# ON THE SPARSITY OF LASSO MINIMIZERS IN SPARSE DATA RECOVERY

SIMON FOUCART, EITAN TADMOR, AND MING ZHONG

# Dedicated to Ron DeVore with friendship and admiration

ABSTRACT. We present a detailed analysis of the unconstrained  $\ell_1$ -weighted LASSO method for recovery of sparse data from its observation by randomly generated matrices, satisfying the Restricted Isometry Property (RIP) with constant  $\delta < 1$ , and subject to negligible measurement and compressibility errors. We prove that if the data is k-sparse, then the size of support of the LASSO minimizer, s, maintains a comparable sparsity,  $s \leqslant C_\delta k$ . For example, if  $\delta = 0.7$  then s < 11k and a slightly smaller  $\delta = 0.4$  yields s < 4k. We also derive new  $\ell_2/\ell_1$  error bounds which highlight precise dependence on k and on the LASSO parameter k, before the error is driven below the scale of negligible measurement/ and compressibility errors.

#### **CONTENTS**

1. Introduction	1
1.1. Statement of main results	2
2. The Robust Null Space Property	4
3. On the sparsity of the unconstrained LASSO minimizer	5
3.1. Bounds of the sparsity	6
3.2. Numerical simulations	8
4. Error bounds	10
4.1. $\ell_2$ -error bounds	10
4.2. $\ell_1$ -error bound	13
4.3. Numerical simulations	15
References	16

## 1. Introduction

In 2006, the pioneering works of Candès, Romberg and Tao [13, 14] and of Donoho [21] suggested the framework of a constrained  $\ell_1$ -method to recover a sparse unknown  $\boldsymbol{x}_* \in \mathbb{R}^N$  from its observation  $\boldsymbol{y}_* = A\boldsymbol{x}_* \in \mathbb{R}^{m1}$ . The key point is that one can design observing matrices  $A \in \mathbb{R}^{m \times N}$ 

Date: March 16, 2022.

<sup>2020</sup> Mathematics Subject Classification. 94A12, 94A20.

Key words and phrases. Inverse problems, data recovery, compressive sensing, basis pursuit de-noising method, robust null space property.

Research was supported in part by NSF grants CCF-1934904, DMS-2053172 (SF), and DMS16-13911 (ET) and ONR grants N00014-2012787 (SF) and N00014-2112773 (ET).

<sup>&</sup>lt;sup>1</sup>Earlier announcement of the works was presented in the workshops, organized together with Ron DeVore, which can be found at home.cscamm.umd.edu/programs/srs05/candes\_srs05.pdf and home.cscamm.umd.edu/programs/srs05/donoho\_srs05.htm.

with a relatively small number of observations,  $m \ll N$ , such that a constrained  $\ell_1$ -method — also known as Basis Pursuit (BP) in [18, 16, 17] — finds a sparse solution as a minimizer of

$$oldsymbol{x}_{BP}\coloneqq rg\min_{oldsymbol{x}\in\mathbb{R}^N} ig\{|oldsymbol{x}|_1 \ ig| \ Aoldsymbol{x}=oldsymbol{y}_*\}, \qquad A\in\mathbb{R}^{m imes N}, \ m\ll N.$$

This is closely related to the well-known LASSO algorithm introduced in 1996 in the statistics literature [39],  $\underset{|\boldsymbol{x}|_1 \leq \delta}{\arg\min} \left\{ |\boldsymbol{y}_*^{\epsilon} - A\boldsymbol{x}|_2^2 \right\}$ , which can be viewed as an  $\ell_1$ -penalty relaxation of a least squares subject to (possibly noisy) observation  $\boldsymbol{y}_*^{\epsilon}$ .

The BP minimizer,  $x_{BP}$ , recovers the sparse  $x_*$  when the observing matrix A satisfies an appropriate recoverability condition; we mention here the Restricted Isometry Property (RIP) introduced in [13], the  $\ell_1$ -Coherence discussed in [40, 27, 22, 23], the restricted eigenvalue condition [6, §3], or the Null Space Property (NSP) of DeVore and his co-authors [19, 20], and related Robust Null Sparse Property (RNSP) of [26]. Important classes of such observing matrices with desired sparse recoverability conditions are randomly generated, e.g., [26, §9].

1.1. **Statement of main results.** Throughout the paper we will be using the two notions of sparsity and compressibility. A vector  $x \in \mathbb{R}^N$  is *sparse* if

$$s_{x} := |x|_{0} \ll N.$$

In applications, sparsity is often difficult to acquire, and clean observations are not always available, since the observation process is inevitably and easily corrupted by errors — human and/or machine measurement errors. We turn our attention to the recovery of compressible unknown from its *noisy* observations. A vector  $\boldsymbol{x} \in \mathbb{R}^N$  is *compressible* of order k, or simply k-compressible, if its content is faithfully captured by a k-sparse vector — specifically, if its  $\ell_1$ -distance to the set of all k-sparse vectors,

(1.1) 
$$\sigma_k(\boldsymbol{x}) := \inf_{\boldsymbol{z} \in \mathbb{R}^N} \left\{ |\boldsymbol{x} - \boldsymbol{z}|_1 : |\boldsymbol{z}|_0 \leqslant k \right\},$$

is small relative to  $|x|_1$ . We note that  $\sigma_k(x)$  is realized by a (not necessarily unique) vector, denoted x(k), whose non-zero entries are the k largest of x in absolute value.

Let  $x_*$  be a compressible unknown of order k so that  $\sigma_k(x_*) \ll |x_*|_1$ , and assume we only have access to its measured observation  $y_*^\epsilon = Ax_* + \varepsilon$ . The term  $\varepsilon$  is the measurement error caused by a number of factors which are assumed statistically independent of the unknown  $x_*$  and the observing operator A. The details of  $\varepsilon$  remain untraceable except for its size which is assumed to be negligibly small. In this case, one should not expect an exact recovery of a sparse  $x_*$ , but instead, accept an approximate solution,  $y_*^\epsilon = Ax_*(k) + \varepsilon'$ , where  $\varepsilon' = A(x_* - x_*(k)) + \varepsilon$  is adapted to the small scale built into the problem, which consists of two contributions — the small measurement error,  $\epsilon := |\varepsilon|_2$ , and the small compressibility error,  $|x_*|_2 = |x_*|_2 = |x_*|$ 

$$\boldsymbol{y}_*^{\epsilon} = A\boldsymbol{x}_*(k) + \boldsymbol{\varepsilon}', \qquad |\boldsymbol{\varepsilon}'|_2 \leqslant \mu,$$

such that  $\mu = \sigma_k(x_*) + \epsilon$  is much smaller relative to the unknown data,  $\mu \ll |x_*|_1$ . Although the observing operator A is linear, the recovery of  $x_*$  by a direct "solution" of the linear problem  $Ax = y_*^{\epsilon}$  is ill-posed, unless additional conditions on A and  $x_*$  are enforced so that the unknown

Given  $x \in \mathbb{R}^N$  we let  $|x|_p$  denote its  $\ell_p$ -norm, with the usual conventional limiting cases of  $p = \infty$  and p = 0, where  $|x|_{\infty} := \max_{1 \leqslant i \leqslant N} |x_i|$ , and respectively  $|x|_0 := |\operatorname{supp}(x)|$  where  $|\cdot|$  is the cardinality of a finite set.

<sup>&</sup>lt;sup>3</sup>The columns of A are assumed  $\ell_2$ -normalized so that  $|A|_{1\to 2}=1$ .

object  $x_*$ , or at least a faithful approximation of it, is recovered by solving an augmented *well-posed* regularized minimization problem. On the way, the original linear problem is replaced by a nonlinear procedure. To capture the compressible information of  $x_*$  from its noisy observation  $y_*^{\epsilon}$ , we seek a minimizer of the unconstrained  $\ell_1$ -regularized Least Squares problem,

(1.2) 
$$\boldsymbol{x}_{\lambda} \coloneqq \underset{\boldsymbol{x} \in \mathbb{R}^{N}}{\min} \left\{ \lambda |\boldsymbol{x}|_{1} + \frac{1}{2} |\boldsymbol{y}_{*}^{\epsilon} - A\boldsymbol{x}|_{2}^{2} \right\}, \qquad A \in \mathbb{R}^{m \times N}, \ m \ll N.$$

The unconstrained variational statement (1.2) falls under the general class of Tikhonov regularization. The distinctive feature is the  $\ell_1$ -regularization, leading to an approximate decomposition of the basis pursuit of Chen & Donoho [18],  $\boldsymbol{y}_*^{\epsilon} = A\boldsymbol{x}_{\lambda} + \boldsymbol{r}_{\lambda}$  with (hopefully) small residual,  $\boldsymbol{r}_{\lambda} = \boldsymbol{y}_*^{\epsilon} - A\boldsymbol{x}_{\lambda}$ , depending on a parameter  $\lambda$ . This version of  $\ell_1$ -regularization, called "Basis Pursuit De-Noising" in [16], which became known as the unconstrained  $\ell_1$ -weighted LASSO, is the main focus of our work. As noted in the 1996 thesis [16], the work on this version of BP was motivated by a series of ideas using  $\ell_0/\ell_1$ -based regularization that appeared in early 1990s, primarily the empirical atomic decomposition of Donoho and Johnstone [24],  $\underset{\boldsymbol{x} \in \mathbb{R}^N}{\operatorname{arg min}} \left\{ \lambda |\boldsymbol{x}|_0 + \frac{1}{2} |\boldsymbol{y}_*^{\epsilon} - A\boldsymbol{x}|_2^2 \right\}$ , the multi-scale edge representation with wavelets of Hwang and Mallat [29] and the TV-based denoising method of ROF [31],  $\underset{\boldsymbol{x} \in \mathbb{R}^N}{\operatorname{arg min}} \left\{ \lambda |\boldsymbol{x}|_{TV} + \frac{1}{2} |\boldsymbol{y}_*^{\epsilon} - A\boldsymbol{x}|_2^2 \right\}$ . These works were later further explored as the Lagrangian formulation of the quadratically constrained Basis Pursuit de-noising [17, 14] and the noise-aware  $\ell_1$ -minimization [26].

