.

In the sparse setting, the parameter $\boldsymbol{\beta}^*$ we seek to recover is also k -sparse, and a key goal is to provide recovery guarantees with sample complexity scaling with k , and *sublinearly* with d . Without

outliers, by now classical results (e.g., [18]) show that $n = \Omega(k \log d)$ samples from a i.i.d sub-Gaussian distribution are enough to give recovery guarantees on β^* with and without additive noise. These strong assumptions on the probabilistic distribution are necessary, since in the worst case, sparse recovery is known to be NP-hard [3, 48].

Sparsity recovery with a constant fraction of arbitrary corruption is fundamentally hard. For instance, to the best of our knowledge, there’s no previous work can provide exact recovery for sparse linear equations with arbitrary corruption in polynomial time. In contrast, a simple exhaustive search can easily enumerate the samples and recover the sparse parameter in exponential time.

In this work, we seek to give an efficient, sample-complexity optimal algorithm that recovers β^* to within accuracy depending on ϵ (the fraction of outliers). In the case of no additive noise, we are interested in algorithms that can guarantee exact recovery, independent of ϵ .

1.1 Related work

The last 10 years have seen a resurgence in interest in robust statistics, including the problem of resilience to outliers in the data. Important problems attacked have included PCA [32, 47, 46, 33, 13], and more recently robust regression (as in this paper) [38, 15, 17, 31] and robust mean estimation [13, 33, 1], among others. We focus now on the recent work most related to the present paper.

Robust regression. Earlier work in robust regression considers corruption only in the output, and shows that algorithms nearly as efficient as for regression without outliers succeeds in parameter recovery, even with a constant fraction of outliers [34, 37, 5, 6, 30]. Yet these algorithms (and their analysis) focus on corruption in y , and do not seem to extend to the setting of corrupted covariates – the setting of this work. In the low dimensional setting, there has been remarkable recent progress. The work in [31] shows that the Sum of Squares (SOS) based semidefinite hierarchy can be used for solving robust regression. Essentially concurrent to the SOS work, [12, 24, 38, 15] use robust gradient descent for empirical risk minimization, by using robust mean estimation as a subroutine to compute robust gradients at each iteration. [17] uses filtering algorithm [13] for robust regression. Computationally, these latter works scale better than the algorithms in [31], as although the Sum of Squares SDP framework gives polynomial time algorithms, they are often not practical [25].

Much less appears to be known in the high-dimensional regime. Exponential time algorithm, such as [21, 29], optimizes Tukey depth [42, 10]. Their results reveal that handling a constant fraction of outliers ($\epsilon = \text{const.}$) is actually minimax-optimal. Work in [11] first provided a polynomial time algorithm for this problem. They show that replacing the standard inner product in Matching Pursuit with a trimmed version, one can recover from an ϵ -fraction of outliers, with $\epsilon = O(1/\sqrt{k})$. Very recently, [36] considered more general sparsity constrained M -estimation by using a trimmed estimator in each step of gradient descent, yet the robustness guarantee $\epsilon = O(1/\sqrt{k})$ is still sub-optimal. Another approach follows as a byproduct of a recent algorithm for robust sparse mean estimation, in [1]. However, their error guarantee scales with $\|\beta^*\|_2$, and moreover, does not provide *exact recovery* in the adversarial corruption case without stochastic noise (i.e., noise variance $\sigma^2 = 0$). We note that this is an inevitable consequence of their approach, as they directly use sparse mean estimation on $\{y_i \mathbf{x}_i\}$, rather than considering Maximum Likelihood Estimation.

Robust mean estimation. The idea in [38, 15, 17] is to leverage recent breakthroughs in robust mean estimation. Very recently, [13, 33] provided the first robust mean estimation algorithms that can handle a constant fraction of outliers (though [33] incurs a small (logarithmic) dependence in the dimension). Following their work, [1] extended the ellipsoid algorithm from [13] to robust sparse mean estimation in high dimensions. They show that k -sparse mean estimation in \mathbb{R}^d with a constant fraction of outliers can be done with $n = \Omega(k^2 \log(d))$ samples. The k^2 term appears to be necessary, as $n = \Omega(k^2)$ follows from an oracle-based lower bound [16].

1.2 Main contributions

- Our result is a robust variant of Iterative Hard Thresholding (IHT) [8]. We provide a deterministic stability result showing that IHT works with any robust sparse mean estimation algorithm. We show our robust IHT does not accumulate the errors of a (any) robust sparse mean estimation subroutine for computing the gradient. Specifically, robust IHT produces a final solution whose error is orderwise the same as the error guaranteed by an single use of the robust mean estimation subroutine. We refer to [Definition 2.1](#) and [Theorem 2.1](#) for the precise statement. Thus our result can be viewed as a meta-theorem that can be coupled with any robust sparse mean estimator.
- Coupling robust IHT with a robust sparse mean estimation subroutine based on a version of the ellipsoid algorithm given and analyzed in [1], our results [Corollary 3.1](#) show that given ϵ -corrupted sparse regression samples with identity covariance, we recover β^* within additive error $O(\sigma\epsilon)$ (which is minimax optimal [21]). The proof of the ellipsoid algorithm’s performance in [1] hinges on obtaining an upper bound on the sparse operator norm (their Lemmas A.2 and A.3). As we show (see [Appendix B](#)), the statement of Lemma A.3 seems to be incorrect, and the general approach of upper bounding the sparse operator norm may not work. Nevertheless, the algorithm performance they claim is correct, as we show through a different avenue (see [Lemma D.3](#) in [Appendix D.3](#)).

Using this ellipsoid algorithm, In particular, we obtain exact recovery if either the fraction of outliers goes to zero (this is just ordinary sparse regression), or in the presence of a constant fraction of outliers but with the additive noise term going to zero (this is the case of robust sparse linear equations). To the best of our knowledge, this is the first result that shows exact recovery for robust sparse linear equations with a constant fraction of outliers. This is the content of [Section 3](#).

- For robust sparse regression with *unknown covariance matrix*, we consider the wide class of sparse covariance matrices [7]. We then prove a result that may be of interest in its own right: we provide a novel robust sparse mean estimation algorithm that is based on a filtering algorithm for sequentially screening and removing potential outliers. We show that the filtering algorithm is flexible enough to deal with unknown covariance, whereas the ellipsoid algorithm cannot. It also runs a factor of $O(d^2)$ faster than the ellipsoid algorithm. If the covariance matrix is sufficiently sparse, our filtering algorithm gives a robust sparse mean estimation algorithm, that can then be coupled with our meta-theorem. Together, these two guarantee recovery of β^* within an additive error of $O(\sigma\sqrt{\epsilon})$. In the case of unknown covariance, this is the best (and in fact, only) result we are aware of for robust sparse regression. We note that it can be applied to the case of known and identity covariance, though it is weaker than the optimal results we obtain using the computationally more expensive ellipsoid algorithm. Nevertheless, in both cases (unknown sparse, or known identity) the result is strong enough to guarantee exact recovery when either σ or ϵ goes to zero. We demonstrate the practical effectiveness of our filtering algorithm in [Appendix H](#). This is the content of [Section 4](#) and [Section 5](#).

1.3 Setup, Notation and Outline

In this subsection, we formally define the corruption model and the sparse regression model. We first introduce the ϵ -corrupted samples described above:

Definition 1.1 (ϵ -corrupted samples). *Let $\{z_i, i \in \mathcal{G}\}$ be i.i.d. observations follow from a distribution P . The ϵ -corrupted samples $\{z_i, i \in \mathcal{S}\}$ are generated by the following process: an adversary chooses an arbitrary ϵ -fraction of the samples in \mathcal{G} and modifies them with arbitrary values. After the corruption, we use \mathcal{S} to denote the observations, and use $\mathcal{B} = \mathcal{S} \setminus \mathcal{G}$ to denote the corruptions.*

The parameter ϵ represents the fraction of outliers. Throughout, we assume that it is a (small) constant, *independent of dimension or other problem parameters*. Furthermore, we assume that the distribution P is the standard Gaussian-design AWGN linear model.

Model 1.1. The observations $\{\mathbf{z}_i = (y_i, \mathbf{x}_i), i \in \mathcal{G}\}$ follow from the linear model $y_i = \mathbf{x}_i^\top \boldsymbol{\beta}^* + \xi_i$, where $\boldsymbol{\beta}^* \in \mathbb{R}^d$ is the model parameter, and assumed to be k -sparse. We assume that $\mathbf{x}_i \sim \mathcal{N}(0, \boldsymbol{\Sigma})$ and $\xi_i \sim \mathcal{N}(0, \sigma^2)$, where $\boldsymbol{\Sigma}$ is the normalized covariance matrix with $\boldsymbol{\Sigma}_{jj} \leq 1$ for all $j \in [d]$. We denote μ_α as the smallest eigenvalue of $\boldsymbol{\Sigma}$, and μ_β as its largest eigenvalue. They are assumed to be universal constants in this paper, and we denote the constant $c_\kappa = \mu_\beta / \mu_\alpha$.

As in [1], we pre-process by removing ‘‘obvious’’ outliers; we henceforth assume that all authentic and corrupted points are within a radius bounded by a polynomial in n , d and $1/\epsilon$.

Notation. We denote the hard thresholding operator of sparsity k' by $\mathsf{P}_{k'}$. We define the k -sparse operator norm as $\|M\|_{\tilde{k}, \text{op}} = \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} |\mathbf{v}^\top M \mathbf{v}|$, where M is not required to be positive-semidefinite (p.s.d.). We use trace inner product $\langle A, B \rangle$ to denote $\text{Tr}(A^\top B)$. We use $\mathbb{E}_{i \in_u \mathcal{S}}$ to denote the expectation operator obtained by the uniform distribution over all samples i in a set \mathcal{S} . Finally, we use the notation $\tilde{O}(\cdot)$ to hide the dependency on $\text{poly} \log(1/\epsilon)$, and $\tilde{\Omega}(\cdot)$ to hide the dependency on $\text{poly} \log(k)$ in our bounds.

2 Hard thresholding with robust gradient estimation

In this section, we present our method of using robust sparse gradient updates in IHT. We then show statistical recovery guarantees given any accurate robust sparse gradient estimation, which is formally defined in [Definition 2.1](#).

We define the notation for the stochastic gradient \mathbf{g}_i corresponding to the i^{th} point \mathbf{z}_i , and the population gradient for $\mathbf{z}_i \sim P$ based on [Model 1.1](#), $\mathbf{g}_i^t = \mathbf{x}_i (\mathbf{x}_i^\top \boldsymbol{\beta}^t - y_i)$, and $\mathbf{G}^t = \mathbb{E}_{\mathbf{z}_i \sim P} (\mathbf{g}_i^t)$, where P is the distribution of the authentic points. Since $\mathbb{E}_{\mathbf{z}_i \sim P} (\mathbf{x}_i \mathbf{x}_i^\top) = \boldsymbol{\Sigma}$, the population mean of all authentic gradients is given by $\mathbf{G}^t = \mathbb{E}_{\mathbf{z}_i \sim P} (\mathbf{x}_i \mathbf{x}_i^\top (\boldsymbol{\beta}^t - \boldsymbol{\beta}^*)) = \boldsymbol{\Sigma}(\boldsymbol{\beta}^t - \boldsymbol{\beta}^*)$.

In the *uncorrupted case* where all samples $\{\mathbf{z}_i, i \in \mathcal{G}\}$ follow from [Model 1.1](#), a single iteration of IHT updates $\boldsymbol{\beta}^t$ via $\boldsymbol{\beta}^{t+1} = \mathsf{P}_{k'}(\boldsymbol{\beta}^t - \mathbb{E}_{i \in_u \mathcal{G}} \mathbf{g}_i^t)$. Here, the hard thresholding operator $\mathsf{P}_{k'}$ selects the k' largest elements in magnitude, and the parameter k' is proportional to k (specified in [Theorem 2.1](#)). However, given ϵ -corrupted samples $\{\mathbf{z}_i, i \in \mathcal{S}\}$ according to [Definition 1.1](#), the IHT update based on empirical average of all gradient samples $\{\mathbf{g}_i, i \in \mathcal{S}\}$ can be arbitrarily bad.

The key goal in this paper is to find a robust estimate $\hat{\mathbf{G}}^t$ to replace \mathbf{G}^t in each step of IHT, with sample complexity *sub-linear* in the dimension d . For instance, we consider robust sparse regression with $\boldsymbol{\Sigma} = \mathbf{I}_d$. Then, $\mathbf{G}^t = \boldsymbol{\beta}^t - \boldsymbol{\beta}^*$ is guaranteed to be $(k' + k)$ -sparse in each iteration of IHT. In this case, given ϵ -corrupted samples, we can use a robust sparse mean estimator to recover the unknown true \mathbf{G}^t from $\{\mathbf{g}_i^t\}_{i=1}^{|\mathcal{S}|}$, with sub-linear sample complexity.

More generally, we propose Robust Sparse Gradient Estimator (RSGE) for gradient estimation given ϵ -corrupted samples, as defined in [Definition 2.1](#), which guarantees that the deviation between the robust estimate $\hat{\mathbf{G}}(\boldsymbol{\beta})$ and true $\mathbf{G}(\boldsymbol{\beta})$, with sample complexity $n \ll d$. For a fixed k -sparse parameter $\boldsymbol{\beta}$, we drop the superscript t without abuse of notation, and use \mathbf{g}_i in place of \mathbf{g}_i^t , and \mathbf{G} in place of \mathbf{G}^t ; $\mathbf{G}(\boldsymbol{\beta})$ denotes the population gradient over the authentic samples’ distribution P , at the point $\boldsymbol{\beta}$.

Definition 2.1 ($\psi(\epsilon)$ -RSGE). Given $n(k, d, \epsilon, \nu)$ ϵ -corrupted samples $\{\mathbf{z}_i\}_{i=1}^n$ from [Model 1.1](#), we call $\hat{\mathbf{G}}(\boldsymbol{\beta})$ a $\psi(\epsilon)$ -RSGE, if given $\{\mathbf{z}_i\}_{i=1}^n$, $\hat{\mathbf{G}}(\boldsymbol{\beta})$ guarantees $\|\hat{\mathbf{G}}(\boldsymbol{\beta}) - \mathbf{G}(\boldsymbol{\beta})\|_2^2 \leq \alpha(\epsilon) \|\mathbf{G}(\boldsymbol{\beta})\|_2^2 + \psi(\epsilon)$, with probability at least $1 - \nu$.

Here, we use $n(k, d, \epsilon, \nu)$ to denote the sample complexity as a function of (k, d, ϵ, ν) , and note that the definition of RSGE does not require $\boldsymbol{\Sigma}$ to be identity matrix. The parameters $\alpha(\epsilon)$ and $\psi(\epsilon)$ will be specified by concrete robust sparse mean estimators in subsequent sections. Equipped with [Definition 2.1](#), we propose [Algorithm 1](#), which takes any RSGE as a subroutine in line 7, and runs a robust variant of IHT with the estimated sparse gradient $\hat{\mathbf{G}}^t$ at each iteration in line 8.¹

¹Our results require sample splitting to maintain independence between subsequent iterations, though we believe this is an artifact of our analysis. Similar approach has been used in [2, 38] for theoretical analysis. We do not use sample splitting technique in the experiments.

Algorithm 1 Robust sparse regression with RSGE

- 1: **Input:** Data samples $\{y_i, \mathbf{x}_i\}_{i=1}^N$, RSGE subroutine.
 - 2: **Output:** The estimation $\hat{\beta}$.
 - 3: **Parameters:** Hard thresholding parameter k' .
-
- 4: Split samples into T subsets of size n . Initialize with $\beta^0 = \mathbf{0}$.
 - 5: **for** $t = 0$ to $T - 1$, **do**
 - 6: At current β^t , calculate all gradients for current n samples: $\mathbf{g}_i^t = \mathbf{x}_i (\mathbf{x}_i^\top \beta^t - y_i)$, $i \in [n]$.
 - 7: The initial input set is $\{\mathbf{g}_i^t\}_{i=1}^n$. We use a RSGE to get $\hat{\mathbf{G}}^t$.
 - 8: Update the parameter: $\beta^{t+1} = \text{P}_{k'}(\beta^t - \eta \hat{\mathbf{G}}^t)$.
 - 9: **end for**
 - 10: Output the estimation $\hat{\beta} = \beta^T$.
-

2.1 Global linear convergence and parameter recovery guarantees

In each single IHT update step, RSGE introduces a controlled amount of error. [Theorem 2.1](#) gives a global linear convergence guarantee for [Algorithm 1](#) by showing that IHT does not accumulate too much error. In particular, we are able to recover β^* within error $O(\sqrt{\psi(\epsilon)})$ given any $\psi(\epsilon)$ -RSGE subroutine. We give the proof of [Theorem 2.1](#) in [Appendix A](#). The hyper-parameter $k' = c_\kappa^2 k$ guarantees global linear convergence of IHT when $c_\kappa > 1$ (when $\Sigma \neq \mathbf{I}_d$). This setup has been used in [\[28, 40\]](#), and is proved to be necessary in [\[35\]](#). Note that [Theorem 2.1](#) is a deterministic stability result in nature, and we obtain probabilistic results by certifying the RSGE condition.

