

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

---

# Robust Regression via Hard Thresholding

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

We study the problem of Robust Least Squares Regression (RLSR) where several response variables can be adversarially corrupted. More specifically, for a data matrix  $X \in \mathbb{R}^{p \times n}$  and an underlying model  $\mathbf{w}^*$ , the response vector is generated as  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b}$  where  $\mathbf{b} \in \mathbb{R}^n$  is the corruption vector supported over at most  $C \cdot n$  coordinates. Existing exact recovery results for RLSR focus solely on  $L_1$ -penalty based convex formulations and impose relatively strict model assumptions such as requiring the corruptions  $\mathbf{b}$  to be selected independently of  $X$ .

In this work, we study a simple hard-thresholding algorithm called TORRENT which, under mild conditions on  $X$ , can recover  $\mathbf{w}^*$  exactly even if  $\mathbf{b}$  corrupts the response variables in an *adversarial* manner, i.e. both the support and entries of  $\mathbf{b}$  are selected adversarially after observing  $X$  and  $\mathbf{w}^*$ . Our results hold under *deterministic* assumptions which are satisfied if  $X$  is sampled from any sub-Gaussian distribution. Finally unlike existing results that apply only to a fixed  $\mathbf{w}^*$ , generated independently of  $X$ , our results are *universal* and hold for any  $\mathbf{w}^* \in \mathbb{R}^p$ .

Next, we propose gradient descent-based extensions of TORRENT that can scale efficiently to large scale problems, such as high dimensional sparse recovery. and prove similar recovery guarantees for these extensions. Empirically we find TORRENT, and more so its extensions, offering significantly faster recovery than the state-of-the-art  $L_1$  solvers. For instance, even on moderate-sized datasets (with  $p = 50K$ ) with around 40% corrupted responses, a variant of our proposed method called TORRENT-HYB is more than  $20\times$  faster than the best  $L_1$  solver.

*“If among these errors are some which appear too large to be admissible, then those equations which produced these errors will be rejected, as coming from too faulty experiments, and the unknowns will be determined by means of the other equations, which will then give much smaller errors.”*

A. M. Legendre, *On the Method of Least Squares*. 1805.

## 1 Introduction

Robust Least Squares Regression (RLSR) addresses the problem of learning a reliable set of regression coefficients in the presence of several arbitrary corruptions in the *response* vector. Owing to the wide-applicability of regression, RLSR features as a critical component of several important real-world applications in a variety of domains such as signal processing [1], economics [2], computer vision [3, 4], and astronomy [2].

Given a data matrix  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n]$  with  $n$  data points in  $\mathbb{R}^p$  and the corresponding response vector  $\mathbf{y} \in \mathbb{R}^n$ , the goal of RLSR is to learn a  $\hat{\mathbf{w}}$  such that,

$$(\hat{\mathbf{w}}, \hat{S}) = \arg \min_{\substack{\mathbf{w} \in \mathbb{R}^p \\ S \subset [n]: |S| \geq (1-\beta) \cdot n}} \sum_{i \in S} (y_i - \mathbf{x}_i^T \mathbf{w})^2, \quad (1)$$

054 That is, we wish to simultaneously determine the set of corruption free points  $\hat{S}$  and also estimate  
 055 the best model parameters over the set of clean points. However, the optimization problem given  
 056 above is non-convex (jointly in  $\mathbf{w}$  and  $S$ ) in general and might not directly admit efficient solutions.  
 057 Indeed there exist reformulations of this problem that are known to be NP-hard to optimize [1].

058 To address this problem, most existing methods with provable guarantees assume that the obser-  
 059 vations are obtained from some generative model. A commonly adopted model is the following  
 060

$$061 \mathbf{y} = X^T \mathbf{w}^* + \mathbf{b}, \quad (2)$$

062 where  $\mathbf{w}^* \in \mathbb{R}^p$  is the *true* model vector that we wish to estimate and  $\mathbf{b} \in \mathbb{R}^n$  is the corruption  
 063 vector that can have arbitrary values. A common assumption is that the corruption vector is *sparse*  
 064 supported i.e.  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$  for some  $\alpha > 0$ .

065 Recently, [4] and [5] obtained a surprising result which shows that one can recover  $\mathbf{w}^*$  *exactly* even  
 066 when  $\alpha \lesssim 1$ , i.e., when almost all the points are corrupted, by solving an  $L_1$ -penalty based convex  
 067 optimization problem:  $\min_{\mathbf{w}, \mathbf{b}} \|\mathbf{w}\|_1 + \lambda \|\mathbf{b}\|_1$ , s.t.,  $X^T \mathbf{w} + \mathbf{b} = \mathbf{y}$ . However, these results require  
 068 the corruption vector  $\mathbf{b}$  to be selected oblivious of  $X$  and  $\mathbf{w}^*$ . Moreover, the results impose severe  
 069 restrictions on the data distribution, requiring that the data be either sampled from an isotropic  
 070 Gaussian ensemble [4], or row-sampled from an incoherent orthogonal matrix [5]. Finally, these  
 071 results hold only for a fixed  $\mathbf{w}^*$  and are not universal in general.

072 In contrast, [6] studied RLSR with less stringent assumptions, allowing arbitrary corruptions in  
 073 response variables as well as in the data matrix  $X$ , and proposed a trimmed inner product based  
 074 algorithm for the problem. However, their recovery guarantees are significantly weaker. Firstly,  
 075 they are able to recover  $\mathbf{w}^*$  only upto an additive error  $\alpha\sqrt{p}$  (or  $\alpha\sqrt{s}$  if  $\mathbf{w}^*$  is  $s$ -sparse). Hence, they  
 076 require  $\alpha \leq 1/\sqrt{p}$  just to claim a non-trivial bound. Note that this amounts to being able to tolerate  
 077 only a vanishing fraction of corruptions. More importantly, even with  $n \rightarrow \infty$  and extremely small  
 078  $\alpha$  they are unable to guarantee exact recovery of  $\mathbf{w}^*$ . A similar result was obtained by [7], albeit  
 079 using a sub-sampling based algorithm with stronger assumptions on  $\mathbf{b}$ .

080 In this paper, we focus on a simple and natural thresholding based algorithm for RLSR. At a high  
 081 level, at each step  $t$ , our algorithm alternately estimates an *active set*  $S_t$  of “clean” points and then  
 082 updates the model to obtain  $\mathbf{w}^{t+1}$  by minimizing the least squares error on the active set. This  
 083 intuitive algorithm seems to embody a long standing heuristic first proposed by Legendre [8] over  
 084 two centuries ago (see introductory quotation in this paper) that has been adopted in later literature  
 085 [9, 10] as well. However, to the best of our knowledge, this technique has never been rigorously  
 086 analyzed before in non-asymptotic settings, despite its appealing simplicity.

087 **Our Contributions:** The main contribution of this paper is an exact recovery guarantee for the  
 088 thresholding algorithm mentioned above that we refer to as TORRENT-FC (see Algorithm 1). We  
 089 provide our guarantees in the model given in 2 where the corruptions  $\mathbf{b}$  are selected *adversarially*  
 090 but restricted to have at most  $\alpha \cdot n$  non-zero entries where  $\alpha < 1/2$  is a global constant dependent  
 091 only on  $X^1$ . Under *deterministic* conditions on  $X$ , namely the subset strong convexity (SSC) and  
 092 smoothness (SSS) properties (see Definition 1), we guarantee that TORRENT-FC converges at a  
 093 *geometric* rate and recovers  $\mathbf{w}^*$  exactly. We further show that these properties (SSC and SSS) are  
 094 satisfied w.h.p. if a) the data  $X$  is sampled from a sub-Gaussian distribution and, b)  $n \geq p \log p$ .

095 We would like to stress three key advantages of our result over the results of [4, 5]: a) we allow  $\mathbf{b}$   
 096 to be adversarial, i.e., both support and values of  $\mathbf{b}$  to be selected adversarially based on  $X$  and  $\mathbf{w}^*$ ,  
 097 b) we make assumptions on data that are natural, as well as significantly less restrictive than what  
 098 existing methods make, and c) our analysis admits universal guarantees, i.e., holds for *any*  $\mathbf{w}^*$ .

099 We would also like to stress that while hard-thresholding based methods have been studied rigor-  
 100 ously for the sparse-recovery problem [11, 12], hard-thresholding has not been studied formally  
 101 for the robust regression problem. Moreover, the two problems are completely different and hence  
 102 techniques from sparse-recovery analysis do not extend to robust regression.

103 Despite its simplicity, TORRENT-FC does not scale very well to datasets with large  $p$  as it solves  
 104 least squares problems at each iteration. We address this issue by designing a gradient descent  
 105

106 <sup>1</sup>Note that for an adaptive adversary, as is the case in our work, recovery cannot be guaranteed for  $\alpha \geq 1/2$   
 107 since the adversary can introduce corruptions as  $\mathbf{b}_i = \mathbf{x}_i^T (\tilde{\mathbf{w}} - \mathbf{w}^*)$  for an adversarially chosen model  $\tilde{\mathbf{w}}$ . This  
 would make it impossible for any algorithm to distinguish between  $\mathbf{w}^*$  and  $\tilde{\mathbf{w}}$  thus making recovery impossible.

108 based algorithm (TORRENT-GD), and a hybrid algorithm (TORRENT-Hyb), both of which enjoy a  
 109 geometric rate of convergence and can recover  $\mathbf{w}^*$  under the model assumptions mentioned above.  
 110 We also propose extensions of TORRENT for the RLSR problem in the sparse regression setting  
 111 where  $p \gg n$  but  $\|\mathbf{w}^*\|_0 = s^* \ll p$ . Our algorithm TORRENT-HD is based on TORRENT-FC but  
 112 uses the Iterative Hard Thresholding (IHT) algorithm, a popular algorithm for sparse regression. As  
 113 before, we show that TORRENT-HD also converges geometrically to  $\mathbf{w}^*$  if a) the corruption index  $\alpha$   
 114 is less than some constant  $C$ , b)  $X$  is sampled from a sub-Gaussian distribution and, c)  $n \geq s^* \log p$ .

115 Finally, we experimentally evaluate existing  $L_1$ -based algorithms and our hard thresholding-based  
 116 algorithms. The results demonstrate that our proposed algorithms (TORRENT-(FC/GD/HYB)) can  
 117 be significantly faster than the best  $L_1$  solvers, exhibit better recovery properties, as well as be more  
 118 robust to dense white noise. For instance, on a problem with  $50K$  dimensions and 40% corruption,  
 119 TORRENT-HYB was found to be  $20\times$  faster than  $L_1$  solvers, as well as achieve lower error rates.

120 **Paper Organization:** We give a formal definition of the RLSR problem in the next section. We then  
 121 introduce our family of algorithms in Section 3 and prove their convergence guarantees in Section 4.  
 122 We present extensions to sparse robust regression in Section 5 and empirical results in Section 6.

## 124 2 Problem Formulation

125  
 126 Given a set of data points  $X = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n]$ , where  $\mathbf{x}_i \in \mathbb{R}^p$  and the corresponding response  
 127 vector  $\mathbf{y} \in \mathbb{R}^n$ , the goal is to recover a parameter vector  $\mathbf{w}^*$  which solves the RLSR problem (1).  
 128 We assume that the response vector  $\mathbf{y}$  is generated using the following model:

$$129 \mathbf{y} = \mathbf{y}^* + \mathbf{b} + \boldsymbol{\varepsilon}, \text{ where } \mathbf{y}^* = X^\top \mathbf{w}^*.$$

130  
 131 Hence, in the above model, (1) reduces to estimating  $\mathbf{w}^*$ . We allow the model  $\mathbf{w}^*$  representing the  
 132 regressor, to be chosen in an adaptive manner *after* the data features have been generated.

133 The above model allows two kinds of perturbations to  $y_i$  – dense but bounded noise  $\varepsilon_i$  (e.g. white  
 134 noise  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ ,  $\sigma \geq 0$ ), as well as potentially unbounded corruptions  $b_i$  – to be introduced  
 135 by an adversary. The only requirement we enforce is that the gross corruptions be sparse.  $\boldsymbol{\varepsilon}$  shall  
 136 represent the dense noise vector, for example  $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \cdot I_{n \times n})$ , and  $\mathbf{b}$ , the corruption vector such  
 137 that  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$  for some *corruption index*  $\alpha > 0$ . We shall use the notation  $S_* = \text{supp}(\mathbf{b}) \subseteq [n]$   
 138 to denote the set of “clean” points, i.e. points that have not faced unbounded corruptions. We consider  
 139 adaptive adversaries that are able to view the generated data points  $\mathbf{x}_i$ , as well as the clean responses  
 140  $y_i^*$  and dense noise values  $\varepsilon_i$  before deciding which locations to corrupt and by what amount.

141 We denote the unit sphere in  $p$  dimensions using  $S^{p-1}$ . For any  $\gamma \in (0, 1]$ , we let  $\mathcal{S}_\gamma =$   
 142  $\{S \subset [n] : |S| = \gamma \cdot n\}$  denote the set of all subsets of size  $\gamma \cdot n$ . For any set  $S$ , we let  $X_S :=$   
 143  $[\mathbf{x}_i]_{i \in S} \in \mathbb{R}^{p \times |S|}$  denote the matrix whose columns are composed of points in that set. Also, for  
 144 any vector  $\mathbf{v} \in \mathbb{R}^n$  we use the notation  $\mathbf{v}_S$  to denote the  $|S|$ -dimensional vector consisting of those  
 145 components that are in  $S$ . We use  $\lambda_{\min}(X)$  and  $\lambda_{\max}(X)$  to denote, respectively, the smallest and  
 146 largest eigenvalues of a square symmetric matrix  $X$ . We now introduce two properties, namely,  
 147 *Subset Strong Convexity* and *Subset Strong Smoothness*, which are key to our analyses.

148 **Definition 1** (SSC and SSS Properties). *A matrix  $X \in \mathbb{R}^{p \times n}$  satisfies the Subset Strong Convexity*  
 149 *Property (resp. Subset Strong Smoothness Property) at level  $\gamma$  with strong convexity constant  $\lambda_\gamma$*   
 150 *(resp. strong smoothness constant  $\Lambda_\gamma$ ) if the following holds:*

$$151 \lambda_\gamma \leq \min_{S \in \mathcal{S}_\gamma} \lambda_{\min}(X_S X_S^\top) \leq \max_{S \in \mathcal{S}_\gamma} \lambda_{\max}(X_S X_S^\top) \leq \Lambda_\gamma.$$

152  
 153 *Remark 1.* We note that the uniformity enforced in the definitions of the SSC and SSS properties is  
 154 not for the sake of convenience but rather a necessity. Indeed, a uniform bound is required in face of  
 155 an adversary which can perform corruptions *after* data and response variables have been generated,  
 156 and choose to corrupt precisely that set of points where the SSC and SSS parameters are the worst.

