# On U-Statistics and Compressed Sensing I: Non-Asymptotic Average-Case Analysis

Fabian Lim* and Vladimir Marko Stojanovic

***Abstract*—Hoeffding's U-statistics model combinatorial-type matrix parameters (appearing in CS theory) in a natural way. This paper proposes using these statistics for analyzing random compressed sensing matrices, in the non-asymptotic regime (relevant to practice). The aim is to address certain pessimisms of "worst-case" *restricted isometry* analyses, as observed by both Blanchard & Dossal, et. al.**

**We show how U-statistics can obtain "average-case" analyses, by relating to *statistical restricted isometry property* (StRIP) type recovery guarantees. However unlike standard StRIP, random signal models are not required; the analysis here holds in the *almost sure* (probabilistic) sense. For Gaussian/bounded entry matrices, we show that both $\ell_1$-minimization and LASSO essentially require on the order of $k \cdot [\log((n-k)/u) + \sqrt{2(k/n)\log(n/k)}]$ measurements to respectively recover at least $1-5u$ fraction, and $1-4u$ fraction, of the signals. Noisy conditions are considered. Empirical evidence suggests our analysis to compare well to Donoho & Tanner's recent large deviation bounds for $\ell_0/\ell_1$-equivalence, in the regime of block lengths $1000 \sim 3000$ with high undersampling ($50 \sim 150$ measurements); similar system sizes are found in recent CS implementation.**

**In this work, it is assumed throughout that matrix columns are independently sampled.**

***Index Terms*—approximation, compressed sensing, satistics, random matrices**

## I. INTRODUCTION

Compressed sensing (CS) analysis involves relatively recent results from random matrix theory [1], whereby recovery guarantees are framed in the context of matrix parameters known as *restricted isometry constants*. Other matrix parameters are also often studied in CS. Earlier work on sparse approximation considered a matrix parameter known as *mutual coherence* [2]–[4]. Fuchs' work on *Karush-Kuhn-Tucker (KKT)* conditions for sparsity pattern recovery considered a parameter involving a matrix *pseudoinverse* [5], re-occurring in recent work [4], [6], [7]. Finally, the *null-space property* [8]–[10] is gaining recent popularity - being the parameter closest related to the fundamental compression limit dictated by *Gel'fand widths*. All above parameters share a similar feature, that is they are defined over subsets of a certain fixed size $k$. This combinatorial nature makes them difficult to evaluate, even for moderate block lengths $n$. Most CS work therefore involve some form of randomization to help the analysis.

While the celebrated $k \log(n/k)$ result was initially approached via asymptotics, *e.g.*, [1], [11]–[13], implementations

F. Lim and V. M. Stojanovic are with the Research Laboratory of Electronics, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA 02139. Email: {flim,vlada}@mit.edu. This work was supported by NSF grant ECCS-1128226.

require finite block sizes. Hence, non-asymptotic analyses are more application relevant. In the same practical aspect, recent work deals with non-asymptotic analysis of *deterministic* CS matrices, see [4], [7], [14], [15]. On the other hand certain situations may not allow control over the sampling process, whereby the sampling may be inherently random, *e.g.*, prediction of clinical outcomes of various tumors based on gene expressions [6]. Random sampling has certain desirable simplicity/efficiency features - see [16] on data acquisition in the distributed sensor setting. Also recent hardware implementations point out energy/complexity-cost benefits of implementing *pseudo-random binary sequences* [17]–[19]; these sequences mimic statistical behavior. Non-asymptotic analysis is particularly valuable, when random samples are costly to acquire. For example, each clinical trial could be expensive to conduct an excessive number of times. In the systems setting, the application could be running on a tight energy budget - whereby processing/communication costs depend on the number of samples acquired.

This work is inspired by the *statistical* notion of the restricted isometry property (StRIP), initially developed for deterministic CS analysis [14], [15]. The idea is to relax the analysis, by allowing sampling matrix parameters (that guarantee signal recovery) to be satisfied for a *fraction* of subsets. Our interest is in "average-case" notions in the context of randomized sampling, reason being that certain pessimisms of "worst-case" restricted isometry analyses have been observed in past works [13], [20], [21]. On the other hand in [22], Donoho & Tanner remarked on potential benefits of the above "average-case" notion, recently pursued in an adaptation of a previous asymptotic result [23]. In the multichannel setting, "average-case" notions are employed to make analysis more tractable [24], [25]. In [26] a simple "thresholding" algorithm is analyzed via an "average" coherence parameter. However the works in this respect are few, most random analyses are of the "worst-case" type, see [12], [13], [21], [27]. We investigate the unexplored, with the aim of providing new insights and obtaining new/improved results for the "average-case".

Here we consider a random analysis tool that is well-suited to the CS context, yet seemingly left untouched in the literature. Our approach differs from that of deterministic matrices, where "average-case" analysis is typically made accessible via mutual coherence, see [14], [15], [18]. For random matrices, we propose an alternative approach via *U-statistics*, which do not require random signal models typically introduced in StRIP analysis, see [14], [25], [26]; here, the results are stated in the *almost sure* sense. U-statistics apply naturally to various kinds of non-asymptotic CS analyses,

since they are designed for combinatorial-type parameters. Also, they have a natural "average-case" interpretation, which we apply to recent recovery guarantees that share the same "average-case" characteristic. Finally thanks to the wealth of U-statistical literature, the theory developed here is open to other extensions, *e.g.*, in related work [28] we demonstrate how U-statistics may also perform "worst-case" analysis.

**Contributions**: "Average-case" analyses are developed based on U-statistics, which are i) empirically observed to have good potential for predicting CS recovery in non-asymptotic regimes, and ii) theoretically obtain measurement rates that incorporate a non-zero failure rate (similar to the $k \log(n/k)$ rate from "worst-case" analyses). We utilize a U-statistical large deviation concentration theorem, under the assumption that the matrix columns are independently sampled. The large deviation error bound holds *almost surely* (Theorem 1). No random signal model is needed, and the error is of the order $(n/k)^{-1} \log(n/k)$, whereby $k$ is the U-statistic kernel size (and $k$ also equals sparsity level). Gaussian/bounded entry matrices are considered. For concreteness, we connect with StRIP-type guarantees (from [6], [7]) to study the fraction of recoverable signals (*i.e.*, "average-case" recovery) of: i) $\ell_1$-*minimization* and ii) *least absolute shrinkage and selection operator (LASSO)*, under noisy conditions. For both these algorithms we show $\text{const} \cdot k[\log((n-k)/u) + \sqrt{2(k/n) \log(n/k)}]$ measurements are essentially required, to respectively recover at least $1 - 5u$ fraction (Theorem 2), and $1 - 4u$ fraction (Theorem 3), of possible signals. This is improved to $1 - 3u$ fraction for the noiseless case. Here $\text{const} = \max(4/(a_1 a_2)^2, 2c_1/(0.29 - a_1)^2)$ for to be specified constants $a_1, a_2, c_1$, where $c_1$ depends on the distribution of matrix entries. Note that the term $\sqrt{2(k/n) \log(n/k)}$ is at most 1 and vanishes with small $k/n$. Empirical evidence suggests that our approach compares well with recent results from Donoho & Tanner [23] - improvement is suggested for system sizes found in implementations [17], with large undersampling (*i.e.*, $m = 50 \sim 100$ and $n = 1000 \sim 3000$). The large deviation analysis here does show some pessimism in the size of $\text{const}$ above, whereby $\text{const} \geq 4$ (we conjecture possible improvement). For Gaussian/Bernoulli matrices, we find $\text{const} \approx 1.8$ to be inherently smaller, *e.g.*, for $k = 4$ this predicts recovery of $1 \times 10^{-6}$ fraction with 153 measurements - empirically $m = 150$.

**Note**: StRIP-type guarantees [6], [7] seem to work well, by simply *not* placing restrictive conditions on the maximum eigenvalues of the size-$k$ submatrices. Our theory applies fairly well for various considered system sizes $k, m, n$ (*e.g.*, Figure 4), however in *noisy* situations, a (relatively small) factor of $\sqrt{k}$ losses is seen without making certain maximum eigenvalue assumptions. For $\ell_1$-recovery, the estimation error is now bounded by a $\sqrt{k}$ factor of its best $k$-term approximation error (both errors measured using the $\ell_1$-norm). For LASSO, the the non-zero signal magnitudes must now be bounded below by a factor $\sqrt{2k \log n}$ (with respect to noise standard deviation), as opposed to $\sqrt{2 \log n}$ in [6]. These losses occur not because of StRIP analyses, but because of the estimation techniques employed here.

**Organization**: We begin with relevant background on CS

in Section II. In Section III we present a general U-statistical theorem for large-deviation ("average-case") behavior. In Section IV the U-statistical machinery is applied to StRIP-type "average-case" recovery. We conclude in Section V.

**Notation**: The set of real numbers is denoted $\mathbb{R}$. Deterministic quantities are denoted using $a, \mathbf{a}$, or $\mathbf{A}$, where bold fonts denote vectors (*i.e.*, $\mathbf{a}$) or matrices (*i.e.*, $\mathbf{A}$). Random quantities are denoted using upper-case *italics*, where $A$ is a random variable (RV), and $\boldsymbol{A}$ a random vector/matrix. Let $\Pr\{A \leq a\}$ denote the probability that event $\{A \leq a\}$ occurs. Sets are denoted using braces, *e.g.*, $\{1, 2, \cdots\}$. The notation $\mathbb{E}$ denotes expectation. The notation $i, j, \ell, \omega$ is used for indexing. We let $|| \cdot ||_p$ denote the $\ell_p$-norm for $p = 1$ and 2.

## II. PRELIMINARIES

### A. Compressed Sensing (CS) Theory

A vector $\mathbf{a}$ is said to be $k$-sparse, if at most $k$ vector coefficients are non-zero (*i.e.*, its $\ell_0$-distance satisfies $||\mathbf{a}||_0 \leq k$). Let $n$ be a positive integer that denotes block length, and let $\boldsymbol{\alpha} = [\alpha_1, \alpha_2, \cdots, \alpha_n]^T$ denote a length-$n$ signal vector with signal coefficients $\alpha_i$. The *best $k$-term approximation* $\overline{\boldsymbol{\alpha}}_k$ of $\boldsymbol{\alpha}$, is obtained by finding the $k$-sparse vector $\overline{\boldsymbol{\alpha}}_k$ that has minimal approximation error $||\overline{\boldsymbol{\alpha}}_k - \boldsymbol{\alpha}||_2$.

Let $\boldsymbol{\Phi}$ denote an $m \times n$ CS sampling matrix, where $m < n$. The length-$m$ **measurement vector** denoted $\mathbf{b} = [b_1, b_2, \cdots, b_m]^T$ of some length-$n$ signal $\boldsymbol{\alpha}$, is formed as $\mathbf{b} = \boldsymbol{\Phi}\boldsymbol{\alpha}$. Recovering $\boldsymbol{\alpha}$ from $\mathbf{b}$ is challenging as $\boldsymbol{\Phi}$ possesses a *non-trivial null-space*. We typically recover $\boldsymbol{\alpha}$ by solving the (convex) $\ell_1$-**minimization** problem

$$\min_{\tilde{\boldsymbol{\alpha}} \in \mathbb{R}^n} ||\tilde{\boldsymbol{\alpha}}||_1 \quad \text{s. t. } ||\tilde{\mathbf{b}} - \boldsymbol{\Phi}\tilde{\boldsymbol{\alpha}}||_2 \leq \epsilon. \quad (1)$$

The vector $\tilde{\mathbf{b}}$ is a *noisy* version of the original measurements $\mathbf{b}$, and here $\epsilon$ bounds the noise error, *i.e.*, $\epsilon \geq ||\tilde{\mathbf{b}} - \mathbf{b}||_2$. Recovery conditions have been considered in many flavors [2], [3], [11], [22], [23], and mostly rely on studying parameters of the sampling matrix $\boldsymbol{\Phi}$.

For $k \leq n$, the $k$-th **restricted isometry constant** $\delta_k$ of an $m \times n$ matrix $\boldsymbol{\Phi}$, equals the smallest constant that satisfies

$$(1 - \delta_k)||\boldsymbol{\alpha}||_2^2 \leq ||\boldsymbol{\Phi}\boldsymbol{\alpha}||_2^2 \leq (1 + \delta_k)||\boldsymbol{\alpha}||_2^2, \quad (2)$$

for any $k$-sparse $\boldsymbol{\alpha}$ in $\mathbb{R}^n$. The following well-known recovery guarantee is stated w.r.t. $\delta_k$ in (2).

**Theorem A, *c.f.*, [29]** *Let $\boldsymbol{\Phi}$ be the sensing matrix. Let $\boldsymbol{\alpha}$ denote the signal vector. Let $\mathbf{b}$ be the measurements,* i.e., $\mathbf{b} = \boldsymbol{\Phi}\boldsymbol{\alpha}$. *Assume that the $(2k)$-th restricted isometry constant $\delta_{2k}$ of $\boldsymbol{\Phi}$ satisfies $\delta_{2k} < \sqrt{2} - 1$, and further assume that the noisy version $\tilde{\mathbf{b}}$ of $\mathbf{b}$ satisfies $||\tilde{\mathbf{b}} - \mathbf{b}||_2 \leq \epsilon$. Let $\overline{\boldsymbol{\alpha}}_k$ denote the best-$k$ approximation to $\boldsymbol{\alpha}$. Then the $\ell_1$-minimum solution $\boldsymbol{\alpha}^*$ to (1) satisfies*

$$||\boldsymbol{\alpha}^* - \boldsymbol{\alpha}||_1 \leq c_1 ||\boldsymbol{\alpha} - \overline{\boldsymbol{\alpha}}_k||_1 + c_2 \epsilon,$$

*for small constants $c_1 = 4\sqrt{1 + \delta_{2k}}/(1 - \delta_{2k}(1 + \sqrt{2}))$ and $c_2 = 2(\delta_{2k}(1 - \sqrt{2}) - 1)/(\delta_{2k}(1 + \sqrt{2}) - 1)$.*

Theorem A is very powerful, on condition that we know the constants $\delta_k$. But because of their combinatoric nature, computing the restricted isometry constants $\delta_k$ is NP-Hard [13].

Let $\mathcal{S}$ denote a size-$k$ subset of indices. Let $\boldsymbol{\Phi}_{\mathcal{S}}$ denote the size $m \times k$ submatrix of $\boldsymbol{\Phi}$, indexed on (column indices) in $\mathcal{S}$. Let $\sigma^2_{\max}(\boldsymbol{\Phi}_{\mathcal{S}})$ and $\sigma^2_{\min}(\boldsymbol{\Phi}_{\mathcal{S}})$ respectively denote the minimum and maximum, *squared-singular values* of $\boldsymbol{\Phi}_{\mathcal{S}}$. Then from (2) if the columns $\boldsymbol{\phi}_i$ of $\boldsymbol{\Phi}$ are properly normalized, *i.e.*, if $||\boldsymbol{\phi}_i||_2 = 1$, we deduce that $\delta_k$ is the smallest constant in $\mathbb{R}$ that satisfies

$$\delta_k \geq \max(\sigma^2_{\max}(\boldsymbol{\Phi}_{\mathcal{S}}) - 1, 1 - \sigma^2_{\min}(\boldsymbol{\Phi}_{\mathcal{S}})), \qquad (3)$$

for all $\binom{n}{k}$ size-$k$ subsets $\mathcal{S}$. For large $n$, the number $\binom{n}{k}$ is huge. Fortunately $\delta_k$ need not be explicitly computed, if we can estimate it after incorporating *randomization* [1], [11].

Recovery guarantee Theorem A involves "worst-case" analysis. If the inequality (3) is violated for *any* one submatrix $\boldsymbol{\Phi}_{\mathcal{S}}$, then the *whole* matrix $\boldsymbol{\Phi}$ is deemed to have restricted isometry constant larger than $\delta_k$. A common complaint of such "worst-case" analyses is pessimism, *e.g.*, in [20] it is found that for $n = 4000$ and $m = 1000$, the restricted isometry property is not even satisfied for sparsity $k = 5$. This motivates the "average-case" analysis investigated here, where the recovery guarantee is relaxed to hold for a large "fraction" of signals (useful in applications that do not demand all possible signals to be completely recovered). We draw ideas from the statistical StRIP notion used in deterministic CS, which only require "most" of the submatrices $\boldsymbol{\Phi}_{\mathcal{S}}$ to satisfy some properties.

