

NONCONVEX ROBUST LOW-RANK MATRIX RECOVERY*

XIAO LI[†], ZHIHUI ZHU[‡], ANTHONY MAN-CHO SO[§], AND RENÉ VIDAL[‡]

Abstract. In this paper we study the problem of recovering a low-rank matrix from a number of random linear measurements that are corrupted by outliers taking arbitrary values. We consider a nonsmooth nonconvex formulation of the problem, in which we explicitly enforce the low-rank property of the solution by using a factored representation of the matrix variable and employ an ℓ_1 -loss function to robustify the solution against outliers. We show that even when a constant fraction (which can be up to almost half) of the measurements are arbitrarily corrupted, as long as certain measurement operators arising from the measurement model satisfy the so-called ℓ_1/ℓ_2 -restricted isometry property, the ground-truth matrix can be exactly recovered from any global minimum of the resulting optimization problem. Furthermore, we show that the objective function of the optimization problem is sharp and weakly convex. Consequently, a subgradient Method (SubGM) with geometrically diminishing step sizes will converge linearly to the ground-truth matrix when suitably initialized. We demonstrate the efficacy of the SubGM for the nonconvex robust low-rank matrix recovery problem with various numerical experiments.

Key words. robust low-rank matrix recovery, sharpness, weak convexity, subgradient method, robust PCA

AMS subject classifications. 65K10, 90C26, 68Q25, 68W40, 62B10.

1. Introduction. Low-rank matrices are ubiquitous in computer vision [8, 23], machine learning [40], and signal processing [13] applications. One fundamental computational task is to recover a low-rank matrix $\mathbf{X}^* \in \mathbb{R}^{n_1 \times n_2}$ from a small number of linear measurements

$$(1.1) \quad \mathbf{y} = \mathcal{A}(\mathbf{X}^*),$$

where $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$ is a known linear operator. Such a task arises in quantum tomography [1], face recognition [8], linear system identification [18], collaborative filtering [10], etc. We refer the interested reader to [13, 53] for more detailed discussions.

Although in many interesting scenarios the number of linear measurements m is much smaller than $n_1 n_2$, the low-rank property of \mathbf{X}^* suggests that its degrees of freedom can also be much smaller than $n_1 n_2$, thus making the task of recovering \mathbf{X}^* possible. This has been demonstrated in, e.g., [10], where a nuclear norm minimization approach for recovering a low-rank matrix from random linear measurements is studied. Despite the strong theoretical guarantees of such approach (see also [21]), most existing methods for solving the nuclear norm minimization problem do not scale well with the problem size (i.e., n_1 , n_2 , and m). To overcome this computational bottleneck, one approach is to enforce the low-rank property explicitly by using a factored representation of the matrix variable in the optimization formulation. Such

*Submitted to the editors December 12, 2019. The first and second authors contributed equally to this paper.

Funding: X. Li was partially supported by the Hong Kong Research Grants Council (RGC) General Research Fund (GRF) Project CUHK 14210617. Z. Zhu and R. Vidal were partially supported by NSF Grant 1704458. A. M.-C. So was partially supported by the Hong Kong Research Grants Council (RGC) General Research Fund (GRF) Project CUHK 14208117.

[†]Department of Electronic Engineering, The Chinese University of Hong Kong. (xli@ee.cuhk.edu.hk, <http://www.ee.cuhk.edu.hk/~xli/>).

[‡]Center for Imaging Science, Mathematical Institute for Data Science, Johns Hopkins University. (zzhu29@jhu.edu, <http://cis.jhu.edu/~zhihui/>; rvidal@jhu.edu, <http://cis.jhu.edu/~rvidal/>).

[§]Department of Systems Engineering and Engineering Management, The Chinese University of Hong Kong. (manchoso@se.cuhk.edu.hk, <http://www.se.cuhk.edu.hk/~manchoso/>).

37 an approach has already been explored in some early works on low-rank semidefinite
 38 programming (see, e.g., [5, 6] and the references therein) but has gained renewed
 39 interest lately in the study of low-rank matrix recovery problems. For the purpose
 40 of illustration, let us first consider the case where the ground-truth matrix \mathbf{X}^* is
 41 symmetric positive semidefinite with rank r . Instead of optimizing, say, an ℓ_2 -loss
 42 function involving an $n \times n$ symmetric positive semidefinite matrix variable \mathbf{X} with
 43 either a constraint or a regularization term controlling the rank of \mathbf{X} , we consider the
 44 factorization $\mathbf{X} = \mathbf{U}\mathbf{U}^T$ and optimize the loss function over the $n \times r$ matrix variable
 45 \mathbf{U} :

$$46 \quad (1.2) \quad \underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\text{minimize}} \left\{ \xi(\mathbf{U}) := \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^T)\|_2^2 \right\}.$$

48 There are two obvious advantages with the formulation (1.2). First, the recovered
 49 matrix will automatically satisfy the rank and positive semidefinite constraints. Se-
 50 cond, when the rank of the ground-truth matrix is small, the size of the variable \mathbf{U}
 51 can be much smaller than that of \mathbf{X} . Although the quadratic nature of $\mathbf{U}\mathbf{U}^T$ renders
 52 the objective function ξ in (1.2) nonconvex, recent advances in the analysis of the
 53 landscapes of structured nonconvex functions allow one to show that when the linear
 54 measurement operator \mathcal{A} satisfies certain restricted isometry property (RIP), local
 55 search algorithms (such as gradient descent) are guaranteed to find a global mini-
 56 mum of (1.2) and exactly recover the underlying low-rank matrix \mathbf{X}^* [4, 19, 35, 41, 52].
 57 Moreover, it was shown in [42, 50] that (1.2) satisfies an error bound condition, indi-
 58 cating that simple gradient descent with an appropriate initialization will converge to
 59 a global minimum at a linear rate; see [12] for a comprehensive review.

60 **1.1. Our Goal and Main Results.** In this paper, we consider the *robust low-*
 61 *rank matrix recovery problem*, in which the measurements are corrupted by *outliers*.
 62 Specifically, we assume that

$$63 \quad (1.3) \quad \mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \mathbf{s}^*,$$

64 where $\mathbf{s}^* \in \mathbb{R}^m$ is an outlier vector such that a small fraction of its entries (the
 65 outliers) have an arbitrary magnitude and the remaining entries are zero. Moreover,
 66 the set of nonzero entries is assumed to be unknown. Outliers are prevalent in the
 67 context of sensor calibration [31] (because of sensor failure), face recognition [16] (due
 68 to self-shadowing, specularities, or saturations in brightness), video surveillance [26]
 69 (where the foreground objects are modeled as outliers), etc.

70 It is well known that the ℓ_2 -loss function is sensitive to outliers, thus rende-
 71 ring (1.2) ineffective for recovering the underlying low-rank matrix. As illustrated in
 72 the top row of Figure 1, the global minima of ξ in (1.2) are perturbed away from the
 73 underlying low-rank matrix because of the outliers, and a larger fraction of outliers
 74 leads to a larger perturbation. By contrast, the ℓ_1 -loss function is more robust against
 75 outliers and has been widely utilized for outlier detection [8, 24, 31]. This motivates us
 76 to adopt the ℓ_1 -loss function together with the factored representation of the matrix
 77 variable to tackle the robust low-rank matrix recovery problem:

$$78 \quad (1.4) \quad \underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\text{minimize}} \left\{ f(\mathbf{U}) := \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^T)\|_1 \right\}.$$

79 The robustness of the ℓ_1 -loss function against outliers can be seen from the bottom row
 80 of Figure 1, where the global minima of (1.4) correspond precisely to the underlying

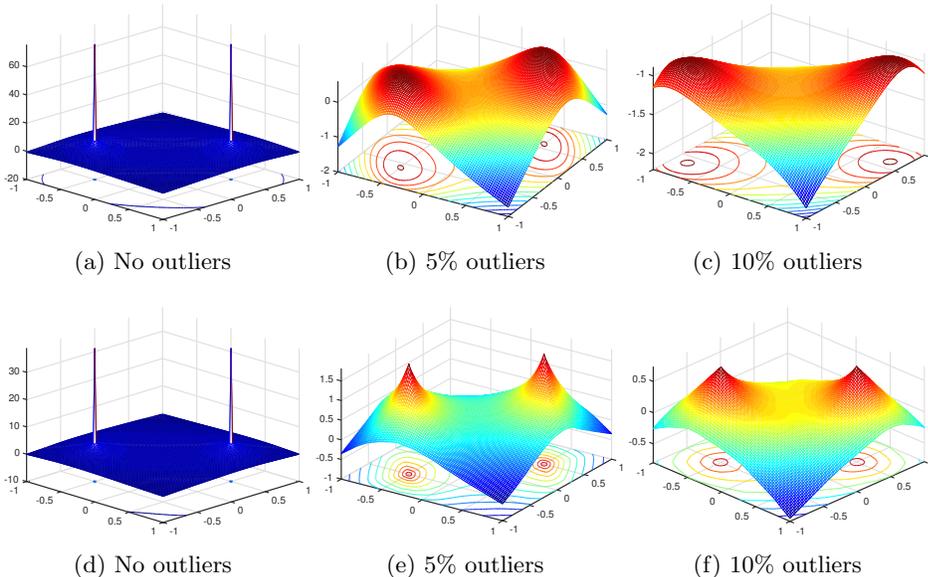


Fig. 1: Landscapes of the objective functions $\mathbf{U} \mapsto \xi(\mathbf{U}) = \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^T)\|_2^2$ (top row) and $\mathbf{U} \mapsto f(\mathbf{U}) = \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^T)\|_1$ (bottom row) for low-rank matrix recovery with different percentages of outliers in the measurement vector \mathbf{y} (1.3). Here, the ground-truth matrix \mathbf{X}^* is given by $\mathbf{X}^* = \mathbf{U}^* \mathbf{U}^{*\top}$ with $\mathbf{U}^* = [0.5 \ 0.5]^\top$ and 40 measurements are taken to form \mathbf{y} . For display purpose, we plot $-\log(\xi(\mathbf{U}))$ and $-\log(f(\mathbf{U}))$ instead of $\xi(\mathbf{U})$ and $f(\mathbf{U})$.

81 low-rank matrix \mathbf{X}^* even in the presence of outliers. However, compared with (1.2),
 82 the exact recovery property of (1.4) (i.e., when the global minima of (1.4) yield the
 83 ground-truth matrix \mathbf{X}^*) and the convergence behavior of local search algorithms for
 84 solving (1.4) are much less understood. This stems in part from the fact that (1.4) is a
 85 nonsmooth nonconvex optimization problem, but most of the algorithmic and analysis
 86 techniques developed in the recent literature on structured nonconvex optimization
 87 problems apply only to the smooth setting.

88 In view of the above discussion, we aim to (i) provide conditions in terms of the
 89 number of linear measurements m and the fraction of outliers that can guarantee the
 90 exact recovery property of (1.4) and (ii) design a first-order method to solve (1.4) and
 91 establish guarantees on its convergence performance. To achieve (i), we utilize the
 92 notion of ℓ_1/ℓ_2 -restricted isometry property (ℓ_1/ℓ_2 -RIP), which has been introduced
 93 previously in the context of low-rank matrix recovery [46, 48] and covariance estima-
 94 tion [11]. We show that if the fraction of outliers is slightly less than $\frac{1}{2}$, then as long
 95 as the measurement operator \mathcal{A} and its restriction \mathcal{A}_{Ω^c} onto the complement of the
 96 support set Ω of the outlier vector \mathbf{s}^* possess the ℓ_1/ℓ_2 -RIP, any global minimum \mathbf{U}^*
 97 of (1.4) must satisfy $\mathbf{U}^* \mathbf{U}^{*\top} = \mathbf{X}^*$. To tackle (ii), we propose to use a subgradient
 98 method (SubGM) to solve (1.4). As a key step in the convergence analysis of the
 99 SubGM, we show that under the aforementioned setting for the fraction of outliers
 100 and the ℓ_1/ℓ_2 -RIP of the operators \mathcal{A} and \mathcal{A}_{Ω^c} , the objective function f in (1.4) is
 101 *sharp* (see Definition 1) and *weakly convex* (see Definition 2). Consequently, we can
 102 apply (a slight variant of) the analysis framework in [14] to show that when initialized
 103 close to the set of global minima of (1.4), the SubGM with geometrically diminishing
 104 step sizes will converge R -linearly to a global minimum. To the best of our knowledge,

105 this is the first time an exact recovery condition (i.e., the ℓ_1/ℓ_2 -RIP of \mathcal{A} and \mathcal{A}_{Ω^c}) for
 106 the optimization formulation (1.4) is shown to also imply its regularity (i.e., sharpness
 107 and weak convexity). We summarize the above results in the following theorem:

108 **THEOREM 1** (informal; see [Theorem 3](#) for the formal statement). *Consider the*
 109 *measurement model (1.3), where the ground-truth matrix \mathbf{X}^* is symmetric positive*
 110 *semidefinite with rank r . Suppose that the fraction of outliers is less than half and both*
 111 *operators \mathcal{A} and \mathcal{A}_{Ω^c} possess the ℓ_1/ℓ_2 -RIP (see [Subsection 3.1](#) and [Subsection 3.2](#)).*
 112 *Then, every global minimum of (1.4) corresponds to the ground-truth matrix \mathbf{X}^**
 113 *and the objective function f is sharp (see [Definition 1](#)) and weakly convex (see [De-](#)*
 114 *inition 2). Consequently, when applied to (1.4), the SubGM with an appropriate*
 115 *initialization will converge to the ground-truth matrix \mathbf{X}^* at a linear rate.*

116 Before we proceed, several remarks are in order. First, for various random measu-
 117 rement operators \mathcal{A} , such as sub-Gaussian measurement operators and the quadratic
 118 measurement operators in [11], as long as the number of measurements is sufficiently
 119 large, the operators \mathcal{A} and \mathcal{A}_{Ω^c} will possess the ℓ_1/ℓ_2 -RIP with high probability. This
 120 is the case, for instance, when \mathcal{A} is a Gaussian measurement operator with $m \gtrsim nr$
 121 measurements.¹ In particular, when combined with [Theorem 1](#), we see that the low-
 122 rank matrix \mathbf{X}^* in (1.3) can be recovered using an information-theoretically optimal
 123 number of measurements. Second, although at first glance (1.4) seems to be more
 124 difficult to solve than (1.2) because of nonsmoothness, [Theorem 1](#) implies that (1.4)
 125 can be solved *as efficiently as* its smooth counterpart (1.2), in the sense that both
 126 can be solved by first-order methods that have a linear convergence guarantee.

127 Although [Theorem 1](#) is concerned with the setting where \mathbf{X}^* is symmetric positive
 128 semidefinite, it can be extended to the general setting where \mathbf{X}^* is a rank- r $n_1 \times n_2$ ma-
 129 trix. Specifically, by using the factorization $\mathbf{X} = \mathbf{U}\mathbf{V}^T$ with $\mathbf{U} \in \mathbb{R}^{n_1 \times r}$, $\mathbf{V} \in \mathbb{R}^{n_2 \times r}$
 130 and utilizing the nonsmooth regularizer $\|\mathbf{U}^T\mathbf{U} - \mathbf{V}^T\mathbf{V}\|_F$ (or $\|\mathbf{U}^T\mathbf{U} - \mathbf{V}^T\mathbf{V}\|_1$) to
 131 account for the ambiguities in the factorization caused by invertible transformations,
 132 we formulate the general robust low-rank matrix recovery problem as follows:

$$133 \quad (1.5) \quad \underset{\mathbf{U} \in \mathbb{R}^{n_1 \times r}, \mathbf{V} \in \mathbb{R}^{n_2 \times r}}{\text{minimize}} \quad \left\{ g(\mathbf{U}, \mathbf{V}) := \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^T)\|_1 + \lambda \|\mathbf{U}^T\mathbf{U} - \mathbf{V}^T\mathbf{V}\|_F \right\}.$$

134 Here, $\lambda > 0$ is a regularization parameter. We remark that the regularizer used in
 135 the above formulation is motivated by but different from that used in [35, 42, 52]. The
 136 latter, which is given by $\|\mathbf{U}^T\mathbf{U} - \mathbf{V}^T\mathbf{V}\|_F^2$, is smooth but is not as well suited for
 137 robustifying the solution against outliers. In [Section 4](#) we show that all the results
 138 established for (1.4) in [Theorem 1](#) carry over to (1.5) for any $\lambda > 0$ (but the choice
 139 of λ affects the sharpness and weak convexity parameters; see the discussion after
 140 [Proposition 6](#)).

