

SHARPNESS AND CONDITIONING OF NONSMOOTH CONVEX FORMULATIONS IN STATISTICAL SIGNAL RECOVERY*

LIJUN DING[†] AND ALEX L. WANG[‡]

Abstract. We study a sample complexity vs. conditioning tradeoff in modern signal recovery problems (including sparse recovery, low-rank matrix sensing, covariance estimation, and abstract phase retrieval), where convex optimization problems are built from sampled observations. We begin by introducing a set of condition numbers related to sharpness in the ℓ_1 or Schatten-1 norm of nonsmooth formulations for these problems. Then, we show that these condition numbers become dimension *independent* constants in each of the example signal recovery problems once the sample size exceeds some constant multiple of the recovery threshold. Structurally, this result ensures that the inaccuracy in the recovered signal due to both observation noise and optimization error is controlled. Algorithmically, such a result ensures that a new restarted mirror descent method achieves nearly-dimension-independent linear convergence to the signal in terms of iterations. This new first-order method is general and applies to any sharp convex function in an ℓ_p or Schatten- p norm for $p \in [1, 2]$.

Key words. Sharp formulation, signal recovery, restarted mirror descent, sparse recovery, matrix sensing, phase retrieval

MSC codes. 90C30, 90C06, 90C22

1. Introduction. This paper studies a *sample complexity vs. convex optimization conditioning* tradeoff in modern signal recovery problems such as sparse recovery [15], low-rank matrix sensing [9, 44], phase retrieval [14], and covariance estimation [20]. To illustrate the setup, the challenges, and our results concretely, we will only consider the abstract phase retrieval problem (henceforth phase retrieval) in the introduction. The general setup and results follow a similar pattern and are deferred to later sections.

Phase retrieval asks us to recover an unknown rank-one positive semidefinite (PSD) matrix $X^\natural \in \mathbb{S}^n$ given a vector of m measurements $b = \mathcal{A}(X^\natural)$ where the linear map $\mathcal{A} : \mathbb{S}^n \rightarrow \mathbb{R}^m$ is defined by $\mathcal{A}(X)_i = g_i^\top X g_i$ for i.i.d. $g_i \sim N(0, I_n/m)$. Here and throughout, \mathbb{S}^n and \mathbb{S}_+^n denote the vector space of symmetric matrices and the cone of PSD matrices respectively. Concretely, we are asked to

recover $X^\natural \in \mathbb{S}^n$ from (b, \mathcal{A}) where $b = \mathcal{A}(X^\natural) \in \mathbb{R}^m$.

In [10], Candès and Li established an *exact recovery* result: If $m = \Theta(n)$, then, with high probability (w.h.p.)¹, X^\natural is the unique optimizer for the convex optimization problem

$$(1.1) \quad \min_{X \in \mathbb{S}^n} \left\{ \text{tr}(X) : \begin{array}{l} \mathcal{A}(X) = b \\ X \in \mathbb{S}_+^n \end{array} \right\}.$$

Here, the choice of the objective function $\text{tr}(X)$ and the constraint $X \in \mathbb{S}_+^n$ reflect prior knowledge that the true signal is low rank and PSD respectively.

*Submitted to the editors January 25, 2026.

Funding: This work was completed in part while A. L. Wang was supported by the Dutch Scientific Council (NWO) grant OCENW.GROOT.2019.015 (OPTIMAL). L. Ding was supported by an NSF TRIPODS grant to the Institute for Foundations of Data Science (NSF DMS-2023239).

[†]Department of Mathematics, University of California San Diego, La Jolla, CA (l2ding@ucsd.edu, <https://www.lijunding.net/>).

[‡]Daniels School of Business, Purdue University, West Lafayette, IN (wang5984@purdue.edu, <https://web.ics.purdue.edu/~wang5984/>)

¹An event happens w.h.p. if the probability is larger than $1 - c_1 \exp(-c_2 m)$ for some numerical constants c_1, c_2 .

37 *Issues and challenges.* Unfortunately, convexity and uniqueness of solution are
 38 not the end of the story as (1.1) may not be well-conditioned. In such a situation,
 39 the recovery process may suffer from:

- 40 • *measurement error:* if the observation b is corrupted by noise, the optimal
 41 solution of (1.1) may deviate greatly from X^\natural or may even cease to exist;
- 42 • *optimization error:* the problem (1.1) may admit near-optimal and near-
 43 feasible solutions in terms of function value and constraint violation that are
 44 far from the true signal X^\natural ;
- 45 • *slow algorithm convergence:* the iteration complexity of existing first-order
 46 algorithms for solving (1.1) may depend polynomially on the dimension due
 47 to polynomially poor conditioning of (1.1).² Such dependence weakens the
 48 applicability of existing first-order methods on large-scale instances of (1.1).

49 *Our contribution.* In this paper, we study a notion of conditioning related to
 50 *sharpness* for (1.1) based on an *unconstrained nonsmooth reformulation* of (1.1). This
 51 notion of conditioning is measured in the Schatten-1 norm and allows us to quanti-
 52 tatively control each of the above issues. To measure the conditioning of (1.1), we
 53 introduce the following penalized nonsmooth version of (1.1):

$$54 \quad (1.2) \quad \min_{X \in \mathbb{S}^n} F_{r,\ell}(X) := \text{tr}(X) + r \|\mathcal{A}(X) - b\|_1 + \ell \text{dist}_1(X, \mathbb{S}_+^n).$$

55 Here, the distance $\text{dist}_1(X, \mathbb{S}_+^n) := \inf_{Y \in \mathbb{S}_+^n} \|Y - X\|_1$ is measured in the Schatten-1
 56 norm, i.e., the sum of the singular values, and $r, \ell \geq 0$ are penalty parameters on the
 57 violations of $\mathcal{A}(X) = b$ and $X \in \mathbb{S}_+^n$ respectively. To clean up the exposition, we will
 58 understand “ ℓ_p norm” of a matrix to mean its Schatten- p norm and “ ℓ_p norm” of a
 59 vector to mean its usual ℓ_p norm. We will also abuse terminology and refer to the rank
 60 of a matrix or the size of the support of a vector as its “ ℓ_0 norm”. We emphasize that
 61 we penalize the violation of $\mathcal{A}(X) = b$ by $\|\mathcal{A}(X) - b\|_1$ and not the more common
 62 least-squares error $\|\mathcal{A}(X) - b\|_2^2$ and we use the ℓ_1 norm in the distance function.
 63 The term $\|\mathcal{A}(X) - b\|_1$ is sometimes referred to as the least absolute deviation error
 64 term [43] and is used as a robust version of the least-squares error. Although this
 65 choice leads to a nonsmooth problem, it also opens the possibility for the function
 66 (1.2) to be μ -sharp around X^\natural . Specifically, we will show that $F_{r,\ell}$ is μ -sharp around
 67 X^\natural w.r.t. the ℓ_1 norm, that is, for all $X \in \mathbb{S}^n$

$$68 \quad (1.3) \quad \underbrace{\text{tr}(X) + r \|\mathcal{A}(X) - b\|_1 + \ell \text{dist}_1(X, \mathbb{S}_+^n)}_{F_{r,\ell}(X)} - \underbrace{\text{tr}(X^\natural)}_{F_{r,\ell}(X^\natural)} \geq \mu \|X - X^\natural\|_1.$$

69 When such a bound holds, we will think of the parameters (μ, r, ℓ) as partially de-
 70 scribing the conditioning of (1.1). Finally, in order to make our notion of conditioning
 71 “aware” of scaling, we additionally track the Lipschitz constant, L , of (1.2) with re-
 72 spect to the ℓ_1 norm:

$$73 \quad (1.4) \quad |F_{r,\ell}(X) - F_{r,\ell}(Y)| \leq L \|X - Y\|_1, \quad \text{for all } X, Y \in \mathbb{S}^n.$$

74 We use the set of numbers (μ, r, ℓ, L) to measure the conditioning of (1.1). The
 75 sharpness parameter μ in some sense will control absolute notions of error, while the

²These polynomial factors may be even worse when measuring error in terms of distance to the true solution as opposed to function value.

condition number $\kappa := L/\mu$ will control relative notions of error and convergence rates of first-order methods.³

We now describe our main results in the setting of phase retrieval:

- *Conditioning*: It is well-known that restricted isometry property (RIP) constants (see Definition 3.2) can be used to prove exact recovery [15, 44]. Our first contribution is a novel connection between RIP constants and sharpness in the sense of (1.3). By adapting known proofs of exact recovery based on restricted isometry, we can describe the parameters (μ, r, ℓ, L) in terms of RIP constants. In particular, we show that there exist $\mu, r, \ell, L = \Theta(1)$ and $m = \Theta(n)$, so that w.h.p., the function $F_{r,\ell}$ is μ -sharp and L -Lipschitz. Almost identical results on dimension-independent conditioning for the other signal recovery problems are also presented in Theorem 4.1.
- *Optimization error and error bound*: In practice, we are unable to solve (1.1) exactly and must resort to numerical optimization. Using μ -sharpness, i.e., (1.3), we have that any ϵ -suboptimal solution⁴ X to (1.1) or (1.2) satisfies $\|X - X^\natural\|_1 \leq \epsilon/\mu$. Using the estimate $\mu = \Theta(1)$, we are able to control how optimization error affects the distance to solution, $\|X - X^\natural\|_1$, a form of recovery error. Thus, μ -sharpness can be viewed as a particular form of an error bound, see (1.5) for a general definition. Compared to approaches [7, 22, 24, 51] based on Hoffman’s error bound [33] or other regularity conditions, e.g., strict complementarity [1], our approach based on the RIP yields a much clearer relationship between the parameters of the error bound, e.g., μ, r, ℓ , and the underlying dimension n . See the paragraph containing (1.5) in Section 1.1 for detailed discussions.
- *Measurement error and sensitivity of solutions*: Suppose that instead of observing $b = \mathcal{A}(X^\natural)$, we observe $\tilde{b} = \mathcal{A}(X^\natural) + \delta$ with a (possibly adversarial) noise term. It is known that (1.1) actually admits a unique *feasible* solution X^\natural w.h.p. [10], thus (1.1) with b replaced by \tilde{b} is likely to be infeasible. On the other hand, given any optimizer \tilde{X} of (1.2) with \tilde{b} in place of b , we may use sharpness in the noiseless case to bound $\|\tilde{X} - X^\natural\|_1 \leq O(\|\delta\|_1/\mu)$ in the noisy case (see Proposition 5.1).
In a similar vein, we may consider a setting where b is corrupted by a sparse vector δ . Sharpness in the noiseless case tells us that the optimizer of (1.2) is *unchanged* if $|\text{supp}(\delta)|/m = O(\mu)$, i.e., if up to an $O(\mu)$ fraction of the entries of b are corrupted (see Proposition 5.6 and Example 5.7 for a formal statement).⁵
Thus, the parameter μ allows us to control how measurement error affects the recovery process. Our estimate $\mu = \Theta(1)$ shows that the optimizer is insensitive to noise in b .
- *Algorithms*: Since (1.1) has well-conditioned sharpness in the ℓ_1 norm, its sharpness in the standard ℓ_2 norm is necessarily ill-conditioned (see Proposition 6.2). This ill-conditioning causes both the theoretical performance guar-

³One may replace the ℓ_1 norm with other norms in the right hand side of these definitions (1.3) and (1.4). The benefit of the ℓ_1 norm, as we will discuss, is that the conditioning in our applications can be shown to be *dimension-independent*.

⁴Given an optimization problem $\min_{x \in \mathcal{X}} f(x)$ with minimum f^* , a feasible region \mathcal{X} in some finite dimensional Euclidean space, and an objective function f , a vector x is an ϵ -suboptimal solution to $\min_{x \in \mathcal{X}} f(x)$ if $f(x) - f^* \leq \epsilon$ and $x \in \mathcal{X}$.

⁵Naturally, the allowed fraction of corrupted entries will satisfy $O(\mu) < 1/2$.

118 antees and the numeric performance of standard sharpness-aware first-order
 119 methods, which are based on the ℓ_2 norm, to deteriorate as the dimension
 120 increases (see Figure 1a). Thus, to solve large-scale instances of (1.1), we de-
 121 velop a new restarted mirror descent algorithm (RMD), which can be applied
 122 to *general* convex functions that are μ -sharp and L -Lipschitz in terms of an
 123 ℓ_p norm with $p \in [1, 2]$ (see Section 6). We may apply this algorithm to (1.2)
 124 with $p = 1$ to produce an ϵ -suboptimal solution in

$$125 \quad O\left(\kappa^2 \log(n) \log\left(\frac{1}{\epsilon}\right)\right)$$

126 iterations. RMD's *logarithmic* dependence on n contrasts with the *linear*
 127 dependence on n that an ℓ_2 -based algorithm would incur. Furthermore, each
 128 iteration of RMD can be implemented in time comparable to a single iteration
 129 of subgradient descent on (1.2).⁶ These convergence guarantees also hold
 130 after appropriate modifications in the presence of corruption or noise (see
 131 Propositions 5.1 and 5.6) with guarantees depending on κ .

132 Our results on optimization error, measurement error, and algorithms above de-
 133 pend solely on the conditions (1.3) and (1.4) and are stated in terms of paramters
 134 (μ, r, ℓ, L) . For a general problem, it may be that these conditions do not hold for
 135 *any* choice of these parameters. However, we will show that under standard statisti-
 136 cal assumptions, not only do (1.3) and (1.4) hold for some choice of (μ, r, ℓ, L) , but
 137 that these parameters are well-controlled. We emphasize again that the above results
 138 continue to hold similarly for other modern signal recovery problems.

139 **1.1. Related work.** To better position our work in the literature, in this section,
 140 we discuss related work on error bounds and conditioning, sharpness and conditioning
 141 in statistical recovery, and first-order methods for minimizing sharp functions.

142 *Establishing error bounds and conditioning.* The sharpness condition in (1.3) can
 143 be viewed as an error bound for (1.1). Error bounds and conditioning are central to
 144 optimization problems both structurally and algorithmically [1, 24, 34, 35, 38, 55].
 145 Consider the following general error bound for (1.1):

$$146 \quad (1.5) \quad \text{tr}(X) + r \| \mathcal{A}(X) - b \| + \ell \text{dist}(X, \mathbb{S}_+^n) - \text{tr}(X^\natural) \geq \mu \| X - X^\natural \|^s, \text{ for all } X \in \mathbb{S}^n,$$

147 where, in addition to the r and ℓ introduced in (1.3), we also have the power s , and
 148 general norms throughout. Most methods for establishing error bounds either require
 149 that (1.1) (or the more general problem (P) introduced in later sections) is linear or its
 150 subdifferential graph is linear [7, 24, 33], or that regularity conditions such as Slater's
 151 condition [5, 54], strict complementarity [22, 24, 51], or nondegeneracy conditions [37]
 152 hold. Additional work has studied error bounds in generic semialgebraic settings [23].
 153 Results of this type generally place an emphasis on establishing the power s (usually
 154 to 1 or 2) and provide only weak bounds on μ , r , and ℓ , that potentially depend on the
 155 dimension polynomially or exponentially. This is natural as the power s determines
 156 the linear or sublinear convergence rate [2, 34]. On the other hand, if μ , r , or ℓ depend
 157 polynomially (or even exponentially) on the dimension, then first-order algorithms
 158 may still be doomed to dimension-dependent convergence rates. This work connects
 159 the quantities μ , r , and ℓ to the quantitatively better-understood restricted isometry

⁶In the phase retrieval case, RMD requires performing two eigenvalue value decompositions per iteration, whereas subgradient descent requires performing one eigenvalue decomposition per iteration.