Since  $\lambda > 0$  controls the distance between  $Ax_{\lambda}$  and  $y_{*}$ , the parameter  $\lambda$  can be interpreted as a regularization *scale*. In a subsequent work, [37], we pursue a *multi-scale* generalization based on a ladder of hierarchical scales constructed by the Hierarchical Decomposition (HD) method [34, 35, 36, 33]. The goal of this work is to analyze the sparsity behavior of the *mono-scale* LASSO (1.2), observed by a sub-class of RIP matrices satisfying the Robust Null Space Property (RNSP) which is discussed in section 2. Our main results, outlined and proved in section 3, are summarized in the following. Our results involve three main parameters: the Restricted Isometry

Constant (RIC) in (2.2) below,  $\delta = \delta_k < 1$ , the related RNSP constant,  $\beta_{\delta} = \frac{\sqrt{1+\delta}}{\sqrt{1-\delta^2}-\delta/4}$ , depending on the RIC  $\delta$ , and the small scale of compressibility+measurement,  $\mu = \sigma_k(\boldsymbol{x}_*) + \epsilon$ .

**Theorem 1.1** (Main result). Let  $x_*$  be k-compressible, and let  $y_*^{\epsilon} = Ax_* + \varepsilon$  be its observation with observing matrix A satisfying the RIP (2.2) with constant  $\delta$  large enough,  $\delta > \delta_t$ , such that (3.7) below holds. Let  $x_{\lambda}$  be the LASSO minimizer (1.2).

(i) (**Sparsity**). The sparsity of the LASSO minimizer,  $s_{\lambda} = s_{x_{\lambda}}$ , does not exceed

$$s_{\lambda} < (1+\delta) \Big( \beta_{\delta} \sqrt{k} + \frac{2\mu}{\lambda} \Big)^2.$$

(ii) ( $\ell_2$ -error bound). The following  $\ell_2$ -error bound holds

$$\frac{1}{\sqrt{1+\delta}} \left( \frac{\sqrt{s_{\lambda}}\lambda}{\sqrt{1+\delta}} - \mu \right) \leqslant |\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{2} \leqslant \frac{1}{\sqrt{1-\delta}} \left( \beta_{\delta} \sqrt{k}\lambda + 3\mu \right).$$

(iii) ( $\ell_1$ -error bound). The following  $\ell_1$ -error bound holds

$$|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{1} < \frac{\sqrt{1+\delta}}{\sqrt{1-\delta}} \frac{1}{\lambda} ((\beta_{\delta} + 1/2)\sqrt{k}\lambda + 2\mu)^{2}.$$

We interpret these bounds as follows. Set  $\theta = 2\mu/\sqrt{k}\lambda$ , then (i) reads

$$s_{\lambda} \leqslant \chi^2 k, \qquad \chi = \sqrt{1+\delta}(\beta_{\delta} + \theta).$$

Thus, if  $\theta \leqslant 1$  — namely, as long as  $\lambda$  does not get exceedingly small so that  $\lambda \geqslant 2\mu/\sqrt{k}$ , then the sparsity of  $\boldsymbol{x}_{\lambda}$  is comparable to the sparsity of  $\boldsymbol{x}_{*}$ . Furthermore, in (ii) we have the  $\ell_{2}$ -error bound of order  $\lesssim \sqrt{k}\lambda + \mu$  and in (iii), an  $\ell_{1}$ -error bound of order  $\lesssim k\lambda + \sqrt{k}\mu$ .

We conclude with a few comments on theorem 1.1. The sparsity bound in (i) with RIC  $\delta=0.7$  yields  $s_{\lambda}<11k$ , while a slightly smaller RIC  $\delta=0.4$  yields  $s_{\lambda}<4k$ . This should be compared with the sparsity bounds in [5, Theorem 3] and [32]. The  $\ell_2$ -upper bound on the right of (ii) is not new; here we recover the  $\ell_2$ -bound, derived under appropriate assumptions, in [11], [6, Theorem 7.1] and [30, 38, 28]. This should be contrasted with the  $\ell_2$ -error lower-bound on the left, derived in section 4.1 (see figure 4.1). Indeed, this  $\ell_2$  lower-bound is the essential ingredient in our proof of the sparsity bound in (i). Finally, the  $\ell_1$ -bound in (iii) with RIC  $\delta=0.7$  yields  $|x_{\lambda}-x_*(k)|_1<16.97k\lambda+24.04\sqrt{k\mu}$ . Here, the linear decay with  $\lambda$  is not new and can be found for example, under various assumptions, in [6, Theorem 7.1] and [9, Theorem 6.1].

### 2. THE ROBUST NULL SPACE PROPERTY

Optimality of the minimizer. The variational problem (1.2) admits a minimizer,  $x_{\lambda}$ , and at least for certain relevant classes of full row rank A's, the minimizer is unique, [41]. The minimizer is completely characterized by its residual,  $r_{\lambda} := y_*^{\epsilon} - Ax_{\lambda}$  (to simplify notations we suppress the dependence of  $r_{\lambda}$  on  $\epsilon$ ). We summarize the results from [35, §2.1],[33, Appendix] where we distinguish between two cases.

- (i) If  $\lambda \geqslant \lambda_{\infty} := |A^{\top} y_{*}^{\epsilon}|_{\infty}$  then (1.2) admits only the trivial minimizer  $x_{\lambda} \equiv 0$ . In this case,  $\lambda$  is too large to extract the compressibility information in  $y_{*}^{\epsilon}$ .
- (ii) If  $\lambda < \lambda_{\infty} = |A^{\top} \boldsymbol{y}_{*}^{\epsilon}|_{\infty}$  then (1.2) admits a non-trivial minimizer,  $\boldsymbol{x}_{\lambda}$ , with the corresponding residual,  $\boldsymbol{r}_{\lambda} = \boldsymbol{y}_{*}^{\epsilon} A\boldsymbol{x}_{\lambda}$ , such that  $(\boldsymbol{x}_{\lambda}, \boldsymbol{r}_{\lambda})$  forms an *extremal pair* in the sense that

(2.1) 
$$\langle A\boldsymbol{x}_{\lambda}, \boldsymbol{r}_{\lambda} \rangle = \lambda |\boldsymbol{x}_{\lambda}|_{1} \text{ and } |A^{\top}\boldsymbol{r}_{\lambda}|_{\infty} = \lambda.$$

To proceed we will need the following notations. The restriction of a vector  $\boldsymbol{w} \in \mathbb{R}^N$  on an index set  $\mathcal{K} \subset \{1,2,\ldots,N\}$  of size  $k=|\mathcal{K}|$  is denoted  $\boldsymbol{w}_{\mathcal{K}}:=\{w_i,\ i\in\mathcal{K}\}\in\mathbb{R}^k$ . Similarly, given a matrix  $W\in\mathbb{R}^{m\times N}$  with columns  $\boldsymbol{w}_1,\boldsymbol{w}_2,\ldots$ , its restriction on an index set  $\mathcal{K}$  of size  $k=|\mathcal{K}|$  consists of the k columns  $W_{\mathcal{K}}:=\operatorname{col}\{\boldsymbol{w}_i,i\in\mathcal{K}\}$ . The size of W can be measured by its induced matrix norm,  $\|W\|_p=\sup_{|\boldsymbol{w}|_p=1}|W\boldsymbol{w}|_p$ . The signum vector is defined component-wise,

$$\operatorname{sgn}(\boldsymbol{w})_i = \operatorname{sgn}(w_i)$$
, in terms of the usual signum function  $\operatorname{sgn}(w) = \left\{ \begin{array}{cc} -1, & w < 0 \\ 1, & w > 0 \end{array} \right\}$  for  $w \neq 0$ .