141 **1.2. Related Work.** By analyzing the optimization geometry, recent works [4,
 142 19, 28, 35, 42] have shown that many local search algorithms with either an appropri-
 143 ate initialization or a random initialization can provably solve the low-rank matrix
 144 recovery problem (1.2) when the measurement operator \mathcal{A} satisfies the RIP. In par-
 145 ticular, gradient descent with an appropriate initialization is shown to converge to
 146 a global optimum at a linear rate [42, 51], while quadratic convergence is establis-
 147 hed for the cubic regularization method [47]. Key to these results is certain error
 148 bound conditions, which elucidate the regularity properties of the underlying opti-
 149 mization problem. Recently, the above results have been extended to cover general

¹See [Subsection 1.3](#) for the meaning of the notation \gtrsim .

150 smooth low-rank matrix optimization problems whose objective functions satisfy the
 151 restricted strong convexity and smoothness properties [27, 51, 52].

152 For the robust low-rank matrix recovery problem, existing solution methods can
 153 be classified into two categories. The first is based on the convex approach [8, 25,
 154 31]. Although such approach enjoys strong statistical guarantees, it is computational
 155 expensive and thus not scalable to practical problems. The second category is based
 156 on the nonconvex approach. This includes the alternating minimization methods
 157 [22, 33, 45, 49], which typically use projected gradient descent for low-rank matrix
 158 recovery and thresholding-based truncation for identification of outliers. However,
 159 these methods typically require performing an SVD in each iteration for projection
 160 onto the set of low-rank matrices. Recently, a median-truncated gradient descent
 161 method has been proposed in [30] to tackle (1.2), where the gradient is modified to
 162 alleviate the effect of outliers. The median-truncated gradient descent is shown to
 163 have a local linear convergence rate [30], but such guarantee requires $m \gtrsim nr \log n$
 164 measurements. Moreover, the maximum number of outliers that can be tolerated is
 165 not explicitly given. By contrast, our result only requires $m \gtrsim nr$ measurements
 166 (which matches the optimal information-theoretic bound) and explicitly bounds the
 167 fraction of outliers that can be present. We also note that a SubGM has been proposed
 168 in [31] for solving (1.4) in the setting where \mathcal{A} is a certain quadratic measurement
 169 operator. As reported in [31], the SubGM exhibits excellent empirical performance
 170 in terms of both computational efficiency and accuracy. In this paper, we provide
 171 a rigorous justification for the empirical success of the SubGM, thus answering a
 172 question that is left open in [31].

173 Finally, we remark that our work is closely related to the recent works [2, 14, 15, 54]
 174 on subgradient methods for nonsmooth nonconvex optimization. A projected subgra-
 175 dient method is proven to converge linearly for the robust subspace recovery pro-
 176 blem [54] and sublinearly for orthonormal dictionary learning [2]. It is shown in [14, 15]
 177 that if the optimization problem at hand is sharp (see Definition 1) and weakly con-
 178 vex (see Definition 2), various subgradient methods for solving it will converge at a
 179 linear rate. Currently, only a few applications are known to give rise to sharp and
 180 weakly convex optimization problems, such as robust phase retrieval [15, 17] and ro-
 181 bust covariance estimation with quadratic sampling [14]. Thus, our result expands
 182 the repertoire of optimization problems that are sharp and weakly convex and contri-
 183 butes to the growing literature on the geometry of structured nonsmooth nonconvex
 184 optimization problems.

185 **1.3. Notation.** Let us introduce the notations used in this paper. Finite-
 186 dimensional vectors and matrices are indicated by bold characters. The symbols \mathbf{I} and
 187 $\mathbf{0}$ represent the identity matrix and zero matrix/vector, respectively. The set of $r \times r$
 188 orthogonal matrices is denoted by $\mathcal{O}_r := \{\mathbf{R} \in \mathbb{R}^{r \times r} : \mathbf{R}^T \mathbf{R} = \mathbf{I}\}$. The subdifferential
 189 of the absolute value function $|\cdot|$ is denoted by Sign ; i.e., $\text{Sign}(a) := \begin{cases} a/|a|, & a \neq 0, \\ [-1, 1], & a = 0. \end{cases}$
 190 We use $\text{Sign}(\mathbf{A})$ to denote the matrix obtained by applying the Sign function to each
 191 element of the matrix \mathbf{A} . Furthermore, we use $\|\mathbf{A}\|_F$ to denote the Frobenius norm
 192 of the matrix \mathbf{A} and $\|\mathbf{a}\|$ to denote the ℓ_2 -norm of the vector \mathbf{a} . Finally, we use $x \lesssim y$
 193 (resp. $x \gtrsim y$) to indicate that $x \leq cy$ (resp. $x \geq cy$) for some universal constant $c > 0$.

194 **2. Problem Setup and Preliminaries.** Consider the general optimization
 195 problem

$$196 \quad (2.1) \quad \inf_{\mathbf{x} \in \mathbb{R}^n} h(\mathbf{x}),$$

197 where $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is a lower semi-continuous, possibly nonsmooth and nonconvex,
198 function. Let h^* denote the optimal value of (2.1) and

$$199 \quad \mathcal{X} := \{\mathbf{z} \in \mathbb{R}^n : h(\mathbf{z}) \leq h(\mathbf{x}), \forall \mathbf{x} \in \mathbb{R}^n\}$$

200 denote the set of global minima of h . We assume that $\mathcal{X} \neq \emptyset$. Given any $\mathbf{x} \in \mathbb{R}^n$, the
201 distance between \mathbf{x} and \mathcal{X} is defined as

$$202 \quad \text{dist}(\mathbf{x}, \mathcal{X}) := \inf_{\mathbf{z} \in \mathcal{X}} \|\mathbf{x} - \mathbf{z}\|.$$

203 Since h can be nonsmooth, we utilize tools from generalized differentiation to for-
204 mulate the optimality condition of (2.1). The (Fréchet) subdifferential of h at \mathbf{x}
205 is defined as

$$206 \quad (2.2) \quad \partial h(\mathbf{x}) := \left\{ \mathbf{d} \in \mathbb{R}^n : \liminf_{\mathbf{y} \rightarrow \mathbf{x}} \frac{h(\mathbf{y}) - h(\mathbf{x}) - \langle \mathbf{d}, \mathbf{y} - \mathbf{x} \rangle}{\|\mathbf{y} - \mathbf{x}\|} \geq 0 \right\},$$

207 where each $\mathbf{d} \in \partial h(\mathbf{x})$ is called a subgradient of h at \mathbf{x} . We say that \mathbf{x} is a critical
208 point of h if $\mathbf{0} \in \partial h(\mathbf{x})$.

209 **2.1. Sharpness and Weak Convexity.** Since our goal is to consider a set of
210 problems that can be solved by the SubGM with a linear rate of convergence, let us
211 introduce two regularity notions for h that are central to our study.

212 **DEFINITION 1** (sharpness; cf. [7]). *We say that $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is sharp with para-*
213 *meter $\alpha > 0$ if*

$$214 \quad h(\mathbf{x}) - h^* \geq \alpha \text{dist}(\mathbf{x}, \mathcal{X})$$

215 *for all $\mathbf{x} \in \mathbb{R}^n$.*

216 **DEFINITION 2** (weak convexity; see, e.g., [44]). *We say that $h : \mathbb{R}^n \rightarrow \mathbb{R}$ is weakly*
217 *convex with parameter $\tau \geq 0$ if $\mathbf{x} \mapsto h(\mathbf{x}) + \frac{\tau}{2}\|\mathbf{x}\|^2$ is convex.*

218 Suppose that h is sharp and weakly convex with parameters $\alpha > 0$ and $\tau \geq 0$,
219 respectively. It is known that for any $\mathbf{x} \notin \mathcal{X}$ with $\text{dist}(\mathbf{x}, \mathcal{X}) < \frac{2\alpha}{\tau}$, we have $\mathbf{0} \notin \partial h(\mathbf{x})$;
220 i.e., \mathbf{x} is not a critical point of h [14, Lemma 3.1]. This suggests the possibility of
221 finding a global minimum of h by initializing local search algorithms with a point
222 that is close to \mathcal{X} . To explore such possibility, let us consider using the SubGM in
223 [Algorithm 2.1](#) to solve the nonsmooth nonconvex optimization problem (2.1).

Algorithm 2.1 Subgradient Method (SubGM) for Solving (2.1)

Initialization: set \mathbf{x}_0 and μ_0 ;

- 1: **for** $k = 0, 1, \dots$ **do**
 - 2: compute a subgradient $\mathbf{d}_k \in \partial h(\mathbf{x}_k)$;
 - 3: update the step size μ_k according to a certain rule;
 - 4: update $\mathbf{x}_{k+1} = \mathbf{x}_k - \mu_k \mathbf{d}_k$;
 - 5: **end for**
-

224 **2.2. Convergence of SubGM for Sharp Weakly Convex Functions.** Un-
225 like gradient descent, the SubGM with a constant step size may not converge to
226 a critical point of a nonsmooth function in general, even when the function is con-
227 vex [3, 32, 38]. To ensure the convergence of the SubGM, a set of diminishing step sizes
228 is generally needed [20, 38]. As it turns out, for a sharp weakly convex function h ,
229 the SubGM with step sizes that are diminishing at a geometric rate can still be shown
230 to converge linearly to a global minimum when initialized close to \mathcal{X} . Specifically, let

$$231 \quad (2.3) \quad \kappa := \sup \left\{ \|\mathbf{d}\| : \mathbf{d} \in \partial h(\mathbf{x}), \text{dist}(\mathbf{x}, \mathcal{X}) < \frac{2\alpha}{\tau} \right\},$$

232 which can be shown to satisfy $\kappa \geq \alpha$; cf. [14, Lemma 3.2]. Then, we have the following
233 result:

234 **THEOREM 2** (local linear convergence of SubGM). *Suppose that the function*
235 *$h : \mathbb{R}^n \rightarrow \mathbb{R}$ is sharp and weakly convex with parameters $\alpha > 0$ and $\tau \geq 0$, re-*
236 *spectively. Suppose further that the SubGM in Algorithm 2.1 is initialized with a*
237 *point \mathbf{x}_0 satisfying $\text{dist}(\mathbf{x}_0, \mathcal{X}) < \frac{2\alpha}{\tau}$ and uses the geometrically diminishing step sizes*
238 *$\mu_k = \rho^k \mu_0$, where the initial step size μ_0 satisfies*

$$239 \quad (2.4) \quad \mu_0 \leq \frac{\alpha^2}{2\tau\kappa^2} \left(1 - \left(\max \left\{ \frac{\tau}{\alpha} \text{dist}(\mathbf{x}_0, \mathcal{X}) - 1, 0 \right\} \right)^2 \right)$$

240 and the decay rate ρ satisfies

$$241 \quad (2.5) \quad 1 > \rho \geq \underline{\rho} := \sqrt{1 - \left(\frac{2\alpha}{\overline{\text{dist}}_0} - \tau \right) \mu_0 + \frac{\kappa^2}{\overline{\text{dist}}_0^2} \mu_0^2}$$

242 with

$$243 \quad (2.6) \quad \overline{\text{dist}}_0 = \max \left\{ \text{dist}(\mathbf{x}_0, \mathcal{X}), \mu_0 \frac{\max\{\kappa^2, 2\alpha^2\}}{\alpha} \right\}.$$

244 Then, the iterates $\{\mathbf{x}_k\}_{k \geq 0}$ generated by the SubGM will converge linearly to a point
245 in \mathcal{X} :

$$246 \quad \text{dist}(\mathbf{x}_k, \mathcal{X}) \leq \rho^k \overline{\text{dist}}_0, \quad \forall k \geq 0.$$

247 We note that a similar result has been established in [14, Corollary 6.1]. Neverthe-
248 less, compared with [14, Corollary 6.1], which requires $\frac{\alpha}{\kappa} \leq \sqrt{\frac{1}{2-\gamma}}$ and $\text{dist}(\mathbf{x}_0, \mathcal{X}) \leq$
249 $\frac{\gamma\alpha}{\tau}$ for some $\gamma \in (0, 1)$, Theorem 2 is less restrictive and allows the larger initialization
250 region $\text{dist}(\mathbf{x}_0, \mathcal{X}) < \frac{2\alpha}{\tau}$. In particular, as $\frac{\alpha}{\kappa}$ tends to 1, so does γ , and the decay rate
251 ρ in [14, Corollary 6.1] approaches 1. Thus, one can no longer use [14, Corollary 6.1]
252 to conclude that the SubGM converges linearly when $\frac{\alpha}{\kappa} = 1$. By contrast, the linear
253 convergence result in Theorem 2 is still valid in this case. Theorem 2 can be proven
254 by refining the arguments in the proof of [14, Theorem 6.1]. We refer the reader to
255 the companion technical report [29] of this paper for details.

256 Before we proceed, it is worth elaborating on the implication of Theorem 2 when
257 h is convex. In this case, we can take $\tau = 0$, which, in view of (2.4), shows that
258 μ_0 can be arbitrarily chosen. If we choose $\mu_0 \geq \frac{\alpha \text{dist}(\mathbf{x}_0, \mathcal{X})}{\max\{\kappa^2, 2\alpha^2\}}$, then by (2.6) we have
259 $\overline{\text{dist}}_0 = \mu_0 \frac{\max\{\kappa^2, 2\alpha^2\}}{\alpha}$, which implies that the decay rate $\underline{\rho}$ satisfies

$$260 \quad \underline{\rho} = \sqrt{1 - \frac{2\alpha^2}{\max\{\kappa^2, 2\alpha^2\}} + \frac{\kappa^2\alpha^2}{(\max\{\kappa^2, 2\alpha^2\})^2}} = \begin{cases} \sqrt{1 - \frac{\alpha^2}{\kappa^2}}, & \kappa^2 \geq 2\alpha^2, \\ \frac{\kappa}{2\alpha}, & \kappa^2 < 2\alpha^2. \end{cases}$$

261 In particular, this is in line with the results in [20, Theorem 4.4].

262 **3. Nonconvex Robust Low-Rank Matrix Recovery: Symmetric Posi-**
263 **tive Semidefinite (PSD) Case.** In the last section we saw that the SubGM with
264 suitable initialization and step sizes converges linearly to a global minimum of a sharp
265 weakly convex function. Naturally, it is of interest to identify concrete problems that
266 possess these two regularity properties. In this section we focus on the robust low-rank
267 matrix recovery problem (1.4) and establish, for the first time, a connection between
268 the exact recovery condition of ℓ_1/ℓ_2 -RIP and the regularity properties of sharpness
269 and weak convexity of the objective function f in (1.4). Specifically, we first show that

270 if the fraction of outliers is slightly less than $\frac{1}{2}$ and certain measurement operators
 271 arising from the measurement model (1.3) possess the ℓ_1/ℓ_2 -RIP, then the sharpness
 272 condition in Definition 1 holds for (1.4). Consequently, all global minima of (1.4)
 273 lead to the exact recovery of the ground-truth matrix \mathbf{X}^* . We then show that (1.4)
 274 also satisfies the weak convexity condition in Definition 2. Hence, by the convergence
 275 result (Theorem 2) in the last section, we conclude that the SubGM can be utilized
 276 to find a global minimum of (1.4) efficiently.