## 158 3 TORRENT: Thresholding Operator-based Robust Regression Method

159  
 160 We now present TORRENT, a Thresholding Operator-based Robust RegrEsson meThod for per-  
 161 forming robust regression at scale. Key to our algorithms is the *Hard Thresholding Operator* which  
 we define below.

162  
163  
164  
165  
166  
167  
168  
169  
170  
171  
172  
173  
174  
175  
176  
177  
178  
179  
180  
181  
182  
183  
184  
185  
186  
187  
188  
189  
190  
191  
192  
193  
194  
195  
196  
197  
198  
199  
200  
201  
202  
203  
204  
205  
206  
207  
208  
209  
210  
211  
212  
213  
214  
215

---

**Algorithm 1** TORRENT: Thresholding Operator-based Robust RegrEsson meThod

---

**Input:** Training data  $\{\mathbf{x}_i, y_i\}, i = 1 \dots n$ , step length  $\eta$ , thresholding parameter  $\beta$ , tolerance  $\epsilon$

- 1:  $\mathbf{w}^0 \leftarrow \mathbf{0}, S_0 = [n], t \leftarrow 0, \mathbf{r}^0 \leftarrow \mathbf{y}$
- 2: **while**  $\|\mathbf{r}_{S_t}^t\|_2 > \epsilon$  **do**
- 3:    $\mathbf{w}^{t+1} \leftarrow \text{UPDATE}(\mathbf{w}^t, S_t, \eta, \mathbf{r}^t, S_{t-1})$
- 4:    $r_i^{t+1} \leftarrow (y_i - \langle \mathbf{w}^{t+1}, \mathbf{x}_i \rangle)$
- 5:    $S_{t+1} \leftarrow \text{HT}(\mathbf{r}^{t+1}, (1 - \beta)n)$
- 6:    $t \leftarrow t + 1$
- 7: **end while**
- 8: **return**  $\mathbf{w}^t$

---

**Algorithm 2** TORRENT-FC

---

**Input:** Current model  $\mathbf{w}$ , current active set  $S$

- 1: **return**  $\arg \min_{\mathbf{w}} \sum_{i \in S} (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2$

---



---

**Algorithm 3** TORRENT-GD

---

**Input:** Current model  $\mathbf{w}$ , current active set  $S$ , step size  $\eta$

- 1:  $\mathbf{g} \leftarrow X_S(X_S^\top \mathbf{w} - \mathbf{y}_S)$
- 2: **return**  $\mathbf{w} - \eta \cdot \mathbf{g}$

---



---

**Algorithm 4** TORRENT-HYB

---

**Input:** Current model  $\mathbf{w}$ , current active set  $S$ , step size  $\eta$ , current residuals  $\mathbf{r}$ , previous active set  $S'$

- 1: // Use the GD update if the active set  $S$  is changing a lot
- 2: **if**  $|S \setminus S'| > \Delta$  **then**
- 3:    $\mathbf{w}' \leftarrow \text{UPDATE-GD}(\mathbf{w}, S, \eta, \mathbf{r}, S')$
- 4: **else**
- 5: // If stable, use the FC update
- 6:    $\mathbf{w}' \leftarrow \text{UPDATE-FC}(\mathbf{w}, S)$
- 7: **end if**
- 8: **return**  $\mathbf{w}'$

---

**Definition 2** (Hard Thresholding Operator). For any vector  $\mathbf{v} \in \mathbb{R}^n$ , let  $\sigma_{\mathbf{v}} \in S_n$  be the permutation that orders elements of  $\mathbf{v}$  in ascending order of their magnitudes i.e.  $|\mathbf{v}_{\sigma_{\mathbf{v}}(1)}| \leq |\mathbf{v}_{\sigma_{\mathbf{v}}(2)}| \leq \dots \leq |\mathbf{v}_{\sigma_{\mathbf{v}}(n)}|$ . Then for any  $k \leq n$ , we define the hard thresholding operator as

$$\text{HT}(\mathbf{v}; k) = \{i \in [n] : \sigma_{\mathbf{v}}^{-1}(i) \leq k\}$$

Using this operator, we present our algorithm TORRENT (Algorithm 1) for robust regression. TORRENT follows a most natural iterative strategy of, alternately, estimating an *active set* of points which have the least residual error on the current regressor, and then updating the regressor to provide a better fit on this active set. We offer three variants of our algorithm, based on how aggressively the algorithm tries to fit the regressor to the current active set.

We first propose a fully corrective algorithm TORRENT-FC (Algorithm 2) that performs a fully corrective least squares regression step in an effort to minimize the regression error on the active set. This algorithm makes significant progress in each step, but at a cost of more expensive updates. To address this, we then propose a milder, gradient descent-based variant TORRENT-GD (Algorithm 3) that performs a much cheaper update of taking a single step in the direction of the gradient of the objective function on the active set. This reduces the regression error on the active set but does not minimize it. This turns out to be beneficial in situations where dense noise is present along with sparse corruptions since it prevents the algorithm from overfitting to the current active set.

Both the algorithms proposed above have their pros and cons – the FC algorithm provides significant improvements with each step, but is expensive to execute whereas the GD variant, although efficient in executing each step, offers slower progress. To get the best of both these algorithms, we propose a third, hybrid variant TORRENT-HYB (Algorithm 4) that adaptively selects either the FC or the GD update depending on whether the active set is stable across iterations or not.

In the next section we show that this hard thresholding-based strategy offers a linear convergence rate for the algorithm in all its three variations. We shall also demonstrate the applicability of this technique to high dimensional sparse recovery settings in a subsequent section.

## 4 Convergence Guarantees

For the sake of ease of exposition, we will first present our convergence analyses for cases where dense noise is not present i.e.  $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$  and will handle cases with dense noise *and* sparse corruptions later. We first analyze the fully corrective TORRENT-FC algorithm. The convergence proof in this case relies on the optimality of the two steps carried out by the algorithm, the fully corrective step that selects the best regressor on the active set, and the hard thresholding step that discovers a new active set by selecting points with the least residual error on the current regressor.

**Theorem 3.** Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix and  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b}$  be the corrupted output with  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Let Algorithm 2 be executed on this data with the thresholding parameter set to  $\beta \geq \alpha$ . Let  $\Sigma_0$  be an invertible matrix such that  $\tilde{X} = \Sigma_0^{-1/2} X$  satisfies the SSC and SSS properties at level  $\gamma$  with constants  $\lambda_\gamma$  and  $\Lambda_\gamma$  respectively (see Definition 1). If the data satisfies  $\frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < 1$ , then after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 2 obtains an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$ .

*Proof (Sketch).* Let  $\mathbf{r}^t = \mathbf{y} - X^T \mathbf{w}^t$  be the vector of residuals at time  $t$  and  $C_t = X_{S_t} X_{S_t}^\top$ . Also let  $S_* = \text{supp}(\mathbf{b})$  be the set of uncorrupted points. The fully corrective step ensures that

$$\mathbf{w}^{t+1} = C_t^{-1} X_{S_t} \mathbf{y}_{S_t} = C_t^{-1} X_{S_t} (X_{S_t}^\top \mathbf{w}^* + \mathbf{b}_{S_t}) = \mathbf{w}^* + C_t^{-1} X_{S_t} \mathbf{b}_{S_t},$$

whereas the hard thresholding step ensures that  $\|\mathbf{r}_{S_{t+1}}^{t+1}\|_2^2 \leq \|\mathbf{r}_{S_*}^{t+1}\|_2^2$ . Combining the two gives us

$$\begin{aligned} \|\mathbf{b}_{S_{t+1}}\|_2^2 &\leq \left\| X_{S_* \setminus S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \\ &\stackrel{\zeta_1}{\leq} \left\| \tilde{X}_{S_* \setminus S_{t+1}}^\top \left( \tilde{X}_{S_t} \tilde{X}_{S_t}^\top \right)^{-1} \tilde{X}_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top \tilde{X}_{S_{t+1}}^\top \left( \tilde{X}_{S_t} \tilde{X}_{S_t}^\top \right)^{-1} \tilde{X}_{S_t} \mathbf{b}_{S_t} \\ &\stackrel{\zeta_2}{\leq} \frac{\Lambda_\beta^2}{\lambda_{1-\beta}^2} \cdot \|\mathbf{b}_{S_t}\|_2^2 + 2 \cdot \frac{\Lambda_\beta}{\lambda_{1-\beta}} \cdot \|\mathbf{b}_{S_t}\|_2 \|\mathbf{b}_{S_{t+1}}\|_2, \end{aligned}$$

where  $\zeta_1$  follows from setting  $\tilde{X} = \Sigma_0^{-1/2} X$  and  $X_{S'}^\top C_t^{-1} X_{S'} = \tilde{X}_{S'}^\top (\tilde{X}_{S_t} \tilde{X}_{S_t}^\top)^{-1} \tilde{X}_{S'}$  and  $\zeta_2$  follows from the SSC and SSS properties,  $\|\mathbf{b}_{S_t}\|_0 \leq \|\mathbf{b}\|_0 \leq \beta \cdot n$  and  $|S_* \setminus S_{t+1}| \leq \beta \cdot n$ . Solving the quadratic equation and performing other manipulations gives us the claimed result.  $\square$

Theorem 3 relies on a deterministic (*fixed design*) assumption, specifically  $\frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < 1$  in order to guarantee convergence. We can show that a large class of random designs, including Gaussian and sub-Gaussian designs actually satisfy this requirement. That is to say, data generated from these distributions satisfy the SSC and SSS conditions such that  $\frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < 1$  with high probability. Theorem 4 explicates this for the class of Gaussian designs.

**Theorem 4.** Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix with each  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$ . Let  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b}$  and  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Also, let  $\alpha \leq \beta < \frac{1}{65}$  and  $n \geq \Omega\left(p + \log \frac{1}{\delta}\right)$ . Then, with probability at least  $1 - \delta$ , the data satisfies  $\frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < \frac{9}{10}$ . More specifically, after  $T \geq 10 \log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)$  iterations of Algorithm 1 with the thresholding parameter set to  $\beta$ , we have  $\|\mathbf{w}^T - \mathbf{w}^*\| \leq \epsilon$ .

*Remark 2.* Note that Theorem 4 provides rates that are independent of the condition number  $\frac{\lambda_{\max}(\Sigma)}{\lambda_{\min}(\Sigma)}$  of the distribution. We also note that results similar to Theorem 4 can be proven for the larger class of sub-Gaussian distributions. We refer the reader to Section G for the same.

*Remark 3.* We remind the reader that our analyses can readily accommodate dense noise in addition to sparse unbounded corruptions. We direct the reader to Appendix A which presents convergence proofs for our algorithms in these settings.

*Remark 4.* We would like to point out that the design requirements made by our analyses are very mild when compared to existing literature. Indeed, the work of [4] assumes the *Bouquet Model* where distributions are restricted to be isotropic Gaussians whereas the work of [5] assumes a more stringent model of sub-orthonormal matrices, something that even Gaussian designs do not satisfy. Our analyses, on the other hand, hold for the general class of sub-Gaussian distributions.

We now analyze the TORRENT-GD algorithm which performs cheaper, gradient-style updates on the active set. We will show that this method nevertheless enjoys a linear rate of convergence.

**Theorem 5.** Let the data settings be as stated in Theorem 3 and let Algorithm 3 be executed on this data with the thresholding parameter set to  $\beta \geq \alpha$  and the step length set to  $\eta = \frac{1}{\Lambda_{1-\beta}}$ . If the data

satisfies  $\max\{\eta\sqrt{\Lambda_\beta}, 1 - \eta\lambda_{1-\beta}\} \leq \frac{1}{4}$ , then after  $t = \mathcal{O}\left(\log\left(\frac{\|\mathbf{b}\|_2}{\sqrt{n}\epsilon}\right)\right)$  iterations, Algorithm 1 obtains an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$ .

Similar to TORRENT-FC, the assumptions made by the TORRENT-GD algorithm are also satisfied by the class of sub-Gaussian distributions. The proof of Theorem 5, given in Appendix D, details these arguments. Given the convergence analyses for TORRENT-FC and GD, we now move on to provide a convergence analysis for the hybrid TORRENT-HYB algorithm which interleaves FC and GD steps. Since the exact interleaving adopted by the algorithm depends on the data, and not known in advance, this poses a problem. We address this problem by giving below a uniform convergence guarantee, one that applies to every interleaving of the FC and GD update steps.

**Theorem 6.** *Suppose Algorithm 4 is executed on data that allows Algorithms 2 and 3 a convergence rate of  $\eta_{\text{FC}}$  and  $\eta_{\text{GD}}$  respectively. Suppose we have  $2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} < 1$ . Then for any interleavings of the FC and GD steps that the policy may enforce, after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 4 ensures an  $\epsilon$ -optimal solution i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\| \leq \epsilon$ .*

We point out to the reader that the assumption made by Theorem 6 i.e.  $2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} < 1$  is readily satisfied by random sub-Gaussian designs, albeit at the cost of reducing the noise tolerance limit. As we shall see, TORRENT-HYB offers attractive convergence properties, merging the fast convergence rates of the FC step, as well as the speed and protection against overfitting provided by the GD step.

## 5 High-dimensional Robust Regression

In this section, we extend our approach to the robust high-dimensional sparse recovery setting. As before, we assume that the response vector  $\mathbf{y}$  is obtained as:  $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$ , where  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . However, this time, we also assume that  $\mathbf{w}^*$  is  $s^*$ -sparse i.e.  $\|\mathbf{w}^*\|_0 \leq s^*$ . As before, we shall neglect white/dense noise for the sake of simplicity. We reiterate that it is not possible to use existing results from sparse recovery (such as [11, 12]) directly to solve this problem.

Our objective would be to recover a sparse model  $\hat{\mathbf{w}}$  so that  $\|\hat{\mathbf{w}} - \mathbf{w}^*\|_2 \leq \epsilon$ . The challenge here is to forgo a sample complexity of  $n \gtrsim p$  and instead, perform recovery with  $n \sim s^* \log p$  samples alone. For this setting, we modify the FC update step of TORRENT-FC method to the following:

$$\mathbf{w}^{t+1} \leftarrow \inf_{\|\mathbf{w}\|_0 \leq s} \sum_{i \in S_t} (y_i - \langle \mathbf{w}, \mathbf{x}_i \rangle)^2, \quad (3)$$

for some target sparsity level  $s \ll p$ . We refer to this modified algorithm as TORRENT-HD. Assuming  $X$  satisfies the RSC/RSS properties (defined below), (3) can be solved efficiently using results from sparse recovery (for example the IHT algorithm [11, 13] analyzed in [12]).