**Restricted Isometry Poperty (RIP)**. A matrix A satisfies the Restricted Isometry Property (RIP) of order k with Restricted Isometry Constant (RIC)  $\delta_k < 1$  if the following holds, [15, 21, 13, 7],

(2.2) 
$$(1 - \delta_k) |\mathbf{x}|_2^2 \leqslant |A\mathbf{x}|_2^2 \leqslant (1 + \delta_k) |\mathbf{x}|_2^2, \qquad \forall |\mathbf{x}|_0 \leqslant k.$$

Throughout the paper we adopt the usual assumption that  $\delta_k$  is measured for A's with  $\ell^2$ -normalized  $columns^4$ . There are two classes of matrices  $A \in \mathbb{R}^{m \times N}$  satisfying the RIP of order k: deterministic

<sup>&</sup>lt;sup>4</sup>The RIP of A asserts that for any subset of its k columns,  $\{a_i\}_{i\in\mathcal{K}}$ , the entries  $|\langle a_i,a_j\rangle|_{i\neq j}\lesssim \delta_k$  while  $|a_i|_2^2=1+\epsilon_i$  such that  $|\epsilon_i|\lesssim \delta_k$ . Therefore, one can always re-normalize the columns of A by a factor  $\lesssim (1-\delta_k)^{-1/2}$  yielding a new RIP matrix with  $\ell_2$ -normalized columns and with possibly slightly larger RIP constant  $\delta_k'\lesssim \delta_k/(1-\delta_k)$ .

A's with number of observations  $m \gtrsim k^2$  (the quadratic bottleneck is lessened in [8]); and a large class of randomly generated A's for which the restriction on the number of observations can be further lessened to having only m observations, [26, §9.4]

(2.3) 
$$m \sim Const \cdot \delta^{-2} k \ln \left( \frac{eN}{k} \right).$$

Candès proved the exactness of the constrained BP for RIP matrices with  $\delta < \sqrt{2} - 1$ , [12]. Further refinements were reported in [25] before the definitve result of [10].

**Robust Null Space Property (RNSP)**. A crucial step in quantifying the recovery error of  $x_*$  using (1.2) is to enforce a recoverability condition on the observing matrix A. This brings us to the Robust Null Sparse Property (RNSP) introduced in [26, §4.3]. A matrix  $A \in \mathbb{R}^{m \times N}$  satisfies the RNSP of order k with constants  $0 < \rho < 1$  and  $\tau > 0$ , if for all  $\mathcal{K} \subset \{1, 2, \dots, N\}$  of size  $|\mathcal{K}| \leq k$ , there holds

$$|\boldsymbol{x}_{\mathcal{K}}|_{1} \leqslant \rho |\boldsymbol{x}_{\mathcal{K}^{c}}|_{1} + \tau |A\boldsymbol{x}|_{2}, \qquad \forall \boldsymbol{x} \in \mathbb{R}^{N}.$$

We refer to the "RNSP $_{\rho,\tau}$  of order k", and unless needed, we suppress the dependence of  $(\rho,\tau)$  on k. In particular, given a k-sparse  $\boldsymbol{v}$  and any  $\boldsymbol{u}$ , we apply (2.4) to  $\boldsymbol{x} = \boldsymbol{u} - \boldsymbol{v}$  with  $\mathcal{K} = \operatorname{supp}(\boldsymbol{v})$ , where  $|\boldsymbol{x}_{\mathcal{K}}|_1 - |\boldsymbol{x}_{\mathcal{K}^c}|_1 \geqslant |\boldsymbol{v}|_1 - |\boldsymbol{u}|_1$  yields the following useful consequence of RNSP.

**Lemma 2.1.** If  $A \in \mathbb{R}^{m \times N}$  satisfies the RNSP<sub> $\rho,\tau$ </sub> of order k, then for all k-sparse v's and any u,

$$|\boldsymbol{v}|_1 - |\boldsymbol{u}|_1 \leqslant \tau |A(\boldsymbol{u} - \boldsymbol{v})|_2, \qquad |\operatorname{supp}(\boldsymbol{v})| \leqslant k.$$

As an example for the class of observation matrices satisfying the RNSP<sub> $\rho,\tau$ </sub> of order k, we mention the class of randomly generated RIP matrices with RICs  $\delta = \delta_{2k}$ , [26, Theorem 6.13],

(2.6) 
$$\rho = \frac{\delta}{\sqrt{1 - \delta^2 - \delta/4}} \quad \text{and} \quad \tau = \beta \sqrt{k}, \quad \beta := \frac{\sqrt{1 + \delta}}{\sqrt{1 - \delta^2 - \delta/4}}, \qquad \delta = \delta_{2k}.$$

These RNSP parameters,  $(\rho, \beta)$ , are dictated as increasing functions of the RIC  $\delta < 1$ . A smaller  $\delta$  requires an increased number of observations. All proofs invoke different classes of observing matrices which are randomly generated so that they satisfy a desirable observing properties—RIP, RNSP, or Constrained Minimal Singular Values (CMSV) property. Accordingly, the error statements are probabilistic in nature, referring to the ensemble of these observations.

### 3. On the sparsity of the unconstrained LASSO minimizer

We analyze the sparsity and  $\ell_1/\ell_2$ -error bounds of the minimizer (1.2) in recovering  $\boldsymbol{x}_*(k)$  from the observation  $\boldsymbol{y}_*^{\epsilon} = A\boldsymbol{x}_* + \boldsymbol{\varepsilon}$ , with small measurement error,  $\epsilon = |\boldsymbol{\varepsilon}|_2$ , and — assuming that  $\boldsymbol{x}_*$  is k-compressible — with small compressibility error,  $\sigma_k(\boldsymbol{x}_*) = |\boldsymbol{x}_* - \boldsymbol{x}_*(k)|_1$ . Set

$$\mu := \sigma_k(\boldsymbol{x}_*) + \epsilon.$$

Clearly, since the exact solution is observed up to  $\ell_2$  residual error of order  $|\boldsymbol{y}_*^{\epsilon} - A\boldsymbol{x}_*|_2 \leqslant \mu$ , we do not have much to say when the computed residual error  $|\boldsymbol{r}_{\lambda}|_2$  is of order  $\mu$  and we will therefore limit ourselves to the parametric regime where  $|\boldsymbol{r}_{\lambda}|_2 \gg \mu$ . Below we show that  $|\boldsymbol{r}_{\lambda}|_2 \sim \lambda \sqrt{k}$  and therefore throughout the paper we make the assumption

(3.1) 
$$\theta := \frac{2\mu}{\lambda\sqrt{k}} \leqslant 1, \qquad \mu = \sigma_k(\boldsymbol{x}_*) + \epsilon.$$

Thus, we assume the LASSO weight,  $\lambda$ , does not get exceedingly small,  $\lambda \geqslant 2\mu/\sqrt{k}$ . In concrete examples demonstrating the sparsity and error bounds reported below we use  $\theta = 0.1$ , corresponding to  $\lambda \geqslant 20\mu/\sqrt{k}$ .

**Lemma 3.1** (The re-scaled residual — an upper-bound). Fix  $\lambda < \lambda_{\infty} := |A^{\top} \boldsymbol{y}_{*}^{\epsilon}|_{\infty}$ . Let  $\boldsymbol{y}_{*}^{\epsilon} = A\boldsymbol{x}_{*} + \varepsilon$  be the observation of a k-compressible unknown  $\boldsymbol{x}_{*} \in \mathbb{R}^{N}$ , observed by  $A \in \mathbb{R}^{m \times N}$  satisfying the  $RNSP_{\rho,\tau}$  of order k, (2.6). Let  $\mu$  denote the small scale of k-compressiblity and measurement errors, see (3.1). Then the residual of the LASSO (1.2),  $\boldsymbol{r}_{\lambda} = \boldsymbol{y}_{*}^{\epsilon} - A\boldsymbol{x}_{\lambda}$ , satisfies

(3.2) 
$$\frac{|\mathbf{r}_{\lambda}|_{2}}{\lambda} \leqslant (\beta_{\delta} + \theta)\sqrt{k}, \qquad \beta_{\delta} = \frac{\sqrt{1+\delta}}{\sqrt{1-\delta^{2}} - \delta/4}.$$

*Proof.* Clearly,  $|A(\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k))|_{2} \leq |\boldsymbol{r}_{\lambda}|_{2} + \mu$ . Using (2.5) with the k-sparse  $\boldsymbol{v} = \boldsymbol{x}_{*}(k)$  and  $\boldsymbol{u} = \boldsymbol{x}_{\lambda}$  yields

(3.3) 
$$|x_*(k)|_1 - |x_\lambda|_1 \leqslant \tau |A(x_\lambda - x_*(k))|_2 \leqslant \tau |r_\lambda|_2 + \tau \mu.$$

Next, a lower-bound for the quantity on the left follows. Recall that  $x_{\lambda}$ , being the LASSO minimizer (1.2), is characterized by the extremal property that its scaled residual  $z = \frac{r_{\lambda}}{\lambda}$  satisfies (2.1),

Hence

$$|\boldsymbol{x}_{*}(k)|_{1} - |\boldsymbol{x}_{\lambda}|_{1} \geqslant \langle \boldsymbol{x}_{*}(k), A^{\top} \boldsymbol{z} \rangle - \langle A \boldsymbol{x}_{\lambda}, \boldsymbol{z} \rangle = \langle A \boldsymbol{x}_{*}(k) - A \boldsymbol{x}_{\lambda}, \boldsymbol{z} \rangle$$

$$= \langle \boldsymbol{r}_{\lambda}, \boldsymbol{z} \rangle + \langle A \boldsymbol{x}_{*}(k) - \boldsymbol{y}_{*}^{\epsilon}, \boldsymbol{z} \rangle$$

$$\geqslant \frac{|\boldsymbol{r}_{\lambda}|_{2}^{2}}{\lambda} - |A \boldsymbol{x}_{*}(k) - \boldsymbol{y}_{*}^{\epsilon}|_{2} \frac{|\boldsymbol{r}_{\lambda}|_{2}}{\lambda}.$$