160 property (RIP). This connection gives precise estimates of the conditioning (μ, r, ℓ, L)
 161 in our settings and shows that these quantities are in fact dimension independent.

162 *Sharpness and well conditioning in statistical recovery.* A series of works [18, 19,
 163 26, 36] studied the local landscape of *nonconvex* nonsmooth formulations in statistical
 164 recovery near the signal, particularly in low-rank matrix recovery. They showed that
 165 near a factored form of the ground truth X^\natural , the natural loss function is sharp, and
 166 the local condition number (in terms of the Frobenius or Schatten-2 norm) does not
 167 depend on the dimension directly. Hence, they can apply off-the-shelf algorithms,
 168 such as subgradient and prox-linear methods [21, 25] to achieve quick convergence to
 169 a factored form of X^\natural . Note that results of this type are *local*, and the algorithms
 170 require specific model knowledge to construct an appropriate initializer to achieve
 171 provable guarantees. In contrast, the methods developed here are convex optimization
 172 methods so they are not as sensitive to poor initialization.

173 *First-order methods for minimizing sharp functions.* Early work in this area es-
 174 tablished linear convergence of the subgradient method (with appropriate step sizes)
 175 for μ -sharp L -Lipschitz functions in the ℓ_2 norm [28, 32, 42]. These methods and their
 176 proofs are adapted to the ℓ_2 norm and incur a *linear dependence on the dimension*
 177 when applied to μ -sharp L -Lipschitz functions in the ℓ_1 norm.

178 Similar guarantees can be derived as a consequence of restarting schemes. Perhaps
 179 the earliest work discussing restarting schemes is [39], where a restarting scheme
 180 was developed for convex Hölder-smooth minimization in ℓ_p spaces. More recent
 181 work [38, 47, 48, 53] has used restarting schemes to accelerate variants of (sub)gradient
 182 descent in the presence of different growth conditions, e.g., a local guarantee of the
 183 form⁷

$$184 \quad f(x) - \min_x f(x) \geq \mu \operatorname{dist}(x, \mathcal{X})^s$$

185 for all x close enough to the set of minimizers, \mathcal{X} of a function f . Here, $\operatorname{dist}(\cdot, \cdot)$ is
 186 typically measured in the ℓ_2 norm and $s \in [1, \infty)$ captures the Hölderian growth of
 187 f . Taking $s = 1$ recovers the usual definition of sharpness. In this setting, restarted
 188 subgradient descent achieves a linear convergence rate of $O(\kappa^2 \log(1/\epsilon))$ [53] where
 189 the condition number κ is measured in the ℓ_2 norm. Adaptive versions of these restart
 190 schemes have also been developed [45, 48] that do not require the sharpness parameter
 191 μ to be specified but incur an additional $\log(1/\epsilon)$ factor in the convergence rate.

192 We emphasize that much of the preceding literature on first-order methods for
 193 sharp convex functions focuses on the ℓ_2 setting. The only exception we know of is the
 194 work of [48], which develops a restarted mirror descent method for functions satisfying
 195 a property inspired by sharpness in the ℓ_2 setting. Roulet and d’Aspremont [48] show
 196 that their method can be used to achieve $O(\kappa^2 \log(n) \log(1/\epsilon))$ convergence on this
 197 particular function class. Unfortunately, in the ℓ_p setting with $p \in [1, 2)$, their proxy
 198 for sharpness is quite restrictive and does not allow solutions with sparseness or low-
 199 rank as is common in signal recovery. Our work in section 6 is closely related to [48]
 200 but uses a different specification of the restarting mechanism. This change allows us
 201 to use the natural definition of sharpness in an ℓ_p norm for $p \in [1, 2]$ (see Remark 6.7
 202 for a thorough comparison).

203 **1.2. Outline and notation.**

204 *Outline.* We begin, in Section 2, by setting up notation and defining sharpness
 205 in the ℓ_1 norm for an abstract signal recovery problem. We will quantify sharpness

⁷Other forms of the necessary “growth” are possible in different settings.

206 and conditioning by a set of parameters (μ, r, ℓ, L) . In Sections 3 and 4, we establish
 207 conditioning and sharpness under the Gaussian sensing models for the problems of
 208 sparse recovery, low-rank matrix and bilinear sensing, and covariance estimation in
 209 the absence of noise. Our proof relies on the well-known restricted isometry property
 210 (RIP). In section 5, we investigate how sharp instances of the abstract signal recovery
 211 problem behave in the presence of noise. In section 6, we describe and analyze an
 212 algorithm restarted mirror descent (RMD) that has linear and nearly dimension-
 213 independent convergence rates when applied to sharp functions in any ℓ_p or Schatten- p
 214 norms for $p \in [1, 2]$ (see Assumption 6.1). In conjunction with the results of Section 5,
 215 we may apply RMD to obtain linear nearly dimension-independent convergence rates
 216 even in the presence of noise for these statistical recovery problems.

217 *A note on ℓ_p and Schatten- p norms.* Let $p \in [1, \infty]$. We will overload notation
 218 and use $\|x\|_p$ for the standard ℓ_p norm when $x \in \mathbb{R}^n$ and the Schatten- p norm when
 219 $x \in \mathbb{R}^{n \times N}$. Recall that the nuclear norm is equivalent to the Schatten-1 norm,
 220 the Frobenius norm is equivalent to the Schatten-2 norm, and the spectral norm is
 221 equivalent to the Schatten- ∞ norm. To streamline the text, we will also take “ ℓ_p
 222 norm” of a matrix to mean its Schatten- p norm.

223 ℓ_0 “norm”. Finally, we abuse notation and write $\|x\|_0$ for the size of the support
 224 of a vector or the rank of a matrix. We caution that this is *not* a norm.

225 **2. Preliminaries.** In this section, we describe an abstract signal recovery prob-
 226 lem and the convex optimization approach. Then, we define our measures of sharpness
 227 and conditioning.

228 **2.1. Signal recovery via convex optimization.** Let V be the space of the
 229 signal, i.e., $V = \mathbb{R}^n$, $\mathbb{R}^{n \times N}$, or \mathbb{S}^n . We will equip V with the ℓ_1 norm. Again, “ ℓ_1 ”
 230 means the standard ℓ_1 norm if $V = \mathbb{R}^n$ and the Schatten-1 norm if $V = \mathbb{R}^{n \times N}$ or
 231 $V = \mathbb{S}^n$. Let \mathbb{R}^m be the space of the measurements equipped with the ℓ_p norm for
 232 some $p \in [1, 2]$ to be determined.

233 In the abstract signal recovery problem, we aim to recover an unknown element
 234 $x^\natural \in V$, called the signal, from a linear sensing operator, $\mathcal{A} : V \rightarrow \mathbb{R}^m$, and a vector
 235 of measurements, $b = \mathcal{A}(x^\natural) \in \mathbb{R}^m$. The convex optimization approach solves the
 236 following generalization of (1.1),

$$237 \quad (\text{P}) \quad \min_{x \in V} \left\{ f(x) : \begin{array}{l} \mathcal{A}(x) = b \\ x \in \mathcal{K} \end{array} \right\},$$

238 to recover x^\natural . Here, the function $f : V \rightarrow \mathbb{R}$ is convex and $\mathcal{K} \subseteq V$ is a closed convex
 239 cone (possibly all of V). They reflect our prior knowledge of x^\natural , e.g., if x^\natural is sparse,
 240 we set $f(x) = \|x\|_1$.

241 Next, we define the sharpness and conditioning of (P).

242 **2.2. Definition of sharpness and conditioning.** Given $x \in V$, and a closed
 243 nonempty set $\mathcal{X} \subseteq V$, let

$$244 \quad \text{dist}_1(x, \mathcal{X}) := \min_{\bar{x} \in \mathcal{X}} \|\bar{x} - x\|_1.$$

245 The following definition specializes the notion of sharpness for a convex function
 246 (also referred to as weak sharpness in the literature [30]) to the ℓ_1 setting.

247 **DEFINITION 2.1 (Sharpness).** *Let $\phi : V \rightarrow \mathbb{R}$ be a convex function, $\mathcal{X} \subseteq V$ a*
 248 *nonempty convex set, and let $\mu > 0$. We say that ϕ is μ -sharp around \mathcal{X} w.r.t. the ℓ_1*

249 norm if $\mathcal{X} = \arg \min \phi$ and

$$250 \quad \phi(\bar{x}) - \min_{x \in V} \phi(x) \geq \mu \operatorname{dist}_1(\bar{x}, \mathcal{X}), \quad \forall \bar{x} \in V.$$

251 We extend the notion of sharpness to a problem of the form (P) as follows.

252 **DEFINITION 2.2** (Sharpness of (P)). *Consider a problem of the form (P). Let*
 253 $x^\natural \in V$, $\mu > 0$, and $r, \ell \geq 0$. *We say that (P) is (μ, r, ℓ) sharp around x^\natural w.r.t. the*
 254 ℓ_1 *norm if x^\natural is feasible in (P) and*

$$255 \quad (2.1) \quad F_{r,\ell}(x) := f(x) + r \|\mathcal{A}(x) - b\|_p + \ell \operatorname{dist}_1(x, \mathcal{K})$$

256 *is μ -sharp around x^\natural . Equivalently, if x^\natural is feasible in (P) and*

$$257 \quad (2.2) \quad \|x - x^\natural\|_1 \leq \frac{1}{\mu} \left(f(x) - f(x^\natural) + r \|\mathcal{A}(x) - b\|_p + \ell \operatorname{dist}_1(x, \mathcal{K}) \right), \quad \forall x \in V.$$

258 *If $\mathcal{K} = V$, then ℓ is inconsequential, so we will simply say (P) is (μ, r) sharp around*
 259 x^\natural *w.r.t. the ℓ_1 norm.*

260 **Remark 2.3.** Suppose (P) is (μ, r, ℓ) -sharp around x^\natural w.r.t. the ℓ_1 norm for some
 261 $\mu > 0$, $r, \ell \geq 0$. Note that x^\natural is necessarily the unique optimal solution of (P).
 262 Moreover, we have an error bound: If we numerically solve (P) and produce $\tilde{x} \in V$,
 263 then we may bound the distance $\|\tilde{x} - x^\natural\|_1$ according to (2.2). This bound holds even
 264 if $f(\tilde{x}) < f(x^\natural)$, which could happen as \tilde{x} is not necessarily feasible for (P).

265 Recall the standard definition of the Lipschitz constant of a function specialized
 266 to the ℓ_1 norm.

267 **DEFINITION 2.4** (Lipschitz constant). *We say that a convex function $\phi : V \rightarrow \mathbb{R}$*
 268 *is L -Lipschitz w.r.t. the ℓ_1 norm if*

$$269 \quad |\phi(x) - \phi(y)| \leq L \|x - y\|_1, \quad \forall x, y \in V.$$

270 We are now ready to define the conditioning of (P). Note that the following
 271 definition depends implicitly on p , which we treat as fixed.

272 **DEFINITION 2.5** (Conditioning and condition number). *Consider a problem of the*
 273 *form (P) and suppose $x^\natural \in V$ is its unique optimizer. Suppose for some $r, \ell \geq 0$ that*
 274 $F_{r,\ell}$ *is μ -sharp around x^\natural and L -Lipschitz w.r.t. the ℓ_1 norm, with $\mu, L > 0$. The*
 275 *conditioning of (P) is measured by (μ, r, ℓ, L) and the condition number is $\kappa = \frac{L}{\mu}$.*

276 It is possible for a problem of the form (P) to not have well-defined conditioning
 277 according to this definition. For example, if x^\natural is *not* the unique minimizer of (P), then
 278 its conditioning is not well-defined. Nevertheless, we will show that under standard
 279 statistical assumptions, many interesting signal recovery problems have well-defined
 280 conditioning under this definition.

281 **3. RIP-based sharpness and conditioning in signal recovery.** This section
 282 shows that under bounds on restricted isometry (defined below), sharpness in the
 283 ℓ_1 norm holds for three archetypal modern signal recovery problems: sparse vector
 284 recovery, low rank matrix sensing, and covariance estimation. Phase retrieval can also
 285 be thought of as a special case of covariance estimation. The signal to be recovered,
 286 the optimization formulation in (P), and the choices of norms are the following.

287 **DEFINITION 3.1.** *We consider the following signal recovery problems that can be*
 288 *formulated as (P) and the corresponding choices for V , f , and \mathcal{K} .*

- 289 • *Sparse vector recovery*: recover a k -sparse vector $x^\natural \in V = \mathbb{R}^n$. In the
290 formulation (P), the objective is $f(x) = \|x\|_1$ and the set $\mathcal{K} = \mathbb{R}^n$.
- 291 • *Low rank matrix sensing*: recover a rank k matrix $X^\natural \in V = \mathbb{R}^{n \times N}$. In the
292 formulation (P), the objective is $f(X) = \|X\|_1$ and the set $\mathcal{K} = \mathbb{R}^n$.
- 293 • *Covariance estimation*: recover a rank k PSD matrix $X^\natural \in V = \mathbb{S}^n$. In the
294 formulation (P), the objective is the trace $f(X) = \text{tr}(X)$ and the set \mathcal{K} is the
295 set of PSD matrices \mathbb{S}_+^n .

296 The norm, $\|\cdot\|_p$ on the measurement space \mathbb{R}^m is, as yet, unspecified in any of
297 the three problems above. We have some freedom in the choice of p as long as crucial
298 bounds on the restricted isometry property (RIP) hold in the ℓ_p norm.

299 **DEFINITION 3.2** (Restricted isometry property). *Let k' be a positive integer and*
300 *$\mathcal{A} : V \rightarrow \mathbb{R}^m$ be a linear operator. We will let $\text{RIP}_{k'}^-(\mathcal{A})$ and $\text{RIP}_{k'}^+(\mathcal{A})$ denote any*
301 *valid uniform lower bound and upper bound on the quantity $\|\mathcal{A}(x)\|_p / \|x\|_2$ as $x \in V$*
302 *ranges over all elements in V with $\|x\|_0 \leq k'$.*

303 Our main result of this section is the following:

304 **PROPOSITION 3.3.** *Consider one of the three problems defined in Definition 3.1.*
305 *Let $c = 1$ for sparse vector recovery and $c = 2$ for the other two cases. Let $k' > 0$ and*
306 *$\epsilon > 0$ and suppose $\text{RIP}_{k'}^+(\mathcal{A}) \geq 1$ and*

$$307 \quad \sqrt{\frac{k'}{ck}} \left(\frac{\text{RIP}_{ck+k'}^-(\mathcal{A})}{\text{RIP}_{k'}^+(\mathcal{A})} \right) \geq 1 + \epsilon.$$

308 *Then, (P) is $\left(\frac{\epsilon}{2+\epsilon}, \sqrt{k'}, 2\right)$ sharp around x^\natural w.r.t. ℓ_1 . Furthermore, the Lipschitz*
309 *constant L w.r.t. ℓ_1 is no more than $3 + \sqrt{k'} \text{RIP}_1^+(\mathcal{A})$.*

310 To put this result in context, we will see in Section 4 that the premise of Propo-
311 sition 3.3 holds once the sample size m is a small multiple of the recovery threshold
312 w.h.p. with $\epsilon = 2$ and $k' = O(k)$, and $\text{RIP}_1^+(\mathcal{A}) = O(1)$. Hence, Problem (P) is
313 $\left(\frac{1}{2}, O(\sqrt{k}), 2\right)$ sharp around x^\natural with a Lipschitz constant no more than $O(\sqrt{k})$
314 w.h.p.