**Definition 7** (RSC and RSS Properties). *A matrix  $X \in \mathbb{R}^{p \times n}$  will be said to satisfy the Restricted Strong Convexity Property (resp. Restricted Strong Smoothness Property) at level  $s = s_1 + s_2$  with strong convexity constant  $\alpha_{s_1+s_2}$  (resp. strong smoothness constant  $L_{s_1+s_2}$ ) if the following holds for all  $\|\mathbf{w}_1\|_0 \leq s_1$  and  $\|\mathbf{w}_2\|_0 \leq s_2$ :*

$$\alpha_s \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2 \leq \|X^\top (\mathbf{w}_1 - \mathbf{w}_2)\|_2^2 \leq L_s \|\mathbf{w}_1 - \mathbf{w}_2\|_2^2$$

For our results, we shall require the subset versions of both these properties.

**Definition 8** (SRSC and SRSS Properties). *A matrix  $X \in \mathbb{R}^{p \times n}$  will be said to satisfy the Subset Restricted Strong Convexity (resp. Subset Restricted Strong Smoothness) Property at level  $(\gamma, s)$  with strong convexity constant  $\alpha_{(\gamma, s)}$  (resp. strong smoothness constant  $L_{(\gamma, s)}$ ) if for all subsets  $S \in \mathcal{S}_\gamma$ , the matrix  $X_S$  satisfies the RSC (resp. RSS) property at level  $s$  with constant  $\alpha_s$  (resp.  $L_s$ ).*

We now state the convergence result for the TORRENT-HD algorithm.

**Theorem 9.** *Let  $X \in \mathbb{R}^{p \times n}$  be the given data matrix and  $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$  be the corrupted output with  $\|\mathbf{w}^*\|_0 \leq s^*$  and  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Let  $\Sigma_0$  be an invertible matrix such that  $\Sigma_0^{-1/2} X$  satisfies the SRSC and SRSS properties at level  $(\gamma, 2s + s^*)$  with constants  $\alpha_{(\gamma, 2s + s^*)}$  and  $L_{(\gamma, 2s + s^*)}$  respectively (see Definition 8). Let Algorithm 2 be executed on this data with the TORRENT-HD update, thresholding parameter set to  $\beta \geq \alpha$ , and  $s \geq 32 \left(\frac{L_{(1-\beta, 2s + s^*)}}{\alpha_{(1-\beta, 2s + s^*)}}\right)$ .*

324  
325  
326  
327  
328  
329  
330  
331  
332  
333  
334  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377

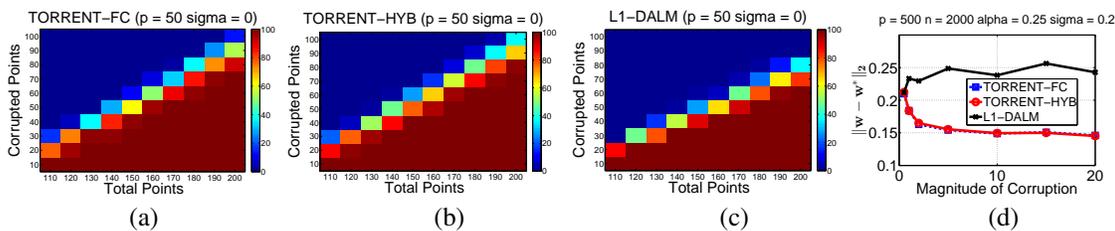


Figure 1: (a), (b) and (c) Phase-transition diagrams depicting the recovery properties of the TORRENT-FC, TORRENT-HYB and  $L_1$  algorithms. The colors red and blue represent a high and low probability of success resp. A method is considered successful in an experiment if it recovers  $w^*$  upto a  $10^{-4}$  relative error. Both variants of TORRENT can be seen to recover  $w^*$  in presence of larger number of corruptions than the  $L_1$  solver. (d) Variation in recovery error with the magnitude of corruption. As the corruption is increased, TORRENT-FC and TORRENT-HYB show improved performance while the problem becomes more difficult for the  $L_1$  solver.

If  $X$  also satisfies  $\frac{4L_{(\beta, s+s^*)}}{\alpha_{(1-\beta, s+s^*)}} < 1$ , then after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|b\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 2 obtains an  $\epsilon$ -accurate solution  $w^t$  i.e.  $\|w^t - w^*\|_2 \leq \epsilon$ .

In particular, if  $X$  is sampled from a Gaussian distribution  $\mathcal{N}(0, \Sigma)$  and  $n \geq \Omega\left(s^* \cdot \frac{\lambda_{\max}(\Sigma)}{\lambda_{\min}(\Sigma)} \log p\right)$ , then for all values of  $\alpha \leq \beta < \frac{1}{65}$ , we can guarantee  $\|w^t - w^*\|_2 \leq \epsilon$  after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|b\|_2}{\epsilon}\right)\right)$  iterations of the algorithm (w.p.  $\geq 1 - 1/n^{10}$ ).

*Remark 5.* The sample complexity required by Theorem 9 is identical to the one required by analyses for high dimensional sparse recovery [12], save constants. Also note that TORRENT-HD can tolerate the same corruption index as TORRENT-FC.

## 6 Experiments

Several numerical simulations were carried out on linear regression problems in low-dimensional, as well as sparse high-dimensional settings. The experiments show that TORRENT not only offers statistically better recovery properties as compared to  $L_1$ -style approaches, but that it can be more than an order of magnitude faster as well.

**Data:** For the low dimensional setting, the regressor  $w^* \in \mathbb{R}^p$  was chosen to be a random unit norm vector. Data was sampled as  $x_i \sim \mathcal{N}(0, I_p)$  and response variables were generated as  $y_i^* = \langle w^*, x_i \rangle$ . The set of corrupted points  $\bar{S}_*$  was selected as a uniformly random  $(\alpha n)$ -sized subset of  $[n]$  and the corruptions were set to  $b_i \sim U(-5 \|y^*\|_\infty, 5 \|y^*\|_\infty)$  for  $i \in \bar{S}_*$ . The corrupted responses were then generated as  $y_i = y_i^* + b_i + \varepsilon_i$  where  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ . For the sparse high-dimensional setting,  $\text{supp}(w^*)$  was selected to be a random  $s^*$ -sized subset of  $[p]$ . Phase-transition diagrams (Figure 1) were generated by repeating each experiment 100 times. For all other plots, each experiment was run over 20 random instances of the data and the plots were drawn to depict the mean results.

**Algorithms:** We compared various variants of our algorithm TORRENT to the regularized  $L_1$  algorithm for robust regression [4, 5]. Note that the  $L_1$  problem can be written as  $\min_z \|z\|_1$  s.t.  $Az = y$ , where  $A = [X^T \frac{1}{\lambda} I_{m \times m}]$  and  $z^* = [w^{*T} \lambda b^T]^T$ . We used the Dual Augmented Lagrange Multiplier (DALM)  $L_1$  solver implemented by [14] to solve the  $L_1$  problem. We ran a fine tuned grid search over the  $\lambda$  parameter for the  $L_1$  solver and quoted the best results obtained from the search. In the low-dimensional setting, we compared the recovery properties of TORRENT-FC (Algorithm 2) and TORRENT-HYB (Algorithm 4) with the DALM- $L_1$  solver, while for the high-dimensional case, we compared TORRENT-HD against the DALM- $L_1$  solver. Both the  $L_1$  solver, as well as our methods, were implemented in Matlab and were run on a single core 2.4GHz machine with 8 GB RAM.

**Choice of  $L_1$ -solver:** An extensive comparative study of various  $L_1$  minimization algorithms was performed by [14] who showed that the DALM and Homotopy solvers outperform other counterparts both in terms of recovery properties, and timings. We extended their study to our observation model and found the DALM solver to be significantly better than the other  $L_1$  solvers; see Figure 3 in the appendix. We also observed, similar to [14], that the Approximate Message Passing (AMP) solver diverges on our problem as the input matrix to the  $L_1$  solver is a non-Gaussian matrix  $A = [X^T \frac{1}{\lambda} I]$ .

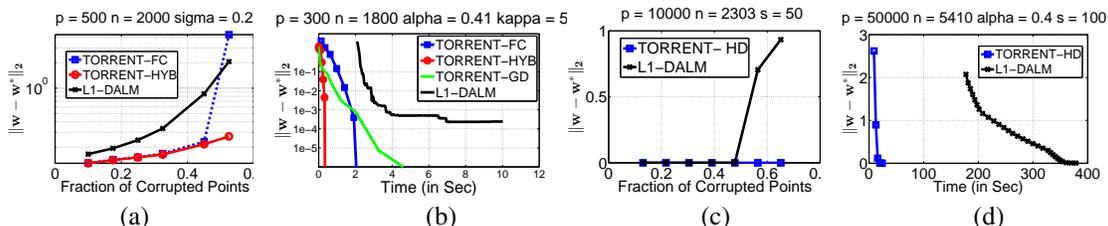


Figure 2: In low-dimensional (a,b), as well as sparse high dimensional (c,d) settings, TORRENT offers better recovery as the fraction of corrupted points  $\alpha$  is varied. In terms of runtime, TORRENT is an order of magnitude faster than  $L_1$  solvers in both settings. In the low-dim. setting, TORRENT-HYB is the fastest of all the variants.

**Evaluation Metric:** We measure the performance of various algorithms using the standard  $L_2$  error:  $r_{\hat{\mathbf{w}}} = \|\hat{\mathbf{w}} - \mathbf{w}^*\|_2$ . For the phase-transition plots (Figure 1), we deemed an algorithm successful on an instance if it obtained a model  $\hat{\mathbf{w}}$  with error  $r_{\hat{\mathbf{w}}} < 10^{-4} \cdot \|\mathbf{w}^*\|_2$ . We also measured the CPU time required by each of the methods, so as to compare their scalability.

## 6.1 Low Dimensional Results

**Recovery Property:** The phase-transition plots presented in Figure 1 represent our recovery experiments in graphical form. Both the fully-corrective and hybrid variants of TORRENT show better recovery properties than the  $L_1$ -minimization approach, indicated by the number of runs in which the algorithm was able to correctly recover  $\mathbf{w}^*$  out of a 100 runs. Figure 2 shows the variation in recovery error as a function of  $\alpha$  in the presence of white noise and exhibits the superiority of TORRENT-FC and TORRENT-HYB over  $L_1$ -DALM. Here again, TORRENT-FC and TORRENT-HYB achieve significantly lesser recovery error than  $L_1$ -DALM for all  $\alpha \leq 0.5$ . Figure 3 in the appendix show that the variations of  $\|\hat{\mathbf{w}} - \mathbf{w}^*\|_2$  with varying  $p, \sigma$  and  $n$  follow a similar trend with TORRENT having significantly lower recovery error in comparison to the  $L_1$  approach.

Figure 1(d) brings out an interesting trend in the recovery property of TORRENT. As we increase the magnitude of corruption from  $U(-\|\mathbf{y}^*\|_\infty, \|\mathbf{y}^*\|_\infty)$  to  $U(-20\|\mathbf{y}^*\|_\infty, 20\|\mathbf{y}^*\|_\infty)$ , the recovery error for TORRENT-HYB and TORRENT-FC decreases as expected since it becomes easier to identify the grossly corrupted points. However the  $L_1$ -solver was unable to exploit this observation and in fact exhibited an increase in recovery error.

**Run Time:** In order to ascertain the recovery guarantees for TORRENT on ill-conditioned problems, we performed an experiment where data was sampled as  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$  where  $\text{diag}(\Sigma) \sim U(0, 5)$ . Figure 2 plots the recovery error as a function of time. TORRENT-HYB was able to correctly recover  $\mathbf{w}^*$  about  $50\times$  faster than  $L_1$ -DALM which spent a considerable amount of time pre-processing the data matrix  $X$ . Even after allowing the  $L_1$  algorithm to run for 500 iterations, it was unable to reach the desired residual error of  $10^{-4}$ . Figure 2 also shows that our TORRENT-HYB algorithm is able to converge to the optimal solution much faster than TORRENT-FC or TORRENT-GD. This is because TORRENT-FC solves a least square problem at each step and thus, even though it requires significantly fewer iterations to converge, each iteration in itself is very expensive. While each iteration of TORRENT-GD is cheap, it is still limited by the slow  $\mathcal{O}\left(\left(1 - \frac{1}{\kappa}\right)^t\right)$  convergence rate of the gradient descent algorithm, where  $\kappa$  is the condition number of the covariance matrix. TORRENT-HYB, on the other hand, is able to combine the strengths of both the methods to achieve faster convergence.

## 6.2 High Dimensional Results

**Recovery Property:** Figure 2 shows the variation in recovery error in the high-dimensional setting as the number of corrupted points was varied. For these experiments,  $n$  was set to  $5s^* \log(p)$  and the fraction of corrupted points  $\alpha$  was varied from 0.1 to 0.7. While  $L_1$ -DALM fails to recover  $\mathbf{w}^*$  for  $\alpha > 0.5$ , TORRENT-HD offers perfect recovery even for  $\alpha$  values upto 0.7.

**Run Time:** Figure 2 shows the variation in recovery error as a function of run time in this setting.  $L_1$ -DALM was found to be an order of magnitude slower than TORRENT-HD, making it infeasible for sparse high-dimensional settings. One key reason for this is that the  $L_1$ -DALM solver is significantly slower in identifying the set of clean points. For instance, whereas TORRENT-HD was able to identify the clean set of points in only 5 iterations, it took  $L_1$  around 250 iterations to do the same.