Theorem 2.1 (Meta-theorem). *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from [Model 1.1](#). [Algorithm 1](#), with $\psi(\epsilon)$ -RSGE defined in [Definition 2.1](#), with step size $\eta = 1/\mu_\beta$ outputs $\hat{\beta}$, such that $\|\hat{\beta} - \beta^*\|_2 = O(\sqrt{\psi(\epsilon)})$, with probability at least $1 - \nu$, by setting $k' = c_\kappa^2 k$ and $T = \Theta(\log(\|\beta^*\|_2 / \sqrt{\psi(\epsilon)}))$. The sample complexity is $N(k, d, \epsilon, \nu) = n(k, d, \epsilon, \nu/T)T$.*

3 Robust sparse regression with near-optimal guarantee

In this section, we provide near optimal statistical guarantee for robust sparse regression when the covariance matrix is identity. Under the assumption $\Sigma = \mathbf{I}_d$, [\[1\]](#) proposes a robust sparse regression estimator based on robust sparse mean estimation on $\{y_i \mathbf{x}_i, i \in \mathcal{S}\}$, leveraging the fact that $\mathbb{E}_{\mathbf{z}_i \sim P}(y_i \mathbf{x}_i) = \beta^*$. With sample complexity $N = \Omega(\frac{k^2 \log(d/\nu)}{\epsilon^2})$, this algorithm produces a $\tilde{\beta}$ such that $\|\tilde{\beta} - \beta^*\|_2^2 = \tilde{O}(\epsilon^2(\|\beta^*\|_2^2 + \sigma^2))$, with probability at least $1 - \nu$. Using [Theorem 2.1](#), we show that we can obtain significantly stronger statistical guarantees which are statistically optimal; in particular, our guarantees are independent of $\|\beta^*\|_2$ and yield exact recovery when $\sigma = 0$.

3.1 RSGE via the ellipsoid algorithm

More specifically, the ellipsoid-based robust sparse mean estimation algorithm [\[1\]](#) deals with outliers by trying to optimize the set of weights $\{w_i, i \in \mathcal{S}\}$ on each of the samples in \mathbb{R}^d – ideally outliers would receive lower weight and hence their impact would be minimized. Since the set of weights is convex, this can be approached using a separation oracle [Algorithm 2](#). The [Algorithm 2](#) depends on a convex relaxation of Sparse PCA, and the hard thresholding parameter is $k = k' + k$, as the population mean of all authentic gradient samples \mathbf{G}^t is guaranteed to be $(k' + k)$ -sparse. In line 4 and 5, we calculate the weighted mean and covariance based on a hard thresholding operator. In line 6 of [Algorithm 2](#), with each call to the relaxation of Sparse PCA, we obtain an optimal value, λ^* , and optimal solution, \mathbf{H}^* , to the problem:

$$\lambda^* = \max_{\mathbf{H}} \text{Tr} \left(\left(\hat{\Sigma} - F(\hat{\mathbf{G}}) \right) \cdot \mathbf{H} \right), \quad \text{subject to } \mathbf{H} \succcurlyeq 0, \|\mathbf{H}\|_{1,1} \leq \tilde{k}, \text{Tr}(\mathbf{H}) = 1. \quad (1)$$

Algorithm 2 Separation oracle for robust sparse estimation [1]

- 1: **Input:** Weights from the previous iteration $\{w_i, i \in \mathcal{S}\}$, gradient samples $\{\mathbf{g}_i, i \in \mathcal{S}\}$.
 - 2: **Output:** Weight $\{w'_i, i \in \mathcal{S}\}$
 - 3: **Parameters:** Hard thresholding parameter \tilde{k} , parameter ρ_{sep} .
-
- 4: Compute the weighted sample mean $\tilde{\mathbf{G}} = \sum_{i \in \mathcal{S}} w_i \mathbf{g}_i$, and $\hat{\mathbf{G}} = \text{P}_{2\tilde{k}}(\tilde{\mathbf{G}})$.
 - 5: Compute the weighted sample covariance matrix $\hat{\Sigma} = \sum_{i \in \mathcal{S}} w_i (\mathbf{g}_i - \hat{\mathbf{G}}) (\mathbf{g}_i - \hat{\mathbf{G}})^\top$.
 - 6: Solve: $\max_{\mathbf{H}} \text{Tr} \left((\hat{\Sigma} - F(\hat{\mathbf{G}})) \cdot \mathbf{H} \right)$, subject to $\mathbf{H} \succcurlyeq 0, \|\mathbf{H}\|_{1,1} \leq \tilde{k}, \text{Tr}(\mathbf{H}) = 1$.
Let λ^* be the optimal value, and \mathbf{H}^* be the corresponding solution.
 - 7: **if** $\lambda^* \leq \rho_{\text{sep}}$, **then return** “Yes”.
 - 8: **return** The hyperplane: $\ell(w') = \left\langle \left(\sum_{i \in \mathcal{S}} w'_i (\mathbf{g}_i - \hat{\mathbf{G}}) (\mathbf{g}_i - \hat{\mathbf{G}})^\top - F(\hat{\mathbf{G}}) \right), \mathbf{H}^* \right\rangle - \lambda^*$.
-

Here, $\hat{\mathbf{G}}, \hat{\Sigma}$ are weighted first and second order moment estimates from ϵ -corrupted samples, and $F: \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d}$ is a function with closed-form

$$F(\hat{\mathbf{G}}) = \|\hat{\mathbf{G}}\|_2^2 \mathbf{I}_d + \hat{\mathbf{G}} \hat{\mathbf{G}}^\top + \sigma^2 \mathbf{I}_d. \quad (2)$$

For eq. (2), given the population mean \mathbf{G} , we have $F(\mathbf{G}) = \mathbb{E}_{\mathbf{z}_i \sim P}((\mathbf{g}_i - \mathbf{G})(\mathbf{g}_i - \mathbf{G})^\top)$, which calculates the underlying true covariance matrix. We provide more details about the calculation of $F(\cdot)$, as well as some smoothness properties, in Appendix C.

The key component in the separation oracle Algorithm 2 is to use convex relaxation of Sparse PCA eq. (1). This idea generalizes existing work on using PCA to detect outliers in low dimensional robust mean estimation [13, 33]. To gain some intuition for eq. (1), if \mathbf{g}_i is an outlier, then the optimal solution of eq. (1), \mathbf{H}^* , may detect the direction of this outlier. And this outlier will be down-weighted in the output of Algorithm 2 by the separating hyperplane. Finally, Algorithm 2 will terminate with $\lambda^* \leq \rho_{\text{sep}}$ (line 7) and output the robust sparse mean estimation of the gradients $\hat{\mathbf{G}}$.

Indeed, the ellipsoid-algorithm-based robust sparse mean estimator gives a RSGE, which we can combine with Theorem 2.1 to obtain stronger results. We state these as Corollary 3.1. We note again that the analysis in [1] has a flaw. Their Lemma A.3 is incorrect, as our counterexample in Appendix B demonstrates. We provide a correct route of analysis in Lemma D.3 of Appendix D.

3.2 Near-optimal statistical guarantees

Corollary 3.1. *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from Model 1.1 with $\Sigma = \mathbf{I}_d$. By setting $\tilde{k} = k' + k$, if we use the ellipsoid algorithm for robust sparse gradient estimation with $\rho_{\text{sep}} = \Theta(\epsilon(\|\mathbf{G}^t\|_2^2 + \sigma^2))$, it requires $N(k, d, \epsilon, \nu) = \Omega\left(\frac{k^2 \log(dT/\nu)}{\epsilon^2}\right)T$ samples, and guarantees $\psi(\epsilon) = \tilde{O}(\epsilon^2 \sigma^2)$. Hence, Algorithm 1 outputs $\hat{\beta}$, such that $\|\hat{\beta} - \beta^*\|_2 = \tilde{O}(\sigma \epsilon)$, with probability at least $1 - \nu$, by setting $T = \Theta\left(\log\left(\frac{\|\beta^*\|_2}{\epsilon \sigma}\right)\right)$.*

For a desired error level $\epsilon' \geq \epsilon$, we only require sample complexity $N(k, d, \epsilon, \nu) = \Omega\left(\frac{k^2 \log(dT/\nu)}{\epsilon'^2}\right)T$. Hence, we can achieve statistical error $\tilde{O}(\sigma(\sqrt{k^2 \log(d)/N} \vee \epsilon))$. Our error bound is nearly optimal compared to the information-theoretically optimal $O(\sigma(\sqrt{k \log(d)/N} \vee \epsilon))$ in [21], as the k^2 term is necessary by an oracle-based SQ lower bound [16].

Proof sketch of Corollary 3.1 The key to the proof relies on showing that λ^* controls the quality of the weights of the current iteration, i.e., small λ^* means good weights and thus a good current solution. Showing this relies on using λ^* to control $\hat{\Sigma} - F(\hat{\mathbf{G}})$. Lemma A.3 in [1] claims that $\lambda^* \geq \|\hat{\Sigma} - F(\hat{\mathbf{G}})\|_{\tilde{k}, \text{op}}$. As we show in Appendix B, however, this need not hold. This is because the trace norm maximization eq. (1) is *not* a valid convex relaxation for the \tilde{k} -sparse operator norm when the term $\hat{\Sigma} - F(\hat{\mathbf{G}})$ is not

Algorithm 3 RSGE via filtering

- 1: **Input:** A set \mathcal{S}_{in} .
 - 2: **Output:** A set \mathcal{S}_{out} or sparse mean vector $\widehat{\mathbf{G}}$.
 - 3: **Parameters:** Hard thresholding parameter \tilde{k} , parameter ρ_{sep} .
-
- 4: Compute the sample mean $\widetilde{\mathbf{G}} = \mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}}(\mathbf{g}_i)$, and $\widehat{\mathbf{G}} = \text{P}_{2\tilde{k}}(\widetilde{\mathbf{G}})$.
 - 5: Compute the sample covariance matrix $\widehat{\Sigma} = \mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}}(\mathbf{g}_i - \widehat{\mathbf{G}})(\mathbf{g}_i - \widehat{\mathbf{G}})^\top$.
 - 6: Solve the following convex program:

$$\max_{\mathbf{H}} \text{Tr}(\widehat{\Sigma} \cdot \mathbf{H}), \quad \text{subject to } \mathbf{H} \succcurlyeq 0, \|\mathbf{H}\|_{1,1} \leq \tilde{k}, \text{Tr}(\mathbf{H}) = 1. \quad (3)$$

Let λ^* be the optimal value, and \mathbf{H}^* be the corresponding solution.

- 7: **if** $\lambda^* \leq \rho_{\text{sep}}$, **then return** with $\widehat{\mathbf{G}}$.
- 8: Calculate projection score for each $i \in \mathcal{S}_{\text{in}}$:

$$\tau_i = \text{Tr}(\mathbf{H}^* \cdot (\mathbf{g}_i - \widehat{\mathbf{G}})(\mathbf{g}_i - \widehat{\mathbf{G}})^\top).$$

- 9: Randomly remove a sample r from \mathcal{S}_{in} according to

$$\Pr(\mathbf{g}_i \text{ is removed}) = \frac{\tau_i}{\sum_{i \in \mathcal{S}_{\text{in}}} \tau_i}. \quad (4)$$

- 10: **return** the set $\mathcal{S}_{\text{out}} = \mathcal{S}_{\text{in}} \setminus \{r\}$.
-

p.s.d. (which indeed it need not be). We provide a different line of analysis in [Lemma D.3](#), essentially showing that even without the claimed (incorrect) bound, λ^* can still provide the control we need. With the corrected analysis for λ^* , the ellipsoid algorithm guarantees $\|\widehat{\mathbf{G}} - \mathbf{G}\|_2^2 = \tilde{O}(\epsilon^2(\|\boldsymbol{\beta} - \boldsymbol{\beta}^*\|_2^2 + \sigma^2))$ with probability at least $1 - \nu$. Therefore, the algorithm provides an $\tilde{O}(\epsilon^2\sigma^2)$ -RSGE.

4 Robust sparse mean estimation via filtering

From a computational viewpoint, the time complexity of [Algorithm 1](#) depends on the RSGE in each iterate. The time complexity of the ellipsoid algorithm is indeed polynomial in the dimension, but it requires $O(d^2)$ calls to a relaxation of Sparse PCA ([\[9\]](#)). In this section, we introduce a faster algorithm as a RSGE, which only requires $O(n)$ calls of Sparse PCA (recall that n only scales with $k^2 \log d$). Importantly, this RSGE is flexible enough to deal with unknown covariance matrix, yet the ellipsoid algorithm cannot. Before we move to the result for unknown covariance matrix in [Section 5](#), we first introduce [Algorithm 3](#) and analyze its performance when the covariance is identity. These supporting Lemmas will be later used in the unknown case.

Our proposed RSGE ([Algorithm 3](#)) attempts to remove one outlier at each iteration, as long as a good solution has not already been identified. It first estimates the gradient $\widehat{\mathbf{G}}$ by hard thresholding (line 4) and then estimates the corresponding sample covariance matrix $\widehat{\Sigma}$ (line 5). By solving (a relaxation of) Sparse PCA, we obtain a scalar λ^* as well as a matrix \mathbf{H}^* . If λ^* is smaller than the predetermined threshold ρ_{sep} , we have a certificate that the effect of the outliers is well-controlled (specified in [eq. \(5\)](#)). Otherwise, we compute a score for each sample based on \mathbf{H}^* , and discard one of the samples according to a probability distribution where each sample's probability of being discarded is proportional to the score we have computed ². [Algorithm 3](#) can be used for other robust sparse

²Although we remove one sample in [Algorithm 3](#), our theoretical analysis naturally extend to removing constant number of outliers. This speeds up the algorithm in practice, yet shares the same computational complexity

functional estimation problems (e.g., robust sparse mean estimation for $\mathcal{N}(\boldsymbol{\mu}, \mathbf{I}_d)$, where $\boldsymbol{\mu} \in \mathbb{R}^d$ is k -sparse). To use [Algorithm 3](#) as a RSGE given n gradient samples (denoted as \mathcal{S}_{in}), we call [Algorithm 3](#) repeatedly on \mathcal{S}_{in} and then on its output, \mathcal{S}_{out} , until it returns a robust estimator $\widehat{\mathbf{G}}$. The next theorem provides guarantees on this iterative application of [Algorithm 3](#).

Theorem 4.1. *Suppose we observe $n = \Omega\left(\frac{k^2 \log(d/\nu)}{\epsilon}\right)$ ϵ -corrupted samples from [Model 1.1](#) with $\boldsymbol{\Sigma} = \mathbf{I}_d$. Let \mathcal{S}_{in} be an ϵ -corrupted set of gradient samples $\{\mathbf{g}_i^t\}_{i=1}^n$. By setting $\tilde{k} = k' + k$, if we run [Algorithm 3](#) iteratively with initial set \mathcal{S}_{in} , and subsequently on \mathcal{S}_{out} , and use $\rho_{\text{sep}} = C_\gamma(\|\mathbf{G}^t\|_2^2 + \sigma^2)$,³ then this repeated use of [Algorithm 3](#) will stop after at most $\frac{1.1\gamma}{\gamma-1}\epsilon n$ iterations, and output $\widehat{\mathbf{G}}^t$, such that $\|\widehat{\mathbf{G}}^t - \mathbf{G}^t\|_2^2 = \tilde{O}(\epsilon(\|\mathbf{G}^t\|_2^2 + \sigma^2))$, with probability at least $1 - \nu - \exp(-\Theta(\epsilon n))$. Here, C_γ is a constant depending on γ , where $\gamma \geq 4$ is a constant.*

Thus, [Theorem 4.1](#) shows that with high probability, [Algorithm 3](#) provides a Robust Sparse Gradient Estimator where $\psi(\epsilon) = \tilde{O}(\epsilon\sigma^2)$. For example, we can take $\nu = d^{-\Theta(1)}$. Combining now with [Theorem 2.1](#), we obtain an error guarantee for robust sparse regression.

Corollary 4.1. *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from [Model 1.1](#) with $\boldsymbol{\Sigma} = \mathbf{I}_d$. Under the same setting as [Theorem 4.1](#), if we use [Algorithm 3](#) for robust sparse gradient estimation, it requires $N(k, d, \epsilon, \nu) = \Omega\left(\frac{k^2 \log(dT/\nu)}{\epsilon}\right) T$ samples, and $T = \Theta\left(\log\left(\frac{\|\boldsymbol{\beta}^*\|_2}{\sigma\sqrt{\epsilon}}\right)\right)$, then we have $\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2 = \tilde{O}(\sigma\sqrt{\epsilon})$ with probability at least $1 - \nu - T \exp(-\Theta(\epsilon n))$.*

Similar to [Section 3](#), we can achieve statistical error $\tilde{O}(\sigma(\sqrt{k^2 \log(d)/N} \vee \sqrt{\epsilon}))$. The scaling of ϵ in [Corollary 4.1](#) is $\tilde{O}(\sqrt{\epsilon})$. These guarantees are worse than $\tilde{O}(\epsilon)$ achieved by ellipsoid methods. Nevertheless, this result is strong enough to guarantee exact recovery when either σ or ϵ goes to zero. The simulation of robust estimation for the filtering algorithm is in [Appendix H](#).