Now, assumption (3.1) and the fact that  $|A|_{1\to 2} \le 1$  imply,

$$|Ax_*(k) - y_*^{\epsilon}|_2 \le |Ax_* - y_*^{\epsilon}|_2 + |A(x_*(k) - x_*)|_2 \le \epsilon + \sigma_k(x_*) = \mu$$

and we end with the desired lower-bound

$$|\boldsymbol{x}_{*}(k)|_{1} - |\boldsymbol{x}_{\lambda}|_{1} \geqslant \frac{|\boldsymbol{r}_{\lambda}|_{2}^{2}}{\lambda} - \mu \frac{|\boldsymbol{r}_{\lambda}|_{2}}{\lambda}.$$

Combining (3.3) and (3.4) we conclude that  $|z|_2 = \frac{|r_{\lambda}|_2}{\lambda}$  satisfies the quadratic inequality,  $|z|_2^2 \le \left(\tau + \frac{\mu}{\lambda}\right)|z|_2 + \tau \frac{\mu}{\lambda}$ , and therefore

(3.5) 
$$\frac{|\boldsymbol{r}_{\lambda}|_{2}}{\lambda} = |\boldsymbol{z}|_{2} < \tau + \frac{2\mu}{\lambda} = (\beta_{\delta} + \theta)\sqrt{k},$$
 proving (3.2).

3.1. **Bounds of the sparsity.** We now come to the main point of the lower-bound on the scaled residual in terms of the size of the support of  $\boldsymbol{x}_{\lambda}$ ,  $\frac{|\boldsymbol{r}_{\lambda}|_2}{\lambda} \gtrsim \sqrt{s_{\lambda}}$ . Fix  $\lambda < \lambda_{\infty} := |A^{\top}\boldsymbol{y}_{*}^{\epsilon}|_{\infty}$ . Recall that if  $\boldsymbol{x}_{\lambda}$  is the LASSO minimizer (1.2) then by the extremal property (2.1), the scaled residual  $\boldsymbol{z} = \frac{\boldsymbol{r}_{\lambda}}{\lambda}$  satisfies the two properties  $\langle A\boldsymbol{x}_{\lambda},\boldsymbol{z}\rangle = |\boldsymbol{x}_{\lambda}|_{1}$  and  $|A^{\top}\boldsymbol{z}|_{\infty} = 1$ . Thus, the extremal  $\boldsymbol{x}_{\lambda}$  with support  $\mathcal{S} = \sup(\boldsymbol{x}_{\lambda})$  of size  $s_{\lambda} = |\boldsymbol{x}_{\lambda}|_{0}$ , is identified by a re-scaled residual satisfying

(3.6) 
$$(A^{\top} \boldsymbol{z})_{\mathcal{S}} = \mathbf{sgn}(\boldsymbol{x}_{\lambda,\mathcal{S}}), \qquad \boldsymbol{z} = \frac{\boldsymbol{r}_{\lambda}}{\lambda}, \quad \mathcal{S} = \mathrm{supp}(\boldsymbol{x}_{\lambda}).$$

Fix the integer t,

(3.7) 
$$t := [(1+\delta)(\beta_{\delta} + \theta)^{2}k] + 1 \text{ with constant } \delta > \delta_{t}.$$

Since the RIC  $\delta_t$  is increasing with the order t, there is no need to trace a precise fixed point associated with (3.7),  $t = [(1 + \delta_t)(\beta_{\delta_t} + \theta)^2 k] + 1$ . Instead, we can use a priori bounds of  $\delta_t$ ; for example, if we restrict ourselves to the range  $\delta < 0.7$ , we can set the integer upper bound t = 11k. Below, we demonstrate refined versions of this bound.

We claim that

(3.8) 
$$s_{\lambda} < t = [(1+\delta)(\beta_{\delta} + \theta)^{2}k] + 1.$$

To this end we proceed by contradiction. Assume  $s_{\lambda} \geqslant t$ . Then the support of  $\boldsymbol{x}_{\lambda}$  has a subset  $\mathcal{T}$  of size t for which the extremal property (3.6) reads  $(A^{\top}\boldsymbol{z})_{\mathcal{T}} = \operatorname{sgn}(\boldsymbol{x}_{\lambda,\mathcal{T}})$ , and the RIP (2.2) for such set  $\mathcal{T}$  implies

$$(3.9) |A(A^{\top} z)_{\mathcal{T}}|_{2}^{2} \leqslant (1 + \delta_{t})|(A^{\top} z)_{\mathcal{T}}|_{2}^{2} = (1 + \delta_{t})|\operatorname{sgn}(x_{\lambda,\mathcal{T}})|_{2}^{2} = (1 + \delta_{t})t.$$

On the other hand, we have

$$|A(A^{\top}\boldsymbol{z})_{\mathcal{T}}|_{2}^{2} \geqslant \frac{1}{|\boldsymbol{z}|_{2}^{2}} \langle A(A^{\top}\boldsymbol{z})_{\mathcal{T}}, \boldsymbol{z} \rangle^{2} = \frac{1}{|\boldsymbol{z}|_{2}^{2}} \langle (A^{\top}\boldsymbol{z})_{\mathcal{T}}, A^{\top}\boldsymbol{z} \rangle^{2} = \frac{1}{|\boldsymbol{z}|_{2}^{2}} |(A^{\top}\boldsymbol{z})_{\mathcal{T}}|_{2}^{4} = \frac{t^{2}}{|\boldsymbol{z}|_{2}^{2}}.$$

The last two inequalities followed by Lemma 3.1 imply  $t \leq (1 + \delta_t)|z|_2^2 < (1 + \delta)(\beta_\delta + \theta)^2 k$ , which contradicts the definition of t,

$$t = [(1+\delta)(\beta_{\delta} + \theta)^{2}k] + 1 \geqslant (1+\delta_{t})(\beta_{\delta} + \theta)^{2}k.$$

Thus, (3.8) holds.

In fact, a refined statement follows. Now that we know  $|\mathcal{S}| \leq [(1+\delta)(\beta_{\delta}+\theta)^2k]$  we can argue along the same line as above with  $\mathcal{T} = \mathcal{S}$ , obtaining  $s_{\lambda} \leq (1+\delta)|z|_2^2$ .

Lemma 3.2 (The re-scaled residual — a lower-bound). Fix  $\lambda < \lambda_{\infty} := |A^{\top} y_{*}^{\epsilon}|_{\infty}$  and let  $x_{\lambda}$  be the  $s_{\lambda}$ -sparse minimizer of the corresponding LASSO (1.2), observed with RIC  $\delta$  such that (3.7) holds. Then the residual,  $r_{\lambda} = y_{*}^{\epsilon} - Ax_{\lambda}$ , satisfies

$$\frac{|\boldsymbol{r}_{\lambda}|_{2}^{2}}{\lambda^{2}} \geqslant \frac{s_{\lambda}}{1+\delta}.$$

Combining the lower- and upper-bounds of  $rac{|r_{\lambda}|_2}{\lambda}$  we conclude the following.

**Theorem 3.3** (Sparsity bound). Fix  $\lambda < \lambda_{\infty} := |A^{\top} y_{*}^{\epsilon}|_{\infty}$  and let  $x_{\lambda}$  be the  $s_{\lambda}$ -sparse minimizer of the corresponding LASSO (1.2), observed with RIC  $\delta$  such that (3.7) holds. Then

$$(3.11) \frac{s_{\lambda}}{1+\delta} \leqslant \frac{|\boldsymbol{r}_{\lambda}|_{2}^{2}}{\lambda^{2}} \leqslant (\beta_{\delta}+\theta)^{2}k, \delta > \delta_{t}, t = [(1+\delta_{t})(\beta_{\delta_{t}}+\theta)^{2}k] + 1.$$

In particular, we recover (3.8),  $s_{\lambda} \leq [(1+\delta)(\beta_{\delta}+\theta)^2k]$ .

We demonstrate the application of corollary 3.3 for different choices of RICs. In all cases, we set  $\theta=0.1$ . We begin with the RIC  $\delta=0.7$ , obtaining  $(\beta_{\delta}+\theta)=2.52 \leadsto s_{\lambda} \leqslant (1+\delta)(\beta_{\delta}+\theta)^2 k < 11k$ . Thus, with t=11k we require  $\delta_{11k}<0.7$  which in turn, by (2.3), set the number of required observations

$$s_{\lambda} < 11k$$
:  $m \approx Const. \frac{11k}{0.7^2} \ln(e^N/k) \approx Const. 22.4 k \ln(e^N/k)$ .

For a second example we choose a smaller RIC  $\delta = 0.4$  and  $\theta = 0.1$ . Recall, that a smaller  $\delta$  requires more observations yet in the number of observations in the present context depends on  $\delta_t$ .

In this case  $(\beta_{\delta} + \theta) = 1.55 \rightsquigarrow s_{\lambda} \leqslant (1 + \delta)(\beta_{\delta} + \theta)^2 k = 3.36k < 4k$ . This requires a slightly larger number of observations (or at least a smaller bound (2.3))

$$s_{\lambda} < 4k$$
:  $m \approx Const. \frac{4k}{0.4^2} \ln(e^N/k) \approx Const. 25 k \ln(e^N/k)$ .