315 The rest of the section consists of proofs of Proposition 3.3 for the three different
316 settings. We start with the proof of sharpness of matrix sensing and discuss how
317 this proof can be modified for the other two settings. We then bound the Lipschitz
318 constant for all three settings.

319 **3.1. Proof of sharpness in matrix sensing.** This section proves Proposi-
320 tion 3.3 for the setting of low-rank matrix sensing and uses many elements of the
321 proofs for exact recovery [15, 44].

322 Let $\Delta \in \mathbb{R}^{n \times N}$ be arbitrary. Our goal is to show that

$$323 \quad (3.1) \quad \|X^\natural + \Delta\|_1 + \sqrt{k'} \|\mathcal{A}(\Delta)\|_p - \|X^\natural\|_1 \geq \frac{\epsilon}{2+\epsilon} \|\Delta\|_1.$$

324 *Change of basis.* Suppose $U_1 \in \mathbb{R}^{n \times n}$ and $U_2 \in \mathbb{R}^{N \times N}$ are orthonormal changes
325 of bases. Note that the inequality (3.1) is invariant upon replacing (X, Δ) with
326 $(U_1 X U_2^\top, U_1 \Delta U_2^\top)$ and replacing \mathcal{A} with $\tilde{\mathcal{A}} : \tilde{\Delta} \mapsto \mathcal{A}(U_1^\top \tilde{\Delta} U_2)$. That is:

$$327 \quad (3.1) \iff$$

$$328 \quad \left\| U_1 X^\natural U_2^\top + U_1 \Delta U_2^\top \right\|_1 + \sqrt{k'} \left\| \tilde{\mathcal{A}}(U_1 \Delta U_2^\top) \right\|_p - \left\| U_1 X^\natural U_2^\top \right\|_1 \geq \frac{\epsilon}{2+\epsilon} \left\| U_1 \Delta U_2^\top \right\|_1.$$

329 Furthermore, $\text{RIP}_{k'}^-(\mathcal{A}) = \text{RIP}_{k'}^-(\tilde{\mathcal{A}})$ and $\text{RIP}_{k'}^+(\mathcal{A}) = \text{RIP}_{k'}^+(\tilde{\mathcal{A}})$ for any k' . Thus,
 330 without loss of generality, we may assume that X^\natural and Δ have the following form:

$$331 \quad X^\natural = \begin{pmatrix} X_{1,1}^\natural & 0_{k \times (N-k)} \\ 0_{(n-k) \times k} & 0_{(n-k) \times (N-k)} \end{pmatrix} \quad \text{and} \quad \Delta = \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & \Delta_{2,2} \end{pmatrix}$$

332 where $X_{1,1}^\natural$ is a nonnegative $k \times k$ diagonal matrix and $\Delta_{2,2} \in \mathbb{R}^{(n-k) \times (N-k)}$ is diagonal
 333 with diagonal entries that are nonnegative and nonincreasing in magnitude.

334 *Partitioning of the difference Δ .* Let $\delta_{k^\perp} := \text{diag}(\Delta_{2,2}) \in \mathbb{R}^{\min(n,N)-k}$. The rest
 335 of the proof is based on the idea in [15, 44]. We decompose $\delta_{k^\perp} = \sigma_1 + \dots + \sigma_t$ where
 336 each $\sigma_i \in \mathbb{R}^{\min(n,N)-k}$. Specifically, let σ_1 extract the first k' coordinates of δ_{k^\perp} as its
 337 first k' coordinates and be zero for other coordinates. We then let each subsequent σ_i
 338 extract the next k' coordinates of δ_{k^\perp} as its $(i-1)k' + 1$ -th to ik' -th coordinates and
 339 be 0 for other coordinates. Finally, σ_t may have fewer than k' nonzero coordinates of
 340 δ_{k^\perp} and its first $(t-1)k'$ coordinates are all zeros. Let Σ_i denote the matrix of size
 341 $n \times N$ with $\text{Diag}(\sigma_i)$ in its bottom-right $(n-k) \times (N-k)$ block and zeros elsewhere.

342 *Lower bound on $\|\Delta_{2,2}\|_1$.* Since $\sigma_1, \dots, \sigma_t$ have disjoint supports, we can lower
 343 bound $\|\Delta_{2,2}\|_1 = \|\delta_{k^\perp}\|_1$ by

$$344 \quad \|\Delta_{2,2}\|_1 = \|\delta_{k^\perp}\|_1 \geq \sum_{i=1}^{t-1} \|\sigma_i\|_1 = k' \sum_{i=1}^{t-1} \frac{\|\sigma_i\|_1}{k'}.$$

345 Using δ_{k^\perp} is nonincreasing and each σ_i is k' sparse, we further have

$$346 \quad \|\Delta_{2,2}\|_1 \geq k' \sum_{i=2}^t \|\sigma_i\|_\infty \geq \sqrt{k'} \sum_{i=2}^t \|\sigma_i\|_2.$$

347 Next, we use the RIP to obtain

$$348 \quad \|\Delta_{2,2}\|_1 \geq \frac{\sqrt{k'}}{\text{RIP}_{k'}^+} \sum_{i=2}^t \|\mathcal{A}(\Sigma_i)\|_p$$

$$349 \quad \stackrel{(a)}{\geq} \frac{\sqrt{k'}}{\text{RIP}_{k'}^+} \left(\left\| \mathcal{A} \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,2} & 0 \end{pmatrix} + \mathcal{A}(\Sigma_1) \right\|_p - \|\mathcal{A}(\Delta)\|_p \right),$$

350 where (a) is due to triangle inequality. Using the RIP condition again and $\text{RIP}_{k'}^+ \geq 1$,
 351 we have

$$352 \quad \|\Delta_{2,2}\|_1 \geq \sqrt{k'} \frac{\text{RIP}_{2k+k'}^-}{\text{RIP}_{k'}^+} \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} + \Sigma_1 \right\|_2 - \sqrt{k'} \|\mathcal{A}(\Delta)\|_p$$

$$353 \quad \geq \sqrt{k'} \frac{\text{RIP}_{2k+k'}^-}{\text{RIP}_{k'}^+} \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_2 - \sqrt{k'} \|\mathcal{A}(\Delta)\|_p.$$

354 Lastly, since $\begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix}$ has rank no more than $2k$, we have

$$355 \quad \|\Delta_{2,2}\|_1 \geq \sqrt{\frac{k'}{2k}} \frac{\text{RIP}_{2k+k'}^-}{\text{RIP}_{k'}^+} \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_1 - \sqrt{k'} \|\mathcal{A}(\Delta)\|_p$$

$$356 \quad \geq (1 + \epsilon) \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_1 - \sqrt{k'} \|\mathcal{A}(\Delta)\|_p.$$

357 *Putting things together.* We are now ready to prove sharpness:

$$\begin{aligned}
358 & \|X^\natural + \Delta\|_1 + \sqrt{k'} \|\mathcal{A}(\Delta)\|_p - \|X^\natural\|_1 \\
359 & \geq \left\| \begin{pmatrix} X_{1,1}^\natural & 0 \\ 0 & \Delta_{2,2} \end{pmatrix} \right\|_1 - \|X^\natural\|_1 - \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_1 \\
360 & \quad + \sqrt{k'} \|\mathcal{A}(\Delta)\|_p \\
361 & \geq \|\delta_{k^\perp}\|_1 - \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_1 + \sqrt{k'} \|\mathcal{A}(\Delta)\|_p \quad (\text{disjoint supports}) \\
362 & = \left(\frac{2}{2+\epsilon} + \frac{\epsilon}{2+\epsilon} \right) \|\delta_{k^\perp}\|_1 - \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_1 \\
363 & \quad + \sqrt{k'} \|\mathcal{A}(\Delta)\|_p \\
364 & \geq \left(\frac{2}{2+\epsilon} (1+\epsilon) - 1 \right) \left\| \begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} \right\|_1 + \frac{\epsilon}{2+\epsilon} \|\delta_{k^\perp}\|_1 \quad (\text{bound on } \|\delta_{k^\perp}\|_1) \\
365 & \quad + \left(1 - \frac{2}{2+\epsilon} \right) \sqrt{k'} \|\mathcal{A}(\Delta)\|_p \\
366 & \geq \frac{\epsilon}{2+\epsilon} \|\Delta\|_1.
\end{aligned}$$

367 This proves the claim (3.1) as $\Delta \in \mathbb{R}^{n \times N}$ was arbitrary.

368 **3.2. Sharpness in sparse vector recovery.** One can follow the proof for ma-
369 trix sensing almost verbatim by treating the signal x^\natural and the sensing vectors in \mathcal{A} as
370 diagonal matrices. The change-of-bases matrices may be picked as any permutation
371 matrix $U_1 = U_2$. The change from $c = 2$ to $c = 1$ is justified by noting there is no
372 extra off-diagonal term for the matrices considered in the argument. Thus, the matrix
373 $\begin{pmatrix} \Delta_{1,1} & \Delta_{1,2} \\ \Delta_{2,1} & 0 \end{pmatrix} = \begin{pmatrix} \Delta_{1,1} & \\ & 0 \end{pmatrix}$ has rank bounded by k instead of $2k$.

374 **3.3. Sharpness in covariance estimation.** We may repeat the proof of Sub-
375 section 3.1 verbatim after replacing $V = \mathbb{R}^{n \times N}$ by $V = \mathbb{S}^n$. The change-of-bases
376 matrices may be picked as any orthonormal matrix $U_1 = U_2$. We deduce that, under
377 the assumptions of this proposition,

$$378 \quad \|X\|_1 + \sqrt{k'} \|\mathcal{A}(X) - b\|_p$$

379 is $\frac{\epsilon}{2+\epsilon}$ sharp around X^\natural . Next, note that $\text{tr}(X) + 2 \text{dist}_1(X, \mathbb{S}_+^n) \geq \|X\|_1$. We conclude
380 that

$$381 \quad \text{tr}(X) + 2 \text{dist}_1(X, \mathbb{S}_+^n) + \sqrt{k'} \|\mathcal{A}(X) - b\|_p$$

382 is $\frac{\epsilon}{2+\epsilon}$ sharp around X^\natural .

383 **3.4. Proof of Lipschitz continuity.** We need to show that

$$384 \quad (3.2) \quad f(x) + \sqrt{k'} \|\mathcal{A}(x) - b\|_p + 2 \text{dist}_1(x, K)$$

385 is $3 + \sqrt{k'} \text{RIP}_1^+(\mathcal{A})$ Lipschitz with respect to the ℓ_1 norm.

386 By the triangle inequality, the objective functions $f(x) = \|x\|_1$ and $f(X) = \text{tr}(X)$
387 are both 1-Lipschitz in terms of the ℓ_1 norm. In the covariance estimation setting,
388 the distance function $\text{dist}_1(X, \mathbb{S}_+^n)$ is also 1-Lipschitz in terms of the ℓ_1 norm by the
389 triangle inequality.

390 Next, we show the function $\|\mathcal{A}(x) - b\|_p$ is RIP_1^+ Lipschitz. Indeed, for any
 391 $x_1, x_2 \in V$, let $\Delta = x_1 - x_2 = \sum_{i=1}^n \Delta_i$, where $\|\Delta_i\|_0 \leq 1$. We have

$$\begin{aligned}
 & \left| \|\mathcal{A}(x_1) - b\|_p - \|\mathcal{A}(x_2) - b\|_p \right| \\
 & \leq \|\mathcal{A}(x_1 - x_2)\|_p = \left\| \mathcal{A} \left(\sum_{i=1}^n \Delta_i \right) \right\|_p \leq \sum_{i=1}^n \|\mathcal{A}(\Delta_i)\|_p \\
 & \stackrel{(a)}{\leq} \sum_{i=1}^n \text{RIP}_1^+ \|\Delta_i\|_2 \stackrel{(b)}{=} \sum_{i=1}^n \text{RIP}_1^+ \|\Delta_i\|_1 = \text{RIP}_1^+ \|\Delta\|_1.
 \end{aligned}$$

395 Here step (a) and (b) are due to the RIP and the fact that $\|\Delta_i\|_0 \leq 1$ for all i .

396 Since $f(x)$ and $\text{dist}_1(x, K)$ are 1 Lipschitz and $\|\mathcal{A}(x) - b\|_p$ is RIP_1^+ Lipschitz,
 397 our proof that the function in (3.2) is $(3 + \sqrt{k'} \text{RIP}_1^+(\mathcal{A}))$ -Lipschitz with respect to
 398 the ℓ_1 norm is complete.

399 **4. Conditioning, sensing models, and sample complexity.** In this section,
 400 we describe different sensing models of \mathcal{A} and show that once m exceeds certain
 401 thresholds, Problem (P) is well-conditioned in terms of ℓ_1 norm (for certain choices
 402 of the norm ℓ_p on the measurement space).

403 To prove this result, we first describe the precise sensing model in Subsection 4.1.
 404 Then we collect bounds on the RIP from the literature that ensure that the premise
 405 of Proposition 3.3 is satisfied. Due to an additional technicality in one of the sensing
 406 models (labeled ‘‘Covariance Estimation I’’ below), we give a separate proof of its
 407 conditioning in Subsection 4.2.

408 We will describe each sensing map \mathcal{A} by specifying its adjoint map \mathcal{A}^* . More
 409 precisely, we equip \mathbb{S}^n and $\mathbb{R}^{n \times N}$ with the trace inner product and equip \mathbb{R}^n and \mathbb{R}^m
 410 with the dot product. For each $i = 1, \dots, m$, let e_i denote the i -th standard basis of
 411 \mathbb{R}^m . We will specify the distribution of $\mathcal{A}^*(e_i) \in V$ and then recover $\mathcal{A} = (\mathcal{A}^*)^*$.

412 A summary of the sensing model, the thresholds, and the norms on \mathbb{R}^m is de-
 413 scribed in Table 1. Our main result of this section is the following theorem.

414 **THEOREM 4.1.** *Suppose the norm on \mathbb{R}^m and the sampling map \mathcal{A} are described*
 415 *according to one of the scenarios in Table 1 and the signal x^\natural satisfies $\|x\|_0 \leq k$.*
 416 *If $m \geq CT(n, N, k)$, where $T(n, N, k)$ is defined in Table 1 and C is a numerical*
 417 *constant, then there are numerical constants $c_1, c_2 > 0$ such that w.h.p. the optimiza-*
 418 *tion problem (P) is $(\frac{1}{2}, \sqrt{c_1 k}, 2)$ -sharp around x^\natural and $F_{\sqrt{c_1 k}, 2}$ has a Lipschitz constant*
 419 *bounded by $\sqrt{c_2 k}$. Consequently, (P) has a condition number κ bounded by $2\sqrt{c_2 k}$.*

420 **4.1. Sensing model, sample complexity, and proof via RIP.** In this sec-
 421 tion, we first describe the sensing models and the thresholds $T(n, N, k)$ on m so that
 422 bounds on RIP hold for \mathcal{A} when \mathbb{R}^m is equipped with either the ℓ_1 or ℓ_2 norms. Then,
 423 we prove Theorem 4.1 by verifying the premise of Proposition 3.3.