The key step in [Algorithm 3](#) is outlier removal [eq. \(4\)](#) based on the solution of Sparse PCA's convex relaxation [eq. \(3\)](#). We describe the outlier removal below, and then give the proofs in [Appendix E](#) and [Appendix F](#).

Outlier removal guarantees in [Algorithm 3](#). We denote samples in the input set \mathcal{S}_{in} as \mathbf{g}_i . This input set \mathcal{S}_{in} can be partitioned into two parts: $\mathcal{S}_{\text{good}} = \{i : i \in \mathcal{G} \text{ and } i \in \mathcal{S}_{\text{in}}\}$, and $\mathcal{S}_{\text{bad}} = \{i : i \in \mathcal{B} \text{ and } i \in \mathcal{S}_{\text{in}}\}$. [Lemma 4.1](#) shows that [Algorithm 3](#) can return a guaranteed gradient estimate, or the outlier removal step [eq. \(4\)](#) is likely to discard an outlier. The guarantee on the outlier removal step [eq. \(4\)](#) hinges on the fact that if $\sum_{i \in \mathcal{S}_{\text{good}}} \tau_i$ is less than $\sum_{i \in \mathcal{S}_{\text{bad}}} \tau_i$, we can show [eq. \(4\)](#) is likely to remove an outlier.

Lemma 4.1. *Suppose we observe $n = \Omega\left(\frac{k^2 \log(d/\nu)}{\epsilon}\right)$ ϵ -corrupted samples from [Model 1.1](#) with $\boldsymbol{\Sigma} = \mathbf{I}_d$. Let \mathcal{S}_{in} be an ϵ -corrupted set $\{\mathbf{g}_i^t\}_{i=1}^n$. [Algorithm 3](#) computes λ^* that satisfies*

$$\lambda^* \geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in \mathcal{S}_{\text{in}}} (\mathbf{g}_i - \widehat{\mathbf{G}}) (\mathbf{g}_i - \widehat{\mathbf{G}})^\top \right) \mathbf{v}. \quad (5)$$

If $\lambda^* \geq \rho_{\text{sep}} = C_\gamma(\|\mathbf{G}^t\|_2^2 + \sigma^2)$, then with probability at least $1 - \nu$, we have $\sum_{i \in \mathcal{S}_{\text{good}}} \tau_i \leq \frac{1}{\gamma} \sum_{i \in \mathcal{S}_{\text{in}}} \tau_i$, where τ_i is defined in [line 8](#), C_γ is a constant depending on γ , and $\gamma \geq 4$ is a constant.

The proofs are collected in [Appendix E](#). In a nutshell, [eq. \(5\)](#) is a natural convex relaxation for the sparsity constraint $\{\mathbf{v} : \|\mathbf{v}\|_2 = 1, \|\mathbf{v}\|_0 \leq \tilde{k}\}$. On the other hand, when $\lambda^* \geq \rho_{\text{sep}}$, the contribution of $\sum_{i \in \mathcal{S}_{\text{good}}} \tau_i$ is relatively small, which can be obtained through concentration inequalities for the samples in $\mathcal{S}_{\text{good}}$. Based on [Lemma 4.1](#), if $\lambda^* \leq \rho_{\text{sep}}$, then the RHS of [eq. \(5\)](#) is bounded, leading to the error guarantee of $\|\widehat{\mathbf{G}}^t - \mathbf{G}^t\|_2^2$. On the other hand, if $\lambda^* \geq \rho_{\text{sep}}$, we can show that [eq. \(4\)](#) is more likely to throw out samples of \mathcal{S}_{bad} rather than $\mathcal{S}_{\text{good}}$. Iteratively applying [Algorithm 3](#) on the remaining samples, we can remove those outliers with large effect, and keep the remaining outliers' effect well-controlled. This leads to the final bounds in [Theorem 4.1](#).

³Similar to [[13](#), [14](#), [15](#), [1](#)], our results seem to require this side information.

5 Robust sparse regression with unknown covariance

In this section, we consider robust sparse regression with unknown covariance matrix Σ , which has additional sparsity structure. Formally, we define the sparse covariance matrices as follows:

Model 5.1 (Sparse covariance matrices). *In Model 1.1, the authentic covariates $\{\mathbf{x}_i, i \in \mathcal{G}\}$ are drawn from $\mathcal{N}(0, \Sigma)$. We assume that each row and column of Σ is r -sparse, but the positions of the non-zero entries are unknown.*

Model 5.1 is widely studied in high dimensional statistics [7, 20, 45]. Under Model 5.1, for the population gradient $\mathbf{G}^t = \mathbb{E}_P(\mathbf{x}_i \mathbf{x}_i^\top (\beta^t - \beta^*)) = \Sigma \omega^t$, where we use ω^t to denote the $(k' + k)$ -sparse vector $\beta^t - \beta^*$, we can guarantee the $\|\mathbf{G}^t\|_0 = \|\Sigma \omega^t\|_0 \leq r(k' + k)$. Hence, we can use the filtering algorithm (Algorithm 3) with $\tilde{k} = r(k' + k)$ as a RSGE for robust sparse regression with unknown Σ . When the covariance is unknown, we cannot evaluate $F(\cdot)$ a priori, thus the ellipsoid algorithm is not applicable to this case. And we provide error guarantees as follows.

Theorem 5.1. *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from Model 1.1, where the covariates \mathbf{x}_i 's follow from Model 5.1. If we use Algorithm 3 for robust sparse gradient estimation, it requires $\tilde{\Omega}\left(\frac{r^2 k^2 \log(dT/\nu)}{\epsilon}\right) T$ samples, and $T = \Theta\left(\log\left(\frac{\|\beta^*\|_2}{\sigma\sqrt{\epsilon}}\right)\right)$, then, we have $\|\hat{\beta} - \beta^*\|_2 = \tilde{O}(\sigma\sqrt{\epsilon})$, with probability at least $1 - \nu - T \exp(-\Theta(\epsilon n))$.*

The proof of Theorem 5.1 is collected in Appendix G, and the main technique hinges on previous analysis for the identity covariance case (Theorem 4.1 and Lemma 4.1). In the case of unknown covariance, this is the best (and in fact, only) recovery guarantee we are aware of for robust sparse regression. We show the performance of robust estimation using our filtering algorithm with unknown covariance in Appendix H, and we observe same linear convergence as Section 4.

6 Acknowledgments

The authors would like to thank Simon S. Du for helpful discussions.

References

- [1] Sivaraman Balakrishnan, Simon S. Du, Jerry Li, and Aarti Singh. Computationally efficient robust sparse estimation in high dimensions. In *Proceedings of the 2017 Conference on Learning Theory*, 2017.
- [2] Sivaraman Balakrishnan, Martin J Wainwright, and Bin Yu. Statistical guarantees for the em algorithm: From population to sample-based analysis. *The Annals of Statistics*, 45(1):77–120, 2017.
- [3] Afonso S. Bandeira, Edgar Dobriban, Dustin G. Mixon, and William F. Sawin. Certifying the restricted isometry property is hard. *IEEE Transactions on Information Theory*, 59(6):3448–3450, 2013.
- [4] Boaz Barak and David Steurer. Proofs, beliefs, and algorithms through the lens of sum-of-squares. *Course notes: <http://www.sumofsquares.org/public/index.html>*, 2016.
- [5] Kush Bhatia, Prateek Jain, and Purushottam Kar. Robust regression via hard thresholding. In *Advances in Neural Information Processing Systems*, pages 721–729, 2015.
- [6] Kush Bhatia, Prateek Jain, and Purushottam Kar. Consistent robust regression. In *Advances in Neural Information Processing Systems*, pages 2107–2116, 2017.
- [7] Peter J Bickel and Elizaveta Levina. Covariance regularization by thresholding. *The Annals of Statistics*, 36(6):2577–2604, 2008.
- [8] Thomas Blumensath and Mike E Davies. Iterative hard thresholding for compressed sensing. *Applied and computational harmonic analysis*, 27(3):265–274, 2009.
- [9] Sébastien Bubeck. Convex optimization: Algorithms and complexity. *Foundations and Trends® in Machine Learning*, 8(3-4):231–357, 2015.
- [10] Mengjie Chen, Chao Gao, and Zhao Ren. Robust covariance and scatter matrix estimation under hubers contamination model. *Ann. Statist.*, 46(5):1932–1960, 10 2018.
- [11] Yudong Chen, Constantine Caramanis, and Shie Mannor. Robust sparse regression under adversarial corruption. In *International Conference on Machine Learning*, pages 774–782, 2013.
- [12] Yudong Chen, Lili Su, and Jiaming Xu. Distributed statistical machine learning in adversarial settings: Byzantine gradient descent. *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, 1(2):44, 2017.
- [13] Ilias Diakonikolas, Gautam Kamath, Daniel M Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Robust estimators in high dimensions without the computational intractability. In *Foundations of Computer Science (FOCS), 2016 IEEE 57th Annual Symposium on*, pages 655–664. IEEE, 2016.
- [14] Ilias Diakonikolas, Gautam Kamath, Daniel M. Kane, Jerry Li, Ankur Moitra, and Alistair Stewart. Being robust (in high dimensions) can be practical. In *International Conference on Machine Learning*, pages 999–1008, 2017.
- [15] Ilias Diakonikolas, Gautam Kamath, Daniel M Kane, Jerry Li, Jacob Steinhardt, and Alistair Stewart. Sever: A robust meta-algorithm for stochastic optimization. *arXiv preprint [arXiv:1803.02815](https://arxiv.org/abs/1803.02815)*, 2018.
- [16] Ilias Diakonikolas, Daniel M Kane, and Alistair Stewart. Statistical query lower bounds for robust estimation of high-dimensional gaussians and gaussian mixtures. In *Foundations of Computer Science (FOCS), 2017 IEEE 58th Annual Symposium on*, pages 73–84. IEEE, 2017.
- [17] Ilias Diakonikolas, Weihao Kong, and Alistair Stewart. Efficient algorithms and lower bounds for robust linear regression. In *Proceedings of the Thirtieth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2745–2754. SIAM, 2019.
- [18] David L Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- [19] Alexandre dAspremont, Francis Bach, and Laurent El Ghaoui. Optimal solutions for sparse principal component analysis. *Journal of Machine Learning Research*, 9(Jul):1269–1294, 2008.

- [20] Nouredine El Karoui. Operator norm consistent estimation of large-dimensional sparse covariance matrices. *The Annals of Statistics*, 36(6):2717–2756, 2008.
- [21] Chao Gao. Robust regression via multivariate regression depth. *arXiv preprint arXiv:1702.04656*, 2017.
- [22] Michael Grant, Stephen Boyd, and Yinyu Ye. Cvx: Matlab software for disciplined convex programming, 2008.
- [23] Frank R Hampel, Elvezio M Ronchetti, Peter J Rousseeuw, and Werner A Stahel. *Robust statistics: the approach based on influence functions*, volume 196. John Wiley & Sons, 2011.
- [24] Matthew J Holland and Kazushi Ikeda. Efficient learning with robust gradient descent. *arXiv preprint arXiv:1706.00182*, 2017.
- [25] Samuel B Hopkins, Tselil Schramm, Jonathan Shi, and David Steurer. Fast spectral algorithms from sum-of-squares proofs: tensor decomposition and planted sparse vectors. In *Proceedings of the forty-eighth annual ACM symposium on Theory of Computing*, pages 178–191. ACM, 2016.
- [26] Peter J Huber. Robust estimation of a location parameter. *The annals of mathematical statistics*, pages 73–101, 1964.
- [27] Peter J Huber. Robust statistics. In *International Encyclopedia of Statistical Science*, pages 1248–1251. Springer, 2011.
- [28] Prateek Jain, Ambuj Tewari, and Purushottam Kar. On iterative hard thresholding methods for high-dimensional m-estimation. In *Advances in Neural Information Processing Systems*, pages 685–693, 2014.
- [29] David S Johnson and Franco P Preparata. The densest hemisphere problem. *Theoretical Computer Science*, 6(1):93–107, 1978.
- [30] Sushrut Karmalkar and Eric Price. Compressed sensing with adversarial sparse noise via l1 regression. *arXiv preprint arXiv:1809.08055*, 2018.
- [31] Adam Klivans, Pravesh K. Kothari, and Raghu Meka. Efficient Algorithms for Outlier-Robust Regression. *arXiv preprint arXiv:1803.03241*, 2018.
- [32] Adam R Klivans, Philip M Long, and Rocco A Servedio. Learning halfspaces with malicious noise. *Journal of Machine Learning Research*, 10(Dec):2715–2740, 2009.
- [33] Kevin A Lai, Anup B Rao, and Santosh Vempala. Agnostic estimation of mean and covariance. In *Foundations of Computer Science (FOCS), 2016 IEEE 57th Annual Symposium on*, pages 665–674. IEEE, 2016.
- [34] Xiaodong Li. Compressed sensing and matrix completion with constant proportion of corruptions. *Constructive Approximation*, 37(1):73–99, 2013.
- [35] Haoyang Liu and Rina Foygel Barber. Between hard and soft thresholding: optimal iterative thresholding algorithms. *arXiv preprint arXiv:1804.08841*, 2018.
- [36] Liu Liu, Tianyang Li, and Constantine Caramanis. High dimensional robust estimation of sparse models via trimmed hard thresholding. *arXiv preprint arXiv:1901.08237*, 2019.
- [37] Nam H Nguyen and Trac D Tran. Exact recoverability from dense corrupted observations via l1-minimization. *IEEE transactions on information theory*, 59(4):2017–2035, 2013.
- [38] Adarsh Prasad, Arun Sai Suggala, Sivaraman Balakrishnan, and Pradeep Ravikumar. Robust estimation via robust gradient estimation. *arXiv preprint arXiv:1802.06485*, 2018.
- [39] Peter J Rousseeuw and Annick M Leroy. *Robust regression and outlier detection*, volume 589. John wiley & sons, 2005.
- [40] Jie Shen and Ping Li. A tight bound of hard thresholding. *The Journal of Machine Learning Research*, 18(1):7650–7691, 2017.

- [41] Charles M Stein. Estimation of the mean of a multivariate normal distribution. *The annals of Statistics*, pages 1135–1151, 1981.
- [42] John W Tukey. Mathematics and the picturing of data. In *Proceedings of the International Congress of Mathematicians, Vancouver, 1975*, volume 2, pages 523–531, 1975.
- [43] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [44] Vincent Q Vu, Juhee Cho, Jing Lei, and Karl Rohe. Fantope projection and selection: A near-optimal convex relaxation of sparse PCA. In *Advances in neural information processing systems*, pages 2670–2678, 2013.
- [45] Martin Wainwright. *High-dimensional statistics: A non-asymptotic viewpoint*. Cambridge University Press, 2019.
- [46] Huan Xu, Constantine Caramanis, and Shie Mannor. Outlier-robust PCA: the high-dimensional case. *IEEE transactions on information theory*, 59(1):546–572, 2013.
- [47] Huan Xu, Constantine Caramanis, and Sujay Sanghavi. Robust PCA via outlier pursuit. *IEEE Transactions on Information Theory*, 58(5):3047–3064, 2012.
- [48] Yuchen Zhang, Martin J Wainwright, and Michael I Jordan. Lower bounds on the performance of polynomial-time algorithms for sparse linear regression. In *Conference on Learning Theory*, pages 921–948, 2014.

Notation. In proofs, we let \otimes denote the Kronecker product, and for a vector \mathbf{u} , we denote the outer product by $\mathbf{u}^{\otimes 2} = \mathbf{u}\mathbf{u}^\top$. We define the infinity norm for a matrix M as $\|M\|_\infty = \max_{i,j} |M_{ij}|$. Given index set \mathcal{J} , $\mathbf{v}^\mathcal{J}$ is the vector restricted to indices \mathcal{J} . Similarly, $M^{\mathcal{J}\mathcal{J}}$ is the sub-matrix on indices $\mathcal{J} \times \mathcal{J}$. we use $\{C_j\}_{j=0}^3$ to denote constants that are independent of dimension, but whose value can change from line to line.

A Proofs for the meta-theorem

In this section, we prove the global linear convergence guarantee given the [Definition 2.1](#). In each iteration of [Algorithm 1](#), we use $\widehat{\mathbf{G}}^t$ to update

$$\boldsymbol{\beta}^{t+1} = \text{P}_{k'}(\boldsymbol{\beta}^t - \eta \widehat{\mathbf{G}}^t),$$

where $\eta = 1/\mu_\beta$ is a fixed step size. Given the condition $\|\widehat{\mathbf{G}}(\boldsymbol{\beta}) - \mathbf{G}(\boldsymbol{\beta})\|_2^2 \leq \alpha(\epsilon)\|\mathbf{G}(\boldsymbol{\beta})\|_2^2 + \psi(\epsilon)$ in RSGE's definition, we show that [Algorithm 1](#) linearly converges to a neighborhood around $\boldsymbol{\beta}^*$ with error at most $O(\sqrt{\psi(\epsilon)})$.