277 To begin, let us collect some preparatory results. Let $\mathbf{X}^* = \mathbf{U}^* \mathbf{U}^{*\top}$ be a
 278 factorization of \mathbf{X}^* , where $\mathbf{U}^* \in \mathbb{R}^{n \times r}$. Note that for any $\mathbf{R} \in \mathcal{O}_r$, we have
 279 $\mathbf{X}^* = \mathbf{U}^* \mathbf{R} (\mathbf{U}^* \mathbf{R})^\top$. Thus, all elements in the set

$$280 \quad \mathcal{U} := \{\mathbf{U}^* \mathbf{R} : \mathbf{R} \in \mathcal{O}_r\}$$

281 are valid factors of \mathbf{X}^* . Furthermore, it is clear that the function f in (1.4) is constant
 282 on the set \mathcal{U} . The following result connects $\text{dist}(\mathbf{U}, \mathcal{U})$ and the distance between $\mathbf{U} \mathbf{U}^\top$
 283 and $\mathbf{U}^* \mathbf{U}^{*\top}$ for any given $\mathbf{U} \in \mathbb{R}^{n \times r}$:

284 LEMMA 1 ([42, Lemma 5.4]). *Given any $\mathbf{U}^* \in \mathbb{R}^{n \times r}$, define $\mathbf{X}^* = \mathbf{U}^* \mathbf{U}^{*\top}$.
 285 Then, for any $\mathbf{U} \in \mathbb{R}^{n \times r}$, we have*

$$286 \quad 2 \left(\sqrt{2} - 1 \right) \sigma_r^2(\mathbf{X}^*) \text{dist}^2(\mathbf{U}, \mathcal{U}) \leq \|\mathbf{U} \mathbf{U}^\top - \mathbf{U}^* \mathbf{U}^{*\top}\|_F^2,$$

287 where σ_r denotes the r -th largest singular value.

288 **3.1. ℓ_1/ℓ_2 -Restricted Isometry Property.** Since the ℓ_1/ℓ_2 -RIP [11, 46, 48]
 289 of the linear measurement operator $\mathcal{A} : \mathbb{R}^{n \times n} \rightarrow \mathbb{R}^m$ in (1.4) plays an impor-
 290 tant role in our subsequent analysis, let us first provide a condition under which
 291 \mathcal{A} will possess such property. Recall that \mathcal{A} can be specified by a collection of
 292 m $n \times n$ matrices $\mathbf{A}_1, \dots, \mathbf{A}_m$. In other words, given any $\mathbf{X} \in \mathbb{R}^{n \times n}$, we have
 293 $\mathcal{A}(\mathbf{X}) = (\langle \mathbf{A}_1, \mathbf{X} \rangle, \dots, \langle \mathbf{A}_m, \mathbf{X} \rangle)$. We now show that if $\mathbf{A}_1, \dots, \mathbf{A}_m$ have indepen-
 294 dent and identically distributed (*i.i.d.*) standard Gaussian entries, then \mathcal{A} will possess
 295 the ℓ_1/ℓ_2 -RIP with high probability.

296 PROPOSITION 1 (ℓ_1/ℓ_2 -RIP of Gaussian measurement operators). *Let $r \geq 1$ be
 297 given. Suppose that $m \gtrsim nr$ and the matrices $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n \times n}$ defining the
 298 linear measurement operator \mathcal{A} have *i.i.d.* standard Gaussian entries. Then, for
 299 any $0 < \delta < \sqrt{\frac{2}{\pi}}$, there exists a universal constant $c > 0$ such that with probability
 300 exceeding $1 - \exp(-c\delta^2 m)$, \mathcal{A} will possess the ℓ_1/ℓ_2 -RIP; *i.e.*, the inequalities*

$$301 \quad (3.1) \quad \left(\sqrt{\frac{2}{\pi}} - \delta \right) \|\mathbf{X}\|_F \leq \frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 \leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{X}\|_F$$

302 hold for any rank- $2r$ matrix $\mathbf{X} \in \mathbb{R}^{n \times n}$.

303 The proof of Proposition 1 is given in Appendix A. It is worth noting that simi-
 304 lar ℓ_1/ℓ_2 -RIPs hold for other types of measurement operators such as the quadratic
 305 measurement operators in [11] and those defined by sub-Gaussian matrices. Thus,
 306 although our results are stated for Gaussian measurement operators, they can be
 307 readily extended to cover other measurement operators that possess similar RIPs.

308 **3.2. Sharpness and Exact Recovery.** Assuming that the linear measurement
 309 operator \mathcal{A} possesses the ℓ_1/ℓ_2 -RIP (3.1), our first goal is to identify further conditions
 310 on the measurement model (1.3) so that any global minimum \mathbf{U}^* of (1.4) can be used

311 to recover the ground-truth matrix \mathbf{X}^* via $\mathbf{U}^*\mathbf{U}^{*\top} = \mathbf{X}^*$. Towards that end, let
 312 $\Omega \subseteq \{1, \dots, m\}$ denote the support of the outlier vector \mathbf{s}^* and $\Omega^c = \{1, \dots, m\} \setminus \Omega$.
 313 Furthermore, let $p = \frac{|\Omega|}{m}$ be the fraction of outliers in \mathbf{y} . Throughout, we do not make
 314 any assumption on the location of the non-zero entries of \mathbf{s}^* . Instead, we assume that
 315 \mathcal{A}_{Ω^c} , the linear operator defined by the matrices in $\{\mathbf{A}_i : i \in \Omega^c\}$, also possesses the
 316 ℓ_1/ℓ_2 -RIP; i.e., we have

$$317 \quad (3.2) \quad \left(\sqrt{\frac{2}{\pi}} - \delta \right) \|\mathbf{X}\|_F \leq \frac{1}{m(1-p)} \|\mathcal{A}(\mathbf{X})\|_{\Omega^c} \leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{X}\|_F$$

318 for any rank- $2r$ matrix \mathbf{X} . When each \mathbf{A}_i is generated with *i.i.d.* standard Gaussian
 319 entries, [Proposition 1](#) implies that \mathcal{A}_{Ω^c} will satisfy (3.2) with high probability as long
 320 as p is a constant. This follows from the fact that $|\Omega^c| = (1-p)m \gtrsim nr$ if $m \gtrsim nr$.

321 **PROPOSITION 2** (sharpness and exact recovery with outliers: PSD case). *Let $0 <$*
 322 *$\delta < \frac{1}{3}\sqrt{\frac{2}{\pi}}$ be given. Suppose that the fraction of outliers p satisfies*

$$323 \quad (3.3) \quad p < \frac{1}{2} - \frac{\delta}{\sqrt{2/\pi} - \delta},$$

324 *and that the linear operators \mathcal{A} and \mathcal{A}_{Ω^c} possess the ℓ_1/ℓ_2 -RIP (3.1) and (3.2), re-*
 325 *spectively. Then, the objective function f in (1.4) satisfies*

$$326 \quad f(\mathbf{U}) - f(\mathbf{U}^*) \geq \alpha \operatorname{dist}(\mathbf{U}, \mathcal{U})$$

327 *for any $\mathbf{U} \in \mathbb{R}^{n \times r}$, where*

$$328 \quad (3.4) \quad \alpha = \sqrt{2(\sqrt{2} - 1)} \left(2(1-p) \left(\sqrt{\frac{2}{\pi}} - \delta \right) - \left(\sqrt{\frac{2}{\pi}} + \delta \right) \right) \sigma_r(\mathbf{X}^*) > 0.$$

329 *In particular, the set \mathcal{U} is precisely the set of global minima of (1.4) and the objective*
 330 *function f is sharp with parameter $\alpha > 0$.*

331 *Proof of Proposition 2.* Using (1.3) and (1.4), we compute

$$\begin{aligned} f(\mathbf{U}) - f(\mathbf{U}^*) &= \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*) - \mathbf{s}^*\|_1 - \frac{1}{m} \|\mathbf{s}^*\|_1 \\ &= \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*)\|_{\Omega^c} + \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*)\|_{\Omega} - \mathbf{s}^*\|_1 - \frac{1}{m} \|\mathbf{s}^*\|_1 \\ &\geq \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*)\|_{\Omega^c} - \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*)\|_{\Omega} \\ 332 &= \frac{2}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{X}^*)\|_{\Omega^c} - \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^\top - \mathbf{U}^*\mathbf{U}^{*\top})\|_1 \\ &\geq \left(2(1-p) \left(\sqrt{\frac{2}{\pi}} - \delta \right) - \left(\sqrt{\frac{2}{\pi}} + \delta \right) \right) \|\mathbf{U}^*\mathbf{U}^{*\top} - \mathbf{U}\mathbf{U}^\top\|_F \\ &\geq \alpha \operatorname{dist}(\mathbf{U}, \mathcal{U}), \end{aligned}$$

333 where the second inequality follows from the ℓ_1/ℓ_2 -RIP of \mathcal{A} and \mathcal{A}_{Ω^c} and the last
 334 inequality follows from [Lemma 1](#). The characterization of the set of global minima
 335 of (1.4) follows immediately from the above inequality and the choice of p in (3.3). \square

336 One interesting consequence of [Proposition 2](#) is that for the robust low-rank ma-
 337 trix recovery problem [\(1.4\)](#), the sharpness condition (which characterizes the geometry
 338 of the optimization problem around the set of global minima) coincides with the exact
 339 recovery property (which is of statistical nature). Moreover, condition [\(3.3\)](#) suggests
 340 that the smaller δ is, the higher the outlier ratio p can be. On the other hand, given an
 341 outlier ratio p , condition [\(3.3\)](#) requires that $\delta < \sqrt{\frac{2}{\pi} - \frac{\sqrt{2/\pi}}{3/2-p}}$, which indirectly imposes
 342 a condition on the number of measurements m . Indeed, [Proposition 1](#) implies that in
 343 order for a Gaussian measurement operator \mathcal{A} to possess the ℓ_1/ℓ_2 -RIP with positive
 344 probability, we need $m \gtrsim nr / (\sqrt{\frac{2}{\pi}} - \frac{\sqrt{2/\pi}}{3/2-p})^2$ measurements. Putting it another way,
 345 the larger the number of measurements m is, the higher the outlier ratio p can be.
 346 We shall elaborate on this point with experiments in [Section 5](#).

347 **3.3. Weak Convexity.** In the last subsection we established the sharpness of
 348 [\(1.4\)](#) and showed that any of its global minimum will lead to the exact recovery
 349 of the ground-truth matrix \mathbf{X}^* , even when the fraction of outliers is up to almost
 350 $\frac{1}{2}$. In this subsection we further establish the weak convexity of [\(1.4\)](#), thus opening
 351 up the possibility of using the machinery developed in [Section 2](#) to obtain provable
 352 convergence guarantees for the SubGM when it is applied to solve [\(1.4\)](#). Towards
 353 that end, we note that the ℓ_1 -norm, being a convex function, is subdifferentially
 354 regular [[37](#), Example 7.27] (see [[37](#), Definition 7.25] for the definition of subdifferential
 355 regularity). Hence, by the chain rule for subdifferentials of subdifferentially regular
 356 functions [[37](#), Corollary 8.11 and Theorem 10.6], we have
 (3.5)

$$357 \quad \partial f(\mathbf{U}) = \frac{1}{m} \left[(\mathcal{A}^* (\text{Sign}(\mathcal{A}(\mathbf{U}\mathbf{U}^T) - \mathbf{y})))^T \mathbf{U} + \mathcal{A}^* (\text{Sign}(\mathcal{A}(\mathbf{U}\mathbf{U}^T) - \mathbf{y})) \mathbf{U} \right].$$

358 We are now ready to prove the following result. Note that the weak convexity para-
 359 meter τ in [\(3.6\)](#) is independent of the fraction of outliers.

360 **PROPOSITION 3** (weak convexity: PSD case). *Suppose that the measurement*
 361 *operator \mathcal{A} satisfies the ℓ_1/ℓ_2 -RIP [\(3.1\)](#). Then, the objective function f in [\(1.4\)](#) is*
 362 *weakly convex with parameter*

$$363 \quad (3.6) \quad \tau = 2 \left(\sqrt{\frac{2}{\pi}} + \delta \right).$$

364 *Proof of Proposition 3.* For any $\mathbf{U}', \mathbf{U} \in \mathbb{R}^{n \times r}$, let $\Delta = \mathbf{U}' - \mathbf{U}$. Then, we have

$$\begin{aligned} 365 \quad f(\mathbf{U}') &= \frac{1}{m} \|\mathcal{A}(\mathbf{U}'\mathbf{U}'^T - \mathbf{X}^*) - \mathbf{s}^*\|_1 \\ 366 \quad &= \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^T - \mathbf{X}^* + \mathbf{U}\Delta^T + \Delta\mathbf{U}^T + \Delta\Delta^T) - \mathbf{s}^*\|_1 \\ 367 \quad &\geq \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^T - \mathbf{X}^* + \mathbf{U}\Delta^T + \Delta\mathbf{U}^T) - \mathbf{s}^*\|_1 - \frac{1}{m} \|\mathcal{A}(\Delta\Delta^T)\|_1 \\ 368 \quad &\geq \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{U}^T - \mathbf{X}^* + \mathbf{U}\Delta^T + \Delta\mathbf{U}^T) - \mathbf{s}^*\|_1 - \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\Delta\Delta^T\|_F \\ 369 \quad &\geq f(\mathbf{U}) + \frac{1}{m} \langle \mathbf{d}, \mathcal{A}(\mathbf{U}\Delta^T + \Delta\mathbf{U}^T) \rangle - \frac{\tau}{2} \|\Delta\|_F^2 \end{aligned}$$

371 for any $\mathbf{d} \in \text{Sign}(\mathcal{A}(\mathbf{U}\mathbf{U}^T) - \mathbf{y})$, where the second inequality follows from the ℓ_1/ℓ_2 -
 372 RIP of \mathcal{A} and the last inequality is due to the convexity of the ℓ_1 -norm and $\|\Delta\Delta^T\|_F \leq$

373 $\|\Delta\|_F^2$. Substituting (3.5) into the above equation gives

$$374 \quad f(\mathbf{U}') \geq f(\mathbf{U}) + \langle \mathbf{D}, \mathbf{U}' - \mathbf{U} \rangle - \frac{\tau}{2} \|\mathbf{U}' - \mathbf{U}\|_F^2, \quad \forall \mathbf{D} \in \partial f(\mathbf{U}).$$

375 This completes the proof. \square

376 **3.4. Putting Everything Together.** With the results in Subsection 3.2 and
 377 Subsection 3.3 in place, in order to show that the SubGM enjoys the convergence
 378 guarantees in Theorem 2 when applied to the robust low-rank matrix recovery pro-
 379 blem (1.4), it remains to determine κ , the bound on the norm of any subgradient of
 380 f in a neighborhood of \mathbf{U} ; see (2.3). This is established by the following result:

381 PROPOSITION 4 (bound on subgradient norm: PSD case). *Suppose that the mea-
 382 surement operator \mathcal{A} satisfies the ℓ_1/ℓ_2 -RIP (3.1). Then, for any $\mathbf{U} \in \mathbb{R}^{n \times r}$ satisfying
 383 $\text{dist}(\mathbf{U}, \mathcal{U}) \leq \frac{2\alpha}{\tau}$, we have*

$$384 \quad (3.7) \quad \|\mathbf{D}\|_F \leq \kappa = 2 \left(\sqrt{\frac{2}{\pi}} + \delta \right) \left(\|\mathbf{U}^*\|_F + \frac{2\alpha}{\tau} \right), \quad \forall \mathbf{D} \in \partial f(\mathbf{U}).$$