424 Recall that $\text{RIP}_{k'}^-(\mathcal{A})$ and $\text{RIP}_{k'}^+(\mathcal{A})$ are any uniform lower and upper bound on
 425 $\|\mathcal{A}x\|_p / \|x\|_2$ as x ranges over elements of V with support or rank bounded by k' .
 426 They will be set to numerical constants c_1, c_2 below and could differ for different
 427 sensing models. The norm of \mathbb{R}^m will be either $\ell_p = \ell_1$ or $\ell_p = \ell_2$.

428 *Sparse vector recovery.* For sparse vector recovery, the measurements are of the
 429 form

$$430 \quad \mathcal{A}^*(e_i) = a_i, \quad i = 1, \dots, m,$$

Task	$\mathcal{A}^*(e_i)$	$T(n, N, k)$	ℓ_p norm on \mathbb{R}^m
Sparse vector recovery	$a_i \in \mathbb{R}^n$	$k \log(n/k)$	ℓ_1 or ℓ_2
Low rank matrix sensing I	$A_i \in \mathbb{R}^{n \times N}$	$\max\{n, N\}k$	ℓ_1 or ℓ_2
Low rank matrix sensing II	$a_i b_i^\top \in \mathbb{R}^{n \times N}$	$\max\{n, N\}k$	ℓ_1
Covariance estimation I	$a_i a_i^\top \in \mathbb{S}^n$	nk	ℓ_1
Covariance estimation II	$a_i a_i^\top - b_i b_i^\top \in \mathbb{S}^n$	nk	ℓ_1

Table 1: Description of different statistical signal recovery tasks and the thresholds for well-conditioning. The entries of A_i , a_i , and b_i are i.i.d. Gaussian random variables with appropriate scaling (see Section 4.1). The conditioning of (P) is measured by $(\frac{1}{2}, \sqrt{c_1 k}, 2, \sqrt{c_2 k})$ where $c_1, c_2 > 0$ are numerical constants.

431 where $a_i \stackrel{\text{i.i.d.}}{\sim} N(0, I/m)$ when $\ell_p = \ell_2$ and $a_i \stackrel{\text{i.i.d.}}{\sim} N(0, I/m^2)$ when $\ell_p = \ell_1$. In
432 both settings, we may set w.h.p. $\text{RIP}_{k'}^+(\mathcal{A}) = c_2$ and $\text{RIP}_{k'}^-(\mathcal{A}) = c_1$ as long as
433 $m \gtrsim k' \log(\frac{n}{k'})$,⁸ where c_1 and c_2 are constants independent of k' satisfying $c_2/c_1 \leq 1.1$
434 and $c_2 \geq 1$. This fact is proved in the ℓ_2 setting following [13, 27] using results on
435 singular values of random matrices; a simple proof can be found in [4, Theorem 5.2].
436 This fact can be proved in the ℓ_1 setting using [41, Lemma 2.1] and [49, Lemma 4.4].
437 ⁹

438 *Matrix sensing I.* For this scenario of matrix sensing, the measurements are of
439 the form

$$440 \quad \mathcal{A}^*(e_i) = A_i, \quad i = 1, \dots, m,$$

441 where each measurement matrix $A_i \in \mathbb{R}^{n \times N}$ has Gaussian entries. Each entry of each
442 matrix is sampled i.i.d. according to $N(0, 1/m)$ when $\ell_p = \ell_2$ setting, and according
443 to $N(0, 1/m^2)$ when $\ell_p = \ell_1$. In both settings, we may set w.h.p. $\text{RIP}_{k'}^+(\mathcal{A}) = c_2$
444 and $\text{RIP}_{k'}^-(\mathcal{A}) = c_1$ as long as $m \gtrsim k' \max(n, N)$, where c_1 and c_2 are constants
445 independent of k' satisfying $c_2/c_1 \leq 1.1$ and $c_2 \geq 1$. This fact is proved in the ℓ_2
446 setting in [12, Theorem 2.3]. This fact is proved in the ℓ_1 setting in [36, Proposition
447 1].

448 *Matrix sensing II.* For this version of matrix sensing, which is more commonly
449 known as bilinear sensing, the measurements are of the form

$$450 \quad \mathcal{A}^*(e_i) = a_i b_i^\top, \quad i = 1, \dots, m,$$

451 and the measurement vectors $a_i \stackrel{\text{i.i.d.}}{\sim} N(0, I_n/m)$ and $b_i \stackrel{\text{i.i.d.}}{\sim} N(0, I_N/m)$. We equip
452 \mathbb{R}^m with the ℓ_1 norm. In this setting, we may set w.h.p. $\text{RIP}_{k'}^+(\mathcal{A}) = c_2$ and
453 $\text{RIP}_{k'}^-(\mathcal{A}) = c_1$ as long as $m \gtrsim k' \max(n, N)$, where c_1 and c_2 and constants inde-
454 pendent of k' satisfying $c_2/c_1 \leq 4$ and $c_2 \geq 1$. This fact is proved according to [9,
455 Theorem 2.2].

456 *Covariance estimation I.* For this scenario of covariance estimation, the measure-
457 ments are of the form

$$458 \quad \mathcal{A}^*(e_i) = a_i a_i^\top, \quad i = 1, \dots, m,$$

⁸We write $a \gtrsim b$ if $a \geq Cb$ for some numerical constant $C > 0$.

⁹Specifically, [41, Lemma 2.1] provides a RIP bound with a quantity called Gaussian width. One then use [49, Lemma 4.4] to estimate the width and obtain the concrete RIP above.

459 where each measurement vector $a_i \stackrel{\text{i.i.d.}}{\sim} N(0, I_n/m)$. We equip \mathbb{R}^m with the ℓ_1 norm.
 460 The proof of Theorem 4.1 for this setting differs from the proof of Theorem 4.1 for
 461 all other settings. This difference stems from the fact that $\langle \mathcal{A}^*(e_i), X^\natural \rangle$ does not
 462 have zero mean, thus biasing the output vector as discussed in [20, Section III.B].
 463 Attempting to follow the same proof strategy will need $\text{RIP}_{k'}^+(\mathcal{A})$ to be a numerical
 464 constant. However, the quantity $\text{RIP}_{k'}^+(\mathcal{A})$ scales as $\sqrt{k'}$. This prevents us from
 465 applying Proposition 3.3 directly in this setting.

466 We provide a separate proof for Theorem 4.1(Covariance estimation I) in Sub-
 467 section 4.2. The proof in this setting is completed by analyzing the conditioning of a
 468 related model of covariance estimation (Covariance estimation II).

469 *Covariance estimation II.* For this scenario of covariance estimation, the mea-
 470 surements are of the form

$$471 \quad \mathcal{A}^*(e_i) = a_i a_i^\top - b_i b_i^\top, \quad i = 1, \dots, m,$$

472 where each measurement vector pair $a_i, b_i \stackrel{\text{i.i.d.}}{\sim} N(0, I_n/(2m))$. We equip \mathbb{R}^m with the
 473 ℓ_1 norm. In this setting, we may set w.h.p. $\text{RIP}_{k'}^+(\mathcal{A}) = c_2$ and $\text{RIP}_{k'}^-(\mathcal{A}) = c_1$ as
 474 long as $m \gtrsim nk'$. Here, c_1 and c_2 are constants independent of k' satisfying $c_2/c_1 \leq 4$
 475 and $c_2 \geq 1$. To prove this fact, consider the operator $\bar{\mathcal{A}} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ such that
 476 $\bar{\mathcal{A}}^*(e_i) = (a_i - b_i)(a_i + b_i)^\top$. Note that $a_i - b_i$ and $a_i + b_i$ are independent of each other
 477 and $a_i - b_i, a_i + b_i \stackrel{\text{i.i.d.}}{\sim} N(0, I_n/m)$ for $i = 1, \dots, m$. Using [9, Theorem 2.2], we see that
 478 RIP holds for $\bar{\mathcal{A}}$ w.h.p. with $\text{RIP}_{k'}^+(\bar{\mathcal{A}}) = c_2$ and $\text{RIP}_{k'}^-(\bar{\mathcal{A}}) = c_1$ as long as $m \gtrsim nk'$.
 479 The proof is completed by noting that $[\bar{\mathcal{A}}(X)]_i = \langle (a_i - b_i)(a_i + b_i)^\top, X \rangle = [\mathcal{A}(X)]_i$
 480 for any $X \in \mathbb{S}^n$.

481 The Gaussian assumption on the sensing vectors or matrices are assumed for
 482 simplicity. One can relax this condition to sub-Gaussian distributions as done in [4,
 483 Theorem 5.2], [41, Lemma 2.1], and [18, Theorem 6.4].

484 *Proof of Theorem 4.1 for all settings except covariance estimation I.* For all set-
 485 tings except covariance estimation I, as long as $m \gtrsim T(n, N, k')$, we may set w.h.p.
 486 $\text{RIP}_{k'}^+(\mathcal{A}) = c_2$ and $\text{RIP}_{k'}^-(\mathcal{A}) = c_1$. By setting $k' = \sqrt{C_1 k}$ for a large numerical con-
 487 stant C_1 , we see the premise of Proposition 3.3 is satisfied and our proof is complete.

488 **4.2. Conditioning of covariance estimation I.** We provide a separate argu-
 489 ment for Theorem 4.1 in the setting of Covariance estimation I.

490 In this setting, we have that $a_i \sim N(0, I_n/m)$ are i.i.d. We equip the space \mathbb{R}^m
 491 with the ℓ_1 norm. Our goal is to show that (P) is sharp in terms of ℓ_1 norm with
 492 parameters $(\frac{1}{2}, \frac{1}{2}\sqrt{c_1 k}, 2)$ and Lipschitz continuous with Lipschitz constant $L = \sqrt{c_2 k}$
 493 w.h.p. once $m \gtrsim nk$. We will do so by comparing Covariance estimation I with
 494 Covariance estimation II.

495 First, we replace the linear constraint in Problem (P) by $\mathcal{B}(X) = d$ where $\mathcal{B} : \mathbb{S}^n \rightarrow \mathbb{R}^{\lfloor m/2 \rfloor}$, and

$$497 \quad (4.1) \quad \mathcal{B}^*(e_i) = \frac{1}{2} a_{2i-1} a_{2i-1}^\top - \frac{1}{2} a_{2i} a_{2i}^\top \quad d_i = \frac{1}{2} b_{2i-1} - \frac{1}{2} b_{2i}, \quad i = 1, \dots, \lfloor m/2 \rfloor.$$

498 This is distributed as an instance of Covariance estimation II. By Theorem 4.1 (Co-
 499 variance estimation II), we know that once $m \gtrsim nk$, it holds that

$$500 \quad (4.2) \quad \text{tr}(X) + \sqrt{c_1 k} \|\mathcal{B}(X) - d\|_1 + 2 \text{dist}_1(X, \mathbb{S}_+^n) - \text{tr}(X^\natural) \geq \frac{1}{2} \|X - X^\natural\|_1.$$

501 Note that for each $i = 1, \dots, \lfloor m/2 \rfloor$, by the construction (4.1), we have

$$502 \quad |\langle \mathcal{B}^*(e_i), X \rangle - d_i| \leq \frac{1}{2} |\langle \mathcal{A}^*(e_{2i-1}), X \rangle - b_{2i-1}| + \frac{1}{2} |\langle \mathcal{A}^*(e_{2i}), X \rangle - b_{2i}|.$$

503 Combining this fact with (4.2), we see that the function $\text{tr}(X) + \frac{1}{2}\sqrt{c_1 k} \|\mathcal{A}(X) - b\|_1 +$
504 $2 \text{dist}_1(X, \mathbb{S}_+^n)$ is $\frac{1}{2}$ sharp around X^\natural as well.

505 To prove the Lipschitz constant is bounded, we utilize [14, Lemma 3.1], which
506 shows that w.h.p.

$$507 \quad \|\mathcal{A}(X)\|_1 \leq 1.1 \|X\|_1.$$

508 Hence we see $\text{tr}(X) + \frac{1}{2}\sqrt{c_1 k} \|\mathcal{A}(X) - b\|_1 + 2 \text{dist}_1(X, \mathbb{S}_+^n)$ is $3 + \sqrt{c_1 k}$ Lipschitz with
509 respect to the ℓ_1 norm. This completes the proof of Theorem 4.1 (Covariance estimation I).
510

511 **5. Sharp problem formulations in the presence of noise.** In this section,
512 we show that the sharpness of a problem (P) in the noiseless setting $b = \mathcal{A}(x^\natural)$ provides
513 (algorithmically useful) information even in the noisy setting, where b is replaced by
514 $\tilde{b} = \mathcal{A}(x^\natural) + \delta$ with δ small or sparse. We begin with the case where δ is small.

515 **PROPOSITION 5.1.** *Suppose (P) is (μ, r, ℓ) sharp around x^\natural in the ℓ_1 norm. Let*
516 $\delta \in \mathbb{R}^m$ and set $\tilde{b} = b + \delta$.

517 • If \tilde{x} minimizes

$$518 \quad \tilde{F}_{r,\ell}(x) := f(x) + r \left\| \mathcal{A}(x) - \tilde{b} \right\|_1 + \ell \text{dist}_1(x, \mathcal{K}),$$

519 then $\|\tilde{x} - x^\natural\|_1 \leq \frac{2r}{\mu} \|\delta\|_p$.

520 • If \tilde{x} minimizes

$$521 \quad \tilde{F}_{r,\ell}^{\text{thresh}}(x) := \max \left(\tilde{F}_{r,\ell}(x), F_{r,\ell}(x^\natural) + 3r \|\delta\|_p \right),$$

522 then $\|\tilde{x} - x^\natural\|_1 \leq \frac{4r}{\mu} \|\delta\|_p$. Furthermore, $\tilde{F}_{r,\ell}^{\text{thresh}}$ is $\frac{\mu}{2}$ -sharp around its opti-
523 mizers.