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485

## References

- [1] Christoph Studer, Patrick Kuppinger, Graeme Pope, and Helmut Bölcskei. Recovery of Sparsely Corrupted Signals. *IEEE Transaction on Information Theory*, 58(5):3115–3130, 2012.
- [2] Peter J. Rousseeuw and Annick M. Leroy. *Robust Regression and Outlier Detection*. John Wiley and Sons, 1987.
- [3] John Wright, Alan Y. Yang, Arvind Ganesh, S. Shankar Sastry, and Yi Ma. Robust Face Recognition via Sparse Representation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 31(2):210–227, 2009.
- [4] John Wright and Yi Ma. Dense Error Correction via  $\ell^1$  Minimization. *IEEE Transaction on Information Theory*, 56(7):3540–3560, 2010.
- [5] Nam H. Nguyen and Trac D. Tran. Exact recoverability from dense corrupted observations via L1 minimization. *IEEE Transaction on Information Theory*, 59(4):2036–2058, 2013.
- [6] Yudong Chen, Constantine Caramanis, and Shie Mannor. Robust Sparse Regression under Adversarial Corruption. In *30th International Conference on Machine Learning (ICML)*, 2013.
- [7] Brian McWilliams, Gabriel Krummenacher, Mario Lucic, and Joachim M. Buhmann. Fast and Robust Least Squares Estimation in Corrupted Linear Models. In *28th Annual Conference on Neural Information Processing Systems (NIPS)*, 2014.
- [8] Adrien-Marie Legendre (1805). On the Method of Least Squares. In (Translated from the French) D.E. Smith, editor, *A Source Book in Mathematics*, pages 576–579. New York: Dover Publications, 1959.
- [9] Peter J. Rousseeuw. Least Median of Squares Regression. *Journal of the American Statistical Association*, 79(388):871–880, 1984.
- [10] Peter J. Rousseeuw and Katrien Driessen. Computing LTS Regression for Large Data Sets. *Journal of Data Mining and Knowledge Discovery*, 12(1):29–45, 2006.
- [11] Thomas Blumensath and Mike E. Davies. Iterative Hard Thresholding for Compressed Sensing. *Applied and Computational Harmonic Analysis*, 27(3):265–274, 2009.
- [12] Prateek Jain, Ambuj Tewari, and Purushottam Kar. On Iterative Hard Thresholding Methods for High-dimensional M-Estimation. In *28th Annual Conference on Neural Information Processing Systems (NIPS)*, 2014.
- [13] Rahul Garg and Rohit Khandekar. Gradient descent with sparsification: an iterative algorithm for sparse recovery with restricted isometry property. In *26th International Conference on Machine Learning (ICML)*, 2009.
- [14] Allen Y. Yang, Arvind Ganesh, Zihan Zhou, Shankar Sastry, and Yi Ma. A Review of Fast  $\ell_1$ -Minimization Algorithms for Robust Face Recognition. CoRR abs/1007.3753, 2012.
- [15] Beatrice Laurent and Pascal Massart. Adaptive estimation of a quadratic functional by model selection. *The Annals of Statistics*, 28(5):1302–1338, 2000.
- [16] Thomas Blumensath. Sampling and reconstructing signals from a union of linear subspaces. *IEEE Transactions on Information Theory*, 57(7):4660–4671, 2011.
- [17] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. In Y. Eldar and G. Kutyniok, editors, *Compressed Sensing, Theory and Applications*, chapter 5, pages 210–268. Cambridge University Press, 2012.

## A Convergence Guarantees with Dense Noise and Sparse Corruptions

We will now present recovery guarantees for the TORRENT-FC algorithm when both, dense noise, as well as sparse adversarial corruptions are present. Extensions for TORRENT-GD and TORRENT-HYB will follow similarly.

**Theorem 10.** *Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix and  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b} + \boldsymbol{\varepsilon}$  be the corrupted output with sparse corruptions  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$  as well as dense bounded noise  $\boldsymbol{\varepsilon}$ . Let Algorithm 2 be executed on this data with the thresholding parameter set to  $\beta \geq \alpha$ . Let  $\Sigma_0$  be an invertible matrix such that  $\tilde{X} = \Sigma_0^{-1/2} X$  satisfies the SSC and SSS properties at level  $\gamma$  with constants  $\lambda_\gamma$  and  $\Lambda_\gamma$  respectively (see Definition 1). If the data satisfies  $\frac{4\sqrt{\Lambda_\beta}}{\sqrt{\lambda_{1-\beta}}} < 1$ , then after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 2 obtains an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon + C \frac{\|\boldsymbol{\varepsilon}\|_2}{\sqrt{n}}$  for some constant  $C > 0$ .*

*Proof.* We begin by observing that the optimality of the model  $\mathbf{w}^{t+1}$  on the active set  $S_t$  ensures

$$\|\mathbf{y}_{S_t} - X_{S_t}^\top \mathbf{w}^{t+1}\|_2 = \|X_{S_t}^\top (\mathbf{w}^* - \mathbf{w}^{t+1}) + \boldsymbol{\varepsilon}_{S_t} + \mathbf{b}_{S_t}\|_2 \leq \|\mathbf{y}_t - X_{S_t}^\top \mathbf{w}^*\|_2 = \|\boldsymbol{\varepsilon}_{S_t} + \mathbf{b}_{S_t}\|_2,$$

which, upon the application of the triangle inequality, gives us

$$\|X_{S_t}^\top (\mathbf{w}^* - \mathbf{w}^{t+1})\|_2 \leq 2 \|\boldsymbol{\varepsilon}_{S_t} + \mathbf{b}_{S_t}\|_2.$$

Since  $\|X_{S_t}^\top (\mathbf{w}^* - \mathbf{w}^{t+1})\|_2 \geq \sqrt{\lambda_{1-\beta}} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2$ , we get

$$\|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 \leq \frac{2}{\sqrt{\lambda_{1-\beta}}} \|\boldsymbol{\varepsilon}_{S_t} + \mathbf{b}_{S_t}\|_2 \leq \frac{2}{\sqrt{\lambda_{1-\beta}}} (\|\boldsymbol{\varepsilon}\|_2 + \|\mathbf{b}_{S_t}\|_2).$$

The hard thresholding step, on the other hand, guarantees that

$$\begin{aligned} \|X_{S_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1}) + \boldsymbol{\varepsilon}_{S_{t+1}} + \mathbf{b}_{S_{t+1}}\|_2^2 &= \|\mathbf{y}_{S_{t+1}} - X_{S_{t+1}}^\top \mathbf{w}^{t+1}\|_2^2 \\ &\leq \|\mathbf{y}_{S_*} - X_{S_*}^\top \mathbf{w}^{t+1}\|_2^2 \\ &= \|X_{S_*}^\top (\mathbf{w}^* - \mathbf{w}^{t+1}) + \boldsymbol{\varepsilon}_{S_*}\|_2^2. \end{aligned}$$

As before, let  $\text{CR}_{t+1} = S_{t+1} \setminus S_*$  and  $\text{MD}_{t+1} = S_* \setminus S_{t+1}$ . Then we have

$$\|X_{\text{CR}_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1}) + \boldsymbol{\varepsilon}_{\text{CR}_{t+1}} + \mathbf{b}_{\text{CR}_{t+1}}\|_2 \leq \|X_{\text{MD}_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1}) + \boldsymbol{\varepsilon}_{\text{MD}_{t+1}}\|_2.$$

An application of the triangle inequality and the fact that  $\|\mathbf{b}_{\text{CR}_{t+1}}\|_2 = \|\mathbf{b}_{S_{t+1}}\|_2$  gives us

$$\begin{aligned} \|\mathbf{b}_{S_{t+1}}\|_2 &\leq \|X_{\text{MD}_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1})\|_2 + \|X_{\text{CR}_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1})\|_2 + \|\boldsymbol{\varepsilon}_{\text{CR}_{t+1}}\|_2 + \|\boldsymbol{\varepsilon}_{\text{MD}_{t+1}}\|_2 \\ &\leq 2\sqrt{\Lambda_\beta} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 + \sqrt{2} \|\boldsymbol{\varepsilon}\|_2, \\ &= \frac{4\sqrt{\Lambda_\beta}}{\sqrt{\lambda_{1-\beta}}} \|\mathbf{b}_{S_t}\|_2 + \left(\frac{4\sqrt{\Lambda_\beta}}{\sqrt{\lambda_{1-\beta}}} + \sqrt{2}\right) \|\boldsymbol{\varepsilon}\|_2 \\ &\leq \eta \cdot \|\mathbf{b}_{S_t}\|_2 + (1 + \sqrt{2}) \|\boldsymbol{\varepsilon}\|_2, \end{aligned}$$

where the second step uses the fact that  $\max\{|\text{CR}_{t+1}|, |\text{MD}_{t+1}|\} \leq \beta \cdot n$  and the Cauchy-Schwartz inequality, and the last step uses the fact that for sufficiently small  $\beta$ , we have  $\eta := \frac{4\sqrt{\Lambda_\beta}}{\sqrt{\lambda_{1-\beta}}}$ . Using the inequality for  $\|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2$  again gives us

$$\begin{aligned} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 &\leq \frac{2}{\sqrt{\lambda_{1-\beta}}} (\|\boldsymbol{\varepsilon}\|_2 + \|\mathbf{b}_{S_t}\|_2) \\ &\leq \frac{4 + 2\sqrt{2}}{\sqrt{\lambda_{1-\beta}}} \|\boldsymbol{\varepsilon}\|_2 + \frac{2 \cdot \eta^t}{\sqrt{\lambda_{1-\beta}}} \|\mathbf{b}\|_2 \end{aligned}$$

For large enough  $n$  we have  $\sqrt{\lambda_{1-\beta}} \geq \mathcal{O}(\sqrt{n})$ , which completes the proof.  $\square$

Notice that for random Gaussian noise, this result gives the following convergence guarantee.

**Corollary 11.** *Let the data be generated as before with random Gaussian dense noise i.e.  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b} + \varepsilon$  with  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$  and  $\varepsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \cdot I)$ . Let Algorithm 2 be executed on this data with the thresholding parameter set to  $\beta \geq \alpha$ . Let  $\Sigma_0$  be an invertible matrix such that  $\tilde{X} = \Sigma_0^{-1/2} X$  satisfies the SSC and SSS properties at level  $\gamma$  with constants  $\lambda_\gamma$  and  $\Lambda_\gamma$  respectively (see Definition 1). If the data satisfies  $\frac{4\sqrt{\Lambda_\beta}}{\sqrt{\lambda_{1-\beta}}} < 1$ , then after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 2 obtains an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon + 2\sigma C$ , where  $C > 0$  is the constant in Theorem 10.*

*Proof.* Using tail bounds on Chi-squared distributions [15], we get, with probability at least  $1 - \delta$ ,

$$\|\varepsilon\|_2^2 \leq \sigma^2 \left( n + 2\sqrt{n \log \frac{1}{\delta}} + 2 \log \frac{1}{\delta} \right).$$

Thus, for  $n > 4 \log \frac{1}{\delta}$ , we have  $\|\varepsilon\|_2^2 \leq 2\sigma n$  which proves the result.  $\square$

*Remark 6.* We note that the design assumptions made by Theorem 10 (i.e.  $\frac{4\sqrt{\Lambda_\beta}}{\sqrt{\lambda_{1-\beta}}} < 1$ ) are similar to those made by Theorem 3 and would be satisfied with high probability by data sampled from sub-Gaussian distributions (see Appendix G for details).

## B Proof of Theorem 3

*Theorem 3.* Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix and  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b}$  be the corrupted output with  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Let Algorithm 2 be executed on this data with the thresholding parameter set to  $\beta \geq \alpha$ . Let  $\Sigma_0$  be an invertible matrix such that  $\tilde{X} = \Sigma_0^{-1/2} X$  satisfies the SSC and SSS properties at level  $\gamma$  with constants  $\lambda_\gamma$  and  $\Lambda_\gamma$  respectively (see Definition 1). If the data satisfies  $\frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < 1$ , then after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 2 obtains an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$ .

*Proof.* Let  $\mathbf{r}^t = \mathbf{y} - X^T \mathbf{w}^t$  be the vector of residuals at time  $t$  and  $C_t = X_{S_t} X_{S_t}^\top$ . Since  $\lambda_\alpha > 0$  (something which we shall establish later), we get

$$\mathbf{w}^{t+1} = C_t^{-1} X_{S_t} \mathbf{y}_{S_t} = C_t^{-1} X_{S_t} (X_{S_t}^\top \mathbf{w}^* + \mathbf{b}_{S_t}) = \mathbf{w}^* + C_t^{-1} X_{S_t} \mathbf{b}_{S_t}.$$

Thus, for any set  $S \subset [n]$ , we have

$$\mathbf{r}_S^{t+1} = \mathbf{y}_S - X_S^\top \mathbf{w}_{t+1} = \mathbf{b}_S - X_S^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t}$$

This, gives us

$$\begin{aligned} \|\mathbf{b}_{S_{t+1}}\|_2^2 &= \left\| \mathbf{b}_{S_{t+1}} - X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 - \left\| X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \\ &\stackrel{\zeta_1}{\leq} \left\| \mathbf{b}_{S_*} - X_{S_*}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 - \left\| X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \\ &\stackrel{\zeta_2}{\leq} \left\| X_{S_*}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 - \left\| X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \\ &\leq \left\| X_{S_* \setminus S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top X_{S_{t+1}}^\top C_t^{-1} X_{S_t} \mathbf{b}_{S_t} \\ &\stackrel{\zeta_3}{\leq} \left\| \tilde{X}_{S_* \setminus S_{t+1}}^\top \left( \tilde{X}_{S_t} \tilde{X}_{S_t}^\top \right)^{-1} \tilde{X}_{S_t} \mathbf{b}_{S_t} \right\|_2^2 + 2 \cdot \mathbf{b}_{S_{t+1}}^\top \tilde{X}_{S_{t+1}}^\top \left( \tilde{X}_{S_t} \tilde{X}_{S_t}^\top \right)^{-1} \tilde{X}_{S_t} \mathbf{b}_{S_t} \\ &\stackrel{\zeta_4}{\leq} \frac{\Lambda_\beta^2}{\lambda_{1-\beta}^2} \cdot \|\mathbf{b}_{S_t}\|_2^2 + 2 \cdot \frac{\Lambda_\beta}{\lambda_{1-\beta}} \cdot \|\mathbf{b}_{S_t}\|_2 \|\mathbf{b}_{S_{t+1}}\|_2, \end{aligned}$$

where  $\zeta_1$  follows since the hard thresholding step ensures  $\|\mathbf{r}_{S_{t+1}}^{t+1}\|_2^2 \leq \|\mathbf{r}_{S_*}^{t+1}\|_2^2$  (see Claim 19 and use the fact that  $\beta \geq \alpha$ ),  $\zeta_2$  notices the fact that  $\mathbf{b}_{S_*} = \mathbf{0}$ .  $\zeta_3$  follows from setting  $\tilde{X} = \Sigma_0^{-1/2} X$  and  $X_S^\top C_t^{-1} X_{S'} = \tilde{X}_S^\top (\tilde{X}_{S_t} \tilde{X}_{S_t}^\top)^{-1} \tilde{X}_{S'}$ .  $\zeta_4$  follows from the definition of SSC and SSS properties,  $\|\mathbf{b}_{S_t}\|_0 \leq \|\mathbf{b}\|_0 \leq \beta \cdot n$  and  $|S_* \setminus S_{t+1}| \leq \beta \cdot n$ . Solving the quadratic equation gives us

$$\|\mathbf{b}_{S_{t+1}}\|_2 \leq (1 + \sqrt{2}) \cdot \frac{\Lambda_\beta}{\lambda_{1-\beta}} \cdot \|\mathbf{b}_{S_t}\|_2. \quad (4)$$