385 *Proof of Proposition 4.* Recall from (2.2) that

$$386 \quad (3.8) \quad \liminf_{\mathbf{U}' \rightarrow \mathbf{U}} \frac{f(\mathbf{U}') - f(\mathbf{U}) - \langle \mathbf{D}, \mathbf{U}' - \mathbf{U} \rangle}{\|\mathbf{U}' - \mathbf{U}\|_F} \geq 0$$

387 for any $\mathbf{D} \in \partial f(\mathbf{U})$. Now, for any $\mathbf{U}' \in \mathbb{R}^{n \times r}$,

$$\begin{aligned} 388 \quad |f(\mathbf{U}') - f(\mathbf{U})| &= \frac{1}{m} \left| \|\mathbf{y} - \mathcal{A}(\mathbf{U}'\mathbf{U}'^T)\|_1 - \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{U}^T)\|_1 \right| \\ 389 &\leq \frac{1}{m} \|\mathcal{A}(\mathbf{U}'\mathbf{U}'^T - \mathbf{U}\mathbf{U}^T)\|_1 \\ 390 &\leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{U}'\mathbf{U}'^T - \mathbf{U}\mathbf{U}^T\|_F \\ 391 &= \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|(\mathbf{U}' - \mathbf{U})\mathbf{U}^T + \mathbf{U}'(\mathbf{U}' - \mathbf{U})^T\|_F \\ 392 &\leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) (\|\mathbf{U}\| + \|\mathbf{U}'\|) \|\mathbf{U}' - \mathbf{U}\|_F, \\ 393 \end{aligned}$$

394 where the second inequality follows from the ℓ_1/ℓ_2 -RIP of \mathcal{A} . It follows that

$$\begin{aligned} 395 \quad \liminf_{\mathbf{U}' \rightarrow \mathbf{U}} \frac{|f(\mathbf{U}') - f(\mathbf{U})|}{\|\mathbf{U}' - \mathbf{U}\|_F} &\leq \lim_{\mathbf{U}' \rightarrow \mathbf{U}} \frac{(\sqrt{2/\pi} + \delta)(\|\mathbf{U}\| + \|\mathbf{U}'\|) \|\mathbf{U}' - \mathbf{U}\|_F}{\|\mathbf{U}' - \mathbf{U}\|_F} \\ 396 &= 2 \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{U}\|. \\ 397 \end{aligned}$$

398 Upon taking $\mathbf{U}' = \mathbf{U} + t\mathbf{D}$, $t \rightarrow 0$ and invoking (3.8), we get

$$399 \quad \|\mathbf{D}\|_F \leq 2 \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{U}\|, \quad \forall \mathbf{D} \in \partial f(\mathbf{U}).$$

400 To complete the proof, it remains to note that for any $\mathbf{U} \in \mathbb{R}^{n \times r}$ satisfying $\text{dist}(\mathbf{U}, \mathcal{U}) \leq$
 401 $\frac{2\alpha}{\tau}$, where α, τ are given in (3.4), (3.6), respectively, the triangle inequality yields
 402 $\|\mathbf{U}\| \leq \|\mathbf{U}^*\|_F + \frac{2\alpha}{\tau}$. \square

403 By collecting [Proposition 2](#), [Proposition 3](#), and [Proposition 4](#) together and in-
 404 voking [Theorem 2](#), we obtain the following guarantees for the SubGM² when it is
 405 applied to the robust low-rank matrix recovery problem [\(1.4\)](#):

406 **THEOREM 3** (nonconvex robust low-rank matrix recovery: PSD case). *Consider*
 407 *the measurement model [\(1.3\)](#), where \mathbf{X}^* is an $n \times n$ rank- r symmetric positive semi-*
 408 *definite matrix. Let $0 < \delta < \frac{1}{3} \sqrt{\frac{2}{\pi}}$ be given. Suppose that the fraction of outliers p in*
 409 *the measurement vector \mathbf{y} satisfies [\(3.3\)](#), and that the linear operators \mathcal{A} , \mathcal{A}_{Ω^c} possess*
 410 *the ℓ_1/ℓ_2 -RIP [\(3.1\)](#), [\(3.2\)](#), respectively. Let α , τ , and κ be given by [\(3.4\)](#), [\(3.6\)](#), and*
 411 *[\(3.7\)](#), respectively. Under such setting, suppose that we apply the SubGM in [Algo-](#)*
 412 *rithm 2.1 to solve [\(1.4\)](#), where the initial point \mathbf{U}_0 satisfies $\text{dist}(\mathbf{U}_0, \mathcal{U}) < \frac{2\alpha}{\tau}$ and the*
 413 *geometrically diminishing step sizes $\mu_k = \rho^k \mu_0$ are used with μ_0 , ρ satisfying [\(2.4\)](#),*
 414 *[\(2.5\)](#), respectively. Then, the sequence of iterates $\{\mathbf{U}_k\}_{k \geq 0}$ generated by the SubGM*
 415 *will converge to a point in \mathcal{U} at a linear rate:*

$$416 \quad \text{dist}(\mathbf{U}_k, \mathcal{U}) \leq \rho^k \max \left\{ \text{dist}(\mathbf{U}_0, \mathcal{U}), \mu_0 \frac{\max\{\kappa^2, 2\alpha^2\}}{\alpha} \right\}.$$

417 *Moreover, the ground-truth matrix \mathbf{X}^* can be exactly recovered by any point $\mathbf{U}^* \in \mathcal{U}$*
 418 *via $\mathbf{X}^* = \mathbf{U}^* \mathbf{U}^{*\text{T}}$.*

419 We remark that a similar result for the smooth counterpart [\(1.2\)](#) without any out-
 420 liers is established in [[42](#), [Theorem 3.3](#)]. Our [Theorem 3](#) implies that the nonsmooth
 421 problem [\(1.4\)](#) can be solved *as efficiently as* its smooth counterpart [\(1.2\)](#), even in the
 422 presence of a substantial fraction of outliers in the measurement vector.

423 **3.5. Initializing the SubGM.** We now discuss some potential initialization
 424 strategies for the SubGM. A common approach to generating an appropriate ini-
 425 tialization for matrix recovery-type problems is the spectral method. In our context,
 426 this entails simply computing the rank- r approximation of $\frac{1}{m} \mathcal{A}^*(\mathbf{y}) = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{A}_i$,
 427 where \mathcal{A}^* is the adjoint operator of \mathcal{A} . Specifically, let $\mathbf{P} \mathbf{\Pi} \mathbf{Q}^{\text{T}}$ be a rank- r SVD of
 428 $\frac{1}{m} \mathcal{A}^*(\mathbf{y})$, where \mathbf{P} , \mathbf{Q} have orthonormal columns and $\mathbf{\Pi}$ is an $r \times r$ diagonal matrix
 429 with the top r singular values of $\frac{1}{m} \mathcal{A}^*(\mathbf{y})$ along its diagonal. In the symmetric po-
 430 sitive semidefinite case, we may assume without loss of generality that $\mathbf{A}_1, \dots, \mathbf{A}_m$
 431 are symmetric. Then, we can take $\mathbf{U}_0 = \mathbf{P} \mathbf{\Pi}^{1/2}$ as the initialization. The main idea
 432 behind this approach is that when there is no outlier (i.e., $\mathbf{y} = \mathcal{A}(\mathbf{X}^*)$ as in [\(1.1\)](#)),
 433 we have $\frac{1}{m} \mathcal{A}^*(\mathbf{y}) = \frac{1}{m} \mathcal{A}^*(\mathcal{A}(\mathbf{X}^*)) \approx \mathbf{X}^*$ when $\frac{1}{m} \mathcal{A}^* \mathcal{A}$ is close to a unitary operator
 434 for low-rank matrices. Thus, \mathbf{U}_0 is also expected to be close to \mathcal{U} . However, when
 435 the measurements are corrupted by outliers, it is possible that $\frac{1}{m} \mathcal{A}^*(\mathbf{y})$ is perturbed
 436 away from $\frac{1}{m} \mathcal{A}^*(\mathcal{A}(\mathbf{X}^*))$ and thus \mathbf{U}_0 may not be close enough to \mathcal{U} . To mitigate the
 437 influence of outliers, Li et al. [[30](#)] have recently proposed a truncated spectral method
 438 for initialization, in which the spectral method is applied to an operator that is formed
 439 by using those measurements whose absolute values do not deviate too much from the
 440 median of the absolute values of certain sampled measurements; see [Algorithm 3.1](#).
 441 They showed that under appropriate conditions, the truncated spectral method can
 442 output an initialization that satisfies the requirement of [Theorem 3](#).

443 **THEOREM 4** (proximity of initialization to optimal set: PSD case; cf. [[30](#), [The-](#)
 444 [orem 3.3](#)]). *Let $r \geq 1$ be given and set $\bar{c} = \frac{\|\mathbf{X}^*\|_F}{\sqrt{r\sigma_r(\mathbf{X}^*)}}$. Suppose that the matrices*
 445 *$\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n \times n}$ defining the linear measurement operator \mathcal{A} are symmetric and*

²In practice, we can just take $\text{Sign}(0) = 0$ when applying the SubGM to solve [\(1.4\)](#).

Algorithm 3.1 Truncated Spectral Method for Initialization [30]

Input: measurement vector \mathbf{y} ; sensing matrices $\mathbf{A}_1, \dots, \mathbf{A}_m$; threshold $\beta > 0$;

- 1: set $\mathbf{y}_1 = \{y_i\}_{i=1}^{\lfloor m/2 \rfloor}$, $\mathbf{y}_2 = \{y_i\}_{\lfloor m/2 \rfloor + 1}^m$;
- 2: Compute the rank- r SVD of

$$\mathbf{E} = \frac{1}{\lfloor m/2 \rfloor} \sum_{i=1}^{\lfloor m/2 \rfloor} y_i \mathbf{A}_i \mathbb{I}_{\{|y_i| \leq \beta \cdot \text{median}(|\mathbf{y}_2|)\}}$$

and denote it by $\mathbf{P}\mathbf{\Pi}\mathbf{Q}^T$, where

$$\mathbb{I}_{\{|y_i| \leq \beta \cdot \text{median}(|\mathbf{y}_2|)\}} = \begin{cases} 1 & \text{if } |y_i| \leq \beta \cdot \text{median}(|\mathbf{y}_2|), \\ 0 & \text{otherwise;} \end{cases}$$

Output: $\mathbf{U}_0 = \mathbf{P}\mathbf{\Pi}^{1/2}$, $\mathbf{V}_0 = \mathbf{Q}\mathbf{\Pi}^{1/2}$;

446 have *i.i.d.* standard Gaussian entries on and above the diagonal, and that the num-
 447 ber of measurements m satisfies $m \gtrsim \beta^2 \bar{c}^2 n r^2 \log n$, where $\beta = 2 \log(r^{1/4} \bar{c}^{1/2} + 20)$.
 448 Furthermore, suppose that the fraction of outliers p in the measurement vector \mathbf{y}
 449 satisfies $p \lesssim \frac{1}{\sqrt{r\bar{c}}}$. Then, with overwhelming probability, [Algorithm 3.1](#) outputs an ini-
 450 tialization $\mathbf{U}_0 \in \mathbb{R}^{n \times r}$ satisfying $\text{dist}(\mathbf{U}_0, \mathcal{U}) \lesssim \sigma_r(\mathbf{X}^*)$ and hence also the requirement
 451 of [Theorem 3](#) (as $\sigma_r(\mathbf{X}^*)$ is of the same order as $\frac{2\alpha}{\tau}$).

452 Note that the requirements on the number of measurements and the fraction of
 453 outliers that can be tolerated are slightly more stringent than those in [Proposition 1](#)
 454 and [Theorem 3](#). However, as will be illustrated in [Section 5](#), our numerical experi-
 455 ments show that even a randomly initialized SubGM can very efficiently find the
 456 global minimum and hence recover the ground-truth matrix \mathbf{X}^* . A theoretical jus-
 457 tification of such a phenomenon will be the subject of a future study. We suspect
 458 that it may be possible to relax the requirement on the initialization in [Theorem 3](#) or
 459 to show that the SubGM enters the region $\{\mathbf{U} : \text{dist}(\mathbf{U}, \mathcal{U}) < \frac{2\alpha}{\tau}\}$ very quickly even
 460 though the random initialization lies outside of this region.

461 **4. Nonconvex Robust Low-Rank Matrix Recovery: General Case.** In
 462 this section we consider the general setting where \mathbf{X}^* is a rank- r $n_1 \times n_2$ matrix. To
 463 extend the nonsmooth nonconvex formulation (1.4) to this setting, a natural approach
 464 is to use the factorization $\mathbf{X} = \mathbf{U}\mathbf{V}^T$ with $\mathbf{U} \in \mathbb{R}^{n_1 \times r}$ and $\mathbf{V} \in \mathbb{R}^{n_2 \times r}$. However, such
 465 a factorization is ambiguous in the sense that if $\mathbf{X} = \mathbf{U}\mathbf{V}^T$, then $\mathbf{X} = (\mathbf{U}\mathbf{T})(\mathbf{V}\mathbf{T}^{-T})^T$
 466 for any invertible matrix $\mathbf{T} \in \mathbb{R}^{r \times r}$. To address this issue, we introduce the nonsmooth
 467 nonconvex regularizer

$$468 \quad (4.1) \quad \phi(\mathbf{U}, \mathbf{V}) := \|\mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V}\|_F,$$

469 which aims to balance the factors \mathbf{U} and \mathbf{V} , and solve the following regularized
 470 problem:

$$471 \quad (4.2) \quad \underset{\mathbf{U} \in \mathbb{R}^{n_1 \times r}, \mathbf{V} \in \mathbb{R}^{n_2 \times r}}{\text{minimize}} \left\{ g(\mathbf{U}, \mathbf{V}) := \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^T)\|_1 + \lambda \|\mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V}\|_F \right\}.$$

472 Here, $\lambda > 0$ is a regularization parameter. We remark that a similar regularizer,
 473 namely,

$$474 \quad \tilde{\phi}(\mathbf{U}, \mathbf{V}) := \|\mathbf{U}^T \mathbf{U} - \mathbf{V}^T \mathbf{V}\|_F^2,$$

475 has been introduced in [35, 42, 52] to account for the ambiguities caused by in-
 476 vertible transformations when minimizing the squared ℓ_2 -loss function $(\mathbf{U}, \mathbf{V}) \mapsto$
 477 $\frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top)\|_2^2$. However, such a regularizer is not entirely suitable for the ℓ_1 -loss
 478 function, as it is no longer clear that the resulting problem will satisfy the sharpness
 479 condition in Definition 1.