In statistics, a well-known notion of a U-statistic (introduced in the next subsection) is very similar to StRIP. We will show how U-statistics naturally lead to "average-case" analysis.

### B. U-statistics & StRIP

A function $\zeta : \mathbb{R}^{m \times k} \to \mathbb{R}$ is said to be a **kernel**, if for any $\mathbf{A}, \mathbf{A}' \in \mathbb{R}^{m \times k}$, we have $\zeta(\mathbf{A}) = \zeta(\mathbf{A}')$ if matrix $\mathbf{A}'$ can be obtained from $\mathbf{A}$ by *column reordering*. Let $\mathbb{R}_{[0,1]}$ be the set of real numbers bounded below by 0 and above by 1, i.e., $\mathbb{R}_{[0,1]} = \{a \in \mathbb{R} : 0 \leq a \leq 1\}$. U-statistics are associated with functions $g : \mathbb{R}^{m \times k} \times \mathbb{R} \to \mathbb{R}_{[0,1]}$ known as **bounded kernels**. To obtain bounded kernels $g$ from indicator functions, simply use some kernel $\zeta$ and set $g(\mathbf{A}, a) = \mathbb{1}\{\zeta(\mathbf{A}) \leq a\}$ or $g(\mathbf{A}, a) = \mathbb{1}\{\zeta(\mathbf{A}) > a\}$, *e.g.* $\mathbb{1}\{\sigma^2_{\max}(\mathbf{A}) \leq a\}$.

**Definition 1** (Bounded Kernel U-Statistics). *Let $\mathbf{A}$ be a random matrix with $n$ columns. Let $\boldsymbol{\Phi}$ be sampled as $\boldsymbol{\Phi} = \mathbf{A}$. Let $g : \mathbb{R}^{m \times k} \times \mathbb{R} \mapsto \mathbb{R}_{[0,1]}$ be a bounded kernel. For any $a \in \mathbb{R}$, the following quantity*

$$U_n(a) \triangleq \frac{1}{\binom{n}{k}} \sum_{\mathcal{S}} g(\boldsymbol{\Phi}_{\mathcal{S}}, a) \qquad (4)$$

*is a U-statistic of the sampled realization $\boldsymbol{\Phi} = \mathbf{A}$, corresponding to the kernel $g$. In (4), the matrix $\boldsymbol{\Phi}_{\mathcal{S}}$ is the submatrix of $\boldsymbol{\Phi}$ indexed on column indices in $\mathcal{S}$, and the sum takes place over all subsets $\mathcal{S}$ in $\{1, 2, \cdots, n\}$. Note, $0 \leq U_n(a) \leq 1$.*

For $k \leq n$ and positive $u$ where $u \leq 1$, a matrix $\boldsymbol{\Phi}$ has $u$-**StRIP constant** $\delta_k$, if $\delta_k$ is the smallest constant s.t.

$$(1 - \delta_k)||\boldsymbol{\alpha}||^2_2 \leq ||\boldsymbol{\Phi}_S \boldsymbol{\alpha}||^2_2 \leq (1 + \delta_k)||\boldsymbol{\alpha}||^2_2, \qquad (5)$$

for any $\boldsymbol{\alpha} \in \mathbb{R}^k$ and fraction $u$ of size-$k$ subsets $\mathcal{S}$. The difference between (5) and (2) is that $\boldsymbol{\Phi}_S$ is in place of $\boldsymbol{\Phi}$. This StRIP notion coincides with [7]. Consider $\zeta(\mathbf{A}) = \max(\sigma^2_{\max}(\mathbf{A}) - 1, 1 - \sigma^2_{\min}(\mathbf{A}))$ where here $\zeta$ is a kernel.

Obtain a bounded kernel $g$ by setting $g(\mathbf{A}, a) = \mathbb{1}\{\zeta(\mathbf{A}) > a\}$. Construct a U-statistic $U_n(\delta)$ of $\boldsymbol{\Phi}$ the form $U_n(\delta) = \binom{n}{k}^{-1} \sum_{\mathcal{S}} \mathbb{1}\{\zeta(\boldsymbol{\Phi}_{\mathcal{S}}) > \delta\}$. Then if this U-statistic satisfies $U_n(\delta) = 1 - u$, the $u$-StRIP constant $\delta_k$ of $\boldsymbol{\Phi}$ is at most $\delta$, *i.e.*, $\delta_k \leq \delta$.

To exploit apparent similarities between U-statistics and StRIP, we turn to two "average-case" guarantees found in the StRIP literature. In the sequel, the conditions required by these two guarantees, will be analyzed in detail via U-statistics - for now let us recap these guarantees. First, an $\ell_1$-minimization recovery guarantee recently given in [7], is a StRIP-adapted version of the "worst-case" guarantee Theorem A. For any non-square matrix $\mathbf{A}$, let $\mathbf{A}^\dagger$ denote the **Moore-Penrose pseudoinverse**[1]. A vector $\boldsymbol{\beta}$ with entries in $\{-1, 1\}$ is termed a **sign vector**. For $\boldsymbol{\alpha} \in \mathbb{R}^n$, we write $\boldsymbol{\alpha}_{\mathcal{S}}$ for the length-$k$ vector supported on $\mathcal{S}$. Let $\mathcal{S}_c$ denote the complementary set of $\mathcal{S}$, *i.e.*, $\mathcal{S}_c = \{1, 2, \cdots, n\} \setminus \mathcal{S}$. The "average-case" guarantees require us to check conditions on $\boldsymbol{\Phi}$ for fractions of subsets $\mathcal{S}$, or **sign-subset** pairs $(\boldsymbol{\beta}, \mathcal{S})$.

**Theorem B, *c.f.*, Lemma 3, [7]** *Let $\boldsymbol{\Phi}$ be an $m \times n$ sensing matrix. Let $\mathcal{S}$ be a size-$k$ subset, and let $\boldsymbol{\beta} \in \{-1, 1\}^k$. Assume that $\boldsymbol{\Phi}$ satisfies*

- *invertibility: for at least a fraction $1 - u_1$ of subsets $\mathcal{S}$, the condition $\sigma_{\min}(\boldsymbol{\Phi}_{\mathcal{S}}) > 0$ holds.*
- *small projections: for at least a fraction $1 - u_2$ of sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$, the condition*

$$\left| (\boldsymbol{\Phi}^\dagger_{\mathcal{S}} \boldsymbol{\phi}_i)^T \boldsymbol{\beta} \right| \leq a_2 \text{ for every } i \notin \mathcal{S}$$

*holds where we assume the constant $a_2 < 1$.*
- *worst-case projections: for at least a fraction $1 - u_3$ of subsets $\mathcal{S}$, the following condition holds*

$$||\boldsymbol{\Phi}^\dagger_{\mathcal{S}} \boldsymbol{\phi}_i||_1 \leq a_3 \text{ for every } i \notin \mathcal{S}.$$

*Then for a fraction $1 - u_1 - u_2 - u_3$ of sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$, the following error bounds are satisfied*

$$||\boldsymbol{\alpha}^*_{\mathcal{S}} - \boldsymbol{\alpha}_{\mathcal{S}}||_1 \leq \frac{2a_3}{1 - a_2}||\boldsymbol{\alpha} - \overline{\boldsymbol{\alpha}}_k||_1,$$

$$||\boldsymbol{\alpha}^*_{\mathcal{S}_c} - \boldsymbol{\alpha}_{\mathcal{S}_c}||_1 \leq \frac{2}{1 - a_2}||\boldsymbol{\alpha} - \overline{\boldsymbol{\alpha}}_k||_1,$$

*where $\boldsymbol{\alpha}$ is a signal vector that satisfies $\text{sgn}(\boldsymbol{\alpha}_{\mathcal{S}}) = \boldsymbol{\beta}$, and $\overline{\boldsymbol{\alpha}}_k$ is the best-$k$ approximation of $\boldsymbol{\alpha}$ and $\overline{\boldsymbol{\alpha}}_k$ is supported on $\mathcal{S}$, and finally $\boldsymbol{\alpha}^*$ is the solution to (1) where the measurements $\mathbf{b}$ satisfy $\mathbf{b} = \boldsymbol{\Phi}\boldsymbol{\alpha}$.*

For convenience, the proof is provided in Supplementary Material A. The second guarantee is a StRIP-type recovery guarantee for the *LASSO* estimate, based on [6] (also see [7]). Consider recovery from noisy measurements

$$\tilde{\mathbf{b}} = \boldsymbol{\Phi}\boldsymbol{\alpha} + \mathbf{z},$$

here $\mathbf{z}$ is a length-$m$ noise realization vector. We assume that the entries $z_i$ of $\mathbf{z}$, are sampled from a zero-mean Gaussian distribution with variance $c^2_Z$. The LASSO estimate considered in [6], is the optimal solution $\boldsymbol{\alpha}^*$ of the optimization problem

$$\min_{\tilde{\boldsymbol{\alpha}} \in \mathbb{R}^n} \frac{1}{2}||\tilde{\mathbf{b}} - \boldsymbol{\Phi}\tilde{\boldsymbol{\alpha}}||_2 + 2c_Z \cdot \theta_n||\tilde{\boldsymbol{\alpha}}||_1. \qquad (6)$$

---

[1]If $\mathbf{A}$ has full column rank, then $\mathbf{A}^\dagger = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$,

The $\ell_1$-regularization parameter is chosen as a *product* of two terms $c_Z$ and $\theta_n$, where we specify $\theta_n = (1 + a)\sqrt{2 \log n}$ for some positive $a$. What differs from convention is that the regularization depends on the noise standard deviation $c_Z$. We assume $c_Z > 0$, otherwise there will be no $\ell_1$-regularization.

**Theorem C, c.f., [6]** *Let $\boldsymbol{\Phi}$ be the $m \times n$ sensing matrix. Let $\mathcal{S}$ be a size-$k$ subset, and let $\boldsymbol{\beta} \in \{-1, 1\}^k$.*

- *invertability: for at least a fraction $1 - u_1$ of subsets $\mathcal{S}$, the condition $\sigma_{\min}(\boldsymbol{\Phi}_{\mathcal{S}}) > a_1$ holds.*
- *small projections: for at least a fraction $1 - u_2$ of subsets $\mathcal{S}$, same as Theorem B.*
- *invertability projections: for at least a fraction $1 - u_3$ of sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$, the following condition holds*

$$||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\beta}||_\infty \leq a_3.$$

*Let $c_Z$ denote noise standard deviation. Assume Gaussian noise realization $\mathbf{z}$ in measurements $\tilde{\mathbf{b}}$, satisfy*

i) *$||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}^T \mathbf{z}||_\infty \leq (c_Z \sqrt{2 \log n})/a_1$, for the constant $a_1$ in the invertability condition.*

ii) *$||\boldsymbol{\Phi}_{\mathcal{S}_c}^T (\mathbf{I} - \boldsymbol{\Phi}_{\mathcal{S}} \boldsymbol{\Phi}_{\mathcal{S}}^\dagger) \mathbf{z}||_\infty \leq c_Z 2 \sqrt{\log n}$, where $\mathcal{S}_c$ is the complementary set of $\mathcal{S}$.*

*For some positive $a$, assume that constant $a_2$ in the small projections condition, satisfies*

$$(\sqrt{2}(1 + a))^{-1} + a_2 < 1. \quad (7)$$

*Then for a fraction $1 - u_1 - u_2 - u_3$ of sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$, the LASSO estimate $\boldsymbol{\alpha}^*$ from (6) with regularization $\theta_n = (1 + a)\sqrt{2 \log n}$ for the same $a$ above, will successfully recover both signs and supports of $\boldsymbol{\alpha}$, if*

$$|\alpha_i| \geq \left[a_1^{-1} + 2a_3(1 + a)\right] \cdot c_Z \sqrt{2 \log n} \quad \text{for all} \quad i \in \mathcal{S} \quad (8)$$

Because of some differences from [6], we also provide the proof in Supplementary Material A. In [6] it is shown that the noise conditions i) and ii) are satisfied with large probability at least $1 - n^{-1}(2\pi \log n)^{-\frac{1}{2}}$ (see Proposition 4 in Supplementary Material A). Theorem C is often referred to as a *sparsity pattern recovery* result, in the sense that it guarantees recovery of the sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$ belonging to a $k$-sparse signal $\boldsymbol{\alpha}$. Fuchs established some of the earlier important results, see [5], [30], [31].

In Theorems B and C, observe that the *invertability* condition can be easily checked using an U-statistic; simply set the bounded kernel $g$ as $g(\mathbf{A}, a_1) = \mathbb{1}\{\sigma_{\min}(\mathbf{A}) \leq a_1\}$ for some positive $a_1$ and measure the fraction $U_n(a_1) = u_1$. Other conditions require slightly different kernels, to be addressed in upcoming Section IV. But first we first introduce the main U-statistical large deviations theorem (central to our analyses) in the next section.

## III. LARGE DEVIATION THEOREM: "AVERAGE-CASE" BEHAVIOR

Consider two bounded kernels $g$ defined for $\mathbf{A} \in \mathbb{R}^{m \times k}$, corresponding to maximum and minimum squared singular values

$$g(\mathbf{A}, a) = \mathbb{1}\left\{\sigma_{\max}^2(\mathbf{A}) \leq a\right\}, \quad \text{and} \quad (9)$$
$$g(\mathbf{A}, a) = \mathbb{1}\left\{\sigma_{\min}^2(\mathbf{A}) \leq a\right\}. \quad (10)$$



Fig. 1. Gaussian measure. Concentration of U-statistic $U_n(a)$ for squared singular value $\sigma_{\min}^2$ and $\sigma_{\max}^2$ kernels $g$, see (9) and (10). Shown for $m = 25, k = 2$ and two values of $n = 25$ and $100$.

Note that restricted isometry conditions (2) and (5) depend on both $\sigma_{\min}^2$ and $\sigma_{\max}^2$ behaviors, although the conditions in the previous StRIP-recovery guarantees Theorem B are explicitly imposed only on $\sigma_{\min}^2$. See [13], [32] for the different behaviors and implications of these two extremal eigenvalues. In this section we consider two U-statistics, corresponding separately to (9) and (10).

Let $\boldsymbol{A}_i$ denote the $i$-th column of $\boldsymbol{A}$, and assume $\boldsymbol{A}_i$ to be IID. For an bounded kernel $g$, let $p(a)$ denote the expectation $\mathbb{E}g(\boldsymbol{A}_{\mathcal{S}}, a)$, i.e., $p(a) = \mathbb{E}g(\boldsymbol{A}_{\mathcal{S}}, a)$ for any size-$k$ subset $\mathcal{S}$. Since $p(a) = \mathbb{E}U_n(a)$, thus the U-statistic mean $\mathbb{E}U_n(a)$ does not depend on block length $n$.