524 **Remark 5.2.** Suppose $F_{r,\ell}$ is L -Lipschitz with $f(x) = \|x\|_1$, $\mathcal{K} = V$ and $L \geq 1$ as
525 in sparse vector recovery or low-rank matrix sensing. Then we have

$$526 \quad L \|x^\natural\|_1 \geq |F_{r,\ell}(0) - F_{r,\ell}(x^\natural)| = \left| r \|b\|_p - \|x^\natural\|_1 \right| \implies (L+1) \|x^\natural\|_1 \geq r \|b\|_p.$$

527 Hence, combining this inequality with the first item of Proposition 5.1, we have

$$528 \quad \frac{\|\tilde{x} - x^\natural\|_1}{\|x^\natural\|_1} \leq \frac{2L+2}{\mu} \frac{\|\delta\|_p}{\|b\|_p}.$$

529 Thus, we see that indeed the condition number of (P) controls the relative change of
530 the solution to the relative perturbation to the data vector b :

$$531 \quad \underbrace{\frac{\|\tilde{x} - x^\natural\|_1}{\|x^\natural\|_1}}_{\text{Relative change in solution}} \leq \mathcal{O}(\kappa) \cdot \underbrace{\frac{\|\delta\|_p}{\|b\|_p}}_{\text{Relative change in data vector}}.$$

532 This is analogous to the use of condition number in linear equations.

533 *Remark 5.3.* When δ is nonzero, the penalization formulation $\tilde{F}_{r,\ell}$ is not neces-
 534 sarily a sharp function. However, the above proposition asserts that $\tilde{F}_{r,\ell}^{\text{thresh}}(x)$ is still
 535 sharp. Hence, we may hope to apply the methods described in Section 6, which apply
 536 to sharp functions: On the surface, evaluating the function $\tilde{F}_{r,\ell}^{\text{thresh}}(x)$ (and its sub-
 537 gradients) requires access to both x^\natural and $\|\delta\|_p$. While we will not have access to these
 538 quantities in practice, that is of little consequence if we only plan to apply first-order
 539 methods (as we suggest in Section 6). Indeed, any first-order method applied to $\tilde{F}_{r,\ell}$
 540 will behave equivalently to the first-order method applied to $\tilde{F}_{r,\ell}^{\text{thresh}}$ until an iterate
 541 \tilde{x} satisfying $\|\tilde{x} - x^\natural\|_1 \leq \frac{4r}{\mu} \|\delta\|_p$ is found. In particular, the algorithms presented in
 542 Section 6 applied to $\tilde{F}_{r,\ell}$ will converge linearly to such a point with a rate depending
 543 on $\frac{\mu}{2}$. We emphasize that such a procedure is *adaptive* (to the noise level $\|\delta\|_p$) in
 544 both $\|\tilde{x} - x^\natural\|_1$ and the rate of convergence to \tilde{x} .

545 *Proof of Proposition 5.1.* For notational convenience, in this proof we will drop
 546 all subscripts r, ℓ .

547 By the triangle inequality, for all $x \in V$, we have

$$548 \quad \left| F(x) - \tilde{F}(x) \right| = \left| r \|\mathcal{A}(x) - b\|_p - r \|\mathcal{A}(x) - \tilde{b}\|_p \right| \leq r \|\delta\|_p.$$

549 For the first claim, note that by optimality of \tilde{x} in \tilde{F} , we have that

$$550 \quad F(\tilde{x}) \leq \tilde{F}(\tilde{x}) + r \|\delta\|_p \leq \tilde{F}(x^\natural) + r \|\delta\|_p \leq F(x^\natural) + 2r \|\delta\|_p.$$

551 Combining this inequality with μ -sharpness of F around x^\natural proves the first claim.

552 Consider the second claim. Note that $\tilde{F}(x^\natural) \leq F(x^\natural) + r \|\delta\|_p$ so that \tilde{F} achieves
 553 values bounded above by the threshold value of $F(x^\natural) + 3r \|\delta\|_p$. Thus, the set of
 554 minimizers of $\tilde{F}^{\text{thresh}}(x)$ is given by $\mathcal{X} := \left\{ x \in V : \tilde{F}(x) \leq F(x^\natural) + 3r \|\delta\|_p \right\}$. Then,
 555 if \tilde{x} minimizes $\tilde{F}^{\text{thresh}}(x)$, we must have

$$556 \quad F(\tilde{x}) \leq \tilde{F}(\tilde{x}) + r \|\delta\|_p \leq F(x^\natural) + 4r \|\delta\|_p.$$

557 Combining this inequality with μ -sharpness of F around x^\natural shows that $\|\tilde{x} - x^\natural\|_1 \leq$
 558 $\frac{4r}{\mu} \|\delta\|_p$.

559 It remains to show that $\tilde{F}^{\text{thresh}}(x)$ is $\mu/2$ sharp around its optimizers \mathcal{X} . By the
 560 definition of sharpness (see Definition 2.1), the goal is to show that for any $\bar{x} \in V \setminus \mathcal{X}$,
 561 there exists $\tilde{x} \in \mathcal{X}$ satisfying

$$562 \quad \frac{\mu}{2} \|\bar{x} - \tilde{x}\|_1 \leq \tilde{F}^{\text{thresh}}(\bar{x}) - \tilde{F}^{\text{thresh}}(\tilde{x}).$$

563 Note that for any $\bar{x} \in V \setminus \mathcal{X}$ and $\tilde{x} \in \mathcal{X}$, we have $\tilde{F}^{\text{thresh}}(\bar{x}) = \tilde{F}(\bar{x})$ and $\tilde{F}^{\text{thresh}}(\tilde{x}) =$
 564 $F(x^\natural) + 3r \|\delta\|_p$.

565 Set $\tilde{x} = (1 - \alpha)x^\natural + \alpha\bar{x}$ where

$$566 \quad \alpha = \frac{2r \|\delta\|_p}{\tilde{F}(\bar{x}) - F(x^\natural) - r \|\delta\|_p}.$$

567 As $\tilde{F}(\bar{x}) > F(x^\natural) + 3r \|\delta\|_p$, we have that α is well-defined and $\alpha \in [0, 1]$. By convexity

568 of \tilde{F} ,

$$\begin{aligned}
569 \quad \tilde{F}(\tilde{x}) &\leq (1 - \alpha)\tilde{F}(x^\natural) + \alpha\tilde{F}(\bar{x}) \\
570 &\leq (1 - \alpha)(F(x^\natural) + r \|\delta\|_p) + \alpha\tilde{F}(\bar{x}) \\
571 &= F(x^\natural) + r \|\delta\|_p + \left(\frac{2r \|\delta\|_p}{\tilde{F}(\bar{x}) - F(x^\natural) - r \|\delta\|_p} \right) (\tilde{F}(\bar{x}) - F(x^\natural) - r \|\delta\|_p) \\
572 &= F(x^\natural) + 3r \|\delta\|_p.
\end{aligned}$$

573 We deduce that $\tilde{x} \in \mathcal{X}$.

574 Finally, we have

$$\begin{aligned}
575 \quad \frac{\mu}{2} \|\bar{x} - \tilde{x}\|_1 &= (1 - \alpha) \frac{\mu}{2} \|x^\natural - \bar{x}\|_1 \\
576 &\leq \frac{1 - \alpha}{2} (F(\bar{x}) - F(x^\natural)) \\
577 &\leq \frac{1}{2} \left(\frac{\tilde{F}(\bar{x}) - F(x^\natural) - 3r \|\delta\|_p}{\tilde{F}(\bar{x}) - F(x^\natural) - r \|\delta\|_p} \right) (\tilde{F}(\bar{x}) - F(x^\natural) + r \|\delta\|_p) \\
578 &= \frac{1}{2} (\tilde{F}(\bar{x}) - F(x^\natural) - 3r \|\delta\|_p) \left(\frac{\tilde{F}(\bar{x}) - F(x^\natural) - 3r \|\delta\|_p + 4r \|\delta\|_p}{\tilde{F}(\bar{x}) - F(x^\natural) - 3r \|\delta\|_p + 2r \|\delta\|_p} \right) \\
579 &\leq \tilde{F}(\bar{x}) - F(x^\natural) - 3r \|\delta\|_p.
\end{aligned}$$

580 Here, the second line follows by μ -sharpness of F , the third line follows by the defini-
581 tion of α , and the final line follows from the premise that $\tilde{F}(\bar{x}) > F(x^\natural) + 3r \|\delta\|_p$. \square

582 *Example 5.4.* If the norm on the measurement space is $\ell_p = \ell_1$, then Proposi-
583 tion 5.1 is tight up to constants. Indeed, consider the sparse recovery problem below
584 (parameterized by $\mu \in (0, 1]$ and $r > 0$):

$$\begin{aligned}
585 \quad V &= \mathbb{R}^2, \quad x^\natural = e_1 \\
586 \quad \mathcal{A} &= \frac{1}{r} \begin{bmatrix} 1 & 0 \\ \mu & 1 - \mu \end{bmatrix}, \quad b = \frac{1}{r} \begin{bmatrix} 1 \\ \mu \end{bmatrix}.
\end{aligned}$$

587 Note that

$$\begin{aligned}
588 \quad F(x) &:= \|x\|_1 + r \|\mathcal{A}(x) - b\|_1 \\
589 &= |x_1| + |x_2| + |x_1 - 1| + |\mu(x_1 - 1) + (1 - \mu)x_2| \\
590 &\geq 1 + \mu |x_1 - 1| + \mu |x_2| \\
591 &= F(x^\natural) + \mu \|x - x^\natural\|_1.
\end{aligned}$$

592 Thus, F is μ -sharp around x^\natural w.r.t. the ℓ_1 norm. Now, suppose instead of b , we receive
593 $\tilde{b} = b - \delta e_2$ where $\delta \in [0, \mu/r]$. A similar calculation shows that $\tilde{x} := (1 - r\delta/\mu)e_1$ is
594 the minimizer of $\tilde{F}(x)$:

$$\begin{aligned}
595 \quad \tilde{F}(x) &:= \|x\|_1 + r \left\| \mathcal{A}(x) - \tilde{b} \right\|_1 \\
596 &= |x_1| + |x_2| + |x_1 - 1| + |\mu(x_1 - 1) + (1 - \mu)x_2 + r\delta| \\
597 &\geq 1 + \mu \left| x_1 - 1 + \frac{r\delta}{\mu} \right| + \mu |x_2| \\
598 &= \tilde{F}(\tilde{x}) + \mu \|x - \tilde{x}\|_1.
\end{aligned}$$

599 We deduce that \tilde{x} , the minimizer of $\tilde{F}(x)$, satisfies $\|\tilde{x} - x^\natural\|_1 = \frac{r}{\mu} \|\delta\|_1$. This example
 600 can be extended to any $n, m \geq 2$ by padding.

601 Thus, if the only structural property of (\mathcal{A}, b) that we are given is that $\mathcal{F}_{r,\ell}$ is
 602 sharp in the noiseless case and δ may be chosen adversarially, then a recovery error
 603 of $\|x - x^\natural\|_1 = O\left(\frac{r\|\delta\|_1}{\mu}\right)$ is optimal. This matches Proposition 5.1.

604 *Remark 5.5.* We believe the bound in Proposition 5.1 to be *suboptimal* if the error
 605 δ is generated randomly. For example, if we take the norm on the measurement space
 606 to be $\ell_p = \ell_2$ and generate \mathcal{A} according to subsection 3.2, $\delta_i \stackrel{\text{i.i.d.}}{\sim} N(0, \sigma^2/m)$, with m
 607 large enough, then Proposition 5.1 guarantees $\|\tilde{x} - x^\natural\|_1 = O\left(\sigma\sqrt{k}\right)$. One may com-
 608 pare this with the nearly statistically optimal error $\|\tilde{x} - x^\natural\|_2 = O\left(\sigma\sqrt{k\log(n)/m}\right)$
 609 under the same assumptions [6]. It is not clear how to extend the nearly optimal
 610 result to the setting where the measurement space is equipped with the ℓ_p norm (for
 611 $p < 2$) and stochastic noise. Indeed, the argument in Proposition 5.1 considers only
 612 the norm on $\|\delta\|_p$ and is fundamentally unable to control cancellations in dealing with
 613 $\mathcal{A}(x) - \tilde{b}$. On the other hand, the argument in [6] specifically uses Euclidean properties
 614 of the ℓ_2 norm and does not extend to general ℓ_p norms. The choice $p = 1$ specifically
 615 is desirable as this enables robustness to grossly-but-sparsely-corrupted observations
 616 (see the following proposition). We leave the detailed analysis in the stochastic noise
 617 setting for future work.

618 The next proposition shows that if the measurement space is equipped with the
 619 $\ell_p = \ell_1$ norm and (P) is sharp in the noiseless case, then exact recovery continues to be
 620 possible (with linearly convergent algorithms) in the presence of grossly-but-sparsely-
 621 corrupted observations. Since the recovery procedure is invariant under permutations
 622 of the measurement space, we may assume without loss of generality that corruption
 623 occurs on the final portion of the measurement vector. Let $\mathbb{R}^{\tilde{m}}$ denote the measure-
 624 ment space. We decompose $\mathbb{R}^{\tilde{m}} = \mathbb{R}^m \times \mathbb{R}^{m'}$ so that the first m observations are
 625 uncorrupted and the remaining m' coordinates are arbitrarily corrupted.

626 **PROPOSITION 5.6.** *Fix V and let $\mathcal{A} : V \rightarrow \mathbb{R}^m$ and $\mathcal{B} : V \rightarrow \mathbb{R}^{m'}$ denote sensing*
 627 *operators. Equip the measurement space $\mathbb{R}^{\tilde{m}} = \mathbb{R}^m \times \mathbb{R}^{m'}$ with the $\ell_p = \ell_1$ norm.*
 628 *Suppose*

$$629 \quad F_{r,\ell}(x) := f(x) + r \|\mathcal{A}(x) - \mathcal{A}(x^\natural)\|_1 + \ell \text{dist}_1(x, \mathcal{K})$$

630 *is μ -sharp around x^\natural w.r.t. the ℓ_1 norm. Let $\|\mathcal{B}\|_{1 \rightarrow 1}$ denote the operator norm*

$$631 \quad \max_{z \in V: \|z\|_1 \leq 1} \|\mathcal{B}(z)\|_1.$$

632 *If $\|\mathcal{B}\|_{1 \rightarrow 1} < \frac{\mu}{r}$, then the function*

$$633 \quad (5.1) \quad \tilde{F}_{r,\ell}(x) := f(x) + r \left\| \begin{bmatrix} \mathcal{A}(x) \\ \mathcal{B}(x) \end{bmatrix} - \begin{bmatrix} \mathcal{A}(x^\natural) \\ \mathcal{B}(x^\natural) + \delta \end{bmatrix} \right\|_1 + \ell \text{dist}_1(x, \mathcal{K})$$

634 *is $(\mu - r \|\mathcal{B}\|_{1 \rightarrow 1})$ -sharp around x^\natural .*

635 *Proof.* Let $x \in V$. Then,

$$636 \quad \begin{aligned} \tilde{F}_{r,\ell}(x) - \tilde{F}_{r,\ell}(x^\natural) &= F_{r,\ell}(x) - F_{r,\ell}(x^\natural) + r \|\mathcal{B}(x) - \mathcal{B}(x^\natural) - \delta\|_1 - r \|\delta\|_1 \\ 637 \quad &\geq \mu \|x - x^\natural\|_1 - r \|\mathcal{B}(x) - \mathcal{B}(x^\natural)\|_1 \\ 638 \quad &\geq (\mu - r \|\mathcal{B}\|_{1 \rightarrow 1}) \|x - x^\natural\|_1. \end{aligned}$$

639 Here, the first inequality is the triangle inequality and the second inequality is the
640 definition of the operator norm. \square

641 *Example 5.7.* Consider the phase retrieval problem where an α -fraction of the
642 observations are corrupted arbitrarily. Formally, consider the following procedure:
643 Let $V = \mathbb{S}^n$, $m = \lceil (1 - \alpha)\tilde{m} \rceil$ and $m' = \tilde{m} - m$. Fix $X^\natural \in \mathbb{S}_+^n \subseteq V$ to be rank-one
644 and let $\tilde{\mathcal{A}} : V \rightarrow \mathbb{R}^{\tilde{m}}$ be the random linear map

$$645 \quad \tilde{\mathcal{A}}(X)_i = g_i^\top X g_i, \quad \text{where } g_i \sim N(0, I/\tilde{m}).$$

646 Let $\delta \in \mathbb{R}^{\tilde{m}}$ denote an arbitrary vector, chosen possibly adversarially, with only the
647 guarantee that $\text{supp}(\delta) \subseteq [m + 1, \tilde{m}]$. Set $\tilde{b} = \tilde{\mathcal{A}}(X^\natural) + \delta$.