480 To simplify notation, we stack \mathbf{U} and \mathbf{V} together as $\mathbf{W} = [\mathbf{U}^\top \ \mathbf{V}^\top]^\top$ and write
 481 $g(\mathbf{W})$ for $g(\mathbf{U}, \mathbf{V})$. Observe that the regularizer ϕ achieves its minimum value of 0
 482 when \mathbf{U} and \mathbf{V} have the same Gram matrices; i.e., $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V}$. Now, let $\mathbf{X}^* =$
 483 $\Phi \Sigma \Psi^\top$ be a rank- r SVD of \mathbf{X}^* , where $\Phi \in \mathbb{R}^{n_1 \times r}$, $\Psi \in \mathbb{R}^{n_2 \times r}$ have orthonormal
 484 columns and $\Sigma \in \mathbb{R}^{r \times r}$ is a diagonal matrix. Define

$$485 \quad \mathbf{U}^* = \Phi \Sigma^{1/2}, \quad \mathbf{V}^* = \Psi \Sigma^{1/2}, \quad \mathbf{W}^* = [\mathbf{U}^{*\top} \ \mathbf{V}^{*\top}]^\top.$$

486 The orthogonal invariance of g (i.e., $g(\mathbf{W}) = g(\mathbf{W}\mathbf{R})$ for any $\mathbf{R} \in \mathcal{O}_r$) implies that g
 487 is constant on the set

$$488 \quad \mathcal{W} := \{\mathbf{W}^* \mathbf{R} : \mathbf{R} \in \mathcal{O}_r\}.$$

489 **4.1. Sharpness and Exact Recovery.** Our immediate goal is to show that \mathcal{W}
 490 is the set of global minima of (4.2). Towards that end, let $0 < \delta < \frac{1}{3} \sqrt{\frac{2}{\pi}}$ be given.
 491 Suppose that the fraction of outliers p in the measurement vector \mathbf{y} satisfies (3.3),
 492 and that the linear operators $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$ and $\mathcal{A}_{\Omega^c} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^{|\Omega^c|}$ possess
 493 the ℓ_1/ℓ_2 -RIP (3.1) and (3.2), respectively.³ Using the argument in the proof of
 494 Proposition 2, we get

$$495 \quad (4.3) \quad \bar{g}(\mathbf{W}) - \bar{g}(\mathbf{W}^*) \geq \left(2(1-p) \left(\sqrt{\frac{2}{\pi}} - \delta \right) - \left(\sqrt{\frac{2}{\pi}} + \delta \right) \right) \|\mathbf{U}\mathbf{V}^\top - \mathbf{X}^*\|_F,$$

496 where

$$497 \quad \bar{g}(\mathbf{W}) = \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^\top)\|_1.$$

498 In particular, we see that $\bar{g}(\mathbf{W}) > \bar{g}(\mathbf{W}^*)$ whenever $\mathbf{U}\mathbf{V}^\top \neq \mathbf{X}^*$. Since $\mathbf{U}^{*\top} \mathbf{U}^* =$
 499 $\mathbf{V}^{*\top} \mathbf{V}^*$ by construction, we conclude that \mathbf{W}^* is a global minimum of (4.2), as \mathbf{W}^* is
 500 a global minimum of both the first term \bar{g} and the second term ϕ of g . It then follows
 501 from the orthogonal invariance of g that every element in \mathcal{W} is a global minimum
 502 of (4.2). The following result further establishes that \mathcal{W} is exactly the set of global
 503 minima of (4.2) and g is sharp.

504 **PROPOSITION 5** (sharpness and exact recovery with outliers: general case). *Let*
 505 $0 < \delta < \frac{1}{3} \sqrt{\frac{2}{\pi}}$ *be given. Suppose that the fraction of outliers* p *satisfies (3.3), and that*
 506 *the linear operators* \mathcal{A} *and* \mathcal{A}_{Ω^c} *possess the* ℓ_1/ℓ_2 -RIP (3.1) *and (3.2), respectively.*
 507 *Then, the objective function* g *in (4.2) satisfies*

$$508 \quad g(\mathbf{W}) - g(\mathbf{W}^*) \geq \alpha \text{dist}(\mathbf{W}, \mathcal{W})$$

509 *for any* $\mathbf{W} \in \mathbb{R}^{(n_1+n_2) \times r}$, *where*

$$510 \quad (4.4) \quad \alpha = \sqrt{\sqrt{2} - 1} \cdot \min \left\{ 2(1-p) \left(\sqrt{\frac{2}{\pi}} - \delta \right) - \left(\sqrt{\frac{2}{\pi}} + \delta \right), 2\lambda \right\} \cdot \sigma_r(\mathbf{X}^*) > 0.$$

³It can be shown that modulo the constants, the Gaussian measurement operator $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$ will possess the ℓ_1/ℓ_2 -RIPs (3.1) and (3.2) with high probability as long as $m \gtrsim \max\{n_1, n_2\}r$. To avoid any distraction caused by the new constants, we shall simply use the ℓ_1/ℓ_2 -RIPs (3.1) and (3.2) in our derivation.

511 *In particular, the set \mathcal{W} is precisely the set of global minima of (4.2) and the objective*
 512 *function g is sharp with parameter $\alpha > 0$.*

513 *Proof of Proposition 5.* Let $\zeta(p, \delta) = 2(1 - p) \left(\sqrt{\frac{2}{\pi}} - \delta \right) - \left(\sqrt{\frac{2}{\pi}} + \delta \right)$. Since
 514 $\mathbf{U}^{\star\text{T}}\mathbf{U}^{\star} = \mathbf{V}^{\star\text{T}}\mathbf{V}^{\star}$, we have $\phi(\mathbf{W}^{\star}) = 0$ by (4.1) and

$$\begin{aligned}
 515 \quad g(\mathbf{W}) - g(\mathbf{W}^{\star}) &= \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{U}\mathbf{V}^{\text{T}})\|_1 - \frac{1}{m} \|\mathbf{y} - \mathcal{A}(\mathbf{X}^{\star})\|_1 + \lambda \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F \\
 516 \quad &\geq \zeta(p, \delta) \|\mathbf{X}^{\star} - \mathbf{U}\mathbf{V}^{\text{T}}\|_F + \lambda \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F \\
 517 \quad &\geq \min \{ \zeta(p, \delta), 2\lambda \} \left(\|\mathbf{X}^{\star} - \mathbf{U}\mathbf{V}^{\text{T}}\|_F + \frac{1}{2} \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F \right) \\
 518 \quad &\geq \min \{ \zeta(p, \delta), 2\lambda \} \sqrt{\|\mathbf{X}^{\star} - \mathbf{U}\mathbf{V}^{\text{T}}\|_F^2 + \frac{1}{4} \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F^2} \\
 519 \quad &\geq \min \left\{ \frac{\zeta(p, \delta)}{2}, \lambda \right\} \|\mathbf{W}\mathbf{W}^{\text{T}} - \mathbf{W}^{\star}\mathbf{W}^{\star\text{T}}\|_F \\
 520 \quad &\geq \min \left\{ \frac{\zeta(p, \delta)}{2}, \lambda \right\} \sqrt{2(\sqrt{2} - 1)} \sigma_r(\mathbf{W}^{\star}) \text{dist}(\mathbf{W}, \mathcal{W}) \\
 521 \quad &= \min \{ \zeta(p, \delta), 2\lambda \} \sqrt{\sqrt{2} - 1} \sigma_r^{1/2}(\mathbf{X}^{\star}) \text{dist}(\mathbf{W}, \mathcal{W}),
 \end{aligned}$$

523 where the first inequality follows from (4.3), the fourth inequality follows from

$$\begin{aligned}
 524 \quad \|\mathbf{X}^{\star} - \mathbf{U}\mathbf{V}^{\text{T}}\|_F^2 + \frac{1}{4} \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F^2 &= \|\mathbf{U}^{\star}\mathbf{V}^{\star\text{T}} - \mathbf{U}\mathbf{V}^{\text{T}}\|_F^2 + \frac{1}{4} \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F^2 \\
 525 \quad &= \frac{1}{4} \|\mathbf{W}\mathbf{W}^{\text{T}} - \mathbf{W}^{\star}\mathbf{W}^{\star\text{T}}\|_F^2 + \nu(\mathbf{W}) \\
 526
 \end{aligned}$$

527 with

$$\begin{aligned}
 528 \quad \nu(\mathbf{W}) &= \frac{1}{2} \|\mathbf{U}\mathbf{V}^{\text{T}} - \mathbf{U}^{\star}\mathbf{V}^{\star\text{T}}\|_F^2 + \frac{1}{4} \|\mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}\|_F^2 \\
 529 \quad &\quad - \frac{1}{4} \|\mathbf{U}\mathbf{U}^{\text{T}} - \mathbf{U}^{\star}\mathbf{U}^{\star\text{T}}\|_F^2 - \frac{1}{4} \|\mathbf{V}\mathbf{V}^{\text{T}} - \mathbf{V}^{\star}\mathbf{V}^{\star\text{T}}\|_F^2 \\
 530 \quad &= \frac{1}{2} \|\mathbf{U}^{\text{T}}\mathbf{U}^{\star}\|_F^2 + \frac{1}{2} \|\mathbf{V}^{\text{T}}\mathbf{V}^{\star}\|_F^2 - \langle \mathbf{U}\mathbf{V}^{\text{T}}, \mathbf{U}^{\star}\mathbf{V}^{\star\text{T}} \rangle \\
 531 \quad &\quad + \frac{1}{2} \|\mathbf{U}^{\star}\mathbf{V}^{\star\text{T}}\|_F^2 - \frac{1}{4} \|\mathbf{U}^{\star}\mathbf{U}^{\star\text{T}}\|_F^2 - \frac{1}{4} \|\mathbf{V}^{\star}\mathbf{V}^{\star\text{T}}\|_F^2 \\
 532 \quad &= \frac{1}{2} \|\mathbf{U}^{\text{T}}\mathbf{U}^{\star} - \mathbf{V}^{\text{T}}\mathbf{V}^{\star}\|_F^2 + \frac{1}{2} \|\mathbf{U}^{\star}\mathbf{V}^{\star\text{T}}\|_F^2 - \frac{1}{4} \|\mathbf{U}^{\star}\mathbf{U}^{\star\text{T}}\|_F^2 - \frac{1}{4} \|\mathbf{V}^{\star}\mathbf{V}^{\star\text{T}}\|_F^2 \\
 533 \quad &= \frac{1}{2} \|\mathbf{U}^{\text{T}}\mathbf{U}^{\star} - \mathbf{V}^{\text{T}}\mathbf{V}^{\star}\|_F^2 \geq 0 \\
 534
 \end{aligned}$$

535 (recall that $\mathbf{U}^{\star\text{T}}\mathbf{U}^{\star} = \mathbf{V}^{\star\text{T}}\mathbf{V}^{\star}$), the fifth inequality is from Lemma 1, and the last
 536 equality follows from the fact that $\sigma_r(\mathbf{W}^{\star}) = \sqrt{2} \sigma_r^{1/2}(\mathbf{X}^{\star})$. This completes the proof. \square

537 By comparing Proposition 2 and Proposition 5, we see that the fraction of outliers
 538 that can be tolerated for exact recovery is the same in both the symmetric positive
 539 semidefinite and general cases. Moreover, the sharpness parameter α in (4.4) demon-
 540 strates the role that the regularizer ϕ plays: When the regularizer ϕ is absent (which
 541 corresponds to $\lambda = 0$), although every element in \mathcal{W} is still a global minimum of (4.2),
 542 we cannot guarantee that there is no other global minimum. Indeed, when $\lambda = 0$, the

543 pair $(\mathbf{U}^*\mathbf{T}, \mathbf{V}^*\mathbf{T}^{-\text{T}})$ is a global minimum of (4.2) for any invertible matrix $\mathbf{T} \in \mathbb{R}^{r \times r}$.
 544 However, when $\lambda > 0$, the regularizer ϕ ensures that the pair $(\mathbf{U}^*\mathbf{T}, \mathbf{V}^*\mathbf{T}^{-\text{T}})$ is a
 545 global minimum of (4.2) only when $\mathbf{T} \in \mathcal{O}_r$.

546 **4.2. Weak Convexity.** Let us now establish the weak convexity of the objective
 547 function g in (4.2).

548 **PROPOSITION 6** (weak convexity: general case). *Suppose that the measurement*
 549 *operator \mathcal{A} satisfies the ℓ_1/ℓ_2 -RIP (3.1). Then, the objective function g in (4.2) is*
 550 *weakly convex with parameter*

$$551 \quad (4.5) \quad \tau = \sqrt{\frac{2}{\pi}} + \delta + 2\lambda.$$

552 *Proof of Proposition 6.* Since $g = \bar{g} + \lambda\phi$, it suffices to show that \bar{g} and ϕ are
 553 both weakly convex. Similar to (3.5), we apply the chain rule for subdifferentials [37,
 554 Corollary 8.11 and Theorem 10.6] to get

$$555 \quad \partial\bar{g}(\mathbf{W}) = \frac{1}{m} \begin{bmatrix} \mathcal{A}^* (\text{Sign}(\mathcal{A}(\mathbf{U}\mathbf{V}^{\text{T}}) - \mathbf{y})) \mathbf{V} \\ (\mathcal{A}^* (\text{Sign}(\mathcal{A}(\mathbf{U}\mathbf{V}^{\text{T}}) - \mathbf{y})))^{\text{T}} \mathbf{U} \end{bmatrix}.$$

556 Using this and the argument in the proof of Proposition 3, we can show that for any
 557 $\mathbf{W}, \mathbf{W}' \in \mathbb{R}^{(n_1+n_2) \times r}$,

$$558 \quad \bar{g}(\mathbf{W}') \geq \bar{g}(\mathbf{W}) + \langle \mathbf{D}, \mathbf{W}' - \mathbf{W} \rangle - \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|(\mathbf{U}' - \mathbf{U})(\mathbf{V}' - \mathbf{V})^{\text{T}}\|_F$$

$$559 \quad \geq \bar{g}(\mathbf{W}) + \langle \mathbf{D}, \mathbf{W}' - \mathbf{W} \rangle - \left(\frac{\sqrt{2/\pi} + \delta}{2} \right) \|\mathbf{W}' - \mathbf{W}\|_F^2, \quad \forall \mathbf{D} \in \partial\bar{g}(\mathbf{W});$$

561 i.e., the function \bar{g} is weakly convex with parameter $\tau_{\bar{g}} = \sqrt{\frac{2}{\pi}} + \delta$.

562 Next, define the matrices

$$563 \quad \underline{\mathbf{W}} = [\mathbf{U}^{\text{T}} \quad -\mathbf{V}^{\text{T}}]^{\text{T}}, \quad \underline{\mathbf{W}}' = [\mathbf{U}'^{\text{T}} \quad -\mathbf{V}'^{\text{T}}]^{\text{T}}$$

564 and note that $\underline{\mathbf{W}}^{\text{T}}\underline{\mathbf{W}} = \mathbf{U}^{\text{T}}\mathbf{U} - \mathbf{V}^{\text{T}}\mathbf{V}$. Furthermore, define the function $\psi : \mathbb{R}^{r \times r} \rightarrow \mathbb{R}$
 565 by $\psi(\mathbf{C}) = \|\mathbf{C}\|_F$, whose subdifferential is

$$566 \quad \partial\psi(\mathbf{C}) = \begin{cases} \left\{ \frac{\mathbf{C}}{\|\mathbf{C}\|_F} \right\}, & \mathbf{C} \neq \mathbf{0}, \\ \{ \mathbf{B} \in \mathbb{R}^{r \times r} : \|\mathbf{B}\|_F \leq 1 \}, & \mathbf{C} = \mathbf{0}. \end{cases}$$

567 Upon setting $\underline{\Delta} = \underline{\mathbf{W}}' - \underline{\mathbf{W}}$ and $\underline{\Delta} = \underline{\mathbf{W}}' - \underline{\mathbf{W}}$, we compute

$$568 \quad (4.6) \quad \begin{aligned} \phi(\mathbf{W}') &= \|\underline{\mathbf{W}}'^{\text{T}}\underline{\mathbf{W}}'\|_F \\ &= \|\underline{\mathbf{W}}^{\text{T}}\underline{\mathbf{W}} + \underline{\mathbf{W}}^{\text{T}}\underline{\Delta} + \underline{\Delta}^{\text{T}}\underline{\mathbf{W}} + \underline{\Delta}^{\text{T}}\underline{\Delta}\|_F \\ &\geq \|\underline{\mathbf{W}}^{\text{T}}\underline{\mathbf{W}} + \underline{\mathbf{W}}^{\text{T}}\underline{\Delta} + \underline{\Delta}^{\text{T}}\underline{\mathbf{W}}\|_F - \|\underline{\Delta}^{\text{T}}\underline{\Delta}\|_F \\ &\geq \|\underline{\mathbf{W}}^{\text{T}}\underline{\mathbf{W}}\|_F + \left\langle \underline{\Psi}, \underline{\mathbf{W}}^{\text{T}}\underline{\Delta} + \underline{\Delta}^{\text{T}}\underline{\mathbf{W}} \right\rangle - \|\underline{\Delta}^{\text{T}}\underline{\Delta}\|_F, \end{aligned}$$

569 where the last inequality holds for any $\underline{\Psi} \in \partial\psi(\underline{\mathbf{W}}^{\text{T}}\underline{\mathbf{W}})$ due to the convexity of the
 570 Frobenius norm. Since the Frobenius norm is subdifferentially regular [37, Example
 571 7.27], the chain rule for subdifferentials [37, Corollary 8.11 and Theorem 10.6] yields

$$572 \quad (4.7) \quad \partial\phi(\mathbf{W}) = \left\{ \underline{\mathbf{W}}(\underline{\Psi} + \underline{\Psi}^{\text{T}}) : \underline{\Psi} \in \partial\psi(\underline{\mathbf{W}}^{\text{T}}\underline{\mathbf{W}}) \right\}.$$

573 It follows from (4.6) and (4.7) that

$$574 \quad \phi(\mathbf{W}') \geq \phi(\mathbf{W}) + \langle \Phi, \mathbf{W}' - \mathbf{W} \rangle - \|\Delta^T \Delta\|_F \\ 575 \quad \geq \phi(\mathbf{W}) + \langle \Phi, \mathbf{W}' - \mathbf{W} \rangle - \|\mathbf{W}' - \mathbf{W}\|_F^2, \quad \forall \Phi \in \partial\phi(\mathbf{W});$$

577 i.e., the function ϕ is weakly convex with parameter $\tau_\phi = 2$.