**Theorem 1.** *Let $\boldsymbol{A}$ be an $m \times n$ random matrix, whereby the columns $\boldsymbol{A}_i$ are IID. Let $g$ be a bounded bounded kernel that maps $\mathbb{R}^{m \times k} \times \mathbb{R} \to \mathbb{R}_{[0,1]}$ and let $p(a) = \mathbb{E}g(\boldsymbol{A}_{\mathcal{S}}, a) = \mathbb{E}U_n(a)$. Let $U_n(a)$ be a U-statistic of the sampled realization $\boldsymbol{\Phi} = \boldsymbol{A}$ corresponding to the bounded kernel $g$. Then almost surely when $n$ is sufficiently large, the deviation $|U_n(a) - p(a)| \leq \epsilon_n(a)$ is bounded by an error term $\epsilon_n(a)$ that satisfies*

$$\epsilon_n^2(a) = 2p(a)(1 - p(a)) \cdot (n/k)^{-1} \log(n/k). \quad (11)$$

Theorem 1 is shown by piecing together (5.5) in [33] and Lemma 2.1 in [34]. The proof is given in Appendix A. Figure 1 empirically illustrates this concentration result for $g$ in (9) and (10), corresponding to $p(a) = \mathbb{E}g(\boldsymbol{A}_{\mathcal{S}}, a) = \Pr\{\sigma_{\max}^2(\boldsymbol{A}_{\mathcal{S}}) \leq a\}$ and $p(a) = \Pr\{\sigma_{\min}^2(\boldsymbol{A}_{\mathcal{S}}) \leq a\}$. Empirical simulation of restricted isometries is very difficult, thus we chose small values $k = 2$, $m = 25$ and block lengths $n = 25$ and $n = 100$. For $n = 25$ the deviation $|U_{25}(a) - p(a)|$ is very noticeable for all values of $a$ and both $\sigma_{\max}^2$ and $\sigma_{\min}^2$. However for larger $n = 100$, the deviation $|U_{100}(a) - p(a)|$ clearly becomes much smaller. This is predicted by vanishing error $\epsilon_n(a)$ given in Theorem 1, which drops as the ratio $n/k$ increases. In fact if $k$ is kept constant then the error behaves as $\mathcal{O}(n^{-1} \log n)$.

Table I reproduces[2] a sample of (asymptotic) estimates for both $\sigma_{\max}^2$ and $\sigma_{\min}^2$ cases, taken from [21]. These estimates are derived for "worst-case" analysis, under assumption that every entry $A_{ij}$ of $\boldsymbol{A}$ is IID and Gaussian distributed (i.e., $A_{ij}$ is Gaussian with variance $1/m$). Table I presents the estimates

---

[2]We point out that Bah actually defined two separate restricted isometry constants, each corresponding to $\sigma_{\min}^2$ and $\sigma_{\max}^2$ in [21]. In this paper to coincide the presentation with our discussion on squared singular values, their results will be discussed in the domain of $\sigma_{\min}^2$ and $\sigma_{\max}^2$.

Fig. 2. Means $p(a) = \mathbb{E}U_n(a)$ for predicting the concentration of $U_n(a)$. Shown for the Gaussian case, (a) $m = 50$ and (b) $m = 150$.

TABLE I
ASYMPTOTIC LOWER AND UPPER BOUNDS ON "WORST-CASE"
EIGENVALUES, [21]

|  |  | Minimum: $\sigma_{min}^2$ | | | Maximum: $\sigma_{max}^2$ | | |
|---|---|---|---|---|---|---|---|
|  |  | $m/n$ | | | $m/n$ | | |
|  |  | 0.1 | 0.3 | 0.5 | 0.1 | 0.3 | 0.5 |
| $k/m$ | 0.1 | 0.095 | 0.118 | 0.130 | 3.952 | 3.610 | 3.459 |
|  | 0.2 | 0.015 | 0.026 | 0.034 | 5.587 | 4.892 | 4.535 |
|  | 0.3 | 0.003 | 0.006 | 0.010 | 6.939 | 5.806 | 5.361 |

according[3] to fixed ratios $k/m$ and $m/n$. To compare, Figure 2 shows the expectations $p(a) = \mathbb{E}U_n(a)$. The values $p(a)$ are interpreted as fractions, and as $n/k$ becomes large $p(a)$ is approached by $U_n(a)$ within a stipulated error $\epsilon_n$. Figure 2 is empirically obtained, though note that in Gaussian case for $p(a)$ we also have exact expressions [32], [35], and the *Bartlett decomposition* [36], available. Again $p(a)$ is a marginal quantity (*i.e.* does not depend on $n$) and simulation is reasonably feasible. In the spirit of non-asymptotics, we consider relatively small $k, m$ values as compared to other works [20], [21]; these adopted values are nevertheless "practical", in the sense they come an implementation paper [17].

Differences are apparent from comparing "average-case" (Figure 2) and "worst-case" (Table I) behavior. Consider $k/m = 0.3$ where Table I shows for all undersampling ratios $m/n$, the worst-case estimate of $\sigma_{min}^2$ is very small, approximately 0.01. But for fixed $m = 50$ and $m = 150$, Figures 2(a) and (b) show that for respectively $k = 0.3 \cdot (150) = 15$ and $k = 45$, a large fraction of subsets $\mathcal{S}$ seem to have $\sigma_{min}^2(\boldsymbol{\Phi}_{\mathcal{S}})$ lying above 0.1. From Table I, the estimates for $\sigma_{min}^2$ gets worse (*i.e.*, gets smaller) as $m/n$ decreases. But the error $\epsilon_n(a)$ in Theorem 1 vanishes with larger $n/k$. For the other $\sigma_{max}^2$ case, we similarly observe that the values in Table I also appear more "pessimistic".

We emphasize that Theorem 1 holds regardless of distribution. Figure 3 is the counterpart figure for Bernoulli and Uniform cases (*i.e.*, each entry $A_{ij}$ is respectively drawn uniformly from $\{-1/\sqrt{m}, 1/\sqrt{m}\}$, or $\{a \in \mathbb{R} : |a| \le \sqrt{3/m}\}$), shown for $m = 50$. Minute differences are seen when comparing with previous Figure 2. For $k = 3$, we observe the fraction $p(a)$ corresponding to $\sigma_{max}^2$ to be roughly 0.95 in the latter case, whereas in the former we have roughly 0.9 in Figure 3(a), and 0.88 in Figure 3(b).

---

[3]The analysis in [21] was performed for the large limit of $k, m$ and $n$, where both $k/m$ and $m/n$ approach fixed constants.



Fig. 3. Means $p(a) = \mathbb{E}U_n(a)$ for $m = 50$ and the (a) Bernoulli and (b) Uniform cases.

**Remark 1.** *Exponential bounds on* $\mathrm{Pr}\{\min_{\mathcal{S}} \sigma_{min}^2(\boldsymbol{A}_{\mathcal{S}}) < 1 - \delta\}$ *and* $\mathrm{Pr}\{\max_{\mathcal{S}} \sigma_{max}^2(\boldsymbol{A}_{\mathcal{S}}) > 1 + \delta\}$ *for* $\max(\delta, \sqrt{k/m}) < \sqrt{2} - 1$, *see (3), employed in "worst-case" analyses, give the optimal* $m = \mathcal{O}(k \log(n/k))$ *rate, see [1], [12], [37]. However the implicit constants are inherently not too small (i.e., these constants cannot be improved).*

These comparisons motivate "average-case" analysis. Marked out on Figures 2 and 3 are the ranges for which $\sigma_{max}^2$ and $\sigma_{min}^2$ must lie to apply Theorem A ("worst-case" analysis). In the cases shown above, the observations are somewhat disappointing - even for small $k$ values, a substantial fraction of eigenvalues lie outside of the required range. Thankfully, there exist "average-case" guarantees, *e.g.*, previous Theorems B and C, addressed in the next section.

## IV. U-STATISTICS & "AVERAGE-CASE" RECOVERY GUARANTEES

### A. Counting argument using U-statistics

Previously we had explained how the *invertability* conditions required by Theorems B and C naturally relate to U-statistics. We now go on to discuss the other conditions, whereby the relationship may not be immediate. We begin with the *projections* conditions, in particular the *worst-case projections* condition. For given $\boldsymbol{\Phi}$, we need to *upper bound* the fraction of subsets $\mathcal{S}$, for which there exists *at least one* column $\boldsymbol{\phi}_j$ where $j \notin \mathcal{S}$, such that $||\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\phi}_j||_{\infty}$ exceeds some value $a$. To this end, let $\mathcal{R}$ denote a size-$(k+1)$ subset, and $\mathcal{R} \setminus \{j\}$ is the size-$k$ subset excluding the index $j$. Consider the bounded kernel $g : \mathbb{R}^{m \times (k+1)} \times \mathbb{R} \mapsto \mathbb{R}_{[0,1]}$ set as

$$g(\mathbf{A}, a) = \frac{1}{k+1} \sum_{j=1}^{k+1} \mathbb{1}\left\{||\mathbf{A}_{\mathcal{R} \setminus \{j\}}^{\dagger} \mathbf{a}_j||_{\infty} > a\right\}, \quad (12)$$

where here $\mathcal{R} = \{1, 2, \cdots, k+1\}$, and $\mathbf{a}_j$ denotes the $j$-th column of $\mathbf{A}$. Consider the U-statistic with bounded kernel (12). We claim that

$$(n-k) \cdot U_n(a)$$
$$= \frac{n-k}{(k+1)\binom{n}{k+1}} \sum_{\mathcal{R}} \sum_{j \in \mathcal{R}} \mathbb{1}\left\{||\boldsymbol{\Phi}_{\mathcal{R} \setminus \{j\}}^{\dagger} \boldsymbol{\phi}_j||_{\infty} > a\right\},$$
$$= \frac{1}{\binom{n}{k}} \sum_{\mathcal{S}} \sum_{j \notin \mathcal{S}} \mathbb{1}\left\{||\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger} \boldsymbol{\phi}_j||_{\infty} > a\right\},$$

where the summations over $\mathcal{R}$ and $\mathcal{S}$ are over all size-$(k+1)$ subsets, and all size-$k$ subsets, respectively. The first equality

follows from Definition 1 and (12). The second equality requires some manipulation. First the coefficient $\binom{n}{k}^{-1}$ follows from the binomial identity $\binom{n}{k+1} \cdot (k+1) = \binom{n}{k} \cdot (n-k)$. Next for some subset $\mathcal{S}$ and index $j$, write the indicator $\mathbb{1}\left\{ \|\mathbf{\Phi}_{\mathcal{S}}^{\dagger}\phi_j\|_{\infty} > a \right\}$ as $\mathbb{1}_{\mathcal{S},j}$ for brevity's sake. By similar counting that proves the previous binomial identity, we argue $\sum_{\mathcal{R}} \sum_{j \in \mathcal{R}} \mathbb{1}_{\mathcal{R}\setminus\{j\},j} = \sum_{\mathcal{S}} \sum_{j \notin \mathcal{S}} \mathbb{1}_{\mathcal{S},j}$, which then proves the claim. Imagine a grid of "pigeon-holes", indexed by pairs $(\mathcal{S}, j)$, where $j \notin \mathcal{S}$. For each size-$(k+1)$ subset $\mathcal{R}$, we assign $k+1$ indicators $\mathbb{1}_{\mathcal{R}\setminus\{j\},j}$ to $k+1$ pairs $(\mathcal{S}, j)$. No "pigeon-hole" gets assigned more than once. In fact we infer from the binomial identity, that every "pigeon-hole" is in fact assigned exactly once, and argument is complete.

Similarly for the *small projections* condition, we define a different bounded kernel $g : \mathbb{R}^{m \times (k+1)} \times \mathbb{R} \mapsto \mathbb{R}_{[0,1]}$ as

$$g(\mathbf{A}, a) = \frac{1}{2^k(k+1)} \sum_{\ell=1}^{2^k} \sum_{j=1}^{k+1} \mathbb{1}\left\{ \left|(\mathbf{A}_{\mathcal{R}\setminus\{j\}}^{\dagger}\mathbf{a}_j)^T \boldsymbol{\beta}_\ell\right| > a \right\}, \quad (13)$$

where $\mathcal{R} = \{1, 2, \cdots, k+1\}$, and $\mathbf{a}_j$ denotes the $j$-th column of $\mathbf{A}$, and $\boldsymbol{\beta}_1, \boldsymbol{\beta}_2, \cdots, \boldsymbol{\beta}_{2^k}$ enumerate all $2^k$ unique sign-vectors in the set $\{-1, 1\}^k$. By similar arguments as before, we can show for the U-statistic $U_n(a)$ of $\bar{\mathbf{\Phi}}$ corresponding to the bounded kernel (13) satisfies

$$(n-k) \cdot U_n(a) = \frac{1}{2^k \binom{n}{k}} \sum_{\ell=1}^{2^k} \sum_{\mathcal{S}} \sum_{j \notin \mathcal{S}} \mathbb{1}\left\{ \left|(\mathbf{\Phi}_{\mathcal{S}}^{\dagger}\phi_j)^T \boldsymbol{\beta}_\ell\right| > a \right\},$$

For indicators $\mathbb{1}_{\mathcal{S},j}$, note that $\sum_{j \notin \mathcal{S}} \mathbb{1}_{\mathcal{S},j} \geq 1$ if *at least one* indicator satisfying $\mathbb{1}_{\mathcal{S},j} = 1$, and we proved the following.

**Proposition 1.** *Let $U_n(a_3)$ be the U-statistic of $\bar{\mathbf{\Phi}}$, corresponding to the bounded kernel $g(\mathbf{A}, a_3)$ in (12). Then the fraction of subsets $\mathcal{S}$ of size-$k$, for which the worst-case projections condition is violated for some $a_3 \in \mathbb{R}$, is at most $(n-k) \cdot U_n(a_3)$. Similarly if $U_n(a_2)$ corresponds to $g(\mathbf{A}, a_2)$ in (13), the fraction sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$, for which the small projections condition is violated for some $a_2 \in \mathbb{R}$, is at most $(n-k) \cdot U_n(a_2)$.*

Referring back to Theorem B, we point out that the *small projections* condition is more stringent than the *worst-case projections* condition. We mean the following: in the former case, the value $a_2$ must be chosen such that $a_2 < 1$; in the latter case, the value $a_3$ is allowed to be larger than 1, its size only affects the constant $2a_3/(1-a_2)$ appearing in the error estimate $\|\boldsymbol{\alpha}_{\mathcal{S}}^* - \boldsymbol{\alpha}_{\mathcal{S}}\|_1$. In fact if the signal $\boldsymbol{\alpha}$ is $k$-sparse, then $\|\boldsymbol{\alpha} - \overline{\boldsymbol{\alpha}}_k\|_1 = 0$ and the size of $a_3$ is inconsequential, *i.e.*, the *worst-case projections* condition is not required in this special case. In this special case, it is best to set $a_2 = 1 - \epsilon$ for some arbitrarily small $\epsilon$. Theorem B is in fact a stronger version of Fuchs' early work on $\ell_0/\ell_1$-*equivalence* [5]. In the same respect, Donoho & Tanner also produced early seminal results from counting faces of random polytopes [22], [23].

Figure 4 shows empirical evidence, where the $k, m, n$ values are inspired by practical system sizes taken from an implementation paper [17]. These experiments consider $\bar{\mathbf{\Phi}}$ sampled from Gaussian matrices $\mathbf{A}$, *exactly $k$-sparse* signals with non-zero $\alpha_i$ sampled from $\{-1, 1\}$, and uses $\ell_1$-minimization recovery

(1). Figure 4(a) plots simulated (sparsity pattern recovery) results for 3 measurement sizes $m = 50, 100$ and $150$ and block sizes $n \geq 200$ and $n \leq 3000$. For example the contour marked "0.1", delineates the $k, n$ values for which recovery fails for a 0.1 fraction of (random) sparsity patterns (sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$). We examine the U-statistic $U_n(a_2)$ with kernel (13), related to the small projections condition. Since $\mathbf{A}$ has Gaussian distribution, we set $a_2 = 1$ in the kernel $g(\mathbf{A}, a_2)$, as $\Pr\{(\mathbf{A}_i^{\dagger}\mathbf{A}_i)^T \boldsymbol{\beta} = 1\} = 0$ for any $(\boldsymbol{\beta}, \mathcal{S})$ and $j \notin \mathcal{S}$. Figure 4(b) plots the expectation $(n-k) \cdot p(1)$, where $p(1) = \mathbb{E} U_n(1) = \mathbb{E} g(\mathbf{A}_{\mathcal{R}}, 1)$ for any size-$(k+1)$ subset $\mathcal{R}$. Again the contour marked "0.1", delineates the $k, n$ values for which $(n-k) \cdot p(1) = 0.1$. Here the values $p(1)$ are empirical. We observe that both Figures 4(a) and (b) are remarkably close for fractions 0.5 and smaller. Figures 4(c) incorporates the large deviation error $\epsilon_n$ given in Theorem 1 (in doing so, we assume $n$ sufficiently large). The bound is still reasonably tight for fractions $\leq 0.5$. Comparing with recent Donoho & Tanners' (also "average-case") results for $\ell_1$-recovery (for only the noiseless case), taken from [23]. For fractions 0.5 and 0.01, we observe that for system parameters $m = 50$ and $n \leq 1000$ (chosen in hardware implementation [17]), we do not obtain reasonable predictions. For $m = 100$, the bounds [23] work only for very small block lengths $n \leq 300$. The only reasonable case here is $m = 150$, where the bounds [23] perform better than ours only for lengths $n \leq 400$ (*i.e.*, Figure 4(c) shows that for $n = 300$, the large deviation bounds predict a 0.01 fraction of size $k = 5$ unrecoverable sparsity patterns, but [23] predict a 0.01 fraction of size $k = 11$ unrecoverable sparsity patterns).