Let  $\eta := \frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}}$  denote the convergence rate in (4). We shall show below that for a large family of random designs, we have  $\eta < 1$  if  $n \geq \Omega(p + \log \frac{1}{\delta})$ . We now recall from our earlier discussion that  $\mathbf{w}^{t+1} = \mathbf{w}^* + C_t^{-1} X_{S_t} \mathbf{b}_{S_t}$  which gives us

$$\|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 = \|C_t^{-1} X_{S_t} \mathbf{b}_{S_t}\|_2 \leq \frac{\sqrt{\Lambda_\beta}}{\lambda_{1-\beta}} \cdot \|\mathbf{b}_{S_t}\|_2 \leq \eta^t \cdot \frac{\sqrt{\Lambda_\beta}}{\lambda_{1-\beta}} \|\mathbf{b}\|_2 \leq \epsilon,$$

for  $t \geq \log \frac{1}{\eta} \left( \frac{\sqrt{\Lambda_\beta}}{\lambda_{1-\beta}} \cdot \frac{\|\mathbf{b}\|_2}{\epsilon} \right)$ . Noting that  $\frac{\sqrt{\Lambda_\beta}}{\lambda_{1-\beta}} \leq \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$  establishes the convergence result.  $\square$

## C Proof of Theorem 4

*Theorem 4.* Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix with each  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$ . Let  $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$  and  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Also, let  $\alpha \leq \beta < \frac{1}{65}$  and  $n \geq \Omega(p + \log \frac{1}{\delta})$ . Then, with probability at least  $1 - \delta$ , the data satisfies  $\frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < \frac{9}{10}$ . More specifically, after  $T \geq 10 \log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)$  iterations of Algorithm 1 with the thresholding parameter set to  $\beta$ , we have  $\|\mathbf{w}^T - \mathbf{w}^*\| \leq \epsilon$ .

*Proof.* We note that whenever  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma)$  then  $\Sigma^{-1/2} \mathbf{x} \sim \mathcal{N}(\mathbf{0}, I)$ . Thus, Theorem 15 assures us that with probability at least  $1 - \delta$ , the data matrix  $\tilde{X} = \Sigma^{-1/2} X$  satisfies the SSC and SSS properties with the following constants

$$\begin{aligned} \Lambda_\beta &\leq \beta n \left( 1 + 3e \sqrt{6 \log \frac{e}{\beta}} \right) + \mathcal{O} \left( \sqrt{np + n \log \frac{1}{\delta}} \right) \\ \lambda_{1-\beta} &\geq n - \beta n \left( 1 + 3e \sqrt{6 \log \frac{e}{\beta}} \right) - \Omega \left( \sqrt{np + n \log \frac{1}{\delta}} \right) \end{aligned}$$

Thus, the convergence given by Algorithm 1, when invoked with  $\Sigma_0 = \Sigma$ , relies on the quantity  $\eta = \frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}}$  being less than unity. This translates to the requirement  $(1 + \sqrt{2})\Lambda_\beta \leq \lambda_{1-\beta}$ . Using the above bounds translates that requirement to

$$\underbrace{(2 + \sqrt{2})\beta \left( 1 + 3e \sqrt{6 \log \frac{e}{\beta}} \right)}_{(A)} + \underbrace{\mathcal{O} \left( \sqrt{\frac{p}{n} + \frac{1}{n} \log \frac{1}{\delta}} \right)}_{(B)} < 1.$$

For  $n = \Omega(p + \log \frac{1}{\delta})$ , the second quantity (B) can be made as small a constant as necessary. Tackling the first quantity (A) turns out to be more challenging. However, we can show that for all  $\beta < \frac{1}{190}$ , we get  $\eta = \frac{(1+\sqrt{2})\Lambda_\beta}{\lambda_{1-\beta}} < \frac{9}{10}$  which establishes the claimed result. Thus, Algorithm 1 can tolerate a corruption index of upto  $\alpha \leq \frac{1}{190}$ . However, we note that using a more finely tuned setting of the constant  $\epsilon$  in the proof of Theorem 15 and a more careful proof using tight tail inequalities for chi-squared distributions [15], we can achieve a better corruption level tolerance of  $\alpha < \frac{1}{65}$ .  $\square$

## D Proof of Theorem 5

*Theorem 5.* Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix and  $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$  be the corrupted output with  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Let  $X$  satisfy the SSC and SSS properties at level  $\gamma$  with

648 constants  $\lambda_\gamma$  and  $\Lambda_\gamma$  respectively (see Definition 1). Let Algorithm 1 be executed on this data with  
 649 the GD update (Algorithm 3) with the thresholding parameter set to  $\beta \geq \alpha$  and the step length set  
 650 to  $\eta = \frac{1}{\Lambda_{1-\beta}}$ . If the data satisfies  $\max\{\eta\sqrt{\Lambda_\beta}, 1 - \eta\lambda_{1-\beta}\} \leq \frac{1}{4}$ , then after  $t = \mathcal{O}\left(\log\left(\frac{\|b\|_2}{\sqrt{n}} \frac{1}{\epsilon}\right)\right)$   
 651 iterations, Algorithm 1 obtains an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$ .  
 652

653  
 654 *Proof.* Let  $\mathbf{r}^t = \mathbf{y} - X^\top \mathbf{w}^t$  be the vector of residuals at time  $t$  and  $C_t = X_{S_t} X_{S_t}^\top$ . We have

$$655 \quad \mathbf{w}^{t+1} = \mathbf{w}^t + \eta \cdot X_{S_t} \mathbf{r}_{S_t}^t = \mathbf{w}^t + \eta \cdot X_{S_t} (\mathbf{y}_{S_t} - X_{S_t}^\top \mathbf{w}^t)$$

656  
 657 The thresholding step ensures that  $\|\mathbf{r}_{S_{t+1}}^{t+1}\|_2^2 \leq \|\mathbf{r}_{S_*}^{t+1}\|_2^2$  (see Claim 19 and use  $\beta \geq \alpha$ ) which  
 658 implies

$$659 \quad \|\mathbf{r}_{\text{CR}_{t+1}}^{t+1}\|_2^2 \leq \|\mathbf{r}_{\text{MD}_{t+1}}^{t+1}\|_2^2,$$

660 where  $\text{CR}_{t+1} = S_{t+1} \setminus S_*$  are the *corrupted recoveries* and  $\text{MD}_{t+1} = S_* \setminus S_{t+1}$  are the clean points  
 661 missed out from *detection*. Note that  $|\text{CR}_{t+1}| \leq \alpha \cdot n$  and  $|\text{MD}_{t+1}| \leq \beta \cdot n$ . Since  $\mathbf{b}_{S_*} = \mathbf{0}$  and  
 662  $\text{MD}_{t+1} \subseteq S_*$ , we get

$$663 \quad \|\mathbf{b}_{\text{CR}_{t+1}} + X_{\text{CR}_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1})\|_2 \leq \|X_{\text{MD}_{t+1}}^\top (\mathbf{w}^* - \mathbf{w}^{t+1})\|_2$$

664  
 665 Using the SSS conditions and the fact that  $\|\mathbf{b}_{S_{t+1}}\|_2 = \|\mathbf{b}_{S_{t+1} \setminus S_*}\|_2$  gives us

$$666 \quad \|\mathbf{b}_{S_{t+1}}\|_2 = \|\mathbf{b}_{\text{CR}_{t+1}}\|_2 \leq (\sqrt{\Lambda_\alpha} + \sqrt{\Lambda_\beta}) \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 \leq 2\sqrt{\Lambda_\beta} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2$$

667 Now, using the expression for  $\mathbf{w}^{t+1}$  gives us

$$668 \quad \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 \leq \|(I - \eta C_t)(\mathbf{w}^* - \mathbf{w}^t)\|_2 + \eta \|X_{S_t} \mathbf{b}_{S_t}\|_2$$

669 We will bound the two terms on the right hand separately. We can bound the second term easily as

$$670 \quad \eta \|X_{S_t} \mathbf{b}_{S_t}\|_2 \leq \eta \sqrt{\Lambda_\alpha} \|\mathbf{b}_{S_t}\|_2 \leq \eta \sqrt{\Lambda_\beta} \|\mathbf{b}_{S_t}\|_2,$$

671 since  $\|\mathbf{b}_{S_t}\|_0 \leq \alpha \cdot n$ . For the first term we observe that for  $\eta \leq \frac{1}{\Lambda_{1-\beta}}$ , we have

$$672 \quad \|I - \eta C_t\|_2 = \sup_{\mathbf{v} \in S^{p-1}} |1 - \eta \cdot \mathbf{v}^\top C_t \mathbf{v}| = \sup_{\mathbf{v} \in S^{p-1}} \{1 - \eta \cdot \mathbf{v}^\top C_t \mathbf{v}\} \leq 1 - \eta \lambda_{1-\beta},$$

673 which we can use to bound

$$674 \quad \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 \leq (1 - \eta \lambda_{1-\beta}) \|\mathbf{w}^* - \mathbf{w}^t\|_2 + \eta \sqrt{\Lambda_\beta} \|\mathbf{b}_{S_t}\|_2$$

675 This gives us, for  $\eta = \frac{1}{\Lambda_{1-\beta}}$ ,

$$676 \quad \|\mathbf{b}_{S_{t+1}}\|_2 \leq 2\sqrt{\Lambda_\beta} \|\mathbf{w}^* - \mathbf{w}^{t+1}\|_2 \leq 2 \underbrace{\left(1 - \frac{\lambda_{1-\beta}}{\Lambda_{1-\beta}}\right)}_{(P)} \sqrt{\Lambda_\beta} \|\mathbf{w}^* - \mathbf{w}^t\|_2 + 2 \underbrace{\frac{\Lambda_\beta}{\Lambda_{1-\beta}}}_{(Q)} \|\mathbf{b}_{S_t}\|_2.$$

677 For Gaussian designs and small enough  $\beta$ , we can show  $(Q) \leq \frac{1}{4}$  as we did in Theorem 4. To bound  
 678  $(P)$ , we use the lower bound on  $\lambda_{1-\beta}$  given by Theorem 15 and use the following tighter upper  
 679 bound for  $\Lambda_{1-\beta}$ :

$$680 \quad \Lambda_{1-\beta} \leq \left( (1 - \beta) + 3e \sqrt{6\beta(1 - \beta) \log \frac{e}{\beta}} \right) n + \mathcal{O} \left( \sqrt{np + n \log \frac{1}{\delta}} \right)$$

681 The above bound is obtained similarly to the one in Theorem 15 but uses the identity  $\binom{n}{k} = \binom{n}{n-k} \leq$   
 682  $\left(\frac{en}{n-k}\right)^{n-k}$  for values of  $k \geq n/2$  instead. For small enough  $\beta$  and  $n = \Omega(\kappa^2(\Sigma)(p + \log \frac{1}{\delta}))$ ,  
 683 we can then show  $(P) \leq \frac{1}{4}$  as well. Let  $\Psi_t := \sqrt{n} \|\mathbf{w}^* - \mathbf{w}^t\|_2 + \|\mathbf{b}_{S_t}\|$ . Using elementary  
 684 manipulations and the fact that  $\sqrt{\Lambda_\beta} \geq \Omega(\sqrt{n})$ , we can then show that

$$685 \quad \Psi_{t+1} \leq 3/4 \cdot \Psi_t.$$

686 Thus, in  $t = \mathcal{O}\left(\log\left(\left(\|\mathbf{w}^*\|_2 + \frac{\|b\|_2}{\sqrt{n}}\right) \frac{1}{\epsilon}\right)\right)$  iterations of the algorithm, we arrive at an  $\epsilon$ -optimal  
 687 solution i.e.  $\|\mathbf{w}^* - \mathbf{w}^t\|_2 \leq \epsilon$ . A similar argument holds true for sub-Gaussian designs as well.  $\square$   
 688  
 689

## 702 E Proof of Theorem 6

703  
704 *Theorem 6.* Suppose Algorithm 4 is executed on data that allows Algorithms 2 and 3 a convergence  
705 rate of  $\eta_{\text{FC}}$  and  $\eta_{\text{GD}}$  respectively. Suppose we have  $2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} < 1$ . Then for *any* interleavings of the  
706 FC and GD steps that the policy may enforce, after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 4  
707 ensures an  $\epsilon$ -optimal solution i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\| \leq \epsilon$ .  
708

709  
710 *Proof.* Our proof shall essentially show that the FC and GD steps do not undo the progress made by  
711 the other if executed in succession and if  $2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} < 1$ , actually ensure non-trivial progress. Let

$$712 \Psi_t^{\text{FC}} = \|\mathbf{b}_{S_t}\|_2$$

$$713 \Psi_t^{\text{GD}} = \sqrt{n} \|\mathbf{w}^t - \mathbf{w}^*\| + \|\mathbf{b}_{S_t}\|_2$$

714  
715 denote the potential functions used in the analyses of the FC and GD algorithms before. Then we  
716 will show below that if the FC and GD algorithms are executed in steps  $t$  and  $t + 1$  then we have

$$717 \Psi_{t+2}^{\text{FC}} \leq 2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} \cdot \Psi_t^{\text{FC}}$$

718 Alternatively, if the GD and FC algorithms are executed in steps  $t$  and  $t + 1$  respectively, then

$$719 \Psi_{t+2}^{\text{GD}} \leq 2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} \cdot \Psi_t^{\text{GD}}$$

720 Thus, if algorithm executes the FC step at the time step  $t$ , then it would at least ensure  $\Psi_t^{\text{FC}} \leq$   
721  $(2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}})^{t/2} \cdot \Psi_0^{\text{FC}}$  (similarly if the last step is a GD step). Since both the FC and GD algorithms  
722 ensure  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$  for  $t \geq \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$ , the claim would follow.  
723  
724  
725