First, we introduce a supporting Lemma from [\[40\]](#), which bounds the distance between $\text{P}_{k'}(\boldsymbol{\beta}^t - \eta \widehat{\mathbf{G}}^t)$ and $\boldsymbol{\beta}^*$ in each iteration of [Algorithm 1](#).

Lemma A.1 (Theorem 1 in [\[40\]](#)). *Let $\mathbf{z} \in \mathbb{R}^d$ be an arbitrary vector and $\boldsymbol{\beta}^* \in \mathbb{R}^d$ be any k -sparse signal. For any $k' \geq k$, we have the following bound:*

$$\|\text{P}_{k'}(\mathbf{z}) - \boldsymbol{\beta}^*\|_2 \leq \sqrt{\zeta}\|\mathbf{z} - \boldsymbol{\beta}^*\|_2, \quad \zeta = 1 + \frac{\rho + \sqrt{(4+\rho)\rho}}{2}, \quad \rho = \frac{\min\{k, d-k'\}}{k' - k + \min\{k, d-k'\}}.$$

We choose the hard thresholding parameter $k' = kc_\kappa^2 \ll d$, hence $\rho = 1/c_\kappa^2$.

Theorem A.1 ([Theorem 2.1](#)). *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from [Model 1.1](#). [Algorithm 1](#), with $\psi(\epsilon)$ -RSGE defined in [Definition 2.1](#), with step size $\eta = 1/\mu_\beta$ outputs $\widehat{\boldsymbol{\beta}}$, such that*

$$\|\widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^*\|_2 = O(\sqrt{\psi(\epsilon)}),$$

with probability at least $1 - \nu$, by setting $k' = c_\kappa^2 k$ and $T = \Theta(\log(\|\boldsymbol{\beta}^*\|_2/\sqrt{\psi(\epsilon)}))$. The sample complexity is $N(k, d, \epsilon, \nu) = n(k, d, \epsilon, \nu/T)T$.

Proof. By splitting N samples into T sets (each set has sample size n), [Algorithm 1](#) collects a fresh batch of samples with size $n(k, d, \epsilon, \nu/T)$ at each iteration $t \in [T]$. [Definition 2.1](#) shows that for the fixed gradient expectation \mathbf{G}^t , the estimate for the gradient $\widehat{\mathbf{G}}^t$ satisfies:

$$\|\widehat{\mathbf{G}}^t - \mathbf{G}^t\|_2^2 \leq \alpha(\epsilon)\|\mathbf{G}^t\|_2^2 + \psi(\epsilon) \tag{6}$$

with probability at least $1 - \nu/T$, where $\alpha(\epsilon)$ is determined by ϵ .

Letting $\mathbf{z}^t = \boldsymbol{\beta}^t - \eta \widehat{\mathbf{G}}^t$, we study the t -th iteration of [Algorithm 1](#). Based on [Lemma A.1](#), we have

$$\|\boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^*\|_2 \leq \sqrt{\zeta}\|\boldsymbol{\beta}^t - \eta \widehat{\mathbf{G}}^t - \boldsymbol{\beta}^*\|_2$$

$$\begin{aligned}
&= \sqrt{\zeta} \left\| \boldsymbol{\beta}^t - \eta \mathbf{G} - \boldsymbol{\beta}^* + \eta(\mathbf{G} - \widehat{\mathbf{G}}) \right\|_2 \\
&\leq \sqrt{\zeta} \left\| \boldsymbol{\beta}^t - \eta \mathbf{G} - \boldsymbol{\beta}^* \right\|_2 + \sqrt{\zeta} \eta \left\| \mathbf{G} - \widehat{\mathbf{G}} \right\|_2 \\
&\stackrel{(i)}{\leq} \sqrt{\zeta} \left\| (\mathbf{I}_d - \eta \boldsymbol{\Sigma})(\boldsymbol{\beta}^t - \boldsymbol{\beta}^*) \right\|_2 + \sqrt{\zeta} \eta \sqrt{\alpha(\epsilon) \left\| \mathbf{G} \right\|_2^2 + \psi(\epsilon)} \\
&\stackrel{(ii)}{\leq} \sqrt{\zeta} \left\| (\mathbf{I}_d - \eta \boldsymbol{\Sigma})(\boldsymbol{\beta}^t - \boldsymbol{\beta}^*) \right\|_2 + \sqrt{\zeta} \eta \sqrt{\alpha(\epsilon)} \left\| \boldsymbol{\Sigma}(\boldsymbol{\beta}^t - \boldsymbol{\beta}^*) \right\|_2 + \sqrt{\zeta} \eta \sqrt{\psi(\epsilon)}
\end{aligned}$$

where (i) follows from the theoretical guarantee of RSGE, and (ii) follows from the basic inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ for non-negative a, b .

By setting $\eta = 1/\mu_\beta$, we have

$$\begin{aligned}
\left\| \boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^* \right\|_2 &\leq \sqrt{\zeta} \left\| (\mathbf{I}_d - \eta \boldsymbol{\Sigma})(\boldsymbol{\beta}^t - \boldsymbol{\beta}^*) \right\|_2 + \sqrt{\zeta} \eta \sqrt{\alpha(\epsilon)} \left\| \boldsymbol{\Sigma}(\boldsymbol{\beta}^t - \boldsymbol{\beta}^*) \right\|_2 + \sqrt{\zeta} \eta \sqrt{\psi(\epsilon)} \\
&\leq \sqrt{\zeta} \left(1 - \frac{1}{c_\kappa}\right) \left\| \boldsymbol{\beta}^t - \boldsymbol{\beta}^* \right\|_2 + \sqrt{\zeta} \sqrt{\alpha(\epsilon)} \left\| \boldsymbol{\beta}^t - \boldsymbol{\beta}^* \right\|_2 + \sqrt{\zeta} \eta \sqrt{\psi(\epsilon)} \\
&\leq \sqrt{\zeta} \left(1 - \frac{1}{c_\kappa} + \sqrt{\alpha(\epsilon)}\right) \left\| \boldsymbol{\beta}^t - \boldsymbol{\beta}^* \right\|_2 + \sqrt{\zeta} \eta \sqrt{\psi(\epsilon)}
\end{aligned} \tag{7}$$

When ϵ is a small enough constant, we have $\sqrt{\alpha(\epsilon)} \leq \frac{1}{2c_\kappa}$, then

$$\begin{aligned}
\sqrt{\zeta} \left(1 - \frac{1}{c_\kappa} + \sqrt{\alpha(\epsilon)}\right) &\leq \sqrt{\zeta} \left(1 - \frac{1}{2c_\kappa}\right) \\
&\leq \sqrt{1 + \frac{\rho + \sqrt{(4+\rho)\rho}}{2}} \left(1 - \frac{1}{2c_\kappa}\right)
\end{aligned}$$

Plugging in the parameter $\rho = 1/c_\kappa^2$ in [Lemma A.1](#), we have

$$\sqrt{\zeta} \left(1 - \frac{1}{c_\kappa} + \sqrt{\alpha(\epsilon)}\right) \leq 1 - \frac{1}{10c_\kappa}$$

Together with [eq. \(7\)](#), we have the recursion

$$\left\| \boldsymbol{\beta}^{t+1} - \boldsymbol{\beta}^* \right\|_2 \leq \left(1 - \frac{1}{10c_\kappa}\right) \left\| \boldsymbol{\beta}^t - \boldsymbol{\beta}^* \right\|_2 + \sqrt{\zeta} \eta \sqrt{\psi(\epsilon)}.$$

By solving this recursion and using a union bound, we have

$$\left\| \boldsymbol{\beta}^t - \boldsymbol{\beta}^* \right\|_2 \leq \left(1 - \frac{1}{10c_\kappa}\right)^t \left\| \boldsymbol{\beta}^0 - \boldsymbol{\beta}^* \right\|_2 + \frac{\sqrt{\zeta} \eta \sqrt{\psi(\epsilon)}}{1 - \left(1 - \frac{1}{10c_\kappa}\right)} \leq (4\alpha(\epsilon))^t \left\| \boldsymbol{\beta}^* \right\|_2^2 + 10c_\kappa \sqrt{\zeta} \eta \sqrt{\psi(\epsilon)},$$

with probability at least $1 - \nu$.

By the definition of c_κ and η , we have $\left\| \widehat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* \right\|_2 = O\left(\frac{\sqrt{\psi(\epsilon)}}{\mu_\alpha}\right)$ □

B Correcting Lemma A.3 in [1]’s proof

A key part of the proof of the main theorem in [1] is to obtain an upper bound on the k -sparse operator norm. Specifically, their Lemmas A.2 and A.3 aim to show:

$$\lambda^* \geq \left\| \sum_{i=1}^{|\mathcal{S}|} w_i \left(\mathbf{g}_i - \widehat{\mathbf{G}}(w) \right)^{\otimes 2} - F \left(\widehat{\mathbf{G}}(w) \right) \right\|_{\widetilde{\mathbf{k}}, \text{op}} \geq \frac{\left\| \mathbf{P}_{\widetilde{\mathbf{k}}} \left(\widetilde{\Delta}(w) \right) \right\|_2^2}{5\epsilon}, \quad (8)$$

where $\widehat{\mathbf{G}}(w) = \mathbf{P}_{2\widetilde{k}} \left(\sum_{i=1}^{|\mathcal{S}|} w_i \mathbf{g}_i \right)$, $\widetilde{\Delta}(w) = \sum_{i=1}^{|\mathcal{S}|} w_i \mathbf{g}_i - \mathbf{G}^4$, and recall λ^* is the solution to the SDP as given in Algorithm 3.

Lemma A.3 asserts the first inequality above, and Lemma A.2 the second. As we show below, Lemma A.3 cannot be correct. Specifically, the issue is that the quantity inside the second term in eq. (8) may not be positive semidefinite. In this case, the convex optimization problem whose solution is λ^* is not a valid relaxation, and hence the λ^* they obtain need not a valid upper bound. Indeed, we give a simple example below that illustrates precisely this potential issue.

Fortunately, not all is lost – indeed, as our results imply, the main results in [1] is correct. The key is to show that while λ^* does not upper bound the sparse operator norm, it does, however, upper bound the quantity

$$\max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \widetilde{k}} \mathbf{v}^\top \left(\sum_{i=1}^{|\mathcal{S}|} w_i \left(\mathbf{g}_i - \widehat{\mathbf{G}}(w) \right)^{\otimes 2} - F \left(\widehat{\mathbf{G}}(w) \right) \right) \mathbf{v}. \quad (9)$$

We show this in Appendix D. More specifically, in Lemma D.3, we replace the \widetilde{k} -sparse operator norm in the second term of eq. (8) by the term in eq. (9). We show this can be used to complete the proof in Appendix D.4.

We now provide a counterexample that shows the first inequality in (8) cannot hold. The main argument is that the convex relaxation for sparse PCA is a valid upper bound of the sparse operator norm only for positive semidefinite matrices. Specifically, denoting $\mathbf{E} = \widehat{\Sigma}(w) - F(\widehat{\mathbf{G}}(w))$ as the matrix in eq. (9), [1] solves the following convex program:

$$\max_{\mathbf{H}} \text{Tr}(\mathbf{E} \cdot \mathbf{H}), \quad \text{subject to } \mathbf{H} \succcurlyeq 0, \|\mathbf{H}\|_{1,1} \leq k, \text{Tr}(\mathbf{H}) = 1.$$

Since $\widehat{\Sigma}(w) - F(\widehat{\mathbf{G}}(w))$ is no longer a p.s.d. matrix, the trace maximization above may not be a valid convex relaxation, and thus not an upper bound. Let us consider a specific example, in robust sparse mean estimation for $\mathcal{N}(\mu, \mathbf{I}_d)$, where function $F(\cdot)$ is a fixed identity matrix \mathbf{I}_d . We choose $\widetilde{k} = 1$, $\mu = [1, 0]^\top$, and $d = 2$. Suppose we observe data to be $x_1 = [2.5, 0]^\top$, $x_2 = [0, 0]^\top$, and the weights for x_1 and x_2 are the same. Then, we can compute the following matrices as:

$$\widehat{\Sigma} = \begin{bmatrix} 1.5625 & 0 \\ 0 & 0 \end{bmatrix}, F = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \mathbf{E} = \widehat{\Sigma} - F = \begin{bmatrix} 0.5625 & 0 \\ 0 & -1 \end{bmatrix}.$$

⁴The $\{w_i\}$ are weights, and these are defined precisely in Section D, but are not required for the present discussion or counterexample.

It is clear that $\|\widehat{\Sigma} - F\|_{\tilde{k}, \text{op}} = 1$. Solving the convex relaxation $\max_{\mathbf{H}} \text{Tr}(\mathbf{E} \cdot \mathbf{H})$ or $\max_{\mathbf{H}} \text{Tr}(\widehat{\Sigma} \cdot \mathbf{H})$ gives answer $\mathbf{H}^* = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$ and the corresponding $\lambda^* = 0.5625$, which is clearly not an upper bound of $\|\widehat{\Sigma} - F\|_{\tilde{k}, \text{op}}$. Hence $\lambda^* \geq \|\widehat{\Sigma} - F\|_{\tilde{k}, \text{op}}$ cannot hold in general.

C Covariance smoothness properties in robust sparse mean estimation

When the covariance is identity, the ellipsoid algorithm requires a closed form expression of the true covariance function $F(\mathbf{G})$. Indeed, the ellipsoid-based robust sparse mean estimation algorithm uses the covariance structure given by $F(\cdot)$ to detect outliers. The accuracy of robust sparse mean estimation explicitly depends on the properties of $F(\mathbf{G})$. L_{cov} and L_F are two important properties of $F(\mathbf{G})$, related to its smoothness. We first provide a closed-form expression for F , and then define precisely smoothness parameters L_{cov} and L_F , and show how these can be controlled.

Closed form expression of $F(\mathbf{G})$.

Lemma C.1. *Suppose we observe i.i.d. samples $\{\mathbf{z}_i, i \in \mathcal{G}\}$ from the distribution P in [Model 1.1](#) with $\Sigma = \mathbf{I}_d$, we have the covariance of gradient as*

$$\text{Cov}(\mathbf{g}) = \mathbb{E}_{\mathbf{z}_i \sim P} \left((\mathbf{g}_i - \mathbf{G})(\mathbf{g}_i - \mathbf{G})^\top \right) = \|\mathbf{G}\|_2^2 \mathbf{I}_d + \mathbf{G}\mathbf{G}^\top + \sigma^2 \mathbf{I}_d.$$

Proof. Since $\mathbf{g}_i = \mathbf{x}_i (\mathbf{x}_i^\top \boldsymbol{\beta} - y_i)$, and $\mathbf{G} = \mathbb{E}_{\mathbf{z}_i \sim P}(\mathbf{g}_i)$ and $\Sigma = \mathbf{I}_d$, we have

$$\begin{aligned} \mathbb{E}_{\mathbf{z}_i \sim P} \left((\mathbf{g}_i - \mathbf{G})(\mathbf{g}_i - \mathbf{G})^\top \right) &= \mathbb{E}_P \left((\mathbf{x}\mathbf{x}^\top - \mathbf{I}_d) \mathbf{G}\mathbf{G}^\top (\mathbf{x}\mathbf{x}^\top - \mathbf{I}_d) \right) + \sigma^2 \mathbf{I}_d \\ &= \mathbb{E}_P \left(\mathbf{x}\mathbf{x}^\top \mathbf{G}\mathbf{G}^\top \mathbf{x}\mathbf{x}^\top \right) - 2 \mathbb{E}_P \left(\mathbf{x}\mathbf{x}^\top \mathbf{G}\mathbf{G}^\top \right) + \mathbf{G}\mathbf{G}^\top + \sigma^2 \mathbf{I}_d, \end{aligned}$$

where we drop i in \mathbf{x}_i without abuse of notation.

Next, we apply the Stein-type Lemma [41] for $\mathbf{x} \sim \mathcal{N}(0, \mathbf{I}_d)$, and a function $f(\mathbf{x})$ whose second derivative exists:

$$\mathbb{E} \left(f(\mathbf{x}) \mathbf{x}\mathbf{x}^\top \right) = \mathbb{E} \left(f(\mathbf{x}) \right) \mathbf{I}_d + \mathbb{E} \left(\nabla^2 f(\mathbf{x}) \right). \quad (10)$$

By eq. (10), we have

$$\text{Cov}(\mathbf{g}) = \|\mathbf{G}\|_2^2 \mathbf{I}_d + \mathbf{G}\mathbf{G}^\top + \sigma^2 \mathbf{I}_d.$$

□

Smoothness properties of $\|F\|_{\text{op}}$. We first assume

$$L_{\text{cov}} = \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} |\mathbf{v}^\top \text{Cov}(\mathbf{g}) \mathbf{v}|. \quad (11)$$

If we define the functional $F(\cdot)$, such that $F(\widehat{\mathbf{G}}) = \|\widehat{\mathbf{G}}\|_2^2 \mathbf{I}_d + \widehat{\mathbf{G}}\widehat{\mathbf{G}}^\top + \sigma^2 \mathbf{I}_d$, and $F(\mathbf{G}) = \|\mathbf{G}\|_2^2 \mathbf{I}_d + \mathbf{G}\mathbf{G}^\top + \sigma^2 \mathbf{I}_d$, then we assume that there exists L_F satisfying

$$\left\| F(\mathbf{G}) - F(\widehat{\mathbf{G}}) \right\|_{\text{op}} \leq L_F \|\mathbf{G} - \widehat{\mathbf{G}}\|_2 + C \|\mathbf{G} - \widehat{\mathbf{G}}\|_2^2, \quad (12)$$

where C is a universal constant.