The above experiments suggest the deviation error $\epsilon_n(a)$ in Theorem 1 to be over-conservative. Fortunately in the next two subsections (pertaining to U-statistics treastise of $\ell_1$-recovery Theorem B (Section IV-B), and LASSO recovery Theorem C (Subsection IV-C)), this conservative-ness does not show up from a rate standpoint (it only shows up in implicit constants). In fact by empirically "adjusting" these constants, we find good measurement rate predictions (akin to moving from Figure 4(c) to (b)).

### B. Rate analysis for $\ell_1$-recovery (Theorem B)

In "worst-case" analysis, it is well-known that it is sufficient to have measurements $m$ on the order of $k \log(n/k)$, in order to have the restricted isometry constants $\delta_k$ defined by (2), satisfy the conditions in Theorem A. We now go on to show that for "average-case", a similar expression for this rate can be obtained. To this end we require tail bounds on salient quantities. Such bounds have been obtained for the *small projections* condition, see [6], [7], [25], where typically an equiprobable distribution is assumed over the sign-vectors $\boldsymbol{\beta}_\ell$. To our knowledge these techniques were born from considering deterministic matrices. Since $\bar{\mathbf{\Phi}}$ is randomly sampled here, we proceed slightly differently (though essentially using similar ideas) without requiring this random signal model. For simplicity, the bound assumes zero mean matrix entries, either i) Gaussian or ii) bounded.

Fig. 4. Gaussian case. Comparing $(a)$ empirical results for $\ell_1$-minimization recovery, $(b)$ mean parameter $(n-k)\cdot p(1)$ (empirically obtained), and $(c)$ after accounting for large deviations (Thm. 1). We show cases $m = 50, 100$ and $150$. We also compare with Donoho & Tanners' (DT) large deviation bounds [23].

**Proposition 2.** *Let $A$ be an $m \times n$ random matrix, whereby its columns $A_i$ are identically distributed. Assume every entry $A_{ij}$ of $A$ has zero mean, i.e., $\mathbb{E}A_{ij} = 0$. Let every $A_{ij}$ be either i) Gaussian with variance $1/m$, or ii) bounded RVs satisfying $|A_{ij}| \leq 1/\sqrt{m}$. Let the rows $[A_{i1}, A_{i2}, \cdots, A_{in}]$ of $A$ be IID.*

*Let $S$ be a size-$k$ subset, and let index $\omega$ be outside of $S$, i.e., $\omega \notin S$. Then for any sign vector $\beta$ in $\{-1, 1\}^k$, we have*

$$\Pr\left\{\left|(A_S^\dagger A_\omega)^T \beta\right| > a\right\} \leq 2\exp\left(-\frac{ma^2\delta}{2k}\right) + \Pr\{\sigma_{\min}^2(A_S) \leq \delta\} \quad (14)$$

*for any positive $\delta \in \mathbb{R}$.*

*Proof:* For $\tau \in \mathbb{R}$, let $\mathcal{E}(\tau) = \{\beta^T(A_S^T A_S)^\dagger\beta \leq \tau\}$ where $\mathcal{E}(\tau)$ is an probabilistic event. Let $\mathcal{E}_c(\tau)$ denote the complementary event. Bound the probability as

$$\Pr\left\{\left|(A_S^\dagger A_\omega)^T\beta\right| > a\right\} \leq \Pr\left\{\left|(A_S^\dagger A_\omega)^T\beta\right| > a \,\Big|\, \mathcal{E}(\tau)\right\} + \Pr\{\mathcal{E}_c(\tau)\}. \quad (15)$$

We upper bound the first term as follows. Denote constants $c_1, c_2, \cdots, c_m$. For entries $(A_\omega)_i$ of $A_\omega$, consider the sum $\sum_{i=1}^m c_i \cdot (m^{-\frac{1}{2}}A_\omega)_i = \frac{1}{m}\sum_{i=1}^m c_i X_i$, where RVs $X_i$ satisfy $X_i = (\sqrt{m}A_\omega)_i$. By standard arguments (see Supplementary Material B) we have the double-sided bound $\Pr\{|\sum_{i=1}^m c_i X_i| > mt\} \leq 2\exp\left(-(mt)^2/(2\cdot||\mathbf{c}||_2^2)\right)$, where vector $\mathbf{c}$ equals $[c_1, c_2, \cdots, c_m]^T$.

Next write $(A_S^\dagger A_\omega)^T\beta = (\sqrt{m}\cdot\beta^T A_S^\dagger)(m^{-\frac{1}{2}}A_\omega)$. When conditioning on $\beta^T A_S^\dagger$, then $\sqrt{m}\cdot\beta^T A_S^\dagger$ is fixed, say equals some vector $\mathbf{c}$. Put $X_i = (\sqrt{m}A_\omega)_i$ and $X_i$'s are independent (by assumed independence of the rows of $A$). Then use the above bound for $\Pr\{\sum_{i=1}^m c_i X_i > t\}$, set $t = a$ and conclude

$$\Pr\left\{\left|(A_S^\dagger A_\omega)^T\beta\right| > a \,\Big|\, \beta^T A_S^\dagger\right\}$$
$$\leq 2\exp\left(-\frac{(ma)^2}{2m||\beta^T A_S^\dagger||_2^2}\right) = 2\exp\left(-\frac{ma^2}{2\cdot\beta^T(A_S^T A_S)^\dagger\beta}\right), \quad (16)$$

where the last equality follows from the identity $A_S^\dagger(A_S^\dagger)^T = (A_S^T A_S)^\dagger$. Further conclude that the first term in (15) is bounded by $2\exp(-ma^2/(2\tau))$, due to further conditioning on the event $\mathcal{E}(\tau) = \{\beta^T(A_S^T A_S)^\dagger\beta \leq \tau\}$.

To bound the second term, let $\varsigma_{\max}(\mathbf{A})$ denote the maximum eigenvalue of matrix $\mathbf{A}$. Since $A_S^T A_S$ is positive semidefinite, note that $\beta^T(A_S^T A_S)^\dagger\beta$ is upper bounded by $||\beta||_2^2 \cdot \varsigma_{\max}((A_S^T A_S)^\dagger)$, which equals $k \cdot \varsigma_{\max}((A_S^T A_S)^\dagger)$. Furthermore $\varsigma_{\max}((A_S^T A_S)^\dagger) \leq 1/\sigma_{\min}^2(A_S)$, where here $\sigma_{\min}(\mathbf{A})$ is the minimum singular value of $\mathbf{A}$. Thus $\Pr\{\mathcal{E}_c(\tau)\} \leq \Pr\{k/\sigma_{\min}^2(A_S) > \tau\}$. Finally put $\tau = \delta k$ to get $\Pr\{\mathcal{E}_c(\tau)\} \leq \Pr\{\sigma_{\min}^2(A_S) \leq \delta^{-1}\}$. ∎

Proposition 2 is used as follows. First recall that previous Proposition 1 allows us to upper bound the fraction $u_2$ of sign-subset pairs $(\beta, S)$ *failing* the *small projections* condition, with the (scaled) U-statistic $(n-k)\cdot U_n(a_2)$ with kernel $g$ in (13) and $|S| = k$. By Theorem 1 the quantity $(n-k)\cdot U_n(a_2)$ concentrates around $(n-k)\cdot p(a_2)$, where $p(a_2) = \mathbb{E}g(A_\mathcal{R}, a_2)$, where $g$ in (13) is defined for size-$(k+1)$ subsets $\mathcal{R}$. We use Proposition 2 to upper estimate $p(a_2)$ using the RHS of (14). Indeed verify that $p(a_2) = 2^{-k}\sum_\ell \Pr\{|(A_S^\dagger A_\omega)^T\beta_\ell| > a_2\}$ for any $S$ and $\omega \notin S$, and the bound (14) holds for any $\beta = \beta_\ell$. Now $p(a_2)$ is bounded by two terms. By $u_2 \leq (n-k)\cdot U_n(a_2)$, thus to have $u_2$ small, we should have the (scaled) first term $2(n-k)\cdot\exp(-ma_2^2\delta/(2k))$ of (14) to be at most some small fraction $u$. This requires

$$m \geq \text{const}\cdot k\log\left(\frac{n-k}{u}\right) \quad (17)$$

with const $= 2/(a_2^2\delta)$ (and we dropped an insignificant $\log 2$ term). Next, for $m \geq 2k$ and $\delta < (0.29)^2$, we can bound[4] the second term $\Pr\{\sigma_{\min}^2(A_S) \leq \delta\}$ of (14)

---

[4]For $m \geq 2k$, we have $\Pr\{\sigma_{\min}(\mathbf{A}) < c\cdot 0.29 - t\} \leq \Pr\{\sigma_{\min}(\mathbf{A}) < 1 - c\cdot\sqrt{k/m} - t\} \leq \exp(-mt^2/c_1)$ for some constants $c, c_1$, where $\mathbf{A}$ has size $m \times k$ and with proper column normalization. For simplicity we drop the constant $c$ in this paper; one simply needs to add $c$ in appropriate places in the exposition. In particular for the Gaussian and Bernoulli cases $c = 1$, and $c_1 = 2$ and $c_1 = 16$, respectively, see Theorem B, [28].

by $\exp(-m \cdot (0.29 - \sqrt{\delta})^2/c_1)$ where $c_1$ is some constant, see [27], Theorem 5.39. Roughly speaking, $\sigma^2_{\min}(\boldsymbol{A}_\mathcal{S}) \geq 0.29$ with "high probability". Figures 2 and 3 (in the previous Section III) empirically support this fact. Again to have $u_2$ small the second term of (14) must be small. This requires $(n-k) \cdot \exp(-m \cdot (0.29 - \sqrt{\delta})^2/c_1) \leq u$ for some small fraction $u$, in which it suffices to have $m$ satisfy (17) with $\mathrm{const} = c_1/(0.29 - \sqrt{\delta})^2$.

For the *invertability* condition in Theorem B, we also need to upper bound the corresponding fraction $u_1$ of size-$k$ subsets $\mathcal{S}$. We simply use an U-statistic $U_n(a_1)$ with kernel $g(\mathbf{A}, a_1) = \mathbb{1}\{\sigma_{\min}(\mathbf{A}) > a_1\}$ for some positive $a_1$ (see also Theorem C). Here Proposition 1 is not needed. To make $p(a_1)$ small, where $p(a_1) = \mathbb{E}g(\boldsymbol{A}_\mathcal{S}, a_1)$, use the previous bound $p(a_1) \leq \exp(-m \cdot (0.29 - a_1)^2/c_1)$, where we set $a_1 = \sqrt{\delta}$ with $a_1 \leq 0.29$. Clearly $p(a_1)$ cannot exceed some fraction $u$, if $m$ satisfies (17) with $\mathrm{const} = c_1/(0.29 - a_1)^2$.

For the time being consider *exactly* $k$-sparse signals $\boldsymbol{\alpha}$. In this special case the *worst-case projections* condition in Theorem B is superfluous (*i.e.*, with no consequence $a_3$ can be arbitrarily big) - only *invertability* and *small projections* conditions are needed. While we have yet to consider the large deviation error $\epsilon_n(a)$ from Theorem 1, doing so will not drastically change the rate. For $U_n(a)$ with kernel $g$ and $p(a)$, where $p(a) = \mathbb{E}g(\boldsymbol{A}, a)$, almost surely

$$U_n(a) \leq p(a) + \epsilon_n(a) \leq (p(a))^{\frac{1}{2}} + \sqrt{2p(a)\omega^{-1}\log\omega}$$
$$\leq (p(a))^{\frac{1}{2}}\left(1 + \sqrt{2\omega^{-1}\log\omega}\right) \quad (18)$$

where the second inequality follows because $p(a) \leq 1$, and by setting $\omega = n/k$. Taking log of the RHS, we obtain $(1/2)\log p(a) + \log(1 + \sqrt{2\omega^{-1}\log\omega})$. Note $\log(1 + \sqrt{2\omega^{-1}\log\omega}) \leq \sqrt{2\omega^{-1}\log\omega}$, since $\log(1 + \alpha) \leq \alpha$ holds for all positive $\alpha$. For the *small projections* condition, bound $(p(a))^{\frac{1}{2}}$ by the sum of the square-roots of each term in (14). Then to have $u_2 \leq (n-k) \cdot U_n(a_2) \leq 2u$, it follows similarly as before that it suffices that (see Supplementary Material C)

$$m \geq \mathrm{const} \cdot k \left[\log\left(\frac{n-k}{u}\right) + \sqrt{2 \cdot (k/n)\log(n/k)}\right] \quad (19)$$

with $\mathrm{const} = \max(4/(a_2^2\delta), 2c_1/(0.29 - \sqrt{\delta})^2)$ where we had set $\sqrt{\delta} = a_1$ (we dropped an insignificant $\log 2$ term). For *invertability* condition do the same. To have $u_1 = U_n(a_1) \leq u$ it suffices that $m$ satisfies (19) with the same const. Observe that the term $\sqrt{2 \cdot (k/n)\log(n/k)}$ is at most 1, and vanishes with high undersampling (small $k/n$). Hence (17) and (19) are similar from a rate standpoint.

We conclude the following: for exactly $k$-sparse signals the rate (19) suffices to recover at least $1 - 3u$ fraction of sign-subset $(\boldsymbol{\beta}, \mathcal{S})$ pairs. While const in (19) must be at least 4 (recall that Figure 4(c) was somewhat pessimistic), for matrices with Gaussian entries we empirically find that const is inherently smaller, whereby const $\approx 1.8$. This is illustrated in Figure 5, for two fractions 0.1 and 0.01 of unrecoverable sign-subset pairs. We observe good match with simulation results shown in the previous Figure 4(a), and



Fig. 5. Measurement rates predicted by equation (19), with const taken to equal 1.8, required to recover at least $1 - 3u = 0.9$ and 0.99 fractions of sign-subset pairs $(\boldsymbol{\beta}, \mathcal{S})$ (when the signal is exactly $k$-sparse), shown respectively in (a) and (b).

quantities[5] $(n-k) \cdot p(1)$ plotted in Figure 4(b). For example, $m = 150$ suffices for a 0.01 fractional recovery failure, for $n = 300 \sim 1000$ and $k = 6 \sim 7$, and for 0.1 fraction then $k = 7 \sim 10$. We conjecture possible improvment for const.

In the more general setting for *approximately* $k$-sparse signals, we can also have rate (19). To see this, observe that Proposition 2 also delivers an exponential bound for the *worst-case projections* condition, see (12). This is because $\|\boldsymbol{A}_\mathcal{S}^\dagger \boldsymbol{A}_\omega\|_1 = \max_{\ell: 1 \leq \ell \leq 2^k} |(\boldsymbol{A}_\mathcal{S}^\dagger \boldsymbol{A}_\omega)^T \boldsymbol{\beta}_\ell|$, and we take a union bound over $2^k$ terms. Set $a_3 = a_2\sqrt{k}$, where $a_2$ and $a_3$ respectively correspond to *small projections* and *invertability* conditions. Then we proceed similarly as before (see Supplementary Material C) to show[6] that the rate for recovering at least $1 - 5u$ fraction of $(\boldsymbol{\beta}, \mathcal{S})$ pairs suffices to be (19). The following is the main result summarizing the exposition so far.