578 Putting the above results together, we conclude that $g = \bar{g} + \lambda\phi$ is weakly convex
579 with parameter $\tau = \tau_{\bar{g}} + \lambda\tau_\phi$, as desired. \square

580 Unlike the sharpness condition in Proposition 5 that requires $\lambda > 0$, the weak
581 convexity condition in Proposition 6 holds even when $\lambda = 0$. Although the parameters
582 α and τ in (4.4) and (4.5) increase as λ increases from 0, the former becomes constant
583 when $\lambda \geq \frac{2(1-p)(\sqrt{2/\pi-\delta}) - (\sqrt{2/\pi+\delta})}{2}$. In view of Theorem 2, it is desirable to choose
584 λ so that the local linear convergence region $\{\mathbf{x} : \text{dist}(\mathbf{x}, \mathcal{X}) < \frac{2\alpha}{\tau}\}$ of the SubGM is
585 as large as possible. Such consideration suggests that we should set

$$586 \quad \lambda = \frac{2(1-p)(\sqrt{2/\pi-\delta}) - (\sqrt{2/\pi+\delta})}{2}.$$

587 **4.3. Putting Everything Together.** As in Subsection 3.4, before we can in-
588 voke Theorem 2 to establish convergence guarantees for the SubGM when applied to
589 the general robust low-rank matrix recovery problem (4.2), we need to bound the norm
590 of any subgradient of g in a neighborhood of \mathcal{W} . This is achieved by the following
591 result:

592 PROPOSITION 7 (bound on subgradient norm: general case). *Suppose that the*
593 *measurement operator \mathcal{A} satisfies the ℓ_1/ℓ_2 -RIP (3.1). Then, for any $\mathbf{W} \in \mathbb{R}^{(n_1+n_2) \times r}$*
594 *satisfying $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \frac{2\alpha}{\tau}$, we have*

$$595 \quad (4.8) \quad \|\mathbf{D}\|_F \leq \kappa = \max \left\{ \sqrt{\frac{2}{\pi}} + \delta, \lambda \right\} \left(\|\mathbf{W}^*\|_F + \frac{2\alpha}{\tau} \right), \quad \forall \mathbf{D} \in \partial g(\mathbf{W}).$$

597 *Proof of Proposition 7.* Observe that for any $\mathbf{W}, \mathbf{W}' \in \mathbb{R}^{(n_1+n_2) \times r}$,

$$598 \quad |g(\mathbf{W}') - g(\mathbf{W})| \leq |\bar{g}(\mathbf{W}') - \bar{g}(\mathbf{W})| + \lambda |\phi(\mathbf{W}') - \phi(\mathbf{W})| \\ 599 \quad \leq \frac{1}{m} \|\mathcal{A}(\mathbf{U}\mathbf{V}^T - \mathbf{U}'\mathbf{V}'^T)\|_1 + \lambda (\|\mathbf{U}^T\mathbf{U} - \mathbf{U}'^T\mathbf{U}'\|_F + \|\mathbf{V}^T\mathbf{V} - \mathbf{V}'^T\mathbf{V}'\|_F) \\ 600 \quad \leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{U}\mathbf{V}^T - \mathbf{U}'\mathbf{V}'^T\|_F + \lambda (\|\mathbf{U}^T\mathbf{U} - \mathbf{U}'^T\mathbf{U}'\|_F + \|\mathbf{V}^T\mathbf{V} - \mathbf{V}'^T\mathbf{V}'\|_F) \\ 601 \quad \leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) (\|\mathbf{V}\|_F \|\mathbf{U} - \mathbf{U}'\|_F + \|\mathbf{U}'\|_F \|\mathbf{V} - \mathbf{V}'\|_F) \\ 602 \quad \quad + \lambda (\|\mathbf{U}\|_F + \|\mathbf{U}'\|_F) \|\mathbf{U} - \mathbf{U}'\|_F + \lambda (\|\mathbf{V}\|_F + \|\mathbf{V}'\|_F) \|\mathbf{V} - \mathbf{V}'\|_F \\ 603 \quad \leq \max \left\{ \sqrt{\frac{2}{\pi}} + \delta, \lambda \right\} (\|\mathbf{W}\|_F + \|\mathbf{W}'\|_F) \|\mathbf{W} - \mathbf{W}'\|_F, \\ 604$$

605 where the third inequality follows from the ℓ_1/ℓ_2 -RIP (3.1). Thus, similar to the
606 derivation of (3.7), for any $\mathbf{W} \in \mathbb{R}^{(n_1+n_2) \times r}$ satisfying $\text{dist}(\mathbf{W}, \mathcal{W}) \leq \frac{2\alpha}{\tau}$, where α

607 and τ are given in (4.4) and (4.5), respectively, we have

$$\begin{aligned}
608 \quad \|\mathbf{D}\|_F &\leq \max \left\{ \sqrt{\frac{2}{\pi}} + \delta, \lambda \right\} \|\mathbf{W}\|_F \\
609 \quad &\leq \max \left\{ \sqrt{\frac{2}{\pi}} + \delta, \lambda \right\} \left(\|\mathbf{W}^*\|_F + \frac{2\alpha}{\tau} \right), \quad \forall \mathbf{D} \in \partial g(\mathbf{W}). \quad \square \\
610
\end{aligned}$$

611 By collecting Proposition 5, Proposition 6, and Proposition 7 together and invo-
612 king Theorem 2, we obtain the following guarantees when the SubGM is used to solve
613 the general robust low-rank matrix recovery problem (4.2):

614 **THEOREM 5** (nonconvex robust low-rank matrix recovery: general case). *Con-*
615 *sider the measurement model (1.3), where \mathbf{X}^* is an $n_1 \times n_2$ rank- r matrix. Let*
616 *$0 < \delta < \frac{1}{3}\sqrt{\frac{2}{\pi}}$ be given. Suppose that the fraction of outliers p in the measure-*
617 *ment vector \mathbf{y} satisfies (3.3), and that the linear operators \mathcal{A} , \mathcal{A}_{Ω^c} possess the ℓ_1/ℓ_2 -*
618 *RIP (3.1), (3.2), respectively. Let α , τ , and κ be given by (4.4), (4.5), and (4.8),*
619 *respectively. Under such setting, suppose that we apply the SubGM in Algorithm 2.1*
620 *to solve (4.2), where the initial point \mathbf{W}_0 satisfies $\text{dist}(\mathbf{W}_0, \mathcal{W}) < \frac{2\alpha}{\tau}$ and the geome-*
621 *trically diminishing step sizes $\mu_k = \rho^k \mu_0$ are used with μ_0 , ρ satisfying (2.4), (2.5),*
622 *respectively. Then, the sequence of iterates $\{\mathbf{W}_k\}_{k \geq 0}$ generated by the SubGM will*
623 *converge to a point in \mathcal{W} at a linear rate:*

$$624 \quad \text{dist}(\mathbf{W}_k, \mathcal{W}) \leq \rho^k \max \left\{ \text{dist}(\mathbf{W}_0, \mathcal{W}), \mu_0 \frac{\max\{\kappa^2, 2\alpha^2\}}{\alpha} \right\}.$$

625 Moreover, the ground-truth matrix \mathbf{X}^* can be exactly recovered by any point $\mathbf{W}^* \in \mathcal{W}$
626 via $\mathbf{X}^* = \mathbf{U}^* \mathbf{V}^{*\text{T}}$.

627 **4.4. Initializing the SubGM.** In the general case, we can still use the trunca-
628 ted spectral method in Algorithm 3.1 to obtain a good initialization for the SubGM.
629 Specifically, we take $\mathbf{W}_0 = [\mathbf{U}_0^{\text{T}} \quad \mathbf{V}_0^{\text{T}}]^{\text{T}}$ as the initialization, where $\mathbf{U}_0, \mathbf{V}_0$ are the
630 outputs of Algorithm 3.1. Then, we have the following result, which is essentially a
631 restatement of [30, Theorem 3.3]:

632 **THEOREM 6** (proximity of initialization to optimal set: general case). *Let $r \geq 1$*
633 *be given and set $n = n_1 + n_2$, $\bar{c} = \frac{\|\mathbf{X}^*\|_F}{\sqrt{r}\sigma_r(\mathbf{X}^*)}$. Suppose that the matrices $\mathbf{A}_1, \dots, \mathbf{A}_m \in$
634 $\mathbb{R}^{n_1 \times n_2}$ defining the linear measurement operator \mathcal{A} have *i.i.d.* standard Gaussian
635 entries, and that the number of measurements m satisfies $m \gtrsim \beta^2 \bar{c}^2 n r^2 \log n$, where
636 $\beta = 2 \log(r^{1/4} \bar{c}^{1/2} + 20)$. Furthermore, suppose that the fraction of outliers p in
637 the measurement vector \mathbf{y} satisfies $p \lesssim \frac{1}{\sqrt{r\bar{c}}}$. Then, with overwhelming probability,
638 Algorithm 3.1 outputs an initialization $\mathbf{W}_0 \in \mathbb{R}^{(n_1+n_2) \times r}$ satisfying $\text{dist}(\mathbf{W}_0, \mathcal{U}) \lesssim$
639 $\sigma_r(\mathbf{X}^*)$ and hence also the requirement of Theorem 5.*

640 **5. Experiments.** In this section we conduct experiments to illustrate the per-
641 formance of the SubGM when applied to robust low-rank matrix recovery problems.
642 The experiments on synthetic data show that the SubGM can exactly and efficiently
643 recover the underlying low-rank matrix from its linear measurements even in the pre-
644 sence of outliers, thus corroborating the result in Theorem 3.

645 We generate the underlying low-rank matrix $\mathbf{X}^* = \mathbf{U}^* \mathbf{U}^{*\text{T}}$ by generating $\mathbf{U}^* \in$
646 $\mathbb{R}^{n \times r}$ with *i.i.d.* standard Gaussian entries. Similarly, we generate the entries of the m
647 sensing matrices $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n \times n}$ (which define the linear measurement operator

648 \mathcal{A}) in an *i.i.d.* fashion according to the standard Gaussian distribution. To generate
 649 the outlier vector $\mathbf{s}^* \in \mathbb{R}^m$, we first randomly select pm locations. Then, we fill each
 650 of the selected location with an *i.i.d.* mean 0 and variance 100 Gaussian entry, while
 651 the remaining locations are set to 0. Here, p is the ratio of the nonzero elements in \mathbf{s}^* .
 652 According to (1.3), the measurement vector \mathbf{y} is then generated by $\mathbf{y} = \mathcal{A}(\mathbf{X}^*) + \mathbf{s}^*$;
 653 i.e., $y_i = \langle \mathbf{A}_i, \mathbf{X}^* \rangle + s_i^*$ for $i = 1, \dots, m$.

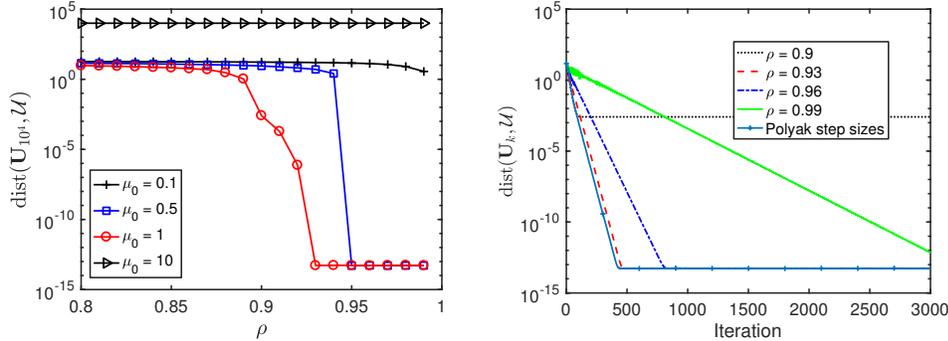
654 To illustrate the performance of the SubGM for recovering the underlying low-
 655 rank matrix \mathbf{X}^* from \mathbf{y} , we first set $n = 50$, $r = 5$, and $p = 0.3$. Throughout the
 656 experiments, we initialize the SubGM with a randomly generated standard Gaussian
 657 vector, as it gives similar practical performance as the one obtained by the truncated
 658 spectral method in Algorithm 3.1. We first run the SubGM for 10^4 iterations using
 659 the geometrically diminishing step sizes $\mu_k = \rho^k \mu_0$, where the initial step size μ_0
 660 and decay rate ρ are selected from $\{0.1, 0.5, 1, 10\}$ and $\{0.80, 0.81, 0.82, \dots, 0.99\}$,
 661 respectively. For each pair of parameters (μ_0, ρ) , we plot the distance of the last
 662 iterate to \mathcal{U} (i.e., $\text{dist}(\mathbf{U}_{10^4}, \mathcal{U})$) in Figure 2a. When the SubGM diverges, we simply
 663 set $\text{dist}(\mathbf{U}_{10^4}, \mathcal{U}) = 10^4$ for the purpose of presenting all results in the same figure.
 664 As observed from Figure 2a, the SubGM diverges when μ_0 is large, say, $\mu_0 = 10$. On
 665 the other hand, it converges to a global minimum when $\mu_0 = 1$, $\rho \in [0.93, 0.99]$ and
 666 $\mu_0 = 0.5$, $\rho \in [0.95, 0.99]$. It is worth noting that the SubGM converges to a global
 667 minimum when $\mu_0 = 1, \rho = 0.93$, but not when $\mu_0 = 0.5, \rho = 0.93$. This is consistent
 668 with Theorem 2, which shows that a larger initial step size μ_0 allows for a smaller
 669 decay rate ρ . Such a phenomenon can also be observed in the case where $\mu_0 = 0.1$,
 670 for which the SubGM fails to find a global minimum even when $\rho \in [0.95, 0.99]$.