Lemma C.2. *Under the same setting as Lemma C.1, we have*

$$L_{\text{cov}} = 2\|\mathbf{G}\|_2^2 + \sigma^2, \text{ and } L_F = 4\|\mathbf{G}\|_2.$$

Proof. L_{cov} is upper bounded by the top eigenvalue of $F(\mathbf{G})$,

$$L_{\text{cov}} \leq \|F(\mathbf{G})\|_2 \leq 2\|\mathbf{G}\|_2^2 + \sigma^2.$$

For the L_F term, we have

$$\begin{aligned} & \left\| F(\mathbf{G}) - F(\widehat{\mathbf{G}}) \right\|_{\text{op}} \\ &= \left\| 2\mathbf{G}^\top (\mathbf{G} - \widehat{\mathbf{G}}) \mathbf{I}_d - \|\mathbf{G} - \widehat{\mathbf{G}}\|_2^2 \mathbf{I}_d + \mathbf{G} (\mathbf{G} - \widehat{\mathbf{G}})^\top + (\mathbf{G} - \widehat{\mathbf{G}}) \mathbf{G}^\top - (\mathbf{G} - \widehat{\mathbf{G}}) (\mathbf{G} - \widehat{\mathbf{G}})^\top \right\|_{\text{op}} \\ &\leq 4\|\mathbf{G}\|_2 \|\mathbf{G} - \widehat{\mathbf{G}}\|_2 + 2\|\mathbf{G} - \widehat{\mathbf{G}}\|_2^2. \end{aligned}$$

Therefore, we can choose $L_F = 4\|\mathbf{G}\|_2$ and $C = 2$. □

D Proofs for the ellipsoid algorithm in robust sparse regression

In this section, we prove guarantees for the ellipsoid algorithm in robust sparse regression. In the theoretical analysis of the ellipsoid algorithm, we use \mathcal{S}_{in} to denote the observations \mathcal{S} , which shares the same notations with Algorithm 3. We first give preliminary definitions of error terms defined on $\mathcal{S}_{\text{good}}$ and \mathcal{S}_{in} , and then prove Lemma D.1. Next, we prove concentration results for gradients of uncorrupted sparse linear regression in Lemma D.2. In Lemma D.3, we provide lower bounds for the \tilde{k} -sparse largest eigenvalue defined in eq. (9). Finally, we prove Corollary 3.1 based on previous Lemmas in Appendix D.4.

D.1 Preliminary definitions and properties related to $\mathcal{S}_{\text{good}}, \mathcal{S}_{\text{bad}}$

Here, we state again the definitions of $\mathcal{S}_{\text{good}}, \mathcal{S}_{\text{bad}}$ and \mathcal{S}_{in} . In Algorithm 3, we denote the input set as \mathcal{S}_{in} , which can be partitioned into two parts: $\mathcal{S}_{\text{good}} = \{i : i \in \mathcal{G} \text{ and } i \in \mathcal{S}_{\text{in}}\}$, and $\mathcal{S}_{\text{bad}} = \{i : i \in \mathcal{B} \text{ and } i \in \mathcal{S}_{\text{in}}\}$. Note that $\mathcal{S}_{\text{in}} = \mathcal{S}_{\text{good}} \cup \mathcal{S}_{\text{bad}}$, and $n = |\mathcal{S}_{\text{in}}|$. For the convenience of our analysis, we

define the following error terms:

$$\begin{aligned}\tilde{\Delta}_{\mathcal{S}_{\text{good}}} &= \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}}(\mathbf{g}_i) - \mathbf{G}, \\ \widehat{\Delta}_{\mathcal{S}_{\text{good}}} &= \mathbf{P}_{2\tilde{k}}(\mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}}(\mathbf{g}_i)) - \mathbf{G}, \\ \tilde{\Delta} &= \mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}}(\mathbf{g}_i) - \mathbf{G}, \\ \widehat{\Delta} &= \mathbf{P}_{2\tilde{k}}(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}}(\mathbf{g}_i)) - \mathbf{G}.\end{aligned}$$

These error terms are defined under a uniform distribution over samples, whereas previous papers using ellipsoid algorithms consider a set of balanced weighted distribution. More specifically, the weights in our setting are defined as:

$$\tilde{w}_i = \frac{1}{n}, \quad \forall i \in \mathcal{S}_{\text{good}} \cup \mathcal{S}_{\text{bad}}.$$

The balanced weighted distribution is defined to satisfy:

$$0 \leq w_i \leq \frac{1}{(1-2\epsilon)n}, \quad \forall i \in \mathcal{S}_{\text{good}} \cup \mathcal{S}_{\text{bad}}, \quad \sum_{i \in \mathcal{S}_{\text{in}}} w_i = 1.$$

Notice that $\sum_{i \in \mathcal{S}_{\text{bad}}} \tilde{w}_i = O(\epsilon)$, and $\sum_{i \in \mathcal{S}_{\text{bad}}} w_i = O\left(\frac{\epsilon}{1-2\epsilon}\right)$ with high probability, which intuitively says that both types of distributions have $O(\epsilon)$ weights over all bad samples. We are interested in considering uniform weighted samples since this formulation helps us analyze the filtering algorithm more conveniently, as we show in the following sections.

We restate the following Lemma which shows the connection of these different error terms.

Lemma D.1 (Lemma A.1 in [1]). *Suppose G is k -sparse. Then we have the following result:*

$$\frac{1}{5} \|\widehat{\Delta}\|_2 \leq \|\mathbf{P}_k(\tilde{\Delta})\|_2 \leq 4 \|\widehat{\Delta}\|_2.$$

D.2 Concentration bounds for gradients in $\mathcal{S}_{\text{good}}$

We first prove concentration bounds for gradients for sparse linear regression in the uncorrupted case. The following is similar to Lemma D.1 in [1].

Lemma D.2. *Suppose we observe i.i.d. gradient samples $\{\mathbf{g}_i, i \in \mathcal{G}\}$ from [Model 1.1](#) with $|\mathcal{G}| = \Omega\left(\frac{k \log(d/\nu)}{\epsilon^2}\right)$. Then, there is a $\delta = \tilde{O}(\epsilon)$, such that with probability at least $1 - \nu$, for any index subset $\mathcal{J} \subset [d]$, $|\mathcal{J}| \leq \tilde{k}$ and for any $\mathcal{G}' \subset \mathcal{G}$, $|\mathcal{G}'| \geq (1-2\epsilon)|\mathcal{G}|$, the following inequalities hold:*

$$\|\mathbb{E}_{i \in_u \mathcal{G}'}(\mathbf{g}_i^{\mathcal{J}}) - \mathbf{G}^{\mathcal{J}}\|_2 \leq \delta (\|\mathbf{G}\|_2 + \sigma), \quad (13)$$

$$\left\| \mathbb{E}_{i \in_u \mathcal{G}'}(\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right\|_{\text{op}} \leq \delta (\|\mathbf{G}\|_2^2 + \sigma^2). \quad (14)$$

Proof. The main difference from their Lemma D.1 is that we consider a uniform distribution over all samples instead of a balanced weighted distribution. Furthermore, eqs. (13) and (14) are the concentration inequalities for the mean and covariance of the collected gradient samples $\{\mathbf{g}_i, i \in \mathcal{G}\}$ in

the good set with the form:

$$\mathbf{g}_i = \mathbf{x}_i \mathbf{x}_i^\top \mathbf{G} - \mathbf{x}_i \xi_i,$$

which is equivalent to their Lemma D.1, where they consider $y_i \mathbf{x}_i = \mathbf{x}_i \mathbf{x}_i^\top \boldsymbol{\beta} + \mathbf{x}_i \xi_i$. Therefore, by setting all weights to $\frac{1}{(1-2\epsilon)|\mathcal{G}|}$ in their Lemma D.1 we obtain the desired concentration properties. \square

D.3 Relationship between the first and second moment of samples in \mathcal{S}_{in}

In this part, we show an important connection between the covariance deviation (the empirical covariance of \mathcal{S}_{in} minus the true covariance of authentic data) and the mean deviation (the empirical mean of \mathcal{S}_{in} minus the true mean of authentic data). When the mean deviation (in ℓ_2 sense) is large, the following Lemma implies that the covariance deviation must also be large. As a result, when the magnitude of the covariance deviation is large, the current set of samples (or the current weights of all samples) needs to be adjusted; when the magnitude of the covariance deviation is small, the average of current sample set (or the weighted sum of samples using current weights) provides a good enough estimate of the model parameter. Moreover, the same principle holds when we use an approximation of the true covariance, which can be efficiently estimated.

Unlike Lemma A.2 in [1], in eq. (17), eq. (18), we provide lower bounds for the \tilde{k} -sparse largest eigenvalue (rigorous definition in eq. (20)), instead of the \tilde{k} -sparse operator norm. As we discussed in Appendix B, λ^* is the convex relaxation of finding the \tilde{k} -sparse largest eigenvalue (instead of the \tilde{k} -sparse operator norm). In the statement of the following Lemma, for the purpose of consistency, we consider the uniform distribution of weights. However, the proof and results can be easily extended to the setting with the balanced distribution of weights. This is due to the similarity between the two types of weight representation, as discussed in Appendix D.1.

Lemma D.3. *Suppose $|\mathcal{S}_{\text{bad}}| \leq 2\epsilon|\mathcal{S}_{\text{in}}|$, $\delta = \Omega(\epsilon)$, and the gradient samples in $\mathcal{S}_{\text{good}}$ satisfy*

$$\left\| \mathbf{P}_{\tilde{k}} \left(\tilde{\Delta}_{\mathcal{S}_{\text{good}}} \right) \right\|_2 \leq c (\|\mathbf{G}\|_2 + \sigma) \delta, \quad (15)$$

$$\left\| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} (\mathbf{g}_i - \mathbf{G})^{\otimes 2} - F(\mathbf{G}) \right\|_{\tilde{k}, \text{op}} \leq c \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) \delta, \quad (16)$$

where c is a constant. If $\left\| \mathbf{P}_{\tilde{k}} \left(\tilde{\Delta} \right) \right\|_2 \geq C_1 (\|\mathbf{G}\|_2 + \sigma) \delta$, where C_1 is a large constant, we have,

$$\max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} (\mathbf{g}_i - \hat{\mathbf{G}})^{\otimes 2} - F(\mathbf{G}) \right) \mathbf{v} \geq \frac{\left\| \mathbf{P}_{\tilde{k}} \left(\tilde{\Delta} \right) \right\|_2^2}{4\epsilon}, \quad (17)$$

$$\max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} (\mathbf{g}_i - \hat{\mathbf{G}})^{\otimes 2} - F(\hat{\mathbf{G}}) \right) \mathbf{v} \geq \frac{\left\| \mathbf{P}_{\tilde{k}} \left(\tilde{\Delta} \right) \right\|_2^2}{5\epsilon}. \quad (18)$$

Proof. We focus on the \tilde{k} -sparse largest eigenvalue (rigorous definition in eq. (20)), which is the correct route of analysis the convex relaxation of Sparse PCA.

Let $\mathcal{J} = \arg \max_{\mathcal{J}' \subset [d], |\mathcal{J}'| \leq \tilde{k}} \left\| \tilde{\Delta}^{\mathcal{J}'} \right\|_2$. Then $\tilde{\Delta}^{\mathcal{J}} = \left\| \mathbf{P}_{\tilde{k}} \left(\tilde{\Delta} \right) \right\|_2 \geq C_1 (\|\mathbf{G}\|_2 + \sigma) \delta$ according to the

assumption. Using $|\mathcal{S}_{\text{in}}|$ to denote the size of \mathcal{S}_{in} , we have a lower bound for the sum over bad samples:

$$\begin{aligned}
\left\| \frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{bad}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}}) \right\|_2 &= \left\| \tilde{\Delta}^{\mathcal{J}} - \frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{good}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}}) \right\|_2 \\
&\geq \left\| \tilde{\Delta}^{\mathcal{J}} \right\|_2 - \left\| \frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{good}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}}) \right\|_2 \\
&\stackrel{(i)}{\geq} \left\| \tilde{\Delta}^{\mathcal{J}} \right\|_2 - c(\|\mathbf{G}\|_2 + \sigma) \delta \\
&\stackrel{(ii)}{\geq} \frac{\left\| \tilde{\Delta}^{\mathcal{J}} \right\|_2}{1.1},
\end{aligned}$$

where (i) follows from eq. (15) and the assumptions; (ii) follows from that we choose C_1 large enough.

By p.s.d.-ness of covariance matrices, we have

$$\frac{1}{|\mathcal{S}_{\text{bad}}|} \sum_{i \in \mathcal{S}_{\text{bad}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}}) (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\top} \succcurlyeq \left(\frac{1}{|\mathcal{S}_{\text{bad}}|} \sum_{i \in \mathcal{S}_{\text{bad}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}}) \right)^{\otimes 2}.$$

Therefore, because $|\mathcal{S}_{\text{bad}}| \leq 2\epsilon|\mathcal{S}_{\text{in}}|$, we have

$$\left\| \frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{bad}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} \right\|_{\text{op}} \geq \frac{\left\| \frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{bad}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}}) \right\|_2^2}{2\epsilon} \geq \frac{\left\| \tilde{\Delta}^{\mathcal{J}} \right\|_2^2}{2.5\epsilon}. \quad (19)$$

With a lower bound of this submatrix of the covariance matrix, we define a vector $\mathbf{v}_0 \in \mathbb{R}^{\tilde{k}}$ as follows:

$$\mathbf{v}_0 = \arg \max_{\|\mathbf{v}\|_2=1} \mathbf{v}^{\top} \left(\sum_{i \in \mathcal{S}_{\text{bad}}} \frac{1}{|\mathcal{S}_{\text{in}}|} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} \right) \mathbf{v}. \quad (20)$$

For this \mathbf{v}_0 , we have

$$\begin{aligned}
&\mathbf{v}_0^{\top} \left(\frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i=1}^{|\mathcal{S}_{\text{in}}|} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right) \mathbf{v}_0 \\
&\geq \mathbf{v}_0^{\top} \left(\frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{bad}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} \right) \mathbf{v}_0 \\
&- \left\| \frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i \in \mathcal{S}_{\text{good}}} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - \frac{|\mathcal{S}_{\text{good}}|}{|\mathcal{S}_{\text{in}}|} F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right\|_{\text{op}} - \left\| \frac{|\mathcal{S}_{\text{bad}}|}{|\mathcal{S}_{\text{in}}|} F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right\|_{\text{op}} \\
&\stackrel{(i)}{\geq} \frac{\left\| \tilde{\Delta}^{\mathcal{J}} \right\|_2^2}{2.5\epsilon} - c(\|\mathbf{G}\|_2^2 + \sigma^2) \delta - 2\epsilon(\|\mathbf{G}\|_2^2 + \sigma^2) \\
&\stackrel{(ii)}{\geq} \frac{\left\| \tilde{\Delta}^{\mathcal{J}} \right\|_2^2}{3\epsilon}, \quad (21)
\end{aligned}$$

where (i) follows from eq. (16) and eq. (19); (ii) follows from the assumption that ϵ is sufficiently small.

Applying eq. (21) on our target $\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} (\mathbf{g}_i - \widehat{\mathbf{G}})^{\otimes 2} - F(\mathbf{G})$, we have

$$\begin{aligned}
& \mathbf{v}_0^\top \left(\frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i=1}^{|\mathcal{S}_{\text{in}}|} (\mathbf{g}_i^{\mathcal{J}} - \widehat{\mathbf{G}}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right) \mathbf{v}_0 \\
&= \mathbf{v}_0^\top \left(\frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i=1}^{|\mathcal{S}_{\text{in}}|} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} - \widehat{\Delta}^{\mathcal{J}} (\widetilde{\Delta}^{\mathcal{J}})^\top - \widetilde{\Delta}^{\mathcal{J}} (\widehat{\Delta}^{\mathcal{J}})^\top + (\widehat{\Delta}^{\mathcal{J}})^{\otimes 2} \right) \mathbf{v}_0 \\
&\stackrel{(i)}{\geq} \mathbf{v}_0^\top \left(\frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i=1}^{|\mathcal{S}_{\text{in}}|} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right) \mathbf{v}_0 - 24 \left(\|\widetilde{\Delta}^{\mathcal{J}}\|_2^2 \right) \\
&\stackrel{(ii)}{\geq} \frac{\|\widetilde{\Delta}^{\mathcal{J}}\|_2^2}{4\epsilon}, \tag{22}
\end{aligned}$$

where (i) follows from Lemma D.1; (ii) follows from eq. (21) and ϵ is sufficiently small. By a construction $\mathbf{v} = (\mathbf{v}_0, \mathbf{0}_{d-\tilde{k}})^\top$, it is easy to see that \mathbf{v}_0 provides a lower bound for the maximum of $\{\mathbf{v} : \|\mathbf{v}\|_2 = 1, \|\mathbf{v}\|_0 \leq \tilde{k}\}$ in eq. (17).