**Theorem 2.** *Let* $\boldsymbol{\Phi}$ *be an* $m \times n$ *matrix, where assume* $n$ *sufficiently large for Theorem 1 to hold. Sample* $\boldsymbol{\Phi} = \boldsymbol{A}$ *whereby the entries* $A_{ij}$ *are IID, and are Gaussian or bounded (as stated in Proposition 2). Then all three conditions in* $\ell_1$-*recovery guarantee Theorem B for* $(\boldsymbol{\beta}, \mathcal{S})$ *with* $|\mathcal{S}| = k$, *with the invertability condition taken as* $\sigma_{\min}(\boldsymbol{\Phi}_\mathcal{S}) \geq a_1$ *with* $a_1 \leq 0.29$. *and with* $a_3 = a_1\sqrt{k}$, *are satisfied for* $u_1 + u_2 + u_3 = 5u$ *for some small fraction* $u$, *if* $m$ *is on the order of (19) with* $\mathrm{const} = \max(4/(a_1a_2)^2, 2c_1/(0.29 - a_1)^2)$, *and* $c_1$ *depends on the distribution of* $A_{ij}$'s. *Note* $\mathrm{const} \geq 4$.

*In the exactly* $k$-sparse case where only the first 2 conditions are required, this improves to $u_1 + u_2 = 3u$.

We end this subsection with two comments on the rate (19) derived here for "average-case" analysis. Firstly (19) is very similar to that of $k\log(n/k)$ for "worst-case" analysis. This

---

[5]Comparing (19) and (17) and the respective expressions for const, dropping const from 4 to 1.8 is akin to ignoring the deviation error $\epsilon_n(a)$. This, and as Figure 4 suggests, the U-statistic "means" $(n-k) \cdot p(1)$ seem to predict recovery remarkably well, with similar rates to (19), and inherent const smaller than that derived here.

[6]We used an assumption that $(n-k)/u$ is suitably larger than 2, see Supplementary Material C.

Fig. 6. Empirical LASSO recovery performance, Bernoullli case. In $(a)$ the non-zero signal magnitudes $|\alpha_i|$ equal 1, and in $(b)$ they are in $\mathbb{R}_{[0,1]}$. Noise variances denoted $c_Z^2$.

justifies the counting employed in previous Subsection IV-A, Proposition 1, and is reassuring since we know that "worst-case" analysis provides the optimal rate [1], [11]. Secondly to have (19) hold for the approximately $k$-sparse case, we lose a factor of $\sqrt{k}$ in the error estimate $||\boldsymbol{\alpha}_S^* - \boldsymbol{\alpha}_S||_1$, as compared to "worst-case" Theorem A. This is because we need to set $a_3 = a_2\sqrt{k}$, as mentioned in the previous paragraph. However, the "average-case" analysis here achieves our primary goal, that is to predict well for system sizes $k, m, n$ when "worst-case" analysis becomes too pessimistic.

### C. Rate analysis for LASSO (Theorem C)

Next we move on to the LASSO estimate of [6]. Recall from (6) that the regularizer depends on the noise standard deviation $c_Z$, and the term $\theta_n = (1+a)\sqrt{2\log n}$ that depends on block length $n$ and some non-negative constant $a$ that we set. This constant $a$ impacts performance [6]. For matrices with Bernoulli entries, Figure 6 shows recovery failure rates for two data sets $m = 50, n = 1000$ and $m = 150, n = 1000$; the sparsity patterns (sign-subset pairs $(\boldsymbol{\beta}, S)$) were chosen at random, and failure rates are shown for various sparsity values $k$, and noises $c_Z$. In Figure 6(a) we set $a = 0$, and in $(b)$ we set $a = 1$. Also, in $(a)$ the non-zero signal magnitudes $|\alpha_i|$ are in $\{1, -1\}$, and in $(b)$ they are in $\mathbb{R}_{[0,1]}$. The performances are clearly different. "Threshold-like" behavior is seen in $(a)$ for both data sets, whereby the performances stay the same for $c_Z$ in the range $5 \times 10^{-2} \sim 1 \times 10^{-4}$, and then catastrophically failing for $c_Z = 1 \times 10^{-1}$. However in $(b)$, for various $c_Z$ the performances seem to be limited by a "noise-floor". We see that in the noiseless limit (more specifically when $c_Z \to 0$),

the performances become the same. In this subsection, we apply U-statistics on the various conditions of Theorem C, in particular the *invertability* and *small projections* conditions have already been discussed in the previous subsection. We account for the observations in Figure 6.

In the noiseless limit, the previously derived rate (19) holds. Here, the regularizer in (6) becomes so small that $a$ (equivalently $\theta_n$) does not matter. As mentioned in [5], LASSO then becomes equivalent to $\ell_1$-minimization (1), hence the (noiseless) performances in Figures 6(a) and $(b)$ are the same. That is, in this special case the rate (19) suffices to recover at least $1 - 3u$ fraction of $(\boldsymbol{\beta}, S)$. To test, take $k = 4$, $n = 3000$, and fraction $1 - 3u = 1 - 6 \times 10^{-6}$, and with const $= 1.8$ gives 153, close to $m$ here which is set to 150.

In the noisy case, we are additionally concerned with the noise conditions i) and ii), conditions (7) and (8), and *invertability projections*. Recall that the noise conditions are satisfied with probability $1 - n^{-1}(2\pi \log n)^{-\frac{1}{2}}$, that goes to 1 superlinearly [6] (Proposition 4, Supplementary Material A). The remaining conditions are influenced by the value $a$ set in the $\theta_n$ regularization term in (6).

In condition (7), the value $a$ sets the maximal value for $a_2$ (when $a = 0$ then $a_2 < 0.2929$, and when $a = 1$ then $a_2 < 0.6464$). This affects the *small projections* condition, to which constant $a_2$ belongs, which in turn affects performance. However from a rate standpoint (19) still holds, only now the value of const (which has the term $4/(a_2^2\delta)$) becomes larger.

In condition (8), the value $a$ affects the size of the term $a_1^{-1} + 2a_3(1 + a)$. The larger $a$ is, the more often (8) fails to satisfy. Here there are two constants $a_1$ and $a_3$. Recall $a_1$ belongs to the *invertability* condition discussed in the previous subsection, which holds with rate (19) with const $= 2c_1/(0.29 - a_1)^2$ and $a_1 \le 0.29$. Consider the case where the non-zero signal magnitudes $|\alpha_i|$ are independently drawn from $\mathbb{R}_{[0,1]}$. Then we observe $(\min_{i \in S} |\alpha_i|) < t$ with probability $1 - (1 - t)^k$ where $t \in \mathbb{R}_{[0,1]}$ and $|S| = k$. For $t$ set equal to the RHS of (8), this gives the probability that condition (8) fails. Figure 6(b) shows good empirical match when setting $a_1 = 0.29$ and $a_3 = 1$, where the dotted curves predict the "error-floors" for various $k$, measurements $m = 50$ and $m = 150$, and noise $c_Z$. In the other case where $|\alpha_i| = 1$ (as in Figure 6(a)), condition (8) remains un-violated as long as $c_Z$ (and $a_1, a_3, n$) allow the RHS to be smaller than 1. Figure 6(a) suggests that for the appropriate choices for $a_1, a_3$, condition (8) is always un-violated when $c_Z \le 5 \times 10^{-2}$, and violated when $c_Z \ge 1 \times 1^{-1}$. For more discussion on noise effects see Supplementary Material D.

The constant $a_3$ belongs to the remaining *invertability projections* condition. The fraction $u_3$ of size-$k$ subsets *failing* the *invertability projections* condition for some $a_3$, can be addressed using U-statistics. Consider the bounded kernel $g : \mathbb{R}^{m \times k} \times \mathbb{R} \to \mathbb{R}_{[0,1]}$, set as

$$g(\mathbf{A}, a) = \frac{1}{2^k} \sum_{\ell=1}^{2^k} \mathbb{1}\left\{(\mathbf{A}^T\mathbf{A})^\dagger \boldsymbol{\beta}_\ell > a\right\} \qquad (20)$$

where $\boldsymbol{\beta}_\ell \in \{-1, 1\}^k$ and $(\mathbf{A}^T\mathbf{A})^\dagger$ is the pseudoinverse of $\mathbf{A}^T\mathbf{A}$. Then $u_3 = U_n(a_3)$, and as before Theorem 1

guarantees the upper bound (18), which depends on $p(a_3)$ where $p(a_3) = \mathbb{E}g(\boldsymbol{A}_{\mathcal{S}}, a_3)$.

We go on to discuss a bound on $p(a_3)$ under some general conditions. In [6], analysis on $p(a_3)$ (see Lemma 3.5) requires $\sigma_{\max}^2(\boldsymbol{A}_{\mathcal{S}}) \leq 1.5$, a condition not explicitly required in Theorem C. Also, empirical evidence suggests not to assume that $\sigma_{\max}^2(\boldsymbol{A}_{\mathcal{S}}) \leq 1.5$. For $m = 150$ and $k = 5$ we see from Figure 6 that (in the noiseless limit) the *failure* rate is on the order of $1 \times 10^{-4}$, but in Figure 2(b) we see $\sigma_{\max}^2(\boldsymbol{A}_{\mathcal{S}}) > 1.5$ occurs with much larger fraction 0.1. Hence we take a different approach. Using ideas behind Bauer's generalization of *Wielandt's inequality* [38], the following proposition allows $\sigma_{\max}^2(\boldsymbol{A}_{\mathcal{S}})$ to arbitrarily exceed 1.5. Also, it does not assume any particular distribution on entries of $\boldsymbol{A}$.

**Proposition 3.** *Let $\mathcal{S}$ be a size-$k$ subset. Assume $k \geq 2$. Let $\boldsymbol{A}_{\mathcal{S}}$ be an $k \times n$ random matrix. Let $\delta_{\min}, \delta_{\max}$ be some positive constants. For any sign vector $\boldsymbol{\beta}$ in $\{-1, 1\}^k$, we have*

$$\Pr\left\{ ||(\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}})^\dagger \boldsymbol{\beta}||_\infty > \frac{(\sqrt{k} + 1) \cdot |\tau_k - 1|}{\delta_{\min}^2 \cdot (\tau_k + 1)} \cdot \right\}$$
$$\leq \Pr\{\mathcal{E}_c(\delta_{\min}, \delta_{\max})\} \qquad (21)$$

*where $\mathcal{E}(\delta_{\min}, \delta_{\max}) = \{\delta_{\min} \leq \sigma_{\min}(\boldsymbol{A}_{\mathcal{S}}) \leq \sigma_{\max}(\boldsymbol{A}_{\mathcal{S}}) \leq \delta_{\max}\}$, and $\mathcal{E}_c(\delta_{\min}, \delta_{\max})$ is the complementary event of $\mathcal{E}(\delta_{\min}, \delta_{\max})$, and the constant $\tau_k$ satisfies*

$$\tau_k = \tau_k(\delta_{\max}, \delta_{\min}) = \left(\frac{\delta_{\max}}{\delta_{\min}}\right)^2 \cdot \frac{1 + k^{-\frac{1}{2}}}{1 - k^{-\frac{1}{2}}}. \qquad (22)$$

We defer the proof for now. If $\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}}$ is "almost" an identity matrix, then we expect $||(\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}})^{-1} \boldsymbol{\beta}||_\infty \approx 1$ for any sign vector $\boldsymbol{\beta}$ (hence our above hueristic whereby we set $a_3 = 1$). Proposition 3 makes a slightly weaker (but relatively general) statement. Now for some appropriately fixed $\delta_{\max}$ and $\delta_{\min}$, we expect $\Pr\{\mathcal{E}_c(\delta_{\min}, \delta_{\max})\}$ in (21) to drop exponentially in $m$. Just as the term $\Pr\{\sigma_{\min}(\boldsymbol{A}_{\mathcal{S}}) \leq \delta_{\min}\}$ in Proposition 2 can be bounded by $\exp(-m \cdot (0.29 - \delta_{\min})^2/c_1)$, we can bound[7] $\Pr\{\sigma_{\max}(\boldsymbol{A}) > \delta_{\max}\} \leq \exp(-m(\delta_{\max} - 1.71)^2/c_1)$ for some $\delta_{\max} \geq 1.71$. Roughly speaking, $\sigma_{\max}(\boldsymbol{A}_{\mathcal{S}}) \leq 1.71$ (or $\sigma_{\max}^2(\boldsymbol{A}_{\mathcal{S}}) \leq 2.92$) with "high probability". We fix $\delta_{\min} = a_1$, where $a_1$ belongs to the *invertability* condition.

So to bound $p(a_3)$, both (20) and Proposition 3 imply $p(a_3) \leq \Pr\{\mathcal{E}_c(\delta_{\min}, \delta_{\max})\}$ for $a_3 = (\sqrt{k} + 1) \cdot |\tau_k - 1|/(\delta_{\min}^2 \cdot (\tau_k + 1))$. Now $\Pr\{\mathcal{E}_c(\delta_{\min}, \delta_{\max})\} \leq 2 \exp(-m \cdot t^2/c_1)$, where we set $t = \delta_{\max} - 1.71 = 0.29 - a_1$ and $\delta_{\min} = a_1$. By (18), the rate (19) suffices to ensure $u_3 = U_n(a_3) \leq u$ for some fraction $u$, with the same const. Thus we proved the other main theorem, similar to Theorem 2.

**Theorem 3.** *Let $\boldsymbol{\Phi}$ be an $m \times n$ matrix, , where assume $n$ sufficiently large for Theorem 1 to hold. Sample $\boldsymbol{\Phi} = \boldsymbol{A}$ whereby the entries $A_{ij}$ are IID, and are Gaussian or bounded (as stated in Proposition 2). Then all three invertability, small projections, and invertability projections conditions in LASSO Theorem C for $(\boldsymbol{\beta}, \mathcal{S})$ with $|\mathcal{S}| = k \geq 2$, with $a_1 \leq 0.29$, with $a_2$ satisfying (7) for some $a$ set in the regularizer $\theta_n$, and with $a_3 = (\sqrt{k} + 1) \cdot |\tau_k - 1|/(a_1^2 \cdot (\tau_k + 1))$ for $\tau_k =$*

[7]For $m \geq 2k$ we have $\Pr\{\sigma_{\max}(\boldsymbol{A}) > 1.71 + t\} \leq \Pr\{\sigma_{\max}(\boldsymbol{A}) > 1 + \sqrt{k/m} + t\} \leq \exp(-mt^2/c_1)$ for some $c_1$, see [27], Theorem 5.39.

$\tau_k(1.42 - a_1, a_1)$ *in (22), are satisfied for $u_1 + u_2 + u_3 = 4u$ for some small fraction $u$, if $m$ is on the order of (19) with* const $= \max(4/(a_1 a_2)^2, 2c_1/(0.29 - a_1)^2)$, *and $c_1$ depends on the distribution of $A_{ij}$'s. Note* const $\geq 4$.

*In the noiseless limit where only the first 2 conditions are required, this improves to $u_1 + u_2 = 3u$.*

**Remark 2.** *We emphasize again that the rate (19) is measured w.r.t. to the three conditions in Theorem 3. The probability for which both noise conditions i) and ii) are satisfied, and for which condition (8) imposed on $\min_{i \in \mathcal{S}} |\alpha_i|$ is satisfied, require additional consideration. For the former the probability is at least $1 - n^{-1}(2\pi \log n)^{-\frac{1}{2}}$, see [6]. For the latter, it has to be derived based on signal statistics, e.g., for $|\alpha_i| \in \mathbb{R}_{[0,1]}$ then $(\min_{i \in \mathcal{S}} |\alpha_i|) > t$ is observed with probability $(1 - t)^k$ with $|\mathcal{S}| = k$.*

Note that the choice for $a_3$ in Theorem 3 implies $||(\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}})^\dagger \boldsymbol{\beta}||_\infty$ is roughly on the order $\sqrt{k}$. Indeed this is true since $\tau_k \geq 1$, and we note $\tau_k = (\delta_{\max}/\delta_{\min})^2 + 2k^{-\frac{1}{2}} + o(k^{-\frac{1}{2}})$, thus $\tau_k \approx (\delta_{\max}/\delta_{\min})^2$ for moderate $k$. Now LASSO recovery also depends on the probability that condition (8) holds. Our choice for $a_3$ causes the RHS of (8) to be roughly of the order $c_Z \sqrt{2k \log n}$. Compare this to [6] (see Theorem 1.3) where it was assumed that $\sigma_{\min}(\boldsymbol{A}_{\mathcal{S}}) \leq 1.5$, they only require $a_3 = 3$, i.e., a factor of $\sqrt{k}$ is lost without this assumption (which was previously argued to be fairly restrictive). To improve Proposition 3, one might additionally assume some specific distributions on $\boldsymbol{A}$. We leave further improvements to future work.