Finally, as a third example we choose an even smaller the RIC  $\delta = 0.26$  and the same  $\theta = 0.1$ . In this case  $(\beta_{\delta} + \theta) = 1.35 \rightsquigarrow s_{\lambda} \leq (1 + \delta)(\beta_{\delta} + \theta)^2 k = 2.28k < 3k$ , and this yields the number of required observations

$$s_{\lambda} < 3k$$
:  $m \approx Const. \frac{3k}{0.26^2} \ln(e^N/k) \approx Const. 44.4 k \ln(e^N/k)$ .

Remark 3.4 (On the threshold parameter  $\chi$ ). Observe that the sparsity bound  $s_{\lambda}$  is uniform in the small scale  $\mu$  throughout the parametric regime assumed in (3.1). Thus, in the range of  $\lambda \gg 2\mu/\sqrt{k}$ , the support of the computed solution  $|\mathbf{x}_{\lambda}|_0$  can grow at most by a fixed factor relative to the k-support of underlying unknown  $\mathbf{x}_*$ , [38, Appendix A]. We write

(3.12) 
$$s_{\lambda} < ([\chi^2] + 1)k, \qquad \chi := \sqrt{1 + \delta}(\beta_{\delta} + \theta) = \frac{1 + \delta}{\sqrt{1 - \delta^2} - \delta/4} + \sqrt{1 + \delta}\theta.$$

We have the theoretical bounds  $[\chi^2] + 1 = 11$  corresponding to  $\delta = 0.7$  and  $[\chi^2] + 1 = 4$  corresponding to  $\delta \approx 0.4$ .

3.2. Numerical simulations. We report here on our simulations of the unconstrained LASSO (1.2), applied to the recovery of k-sparse data,  $\sigma_k = 0$ , that is  $\mu = \epsilon$ , with (k, m, N) = (160, 1024, 4096). We consider different levels of noise  $\epsilon = 10^{-3}, 10^{-2}, 10^{-1}$ , in the corresponding parametric regime (3.1),  $\lambda > 2\mu/\sqrt{k} = 0.16\epsilon$ . The results are obtained by averaging 100 observations using randomly generated RNSP<sub> $\rho,\tau$ </sub> matrices based on Gaussian distributions. A simple proof of the RIP for such matrices can be found in [4]. The results are compared with the sparsity bound of theorem 3.3. We note that our sparsity bound depends in an essential manner on the RICs,  $1 \pm \delta$ , in (2.2). The parametric regime in (2.3) provides only a rough estimate on the range of allowable RICs, and in particular, does not cover the parameters used in the simulations below, [26, §9.4]. A detailed study which traces the sharp RICs can be found in [2, 3], but is beyond the scope of our work. We compare the simulations with our sparsity bound based on the RIC  $\delta = 0.7$ . This is partly motivated by the result of [10] in which the authors prove an exact BP recovery of k-sparse data from the RIP with  $\delta_{tk} < \sqrt{(t-1)/t}$ . In our case, the computation reported in figure 3.1 indicates the actual sparsity  $s_{\lambda} < tk$  with t=2 which is consistent with  $\delta < \sqrt{1/2} \approx 0.7$ . Although the RIC  $\delta = 0.7$  does not provide a tight bound,  $s_{\lambda} < 11k$ , it suffices to detect the correct behavior of the LASSO minimizer, reported in figures 3.1–3.3 and 4.1–4.2.

We record here the corresponding parameters involved in our bounds:

$$\beta_{\delta_{|\delta=0.7}} = \frac{\sqrt{1+\delta}}{\sqrt{1-\delta^2} - \delta/4} = 2.42, \quad \chi_{|\delta=0.7} = \sqrt{1+\delta}(\beta_{\delta} + \theta) = 3.16, \quad \eta_{|\delta=0.7} = \frac{1}{1+\delta} = 0.59$$

Our main result on the sparsity of the LASSO minimizer in theorem 3.3 provides a reasonably accurate information about the behavior of the unconstrained LASSO minimization (1.2). Figure 3.1 shows the behavior of the support,  $s_{\lambda} = |x_{\lambda}|_0$ , starting with  $s_{\lambda} = 0$  for  $\lambda > \lambda_{\infty}$  and monotonically increasing as  $\lambda$  decreases all the way to a critical value,  $\lambda_c \sim 0.11$ , at which point  $s_{\lambda_c}$  reaches its maximal value of 215. This should be compared with our bound  $s_{\lambda} \leq (1+\delta)(\beta_{\delta}+\theta)^2k$ . For  $\delta = 0.7$  we have  $s_{\lambda} \leq 11k$ , which is a rough sparsity bound, relative to the actual  $s_{\lambda} \sim 215$ . A smaller RIC  $\delta \sim 0.2$  yields a tighter sparsity bound  $1.66k \sim 313$ .

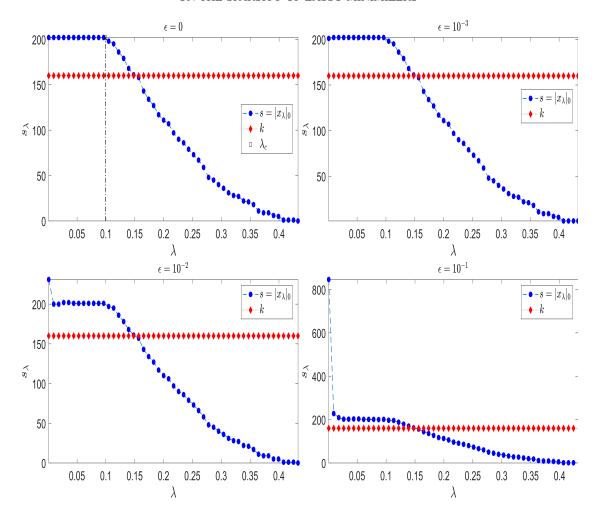


FIGURE 3.1. The support for computed minimizer  $s_{\lambda} = |x_{\lambda}|_0$  of k-sparse data, k = 160, peaks at the threshold value of  $k_{\text{max}} \sim 215$  when  $\lambda$  reaches  $\lambda_c \sim 0.11$ . This should be compared with the rough upper bound  $(1+\delta)(\beta_{\delta}+\theta)^2k \leqslant 11k$  corresponding to the RIC  $\delta=0.7$ , and the more realistic bound 4k corresponding to  $\delta=0.4$ . Observe (lower figures) that for exceedingly small  $\lambda \ll \epsilon$ , there is an additional growth of order  $\frac{\epsilon}{\lambda}$ .

Observe that according to (3.4), the  $\ell_1$ -size of the LASSO minimizer  $x_{\lambda}$  remains smaller than the target  $|x_*(k)|_1$ , Indeed, as long as the residual  $|r_{\lambda}|_2 > \mu$ , then

(3.13) 
$$|x_*(k)|_1 - |x_\lambda|_1 \geqslant (|r_\lambda|_2 - \mu) \frac{|r_\lambda|_2}{\lambda}.$$

This is depicted in figure 3.2: as  $\lambda$  decreases, the ratio  $\frac{|r_{\lambda}|_2}{\lambda} \gtrsim \sqrt{s_{\lambda}}$  is increasing until  $|x_{\lambda}|_1$  reaches its upper bound of  $|x_*(k)|_1$ .

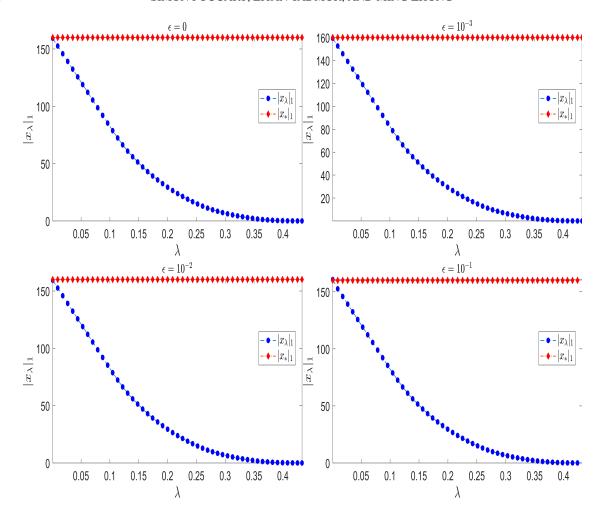


FIGURE 3.2.  $\ell_1$  norm of  $x_{\lambda}$  approaches its upper-bound  $|x_*(k)|_1$  as  $\lambda$  decreases.

Figure 3.3 shows the aptitude of the lower- and upper-bounds of the re-scaled residual (3.11), in capturing the re-scaled residual  $\frac{|{\bm r}_\lambda|_2}{\lambda}$ . Again, the three quantities increase with deceasing  $\lambda$ , until  $\lambda$  reaches the threshold  $\lambda_c$  at which point the re-scaled residual,  $\frac{|{\bm r}_\lambda|_2}{\lambda}$ , peaks at its maximal value  $\sim 27$ , in agreement with the upper-bound (3.2),  $\frac{|{\bm r}_\lambda|_2}{\lambda} < \beta_\delta \sqrt{k} + \frac{2\epsilon}{\lambda} < 30.61 + \frac{2\epsilon}{\lambda}$ .