By eq. (22), we already know that

$$\mathbf{v}_0^\top \left(\frac{1}{|\mathcal{S}_{\text{in}}|} \sum_{i=1}^{|\mathcal{S}_{\text{in}}|} (\mathbf{g}_i^{\mathcal{J}} - \widehat{\mathbf{G}}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right) \mathbf{v}_0 \geq \frac{\|\widetilde{\Delta}^{\mathcal{J}}\|_2^2}{4\epsilon}.$$

By our assumptions on F , we have

$$\begin{aligned}
\|F(\mathbf{G}) - F(\widehat{\mathbf{G}})\|_{\tilde{k}, \text{op}} &\leq L_F \|\widehat{\Delta}\|_2 + C \|\widehat{\Delta}\|_2^2 \\
&\stackrel{(i)}{\leq} 5L_F \|\widetilde{\Delta}^{\mathcal{J}}\|_2 + 5C \|\widetilde{\Delta}^{\mathcal{J}}\|_2^2,
\end{aligned}$$

where (i) follows from Lemma D.1. Since $\delta = \Omega(\epsilon)$, we obtain eq. (18) by using the triangle inequality. \square

D.4 Proof of Corollary 3.1

Equipped with Lemma D.1, Lemma D.2 and Lemma D.3, we can now prove Corollary 3.1.

Corollary D.1 (Corollary 3.1). *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from Model 1.1 with $\Sigma = \mathbf{I}_d$. By setting $\tilde{k} = k' + k$, if we use the ellipsoid algorithm for robust sparse gradient estimation with $\rho_{\text{sep}} = \Theta(\epsilon(\|\mathbf{G}^t\|_2^2 + \sigma^2))$, it requires $N(k, d, \epsilon, \nu) = \Omega\left(\frac{k^2 \log(dT/\nu)}{\epsilon^2}\right)T$ samples, and guarantees $\psi(\epsilon) = \tilde{O}(\epsilon^2 \sigma^2)$. Hence, Algorithm 1 outputs $\widehat{\beta}$, such that*

$$\|\widehat{\beta} - \beta^*\|_2 = \tilde{O}(\sigma\epsilon),$$

with probability at least $1 - \nu$, by setting $T = \Theta\left(\log\left(\frac{\|\beta^*\|_2}{\epsilon\sigma}\right)\right)$.

Proof. We consider only the t -th iteration, and thus omit t in \mathbf{g}_i^t and \mathbf{G}^t . The function $F(\mathbf{G})$ is given by $F(\mathbf{G}) = \|\mathbf{G}\|_2^2 \mathbf{I}_d + \mathbf{G}\mathbf{G}^\top + \sigma^2 \mathbf{I}_d$, as in [Appendix C](#). The accuracy in robust sparse estimation on gradients depends on two parameters for $F(\mathbf{G})$: $L_{\text{cov}} = 2\|\mathbf{G}\|_2^2 + \sigma^2$, and $L_F = 4\|\mathbf{G}\|_2$, which are calculated in [Appendix C](#).

Under the statistical model and the contamination model described in [Theorem 2.1](#), we can set the parameters $\rho_{\text{sep}} = \Theta(\epsilon(\|\mathbf{G}^t\|_2^2 + \sigma^2))$ in [Algorithm 2](#) by the calculation of L_{cov} and L_F .

The ellipsoid algorithm considers all possible sample weights in a convex set and finds the optimal weight for each sample. The algorithm iteratively uses a separation oracle [Algorithm 2](#), which solves the convex relaxation of Sparse PCA at each iteration:

$$\lambda^* = \max_{\mathbf{H}} \text{Tr} \left(\left(\widehat{\Sigma} - F(\widehat{\mathbf{G}}) \right) \cdot \mathbf{H} \right), \quad \text{subject to } \mathbf{H} \succcurlyeq 0, \|\mathbf{H}\|_{1,1} \leq \tilde{k}, \text{Tr}(\mathbf{H}) = 1. \quad (23)$$

To prove the Main Theorem (Theorem 3.1) in [1], the only modification is to replace the lower bound of λ^* in their Lemma A.3.

A weighted version of [Lemma D.3](#) implies that if the mean deviation is large, then

$$\max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\sum_{i=1}^{|\mathcal{S}_{\text{in}}|} w_i (\mathbf{g}_i - \widehat{\mathbf{G}}(w))^{\otimes 2} - F(\widehat{\mathbf{G}}(w)) \right) \mathbf{v} \geq \frac{\left\| \mathbf{P}_{\tilde{k}}(\tilde{\Delta}(w)) \right\|_2^2}{5\epsilon}, \quad (24)$$

where $\widehat{\mathbf{G}}(w) = \mathbf{P}_{2\tilde{k}} \left(\sum_{i=1}^{|\mathcal{S}_{\text{in}}|} w_i \mathbf{g}_i \right)$, and $\tilde{\Delta}(w) = \sum_{i=1}^{|\mathcal{S}_{\text{in}}|} w_i \mathbf{g}_i - \mathbf{G}$. Then, λ^* in the ellipsoid algorithm satisfies

$$\lambda^* \geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\sum_{i=1}^{|\mathcal{S}_{\text{in}}|} w_i (\mathbf{g}_i - \widehat{\mathbf{G}}(w))^{\otimes 2} - F(\widehat{\mathbf{G}}(w)) \right) \mathbf{v}, \quad (25)$$

since λ^* is the solution to the trace norm maximization [eq. \(23\)](#), which is the convex relaxation of finding the \tilde{k} -sparse largest eigenvalue.

Combining [eq. \(24\)](#) and [eq. \(25\)](#), we have

$$\lambda^* \geq \frac{\left\| \mathbf{P}_{\tilde{k}}(\tilde{\Delta}(w)) \right\|_2^2}{5\epsilon}, \quad (26)$$

which recovers the correctness of the separation oracle in the ellipsoid algorithm, and their Main Theorem (Theorem 3.1).

Finally, the ellipsoid algorithm guarantees that, with sample complexity $\Omega \left(\frac{k^2 \log(d/\nu)}{\epsilon^2} \right)$, the estimate $\widehat{\mathbf{G}}$ satisfies

$$\left\| \widehat{\mathbf{G}} - \mathbf{G} \right\|_2^2 = \tilde{O} \left(\epsilon^2 (L_F^2 + L_{\text{cov}}) \right) = \tilde{O} \left(\epsilon^2 \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) \right), \quad (27)$$

with probability at least $1 - \nu$. This exactly gives us a $\tilde{O}(\epsilon^2 \sigma^2)$ -RSGE. Hence, we can apply [eq. \(27\)](#) as the RSGE in [Theorem 2.1](#) to prove [Corollary 3.1](#). \square

E Outlier removal guarantees in the filtering algorithm

In this section, we consider a single iteration of [Algorithm 1](#), and prove [Lemma 4.1](#) at the t -th step. For clarity, we omit the superscript t in both \mathbf{g}_i^t and \mathbf{G}^t .

In order to show guarantees for [Lemma 4.1](#), we leverage previous results [Lemma D.2](#) and [Lemma D.3](#). We state [Lemma E.1](#) as a modification of [Lemma D.2](#) by replacing ϵ by $\sqrt{\epsilon}$, using concentration results in [Lemma D.2](#), and replacing ϵ by $\sqrt{\epsilon}$. We state [Lemma E.2](#) as a modification of [Lemma D.3](#) by replacing $\delta = \Omega(\epsilon)$ with $\delta = \Omega(\sqrt{\epsilon})$, since the results for $\delta = \Omega(\epsilon)$ implies the results for $\delta = \Omega(\sqrt{\epsilon})$.

The reason we modify the above is to prove guarantees for our computationally more efficient RSGE described in [Algorithm 3](#). Our motivation for calculating the score for each sample according to $\tau_i = \text{Tr}(\mathbf{H}^* \cdot (\mathbf{g}_i - \widehat{\mathbf{G}})(\mathbf{g}_i - \widehat{\mathbf{G}})^\top)$ is to make sure that all the scores τ_i are positive (notice that the scores calculated based on the original non-p.s.d matrix may be negative). Based on this, we show that the sum of scores over all bad samples is a large constant (> 1) times larger than the sum of scores over all good samples. When finding an upper bound for $\sum_{i \in \mathcal{S}_{\text{good}}} \tau_i$, we compromise an ϵ factor in the value of λ^* , which results in an $\sqrt{\epsilon}$ factor in the recovery guarantee.

As described above, we immediately have [Lemma E.1](#) and [Lemma E.2](#) given the proofs in [Appendix D](#). Note that we still use the same definitions $\widetilde{\Delta}_{\mathcal{S}_{\text{good}}}$ and $\widetilde{\Delta}$ on set $\mathcal{S}_{\text{good}}$ and \mathcal{S}_{in} respectively as in [Appendix D.1](#).

Lemma E.1. *Suppose we observe i.i.d. gradient samples $\{\mathbf{g}_i, i \in \mathcal{G}\}$ from [Model 1.1](#) with $|\mathcal{G}| = \Omega\left(\frac{k \log(d/\nu)}{\epsilon}\right)$. Then there is a $\delta = \widetilde{O}(\sqrt{\epsilon})$ that with probability at least $1 - \nu$, we have for any subset $\mathcal{J} \subset [d]$, $|\mathcal{J}| \leq \widetilde{k}$, and for any $\mathcal{G}' \subset \mathcal{G}$, $|\mathcal{G}'| \geq (1 - 2\epsilon)|\mathcal{G}|$, the following inequalities hold:*

$$\left\| \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{g}_i^{\mathcal{J}}) - \mathbf{G}^{\mathcal{J}} \right\|_2 \leq \delta (\|\mathbf{G}\|_2 + \sigma), \quad (28)$$

$$\left\| \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right\|_{\text{op}} \leq \delta \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right). \quad (29)$$

Lemma E.2. *Suppose $|\mathcal{S}_{\text{bad}}| \leq 2\epsilon|\mathcal{S}_{\text{in}}|$, $\delta = \Omega(\sqrt{\epsilon})$, and the gradient samples in $\mathcal{S}_{\text{good}}$ satisfy*

$$\left\| \mathbb{P}_{\widetilde{k}} \left(\widetilde{\Delta}_{\mathcal{S}_{\text{good}}} \right) \right\|_2 \leq c (\|\mathbf{G}\|_2 + \sigma) \delta, \quad (30)$$

$$\left\| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} (\mathbf{g}_i - \mathbf{G})^{\otimes 2} - F(\mathbf{G}) \right\|_{\widetilde{k}, \text{op}} \leq c \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) \delta, \quad (31)$$

where c is a constant. If $\left\| \mathbb{P}_{\widetilde{k}} \left(\widetilde{\Delta} \right) \right\|_2 \geq C_1 (\|\mathbf{G}\|_2 + \sigma) \delta$, where C_1 is a constant. Then we have,

$$\max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \widetilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} (\mathbf{g}_i - \widehat{\mathbf{G}})^{\otimes 2} - F(\mathbf{G}) \right) \mathbf{v} \geq \frac{\left\| \mathbb{P}_{\widetilde{k}} \left(\widetilde{\Delta} \right) \right\|_2^2}{4\epsilon}, \quad (32)$$

$$\max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \widetilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} (\mathbf{g}_i - \widehat{\mathbf{G}})^{\otimes 2} - F(\widehat{\mathbf{G}}) \right) \mathbf{v} \geq \frac{\left\| \mathbb{P}_{\widetilde{k}} \left(\widetilde{\Delta} \right) \right\|_2^2}{5\epsilon}. \quad (33)$$

By [Lemma E.1](#), eq. (30) and eq. (31) in [Lemma E.2](#) are satisfied, provided that we have $|\mathcal{G}| = \Omega\left(\frac{k \log(d/\nu)}{\epsilon}\right)$. Now, equipped with [Lemma E.1](#) and [Lemma E.2](#), the effect of good samples can be controlled by concentration inequalities. Based on these, we are ready to prove [Lemma 4.1](#).

Lemma E.3 (Lemma 4.1). *Suppose we observe $n = \Omega\left(\frac{k^2 \log(d/\nu)}{\epsilon}\right)$ ϵ -corrupted samples from [Model 1.1](#) with $\Sigma = \mathbf{I}_d$. Let \mathcal{S}_{in} be an ϵ -corrupted set of gradient samples $\{\mathbf{g}_i^t\}_{i=1}^n$. [Algorithm 3](#) computes λ^* that satisfies*

$$\lambda^* \geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} \right) \mathbf{v}. \quad (34)$$

If $\lambda^* \geq \rho_{\text{sep}} = C_\gamma \left(\|\mathbf{G}^t\|_2^2 + \sigma^2 \right)$, then with probability at least $1 - \nu$, we have

$$\sum_{i \in \mathcal{S}_{\text{good}}} \tau_i \leq \frac{1}{\gamma} \sum_{i \in \mathcal{S}_{\text{in}}} \tau_i, \quad (35)$$

where τ_i is defined in [line 10](#), C_γ is a constant depending on γ , and $\gamma \geq 4$ is a constant.

Proof. Since λ^* is the solution of the convex relaxation of Sparse PCA, we have

$$\begin{aligned} \lambda^* &= \text{Tr} \left(H^* \cdot \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} \right) \right) \\ &\geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} \right) \mathbf{v}. \end{aligned}$$

By [Theorem A.1](#) in [\[1\]](#), we have

$$\begin{aligned} &\text{Tr} \left(H^* \cdot \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} - F(\widehat{\mathbf{G}}) \right) \right) \\ &\leq C \left(\left\| \widehat{\Delta} \right\|_2^2 + \left(L_F + \tilde{k} \left\| \widetilde{\Delta}_{\mathcal{S}_{\text{good}}} \right\|_\infty \right) \left\| \widehat{\Delta} \right\|_2 + \tilde{k} \left\| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \mathbf{G} \right)^{\otimes 2} - F(\mathbf{G}) \right\|_\infty \right), \quad (36) \end{aligned}$$

where C is a constant. Noticing that $\left\| \widetilde{\Delta}_{\mathcal{S}_{\text{good}}} \right\|_\infty$ and $\left\| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \mathbf{G} \right)^{\otimes 2} - F(\mathbf{G}) \right\|_\infty$ are unrelated to $\widehat{\mathbf{G}}$ and only defined on $\mathcal{S}_{\text{good}}$, [\[1\]](#) shows concentration bounds for these two terms, when $n = \Omega\left(\frac{\tilde{k}^2 \log(d/\nu)}{\epsilon}\right)$. Specifically, it showed that with probability at least $1 - \nu$, we have

$$\left\| \widetilde{\Delta}_{\mathcal{S}_{\text{good}}} \right\|_\infty \leq C_1 \left(L_F + \sqrt{L_{\text{cov}}} \right) \sqrt{\epsilon/\tilde{k}} \quad (37)$$

$$\left\| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \mathbf{G} \right)^{\otimes 2} - F(\mathbf{G}) \right\|_\infty \leq C_1 \left(L_F^2 + L_{\text{cov}} \right) \sqrt{\epsilon/\tilde{k}} \quad (38)$$

Now, we focus on the LHS of [eq. \(35\)](#), the sum of scores of points in $\mathcal{S}_{\text{good}}$. By definition, we have

$$\begin{aligned} &\mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \tau_i \\ &= \text{Tr} \left(H^* \cdot \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} \right) \right) \\ &= \text{Tr} \left(H^* \cdot \left(\mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} - F(\widehat{\mathbf{G}}) \right) \right) + \text{Tr} \left(H^* F(\widehat{\mathbf{G}}) \right) \\ &\stackrel{(i)}{\leq} C \left(\left\| \widehat{\Delta} \right\|_2^2 + \left(L_F + \tilde{k} \left\| \widetilde{\Delta}_{\mathcal{S}_{\text{good}}} \right\|_\infty \right) \left\| \widehat{\Delta} \right\|_2 + \tilde{k} \left\| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \left(\mathbf{g}_i - \mathbf{G} \right)^{\otimes 2} - F(\mathbf{G}) \right\|_\infty \right) \end{aligned}$$

$$+ \text{Tr} \left(H^* F \left(\widehat{\mathbf{G}} \right) \right),$$

where (i) follows from eq. (36).

To bound the RHS above, we first bound $\text{Tr} \left(H^* \cdot F \left(\widehat{\mathbf{G}} \right) \right)$. Because of the constraint of the SDP given in eq. (3), H^* belongs to the Fantope \mathcal{F}^1 [44], and thus for any matrix A , we have $\text{Tr} (A \cdot H^*) \leq \|A\|_{\text{op}}$.