*Proof of Proposition 3:* For notational convenience, put $\boldsymbol{X} = (\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}})^\dagger$. Bound the probability

$$\Pr\left\{ ||\boldsymbol{X}\boldsymbol{\beta}||_\infty > a\sqrt{k} \right\} \leq \Pr\left\{ ||\boldsymbol{X}\boldsymbol{\beta}||_\infty > a\sqrt{k} \,\Big|\, \mathcal{E}(\delta_{\min}, \delta_{\max}) \right\}$$
$$+ \Pr\{\mathcal{E}_c(\delta_{\min}, \delta_{\max})\}. \qquad (23)$$

where we take $a$ to mean

$$a = \frac{|\tau_k - 1|}{\tau_k + 1} \cdot \frac{1 + k^{-\frac{1}{2}}}{\sigma_{\min}^2(\boldsymbol{A}_{\mathcal{S}})} \qquad (24)$$

for $\tau_k$ chosen as in (22). We claim that every entry $(\boldsymbol{X}\boldsymbol{\beta})_i$ of $\boldsymbol{X}\boldsymbol{\beta}$ is upper bounded by $a\sqrt{k}$, for $a$ as in (24). Then by definition of $\mathcal{E}(\delta_{\min}, \delta_{\max})$, the first term in (23) equals 0 and we would have proven the bound (21).

Let $\boldsymbol{C}$ denote a $k \times 2$ matrix. The first column $\boldsymbol{C}$ is be a normalized version of $\boldsymbol{\beta}$, more specifically it equals $k^{-\frac{1}{2}}\boldsymbol{\beta}_i$. The second column equals the canonical basis vector $\mathbf{c}_i$, where $\mathbf{c}_i$ is a 0-1 vector whereby $(\mathbf{c}_i)_j = 1$ if and only if $j = i$. Consider the $2 \times 2$ matrix $\boldsymbol{X}'$ that satisfies $\boldsymbol{X}' = \boldsymbol{C}^T \boldsymbol{X} \boldsymbol{C}$. This matrix $\boldsymbol{X}'$ is symmetric (from symmetry of $\boldsymbol{X}$) and $k^{-\frac{1}{2}}(\boldsymbol{X}\boldsymbol{\beta})_i = X'_{1,2} = X'_{2,1}$ (from our construction of $\boldsymbol{C}$). That is the entry $X'_{1,2}$ (and $X'_{2,1}$) of $\boldsymbol{X}'$, correspond to the (scaled) quantity $k^{-\frac{1}{2}}(\boldsymbol{X}\boldsymbol{\beta})_i$ that we want to bound.

Condition on the event $\mathcal{E}_c(\delta_{\min}, \delta_{\max})$, then $\boldsymbol{A}_{\mathcal{S}}$ has rank $k$ and therefore $\boldsymbol{X} = (\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}})^\dagger = (\boldsymbol{A}_{\mathcal{S}}^T \boldsymbol{A}_{\mathcal{S}})^{-1}$. Let $\det(\cdot)$ and $\mathrm{Tr}(\cdot)$ denote determinant and trace. As in [38] equation (11),

we have

$$1 - \frac{X'_{1,2}X'_{1,2}}{X'_{1,1}X'_{2,2}} = \frac{4\det(\boldsymbol{X}')}{(\mathrm{Tr}(\boldsymbol{X}'))^2 - (X'_{1,1} - X'_{2,2})^2}$$

$$\geq \frac{4\varsigma_{\max}(\boldsymbol{X}') \cdot \varsigma_{\min}(\boldsymbol{X}')}{(\mathrm{Tr}(\boldsymbol{X}'))^2} = \frac{4t}{(1+t)^2} \quad (25)$$

where $t = \varsigma_{\max}(\boldsymbol{X}')/\varsigma_{\min}(\boldsymbol{X}')$ and $\varsigma_{\max}$ and $\varsigma_{\min}$ respectively denote the maximum and minimum eigenvalues. Now $t = \varsigma_{\max}(\boldsymbol{X}')/\varsigma_{\min}(\boldsymbol{X}') \geq 1$. If $t = 1$ then $4t/(1+t)^2 = 1$, and for $t \geq 1$ the function $4t/(1+t)^2$ decreases monotonically. We claim that $\tau_k$ in (22) upper bounds $\varsigma_{\max}(\boldsymbol{X}')/\varsigma_{\min}(\boldsymbol{X}')$, and (25) then allows us to produce the following upper bound

$$|X'_{1,2}| \leq \sqrt{X'_{1,1}X'_{2,2} \cdot \left(1 - \frac{4\tau_k}{(1+\tau_k)^2}\right)}$$

$$= \sqrt{X'_{1,1}X'_{2,2}} \cdot \frac{|\tau_k - 1|}{1+\tau_k}. \quad (26)$$

Bound $(X'_{1,1}X'_{2,2})^{\frac{1}{2}}$ by the maximum eigenvalue $\varsigma_{\max}(\boldsymbol{X}')$ of $\boldsymbol{X}'$. Then, further bound $\varsigma_{\max}(\boldsymbol{X}')$ by $(1 + k^{-\frac{1}{2}})/\sigma_{\min}^2(\boldsymbol{A}_\mathcal{S})$, which gives the form (24). This bound is argued as follows. For $k \geq 2$, we have the columns in $\mathbf{C}$ to be *linearly independent*. Since $\boldsymbol{X}' = \mathbf{C}^T\boldsymbol{X}\mathbf{C}$ and $\boldsymbol{X}$ is *positive definite*, it is then clear that $\varsigma_{\max}(\boldsymbol{X}') \leq \varsigma_{\max}(\mathbf{C}^T\mathbf{C}) \cdot \varsigma_{\max}(\boldsymbol{X})$. Now $\mathbf{C}^T\mathbf{C}$ is a $2 \times 2$ matrix with diagonal elements 1, and off-diagonal elements $\pm 1/\sqrt{k}$. Hence $\varsigma_{\max}(\mathbf{C}^T\mathbf{C}) = 1 + k^{-\frac{1}{2}}$. Also $\varsigma_{\max}(\boldsymbol{X}) \leq 1/\sigma_{\min}^2(\boldsymbol{A}_\mathcal{S})$, and the bound follows.

To finish, we show the claim $\tau_k \geq \varsigma_{\max}(\boldsymbol{X}')/\varsigma_{\min}(\boldsymbol{X}')$. By similar arguments as above, it follows that

$$\frac{\varsigma_{\max}(\boldsymbol{X}')}{\varsigma_{\min}(\boldsymbol{X}')} \leq \frac{\varsigma_{\max}(\mathbf{C}^T\mathbf{C})}{\varsigma_{\min}(\mathbf{C}^T\mathbf{C})} \cdot \frac{\varsigma_{\max}(\boldsymbol{X})}{\varsigma_{\min}(\boldsymbol{X})} = \frac{1 + k^{-\frac{1}{2}}}{1 - k^{-\frac{1}{2}}} \cdot \frac{\sigma_{\max}^2(\boldsymbol{A}_\mathcal{S})}{\sigma_{\min}^2(\boldsymbol{A}_\mathcal{S})} \leq \tau_k$$

since $\varsigma_{\min}(\boldsymbol{X}') \geq \varsigma_{\min}(\mathbf{C}^T\mathbf{C}) \cdot \varsigma_{\min}(\boldsymbol{X})$, and $\varsigma_{\min}(\boldsymbol{X}') = 1 - k^{-\frac{1}{2}}$, and $\boldsymbol{X} = (\boldsymbol{A}_\mathcal{S}^T\boldsymbol{A}_\mathcal{S})^{-1}$. We are done. ∎

## V. Conclusion

We take a first look at U-statistical theory for predicting the "average-case" behavior of salient CS matrix parameters. Leveraging on the generality of this theory, we consider two different recovery algorithms i) $\ell_1$-minimization and ii) LASSO. The developed analysis is observed to have good potential for predicting CS recovery, and compares well (empirically) with Donoho & Tanner [23] recent "average-case" analysis for system sizes found in implementations. Measurement rates that incorporate fractional $u$ failure rates, are derived to be on the order of $k[\log((n-k)/u) + \sqrt{2(k/n)\log(n/k)}]$, similar to the known optimal $k\log(n/k)$ rate. Empirical observations suggest possible improvement for const (as opposed to typical "worst-case" analyses whereby implicit constants are known to be inherently large).

There are multiple directions for future work. Firstly while restrictive maximum eigenvalue assumptions are avoided (as StRIP-recovery does not require them), the applied techniques could be fine-tuned. It is desirable to overcome the $\sqrt{k}$ losses observed here for noisy conditions. Secondly, it is interesting to further leverage the general U-statistical techniques to other different recovery algorithms, to try and obtain their good "average-case" analyses. Finally, one might consider

similar U-statistical "average-case" analyses for the case where the sampling matrix columns are dependent, which requires appropriate extensions of Theorem 1.

## Appendix

### A. Proof of Theorem 1

For notational simplicity we shall henceforth drop explicit dependence on $a$ from all three quantities $U_n(a), p(a)$ and $g(\mathbf{A}, a)$ in this appendix subsection. While $U_n$ is made explicit in Definition 1 as a statistic corresponding to the realization $\boldsymbol{\Phi} = \boldsymbol{A}$, this proof considers $U_n$ consisting of random terms $g(\boldsymbol{A}_\mathcal{S})$ for purposes of making probabilistic estimates. Theorem 1 is really a *law of large numbers* result. However even when the columns $\boldsymbol{A}_i$ are assumed to be IID, the terms $g(\boldsymbol{A}_\mathcal{S})$ in $U_n$ depend on each other. As such, the usual techniques for IID sequences do not apply. Aside from large deviation results such as Thm. 1, there exist *strong law* results, see [39]. The following proof is obtained by combining ideas taken from [33] and [34]. We use the following new notation just in this subsection of the appendix. Partition the index set $\{1, 2, \cdots, n\}$ into $\omega_n = \lfloor n/k \rfloor$ subsets denoted $\mathcal{S}_i$ each of size $k$, and a single subset $\mathcal{R}$ of size at most $k$. More specifically, let $\mathcal{S}_i = \{(i-1) \cdot k + 1, (i-1) \cdot k + 2, \cdots, i \cdot k\}$ and let $\mathcal{R} = \{\lfloor n/k \rfloor \cdot k + 1, \lfloor n/k \rfloor \cdot k + 2, \cdots, n\}$. Let $\pi$ denote a *permutation* (bijective) mapping $\{1, 2, \cdots, n\} \rightarrow \{1, 2, \cdots, n\}$. The notation $\pi(\mathcal{S})$ denotes the set of all *images* of each element in $\mathcal{S}$, under the mapping $\pi$. Following Section 5c in [33] we express the U-statistic $U_n$ of $\boldsymbol{A}$ in the form

$$U_n = \frac{1}{n!}\sum_\pi \left(\frac{1}{\omega_n}\sum_{i=1}^{\omega_n} g(\boldsymbol{A}_{\pi(\mathcal{S}_i)})\right), \quad (27)$$

the first summation taken over all $n!$ possible permutations $\pi$ of $\{1, 2, \cdots, n\}$. To verify, observe that any subset $\mathcal{S}$ is counted exactly $\omega_n \cdot k!(n-k)!$ times in the RHS of (27).

Recall $p = \mathbb{E}g(\boldsymbol{A}_\mathcal{S}) = \mathbb{E}U_n$. From the theorem statement let the term $\epsilon_n^2$ equal $cp(1-p) \cdot \omega_n^{-1}\log\omega_n$ where $c > 2$. We show that the probabilities $\Pr\{|U_n - p| > \epsilon_n\}$ for each $n$ are small. For brevity, we shall only explicitly treat the upper tail probability $\Pr\{U_n - p > \epsilon_n\}$, where standard modifications of the below arguments will address the lower tail probability $\Pr\{-U_n + p > \epsilon_n\}$ (see comment in p. 1, [33]). Using the expression (27) for $U_n$, write the probability $\Pr\{U_n - p > \epsilon_n\}$ for any $h > 0$ as

$$\Pr\{U_n - p > \epsilon_n\} \leq \mathbb{E}\exp(h(U_n - p + \epsilon_n))$$

$$= \mathbb{E}\exp\left(\frac{1}{n!}\left(\sum_\pi h(S_\pi - p + \epsilon_n)\right)\right),$$

where here $S_\pi$ is a RV that equals the inner summation in (27), i.e. $S_\pi = \frac{1}{\omega_n}\sum_{i=1}^{\omega_n} g(\boldsymbol{A}_{\pi(\mathcal{S}_i)})$. Using convexity of the function $\exp(\cdot)$ we express

$$\Pr\{U_n - p > \epsilon_n\} \leq \frac{1}{n!}\sum_\pi \mathbb{E}\exp(h(S_\pi - p + \epsilon_n)).$$

Now observe that the RV $S_\pi$ is an average of $\omega_n$ IID terms $g(\boldsymbol{A}_{\pi(\mathcal{S}_i)})$. This is due to the assumption that the columns $\boldsymbol{A}_i$ of $\boldsymbol{A}$ are IID, and also due to the fact that the sets $\pi(\mathcal{S}_i)$ are disjoint (recall sets $\mathcal{S}_i$ are disjoint). Hence for any permutation $\pi$, by this independence we have $\mathbb{E}\exp(hS_\pi) = (\mathbb{E}\exp(h'\cdot g(\boldsymbol{A}_{\pi(\mathcal{S}_1)})))^{\omega_n}$, where the normalization $h' = h/\omega_n$ bears no consequence. The RV $g(\boldsymbol{A}_{\pi(\mathcal{S}_1)})$ is bounded, i.e. $0 \le g(\boldsymbol{A}_{\pi(\mathcal{S}_i)}) \le 1$, and its expectation $\mathbb{E}g(\boldsymbol{A}_{\pi(\mathcal{S}_1)})$ equals $p$. By convexity of $\exp(\cdot)$ again and for all $h > 0$, the inequality $e^{h\alpha} \le e^h\alpha + 1 - \alpha$ holds for all $0 \le \alpha \le 1$. Therefore putting $\alpha = g(\boldsymbol{A}_{\pi(\mathcal{S}_1)})$ we get the inequality $\exp(h\cdot g(\boldsymbol{A}_{\pi(\mathcal{S}_1)})) \le 1+(e^h-1)\cdot g(\boldsymbol{A}_{\pi(\mathcal{S}_1)})$. By the irrelevance of $\pi$ in previous arguments, by putting $\mathbb{E}g(\boldsymbol{A}_{\pi(\mathcal{S}_1)}) = p$

$$\Pr\{U_n - p > \epsilon_n\} \le e^{-h(\epsilon_n+p)}\left(1-p+pe^h\right)^{\omega_n}.$$

We optimize the bound by putting $pe^h = (1-p)(p+\epsilon_n)/(1-p-\epsilon_n)$, see (4.7) in [33], to get

$$\Pr\{U_n - p > \epsilon_n\}$$
$$\le \left((1+\epsilon_n p^{-1})^{p+\epsilon_n}(1-\epsilon_n(1-p)^{-1})^{1-p-\epsilon_n}\right)^{-\omega_n}. \quad (28)$$

Following (2.20) in [34] we use the relation $\log(1+\alpha) = \alpha - \frac{1}{2}\alpha^2 + o(\alpha^2)$ as $\alpha \to 0$, to express the logarithmic exponent on the RHS of (28) as

$$\frac{-\omega_n\epsilon_n^2 \cdot (1+o(1))}{2p(1-p)}.$$

Therefore by the form $\epsilon_n^2 = cp(1-p) \cdot \omega_n^{-1}\log\omega_n$ where $c > 2$, for sufficiently large $n$ we have

$$\Pr\{U_n - p > \epsilon_n\} \le \omega_n^{-c/2} < \omega_n^{-1}$$

which in turn implies $\sum_{n=k}^\infty \Pr\{U_n - p > \epsilon_n\} < \infty$. Repeating similar arguments for the lower tail probability $\Pr\{-U_n+p > \epsilon_n\}$, we eventually prove $\sum_{n=k}^\infty \Pr\{|U_n-p| > \epsilon_n\} < \infty$ which implies the claim.