#### 4. Error bounds

4.1.  $\ell_2$ -error bounds. The sparsity bound (3.11) was derived based on a two-sided  $\ell_2$ -bound of the scaled residual. The latter can be converted into a two-sided  $\ell_2$  error bound of  $|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_*(k)|_2$ . Note that since  $||\boldsymbol{r}_{\lambda}|_2 - |A(\boldsymbol{x}_*(k) - \boldsymbol{x}_{\lambda})|_2| \leq \mu$ , then the upper-bound on  $|\boldsymbol{r}_{\lambda}|_2$ , see (3.5), also bounds the 'observed error'  $A(\boldsymbol{x}_{\lambda} - \boldsymbol{x}_*(k))$ ,

$$(4.1) |A(\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k))|_{2} \leqslant (\beta_{\delta} + \theta)\sqrt{k}\lambda + \mu \leqslant \beta_{\delta}\sqrt{k}\lambda + 3\mu.$$

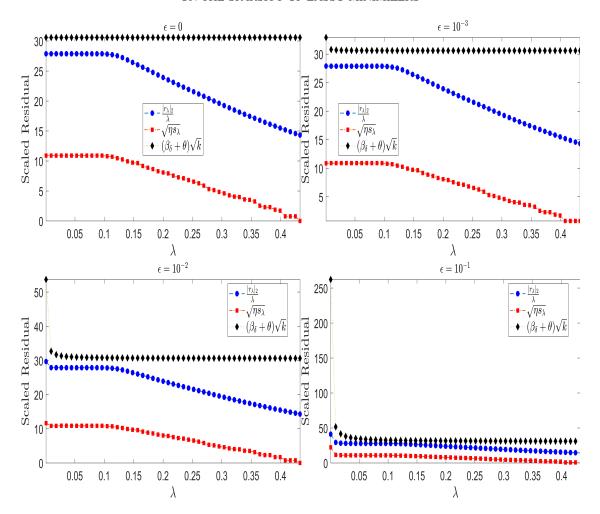


FIGURE 3.3. Re-scaled residual  $\frac{|r_{\lambda}|_2}{\lambda}$  captured between its lower- and upper-bounds (3.10) and (3.2),  $\sqrt{\eta s_{\lambda}} \leqslant \frac{|r_{\lambda}|_2}{\lambda} \leqslant \beta_{\delta} \sqrt{k} + \frac{2\epsilon}{\lambda} \approx 30.61$  with  $(\eta, \beta, \theta) = (0.59, 2.42, 0.1)$  corresponding to  $\delta = 0.7$ . It peaks at a threshold value of 27, independent of the level of noise. When  $\lambda \ll \epsilon$ , there is an additional large term of order  $\frac{2\epsilon}{\lambda}$ .

The sparsity of  $x_{\lambda}$  does not exceed  $s_{\lambda} \leq ([\chi^2] + 1)k$  hence  $x_{\lambda} - x_*(k)$  has sparsity  $([\chi^2] + 2)k$ , and the RIP (2.2) implies the  $\ell_2$ -error upper-bound

$$(4.2) |\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{2} \leqslant \frac{1}{\sqrt{1-\delta}} (\beta_{\delta} \sqrt{k}\lambda + 3\mu), \delta = \delta_{([\chi^{2}]+2)k}.$$

In particular, (4.2) with  $\frac{1}{\sqrt{1-\delta}} \leqslant 1.83$  and  $\beta_{\delta} \leqslant 2.42$  corresponding to  $\delta = 0.7$  yields

$$|x_{\lambda} - x_{*}(k)|_{2} \leq 4.43\sqrt{k}\lambda + 5.48\mu.$$

This recovers a quantitative version of the  $\ell_2$  upper bound proved under additional condition of an incoherence design assumption in [30, Theorem 1], an  $\ell_1$ -CMSV assumption<sup>5</sup> [38], or restricted eigenvalue bound in [28, Theorem 11.1].

<sup>&</sup>lt;sup>5</sup>In fact, we slightly improve the quadratic dependence of the bound in [38, (23)] on the  $\ell_1$ -CMSV constant  $\sim \rho_{4k}^{-2}$ , mentioned in (4.6) below.

The upper-bound (4.2) is sharp in the sense of having a tight  $\ell_2$ -lower bound: since the error  $x_{\lambda} - x_*(k)$  is at most  $([\chi^2] + 2)k$ -sparse, we can use the RIP to translate the lower bound (3.10) into an  $\ell_2$  lower-bound,

$$|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{2} \geqslant \frac{1}{\sqrt{1+\delta}} |A(\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k))|_{2} \geqslant \frac{1}{\sqrt{1+\delta}} (|\boldsymbol{r}_{\lambda}|_{2} - \mu) \geqslant \frac{\sqrt{s_{\lambda}}\lambda}{1+\delta} - \frac{\mu}{\sqrt{1+\delta}}.$$

We summarize these bounds in the following form.

**Theorem 4.1** ( $\ell_2$ -bound). Fix  $\lambda < \lambda_{\infty} := |A^{\top} y_*^{\epsilon}|_{\infty}$  and let  $x_{\lambda}$  be the  $s_{\lambda}$ -sparse minimizer of the corresponding LASSO (1.2), observed with RIP matrix A. Then

$$(4.4) \quad \frac{1}{\sqrt{1+\delta}} \left( \frac{\sqrt{s_{\lambda}}\lambda}{\sqrt{1+\delta}} - \mu \right) \leqslant |\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{2} \leqslant \frac{1}{\sqrt{1-\delta}} (\beta_{\delta}\sqrt{k}\lambda + 3\mu), \qquad \delta = \delta_{([\chi^{2}]+2)k}.$$

Remark 4.2 (Compared with the  $\ell_1$ -entropy bound). The extremal relation  $\langle A_{\mathcal{S}} \boldsymbol{x}_{\lambda,\mathcal{S}}, \boldsymbol{r}_{\lambda} \rangle = \lambda |\boldsymbol{x}_{\lambda,\mathcal{S}}|_1$  and the RIP (2.2) yield  $\lambda |\boldsymbol{x}_{\lambda,S}|_1 \leqslant \sqrt{1+\delta} |\boldsymbol{x}_{\lambda,S}|_2 |\boldsymbol{r}_{\lambda}|_2$ , and hence we end up with a lower-bound involving the  $\ell_1$ -entropy of  $\{\boldsymbol{x}_{\lambda,S}\}$ ,

$$\frac{|\boldsymbol{r}_{\lambda}|_{2}^{2}}{\lambda^{2}} \geqslant \frac{1}{1+\delta} Ent(\boldsymbol{x}_{\lambda,S}) \qquad Ent(\boldsymbol{x}) := \frac{|\boldsymbol{x}|_{1}^{2}}{|\boldsymbol{x}|_{2}^{2}}.$$

This bound is tied to a Null Entropy Property of A [1, §3.2] or the  $\ell_1$ -CMSV constant  $\rho_s(A)$  introduced in [38]<sup>6</sup>

$$(4.6) \qquad \frac{|\mathbf{r}_{\lambda}|_{2}^{2}}{\lambda^{2}} \geqslant \frac{\operatorname{Ent}(\mathbf{x}_{\lambda,S})}{\rho_{s}(A)}, \qquad \rho_{s}(A) := \min_{|\mathbf{x}|_{2}=1} \left\{ |A\mathbf{x}|_{2} : \operatorname{Ent}(\mathbf{x}) \leqslant s \right\}.$$

<sup>&</sup>lt;sup>6</sup>Which is not to be confused with the RNSP parameter in (2.6)

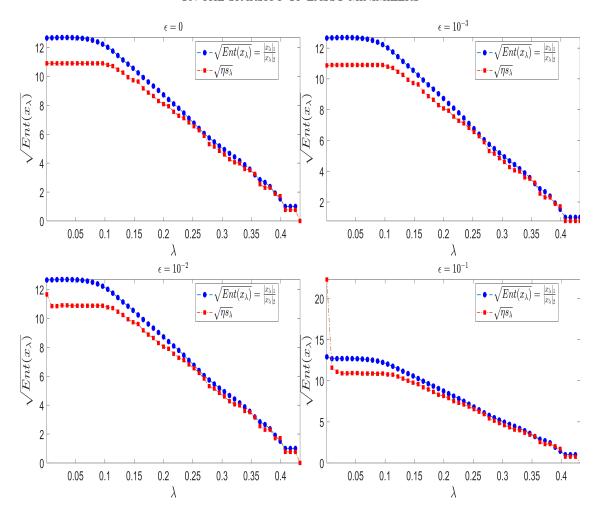


FIGURE 4.1. Lower bounds of the re-scaled residual: (3.10) with  $\eta := \frac{1}{1+\delta} = 0.592$  vs. the  $\ell_1$ -entropy based (4.5).