726 We now prove the two claimed results regarding the two types of interleaving below  
727

### 728 1. FC $\rightarrow$ GD

729 The FC step guarantees  $\|\mathbf{b}_{S_{t+1}}\|_2 \leq \eta_{\text{FC}} \cdot \|\mathbf{b}_{S_t}\|_2$  as well as  $\|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 \leq \eta_{\text{FC}} \cdot \frac{\|\mathbf{b}_{S_t}\|_2}{\sqrt{n}}$ ,  
730 whereas the GD step guarantees  $\Psi_{t+2}^{\text{GD}} \leq \eta_{\text{GD}} \cdot \Psi_{t+1}^{\text{GD}}$ . Together these guarantee

$$731 \sqrt{n} \|\mathbf{w}^{t+2} - \mathbf{w}^*\|_2 + \|\mathbf{b}_{S_{t+2}}\|_2 \leq \eta_{\text{GD}} \cdot \sqrt{n} \|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 + \|\mathbf{b}_{S_{t+1}}\|_2$$

$$732 \leq 2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} \cdot \|\mathbf{b}_{S_t}\|_2$$

733 Since  $\sqrt{n} \|\mathbf{w}^{t+2} - \mathbf{w}^*\|_2 \geq 0$ , this yields the result.  
734

### 735 2. GD $\rightarrow$ FC

736 The GD step guarantees  $\Psi_{t+1}^{\text{GD}} \leq \eta_{\text{GD}} \cdot \Psi_t^{\text{GD}}$  whereas the FC step guarantees  $\|\mathbf{b}_{S_{t+2}}\|_2 \leq$   
737  $\eta_{\text{FC}} \cdot \|\mathbf{b}_{S_{t+1}}\|_2$  as well as  $\|\mathbf{w}^{t+2} - \mathbf{w}^*\|_2 \leq \eta_{\text{FC}} \cdot \frac{\|\mathbf{b}_{S_{t+1}}\|_2}{\sqrt{n}}$ . Together these guarantee

$$738 \sqrt{n} \|\mathbf{w}^{t+2} - \mathbf{w}^*\|_2 + \|\mathbf{b}_{S_{t+2}}\|_2 \leq 2\eta_{\text{FC}} \|\mathbf{b}_{S_{t+1}}\|_2$$

$$739 \leq 2 \cdot \eta_{\text{FC}} \cdot \eta_{\text{GD}} \cdot \Psi_t^{\text{GD}},$$

740 where the second step follows from the GD step guarantee since  $\sqrt{n} \|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 \geq 0$ .  
741  
742

743 This finishes the proof.  $\square$   
744  
745

## 746 F Proof of Theorem 9

747  
748 *Theorem 9.* Let  $X = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{p \times n}$  be the given data matrix and  $\mathbf{y} = X^T \mathbf{w}^* + \mathbf{b}$  be the  
749 corrupted output with  $\|\mathbf{w}^*\|_0 \leq s^*$  and  $\|\mathbf{b}\|_0 \leq \alpha \cdot n$ . Let Algorithm 2 be executed on this data  
750 with the IHT update from [12] and thresholding parameter set to  $\beta \geq \alpha$ . Let  $\Sigma_0$  be an invertible  
751 matrix such that  $\Sigma_0^{-1/2} X$  satisfies the SRSC and SRSS properties at level  $(\gamma, 2s + s^*)$  with constants  
752  $\alpha_{(\gamma, 2s + s^*)}$  and  $L_{(\gamma, 2s + s^*)}$  respectively (see Definition 8) for  $s \geq 32 \left(\frac{L_{(\gamma, 2s + s^*)}}{\alpha_{(\gamma, 2s + s^*)}}\right)$  with  $\gamma = 1 - \beta$ . If  
753  $X$  also satisfies  $\frac{4L_{(\beta, s + s^*)}}{\alpha_{(1 - \beta, s + s^*)}} < 1$ , then after  $t = \mathcal{O}\left(\log\left(\frac{1}{\sqrt{n}} \frac{\|\mathbf{b}\|_2}{\epsilon}\right)\right)$  iterations, Algorithm 2 obtains  
754  
755

756 an  $\epsilon$ -accurate solution  $\mathbf{w}^t$  i.e.  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$ . In particular, if  $X$  is sampled from a Gaussian  
 757 distribution  $\mathcal{N}(\mathbf{0}, \Sigma)$  and  $n \geq \Omega\left((2s + s^*) \log p + \log \frac{1}{\delta}\right)$ , then for all values of  $\alpha \leq \beta < \frac{1}{65}$ , we  
 758 can guarantee recovery as  $\|\mathbf{w}^t - \mathbf{w}^*\|_2 \leq \epsilon$ .  
 759

760 *Proof.* We first begin with the guarantee provided by existing sparse recovery techniques. The  
 761 results of [12], for example, indicate that if the input to the algorithm indeed satisfies the RSC and  
 762 RSS properties at the level  $(1 - \beta, 2s + s^*)$  with constants  $\alpha_{2s+s^*}$  and  $L_{2s+s^*}$  for  $s \geq 32 \left(\frac{L_{2s+s^*}}{\alpha_{2s+s^*}}\right)$ ,  
 763 then in time  $\tau = \mathcal{O}\left(\frac{L_{2s+s^*}}{\alpha_{2s+s^*}} \cdot \log\left(\frac{\|b\|_2}{\rho}\right)\right)$ , the IHT algorithm [12, Algorithm 1] outputs an updated  
 764 model  $\mathbf{w}^{t+1}$  that satisfies  $\|\mathbf{w}^{t+1}\|_0 \leq s$ , as well as  
 765  
 766

$$767 \quad \|\mathbf{y}_{S_t} - X_{S_t}^\top \mathbf{w}^{t+1}\|_2^2 \leq \|\mathbf{y}_{S_t} - X_{S_t}^\top \mathbf{w}^*\|_2^2 + \rho.$$

769 We will set  $\rho$  later. Since the SRSC and SRSS properties ensure the above and  $\mathbf{y} = X^\top \mathbf{w}^* + \mathbf{b}$ , this  
 770 gives us  
 771

$$772 \quad \|X_{S_t}^\top (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2^2 \leq 2(\mathbf{w}^{t+1} - \mathbf{w}^*)^\top X_{S_t}^\top \mathbf{b}_{S_t} + \rho = 2(\mathbf{w}^{t+1} - \mathbf{w}^*)^\top X_{S_t \cap \bar{S}_*}^\top \mathbf{b}_{S_t \cap \bar{S}_*} + \rho,$$

773 since  $\mathbf{b}_S = \mathbf{0}$  for any set  $S \cap \bar{S}_* = \phi$ . We now analyze the two sides separately below using the  
 774 SRSC and SRSS properties below. For any  $S \subset [n]$ , denote  $\tilde{X}_S := \Sigma_0^{-1/2} X$ .  
 775

$$776 \quad \|X_{S_t}^\top (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2^2 = \|\tilde{X}_{S_t}^\top \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2^2 \geq \alpha_{(1-\beta, s+s^*)} \left\| \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*) \right\|_2^2$$

$$777 \quad \|X_{S_t \cap \bar{S}_*}^\top (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2 = \|\tilde{X}_{S_t \cap \bar{S}_*}^\top \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2 \leq \sqrt{L_{(\beta, s+s^*)}} \left\| \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*) \right\|_2.$$

778  
 779 Now, if  $\|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 \geq \epsilon$ , then  $\left\| \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*) \right\|_2 \geq \sqrt{\lambda_{\min}(\Sigma_0)} \cdot \epsilon$ . This give us  
 780

$$781 \quad \left\| \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*) \right\|_2 \leq \frac{2\sqrt{L_{(\beta, s+s^*)}}}{\alpha_{(1-\beta, s+s^*)}} \|\mathbf{b}_{S_t \cap \bar{S}_*}\|_2 + \frac{\rho}{\alpha_{(1-\beta, s+s^*)}}$$

$$782 \quad = \frac{2\sqrt{L_{(\beta, s+s^*)}}}{\alpha_{(1-\beta, s+s^*)}} \|\mathbf{b}_{S_t}\|_2 + \frac{\rho}{\epsilon \cdot \sqrt{\lambda_{\min}(\Sigma_0)} \cdot \alpha_{(1-\beta, s+s^*)}}.$$

783 We note that although we declared the SRSC and SRSS properties for the action of matrices on  
 784 sparse vectors (such as  $\mathbf{w}^* - \mathbf{w}^{t+1}$ ), we instead applied them above to the action of matrices on  
 785 sparse vectors transformed by  $\Sigma_0^{1/2}$  ( $\Sigma_0^{1/2} (\mathbf{w}^* - \mathbf{w}^{t+1})$ ). Since  $\Sigma_0^{1/2} \mathbf{v}$  need not be sparse even if  $\mathbf{v}$   
 786 is sparse, this appears to pose a problem. However, all we need to resolve this is to notice that the  
 787 proof technique of Theorem 18 which would be used to establish the SRSC and SRSS properties,  
 788 holds in general for not just the action of a matrix on the set of sparse vectors, but on vectors in the  
 789 union of any fixed set of low dimensional subspaces.

790 More specifically, we can modify the RSC and RSS properties (and by extension, the SRSC and  
 791 SRSS properties), to requiring that the matrix  $X$  act as an approximate isometry on the following  
 792 set of vectors  $S_{(s, \Sigma_0)}^{p-1} := \left\{ \mathbf{v} : \mathbf{v} = \Sigma_0^{-1/2} \mathbf{v}' \text{ for some } \mathbf{v}' \in S_s^{p-1} \right\}$ . We refer the reader to the work  
 793 of [16] which describes this technique in great detail. Proceeding with the proof, the assurance of  
 794 the thresholding step, as used in the proof of Theorem 5, along with a straightforward application of  
 795 the (modified) SRSS property gives us  
 796

$$801 \quad \|\mathbf{b}_{S_{t+1}}\|_2 \leq \|X_{\text{CR}_{t+1}}^\top (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2 + \|X_{\text{MD}_{t+1}}^\top (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2$$

$$802 \quad = \|\tilde{X}_{\text{CR}_{t+1}}^\top \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2 + \|\tilde{X}_{\text{MD}_{t+1}}^\top \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*)\|_2$$

$$803 \quad \leq 2\sqrt{L_{(\beta, s+s^*)}} \left\| \Sigma_0^{1/2} (\mathbf{w}^{t+1} - \mathbf{w}^*) \right\|_2$$

$$804 \quad \leq \frac{4L_{(\beta, s+s^*)}}{\alpha_{(1-\beta, s+s^*)}} \|\mathbf{b}_{S_t}\|_2 + \frac{2\rho\sqrt{L_{(\beta, s+s^*)}}}{\epsilon \cdot \sqrt{\lambda_{\min}(\Sigma_0)} \cdot \alpha_{(1-\beta, s+s^*)}}$$

Thus, whenever  $\|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 > \epsilon$ , in successive steps,  $\|\mathbf{b}_{S_t}\|_2$  undergoes a linear decrease. Denoting  $\eta := \frac{4L_{(\beta, s+s^*)}}{\alpha_{(1-\beta, s+s^*)}}$ , we get

$$\|\mathbf{b}_{S_{t+1}}\|_2 \leq \eta^t \cdot \|\mathbf{b}\|_2 + \left(\frac{1-\eta^t}{1-\eta}\right) \frac{2\rho\sqrt{L_{(\beta, s+s^*)}}}{\epsilon \cdot \sqrt{\lambda_{\min}(\Sigma_0)} \cdot \alpha_{(1-\beta, s+s^*)}}$$

and using  $\left\|\Sigma_0^{1/2}(\mathbf{w}^t - \mathbf{w}^*)\right\|_2 \geq \sqrt{\lambda_{\min}(\Sigma_0)} \|\mathbf{w}^t - \mathbf{w}^*\|_2$  gives us

$$\begin{aligned} \|\mathbf{w}^{t+1} - \mathbf{w}^*\|_2 &\leq \frac{2\sqrt{L_{(\beta, s+s^*)}}}{\sqrt{\lambda_{\min}(\Sigma_0)} \cdot \alpha_{(1-\beta, s+s^*)}} \|\mathbf{b}_{S_{t+1}}\|_2 + \frac{\rho}{\lambda_{\min}(\Sigma_0) \cdot \alpha_{(1-\beta, s+s^*)}} \\ &\leq \eta^t \frac{2\sqrt{L_{(\beta, s+s^*)}}}{\sqrt{\lambda_{\min}(\Sigma_0)} \cdot \alpha_{(1-\beta, s+s^*)}} \|\mathbf{b}\|_2 + \frac{36\rho}{\epsilon \cdot \lambda_{\min}(\Sigma_0) \cdot \alpha_{(1-\beta, s+s^*)}}, \end{aligned}$$

where we have assumed that  $\frac{4L_{(\beta, s+s^*)}}{\alpha_{(1-\beta, s+s^*)}} < 9/10$ , something that we shall establish below. Note that  $\lambda_{\min}(\Sigma_0) > 0$  since  $\Sigma$  is assumed to be invertible. In the random design settings we shall consider, we also have  $\frac{\sqrt{L_{(\beta, s+s^*)}}}{\sqrt{\lambda_{\min}(\Sigma_0)} \cdot \alpha_{(1-\beta, s+s^*)}} = \mathcal{O}\left(\frac{1}{\sqrt{n}}\right)$ . Then setting  $\rho \leq \frac{1}{72}\epsilon^2 \cdot \lambda_{\min}(\Sigma_0) \cdot \alpha_{(1-\beta, s+s^*)}$  proves the convergence result.

As before, we can use the above result to establish sparse recovery guarantees in the statistical setting for Gaussian and sub-Gaussian design models. If our data matrix  $X$  is generated from a Gaussian distribution  $\mathcal{N}(\mathbf{0}, \Sigma)$  for some invertible  $\Sigma$ , then the results in Theorem 18 can be used to establish that  $\Sigma^{-1/2}X$  satisfies the SRSC and SRSS properties at the required levels and that for  $\alpha < \frac{1}{190}$  and  $n \geq \Omega\left((2s + s^*) \log p + \log \frac{1}{\delta}\right)$ , we have  $\eta = \frac{2L_{(\beta, s+s^*)}}{\alpha_{(1-\beta, s+s^*)}} < 9/10$ .

Thus, the above result can be applied with  $\Sigma_0 = \Sigma$  to get convergence guarantees in the general Gaussian setting. We note that the above analysis can tolerate the same level of corruption as Theorem 4 and thus, we can improve the noise tolerance level to  $\alpha \leq \frac{1}{65}$  here as well. We also note that these results can be readily extended to the sub-Gaussian setting as well.  $\square$

## G Robust Statistical Estimation

This section elaborates on how results on the convergence guarantees of our algorithms can be used to give guarantees for robust statistical estimation problems. We begin with a few definition of sampling models that would be used in our results.