## REFERENCES

[1] E. Candes and T. Tao, "Decoding by linear programming," *IEEE Trans. on Inform. Theory*, vol. 51, pp. 4203–4215, Dec. 2005.

[2] R. Gribonval and M. Nielsen, "Sparse representations in unions of bases," *IEEE Trans. on Inform. Theory*, vol. 49, no. 12, pp. 3320–3325, Dec. 2003.

[3] D. L. Donoho and M. Elad, "Optimally sparse representation in general (nonorthogonal) dictionaries via $\ell_1$ minimization," *Proc. of the Nat. Acad. of Sci. (PNAS)*, vol. 100, no. 5, pp. 2197–2202, Mar. 2003.

[4] J. A. Tropp, "On the conditioning of random subdictionaries," *Applied and Computational Harmonic Analysis*, vol. 25, no. 1, pp. 1–24, Jul. 2008.

[5] J. J. Fuchs, "On sparse representations in arbitrary redundant bases," *IEEE Trans. on Inform. Theory*, vol. 50, no. 6, pp. 1341–1344, Jun. 2004.

[6] E. Candès and Y. Plan, "Near-ideal model selection by ell_1 minimization," *The Annals of Statistics*, vol. 37, no. 5A, pp. 2145–2177, 2009.

[7] A. Mazumdar and A. Barg, "Sparse recovery properties of statistical RIP matrices," in *IEEE 49th Annual Conf.. on Communication, Control, and Computing (Allerton)*, Sep. 2011, pp. 9–12.

[8] Y. Zhang, "Theory of compressive sensing via ell_1-minimization: a non-rip analysis and extensions," *Techical Report, Rice University*, 2008.

[9] A. Cohen, W. Dahmen, and R. DeVore, "Compressed sensing and best k-term approximation," *Journal of the American Mathematical Society*, vol. 22, no. 1, pp. 211–231, Jul. 2008.

[10] A. D'Aspremont and L. El Ghaoui, "Testing the nullspace property using semidefinite programming," *Mathematical Programming*, vol. 127, no. 1, pp. 123–144, Mar. 2011.

[11] D. L. Donoho, "Compressed sensing," *IEEE Trans. on Inform. Theory*, vol. 52, no. 4, pp. 1289 – 1306, Apr. 2006.

[12] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, "A simple proof of the restricted isometry property for random matrices," *Constructive Approximation*, vol. 28, Dec. 2008.

[13] J. D. Blanchard, C. Cartis, and J. Tanner, "Compressed sensing: How sharp is the RIP?" *SIAM Review*, vol. 53, pp. 105–525, Feb. 2011.

[14] R. Calderbank, S. Howard, and S. Jafarpour, "Construction of a large class of deterministic sensing matrices that satisfy a statistical isometry property," *IEEE Journal of Sel. Topics in Signal Proc.*, vol. 4, no. 2, pp. 358–374, Apr. 2010.

[15] L. Gan, C. Long, T. T. Do, and T. D. Tran. (2009) Analysis of the statistical restricted isometry property for deterministic sensing matrices using Stien's method. [Online]. Available: http://www.dsp.rice.edu/files/cs/Gan_StatRIP.pdf

[16] M. Sartipi and R. Fletcher, "Energy-Efficient Data Acquisition in Wireless Sensor Networks Using Compressed Sensing," *2011 Data Compression Conference*, pp. 223–232, Mar. 2011.

[17] F. Chen, A. P. Chandrakasan, and V. Stojanovic, "A signal-agnostic compressed sensing acquisition system for wireless and implantable systems," in *Proceedings of the IEEE Custom Integrated Circuits Conference (CICC '10)*, San Jose, CA, Sep. 2010, pp. 1–4.

[18] M. Mishali and Y. C. Eldar, "Expected RIP: Conditioning of the modulated wideband converter," in *Proceedings of the IEEE Information Theory Workshop (ITW '09)*, Taormina, Sicily, Oct. 2009, pp. 343–347.

[19] K. Kanoun, H. Mamaghanian, N. Khaled, and D. Atienza, "A real-time compressed sensing-based personal EKG monitoring system," in *Proceedings of the IEEE/ACM Design, Automation and Test in Europe Conference (DATE '11)*, Grenoble, France, Mar. 2011, pp. 1–6.

[20] C. Dossal, G. Peyré, and J. Fadili, "A numerical exploration of compressed sampling recovery," *Linear Algebra and its Applications*, vol. 432, no. 7, pp. 1663–1679, Mar. 2010.

[21] B. Bah and J. Tanner, "Improved bounds on restricted isometry constants for gaussian matrices," *SIAM J. Matrix Analysis*, vol. 31, no. 5, pp. 2882–2898, 2010.

[22] D. L. Donoho and J. Tanner, "Neighborliness of randomly-projected simplices in high dimensions," *Proc. of the Nat. Acad. of Sci. (PNAS)*, vol. 102, pp. 9452–9457, Jul. 2005.

[23] ——, "Exponential bounds implying construction of compressed sensing matrices, error-correcting codes, and neighborly polytopes by random sampling," *IEEE Trans. on Inform. Theory*, vol. 56, no. 4, pp. 2002–2016, Apr. 2010.

[24] R. Gribnoval, B. Mailhe, H. Rauhut, K. . Schnass, and P. Vandergheynst, "Average case analysis of multichannel thresholding," in *IEEE. International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, vol. 2. IEEE, Mar. 2007, pp. II–853–II–856.

[25] Y. C. Eldar and H. Rauhut, "Average Case Analysis of Multichannel Sparse Recovery Using Convex Relaxation," *IEEE Transactions on Information Theory*, vol. 56, no. 1, pp. 1–15, Jan. 2009.

[26] M. Golbabaee and P. Vandergheynst, "Average case analysis of sparse recovery with thresholding: New bounds based on average dictionary coherence," in *IEEE. International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Mar. 2008, pp. 3877–3880.

[27] R. Vershynin, "Introduction to the non-asymptotic analysis of random matrices," in *Compressed Sensing, Theory and Applications*, Y. Eldar and G. Kutyniok, Eds. Cambridge University Press, 2012, ch. 5, pp. 210–268.

[28] F. Lim and V. Stojanovic, "On U-statistics and compressed sensing II: Non-asymptotic worst-case analysis," *submitted*, 2012.

[29] E. Candes, "The restricted isometry property and its implications for compressed sensing," *Compte Rendus de l'Academie des Sciences, Paris, Serie I*, vol. 336, pp. 589–592, May 2008.

[30] J. A. Tropp, "Recovery of short, complex linear combinations via $\ell_1$ minimization," *IEEE Trans. on Inform. Theory*, vol. 51, no. 4, pp. 1568–1571, Apr. 2004.

[31] J. J. Fuchs, "Recovery of exact sparse representations in the presence of bounded noise," *IEEE Transactions on Information Theory*, vol. 51, no. 10, pp. 3601–3608, Oct. 2005.

[32] A. Edelman, "Eigenvalues and condition numbers of random matrices," *SIAM J. Matrix Analysis*, vol. 9, no. 4, pp. 543–560, 1988.

[33] W. Hoeffding, "Probability inequalities for sums of bounded random variables," *J. Am. Stat. Assoc.*, vol. 58, no. 301, pp. 13–30, Mar. 1963.

[34] P. K. Sen, "Asymptotic normality of sample quantiles for $m$-dependent processes," *Ann. Math. Statist.*, vol. 39, no. 5, pp. 1724–1730, 1968.

[35] P. Koev and A. Edelman, "The efficient evaluation of the hypergeometric function of a matrix argument," *Mathematics of Computation*, vol. 75, no. 254, pp. 833–846, 2006.

[36] A. Edelman and N. R. Rao, "Random matrix theory," *Acta Numerica*, vol. 14, pp. 233–297, May 2005.

[37] M. Rudelson and R. Vershynin, "Non-asymptotic theory of random matrices : extreme singular values," in *Proceedings of the International Congress of Mathematicians*. New Delhi: Hindustan Book Agency, 2010, pp. 1576–1602.

[38] A. Householder, "The Kantorovich and some related inequalities," *SIAM Review*, vol. 7, no. 4, pp. 463–473, 1965.

[39] R. H. Berk, "Limiting behavior of posterior distribution when the model is incorrect," *Annals of Math. Stat.*, vol. 37, pp. 51–58, 1966.

*A. Proofs of StRIP-type recovery guarantees appearing in Subsection II-B*

In this part of the appendix we provide the proofs for the two StRIP-type recovery guarantees discussed in this paper. The following are proofs for Theorems B and C.

*Proof of Theorem B, c.f., Lemma 3, [7]:* Define $\boldsymbol{\epsilon} \in \mathbb{R}^n$ as $\boldsymbol{\epsilon} = \boldsymbol{\alpha}^* - \boldsymbol{\alpha}$, *i.e.*, $\boldsymbol{\epsilon}$ is the recovery error vector. The proof technique closely follows that of Theorem 1.2, *c.f.*, [29]. Since $\text{sgn}(\boldsymbol{\alpha}_{\mathcal{S}}) = \boldsymbol{\beta}$, we have the inequality

$$||(\boldsymbol{\alpha} + \boldsymbol{\epsilon})_{\mathcal{S}}||_1 \geq ||\boldsymbol{\alpha}_{\mathcal{S}}||_1 + \boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}}. \tag{29}$$

Since $\boldsymbol{\alpha}^*$ solves (1), hence $||\boldsymbol{\alpha}^*||_1 \leq ||\boldsymbol{\alpha}||_1$. Putting $\boldsymbol{\alpha}^* = \boldsymbol{\alpha} + \boldsymbol{\epsilon}$, we have

$$
\begin{aligned}
||\boldsymbol{\alpha}||_1 \geq ||\boldsymbol{\alpha} + \boldsymbol{\epsilon}||_1 = & \ ||(\boldsymbol{\alpha} + \boldsymbol{\epsilon})_{\mathcal{S}}||_1 + ||(\boldsymbol{\alpha} + \boldsymbol{\epsilon})_{\mathcal{S}_c}||_1 \\
\geq & \ ||\boldsymbol{\alpha}_{\mathcal{S}}||_1 + \boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}} + ||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1 - ||\boldsymbol{\alpha}_{\mathcal{S}_c}||_1,
\end{aligned} \tag{30}
$$

where the last step follows the inequality (29), and the triangular inequality. Re-arranging and putting $||\boldsymbol{\alpha}_{\mathcal{S}_c}||_1 = ||\boldsymbol{\alpha}||_1 - ||\boldsymbol{\alpha}_{\mathcal{S}}||_1$ we get

$$||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1 \leq -\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}} + 2||\boldsymbol{\alpha}_{\mathcal{S}_c}||_1. \tag{31}$$

We next bound the term $-\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}}$ with $|\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}}|$, and for now assume that the following claim holds

$$|\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}}| \leq ||\boldsymbol{\beta}^T \boldsymbol{\Phi}_{\mathcal{S}}^{\dagger} \boldsymbol{\Phi}_{\mathcal{S}_c}||_\infty \cdot ||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1. \tag{32}$$

We then proceed to show the bound on $||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1$ (or $||\boldsymbol{\alpha}_{\mathcal{S}}^* - \boldsymbol{\alpha}_{\mathcal{S}}||_1$) to complete the first part of the proof. Using the *small projections* condition, bound $||\boldsymbol{\beta}^T \boldsymbol{\Phi}_{\mathcal{S}}^{\dagger} \boldsymbol{\Phi}_{\mathcal{S}_c}||_\infty \leq a_2$ using some $a_2 < 1$. This gives a upper bound of $a_2 \cdot ||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1$ on $|\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}}|$ in (32). Finally use this in (31) get $||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1 \leq a_2||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1 + 2||\boldsymbol{\alpha}_{\mathcal{S}_c}||_1$, or equivalently $||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1 \leq 2/(1 - a_2) \cdot ||\boldsymbol{\alpha}_{\mathcal{S}_c}||_1$. To show the claim (32), note that $\boldsymbol{\epsilon}$ is in the null-space of $\boldsymbol{\Phi}$, i.e. $\boldsymbol{\Phi}\boldsymbol{\epsilon} = \mathbf{0}$, or equivalently, $\boldsymbol{\Phi}_{\mathcal{S}}\boldsymbol{\epsilon}_{\mathcal{S}} = -\boldsymbol{\Phi}_{\mathcal{S}_c}\boldsymbol{\epsilon}_{\mathcal{S}_c}$. Let $\mathbf{I}$ denote the size-$k$ identity matrix. By the *invertability* condition, the pseudoinverse $\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}$ satisfies $\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}} = \mathbf{I}$. Hence

$$\boldsymbol{\epsilon}_{\mathcal{S}} = -\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}\boldsymbol{\epsilon}_{\mathcal{S}_c}, \tag{33}$$

and take the vector inner product with $\boldsymbol{\beta}$ on both sides to obtain $\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}} = -\boldsymbol{\beta}^T \boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}\boldsymbol{\epsilon}_{\mathcal{S}_c}$. Finally (32) holds by taking absolute value of $\boldsymbol{\beta}^T \boldsymbol{\epsilon}_{\mathcal{S}}$, and writing $|\boldsymbol{\beta}^T \boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}\boldsymbol{\epsilon}_{\mathcal{S}_c}| \leq ||\boldsymbol{\beta}^T \boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}||_\infty \cdot ||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1$.

To second part is to elucidate the bound on $||\boldsymbol{\epsilon}_{\mathcal{S}}||_1$ (or $||\boldsymbol{\alpha}_{\mathcal{S}_c}^* - \boldsymbol{\alpha}_{\mathcal{S}_c}||_1$). Starting from the previous relationship (33) we have $||\boldsymbol{\epsilon}_{\mathcal{S}}||_1 = ||\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1 \leq ||\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}||_\infty \cdot ||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1$. The result then follows by using the *worst-case projections* condition to bound $||\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}\boldsymbol{\Phi}_{\mathcal{S}_c}||_\infty$ by some positive $a_3$, and also bounding $||\boldsymbol{\epsilon}_{\mathcal{S}_c}||_1$ using the bound obtained in the first part of this proof. ∎

For the next two proofs we use the following notation. Let $\mathbf{I}$ denote the identity matrix, and let $\mathbf{P}$ denote a projection matrix onto the column subspace of $\boldsymbol{\Phi}_{\mathcal{S}}$, i.e., $\mathbf{P} = \boldsymbol{\Phi}_{\mathcal{S}}\boldsymbol{\Phi}_{\mathcal{S}}^{\dagger}$. We first address the proof of Proposition 4.

**Proposition 4** (*c.f.*, [6]). *Let $\mathbf{Z}$ be a random noise vector, whose components are IID zero mean Gaussian with variance $c_Z^2$. Assume that the matrix $\boldsymbol{\Phi}$ satisfies $||\boldsymbol{\phi}_i||_2 = 1$*

*for all columns $\boldsymbol{\phi}_i$. Then the realization $\boldsymbol{Z} = \mathbf{z}$ satisfies conditions i) and ii) in Theorem C with probability at least $1 - n^{-1}(2\pi \log n)^{-\frac{1}{2}}$.*

*Proof of Proposition 4, c.f., [6]:* The result will follow by showing i) holds with probability $k \cdot n^{-2}(2\pi \log n)^{-\frac{1}{2}}$, and by showing ii) holds with probability $(n - k) \cdot n^{-2}(2\pi \log n)^{-\frac{1}{2}}$.