Thus, we have

$$\begin{aligned} \text{Tr} \left(H^* \cdot F \left(\widehat{\mathbf{G}} \right) \right) &= \text{Tr} \left(H^* \cdot F \left(\mathbf{G} \right) \right) + \text{Tr} \left(H^* * \left(F \left(\widehat{\mathbf{G}} \right) - F \left(\mathbf{G} \right) \right) \right) \\ &\leq \|F \left(\mathbf{G} \right)\|_{\text{op}} + \left\| F \left(\widehat{\mathbf{G}} \right) - F \left(\mathbf{G} \right) \right\|_{\text{op}} \\ &\stackrel{(i)}{\leq} C_1 \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) + \left\| F \left(\widehat{\mathbf{G}} \right) - F \left(\mathbf{G} \right) \right\|_{\text{op}} \\ &\stackrel{(ii)}{\leq} C_1 \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) + L_F \|\widehat{\Delta}\|_2 + C_2 \|\widehat{\Delta}\|_2^2, \end{aligned} \quad (39)$$

where (i) follows from the expression of $F(G)$ in Appendix C; (ii) from the smoothness of $F(G)$.

By plugging in the concentration guarantees eq. (37) and combining eq. (39), we have

$$\begin{aligned} &\mathbb{E}_{i \in \mathcal{U} \mathcal{S}_{\text{good}}} \tau_i \\ &\leq C_2 \left((L_F^2 + L_{\text{cov}}) \sqrt{\epsilon} + \left((L_F + \sqrt{L_{\text{cov}}}) \sqrt{\epsilon} + L_F \right) \|\widehat{\Delta}\|_2 + \|\widehat{\Delta}\|_2^2 \right) + C_1 \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) \\ &\stackrel{(i)}{\leq} C_2 \left(\|\mathbf{G}\|_2 \|\widehat{\Delta}\|_2 + \|\widehat{\Delta}\|_2^2 \right) + C_1 \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right) \\ &\leq C_1 \left(\|\mathbf{G}\|_2 \|\widehat{\Delta}\|_2 + \|\widehat{\Delta}\|_2^2 + \|\mathbf{G}\|_2^2 + \sigma^2 \right), \end{aligned} \quad (40)$$

where (i) follows from the fact that ϵ is sufficiently small.

On the other hand, we know that: $\mathbb{E}_{i \in \mathcal{U} \mathcal{S}_{\text{in}}} \tau_i = \lambda^*$.

Now, under the condition $\lambda^* \geq \rho_{\text{sep}} = \Theta \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, we consider two cases separately. By separating two cases, we can always show λ^* is very large, and the contribution from good samples is limited.

First, if $\|\widehat{\Delta}\|_2^2 \geq \Theta \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, then in eq. (40), we have

$$\|\widehat{\Delta}\|_2^2 \gtrsim \|\mathbf{G}\|_2 \|\widehat{\Delta}\|_2 \gtrsim \|\mathbf{G}\|_2^2, \quad \text{and} \quad \|\widehat{\Delta}\|_2^2 \gtrsim \sigma^2.$$

Thus, we only need to compare λ^* and $\|\widehat{\Delta}\|_2^2$. By Lemma E.2, we have

$$\begin{aligned} \mathbb{E}_{i \in \mathcal{U} \mathcal{S}_{\text{in}}} \tau_i = \lambda^* &\geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in \mathcal{U} \mathcal{S}_{\text{in}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} \right) \mathbf{v} \\ &\geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq \tilde{k}} \mathbf{v}^\top \left(\mathbb{E}_{i \in \mathcal{U} \mathcal{S}_{\text{in}}} \left(\mathbf{g}_i - \widehat{\mathbf{G}} \right)^{\otimes 2} - F \left(\widehat{\mathbf{G}} \right) \right) \mathbf{v} \\ &\geq \frac{\|\widehat{\Delta}\|_2^2}{\epsilon}. \end{aligned}$$

Hence, by eq. (40), we have $\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} \tau_i \geq \gamma \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \tau_i$, where $\gamma \geq 4$ is a constant.

Second, if $\|\widehat{\Delta}\|_2^2 \leq \Theta \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, then in eq. (40), we have

$$\|\mathbf{G}\|_2^2 \gtrsim \|\mathbf{G}\|_2 \|\widehat{\Delta}\|_2 \gtrsim \|\widehat{\Delta}\|_2^2, \quad \text{or } \sigma^2 \gtrsim \|\widehat{\Delta}\|_2^2.$$

Thus, we only need to compare λ^* and $\max \left(\|\mathbf{G}\|_2^2, \sigma^2 \right)$. Since $\lambda^* \geq C_\gamma \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$ by the condition of Lemma 4.1, we still have $\mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} \tau_i \geq \gamma \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \tau_i$, where $\gamma \geq 4$ is a constant.

Combing all of above, and setting $\rho_{\text{sep}} = C_\gamma \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, we have

$$\sum_{i \in \mathcal{S}_{\text{in}}} \tau_i = |\mathcal{S}_{\text{in}}| \mathbb{E}_{i \in_u \mathcal{S}_{\text{in}}} \tau_i \geq \gamma |\mathcal{S}_{\text{good}}| \mathbb{E}_{i \in_u \mathcal{S}_{\text{good}}} \tau_i = \gamma \sum_{i \in \mathcal{S}_{\text{good}}} \tau_i.$$

□

F RSGE via the filtering algorithm

In this section, we still consider the t -th iteration of Algorithm 1 and prove Theorem 4.1 on t . We omit t in \mathbf{g}_i^t and \mathbf{G}^t .

In the case of $\lambda^* \geq C_\gamma \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, Algorithm 3 iteratively removes one sample according to the probability distribution eq. (4). We denote the steps of this outlier removal procedure as $l = 1, 2, \dots, n$. The first step of proving Theorem 4.1 is to show we can remove a corrupted samples with high probability at each step, which is a result by Lemma 4.1.

Intuitively, if all subsequent steps are i.i.d., we can expect Algorithm 3 to remove outliers within around en iterations, with exponentially high probability. However, the subsequent steps in Algorithm 3 are not independent. To circumvent this challenge we appeal to a martingale argument.

F.1 Supermartingale construction

Let \mathcal{F}^l be the filtration generated by the set of events until iteration l of Algorithm 3. We define the corresponding set $\mathcal{S}_{\text{in}}^l$, $\mathcal{S}_{\text{good}}^l$ and $\mathcal{S}_{\text{bad}}^l$ at the step l . We have that $\mathcal{S}_{\text{in}}^l, \mathcal{S}_{\text{good}}^l, \mathcal{S}_{\text{bad}}^l \in \mathcal{F}^l$, and $|\mathcal{S}_{\text{in}}^l| = n - l$.

We denote a good event \mathcal{E}^l at step l as

$$\sum_{i \in \mathcal{S}_{\text{bad}}^l} \tau_i \leq (\gamma - 1) \sum_{i \in \mathcal{S}_{\text{good}}^l} \tau_i.$$

Then, by the definition of Algorithm 3 and Lemma 4.1, if $\lambda^* \geq C_\gamma \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, \mathcal{E}^l is not true; if \mathcal{E}^l is true, then Algorithm 3 will return a $\widehat{\mathbf{G}}$.

In Lemma F.1, we show that at any step l when \mathcal{E}^l is not true, the random outlier removal procedure removes a corrupted sample with probability at least $(\gamma - 1) / \gamma$.

Lemma F.1. *In each subsequent step l , if \mathcal{E}^l is not true, then we can remove one remaining outlier*

from $\mathcal{S}_{\text{in}}^l$ with probability at least $(\gamma - 1) / \gamma$:

$$\Pr(\text{one sample from } \mathcal{S}_{\text{bad}}^l \text{ is removed} \mid \mathcal{F}_l) \geq \frac{\gamma - 1}{\gamma}.$$

Proof of Lemma F.1. When $\lambda^* \geq C_\gamma (\|\mathbf{G}\|_2^2 + \sigma^2)$, Lemma 4.1 implies

$$\sum_{i \in \mathcal{S}_{\text{bad}}^l} \tau_i \geq (\gamma - 1) \sum_{i \in \mathcal{S}_{\text{good}}^l} \tau_i.$$

Then we randomly remove a sample r from \mathcal{S}_{in} according to

$$\Pr(\mathbf{g}_i \text{ is removed} \mid \mathcal{F}_l) = \frac{\tau_i}{\sum_{i \in \mathcal{S}_{\text{in}}^l} \tau_i}.$$

Finally,

$$\Pr(\text{one sample from } \mathcal{S}_{\text{bad}}^l \text{ is removed} \mid \mathcal{F}_l) = \sum_{i \in \mathcal{S}_{\text{bad}}^l} \frac{\tau_i}{\sum_{i \in \mathcal{S}_{\text{in}}^l} \tau_i} \geq \frac{\gamma - 1}{\gamma}.$$

□

Since subsequent steps for applying Algorithm 3 on \mathcal{S}_{in} are not independent, we need martingale arguments to show the total iterations of applying Algorithm 3 is limited.

We use the martingale technique in [46], by defining $T: T = \min\{l : \mathcal{E}^l \text{ is true}\}$. Based on T , we define a random variable:

$$Y^l = \begin{cases} |\mathcal{S}_{\text{bad}}^{T-1}| + \frac{\gamma-1}{\gamma} (T-1), & \text{if } l \geq T \\ |\mathcal{S}_{\text{bad}}^l| + \frac{\gamma-1}{\gamma} l, & \text{if } l < T \end{cases}$$

Lemma F.2 (Lemma 1 in [46]). $\{Y^l, \mathcal{F}^l\}$ is a supermartingale.

Now, equipped with Lemma F.1 and Lemma F.2, we are ready to prove Theorem 4.1.

F.2 Proof of Theorem 4.1

Theorem F.1 (Theorem 4.1). Suppose we observe $n = \Omega\left(\frac{k^2 \log(d/\nu)}{\epsilon}\right)$ ϵ -corrupted samples from Model 1.1 with $\Sigma = \mathbf{I}_d$. Let \mathcal{S}_{in} be an ϵ -corrupted set of gradient samples $\{\mathbf{g}_i^t\}_{i=1}^n$. By setting $\tilde{k} = k' + k$, if we run Algorithm 3 iteratively with initial set \mathcal{S}_{in} , and subsequently on \mathcal{S}_{out} , and use $\rho_{\text{sep}} = C_\gamma (\|\mathbf{G}^t\|_2^2 + \sigma^2)$, then this repeated use of Algorithm 3 will stop after at most $\frac{1.1\gamma}{\gamma-1} \epsilon n$ iterations, and output $\widehat{\mathbf{G}}^t$, such that

$$\left\| \widehat{\mathbf{G}}^t - \mathbf{G}^t \right\|_2^2 = \tilde{O}\left(\epsilon \left(\|\mathbf{G}^t\|_2^2 + \sigma^2\right)\right),$$

with probability at least $1 - \nu - \exp(-\Theta(\epsilon n))$. Here, C_γ is a constant depending on γ , where $\gamma \geq 4$ is a constant.

Proof. We analyze [Algorithm 3](#) by discussing a series of $\{\mathcal{E}^l\}$.

If \mathcal{E}^l is true, then $\lambda^* \leq \rho_{\text{sep}} = C_\gamma \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$. By [Lemma E.2](#), we have

$$\lambda^* \geq \max_{\|\mathbf{v}\|_2=1, \|\mathbf{v}\|_0 \leq k} \mathbf{v}^\top \left(\mathbb{E}_{i \in \mathcal{u}_S} (\mathbf{g}_i - \widehat{\mathbf{G}})^{\otimes 2} - F(\widehat{\mathbf{G}}) \right) \mathbf{v} \geq \frac{\left\| \mathbb{P}_{\tilde{k}}(\widetilde{\Delta}_S) \right\|_2^2}{5\epsilon}.$$

Plugging in $\lambda^* \leq C_\gamma \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)$, we have

$$\frac{1}{5} \left\| \widehat{\Delta}_S \right\|_2^2 \stackrel{(i)}{\leq} \left\| \mathbb{P}_{\tilde{k}}(\widetilde{\Delta}_S) \right\|_2^2 \leq 5\epsilon\lambda^* \leq O\left(\epsilon \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)\right),$$

where (i) follows from [Lemma D.1](#). Hence, when \mathcal{E}^l is true, [Algorithm 3](#) can return a $\widehat{\mathbf{G}}$, such that $\left\| \widehat{\mathbf{G}} - \mathbf{G} \right\|_2^2 \leq O\left(\epsilon \left(\|\mathbf{G}\|_2^2 + \sigma^2 \right)\right)$.

Then, we only need to show $\bigcup_{l=1}^L \mathcal{E}^l$ is true, where $L = \frac{1.1\gamma}{\gamma-1}\epsilon n$, with high probability. That said, we need to upper bound the probability

$$\Pr\left(\bigcap_{l=1}^L \overline{\mathcal{E}^l}\right) = \Pr(T \geq L) \leq \Pr\left(Y^L \geq \frac{\gamma-1}{\gamma}L\right) = \Pr(Y^L \geq 1.1\epsilon n). \quad (41)$$

Then, we can construct the martingale difference according to [\[46\]](#). Let $D^l = Y^l - Y^{l-1}$, where $Y^0 = \epsilon n$, and

$$\bar{D}^l = D^l - \mathbb{E}(D^l | D^1, \dots, D^{l-1}).$$

Thus $\{\bar{D}^l\}$ is a martingale difference process, and $\mathbb{E}(D^l | D^1, \dots, D^{l-1}) \leq 0$, since $\{Y^l\}$ is a supermartingale. Now, [eq. \(41\)](#) can be viewed as a bound for the sum of the associated martingale difference sequence.

$$Y^l - Y^0 = \sum_{j=1}^l D^j = \sum_{j=1}^l \bar{D}^j + \sum_{j=1}^l \mathbb{E}(D^j | D^1, \dots, D^{j-1}) \leq \sum_{j=1}^l \bar{D}^j.$$

Since we only remove one example from the set $\mathcal{S}_{\text{in}}^l$, we can guarantee $|D^l| \leq 1$ and $|\bar{D}^l| \leq 2$. For these bounded random variables, by applying the Azuma-Hoeffding inequality, we have

$$\begin{aligned} \Pr(Y^L \geq 1.1\epsilon n) &\leq \Pr\left(\sum_{l=1}^L \bar{D}^l \geq 0.1\epsilon n\right) \\ &\leq \exp\left(\frac{-(0.1\epsilon n)^2}{8L}\right). \end{aligned}$$

Plugging in $L = \frac{1.1\gamma}{\gamma-1}\epsilon n$, this probability is upper bounded by $\exp(-\Theta(\epsilon n))$.

Notice that $L = \frac{1.1\gamma}{\gamma-1}\epsilon n \leq 1.5\epsilon n$, by setting $\gamma \geq 4$. Hence, from $l = 1$ to L , we always have $|\mathcal{S}_{\text{bad}}^l| \leq 2\epsilon |\mathcal{S}_{\text{in}}^l|$. Then [Lemma E.1](#) and [Lemma E.2](#) hold and [Lemma 4.1](#) is still valid.

Combining all of the above, we have proven that, with exponentially high probability, [Algorithm 3](#)

returns a $\widehat{\mathbf{G}}$ satisfying $\|\widehat{\mathbf{G}} - \mathbf{G}\|_2^2 \leq O\left(\epsilon\left(\|\mathbf{G}\|_2^2 + \sigma^2\right)\right)$, within $\frac{1.1\gamma}{\gamma-1}\epsilon n$ iterations. \square

G Robust sparse regression with unknown covariance

In this section, we prove the guarantees for RSGE when the covariance matrix Σ is unknown, but each row and column is sparse. In this case, the population mean of all authentic gradients \mathbf{G}^t can be calculated as

$$\mathbf{G}^t = \mathbb{E}_P(\mathbf{g}_i^t) = \mathbb{E}_P(\mathbf{x}_i \mathbf{x}_i^\top (\beta^t - \beta^*)) = \Sigma \omega^t.$$

Therefore, $\mathbf{G}^t = \Sigma \omega^t$ is guaranteed to be $r(k' + k)$ sparse. And we use the filtering algorithm (Algorithm 3) with $\tilde{k} = r(k' + k)$ as a RSGE.

First, we derive the functional $F(\mathbf{G})$ with general covariance matrix Σ , and compute the corresponding L_F, L_{cov} , which has been defined in eq. (11) and eq. (12) for the case $\Sigma = \mathbf{I}_d$ in Appendix C.