Clearly, if  $\mathbf{x}_{\lambda}$  has the sparsity  $s_{\lambda}$  then  $\operatorname{Ent}(\mathbf{x}_{\lambda,S}) \leqslant s_{\lambda}$ . Here we note about the reverse implication, namely — if the reverse inequality holds,  $\operatorname{Ent}(\mathbf{x}_{\lambda,S}) \gtrsim s_{\lambda}$ , then it would yield our sparsity result of lemma 3.2, based on the lower bound  $\frac{|\mathbf{r}_{\lambda}|_2}{\lambda} \gtrsim \frac{\sqrt{s_{\lambda}}}{\rho_{s_{\lambda}}}$ . Theorem 3.3 suggests the lower-entropy bound for the minimizers  $\mathbf{x}_{\lambda}$ . Indeed, figure 4.1 shows a remarkable agreement between the lower bound (3.10) with  $\delta=0.7$  and the  $\ell_1$ -entropy bound (4.5),  $\operatorname{Ent}(\mathbf{x}_{\lambda})$ , at least before the support of  $\mathbf{x}_{\lambda}$  reaches its peak at  $k_{max}$ .

4.2.  $\ell_1$ -error bound. We recall the  $\ell_2$ -bound (4.2) which we express in the form  $|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_*(k)|_2 \le \frac{1}{\sqrt{1-\delta}}(\beta_{\delta} + 3/2\theta)\sqrt{k}\lambda$ . Since  $\boldsymbol{x}_{\lambda} - \boldsymbol{x}_*(k)$  has sparsity of order  $\leqslant k + \chi^2 k$ , we derive the following  $\ell_1$ -bound.

**Theorem 4.3** ( $\ell_1$ -error bound). Fix  $\lambda < \lambda_{\infty} := |A^{\top} \boldsymbol{y}_{*}^{\epsilon}|_{\infty}$  and let  $\boldsymbol{x}_{\lambda}$  be the LASSO minimizer of (1.2), observed with RIP matrix A with RIC  $\delta$  such that (3.7) holds. Then the following  $\ell_1$ -error

bound holds,

$$|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{1} \leqslant \sqrt{\left(1 + (1+\delta)(\beta_{\delta} + \theta)^{2}\right)k}|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{2}$$

$$< \sqrt{1+\delta}(\beta_{\delta} + \frac{1}{2} + \theta)\sqrt{k}\frac{1}{\sqrt{1-\delta}}\left(\beta_{\delta} + \frac{3}{2}\theta\right)\sqrt{k}\lambda$$

$$< \frac{\sqrt{1+\delta}}{\sqrt{1-\delta}}\frac{1}{\lambda}\left(\left(\beta_{\delta} + \frac{1}{2}\right)\sqrt{k}\lambda + 2\mu\right)^{2}.$$

The amplitude of  $k\lambda$  in the  $\ell_1$ -error bound (4.7) is not sharp. For example, with RIC  $\delta=\delta_{11k}<0.7$  we have  $\beta_\delta>2$  in which case, omitting the negligibly small  $\mu^2/\lambda$  terms, one ends up with the improved bound

$$(4.8) |x_{\lambda} - x_{*}(k)|_{1} < \frac{\sqrt{1+\delta}}{\sqrt{1-\delta}} \Big( (\beta_{\delta} + 1/4)^{2} k\lambda + (4\beta_{\delta} + 1)\sqrt{k}\mu \Big) < 16.97k\lambda + 24.04\sqrt{k}\mu.$$

We conclude with an alternative derivation of an  $\ell_1$ -error bound. To this end, we recall the RNSP bound [26, Theorem 4.20], which states that for all  $\mathcal{K} \subset \{1, 2, ..., N\}$  of size  $\leqslant k$  and for any u,  $v \in \mathbb{R}^N$ , the following holds,

$$|\boldsymbol{u} - \boldsymbol{v}|_1 \leqslant \frac{1 + \rho}{1 - \rho} (|\boldsymbol{u}|_1 - |\boldsymbol{v}|_1 + 2|\boldsymbol{v}_{\mathcal{K}^c}|_1) + \frac{2\tau}{1 - \rho} |A(\boldsymbol{u} - \boldsymbol{v})|_2, \qquad |\mathcal{K}| \leqslant k.$$

Using it with  $(\boldsymbol{u}, \boldsymbol{v}) = (\boldsymbol{x}_{\lambda}, \boldsymbol{x}_{*}(k))$  and  $\mathcal{K} = \operatorname{supp}(\boldsymbol{x}_{*}(k))$  yields

$$(4.9) |\mathbf{x}_{\lambda} - \mathbf{x}_{*}(k)|_{1} \leq \frac{1 + \rho}{1 - \rho} (|\mathbf{x}_{\lambda}|_{1} - |\mathbf{x}_{*}(k)|_{1}) + \frac{2\tau}{1 - \rho} |A(\mathbf{x}_{\lambda} - \mathbf{x}_{*}(k))|_{2}.$$

Now, using (3.4) to bound the term inside the first parenthesis on the right, and as before, noting that the second term does not exceed  $|A(x_{\lambda} - x_{*}(k))|_{2} \leq |r_{\lambda}|_{2} + \mu$ , we find

$$|\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{1} \leqslant \frac{1+\rho}{1-\rho} \Big\{ |\boldsymbol{r}_{\lambda}|_{2} \Big( \frac{2\tau}{1+\rho} + \frac{\mu}{\lambda} - \frac{|\boldsymbol{r}_{\lambda}|_{2}}{\lambda} \Big) + \frac{2\tau}{1+\rho} \mu \Big\}.$$

Given the RNSP parameters (2.6),  $2\tau = 2\beta\sqrt{k}$  and  $\frac{\mu}{\lambda} < \frac{\theta}{1+\rho}\sqrt{k}$ , the last bound yields

$$(4.10) |\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{1} \leqslant \frac{1}{1 - \rho} \left\{ |\boldsymbol{r}_{\lambda}|_{2} \left( 2(\beta_{\delta} + \theta)\sqrt{k} - (1 + \rho) \frac{|\boldsymbol{r}_{\lambda}|_{2}}{\lambda} \right) + \beta_{\delta}\theta k \lambda \right\}.$$

Viewed as quadratic in  $\frac{|r_{\lambda}|_2}{\lambda}$ , the first expression on the right admits a maximal value  $\frac{(\beta_{\delta}+\theta)^2}{1+\rho}k\lambda$ , and we finally end up with

$$(4.11) |\boldsymbol{x}_{\lambda} - \boldsymbol{x}_{*}(k)|_{1} \leqslant \frac{1}{1 - \rho^{2}} \left( (\beta_{\delta} + \theta)^{2} k \lambda + (1 + \rho) \beta_{\delta} \theta k \lambda \right) \leqslant \frac{(\beta_{\delta} + 2\theta)^{2}}{1 - \rho^{2}} k \lambda.$$

This recovers the  $\ell_1$ -bound of order  $\mathcal{O}(k\lambda)$  as in (4.7). However, since the  $\ell_1$  bound (4.11) involves the value of  $1/(1-\rho)^2$ , it is therefore limited to the RIC  $\delta < 4/\sqrt{41}$  where  $\rho$  approaches 1.

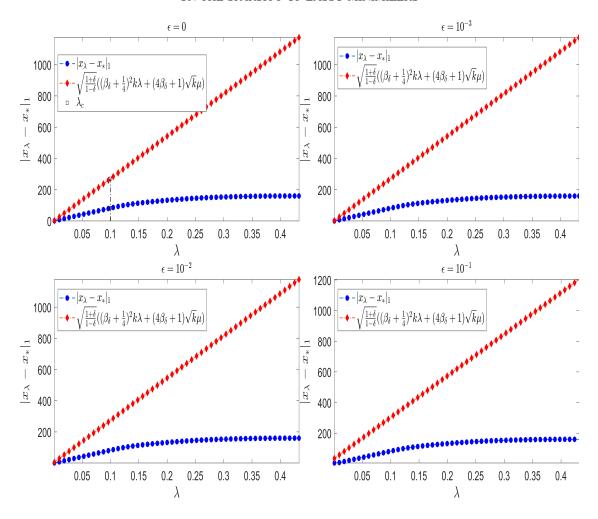


FIGURE 4.2.  $\ell_1$ -error for recovery of sparse data compared with the upper-bound (4.12).

4.3. Numerical simulations. We report on the error behavior in our simulations of the unconstrained LASSO (1.2), applied to the recovery of k-sparse data,  $\sigma_k = 0$ , that is  $\mu = \epsilon$ , with (k, m, N) = (160, 1024, 4096). The results are obtained by averaging 100 observations using randomly generated RNSP $_{\rho,\tau}$  matrices based on Gaussian distributions. We compare the  $\ell_1$ -error with the error bound (4.8)

$$|\mathbf{x}_{\lambda} - \mathbf{x}_{*}(k)|_{1} \leq 16.97 * 160\lambda + 24.04\sqrt{160} \epsilon.$$

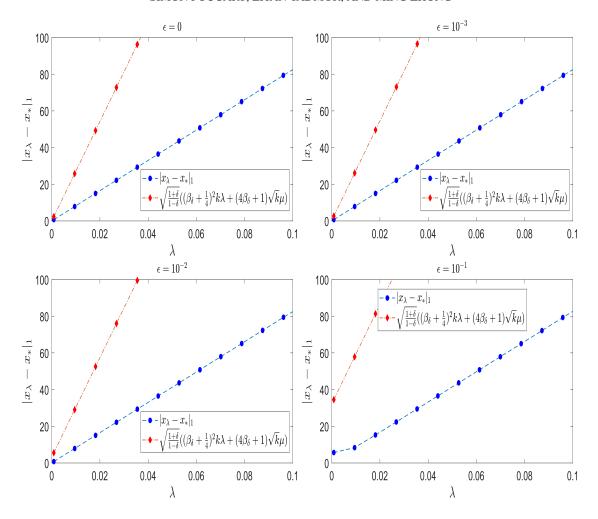


FIGURE 4.3. The  $\ell_1$ -error compared with the upper-bound (4.12) zoomed near  $\lambda = 0$ .