For i), first assume each component of $\boldsymbol{Z}$ has variance 1. Let $\mathbf{c}_i$ denote the $i$-th row of $(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}$, thus we have $||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}} \boldsymbol{Z}||_\infty = \max_i |\mathbf{c}_i^T \boldsymbol{Z}|$. Since $\boldsymbol{Z}$ is Gaussian, thus

$$\Pr\{||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}} \boldsymbol{Z}||_\infty > z\} \le k \cdot \Pr\{|\tilde{Z}| > z\}, \quad (34)$$

where $\tilde{Z}$ is a Gaussian RV with standard deviation at least the $\ell_2$-norm of any row $\mathbf{c}_i$. It remains to then upper bound $||\mathbf{c}_i||_2$ for all $i$, which follows as $||\mathbf{c}_i||_2 \le ||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}||_2$. The spectral norm $||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}||_2$ is at most the reciprocal of the smallest non-zero singular value of $\boldsymbol{\Phi}_{\mathcal{S}}$, and by the *invertability* condition for some positive $a_1$, we have $||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}||_2 \le a_1^{-1}$. Then we let $\tilde{Z}$ in (34) have standard deviation $a_1^{-1}$. Equivalently,

$$\Pr\{||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{Z}||_\infty > z\} \le k \cdot \Pr\{|Z| > a_1 \cdot z\}$$
$$\le 2k \cdot f_Z(a_1 z)/(a_1 z) \quad (35)$$

where $Z$ is a standard normal RV with density function $f_Z(z)$. Generalizing to the case where each component of $\boldsymbol{Z}$ has variance $c_Z$, the upper bound becomes $2k \cdot f_Z((a_1 z)/c_Z)/((a_1 z)/c_Z)$. Put $z = (c_Z\sqrt{2 \log n})/a_1$ to get the claimed probabilistic upper estimate $k \cdot n^{-2}(2\pi \log n)^{-\frac{1}{2}}$.

For ii) we proceed similarly. Observe that for any $i \notin \mathcal{S}$, we have $||\boldsymbol{\phi}_i^T (\mathbf{I} - \mathbf{P})||_2 \le ||\boldsymbol{\phi}_i||_2 = 1$. Then put $z = c_Z 2\sqrt{\log n}$ in case ii) to get the claimed probabilistic upper estimate $(n - k) \cdot n^{-2}(2\pi \log n)^{-\frac{1}{2}}$. ∎

*Proof of Theorem 1.3, c.f., [6]:* We shall show that any signal $\boldsymbol{\alpha}$ with sign $\boldsymbol{\beta}$ and support $\mathcal{S}$, assuming $(\boldsymbol{\beta}, \mathcal{S})$ satisfy all three *invertability*, *small projections*, and *invertability projections* conditions together with (7) and (8), will have both sign and support successfully recovered.

The proof follows by constructing a vector $\boldsymbol{\alpha}'$ from $\boldsymbol{\alpha}$ as follows. Let $\boldsymbol{\epsilon}$ denote the error $\boldsymbol{\epsilon} = \boldsymbol{\alpha}' - \boldsymbol{\alpha}$, and $\boldsymbol{\alpha}'$ is defined by letting $\boldsymbol{\epsilon}$ satisfy

$$\begin{aligned} \boldsymbol{\epsilon}_{\mathcal{S}} &= (\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1}(\boldsymbol{\Phi}_{\mathcal{S}}^T \mathbf{z} - 2c_Z \theta_n \boldsymbol{\beta}), \\ \boldsymbol{\epsilon}_{\mathcal{S}_c} &= \mathbf{0}. \end{aligned} \quad (36)$$

Let us first claim that if (8) holds, then the support of $\boldsymbol{\alpha}'$ equals that of $\boldsymbol{\alpha}$. If this is true, then standard subgradient arguments, see [6], [31], will lead us to conclude that $\boldsymbol{\alpha}'$ must be the unique Lasso (6) solution (i.e., $\boldsymbol{\alpha}' = \boldsymbol{\alpha}^*$) if i) it satisfies

$$\begin{aligned} \boldsymbol{\phi}_i^T(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}') &= 2c_Z \theta_n \cdot \mathrm{sgn}(\alpha_i'), \quad \text{if } i \in \mathcal{S}, \\ |\boldsymbol{\phi}_i^T(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}')| &< 2c_Z \theta_n, \quad \text{if } i \notin \mathcal{S}, \end{aligned} \quad (37)$$

and ii) the submatrix $\boldsymbol{\Phi}_{\mathcal{S}}$ has full column rank. The condition ii) follows from the *invertability* condition, and the latter half of the proof will verify i). Let us first verify the previous claim that both $\boldsymbol{\alpha}'$ and $\boldsymbol{\alpha}$ have exact same supports. In fact, we go further to verify that $\boldsymbol{\alpha}'$ and $\boldsymbol{\alpha}$ also have the same signs. First

check

$$\begin{aligned} ||\boldsymbol{\epsilon}_{\mathcal{S}}||_\infty &\le ||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}^T \mathbf{z}||_\infty + 2\theta_n c_Z \cdot ||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\beta}||_\infty \\ &\le a_1^{-1} c_Z \cdot \sqrt{2 \log n} + 2a_3 c_Z \cdot \theta_n, \quad (38) \end{aligned}$$

where the final inequality follows from noise condition i) from Proposition 4, and the *invertability projections* condition which provides the bound $||(\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\beta}||_\infty \le a_3$ for some positive $a_3$. By assumption (8) and comparing with the above upper estimate for $||\boldsymbol{\epsilon}_{\mathcal{S}}||_\infty$, our claim must hold.

Next we go on to verify $\boldsymbol{\alpha}'$ satisfies (37). We have

$$\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}' = \mathbf{z} - \boldsymbol{\Phi}\boldsymbol{\epsilon} = \mathbf{z} - (\boldsymbol{\Phi}_{\mathcal{S}}^\dagger)^T (\boldsymbol{\Phi}_{\mathcal{S}}^T \mathbf{z} - 2c_Z \theta_n \cdot \boldsymbol{\beta}) \quad (39)$$

where the last equality follows by first writing $\boldsymbol{\Phi}\boldsymbol{\epsilon} = \boldsymbol{\Phi}\boldsymbol{\epsilon}_{\mathcal{S}}$, then substituting (36), and putting $\boldsymbol{\Phi}_{\mathcal{S}}^\dagger = (\boldsymbol{\Phi}_{\mathcal{S}}^T \boldsymbol{\Phi}_{\mathcal{S}})^{-1} \boldsymbol{\Phi}_{\mathcal{S}}^T$. Now because $\boldsymbol{\Phi}_{\mathcal{S}}^\dagger$ is a right inverse of $\boldsymbol{\Phi}_{\mathcal{S}}^T$, by left multiplying the above expression by $\boldsymbol{\Phi}_{\mathcal{S}}^T$ we conclude

$$\boldsymbol{\Phi}_{\mathcal{S}}^T(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}') = 2c_Z \theta_n \cdot \boldsymbol{\beta},$$

which is equivalent to the first set of equations of (36) as we verified before that $\boldsymbol{\beta} = \mathrm{sgn}(\boldsymbol{\alpha}_{\mathcal{S}}')$. For the second set of equations, observe from (39) that

$$\begin{aligned} (\mathbf{I} - \mathbf{P})(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}^*) &= (\mathbf{I} - \mathbf{P})\mathbf{z}, \\ \mathbf{P}(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}^*) &= 2c_Z \theta_n \cdot (\boldsymbol{\Phi}_{\mathcal{S}}^\dagger)^T \boldsymbol{\beta}, \end{aligned}$$

where the first equality follows because $(\mathbf{I} - \mathbf{P})(\boldsymbol{\Phi}_{\mathcal{S}}^\dagger)^T = \mathbf{0}$, and the second equality follows because $\mathbf{P}(\boldsymbol{\Phi}_{\mathcal{S}}^\dagger)^T \boldsymbol{\Phi}_{\mathcal{S}}^T = \mathbf{P}\mathbf{P}^T = \mathbf{P}^2 = \mathbf{P}$. Using the above two identities, we estimate

$$\begin{aligned} &||\boldsymbol{\Phi}_{\mathcal{S}_c}^T(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}^*)||_\infty \\ &\le ||\boldsymbol{\Phi}_{\mathcal{S}_c}^T(\mathbf{I} - \mathbf{P})(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}')||_\infty + ||\boldsymbol{\Phi}_{\mathcal{S}_c}^T \mathbf{P}(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}')||_\infty \\ &= ||\boldsymbol{\Phi}_{\mathcal{S}_c}^T(\mathbf{I} - \mathbf{P})\mathbf{z}||_\infty + 2c_Z \theta_n \cdot ||\boldsymbol{\Phi}_{\mathcal{S}_c}^T(\boldsymbol{\Phi}_{\mathcal{S}}^\dagger)^T \boldsymbol{\beta}||_\infty \\ &\le \frac{c_Z\sqrt{2}\theta_n}{1 + a} + 2c_Z a_2 \cdot \theta_n, \quad (40) \end{aligned}$$

where the upper estimate $(c_Z\sqrt{2}\theta_n)/(1 + a) = c_Z 2\sqrt{\log n}$ follows from noise condition ii) stated in Proposition 4, and $||\boldsymbol{\Phi}_{\mathcal{S}_c}^T(\boldsymbol{\Phi}_{\mathcal{S}}^\dagger)^T \boldsymbol{\beta}||_\infty \le a_2$ follows from the *small projections* property. Finally from assuming (7) we have $\sqrt{2}(1 + a)^{-1} + 2a_2 < 2$, and applying to the last member of (40) proves $||\boldsymbol{\Phi}_{\mathcal{S}_c}^T(\tilde{\mathbf{b}} - \boldsymbol{\Phi}\boldsymbol{\alpha}')||_\infty < 2c_Z \theta_n$, which verifies $\boldsymbol{\alpha}'$ satisfies the second set of equations of (36). Thus we verified $\boldsymbol{\alpha}' = \boldsymbol{\alpha}^*$ which is what we need to complete the proof. ∎

### B. Derivation of standard bounds

In the Gaussian case note $\mathbb{E}X_i^2 = 1$ and $\mathbb{E}X_i = 0$. Then $\sum_{i=1}^m c_i X_i$ is also Gaussian with variance $||\mathbf{c}||_2^2$. Hence by Markov's inequality we have the (single-sided) inequality $\Pr\{\sum_{i=1}^m c_i X_i > t\} \le \exp(-ht + h^2/||\mathbf{c}||_2^2)$ for any $h > 0$. The claim for the Gaussian case will follow by setting $h = t \cdot ||\mathbf{c}||_2^2/2$, and noting that for the other side $\Pr\{-(\sum_{i=1}^m c_i X_i) > t\} = \Pr\{\sum_{i=1}^m c_i X_i > t\}$. For the bounded case, note $|X_i| \le 1$ and $\mathbb{E}X_i = 0$, and the claim follows from Hoeffding's (2.6) in [33].

## C. Derivation of measurement rates

For the *small projections* condition, start from $p(a_2)$ being bounded by the RHS of (14) where $a = a_2$. As before bound $\Pr\{\sigma_{\min}(\boldsymbol{A}_{\mathcal{S}}) \leq a_1\} \leq \exp(-m \cdot (0.29 - a_1)^2/c_1)$, where we had set $\sqrt{\delta} = a_1$. From the identity $\sqrt{\alpha_1} \leq \sqrt{\alpha_2} + \sqrt{\alpha_3}$ for positive quantities $\alpha_i$, it follows from Theorem 1 and (18) that we will have $u_2 \leq (n - k) \cdot U_n(a_3) \leq 2u$, if we enforce

$$\frac{1}{2}\left[\log 2 + \log(n - k) - \frac{m(a_1 a_2)^2}{2k}\right] + t \leq \log u,$$

$$\frac{1}{2}\left[\log(n - k) - \frac{m(0.29 - a_1)^2}{c_1}\right] + t \leq \log u,$$

where $t = \sqrt{2(k/n)\log(n/k)}$. Ignoring the $\log 2$ term, and using $\sqrt{n - k} \leq n - k$, it follows that (19) enforces the two above conditions.

Similarly for the *invertability* condition, to have $u_1 = U_n(a_1) \leq u$ it follows from Theorem 1 and (18) that we need to enforce to second condition above.

For the *worst-case projections* condition, to have $u_3 \leq (n - k) \cdot U_n(a_3) \leq 2u$ we need to enforce

$$\frac{1}{2}\left[(k + 1) \cdot \log 2 + \log(n - k) - \frac{m(a_1 a_3)^2}{2k}\right] + t \leq \log u,$$

$$\frac{1}{2}\left[k\log 2 + \log(n - k) - \frac{m(0.29 - a_1)^2}{c_1}\right] + t \leq \log u.$$

Taking

$$k\log\left(\frac{n - k}{u}\right) \geq (k + 1) \cdot \log 2 + \log\left(\frac{n - k}{u}\right),$$

justifiable for $(n - k)/u$ suitably larger than 2, the rate (19) generously suffices to ensure these 2 conditions.

## D. More on noisy LASSO performance

The aim here is to provide more empirical evidence to support observations made in Figure 6 for more block lengths. Here Figure D.1 shows LASSO performance now for a wider range of $n$. We only consider $m = 150$, and show various recovery failure rates displayed via contoured lines, for various sparsities $k$ and block lengths $n$. Figures D.1(a) and (b) are companion to Figures 6(a) and (b), in that they respectively correspond to cases where the non-zero signal magnitudes $|\alpha_i|$ equal 1 (and $a = 0$), and in $\mathbb{R}_{[0,1]}$ (and $a = 1$). That is, for $n = 1000$, and $k = 4$ and $c_Z = 1 \times 10^{-4}$, we see the recovery failure is approximately $1 \times 10^{-3}$ in both Figure D.1(a) and Figure 6(a).

As mentioned in Subsection IV-C we observe good empirical match when adjusting the term $t = (a_1^{-1} + 2a_3(1 + a)) \cdot c_Z\sqrt{2\log n}$ (on the RHS of (8)) with $a_1 = 0.29$ and $a_3 = 1$. Figure D.1 provides further support. In (a) we show the values of the term $t$ for values $n = 300$ and $n = 3000$. Recall in this case when $t > 1$ condition (8) (and thus recovery) fails. Observe when $c_Z = 5 \times 10^{-2}$ the values of $t$ are very close to 1, and for $c_Z = 1 \times 10^{-1}$ they exceed 1. This matches with our observation in Figure 6(a) that $c_Z = 5 \times 10^{-2}$ is the critical point, beyond which for large $c_Z$ recovery fails catastrophically.

In (b) and (c) we look at the other case where $|\alpha_i| \in \mathbb{R}_{[0,1]}$. Here (c) plots the probability $1 - (1 - t)^k$ that (8) fails. Again the contoured lines delineate a particular fixed value of $1 - (1 - t)^k$ for various $k, n$ values, whereby we set $t = 7.4 \cdot c_Z\sqrt{2\log n}$ (recall we used $a = 1$ here). We observe how closely (c) tracks the noise floor regions in (b) (indicated by shading). More specifically note $t$ really depends on $n$, and the larger the probabilities $1 - (1 - t)^k$ get for various $k, n$ in Figure D.1(c), this probability overwhelms the LASSO recovery rates in Figure D.1(b). This matches with our previous observations in Figure 6(b).

Fig. D.1. Empirical LASSO performance shown for $m = 150$ for range of $k, n$ values. In $(a)$ the non-zero signal magnitudes $|\alpha_i|$ equal 1, and in $(b)$ they are in $\mathbb{R}_{[0,1]}$. In $(c)$ we plot a curve (expression) $1 - (1 - t)^k$ for $t = (3.4 + 2(1 + a)) \cdot c_Z \sqrt{2 \log n}$.