671 In Figure 2b, we fix $\mu_0 = 1$ and plot the convergence behavior of the SubGM
 672 with $\rho \in \{0.9, 0.93, 0.96, 0.99\}$. As observed from the figure, when ρ is not too small
 673 (say, larger than 0.93), the distances $\{\text{dist}(\mathbf{U}_k, \mathcal{U})\}_{k \geq 0}$ converge to 0 at a linear rate,
 674 thus implying that the SubGM with geometrically diminishing step sizes can exactly
 675 recover the underlying low-rank matrix \mathbf{X}^* . We observe that a smaller ρ gives faster
 676 convergence. This corroborates the results in Theorem 2, which guarantee that
 677 $\{\text{dist}(\mathbf{U}_k, \mathcal{U})\}_{k \geq 0}$ decays at the rate $O(\rho^k)$ as long as ρ is not too small (i.e., satisfying
 678 (2.5)). We also consider the SubGM with the Polyak step size rule [36], which, in the
 679 context of (1.4), is given by $\mu_k = \frac{f(\mathbf{U}_k) - f^*}{\|\mathbf{d}_k\|^2}$, where f^* is the optimal value of (1.4) and
 680 $\mathbf{d}_k \in \partial f(\mathbf{U}_k)$ (the method terminates when $\mathbf{d}_k = \mathbf{0}$). The convergence rate of such
 681 method for sharp weakly convex minimization has been analyzed in [14]. We plot
 682 the convergence behavior of the SubGM with the Polyak step size rule in Figure 2b,
 683 which also shows its linear convergence. However, we note that the Polyak step size
 684 rule is generally not easy to implement, as it requires the knowledge of f^* .

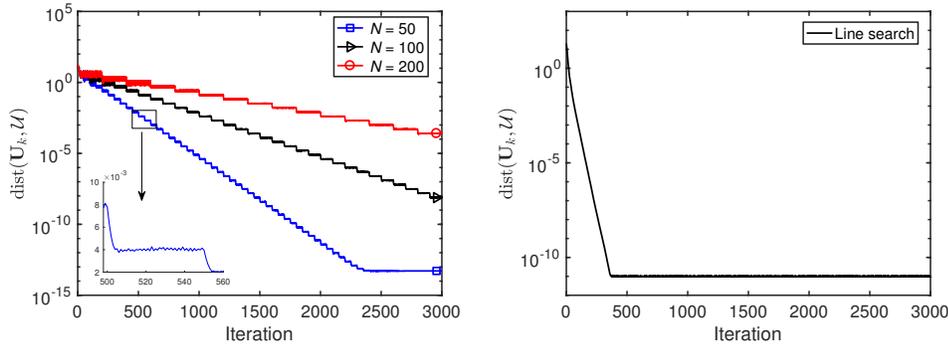
685 Then, we consider the SubGM with piecewise geometrically diminishing step si-
 686 zes, which dates as far back as to the work [39] and has recently been used in [54].
 687 Specifically, we set $\mu_k = \frac{1}{2^{\lfloor k/N \rfloor}}$ with $N \in \{50, 100, 200\}$. Compared to the vanilla
 688 strategy $\mu_k = \rho^k \mu_0$, the piecewise strategy allows for a smaller decay rate ρ (here,
 689 we use $\rho = \frac{1}{2}$) and keeps the same step size for N iterations. As can be seen from
 690 Figure 2c, the method converges at a piecewise linear rate. Nevertheless, we observe
 691 that the piecewise strategy is slightly less efficient than the vanilla one in general.

692 We also consider a modified backtracking line search strategy in [34] to choose the
 693 step size. Although such a strategy is generally designed for smooth problems, it is
 694 empirically used in [54] for a nonsmooth nonconvex optimization problem to achieve
 695 fast convergence. Inspired by the strategy of choosing geometrically diminishing step
 696 sizes, we modify the backtracking line search strategy in [34] by (i) setting $\mu_k = \mu_{k-1}$

697 and (ii) reducing it according to $\mu_k \leftarrow \mu_k \rho$ until the condition $f(\mathbf{U}_k - \mu_k \mathbf{d}_k) >$
 698 $f(\mathbf{U}_k) - \eta \mu_k \|\mathbf{d}_k\|$ is satisfied. We set $\eta = 10^{-3}$, $\rho = 0.85$, $\mu_0 = 1$ and plot the
 699 convergence behavior of the resulting method in Figure 2d. As can be seen from the
 700 figure, the method converges at a linear rate. Moreover, we observe empirically that
 701 the choice of parameters above works for other settings (i.e., different n, r, m, p). We
 702 leave the convergence analysis of the SubGM with backtracking line search as a future
 703 work.



(a) Distance of last iterate to optimal set with $\mu_0 \in \{0.1, 0.5, 1, 10\}$ and $\rho \in \{0.80, 0.81, \dots, 0.99\}$ (b) Convergence of SubGM with geometrically diminishing ($\mu_k = \rho^k$, $\rho \in \{0.90, 0.93, 0.96, 0.99\}$) and Polyak step sizes



(c) Convergence of SubGM with piecewise geometrically diminishing ($\mu_k = \frac{1}{2^{\lfloor k/N \rfloor}}$, $N \in \{50, 100, 200\}$) step sizes (d) Convergence of SubGM with modified backtracking line search ($\eta = 10^{-3}$, $\rho = 0.85$, $\mu_0 = 1$)

Fig. 2: Behavior of SubGM when applied to robust low-rank matrix recovery with $n = 50$, $r = 5$, $m = 5nr$, and $p = 0.3$.

704 Next, we study the performance of the SubGM with geometrically diminishing
 705 step sizes by varying the outlier ratio p and the number of measurements m . In these
 706 experiments we run the SubGM for 2×10^3 iterations with initial step size $\mu_0 = 1$ and
 707 decay rate $\rho = 0.99$. We also conduct experiments on the median-truncated gradient
 708 descent (MTGD) with the setting used in [30]. In particular, we initialize the MTGD
 709 with the truncated spectral method in Algorithm 3.1 and run it for 10^4 iterations.
 710 For each pair of p and m , 10 Monte Carlo trials are carried out, and for each trial

711 we declare the recovery to be successful if the relative reconstruction error satisfies
 712 $\frac{\|\widehat{\mathbf{X}} - \mathbf{X}^*\|_F}{\|\mathbf{X}^*\|_F} \leq 10^{-6}$, where $\widehat{\mathbf{X}}$ is the reconstructed matrix. Figure 3 displays the phase
 713 transition of MTGD and SubGM using the average result of 10 independent trials.
 714 In this figure, white indicates successful recovery while black indicates failure. It is of
 715 interest to observe that when the outlier ratio p is small, both the SubGM and MTGD
 716 can exactly recover the underlying low-rank matrix \mathbf{X}^* even with only $m = 2nr$
 717 measurements. On the other hand, given sufficiently large number of measurements
 718 (say $m = 7nr$), the SubGM is able to exactly recover the ground-truth matrix even
 719 when half of the measurements are corrupted by outliers, while the MTGD fails in
 720 this case. In particular, by comparing Figure 3a with Figure 3b, we observe that the
 721 SubGM is more robust to outliers than MTGD, especially in the case of high outlier
 722 ratio. We also observe from Figure 3 that with more measurements, the robust low-
 723 rank matrix recovery formulation (1.4) can tolerate not only more outliers but also
 724 a higher fraction of outliers. This provides further explanation to the observations
 725 made after the proof of Proposition 2.

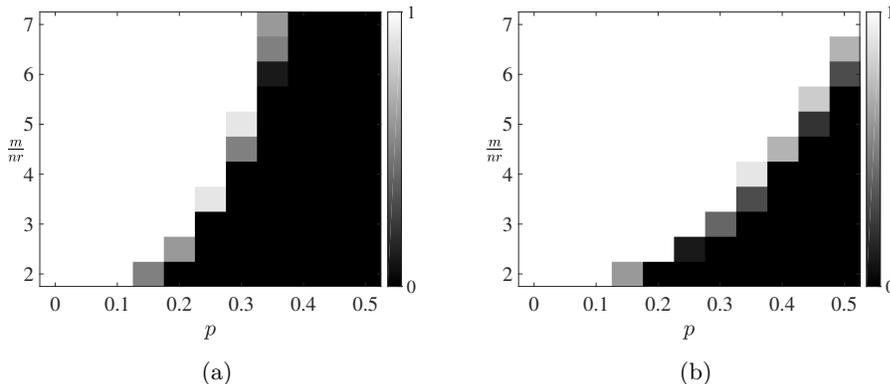


Fig. 3: Phase transition of robust low-rank matrix recovery using (a) median-truncated gradient descent (MTGD) [30] and (b) SubGM. Here, we fix $n = 50$, $r = 5$ and vary the outlier ratio p from 0 to 0.5. In addition, we vary m so that the ratio $\frac{m}{nr}$ varies from 2 to 7. Successful recovery is indicated by white and failure by black. Results are averaged over 10 independent trials.

726 **6. Conclusion.** In this paper we gave a nonsmooth nonconvex formulation of
 727 the problem of recovering a rank- r matrix $\mathbf{X}^* \in \mathbb{R}^{n_1 \times n_2}$ from corrupted linear mea-
 728 surements. The formulation enforces the low-rank property of the solution by using
 729 a factored representation of the matrix variable and employs an ℓ_1 -loss function to
 730 robustify the solution against outliers. We showed that even when close to half of
 731 the measurements are arbitrarily corrupted, as long as certain measurement opera-
 732 tors arising from the measurement model satisfy the ℓ_1/ℓ_2 -RIP, the formulation will
 733 be sharp and weakly convex. Consequently, the ground-truth matrix can be exactly
 734 recovered from any of its global minimum. Moreover, when suitably initialized, the
 735 SubGM with geometrically diminishing step sizes will converge to the ground-truth
 736 matrix at a linear rate.

737 **7. Acknowledgment.** We thank the Associate Editor and two anonymous re-
 738 viewers for their detailed and helpful comments.

- 740 [1] S. AARONSON, *The Learnability of Quantum States*, in Proceedings of the Royal Society of
741 London A: Mathematical, Physical and Engineering Sciences, vol. 463, 2007, pp. 3089–
742 3114.
- 743 [2] Y. BAI, Q. JIANG, AND J. SUN, *Subgradient Descent Learns Orthogonal Dictionaries*, Interna-
744 tional Conference on Learning Representations (ICLR), (2019).
- 745 [3] D. P. BERTSEKAS, *Incremental Gradient, Subgradient, and Proximal Methods for Convex Op-*
746 *timization*, in Optimization for Machine Learning, S. Sra, S. Nowozin, and S. J. Wright,
747 eds., Neural Information Processing Series, MIT Press, Cambridge, Massachusetts, 2012,
748 pp. 85–119.
- 749 [4] S. BHOJANAPALLI, B. NEYSHABUR, AND N. SREBRO, *Global Optimality of Local Search for Low*
750 *Rank Matrix Recovery*, in Advances in Neural Information Processing Systems 29 (NIPS),
751 D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds., 2016, pp. 3873–
752 3881.
- 753 [5] S. BURER AND R. D. MONTEIRO, *A Nonlinear Programming Algorithm for Solving Semidefinite*
754 *Programs via Low-Rank Factorization*, Mathematical Programming, 95 (2003), pp. 329–
755 357.
- 756 [6] S. BURER AND R. D. MONTEIRO, *Local Minima and Convergence in Low-Rank Semidefinite*
757 *Programming*, Mathematical Programming, 103 (2005), pp. 427–444.
- 758 [7] J. V. BURKE AND M. C. FERRIS, *Weak Sharp Minima in Mathematical Programming*, SIAM
759 Journal on Control and Optimization, 31 (1993), pp. 1340–1359.
- 760 [8] E. J. CANDÈS, X. LI, Y. MA, AND J. WRIGHT, *Robust Principal Component Analysis?*, Journal
761 of the ACM, 58 (2011), p. Article 11.
- 762 [9] E. J. CANDÈS AND Y. PLAN, *Tight Oracle Inequalities for Low-Rank Matrix Recovery from*
763 *a Minimal Number of Noisy Random Measurements*, IEEE Transactions on Information
764 Theory, 57 (2011), pp. 2342–2359.
- 765 [10] E. J. CANDÈS AND B. RECHT, *Exact Matrix Completion via Convex Optimization*, Foundations
766 of Computational Mathematics, 9 (2009), pp. 717–772.
- 767 [11] Y. CHEN, Y. CHI, AND A. J. GOLDSMITH, *Exact and Stable Covariance Estimation from Qua-*
768 *dratic Sampling via Convex Programming*, IEEE Transactions on Information Theory, 61
769 (2015), pp. 4034–4059.
- 770 [12] Y. CHI, Y. M. LU, AND Y. CHEN, *Nonconvex Optimization Meets Low-Rank Matrix Factori-*
771 *zation: An Overview*, arXiv preprint arXiv:1809.09573, (2018).
- 772 [13] M. A. DAVENPORT AND J. ROMBERG, *An Overview of Low-Rank Matrix Recovery from In-*
773 *complete Observations*, IEEE Journal of Selected Topics in Signal Processing, 10 (2016),
774 pp. 608–622.
- 775 [14] D. DAVIS, D. DRUSVYATSKIY, K. J. MACPHEE, AND C. PAQUETTE, *Subgradient Methods for*
776 *Sharp Weakly Convex Functions*, Journal of Optimization Theory and Applications, 179
777 (2018), pp. 962–982.
- 778 [15] D. DAVIS, D. DRUSVYATSKIY, AND C. PAQUETTE, *The Nonsmooth Landscape of Phase Retrieval*,
779 arXiv preprint arXiv:1711.03247, (2017).
- 780 [16] F. DE LA TORRE AND M. J. BLACK, *A Framework for Robust Subspace Learning*, International
781 Journal of Computer Vision, 54 (2003), pp. 117–142.
- 782 [17] J. C. DUCHI AND F. RUAN, *Solving (Most) of a Set of Quadratic Equalities: Composite Op-*
783 *timization for Robust Phase Retrieval*, Information and Inference: A Journal of the IMA,
784 (2018), p. iay015, <https://doi.org/10.1093/imaiai/iay015>.
- 785 [18] M. FAZEL, H. HINDI, AND S. BOYD, *Rank Minimization and Applications in System Theory*, in
786 Proceedings of the 2004 American Control Conference, vol. 4, IEEE, 2004, pp. 3273–3278.
- 787 [19] R. GE, J. D. LEE, AND T. MA, *Matrix Completion has No Spurious Local Minima*, in Advan-
788 ces in Neural Information Processing Systems, D. D. Lee, M. Sugiyama, U. V. Luxburg,
789 I. Guyon, and R. Garnett, eds., 2016, pp. 2973–2981.
- 790 [20] J.-L. GOFFIN, *On Convergence Rates of Subgradient Optimization Methods*, Mathematical
791 programming, 13 (1977), pp. 329–347.
- 792 [21] D. GROSS, *Recovering Low-Rank Matrices from Few Coefficients in Any Basis*, IEEE Tran-
793 sactions on Information Theory, 57 (2011), pp. 1548–1566.
- 794 [22] Q. GU, Z. W. WANG, AND H. LIU, *Low-Rank and Sparse Structure Pursuit via Alternating*
795 *Minimization*, in Proceedings of the 19th International Conference on Artificial Intelligence
796 and Statistics (AISTATS 2016), 2016, pp. 600–609.
- 797 [23] B. HAEFFELE, E. YOUNG, AND R. VIDAL, *Structured Low-Rank Matrix Factorization: Op-*
798 *timality, Algorithm, and Applications to Image Processing*, in Proceedings of the 31st
799 International Conference on Machine Learning (ICML 2014), 2014, pp. 2007–2015.