**Acknowledgment**. We are indebted to Wolfgang Dahmen for comments which greatly improved the presentation of material of this paper.

### REFERENCES

- [1] J. ANDERSSON AND J.-O. STRÖMBERG, On the theorem of uniform recovery of random sampling matrices, IEEE Transactions on Information Theory, 60 (2014), pp. 1700–1710.
- [2] B. BAH AND J. TANNER, *Improved bounds on restricted isometry constants for Gaussian matrices*, SIAM Journal on Matrix Analysis and Applications, 31 (2010), pp. 2882–2898.
- [3] ——, Bounds of restricted isometry constants in extreme asymptotics: formulae for Gaussian matrices, Linear Algebra and its Applications, 441 (2014), pp. 88–109.
- [4] R. BARANIUK, M. DAVENPORT, R. DEVORE, AND M. WAKIN, A simple proof of the restricted isometry property for random matrices, Constructive Approximation, 28 (2008), pp. 253–263.
- [5] A. Belloni and V. Chernozhukov, *Least squares after model selection in high-dimensional sparse models*, Bernoulli, 19 (2013), pp. 521–547.
- [6] P. J. BICKEL, Y. RITOV, AND A. B. TSYBAKOV, Simultaneous analysis of Lasso and Dantzig selector, The Annals of Statistics, 37 (2009), pp. 1705–1732.
- [7] J. D. BLANCHARD, C. CARTIS, AND J. TANNER, Compressed sensing: How sharp is the restricted isometry property?, SIAM Review, 53 (2011), pp. 105–125.

- [8] J. BOURGAIN, S. DILWORTH, K. FORD, S. KONYAGIN, AND D. KUTZAROVA, *Explicit constructions of RIP matrices and related problems*, Duke Mathematical Journal, 159 (2011), pp. 145–185.
- [9] P. BÜHLMANN AND S. VAN DE GEER, Statistics for High-Dimensional Data: Methods, Theory and Applications, Springer Science & Business Media, 2011.
- [10] T. T. CAI AND A. ZHANG, Sparse representation of a polytope and recovery of sparse signals and low-rank matrices, IEEE transactions on information theory, 60 (2013), pp. 122–132.
- [11] E. CANDÈS AND T. TAO, *The Dantzig selector: Statistical estimation when p is much larger than n*, The Annals of Statistics, 35 (2007), pp. 2313–2351.
- [12] E. J. CANDÈS, *The restricted isometry property and its implications for compressed sensing*, Comptes Rendus Mathematique, 346 (2008), pp. 589–592.
- [13] E. J. CANDÈS, J. K. ROMBERG, AND T. TAO, *Robust uncertainty principles: exact signal reconstruction from high incomplete frequency information*, IEEE Trans. Infom. Theor., 52 (2006), pp. 489 509.
- [14] —, Stable signal recovery from incomplete and inaccurate measurements, Commun. Pure Appl. Math., 59 (2006), pp. 1207 1233.
- [15] E. J. CANDÈS AND T. TAO, The Dantzig selector: statistical estimation when p is much larger than n, Ann. Stat., 35 (2007), pp. 2313 2351.
- [16] S. S. CHEN, *Basis pursuit*, PhD thesis, Stanford University, 1996.
- [17] S. S. CHEN, D. L. DOHONO, AND M. A. SAUNDERS, *Atomic decomposition by Basis Pursuit*, SIAM J. Sci. Comput., 20 (1999), pp. 33 61.
- [18] S. S. CHEN AND D. L. DONOHO, *Basis pursuit*, in Proceedings of 1994 28th Asilomar Conference on Signals, Systems and Computers, vol. 1, IEEE, 1994, pp. 41–44.
- [19] A. COHEN, W. DAHMEN, AND R. A. DEVORE, Compressed sensing and the best k-term approximation, J. Amer. Math. Soc., 22 (2009), pp. 211 231.
- [20] R. DEVORE, G. PETROVA, AND P. WOJTASZCZYK, *Instance-optimality in probability with an*  $\ell_1$ -minimization decoder, Appl. Comput. Harmon. Anal., 27 (2009), pp. 275–288.
- [21] D. L. DONOHO, *Compressed sensing*, IEEE Trans. Infom. Theor., 52 (2006), pp. 1289 1306.
- [22] D. L. DONOHO AND M. ELAD, Optimally sparse representations in general (non-orthogonal) dictionaries via  $\ell^1$  minimization, Proc. Nat. Acad. Sci., 100 (2003), pp. 2197 2202.
- [23] D. L. DONOHO AND X. Huo, *Uncertainty principles and ideal atomic decompositions*, IEEE Trans. Inform. Theor., 47 (2001), pp. 2845 2862.
- [24] D. L. DONOHO AND I. M. JOHNSTONE, Empirical atomic decomposition, Unpublished manuscript, (1995).
- [25] S. FOUCART, Sparse recovery algorithms: sufficient conditions in terms of restricted isometry constants, in Approximation Theory XIII: San Antonio 2010, Springer, 2012, pp. 65–77.
- [26] S. FOUCART AND H. RAUHUT, A Mathematical Introduction to Compressive Sensing, Applied and Numerical Harmonic Analysis, Birkhäuser, 2013.
- [27] R. GRIBONVAL AND M. NIELSEN, *Sparse representations in unions of bases*, IEEE Trans. Inform. Theor., 49 (2003), pp. 3320 3325.
- [28] T. HASTIE, R. TIBSHIRANI, AND M. WAINWRIGHT, Statistical Learning with Sparsity: the LASSO and generalizations, CRC Press, 2015.
- [29] S. MALLAT AND W. L. HWANG, Singularity detection and processing with wavelets, IEEE Transactions on Information Theory, 38 (1992), pp. 617–643.
- [30] N. MEINSHAUSEN, B. YU, ET AL., Lasso-type recovery of sparse representations for high-dimensional data, The Annals of Statistics, 37 (2009), pp. 246–270.
- [31] L. I. RUDIN, S. OSHER, AND E. FATEMI, *Nonlinear total variation based noise removal algorithms*, Physica D: nonlinear phenomena, 60 (1992), pp. 259–268.
- [32] W. Su, M. Bogdan, and E. Candès, *False discoveries occur early on the Lasso path*, The Annals of Statistics, (2017), pp. 2133–2150.
- [33] E. TADMOR, *Hierarchical construction of boudned solutions in regularity spaces*, Communications in Pure & Applied Mathematics, 69(6) (2015), pp. 1087–1109.
- [34] E. TADMOR, S. NEZZAR, AND L. VESSE, A multiscale image representation using hierarchical ( $BV, L^2$ ) decompositions, Multiscale Model. Simul., 2 (2004), pp. 554 579.
- [35] ——, Multiscale hierarchical decomposition of images with applications to deblurring, denoising and segmentation, Commun. Math. Sci., 6 (2008), pp. 281 307.

- [36] E. TADMOR AND C. TAN, Hierarchical construction of bounded solutions of div U=F in critical regularity spaces, in Nonlinear Partial Differential Equations, H. Holden and K. Karlsen, eds., Abel Symposia 7, Oslo, Sep. 2010, pp. 255 269. 2010 Abel Symposium.
- [37] E. TADMOR AND M. ZHONG, The multi-scale hierarchical reconstruction of ill-posed inverse problems, In preparation, (2021).
- [38] G. TANG AND A. NEHORAI, *Performance analysis of sparse recovery based on constrained minimal singular values*, IEEE Transactions on Signal Processing, 59 (2011), pp. 5734–5745.
- [39] R. TIBSHIRANI, Regression shrinkage and selection via the lasso, J. Roy. Stat. Soc. B, 58 (1996), pp. 267 288.
- [40] J. A. TROPP, *Greed is good: Algorithmic results for sparse approximation*, IEEE Trans. Inform. Theor., 50 (2004), pp. 2231 2242.
- [41] H. ZHANG, W. YIN, AND L. CHENG, Necessary and sufficient conditions on solution uniqueness in  $\ell_1$  minimization, Journal of Optimization Theory and Applications, 164 (2015), pp. 109 122.

DEPARTMENT OF MATHEMATICS AND TEXAS A&M INSTITUTE OF DATA SCIENCE

TEXAS A&M UNIVERSITY, COLLEGE STATION, TX 77843

Email address: foucart@tamu.edu

DEPARTMENT OF MATHEMATICS AND INSTITUTE FOR PHYSICAL SCIENCES & TECHNOLOGY (IPST)

UNIVERSITY OF MARYLAND, COLLEGE PARK, MD 20742

Email address: tadmor@umd.edu

TEXAS A&M INSTITUTE OF DATA SCIENCE
TEXAS A&M UNIVERSITY, COLLEGE STATION, TX 77843

Email address: mingzhong@tamu.edu