648 Let $\mathcal{A} : V \rightarrow \mathbb{R}^m$ and $\mathcal{B} : V \rightarrow \mathbb{R}^{m'}$ denote the decomposition of $\tilde{\mathcal{A}} : V \rightarrow \mathbb{R}^{\tilde{m}}$
649 corresponding to the decomposition $\mathbb{R}^{\tilde{m}} = \mathbb{R}^m \times \mathbb{R}^{m'}$. Slightly abusing notation, we
650 will also let δ denote the restriction of δ to $\mathbb{R}^{m'}$.

651 Our goal is to recover X^\natural from

$$652 \quad \tilde{F}_{r,\ell}(X) := \text{tr}(X) + r \left\| \tilde{\mathcal{A}}(X) - \tilde{b} \right\|_1 + \ell \text{dist}_1(X, \mathbb{S}_+^n).$$

653 Now, suppose that

$$654 \quad \min_{X \in \mathbb{S}_+^n} \left\{ \text{tr}(X) : \begin{array}{l} \mathcal{A}(X) = \mathcal{A}(X^\natural) \\ X \in \mathbb{S}_+^n \end{array} \right\}$$

655 is (μ, r, ℓ) sharp around X^\natural . Proposition 5.6 states that if $\|\mathcal{B}\|_{1 \rightarrow 1} < \frac{\mu}{r}$, then X^\natural is
656 the unique minimizer of $\tilde{F}_{r,\ell}$ despite the corruption δ . By [52, Corollary 5.35], we
657 know that w.h.p., $\|\mathcal{B}\|_{1 \rightarrow 1} \leq \left(\frac{\sqrt{\alpha\tilde{m} + 2\sqrt{n}}}{\sqrt{\tilde{m}}} \right)^2 \leq 2\alpha + 8(n/\tilde{m})$. Combining these bounds,
658 we have that X^\natural is the unique minimizer of the sharp unconstrained minimization
659 problem $\tilde{F}_{r,\ell}$ if $\alpha \lesssim \mu/r$ and $n/\tilde{m} \lesssim \mu/r$. In particular, μ/r controls the fraction of
660 allowed gross corruption in the observation $\tilde{\mathcal{A}}(X^\natural)$.

661 **6. First-order methods for non-Euclidean sharp minimization.** We next
662 turn to algorithms for minimizing the sharp nonsmooth formulations proposed in this
663 paper. Our goal is to take advantage of the conditioning results (Theorem 4.1) and
664 to design a first-order method tailored to sharp Lipschitz functions in the ℓ_1 norm.

665 Throughout this section, we make the following blanket assumption:

666 **ASSUMPTION 6.1.** *Let $V = \mathbb{R}^n$, \mathbb{S}^n , or $\mathbb{R}^{n \times N}$. In the third case, assume without*
667 *loss of generality that $n \leq N$. Assume that $f : V \rightarrow \mathbb{R}$ is a convex function with*
668 *minimizers.*

669 Instead of directly designing algorithms for sharp Lipschitz functions w.r.t. the ℓ_1
670 norm, we will design algorithms for sharp Lipschitz functions w.r.t. the ℓ_p norm for
671 all $p \in (1, 2]$ and view the $\ell_p = \ell_1$ case as a limiting case.

672 Let $p \in [1, 2]$. Recall that a function f is said to be μ -sharp w.r.t. the ℓ_p norm if

$$673 \quad f(x) - f^* \geq \mu \text{dist}_p \left(x, \arg \min_z f(z) \right) \quad \forall x \in V,$$

674 where $f^* = \min_z f(z)$, and L -Lipschitz w.r.t. the ℓ_p norm if

$$675 \quad |f(x) - f(y)| \leq L \|x - y\|_p \quad \forall x, y \in V.$$

676 We caution that the ℓ_p norm in this section is the norm on the signal space V and
 677 not the norm on the measurement space \mathbb{R}^m as in prior sections. In prior sections,
 678 we always fixed $\ell_p = \ell_1$ in the signal space V .

679 Before we discuss our new algorithm, we first point out that sharp functions that
 680 are well-conditioned in ℓ_1 are necessarily poorly conditioned in ℓ_2 .

681 **PROPOSITION 6.2.** *Suppose $f : V \rightarrow \mathbb{R}$ is μ -sharp and L -Lipschitz in the ℓ_1 norm
 682 around $x^\natural \in V$. If f is α -sharp and β -Lipschitz in the ℓ_2 norm around $x^\natural \in V$, then
 683 α and β must satisfy*

$$684 \quad \alpha \leq L, \quad \text{and} \quad \beta \geq \mu\sqrt{n}.$$

685 *Proof.* First, suppose $V = \mathbb{R}^n$. Let $y = x^\natural + (\text{any } 1\text{-sparse vector})$. Then,

$$686 \quad \alpha \|y - x^\natural\|_2 \leq f(y) - f(x^\natural) \leq L \|y - x^\natural\|_1.$$

687 We deduce that $\alpha \leq L$. Similarly, let $z = x^\natural + \mathbf{1}_n$ (the all-ones vector). Then,

$$688 \quad \beta \|z - x^\natural\|_2 \geq f(z) - f(x^\natural) \geq \mu \|z - x^\natural\|_1.$$

689 We deduce that $\beta \geq \mu\sqrt{n}$.

690 The other two settings are proved analogously: take $y = x^\natural + (\text{any rank-1 matrix})$
 691 and $z = x^\natural + I_n$ where I_n is the $n \times n$ identity matrix padded to $n \times N$. \square

692 This proposition tells us that the functions $F_{r,\ell}$ from Theorem 4.1 have constant
 693 sharpness in the ℓ_2 norm and an $\Omega(\sqrt{n})$ Lipschitz constant in the ℓ_2 norm. Thus,
 694 first-order methods for sharp Lipschitz functions in the ℓ_2 norm [32, 42, 53] can at
 695 best guarantee a convergence rate of $O\left(n \log\left(\frac{1}{\epsilon}\right)\right)$. This dependence on n precludes
 696 the use of such algorithms (for example, the Polyak variant of the subgradient method
 697 **Polyak-Subgrad**) on large-scale instances of our signal recovery problems (see Fig-
 698 ure 1a) and motivates our development of an algorithm with almost dimension-free
 699 iteration complexity in the next subsection.

700 **6.1. Restarted Mirror Descent.** We now describe the restarted mirror de-
 701 scent (RMD) algorithm, Algorithm 6.2. This algorithm generalizes similar algo-
 702 rithms for minimizing sharp functions in a Euclidean norm [32, 42, 53] and has nearly
 703 dimension-independent linear convergence rates that depend explicitly on sharpness
 704 (see Theorem 6.6) in an ℓ_p norm ($p \in [1, 2]$). Algorithm 6.2 can be applied to $F_{r,\ell}$,
 705 the sharp exact penalty formulation of (P), or any of its sharp perturbations (see
 706 Propositions 5.1 and 5.6).

707 Recall the mirror descent algorithm and its guarantee [8, Theorem 4.2].

Algorithm 6.1 Mirror Descent

Given $f : V \rightarrow \mathbb{R}$, $\bar{x} \in V$, $\eta > 0$, $T \in \mathbb{N}$, $h : V \rightarrow \mathbb{R}$

- Let $x_0 = \bar{x}$, $\theta_0 = 0 \in V^*$
 - For $t = 1, \dots, T$
 - Let $g_t \in \partial f(x_{t-1})$ and set $\theta_t = \theta_{t-1} - \eta \cdot g_t$
 - Set $x_t = (\nabla h)^{-1}(\theta_t)$
 - Output the x_t minimizing $f(x_t)$ among $t \in [0, T]$.
-

708 Here, V^* is the dual space of V . As a vector space, V^* can be identified with V via
 709 the standard inner product if $V = \mathbb{R}^n$ and via the trace inner product if $V = \mathbb{R}^{n \times N}$
 710 or $V = \mathbb{S}^n$. The norm on V^* is the ℓ_q norm where q is the Hölder dual to p , i.e., the
 711 quantities q and p satisfies that $\frac{1}{p} + \frac{1}{q} = 1$.

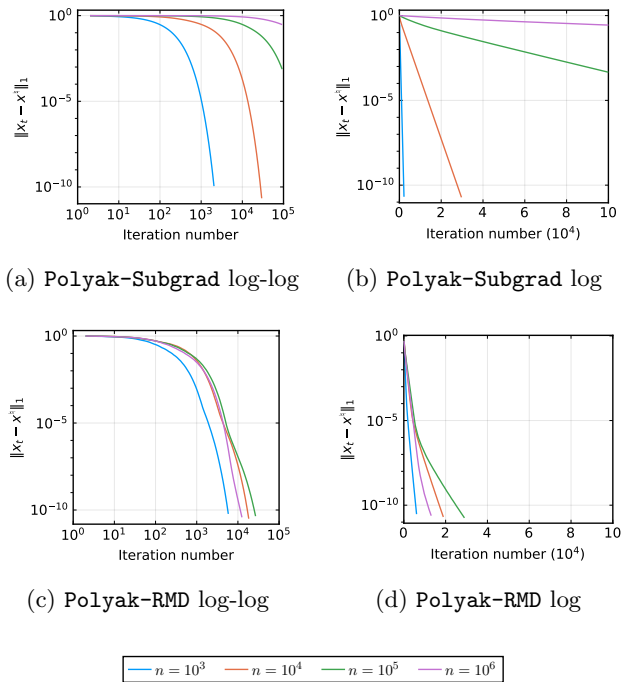


Fig. 1: The convergence of Polyak-Subgrad and Polyak-RMD (described in Section 6.1) on sparse vector recovery for different values of n , the ambient dimension. The support size of x^h is $k = 5$ and the sensing map \mathcal{A} is generated as described in Section 4.1. The number of observations is $m = 2T$ where $T = (2k \log(n/k) + 1.25k + 1)$ is the current best estimate of the statistical threshold for sparse recovery [17]. See https://github.com/alexlihengwang/sharpness_well_conditioning for implementation details. Note Polyak-Subgrad deteriorates as the dimension n increases while Polyak-RMD is relatively insensitive to the dimension in terms of iteration number. Furthermore, one can see in subfigure (d) that Polyak-RMD achieves a linear convergence rate as predicted by Theorem 6.6 and Proposition 6.9.

712 LEMMA 6.3. Let $h : V \rightarrow \mathbb{R}$ be differentiable and σ -strongly convex with respect
713 to the ℓ_p norm. Suppose $f : V \rightarrow \mathbb{R}$ is convex and L -Lipschitz with respect to the ℓ_p
714 norm. Then, the mirror descent algorithm initialized at \bar{x} , run for t iterations, with
715 step-size η , produces \tilde{x} such that

$$716 \quad f(\tilde{x}) - f(x^*) \leq \frac{L^2 \eta}{2\sigma} + \frac{D_h(x^* || \bar{x})}{\eta t}$$

717 where $x^* \in V$ is any minimizer of f . Here, each iteration requires computing a single
718 subgradient of f and applying $(\nabla h)^{-1}$ and arithmetic operations in V and V^* .

719 In the bound above, the quantity $D_h(\cdot || \cdot)$ is a Bregman divergence term. We
720 elaborate on this term for our choice of the map h . Let $p \in (1, 2]$ and define $h_{\bar{x}}(x) :=$
721 $\frac{1}{2} \|x - \bar{x}\|_p^2$ for $\bar{x} \in V$. It is known [3, Theorem 1] that in each of the three setups
722 in Assumption 6.1, that $h_{\bar{x}}$ is differentiable and $(p - 1)$ -strongly convex w.r.t. the ℓ_p
723 norm. Furthermore, for all $x^* \in V$, the Bregman divergence (associated to $h_{\bar{x}}$) of x^*

724 with respect to \bar{x} is

$$725 \quad D_{h_{\bar{x}}}(x^*|\bar{x}) := h_{\bar{x}}(x^*) - h_{\bar{x}}(\bar{x}) - \langle \nabla h_{\bar{x}}(\bar{x}), x^* - \bar{x} \rangle = h_{\bar{x}}(x^*) = \frac{1}{2} \|x^* - \bar{x}\|_p^2.$$

726 Thus, if f is L -Lipschitz w.r.t. the ℓ_p norm, the output $\tilde{x} = \text{mirror}(h_{\bar{x}}, f, L, \bar{x}, t, \eta)$
727 has suboptimality bounded by

$$728 \quad f(\tilde{x}) - f(x^*) \leq \frac{L^2 \eta}{2(p-1)} + \frac{\|x^* - \bar{x}\|_p^2}{2\eta T}.$$

729 This bound holds simultaneously for all minimizers $x^* \in V$ of f .

730 *Remark 6.4.* Consider a setup from Assumption 6.1 and let $p \in (1, 2]$. In each
731 iteration of mirror descent (applied with $h_{\bar{x}}$), we must evaluate $(\nabla h_{\bar{x}})^{-1}$ on some
732 input $\theta \in V^*$. By [46, Corollary 23.5.1], this is equivalent to evaluating $\nabla h_{\bar{x}}^*(\theta)$ where
733 $h_{\bar{x}}^*$ is the convex conjugate of $h_{\bar{x}}$. Thus,

$$734 \quad (\nabla h_{\bar{x}})^{-1}(\theta) = \nabla h_{\bar{x}}^*(\theta) = \bar{x} + \nabla \frac{1}{2} \|\theta\|_q^2 = \bar{x} + \frac{\text{sign}(\theta) \circ |\theta|^{q-1}}{\|\theta\|_q^{q-2}}.$$

735 Here, q is the Hölder dual to p and the expression $\text{sign}(\theta) \circ |\theta|^{q-1}$ is applied entrywise
736 to the entries of θ if θ is a vector and to the singular values of θ if θ is a matrix.

737 *Remark 6.5.* In the matrix setting, computing the mirror map requires perform-
738 ing a singular value decomposition on θ . While this is an expensive operation, we
739 emphasize that this is the same as the cost of simply *computing a subgradient* in all
740 of the matrix recovery applications considered in this paper. For example, in the
741 low-rank matrix sensing applications, computing the subgradient of

$$742 \quad \|X\|_1 + r \|\mathcal{A}(X) - b\|_1$$

743 requires performing an SVD on X , and in covariance estimation, computing the sub-
744 gradient of

$$745 \quad \text{tr}(X) + r \|\mathcal{A}(X) - b\|_1 + \ell \text{dist}_1(X, \mathbb{S}_+^n)$$

746 requires performing an eigenvalue decomposition (EVD) on X , which can be obtained
747 from an SVD of X or computed at a similar cost of an SVD of X . Thus, the work
748 performed in each iteration of mirror descent is comparable to the work performed
749 in each iteration of subgradient descent: two SVDs for mirror descent (one for the
750 subgradient and one for the mirror map) and one SVD for subgradient descent. We
751 expect that this full SVD may be replaced by a partial SVD near an optimal low-rank
752 solution (following [31]) and leave this extension for future work.