**Definition 12.** A random variable  $x \in \mathbb{R}$  is called sub-Gaussian if the following quantity is finite

$$\sup_{p \geq 1} p^{-1/2} (\mathbb{E} |x|^p)^{1/p}.$$

Moreover, the smallest upper bound on this quantity is referred to as the sub-Gaussian norm of  $x$  and denoted as  $\|x\|_{\psi_2}$ .

**Definition 13.** A vector-valued random variable  $\mathbf{x} \in \mathbb{R}^p$  is called sub-Gaussian if its unidimensional marginals  $\langle \mathbf{x}, \mathbf{v} \rangle$  are sub-Gaussian for all  $\mathbf{v} \in S^{p-1}$ . Moreover, its sub-Gaussian norm is defined as follows

$$\|X\|_{\psi_2} := \sup_{\mathbf{v} \in S^{p-1}} \|\langle \mathbf{x}, \mathbf{v} \rangle\|_{\psi_2}$$

We will begin with the analysis of Gaussian designs and then extend our analysis for the class of general sub-Gaussian designs.

**Lemma 14.** Let  $X \in \mathbb{R}^{p \times n}$  be a matrix whose columns are sampled i.i.d from a standard Gaussian distribution i.e.  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, I)$ . Then for any  $\epsilon > 0$ , with probability at least  $1 - \delta$ ,  $X$  satisfies

$$\begin{aligned} s_{\max}(XX^\top) &\leq n + (1 - 2\epsilon)^{-1} \sqrt{cnp + c'n \log \frac{2}{\delta}} \\ s_{\min}(XX^\top) &\geq n - (1 - 2\epsilon)^{-1} \sqrt{cnp + c'n \log \frac{2}{\delta}}, \end{aligned}$$

where  $c = 24e^2 \log \frac{3}{\epsilon}$  and  $c' = 24e^2$ .

864 *Proof.* We will first use the fact that  $X$  is sampled from a standard Gaussian to show that its covari-  
865 ance concentrates around identity. Thus, we first show that with high probability,  
866

$$867 \quad \|XX^\top - nI\|_2 \leq \epsilon_1$$

868 for some  $\epsilon_1 < 1$ . Doing so will automatically establish the following result  
869

$$870 \quad n - \epsilon_1 \leq s_{\min}(XX^\top) \leq s_{\max}(XX^\top) \leq n + \epsilon_1.$$

871 Let  $A := XX^\top - I$ . We will use the technique of covering numbers [17] to establish the above.  
872 Let  $\mathcal{C}^{p-1}(\epsilon) \subset S^{p-1}$  be an  $\epsilon$  cover for  $S^{p-1}$  i.e. for all  $\mathbf{u} \in S^{p-1}$ , there exists at least one  $\mathbf{v} \in \mathcal{C}^{p-1}$   
873 such that  $\|\mathbf{u} - \mathbf{v}\|_2 \leq \epsilon$ . Standard constructions [17, see Lemma 5.2] guarantee such a cover of size  
874 at most  $(1 + \frac{2}{\epsilon})^p \leq (\frac{3}{\epsilon})^p$ . Now for any  $\mathbf{u} \in S^{p-1}$  and  $\mathbf{v} \in \mathcal{C}^{p-1}$  such that  $\|\mathbf{u} - \mathbf{v}\|_2 \leq \epsilon$ , we have  
875

$$876 \quad |\mathbf{u}^\top A\mathbf{u} - \mathbf{v}^\top A\mathbf{v}| \leq |\mathbf{u}^\top A(\mathbf{u} - \mathbf{v})| + |\mathbf{v}^\top A(\mathbf{u} - \mathbf{v})| \leq 2\epsilon \|A\|_2,$$

877 which gives us  
878

$$879 \quad \|XX^\top - nI\|_2 \leq (1 - 2\epsilon)^{-1} \cdot \sup_{\mathbf{v} \in \mathcal{C}^{p-1}(\epsilon)} \left| \|X^\top \mathbf{v}\|_2^2 - n \right|.$$

880 Now for a fixed  $\mathbf{v} \in S^{n-1}$ , the random variable  $\|X^\top \mathbf{v}\|_2^2$  is distributed as a  $\chi^2(n)$  distribution with  
881  $n$  degrees of freedom. Using Lemma 20, we get, for any  $\mu < 1$ ,  
882

$$883 \quad \mathbb{P} \left[ \left| \|X^\top \mathbf{v}\|_2^2 - n \right| \geq \mu n \right] \leq 2 \exp \left( - \min \left\{ \frac{\mu^2 n^2}{24ne^2}, \frac{\mu n}{4\sqrt{3e}} \right\} \right) \leq 2 \exp \left( - \frac{\mu^2 n}{24e^2} \right).$$

884 Setting  $\mu^2 = c \cdot \frac{n}{n} + c' \cdot \frac{\log \frac{2}{\delta}}{n}$ , where  $c = 24e^2 \log \frac{3}{\epsilon}$  and  $c' = 24e^2$ , and taking a union bound over  
885 all  $\mathcal{C}^{p-1}(\epsilon)$ , we get  
886

$$887 \quad \mathbb{P} \left[ \sup_{\mathbf{v} \in \mathcal{C}^{p-1}(\epsilon)} \left| \|X^\top \mathbf{v}\|_2^2 - n \right| \geq \sqrt{cnp + c'n \log \frac{2}{\delta}} \right] \leq 2 \left( \frac{3}{\epsilon} \right)^p \exp \left( - \frac{\mu^2 n}{24e^2} \right) \leq \delta.$$

888 This implies that with probability at least  $1 - \delta$ ,  
889

$$890 \quad \|XX^\top - nI\|_2 \leq (1 - 2\epsilon)^{-1} \sqrt{cnp + c'n \log \frac{2}{\delta}},$$

891 which gives us the claimed bounds on the singular values of  $XX^\top$ .  $\square$   
892

893 **Theorem 15.** Let  $X \in \mathbb{R}^{p \times n}$  be a matrix whose columns are sampled i.i.d from a standard Gaussian  
894 distribution i.e.  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, I)$ . Then for any  $\gamma > 0$ , with probability at least  $1 - \delta$ , the matrix  $X$   
895 satisfies the SSC and SSS properties with constants  
896

$$897 \quad \Lambda_\gamma^{Gauss} \leq \gamma n \left( 1 + 3e \sqrt{6 \log \frac{e}{\gamma}} \right) + \mathcal{O} \left( \sqrt{np + n \log \frac{1}{\delta}} \right)$$

$$898 \quad \lambda_\gamma^{Gauss} \geq n - (1 - \gamma)n \left( 1 + 3e \sqrt{6 \log \frac{e}{1 - \gamma}} \right) - \Omega \left( \sqrt{np + n \log \frac{1}{\delta}} \right).$$

899 *Proof.* For any fixed  $S \in \mathcal{S}_\gamma$ , Lemma 14 guarantees the following bound  
900

$$901 \quad s_{\max}(X_S X_S^\top) \leq \gamma n + (1 - 2\epsilon)^{-1} \sqrt{c\gamma np + c'\gamma n \log \frac{2}{\delta}}.$$

902 Taking a union bound over  $\mathcal{S}_\gamma$  and noting that  $\binom{n}{k} \leq \left( \frac{en}{k} \right)^k$  for all  $1 \leq k \leq n$ , gives us  
903

$$904 \quad \Lambda_\gamma \leq \gamma n + (1 - 2\epsilon)^{-1} \sqrt{c\gamma np + c'\gamma^2 n^2 \log \frac{e}{\gamma} + c'\gamma n \log \frac{2}{\delta}}$$

$$905 \quad \leq \gamma n \left( 1 + (1 - 2\epsilon)^{-1} \sqrt{c' \log \frac{e}{\gamma}} \right) + (1 - 2\epsilon)^{-1} \sqrt{c\gamma np + c'\gamma n \log \frac{2}{\delta}},$$

which finishes the first bound after setting  $\epsilon = 1/6$ . For the second bound, we use the equality

$$X_S X_S^\top = X X^\top - X_{\bar{S}} X_{\bar{S}}^\top,$$

which provides the following bound for  $\lambda_\gamma$

$$\lambda_\gamma \geq s_{\min}(X X^\top) - \sup_{T \in \mathcal{S}_{1-\gamma}} X_T X_T^\top = s_{\min}(X X^\top) - \Lambda_{1-\gamma}.$$

Using Lemma 14 to bound the first quantity and the first part of this theorem to bound the second quantity gives us, with probability at least  $1 - \delta$ ,

$$\lambda_\gamma \geq n - \gamma' n \left( 1 + (1 - 2\epsilon)^{-1} \sqrt{c' \log \frac{e}{\gamma'}} \right) - (1 - 2\epsilon)^{-1} \left( 1 + \sqrt{\gamma'} \right) \sqrt{cnp + c'n \log \frac{2}{\delta}},$$

where  $\gamma' = 1 - \gamma$ . This proves the second bound after setting  $\epsilon = 1/6$ .  $\square$

We now extend our analysis to the class of isotropic subGaussian distributions. We note that this analysis is without loss of generality since for non-isotropic sub-Gaussian distributions, we can simply use the fact that Theorem 3 can admit whitened data for calculation of the SSC and SSS constants as we did for the case of non-isotropic Gaussian distributions.

**Lemma 16.** *Let  $X \in \mathbb{R}^{p \times n}$  be a matrix with columns sampled from some sub-Gaussian distribution with sub-Gaussian norm  $K$  and covariance  $\Sigma$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , each of the following statements holds true:*

$$\begin{aligned} s_{\max}(X X^\top) &\leq \lambda_{\max}(\Sigma) \cdot n + C_K \cdot \sqrt{pn} + t\sqrt{n} \\ s_{\min}(X X^\top) &\geq \lambda_{\min}(\Sigma) \cdot n - C_K \cdot \sqrt{pn} - t\sqrt{n}, \end{aligned}$$

where  $t = \sqrt{\frac{1}{c_K} \log \frac{2}{\delta}}$ , and  $c_K, C_K$  are absolute constants that depend only on the sub-Gaussian norm  $K$  of the distribution.

*Proof.* Since the singular values of a matrix are unchanged upon transposition, we shall prove the above statements for  $X^\top$ . The benefit of this is that we get to work with a matrix with independent rows, so that standard results can be applied. The proof technique used in [17, Theorem 5.39] (see also Remark 5.40 (1) therein) can be used to establish the following result: with probability at least  $1 - \delta$ , with  $t$  set as mentioned in the theorem statement, we have

$$\left\| \frac{1}{n} X X^\top - \Sigma \right\| \leq C_K \sqrt{\frac{p}{n}} + \frac{t}{\sqrt{n}}$$

This implies that for any  $\mathbf{v} \in S^{p-1}$ , we have

$$\left| \frac{1}{n} \|X^\top \mathbf{v}\|_2^2 - \mathbf{v}^\top \Sigma \mathbf{v} \right| = \left| \frac{1}{n} \mathbf{v}^\top X X^\top \mathbf{v} - \mathbf{v}^\top \Sigma \mathbf{v} \right| \leq \left| \frac{1}{n} X X^\top \mathbf{v} - \Sigma \mathbf{v} \right| \leq C_K \sqrt{\frac{p}{n}} + \frac{t}{\sqrt{n}}.$$

The results then follow from elementary manipulations and the fact that the singular values and eigenvalues of real symmetric matrices coincide.  $\square$

**Theorem 17.** *Let  $X \in \mathbb{R}^{p \times n}$  be a matrix with columns sampled from some sub-Gaussian distribution with sub-Gaussian norm  $K$  and covariance  $\Sigma$ . Let  $c_K, C_K$  and  $t$  be fixed to values as required in Lemma 16. Note that  $c_K$  and  $C_K$  are absolute constants depend only on the sub-Gaussian norm  $K$  of the distribution. Let  $\gamma \in (0, 1]$  be some fixed constant. Then, with we have the following:*

$$\Lambda_\gamma^{\text{subGauss}(K, \Sigma)} \leq \left( \lambda_{\max}(\Sigma) \cdot \gamma + \sqrt{\frac{\gamma}{c_K} \log \frac{e}{\gamma}} \right) \cdot n + C_K \cdot \sqrt{\gamma pn} + t\sqrt{n}.$$

Furthermore, fix any  $\epsilon \in (0, 1)$  and let  $\gamma$  be a value in  $(0, 1)$  satisfying the following

$$\gamma > 1 - \min \left\{ \frac{\epsilon \cdot \lambda_{\min}(\Sigma)}{\lambda_{\max}(\Sigma)}, \exp \left( 1 + W_{-1} \left( -\frac{c_K \epsilon^2 \cdot \lambda_{\min}^2(\Sigma)}{e} \right) \right) \right\},$$

where  $W_{-1}(\cdot)$  is the lower branch of the real valued restriction of the Lambert W function. Then we have, with the same confidence,

$$\lambda_\gamma^{\text{subGauss}(K, \Sigma)} \geq (1 - 2\epsilon) \cdot \lambda_{\min}(\Sigma) \cdot n - C_K \left( 1 + \sqrt{1 - \gamma} \right) \sqrt{pn} - 2t\sqrt{n}$$

972 *Proof.* The first result follows from an application of Lemma 16, a union bound over sets in  $\mathcal{S}_\gamma$ , as  
 973 well as the bound  $\binom{n}{k} \leq \left(\frac{en}{k}\right)^k$  for all  $1 \leq k \leq n$  which puts a bound on the number of sparse sets  
 974 as  $\log |\mathcal{S}_\gamma| \leq \gamma \cdot n \log \frac{e}{\gamma}$ .  
 975

976 For the second result, we observe that  $X_S X_S^\top = X X^\top - X_{\bar{S}} X_{\bar{S}}^\top$ , so that  $s_{\min}(X_S X_S^\top) \geq$   
 977  $s_{\min}(X X^\top) - s_{\max}(X_{\bar{S}} X_{\bar{S}}^\top)$ . This gives us  
 978

$$979 \inf_{S \in \mathcal{S}_\gamma} s_{\min}(X_S X_S^\top) \geq s_{\min}(X X^\top) - \sup_{S \in \mathcal{S}_{1-\gamma}} s_{\max}(X_S X_S^\top).$$