Lemma G.1. *Suppose we observe i.i.d. samples $\{\mathbf{z}_i, i \in \mathcal{G}\}$ from the distribution P in Model 1.1 with an unknown Σ , we have the covariance of gradient as*

$$\text{Cov}(\mathbf{g}) := \mathbb{E}_{\mathbf{z}_i \sim P} \left((\mathbf{g}_i - \mathbf{G})(\mathbf{g}_i - \mathbf{G})^\top \right) = \Sigma \left\| \Sigma^{-\frac{1}{2}} \mathbf{G} \right\|_2^2 + \mathbf{G} \mathbf{G}^\top + \sigma^2 \Sigma.$$

Proof. As in the Model 1.1, we draw \mathbf{x} from Gaussian distribution $\mathcal{N}(0, \Sigma)$, the expression of $F(\cdot)$ is given by

$$\begin{aligned} \text{Cov}(\mathbf{g}) &= \mathbb{E} \left((\mathbf{g}_i - \mathbf{G})(\mathbf{g}_i - \mathbf{G})^\top \right) \\ &= \mathbb{E} \left((\mathbf{x} \mathbf{x}^\top - \Sigma) \omega \omega^\top (\mathbf{x} \mathbf{x}^\top - \Sigma) \right) + \sigma^2 \Sigma \\ &\stackrel{(i)}{=} \mathbb{E} \left(\Sigma^{\frac{1}{2}} (\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top - \mathbf{I}_d) \Sigma^{\frac{1}{2}} \omega \omega^\top \Sigma^{\frac{1}{2}} (\tilde{\mathbf{x}} \tilde{\mathbf{x}}^\top - \mathbf{I}_d) \Sigma^{\frac{1}{2}} \right) + \sigma^2 \Sigma \\ &\stackrel{(ii)}{=} \Sigma \left\| \Sigma^{-\frac{1}{2}} \mathbf{G} \right\|_2^2 + \mathbf{G} \mathbf{G}^\top + \sigma^2 \Sigma. \end{aligned}$$

where (i) follows from the re-parameterization $\mathbf{x} = \Sigma^{\frac{1}{2}} \tilde{\mathbf{x}}$, where $\tilde{\mathbf{x}} \sim \mathcal{N}(0, \mathbf{I}_d)$, and (ii) follows from the Stein-type Lemma as in Appendix C. \square

By Lemma G.1, we define the functional $F(\mathbf{G}) = \Sigma \left\| \Sigma^{-\frac{1}{2}} \mathbf{G} \right\|_2^2 + \mathbf{G} \mathbf{G}^\top + \sigma^2 \Sigma$. In Algorithm 3, we do not need to evaluate $F(\cdot)$, but our analysis requires upper bounds for two parameters of $F(\cdot) - L_{\text{cov}}, L_F$ – to control tail bounds. Under the same setting as Lemma G.1, we use similar bounds as Appendix C, based on assumptions in Model 1.1. Hence, we have $L_{\text{cov}} = \Theta(\|\mathbf{G}\|_2^2 + \sigma^2)$, and $L_F = \Theta(\|\mathbf{G}\|_2)$.

Next, we show concentration bounds (Lemma G.2) similar to Lemma E.1, which controls deviation of empirical mean and covariance for all samples in the good set \mathcal{G} .

Lemma G.2. *Suppose we observe i.i.d. gradient samples $\{\mathbf{g}_i, i \in \mathcal{G}\}$ from Model 1.1 with $|\mathcal{G}| = \tilde{\Omega}\left(\frac{\tilde{k} \log(d/\nu)}{\epsilon}\right)$. Then, there is a $\delta = \tilde{O}(\sqrt{\epsilon})$, such that with probability at least $1 - \nu$, for any index*

subset $\mathcal{J} \subset [d]$, $|\mathcal{J}| \leq \tilde{k}$ and for any $\mathcal{G}' \subset \mathcal{G}$, $|\mathcal{G}'| \geq (1 - 2\epsilon)|\mathcal{G}|$, we have

$$\left\| \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{g}_i^{\mathcal{J}}) - \mathbf{G}^{\mathcal{J}} \right\|_2 \leq \delta (\|\mathbf{G}\|_2 + \sigma), \quad (42)$$

$$\left\| \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \right\|_{\text{op}} \leq \delta (\|\mathbf{G}\|_2^2 + \sigma^2). \quad (43)$$

Proof. We prove the concentration inequality for the covariance eq. (43), the bound for mean eq. (42) is similar. For any index subset $\mathcal{J} \subset [d]$, $|\mathcal{J}| \leq \tilde{k}$, we can expand eq. (43) as follows,

$$\begin{aligned} & \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{g}_i^{\mathcal{J}} - \mathbf{G}^{\mathcal{J}})^{\otimes 2} - F(\mathbf{G})^{\mathcal{J}\mathcal{J}} \\ &= \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{x}^{\mathcal{J}} \mathbf{x}^{\mathcal{J}\top} \omega \omega^{\mathcal{J}\top} \mathbf{x}(\mathbf{x}^{\mathcal{J}})^{\top}) - \left(\boldsymbol{\Sigma}^{\mathcal{J}\mathcal{J}} \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \omega \right\|_2^2 + 2\mathbf{G}^{\mathcal{J}} (\mathbf{G}^{\mathcal{J}})^{\top} \right) \end{aligned} \quad (44)$$

$$- \mathbb{E}_{i \in_u \mathcal{G}'} (\mathbf{x} \mathbf{x}^{\top} \omega \omega^{\top} \boldsymbol{\Sigma})^{\mathcal{J}\mathcal{J}} + \mathbf{G}^{\mathcal{J}} (\mathbf{G}^{\mathcal{J}})^{\top} \quad (45)$$

$$+ \mathbb{E}_{i \in_u \mathcal{G}'} \xi_i^2 \mathbf{x}^{\mathcal{J}} (\mathbf{x}^{\mathcal{J}})^{\top} - \sigma^2 \boldsymbol{\Sigma}^{\mathcal{J}\mathcal{J}} \quad (46)$$

Here, we prove the concentration inequality for eq. (44), and the other two terms can be bounded by the same technique. It is sufficient to prove an upper bound for the operator norm as follows

$$\left\| \mathbb{E}_{i \in_u \mathcal{G}'} \mathbf{x}^{\mathcal{J}} (\mathbf{x}^{\mathcal{J}})^{\top} \omega^{\mathcal{J}} (\omega^{\mathcal{J}})^{\top} \mathbf{x}^{\mathcal{J}} (\mathbf{x}^{\mathcal{J}})^{\top} - \left(\boldsymbol{\Sigma}^{\mathcal{J}\mathcal{J}} \left\| \boldsymbol{\Sigma}^{\frac{1}{2}} \omega \right\|_2^2 + 2\mathbf{G}^{\mathcal{J}} (\mathbf{G}^{\mathcal{J}})^{\top} \right) \right\|_{\text{op}} \leq \delta \|\mathbf{G}\|_2^2, \quad (47)$$

where \mathbf{x} is drawn from a Gaussian distribution $\mathcal{N}(0, \boldsymbol{\Sigma})$. Note that the index subset \mathcal{J} reduce the matrix to $\mathbb{R}^{|\mathcal{J}| \times |\mathcal{J}|}$. For the concentration bounds of covariance matrix estimation eq. (47), we have a near identical argument as Lemma 4.5 of [13], by replacing Theorem 5.50 with Theorem 5.44 in [43].

This establishes eq. (47) with sample complexity $n = \tilde{\Omega} \left(\frac{\tilde{k} \log(1/\nu)}{\epsilon} \right)$, with probability at least $1 - \nu$. Next, we take a union bound over all possible subsets $\mathcal{J} \subset [d]$, and this gives concentration results for the covariance eq. (43). Hence we have proved the concentration results for the gradient under the assumption that $\boldsymbol{\Sigma}$ is row/column sparse. \square

Based on Lemma G.2, we have Theorem 5.1, which guarantees the recovery of $\boldsymbol{\beta}^*$ in robust sparse regression with unknown covariance as defined in Model 5.1.

Corollary G.1 (Theorem 5.1). *Suppose we observe $N(k, d, \epsilon, \nu)$ ϵ -corrupted samples from Model 1.1, where the covariates \mathbf{x}_i 's follow from Model 5.1. If we use Algorithm 3 for robust sparse gradient estimation, it requires $\tilde{\Omega} \left(\frac{r^2 k^2 \log(dT/\nu)}{\epsilon} \right) T$ samples, and $T = \Theta \left(\log \left(\frac{\|\boldsymbol{\beta}^*\|_2}{\sigma \sqrt{\epsilon}} \right) \right)$, then, we have*

$$\left\| \hat{\boldsymbol{\beta}} - \boldsymbol{\beta}^* \right\|_2 = \tilde{O}(\sigma \sqrt{\epsilon}), \quad (48)$$

with probability at least $1 - \nu - T \exp(-\Theta(\epsilon n))$.

Proof. With the concentration result Lemma G.2 in hand, the remaining parts share the same theoretical analysis as Appendix E and Appendix F, by replacing $(k' + k)^2$ with $r^2(k' + k)^2 = \Theta(r^2 k^2)$. Hence, we have a result similar to Corollary 4.1, with sample complexity $\tilde{\Omega} \left(\frac{r^2 k^2 \log(dT/\nu)}{\epsilon} \right)$. And this yields Theorem 5.1. \square

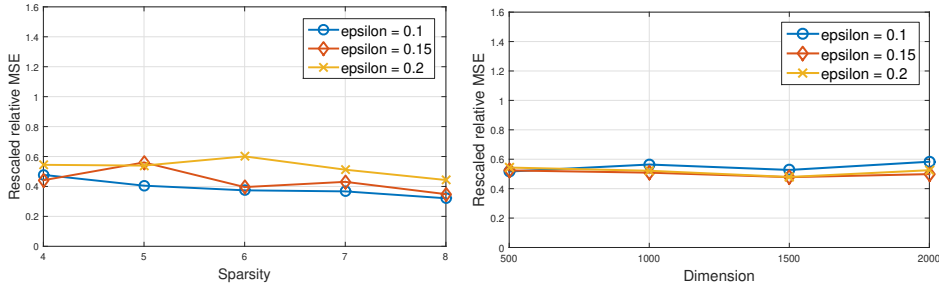


Figure 1: Simulations for Algorithm 3 showing the dependence of relative MSE on sparsity and dimension. For each parameter, we choose corresponding sample complexity $n \propto k^2 \log(d)/\epsilon$. Different curves for $\epsilon \in \{0.1, 0.15, 0.2\}$ are the average of 15 trials. Consistent with the theory, the rescaled relative MSE's are nearly independent of sparsity and dimension. Furthermore, by rescaling for different ϵ , three curves have the same magnitude.

H Numerical results

H.1 Robust sparse mean estimation

We first demonstrate the performance of Algorithm 3 for robust sparse mean estimation, and then move to Algorithm 1 for robust sparse regression. For the robust sparse gradient estimation, we generate samples through $\mathbf{g}_i = \mathbf{x}_i \mathbf{x}_i^\top \mathbf{G} - \mathbf{x}_i \xi_i$, where the unknown true mean \mathbf{G} is k -sparse. The authentic \mathbf{x}_i 's are generated from $\mathcal{N}(0, \mathbf{I}_d)$. We set $\sigma = 0$, since the main part of the error in robust sparse mean estimation is \mathbf{G} . Each entry of \mathbf{G} is either $+1$ or -1 , hence $\|\mathbf{G}\|_2^2 = k$.

The outliers are specially designed: the norm of the outliers is $\|\mathbf{G}\|_2$, and the directions are orthogonal to \mathbf{G} . Through this construction, outliers cannot be easily removed by simple pruning, and the directions of outliers can cause large effects on the estimation of \mathbf{G} . We plot the relative MSE of parameter recovery, defined as $\|\hat{\mathbf{G}} - \mathbf{G}\|_2^2 / \|\mathbf{G}\|_2^2$, with respect to different sparsities and dimensions.

Parameter error vs. sparsity k . We fix the dimension to be $d = 50$. We solve the trace norm maximization in Algorithm 3 using CVX [22]. We solve robust sparse gradient estimation under different levels of outlier fraction ϵ and different sparsity values k .

Parameter error vs. dimension d . We fix $k = 5$. We use a Sparse PCA solver from [19] which is much more efficient for higher dimensions. We run robust sparse gradient estimation Algorithm 3 under different levels of outlier fraction ϵ and different dimensions d .

For each parameter, the corresponding number of samples required for the authentic data is $n \propto k^2 \log(d)/\epsilon$ according to Theorem 4.1. Therefore, we add $\epsilon n / (1 - \epsilon)$ outliers (so that the outliers are an ϵ -fraction of the total samples), and then run Algorithm 3. According to Theorem 4.1, the rescaled relative MSE: $\|\hat{\mathbf{G}} - \mathbf{G}\|_2^2 / (\epsilon \|\mathbf{G}\|_2^2)$ should be independent of the parameters $\{\epsilon, k, d\}$. We plot this in Figure 1, and these plots validate our theorem on the sample complexity in robust sparse mean estimation problems.

H.2 Robust sparse regression with identity covariance

We use Algorithm 1 for robust sparse regression. Similarly as in Appendix H.1, we use Algorithm 3 as our Robust Sparse Gradient Estimator, and leverage the Sparse PCA solver from [19]. In the simulation, we fix $d = 500$, and $k = 5$, hence the corresponding sample complexity is $n \propto 1/\epsilon$. However,

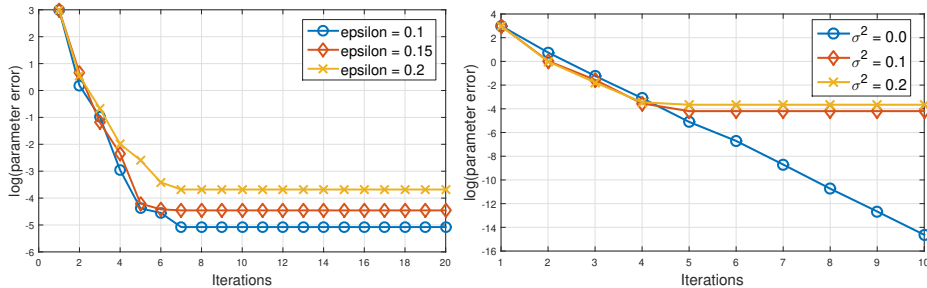


Figure 2: Empirical illustration of the linear convergence of $\log(\|\beta^t - \beta^*\|_2^2)$ vs. iteration counts in the Algorithm 1. In all cases, we fix $k = 5$, $d = 500$, and choose the sample complexity $n \propto 1/\epsilon$. The left plot considers different ϵ with fixed $\sigma^2 = 0.1$. The right plot considers different σ^2 with fixed $\epsilon = 0.1$. As expected, the convergence is linear, and flatten out at the level of the final error.

we do not use the sample splitting technique in the simulations.

The entries of the true parameter β^* are set to be either $+1$ or -1 , hence $\|\beta^*\|_2^2 = k$ is fixed. The authentic \mathbf{x}_i s are generated from $\mathcal{N}(0, \mathbf{I}_d)$, and the authentic $y_i = \mathbf{x}_i^\top \beta^* + \xi_i$ as in Model 1.1. We set the covariates of the outliers as A , where A is a random ± 1 matrix of dimension $\epsilon n / (1 - \epsilon) \times d$, and set the responses of outliers to $-A\beta^*$.

To show the performance of Algorithm 1 under different settings, we use different levels of ϵ and σ in Figure 2, and track the parameter error $\|\beta^t - \beta^*\|_2^2$ of Algorithm 1 in each iteration. Consistent with the theory, the algorithm displays linear convergence, and the error curves flatten out at the level of the final error. Furthermore, Algorithm 1 achieves machine precision when $\sigma^2 = 0$ in the right plot of Figure 2.

H.3 Robust sparse regression with unknown covariance matrix

Following Appendix H.2, we study the empirical performance of robust sparse regression with unknown covariance matrix Σ following from Model 5.1.

We use the same experimental setup as in Appendix H.2, but modify the covariance matrix to be a Toeplitz matrix with a decay $\Sigma_{ij} = \exp^{-(i-j)^2}$. Under this setting, the covariance matrix is sparse, thus follows from Model 5.1. Figure 3 indicates that we have nearly the same performance as the $\Sigma = \mathbf{I}_d$ case.

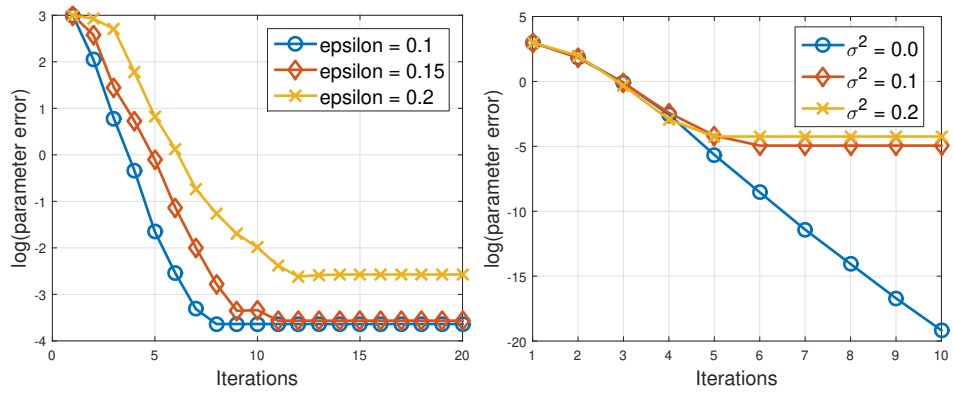


Figure 3: Empirical illustration of the linear convergence of $\log(\|\beta^t - \beta^*\|_2^2)$ vs. iteration counts in the Algorithm 1 with unknown covariance matrix which is a Toeplitz matrix with a decay $\Sigma_{ij} = \exp^{-(i-j)^2}$. The other settings are the same as Figure 2. Even though the covariance matrix is unknown, we observe similar performance in linear convergence as Figure 2.