753 Now consider the following restarted variant of mirror descent where the step size
754 *and* mirror map h update at each restart.

Algorithm 6.2 $\text{RMD}(f, L, x_0, K, \{\eta_k\}, t, p)$

- For $k = 1, 2, \dots, K$
 - Set $x_k \leftarrow \text{mirror}(h_{x_{k-1}}, f, L, x_{k-1}, t, \eta_k)$ where

$$h_{x_{k-1}}(x) := \frac{1}{2} \|x - x_{k-1}\|_p^2$$

- Output x_K
-

755 Readers who prefer a more concrete description of RMD can find one for the setting
756 of sparse recovery in Appendix A.

757 **THEOREM 6.6.** *Consider a setup from Assumption 6.1 and let $p \in (1, 2]$. Suppose*
758 *$f : V \rightarrow \mathbb{R}$ is L -Lipschitz and μ -sharp w.r.t. the ℓ_p norm and $x_0 \in V$ satisfies*
759 *$f(x_0) - f^* \leq \epsilon_0$. Let $\epsilon_k = \epsilon_0 e^{-k/2}$, $K = \lceil 2 \ln(\frac{\epsilon_0}{\epsilon}) \rceil$, $t = \lceil \frac{\epsilon L^2}{\mu^2(p-1)} \rceil$, and $\eta_k = \frac{(p-1)\epsilon_k}{L^2}$.*
760 *Then, each iterate x_k in $\text{RMD}(f, L, x_0, K, \{\eta_k\}, t, p)$ satisfies $f(x_k) - f^* \leq \epsilon_k$. In*
761 *particular, an $0 < \epsilon \leq \epsilon_0$ -optimizer can be computed in*

$$762 \quad O\left(\frac{L^2}{\mu^2(p-1)} \log\left(\frac{\epsilon_0}{\epsilon}\right)\right)$$

763 *total mirror descent steps. If $f : V \rightarrow \mathbb{R}$ is L -Lipschitz and μ -sharp w.r.t. the ℓ_1*
764 *norm, we may apply the above statement with $p = 1 + \frac{1}{\ln n}$, sharpness μ , and Lipschitz*
765 *constant eL w.r.t. the ℓ_p norm.*

766 *Proof.* We first handle the case $p \in (1, 2]$. The claim holds for $k = 0$. Now let $k \geq$
767 1 . By μ -sharpness of f , there exists an optimizer x^* of f , satisfying $\mu \|x_{k-1} - x^*\|_p \leq$
768 $f(x_{k-1}) - f^*$. By Lemma 6.3, we have

$$769 \quad f(x_k) - f^* \leq \frac{L^2 \eta_k}{2(p-1)} + \frac{\|x^* - x_{k-1}\|_p^2}{2\eta_k t}$$

$$770 \quad \leq \frac{L^2 \eta_k}{2(p-1)} + \frac{e\epsilon_k^2}{2\mu^2 \eta_k t}$$

$$771 \quad \leq \frac{L^2 \eta_k}{2(p-1)} + \frac{\epsilon_k^2 (p-1)}{2\eta_k L^2}$$

$$772 \quad = \epsilon_k.$$

773 The setting of L -Lipschitz, μ -sharp convex functions w.r.t. the ℓ_1 norm reduces
774 to the setting of $p = 1 + \frac{1}{\ln(n)}$ by the bounds

$$775 \quad \frac{1}{e} \|w\|_1 \leq \|w\|_p \leq \|w\|_1, \quad \forall w \in V,$$

776 which hold [3] in each of the setups in Assumption 6.1. Specifically, these inequalities
777 imply that f is μ -sharp and eL -Lipschitz w.r.t. the ℓ_p norm. \square

778 *Remark 6.7.* In [48], the authors consider restarted versions of various accelerated
779 first order methods for sharp problems including problems in non-Euclidean spaces.
780 In the non-Euclidean setting, [48] suggests algorithms for functions f that satisfy the
781 following proxy for sharpness:

$$782 \quad (6.1) \quad f(x) - f(x^\natural) \geq \mu \sqrt{D_h(x^\natural \| x)}.$$

783 Here, $h : V \rightarrow \mathbb{R}$ is a *fixed* differentiable and σ -strongly convex function¹⁰ and D_h
784 is the Bregman divergence induced by h . This condition is inspired by sharpness in
785 the ℓ_2 setting. Indeed, if we take $h(x) = \frac{1}{2} \|x\|_2^2$ to be the standard prox setup for a
786 Euclidean space, then

$$787 \quad \mu \sqrt{D_h(x^\natural \| x)} = \frac{\mu}{\sqrt{2}} \|x^\natural - x\|_2.$$

¹⁰The notation of [48] uses a 1-strongly convex function, but this only changes the value of μ by a factor of $\sqrt{\sigma}$.

788 In other words, under the standard Euclidean prox setup, (6.1) is equivalent to $\frac{\mu}{\sqrt{2}}$ -
 789 sharpness w.r.t. the ℓ_2 norm. Assuming that the condition (6.1) holds, calculations al-
 790 most identical to the proof of Theorem 6.6 show that restarted mirror descent with the
 791 mirror map h in every restart produces an ϵ -suboptimal solution in $O\left(\frac{L^2}{\mu^2\sigma} \log\left(\frac{\epsilon_0}{\epsilon}\right)\right)$
 792 mirror descent steps.

793 Unfortunately, (6.1) is *not* implied by sharpness in the ℓ_p norm if $p \in [1, 2)$ and
 794 may be overly restrictive, as illustrated in Example 6.8 (see below). Specifically,
 795 in standard prox setups for mirror descent in ℓ_p norms or Schatten- p norms with
 796 $p \in [1, 2)$, (6.1) cannot hold for any $\mu > 0$ unless x^\natural has full support in the vector case
 797 or full rank in the matrix case.

798 Beyond the condition (6.1), [48] also considers extensions allowing for (possibly
 799 only local) Hölderian growth $f(x) - f(x^\natural) \geq \frac{\mu}{s}(D_h(x^\natural||x))^{s/2}$ for some $s \geq 1$. When
 800 f is nonsmooth and $s > 1$, [48, Corollary 5.4] guarantees an iteration complexity
 801 of $O(\epsilon^{-2\frac{s-1}{s}})$ after entering the local region. This guarantee can be improved to
 802 $O(\epsilon^{-\frac{s-1}{s}})$ for certain non-smooth problems by using smoothing (see [48, Proposition
 803 5.5]). However, neither of these results offer an approach to achieve the $O(\log(1/\epsilon))$
 804 iteration complexity achieved by RMD, unless $s = 1$.

805 *Example 6.8.* Let $V = \mathbb{R}^n$ and assume that $n \geq 2$. Let $p \in [1, 2)$. Let $f(x) :=$
 806 $\|x - x^\natural\|_p$ where $x^\natural = e_1$. Clearly, f is 1-sharp around x^\natural and 1-Lipschitz w.r.t. the
 807 ℓ_p norm.

808 Two standard prox setups for mirror descent [40] in this setting are to take $h(x) =$
 809 $\frac{1}{2}\|x\|_p^2$ (in the unbounded case) or $h(x) = \frac{1}{p}\|x\|_p^p$ (in the bounded case); see [40,
 810 Theorem 2.1] for the corresponding strong convexity parameters. Consider the first
 811 setting, i.e., $h(x) = \frac{1}{2}\|x\|_p^2$. Letting $x_\epsilon = x^\natural + \epsilon e_2$, we have by Bernoulli's inequality
 812 that

$$\begin{aligned} 813 \quad \sqrt{D_h(x^\natural||x_\epsilon)} &= \sqrt{\frac{1}{2}(1 + |\epsilon|^p)^{2/p} - \frac{1}{2} - \langle e_1, \epsilon e_2 \rangle} \\ 814 &= 2^{-1/2} \sqrt{(1 + |\epsilon|^p)^{2/p} - 1} \\ 815 &\geq p^{-1/2} |\epsilon|^{p/2}. \end{aligned}$$

816 Similarly, suppose $h(x) = \frac{1}{p}\|x\|_p^p$. Then,

$$817 \quad \sqrt{D_h(x^\natural||x_\epsilon)} = \sqrt{\frac{1}{p}(1 + |\epsilon|^p) - \frac{1}{p} - \langle e_1, \epsilon e_2 \rangle} = p^{-1/2} |\epsilon|^{p/2}.$$

818 On the other hand, $f(x_\epsilon) - f(x^\natural) \leq |\epsilon|$. Thus, letting $|\epsilon| \rightarrow 0$, we see that

$$819 \quad f(x_\epsilon) - f(x^\natural) \geq \mu \sqrt{D_h(x^\natural||x_\epsilon)}$$

820 cannot hold for *any* $\mu > 0$ in either of the standard prox setups.

821 An almost identical argument shows that a local Hölderian error bound of the
 822 form

$$823 \quad f(x_\epsilon) - f(x^\natural) \geq \frac{\mu}{s} D_h(x^\natural||x_\epsilon)^{s/2} \quad \forall \epsilon \text{ small}$$

824 can only hold with a positive $\mu > 0$ if $s \geq 2/p$. We conclude that the best guar-
 825 antee one could achieve by applying¹¹ [48, Corollary 5.4] or [48, Proposition 5.5]

¹¹This *assumes* that the conditions of [48, Corollary 5.4, Proposition 5.5] hold.

826 naïvely to the objective function $f(x) = \|x - \bar{x}\|_p$ would be an iteration complexity
 827 of $O(\epsilon^{(p-2)/2})$. Heuristically substituting $p \leftarrow 1$ gives $O(\epsilon^{-1/2})$ contrasting with the
 828 $O(\log(1/\epsilon))$ iteration complexity achieved by RMD.

829 The following variant of Algorithm 6.2 replaces knowledge of μ with knowledge of
 830 f^* . This “Polyak” variant achieves the same convergence rate as Algorithm 6.2 with
 831 the optimal choice of μ .

Algorithm 6.3 Polyak-RMD($f, L, x_0, K, \{\eta_k\}, f^*, \{\epsilon_k\}, p$)

- For $k = 1, 2, \dots, K$
 - Run mirror descent $\text{mirror}(h_{x_{k-1}}, f, L, x_{k-1}, \infty, \eta_k)$ until it finds an iterate x_k satisfying

$$f(x_k) - f^* \leq \epsilon_k$$

- Output x_K
-

832 The convergence guarantee for Polyak-RMD (Algorithm 6.3) and its proof are
 833 entirely identical to that of Algorithm 6.2.

834 **PROPOSITION 6.9.** *Consider a setup from Assumption 6.1 and let $p \in (1, 2]$. Suppose $f : V \rightarrow \mathbb{R}$ is L -Lipschitz and μ -sharp w.r.t. the ℓ_p norm and $x_0 \in V$ satisfies*
 835 *$f(x_0) - f^* \leq \epsilon_0$. Let $\epsilon_k = \epsilon_0 e^{-k/2}$, $K = \lceil 2 \ln(\frac{\epsilon_0}{\epsilon}) \rceil$, and $\eta_k = \frac{(p-1)\epsilon_k}{L^2}$. Then, each iter-*
 836 *ate x_k in Polyak-RMD($f, L, x_0, K, \{\eta_k\}, f^*, \{\epsilon_k\}, p$) is computed in at most $\lceil \frac{\epsilon L^2}{\mu^2(p-1)} \rceil$*
 837 *mirror descent steps. In particular, an $0 < \epsilon \leq \epsilon_0$ -minimizer can be computed in*

$$839 \quad O\left(\frac{L^2}{\mu^2(p-1)} \log\left(\frac{\epsilon_0}{\epsilon}\right)\right)$$

840 *total mirror descent steps. If $f : V \rightarrow \mathbb{R}$ is L -Lipschitz and μ -sharp w.r.t. the ℓ_1*
 841 *norm, we may apply the above statement with $p = 1 + \frac{1}{\ln n}$, sharpness μ , and Lipschitz*
 842 *constant eL w.r.t. the ℓ_p norm.*

843 **Remark 6.10.** The basic algorithmic structure of Algorithms 6.2 and 6.3 can also
 844 be extended to an adaptive variant. Specifically, it is possible to design an adaptive
 845 version of Algorithm 6.2 (similar in spirit to [45]) that does not require knowledge of
 846 either f^* or μ but would have a convergence rate of $O\left(\frac{L^2}{\mu^2(p-1)} (\log \frac{\epsilon_0}{\epsilon})^2\right)$.

847 **7. Discussion.** This paper shows that for various statistical signal recovery
 848 tasks, once the sample size is greater than a constant multiple of the recovery thresh-
 849 old, then the convex optimization problem becomes well-conditioned in the sense of
 850 sharpness w.r.t. the ℓ_1 or Schatten-1 norm. In turn, this fact shows the optimiza-
 851 tion problem is robust to measurement error and optimization error. Furthermore,
 852 the newly developed algorithm RMD is able to achieve nearly-dimension-independent
 853 convergence rates.

854 We hope this paper induces interest in the interplay between statistics and convex
 855 optimization, especially in nonsmooth formulations and algorithms for these non-
 856 smooth formulations. In particular, the following three directions might be of interest
 857 for future investigations:

- 858 • **Good conditioning beyond RIP:** This paper considers recovery tasks
 859 where strong RIP bounds have been established. Can we prove conditioning
 860 results for other statistical problems such as matrix completion [16] and phase

861 retrieval with coded diffraction patterns [11] where standard RIP bounds may
862 not hold?

863 • **Adaptive penalty parameters:** Given a concrete statistical model, we
864 can give upper bounds on r and ℓ so that the corresponding function $F_{r,\ell}$
865 is μ -sharp. On the other hand, determining these parameters may be diffi-
866 cult without a precise statistical model. Thus, are there algorithmic ways of
867 estimating or adaptively choosing these two parameters?

868 • **Beyond convexity:** In our current setup, the objective function in (P) is
869 required to be convex. For some applications in signal recovery, e.g. [29, 50],
870 it is more convenient to use a nonconvex objective function. Can we extend
871 the results of the current paper to those settings?

872 **Appendix A. Concrete description of RMD for sparse recovery.**

873 This appendix specializes RMD (Algorithm 6.2) to the sparse recovery setting, while
874 fully unpacking the inner steps of Mirror Descent as well. This appendix does not
875 present any new results beyond those already contained in the main text.

876 Our goal is to minimize

$$877 \text{ (A.1)} \quad f(x) = \|x\|_1 + r \|Ax - b\|_1,$$

878 where $b = Ax^\natural$, over choices $x \in \mathbb{R}^n$. We will assume that $f(x)$ is μ -sharp and
879 L -Lipschitz with respect to the ℓ_1 norm. Then, Theorem 6.6 guarantees that the
880 following algorithm will output a \tilde{x} satisfying $f(\tilde{x}) \leq \epsilon$ after $O\left(\frac{L^2}{\mu^2} \ln(n) \ln\left(\frac{r\|b\|_1}{\epsilon}\right)\right)$
881 iterations. Using μ -sharpness one more time gives $\|x^\natural - \tilde{x}\|_1 \leq \frac{\epsilon}{\mu}$.