981 Using Lemma 16 and the first part of this result gives us  
 982

$$983 \inf_{S \in \mathcal{S}_\gamma} s_{\min}(X_S X_S^\top) \geq \lambda_{\min}(\Sigma) \cdot n - C_K \cdot \sqrt{pn} - t\sqrt{n}$$

$$984 - \left( \lambda_{\max}(\Sigma)(1-\gamma) + \sqrt{\frac{1-\gamma}{c_K} \log \frac{e}{1-\gamma}} \right) n - C_K \sqrt{(1-\gamma)pn} - t\sqrt{n}$$

$$985 = \left( \lambda_{\min}(\Sigma) - \lambda_{\max}(\Sigma)(1-\gamma) - \sqrt{\frac{1-\gamma}{c_K} \log \frac{e}{1-\gamma}} \right) n$$

$$986 - C_K \left( 1 + \sqrt{1-\gamma} \right) \sqrt{pn} - 2t\sqrt{n}$$

$$987 \geq (1-2\epsilon) \cdot \lambda_{\min}(\Sigma) \cdot n - C_K \left( 1 + \sqrt{1-\gamma} \right) \sqrt{pn} - 2t\sqrt{n},$$

988 where the last step follows from the assumptions on  $\gamma$  and by noticing that it suffices to show the  
 989 following two inequalities to establish the last step  
 990

- 991 1.  $\lambda_{\max}(\Sigma)(1-\gamma) \leq \epsilon \cdot \lambda_{\min}(\Sigma)$
- 992 2.  $(1-\gamma) \log \frac{e}{1-\gamma} \leq c_K \epsilon^2 \cdot \lambda_{\min}^2(\Sigma)$

993 The first part gives us the condition  $\gamma > 1 - \frac{\epsilon \cdot \lambda_{\min}(\Sigma)}{\lambda_{\max}(\Sigma)}$  in a straightforward manner. For the second  
 994 part, denote  $v = c_K \epsilon^2 \cdot \lambda_{\min}^2(\Sigma)$ . Note that for  $v \geq 1$ , all values of  $\gamma \in (0, 1]$  satisfy the inequality.  
 995

996 Otherwise we require the use of the Lambert W function (also known as the product logarithm  
 997 function). This function ensures that its value  $W(z)$  for any  $z > -1/e$  satisfies  $z = W(z)e^{W(z)}$ . In  
 998 our case, making a change of variable  $(1-\gamma) = e^\eta$  gives us the inequality  $(\eta-1)e^{\eta-1} \geq -v/e$ .  
 999 Note that since  $v \leq 1$  in this case,  $-v/e \in (-1/e, 0)$  i.e. a valid value for the Lambert W function.  
 1000 However,  $(-1/e, 0)$  is also the region in which the Lambert W function is multi-valued. Taking  
 1001 the worse bound for  $\gamma$  by choosing the lower branch  $W_{-1}(\cdot)$  gives us the second condition  $\gamma \geq$   
 1002  $1 - \exp\left(1 + W_{-1}\left(-\frac{c_K \epsilon^2 \cdot \lambda_{\min}^2(\Sigma)}{e}\right)\right)$ .  $\square$   
 1003

1004 It is important to note that for any  $-1/e \leq z < 0$ , we have  $\exp(1 + W_{-1}(z)) > 0$  which means  
 1005 that the bounds imposed on  $\gamma$  by Theorem 17 always allow a non-zero fraction of the data points  
 1006 to be corrupted in an adversarial manner. However, the exact value of that fraction depends, in  
 1007 a complicated manner, on the sub-Gaussian norm of the underlying distribution, as well as the  
 1008 condition number and the smallest eigenvalue of the second moment of the underlying distribution.  
 1009

1010 We also note that due to the generic nature of the previous analysis, which can handle the entire class  
 1011 of sub-Gaussian distributions, the bounds are not as explicitly stated in terms of universal constants  
 1012 as they are for the standard Gaussian design setting (Theorem 15).  
 1013

1014 We now establish that for a wide family of random designs, the SRSC and SRSS properties are  
 1015 satisfied with high probability as well. For sake of simplicity, we will present our analysis for the  
 1016 standard Gaussian design. However, the results would readily extend to general Gaussian and sub-  
 1017 Gaussian designs using techniques similar to Theorem 17.  
 1018

1019 **Theorem 18.** Let  $X \in \mathbb{R}^{p \times n}$  be a matrix whose columns are sampled i.i.d from a standard Gaussian  
 1020 distribution i.e.  $\mathbf{x}_i \sim \mathcal{N}(\mathbf{0}, I)$ . Then for any  $\gamma > 0$  and  $s \leq p$ , with probability at least  $1 - \delta$ , the  
 1021  
 1022  
 1023  
 1024  
 1025

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

matrix  $X$  satisfies the SRSC and SRSS properties with constants

$$L_{(\gamma,s)}^{Gauss} \leq \gamma n \left( 1 + 3e \sqrt{6 \log \frac{e}{\gamma}} \right) + \tilde{O} \left( \sqrt{ns + n \log \frac{1}{\delta}} \right)$$

$$\alpha_{(\gamma,s)}^{Gauss} \geq n - (1 - \gamma)n \left( 1 + 3e \sqrt{6 \log \frac{e}{1 - \gamma}} \right) - \tilde{\Omega} \left( \sqrt{ns + n \log \frac{1}{\delta}} \right).$$

*Proof.* The proof of this theorem proceeds similarly to that of Theorem 15. Hence, we simply point out the main differences. First, we shall establish, that for any  $\epsilon > 0$ , with probability at least  $1 - \delta$ ,  $X$  satisfies the RSC and RSS properties at level  $s$  with the following constants

$$L_s \leq n + (1 - 2\epsilon)^{-1} \sqrt{bns + b'n \log \frac{2}{\delta}}$$

$$\alpha_s \geq n - (1 - 2\epsilon)^{-1} \sqrt{bns + b'n \log \frac{2}{\delta}},$$

where  $b = 24e^2 \log \frac{3ep}{\epsilon s}$  and  $b' = 24e^2$ . To do so we notice that the only change needed to be made would be in the application of the covering number argument. Instead of applying the union bound over an  $\epsilon$ -cover  $\mathcal{C}^{p-1}$  of  $S^{p-1}$ , we would only have to consider an  $\epsilon$ -cover  $\mathcal{C}_s^{p-1}$  of the set  $S_s^{p-1}$  of all  $s$ -sparse unit vectors in  $p$ -dimensions. A straightforward calculation shows us that

$$|\mathcal{C}_s^{p-1}| \leq \binom{p}{s} \left( 1 + \frac{2}{\epsilon} \right)^s \leq \left( \frac{3ep}{\epsilon s} \right)^s.$$

Thus, setting  $\mu^2 = b \cdot \frac{s}{n} + b' \cdot \frac{\log \frac{2}{\delta}}{n}$ , where  $b = 24e^2 \log \frac{3ep}{\epsilon s}$  and  $b' = 24e^2$ , we get

$$\mathbb{P} \left[ \sup_{\mathbf{v} \in \mathcal{C}_s^{p-1}} \left| \|X\mathbf{v}\|_2^2 - n \right| \geq \sqrt{bns + b'n \log \frac{2}{\delta}} \right] \leq \delta,$$

which establishes the required RSC and RSS constants for  $X$ . Now, moving on to the SRSS constant, it follows simply by applying a union bound over all sets in  $\mathcal{S}_\gamma$  much like in Theorem 15. One can then proceed to bound the SRSC constant in a similar manner.

We note that the nature of the SRSC and SRSS bounds indicate that our TORRENT-FC algorithm in the high dimensional sparse recovery setting has noise tolerance properties, characterized by the largest corruption index  $\alpha$  that can be tolerated, identical to its low dimensional counterpart - something that Theorem 9 states explicitly.  $\square$

## H Supplementary Results

**Claim 19.** Given any vector  $\mathbf{v} \in \mathbb{R}^n$ , let  $\sigma \in S_n$  be defined as the permutation that orders elements of  $\mathbf{v}$  in descending order of their magnitudes i.e.  $|v_{\sigma(1)}| \geq |v_{\sigma(2)}| \geq \dots \geq |v_{\sigma(n)}|$ . For any  $0 < p \leq q \leq 1$ , let  $S_1 \in \mathcal{S}_q$  be an arbitrary set of size  $q \cdot n$  and  $S_2 = \{\sigma(i) : n - p \cdot n + 1 \leq i \leq n\}$ . Then we have  $\|\mathbf{v}_{S_2}\|_2^2 \leq \frac{p}{q} \|\mathbf{v}_{S_1}\|_2^2 \leq \|\mathbf{v}_{S_1}\|_2^2$ .

*Proof.* Let  $S_3 = \{\sigma(i) : n - q \cdot n + 1 \leq i \leq n\}$  and  $S_4 = \{\sigma(i) : n - q \cdot n + 1 \leq i \leq n - p \cdot n\}$ . Clearly, we have  $\|\mathbf{v}_{S_3}\|_2^2 \leq \|\mathbf{v}_{S_1}\|_2^2$  since  $S_3$  contains the smallest  $q \cdot n$  elements (by magnitude). Now we have  $\|\mathbf{v}_{S_3}\|_2^2 = \|\mathbf{v}_{S_2}\|_2^2 + \|\mathbf{v}_{S_4}\|_2^2$ . Moreover, since each element of  $S_4$  is larger in magnitude than every element of  $S_2$ , we have

$$\frac{1}{|S_4|} \|\mathbf{v}_{S_4}\|_2^2 \geq \frac{1}{|S_2|} \|\mathbf{v}_{S_2}\|_2^2.$$

This gives us

$$\|\mathbf{v}_{S_2}\|_2^2 = \|\mathbf{v}_{S_3}\|_2^2 - \|\mathbf{v}_{S_4}\|_2^2 \leq \|\mathbf{v}_{S_3}\|_2^2 - \frac{|S_4|}{|S_2|} \|\mathbf{v}_{S_2}\|_2^2,$$

which upon simple manipulations, gives us the claimed result.  $\square$

**Lemma 20.** Let  $Z$  be distributed according to the chi-squared distribution with  $k$  degrees of freedom i.e.  $Z \sim \chi^2(k)$ . Then for all  $t \geq 0$ ,

$$\mathbb{P}[|Z - k| \geq t] \leq 2 \exp\left(-\min\left\{\frac{t^2}{24ke^2}, \frac{t}{4\sqrt{3}e}\right\}\right)$$

*Proof.* This lemma requires a proof structure that traces several basic results in concentration inequalities for sub-exponential variables [17, Lemma 5.5, 5.15, Proposition 5.17]. The purpose of performing this exercise is to explicate the constants involved so that a crisp bound can be provided on the corruption index that our algorithm can tolerate in the standard Gaussian design case.

We first begin by establishing the sub-exponential norm of a chi-squared random variable with a single degree of freedom. Let  $X \sim \chi^2(1)$ . Then using standard results on the moments of the standard normal distribution gives us, for all  $p \geq 2$ ,

$$(\mathbb{E}|X|^p)^{1/p} = ((2p-1)!!)^{1/p} = \left(\frac{(2p)!}{2^p p!}\right)^{1/p} \leq \frac{\sqrt{3}}{2} p$$

Thus, the sub-exponential norm of  $X$  is upper bounded by  $\sqrt{3}/2$ . By applying the triangle inequality, we obtain, as a corollary, an upper bound on the sub-exponential norm of the centered random variable  $Y = X - 1$  as  $\|Y\|_{\psi_1} \leq 2\|X\|_{\psi_1} \leq \sqrt{3}$ .

Now we bound the moment generating function of the random variable  $Y$ . Noting that  $\mathbb{E}Y = 0$ , we have, for any  $|\lambda| \leq \frac{1}{2\sqrt{3}e}$ ,

$$\mathbb{E} \exp(\lambda Y) = 1 + \sum_{q=2}^{\infty} \frac{\mathbb{E}(\lambda Y)^q}{q!} \leq 1 + \sum_{q=2}^{\infty} \frac{(\sqrt{3}|\lambda|q)^q}{q!} \leq 1 + \sum_{q=2}^{\infty} (\sqrt{3}e|\lambda|)^q \leq 1 + 6e^2\lambda^2 \leq \exp(6e^2\lambda^2).$$

Note that the second step uses the sub-exponentiality of  $Y$ , the third step uses the fact that  $q! \geq (q/e)^q$ , and the fourth step uses the bound on  $|\lambda|$ . Now let  $X_1, X_2, \dots, X_k$  be  $k$  independent random variables distributed as  $\chi^2(1)$ . Then we have  $Z \sim \sum_{i=1}^k X_i$ . Using the exponential Markov's inequality, and the independence of the random variables  $X_i$  gives us

$$\mathbb{P}[Z - k \geq t] = \mathbb{P}\left[e^{\lambda(Z-k)} \geq e^{\lambda t}\right] \leq e^{-\lambda t} \mathbb{E} e^{\lambda(Z-k)} = e^{-\lambda t} \prod_{i=1}^k \mathbb{E} \exp(\lambda(X_i - 1)).$$

For any  $|\lambda| \leq \frac{1}{2\sqrt{3}e}$ , the above bounds on the moment generating function give us

$$\mathbb{P}[Z - k \geq t] \leq e^{-\lambda t} \prod_{i=1}^k \exp(6e^2\lambda^2) = \exp(-\lambda t + 6ke^2\lambda^2).$$

Choosing  $\lambda = \min\left\{\frac{1}{2\sqrt{3}e}, \frac{t}{12ke^2}\right\}$ , we get

$$\mathbb{P}[Z - k \geq t] \leq \exp\left(-\min\left\{\frac{t^2}{24ke^2}, \frac{t}{4\sqrt{3}e}\right\}\right).$$

Repeating this argument gives us the same bound for  $\mathbb{P}[k - Z \geq t]$ . This completes the proof.  $\square$

1134 **I Supplementary Experimental Results**

1135  
1136  
1137  
1138  
1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187

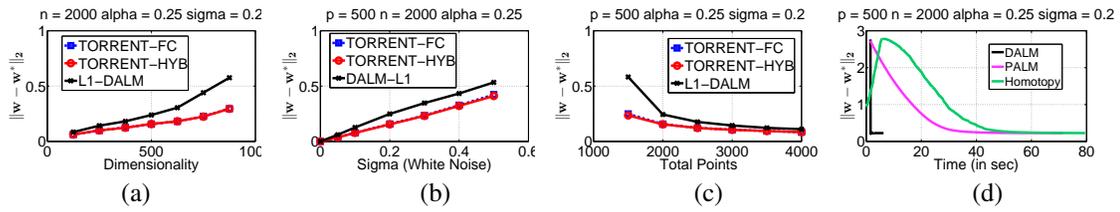


Figure 3: (a), (b), (c) Variation of recovery error with varying  $p, \sigma$  and  $n$ . TORRENT was found to outperform DALM- $L_1$  in all these settings. (d) Recovery error as a function of runtime for various state-of-the-art  $L_1$  solvers as indicated in [14].