Algorithm A.1 RMD-SR($r, A, b, L, \epsilon, \mu$)

- Initialize $\bar{x}_0 = 0 \in \mathbb{R}^n$, $\epsilon_0 = r \|b\|_1$, and $\eta = \frac{\epsilon_0}{\ln(n)L^2}$
 - Let $K = \lceil 2 \ln \frac{\epsilon_0}{\epsilon} \rceil$, $t = \lceil \frac{eL^2 \ln(n)}{\mu^2} \rceil$, $p = 1 + \frac{1}{\ln(n)}$, $q = 1 + \ln(n)$
 - For $k = 1, 2, \dots, K$
 - Initialize $x_0 = \bar{x}_{k-1}$, $\theta_0 = 0 \in \mathbb{R}^n$
 - Set $\eta \leftarrow \eta / \sqrt{e}$
 - For $i = 1, \dots, t$
 - * Let $g_i \in \partial f(x_{i-1})$
 - * Let $\theta_i = \theta_{i-1} - \eta g_i$
 - * Let $x_i = \bar{x}_{k-1} + \text{sign}(\theta_i) \circ |\theta_i|^{q-1} / \|\theta_i\|_q^{q-2}$
 - Set $\bar{x}_k = \arg \min_{x \in \{x_0, \dots, x_t\}} f(x)$
 - Output \bar{x}_K
-

882 One can see that RMD-SR is a double-loop method, where the outer iterates
883 $\bar{x}_0, \dots, \bar{x}_K$ are each computed using an auxiliary sequence of iterates x_0, \dots, x_t . The
884 stepsize, η , that is used in the inner loop (Mirror Descent) decays geometrically with
885 the outer iterations. Within each inner loop, the function that maps gradients θ_i to
886 iterates x_i , gets “re-centered” around the most recent outer iterate \bar{x}_{k-1} . This map
887 raises the magnitude of each entry of θ to a large power and then normalizes so that
888 $\|x_i - \bar{x}_{k-1}\|_p = \|\theta\|_q$. This is exactly the behavior of RMD (Algorithm 6.2) in the
889 general case as well.

890 **Acknowledgments.** Lijun Ding would like to thank Dmitriy Drusvyatskiy and
891 Mateo Diaz for helpful discussions.

References.

- 892
893 [1] F. ALIZADEH, J.-P. A. HAEBERLY, AND M. L. OVERTON, *Complementarity*
894 *and nondegeneracy in semidefinite programming*, Mathematical programming,
895 77 (1997), pp. 111–128.
- 896 [2] H. ATTOUCH, J. BOLTE, P. REDONT, AND A. SOUBEYRAN, *Proximal alternat-*
897 *ing minimization and projection methods for nonconvex problems: An approach*
898 *based on the kurdyka-lojasiewicz inequality*, Mathematics of operations research,
899 35 (2010), pp. 438–457.
- 900 [3] K. BALL, E. A. CARLEN, AND E. H. LIEB, *Sharp uniform convexity and smooth-*
901 *ness inequalities for trace norms*, Inequalities: Selecta of Elliott H. Lieb, (2002),
902 pp. 171–190.
- 903 [4] R. BARANIUK, M. DAVENPORT, R. DEVORE, AND M. WAKIN, *A simple proof*
904 *of the restricted isometry property for random matrices*, Constructive approxi-
905 mation, 28 (2008), pp. 253–263.
- 906 [5] H. H. BAUSCHKE, J. M. BORWEIN, AND W. LI, *Strong conical hull intersection*
907 *property, bounded linear regularity, jameson’s property (g), and error bounds in*
908 *convex optimization*, Mathematical Programming, 86 (1999), pp. 135–160.
- 909 [6] A. BELLONI, V. CHERNOZHUKOV, AND L. WANG, *Square-root lasso: pivotal*
910 *recovery of sparse signals via conic programming*, Biometrika, 98 (2011), pp. 791–
911 806.
- 912 [7] J. BOLTE, T. P. NGUYEN, J. PEYPOUQUET, AND B. W. SUTER, *From er-*
913 *ror bounds to the complexity of first-order descent methods for convex functions*,
914 Mathematical Programming, 165 (2017), pp. 471–507.
- 915 [8] S. BUBECK ET AL., *Convex optimization: Algorithms and complexity*, Founda-
916 tions and Trends® in Machine Learning, 8 (2015), pp. 231–357.
- 917 [9] T. T. CAI AND A. ZHANG, *Rop: Matrix recovery via rank-one projections*, The
918 Annals of Statistics, 43 (2015), pp. 102–138.
- 919 [10] E. J. CANDÈS AND X. LI, *Solving quadratic equations via phaselift when there*
920 *are about as many equations as unknowns*, Foundations of Computational Math-
921 ematics, 14 (2014), pp. 1017–1026.
- 922 [11] E. J. CANDÈS, X. LI, AND M. SOLTANOLKOTABI, *Phase retrieval from coded*
923 *diffraction patterns*, Applied and Computational Harmonic Analysis, 39 (2015),
924 pp. 277–299.
- 925 [12] E. J. CANDÈS AND Y. PLAN, *Tight oracle inequalities for low-rank matrix recov-*
926 *ery from a minimal number of noisy random measurements*, IEEE Transactions
927 on Information Theory, 57 (2011), pp. 2342–2359.
- 928 [13] E. J. CANDÈS, J. K. ROMBERG, AND T. TAO, *Stable signal recovery from in-*
929 *complete and inaccurate measurements*, Communications on Pure and Applied
930 Mathematics: A Journal Issued by the Courant Institute of Mathematical Sci-
931 ences, 59 (2006), pp. 1207–1223.
- 932 [14] E. J. CANDÈS, T. STROHMER, AND V. VORONINSKI, *Phaselift: Exact and*
933 *stable signal recovery from magnitude measurements via convex programming*,
934 Communications on Pure and Applied Mathematics, 66 (2013), pp. 1241–1274.
- 935 [15] E. J. CANDÈS AND T. TAO, *Decoding by linear programming*, IEEE transactions
936 on information theory, 51 (2005), pp. 4203–4215.
- 937 [16] E. J. CANDÈS AND T. TAO, *The power of convex relaxation: Near-optimal ma-*
938 *trix completion*, IEEE Transactions on Information Theory, 56 (2010), pp. 2053–
939 2080.
- 940 [17] V. CHANDRASEKARAN, B. RECHT, P. A. PARRILO, AND A. S. WILLSKY,
941 *The convex geometry of linear inverse problems*, Foundations of Computational

- 942 mathematics, 12 (2012), pp. 805–849.
- 943 [18] V. CHARISOPOULOS, Y. CHEN, D. DAVIS, M. DÍAZ, L. DING, AND D. DRUSVY-
944 ATSKIY, *Low-rank matrix recovery with composite optimization: good condi-*
945 *tioning and rapid convergence*, Foundations of Computational Mathematics, 21
946 (2021), pp. 1505–1593.
- 947 [19] V. CHARISOPOULOS, D. DAVIS, M. DÍAZ, AND D. DRUSVYATSKIY, *Compos-*
948 *ite optimization for robust blind deconvolution*, arXiv preprint arXiv:1901.01624,
949 (2019).
- 950 [20] Y. CHEN, Y. CHI, AND A. J. GOLDSMITH, *Exact and stable covariance esti-*
951 *mation from quadratic sampling via convex programming*, IEEE Transactions on
952 Information Theory, 61 (2015), pp. 4034–4059.
- 953 [21] D. DAVIS, D. DRUSVYATSKIY, K. J. MACPHEE, AND C. PAQUETTE, *Subgradi-*
954 *ent methods for sharp weakly convex functions*, Journal of Optimization Theory
955 and Applications, 179 (2018), pp. 962–982.
- 956 [22] L. DING AND M. UDELL, *A strict complementarity approach to error bound and*
957 *sensitivity of solution of conic programs*, Optimization Letters, (2022), pp. 1–24.
- 958 [23] D. DRUSVYATSKIY, A. D. IOFFE, AND A. S. LEWIS, *Generic minimizing be-*
959 *havior in semialgebraic optimization*, SIAM Journal on Optimization, 26 (2016),
960 pp. 513–534.
- 961 [24] D. DRUSVYATSKIY AND A. S. LEWIS, *Error bounds, quadratic growth, and lin-*
962 *ear convergence of proximal methods*, Mathematics of Operations Research, 43
963 (2018), pp. 919–948.
- 964 [25] D. DRUSVYATSKIY AND C. PAQUETTE, *Efficiency of minimizing compositions*
965 *of convex functions and smooth maps*, Mathematical Programming, 178 (2019),
966 pp. 503–558.
- 967 [26] J. C. DUCHI AND F. RUAN, *Solving (most) of a set of quadratic equalities:*
968 *Composite optimization for robust phase retrieval*, Information and Inference: A
969 Journal of the IMA, 8 (2019), pp. 471–529.
- 970 [27] A. EDELMAN, *Eigenvalues and condition numbers of random matrices*, SIAM
971 journal on matrix analysis and applications, 9 (1988), pp. 543–560.
- 972 [28] I. I. EREMIN, *The relaxation method of solving systems of inequalities with convex*
973 *functions on the left-hand side*, in Doklady Akademii Nauk, vol. 160, Russian
974 Academy of Sciences, 1965, pp. 994–996.
- 975 [29] J. FAN, L. DING, Y. CHEN, AND M. UDELL, *Factor group-sparse regularization*
976 *for efficient low-rank matrix recovery*, Advances in neural information processing
977 Systems, 32 (2019).
- 978 [30] M. C. FERRIS, *Weak sharp minima and exact penalty functions*, tech. report,
979 University of Wisconsin-Madison Department of Computer Sciences, 1988.
- 980 [31] D. GARBER AND A. KAPLAN, *On the efficient implementation of the matrix*
981 *exponentiated gradient algorithm for low-rank matrix optimization*, Mathematics
982 of Operations Research, (2022).
- 983 [32] J.-L. GOFFIN, *On convergence rates of subgradient optimization methods*, Math-
984 ematical programming, 13 (1977), pp. 329–347.
- 985 [33] A. J. HOFFMAN, *On approximate solutions of systems of linear inequalities*, in
986 Selected Papers Of Alan J Hoffman: With Commentary, World Scientific, 2003,
987 pp. 174–176.
- 988 [34] P. R. JOHNSTONE AND P. MOULIN, *Faster subgradient methods for functions*
989 *with hölderian growth*, Mathematical Programming, 180 (2020), pp. 417–450.
- 990 [35] A. S. LEWIS, *Nonsmooth optimization: conditioning, convergence and semi-*
991 *algebraic models*, in Proceedings of the International Congress of Mathematicians,

- 992 Seoul, vol. 4, 2014, pp. 872–895.
- 993 [36] X. LI, Z. ZHU, A. MAN-CHO SO, AND R. VIDAL, *Nonconvex robust low-rank*
994 *matrix recovery*, SIAM Journal on Optimization, 30 (2020), pp. 660–686.
- 995 [37] M. V. NAYAKKANKUPPAM AND M. L. OVERTON, *Conditioning of semidefinite*
996 *programs.*, Mathematical programming, 85 (1999).
- 997 [38] I. NECOARA, Y. NESTEROV, AND F. GLINEUR, *Linear convergence of first order*
998 *methods for non-strongly convex optimization*, Mathematical Programming, 175
999 (2019), pp. 69–107.
- 1000 [39] A. S. NEMIROVSKI AND Y. E. NESTEROV, *Optimal methods of smooth convex*
1001 *minimization*, Zhurnal Vychislitel'noi Matematiki i Matematicheskoi Fiziki, 25
1002 (1985), pp. 356–369.
- 1003 [40] Y. NESTEROV AND A. NEMIROVSKI, *On first-order algorithms for l1/nuclear*
1004 *norm minimization*, Acta Numerica, 22 (2013), pp. 509–575.
- 1005 [41] Y. PLAN AND R. VERSHYNIN, *Dimension reduction by random hyperplane tes-*
1006 *sellations*, Discrete & Computational Geometry, 51 (2014), pp. 438–461.
- 1007 [42] B. T. POLYAK, *Minimization of unsmooth functionals*, USSR Computational
1008 Mathematics and Mathematical Physics, 9 (1969), pp. 14–29.
- 1009 [43] S. PORTNOY AND R. KOENKER, *The gaussian hare and the laplacian tortoise:*
1010 *computability of squared-error versus absolute-error estimators*, Statistical Sci-
1011 ence, 12 (1997), pp. 279–300.
- 1012 [44] B. RECHT, M. FAZEL, AND P. A. PARRILO, *Guaranteed minimum-rank solu-*
1013 *tions of linear matrix equations via nuclear norm minimization*, SIAM review,
1014 52 (2010), pp. 471–501.
- 1015 [45] J. RENEGAR AND B. GRIMMER, *A simple nearly optimal restart scheme for*
1016 *speeding up first-order methods*, Foundations of Computational Mathematics, 22
1017 (2022), pp. 211–256.
- 1018 [46] R. T. ROCKAFELLAR, *Convex analysis*, vol. 11, Princeton university press, 1997.
- 1019 [47] V. ROULET, N. BOUMAL, AND A. D'ASPREMONT, *Renegar's condition number,*
1020 *sharpness and compressed sensing performance*, arXiv preprint arXiv:1506.03295,
1021 (2015).
- 1022 [48] V. ROULET AND A. D'ASPREMONT, *Sharpness, restart, and acceleration*, SIAM
1023 Journal on Optimization, 30 (2020), pp. 262–289.
- 1024 [49] M. RUDELSON AND R. VERSHYNIN, *Sparse reconstruction by convex relaxation:*
1025 *Fourier and gaussian measurements*, in 2006 40th Annual Conference on Infor-
1026 mation Sciences and Systems, IEEE, 2006, pp. 207–212.
- 1027 [50] F. SHANG, J. CHENG, Y. LIU, Z.-Q. LUO, AND Z. LIN, *Bilinear factor matrix*
1028 *norm minimization for robust pca: Algorithms and applications*, IEEE transac-
1029 tions on pattern analysis and machine intelligence, 40 (2017), pp. 2066–2080.
- 1030 [51] J. F. STURM, *Error bounds for linear matrix inequalities*, SIAM Journal on
1031 Optimization, 10 (2000), pp. 1228–1248.
- 1032 [52] R. VERSHYNIN, *Introduction to the non-asymptotic analysis of random matrices*,
1033 Cambridge University Press, 2012, p. 210–268.
- 1034 [53] T. YANG AND Q. LIN, *Rsg: Beating subgradient method without smoothness and*
1035 *strong convexity*, The Journal of Machine Learning Research, 19 (2018), pp. 236–
1036 268.
- 1037 [54] S. ZHANG, *Global error bounds for convex conic problems*, SIAM Journal on
1038 Optimization, 10 (2000), pp. 836–851.
- 1039 [55] Z. ZHOU AND A. M.-C. SO, *A unified approach to error bounds for structured*
1040 *convex optimization problems*, Mathematical Programming, 165 (2017), pp. 689–
1041 728.