- 800 [24] C. JOSZ, Y. OUYANG, R. ZHANG, J. LAVAEI, AND S. SOJOUDI, *A Theory on the Absence of*
801 *Spurious Solutions for Nonconvex and Nonsmooth Optimization*, in Advances in Neural
802 Information Processing Systems 31 (NeurIPS), S. Bengio, H. Wallach, H. Larochelle,
803 K. Grauman, N. Cesa-Bianchi, and R. Garnett, eds., 2018, pp. 2441–2449.
- 804 [25] Q. KE AND T. KANADE, *Robust L_1 Norm Factorization in the Presence of Outliers and Missing*
805 *Data by Alternative Convex Programming*, in Proceedings of the 2005 IEEE Computer
806 Society Conference on Computer Vision and Pattern Recognition (CVPR 2005), vol. 1,
807 IEEE, 2005, pp. 739–746.
- 808 [26] L. LI, W. HUANG, I. Y.-H. GU, AND Q. TIAN, *Statistical Modeling of Complex Backgrounds for*
809 *Foreground Object Detection*, IEEE Transactions on Image Processing, 13 (2004), pp. 1459–
810 1472.
- 811 [27] Q. LI, Z. ZHU, AND G. TANG, *The Non-Convex Geometry of Low-Rank Matrix Optimization,*
812 *Information and Inference: A Journal of the IMA*, (2018), p. iay003, [https://doi.org/10.](https://doi.org/10.1093/imaiai/iay003)
813 [1093/imaiai/iay003](https://doi.org/10.1093/imaiai/iay003).
- 814 [28] X. LI, J. LU, R. ARORA, J. HAUPT, H. LIU, Z. WANG, AND T. ZHAO, *Symmetry, Saddle*
815 *Points, and Global Optimization Landscape of Nonconvex Matrix Factorization*, IEEE
816 Transactions on Information Theory, 65 (2019), pp. 3489–3514.
- 817 [29] X. LI, Z. ZHU, A. M.-C. SO, AND R. VIDAL, *Nonconvex Robust Low-Rank Matrix Recovery.*
818 Companion technical report, available at <https://arxiv.org/abs/1809.09237>, 2018.
- 819 [30] Y. LI, Y. CHI, H. ZHANG, AND Y. LIANG, *Nonconvex Low-Rank Matrix Recovery with Arbitrary*
820 *Outliers via Median-Truncated Gradient Descent*, Information and Inference: A Journal
821 of the IMA, (2019), p. iaz009, <https://doi.org/10.1093/imaiai/iaz009>.
- 822 [31] Y. LI, Y. SUN, AND Y. CHI, *Low-Rank Positive Semidefinite Matrix Recovery from Corrupted*
823 *Rank-One Measurements*, IEEE Transactions on Signal Processing, 65 (2017), pp. 397–408.
- 824 [32] A. NEDIĆ AND D. BERTSEKAS, *Convergence Rate of Incremental Subgradient Algorithms*, in
825 Stochastic Optimization: Algorithms and Applications, S. Uryasev and P. M. Pardalos,
826 eds., vol. 54 of Applied Optimization, Springer Science+Business Media, Dordrecht, 2001.
- 827 [33] P. NETRAPALLI, U. N. NIRANJAN, S. SANGHAVI, A. ANANDKUMAR, AND P. JAIN, *Non-Convex*
828 *Robust PCA*, in Advances in Neural Information Processing Systems 27 (NIPS), Z. Ghahra-
829 mani, M. Welling, C. Cortes, N. D. Lawrence, and K. Q. Weinberger, eds., 2014, pp. 1107–
830 1115.
- 831 [34] J. NOCEDAL AND S. WRIGHT, *Numerical optimization*, Springer Science & Business Media,
832 2006.
- 833 [35] D. PARK, A. KYRILLIDIS, C. CARAMANIS, AND S. SANGHAVI, *Non-Square Matrix Sensing with-*
834 *out Spurious Local Minima via the Burer-Monteiro Approach*, in Proceedings of the 20th
835 International Conference on Artificial Intelligence and Statistics (AISTATS 2017), 2017,
836 pp. 65–74.
- 837 [36] B. T. POLYAK, *Minimization of Unsmooth Functions*, USSR Computational Mathematics and
838 Mathematical Physics, 9 (1969), pp. 14–29.
- 839 [37] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317 of Grundlehren der
840 mathematischen Wissenschaften, Springer-Verlag, Berlin Heidelberg, second ed., 2004.
- 841 [38] N. Z. SHOR, *Minimization Methods for Non-Differentiable Functions*, vol. 3 of Springer Series
842 in Computational Mathematics, Springer-Verlag, Berlin Heidelberg, 1985.
- 843 [39] N. Z. SHOR AND M. B. SHCHEPAKIN, *Algorithms for the Solution of the Two-Stage Problem in*
844 *Stochastic Programming*, Kibernetika, 4 (1968), pp. 56–58.
- 845 [40] N. SREBRO, J. RENNIE, AND T. S. JAAKKOLA, *Maximum-Margin Matrix Factorization*, in
846 Advances in Neural Information Processing Systems 17 (NIPS), L. K. Saul, Y. Weiss, and
847 L. Bottou, eds., 2004, pp. 1329–1336.
- 848 [41] R. SUN AND Z.-Q. LUO, *Guaranteed Matrix Completion via Non-Convex Factorization*, IEEE
849 Transactions on Information Theory, 62 (2016), pp. 6535–6579.
- 850 [42] S. TU, R. BO CZAR, M. SIMCHOWITZ, M. SOLTANOLKOTABI, AND B. RECHT, *Low-Rank Solutions*
851 *of Linear Matrix Equations via Procrustes Flow*, in Proceedings of the 33rd International
852 Conference on Machine Learning (ICML 2016), 2016, pp. 964–973.
- 853 [43] R. VERSHYNIN, *Introduction to the Non-Asymptotic Analysis of Random Matrices*, in Com-
854 pressed Sensing: Theory and Applications, Y. C. Eldar and G. Kutyniok, eds., Cambridge
855 University Press, New York, 2012, pp. 210–268.
- 856 [44] J.-P. VIAL, *Strong and Weak Convexity of Sets and Functions*, Mathematics of Operations
857 Research, 8 (1983), pp. 231–259.
- 858 [45] X. YI, D. PARK, Y. CHEN, AND C. CARAMANIS, *Fast Algorithms for Robust PCA via Gradient*
859 *Descent*, in Advances in Neural Information Processing Systems 29 (NIPS), D. D. Lee,
860 M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, eds., 2016, pp. 4152–4160.
- 861 [46] M.-C. YUE AND A. M.-C. SO, *A Perturbation Inequality for Concave Functions of Singular*

- 862 *Values and Its Applications in Low-Rank Matrix Recovery*, Applied and Computational
863 Harmonic Analysis, 40 (2016), pp. 396–416.
- 864 [47] M.-C. YUE, Z. ZHOU, AND A. M.-C. SO, *On the Quadratic Convergence of the Cubic Regu-*
865 *larization Method under a Local Error Bound Condition*, SIAM Journal on Optimization,
866 29 (2019), pp. 904–932.
- 867 [48] M. ZHANG, Z.-H. HUANG, AND Y. ZHANG, *Restricted p -Isometry Properties of Nonconvex*
868 *Matrix Recovery*, IEEE Transactions on Information Theory, 59 (2013), pp. 4316–4323.
- 869 [49] X. ZHANG, L. WANG, AND Q. GU, *A Unified Framework for Nonconvex Low-Rank plus Sparse*
870 *Matrix Recovery*, in Proceedings of the 21st International Conference on Artificial Intelli-
871 gence and Statistics (AISTATS 2018), 2018, pp. 1097–1107.
- 872 [50] Q. ZHENG AND J. LAFFERTY, *A Convergent Gradient Descent Algorithm for Rank Minimization*
873 *and Semidefinite Programming from Random Linear Measurements*, in Advances in
874 Neural Information Processing Systems 28 (NIPS), C. Cortes, N. D. Lawrence, D. D. Lee,
875 M. Sugiyama, and R. Garnett, eds., 2015, pp. 109–117.
- 876 [51] Z. ZHU, Q. LI, G. TANG, AND M. B. WAKIN, *The Global Optimization Geometry of Low-Rank*
877 *Matrix Optimization*, arXiv preprint arXiv:1703.01256, (2017).
- 878 [52] Z. ZHU, Q. LI, G. TANG, AND M. B. WAKIN, *Global Optimality in Low-Rank Matrix Optimi-*
879 *zation*, IEEE Transactions on Signal Processing, 66 (2018), pp. 3614–3628.
- 880 [53] Z. ZHU, A. M.-C. SO, AND Y. YE, *Fast and Near-Optimal Matrix Completion via Randomized*
881 *Basis Pursuit*, in Fifth International Congress of Chinese Mathematicians, L. Ji, Y. S.
882 Poon, L. Yang, and S.-T. Yau, eds., vol. 51, Part 2 of AMS/IP Studies in Advanced
883 Mathematics, American Mathematical Society and International Press, 2012, pp. 859–882.
- 884 [54] Z. ZHU, Y. WANG, D. ROBINSON, D. NAIMAN, R. VIDAL, AND M. TSAKIRIS, *Dual Principal*
885 *Component Pursuit: Improved Analysis and Efficient Algorithms*, in Advances in Neu-
886 ral Information Processing Systems 31 (NeurIPS), S. Bengio, H. Wallach, H. Larochelle,
887 K. Grauman, N. Cesa-Bianchi, and R. Garnett, eds., 2018, pp. 2171–2181.

888 Appendix A. Proof of Proposition 1.

889 **A.1. Preliminaries.** We say that a random variable X is sub-Gaussian if $\Pr[|X| > t] \leq$

890 $\exp\left(1 - \frac{t^2}{K_1^2}\right)$, $\forall t \geq 0$ for some constant $K_1 > 0$. This is equivalent to

$$891 \quad (\text{A.1}) \quad (\mathbb{E}[|X|^p])^{1/p} \leq K_2 \sqrt{p}, \quad \forall p \geq 1$$

892 for some constant $K_2 > 0$. The constants K_1 and K_2 differ from each other by at
893 most an absolute constant factor; see [43, Lemma 5.5]. The sub-Gaussian norm of a
894 sub-Gaussian random variable X is defined as $\|X\|_{\psi_2} = \sup_{p \geq 1} \{p^{-1/2} \mathbb{E}[|X|^p]^{1/p}\}$.
895 We then have the following Hoeffding-type inequality:

896 **LEMMA 2** ([43, Proposition 5.10]). *Let X_1, \dots, X_m be independent sub-Gaussian*
897 *random variables with $\mathbb{E}[X_i] = 0$ for $i = 1, \dots, m$ and $K = \max_{i \in \{1, \dots, m\}} \|X_i\|_{\psi_2}$.*
898 *Then, for any $t > 0$, we have*

$$899 \quad (\text{A.2}) \quad \Pr\left[\frac{1}{m} \left| \sum_{i=1}^m X_i \right| > t\right] \leq 2 \exp\left(-\frac{cmt^2}{K^2}\right)$$

900 for some constant $c > 0$.

901 We also need the following result on the covering number of the set of low-rank
902 matrices:

903 **LEMMA 3** ([9, Lemma 3.1]). *Let $\mathbb{S}_r = \{\mathbf{X} \in \mathbb{R}^{n \times n} : \|\mathbf{X}\|_F = 1, \text{rank}(\mathbf{X}) \leq r\}$.*
904 *Then, there exists an ϵ -net $\overline{\mathbb{S}}_{r,\epsilon} \subset \mathbb{S}_r$ with respect to the Frobenius norm (i.e., for*
905 *any $\mathbf{X} \in \mathbb{S}_r$, there exists an $\overline{\mathbf{X}} \in \overline{\mathbb{S}}_{r,\epsilon}$ such that $\|\mathbf{X} - \overline{\mathbf{X}}\|_F \leq \epsilon$) satisfying $|\overline{\mathbb{S}}_{r,\epsilon}| \leq$
906 $\left(\frac{9}{\epsilon}\right)^{(2n+1)r}$.*

907 A.2. Isometry Property of a Given Matrix.

908 LEMMA 4. Suppose that the matrices $\mathbf{A}_1, \dots, \mathbf{A}_m \in \mathbb{R}^{n \times n}$ defining the linear mea-
 909 surement operator \mathcal{A} have i.i.d. standard Gaussian entries. Then, for any $\mathbf{X} \in \mathbb{R}^{n \times n}$
 910 and $0 < \delta < 1$, there exists a constant $c_1 > 0$ such that with probability exceeding
 911 $1 - 2 \exp(-c_1 \delta^2 m)$, we have

$$912 \quad (\text{A.3}) \quad \left(\sqrt{\frac{2}{\pi}} - \delta \right) \|\mathbf{X}\|_F \leq \frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 \leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{X}\|_F.$$

913 *Proof of Lemma 4.* Since \mathbf{A}_i has i.i.d. standard Gaussian entries, the random
 914 variable $\langle \mathbf{A}_i, \mathbf{X} \rangle$ is Gaussian with mean zero and variance $\|\mathbf{X}\|_F^2$. It follows that

$$915 \quad (\text{A.4}) \quad \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle] = \sqrt{\frac{2}{\pi}} \|\mathbf{X}\|_F, \quad \mathbb{E}[\|\mathcal{A}(\mathbf{X})\|_1] = m \sqrt{\frac{2}{\pi}} \|\mathbf{X}\|_F.$$

916 Now, let $Z_i = |\langle \mathbf{A}_i, \mathbf{X} \rangle| - \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]$, which satisfies $\mathbb{E}[Z_i] = 0$. We claim that Z_i
 917 is a sub-Gaussian random variable. To establish the claim, it suffices to bound the
 918 sub-Gaussian norm of Z_i . Towards that end, we first observe that $\Pr[|\langle \mathbf{A}_i, \mathbf{X} \rangle| > t] \leq$
 919 $2 \exp\left(-\frac{t^2}{2\|\mathbf{X}\|_F^2}\right)$. Together with (A.4), this implies that for any $t > \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]$,

$$920 \quad \Pr[|Z_i| > t] = \Pr[|\langle \mathbf{A}_i, \mathbf{X} \rangle| > t + \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]] + \Pr[|\langle \mathbf{A}_i, \mathbf{X} \rangle| < -t + \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]] \\
 921 \quad \leq 2 \exp\left(-\frac{(t + \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle])^2}{2\|\mathbf{X}\|_F^2}\right) + \Pr[|\langle \mathbf{A}_i, \mathbf{X} \rangle| < -t + \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]] \\
 922 \quad \leq 2 \exp\left(-\frac{(t + \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle])^2}{2\|\mathbf{X}\|_F^2}\right) \leq \exp\left(1 - \frac{t^2}{\|\mathbf{X}\|_F^2}\right), \\
 923$$

924 where the second inequality follows because $\Pr[|\langle \mathbf{A}_i, \mathbf{X} \rangle| < -t + \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]] = 0$ for
 925 all $t > \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle]$. Since $\exp\left(1 - \frac{t^2}{\|\mathbf{X}\|_F^2}\right) \geq 1$ for all $t \leq \mathbb{E}[\langle \mathbf{A}_i, \mathbf{X} \rangle] = \sqrt{\frac{2}{\pi}} \|\mathbf{X}\|_F$,
 926 we then have $\Pr[|Z_i| > t] \leq \exp\left(1 - \frac{t^2}{\|\mathbf{X}\|_F^2}\right)$, $\forall t \geq 0$. This, together with (A.1),
 927 implies that $(\mathbb{E}[|Z_i|^p])^{1/p} \leq c p^{1/2} \|\mathbf{X}\|_F$, $\forall p \geq 1$, where $c > 0$ is a constant. It follows
 928 that $\|Z_i\|_{\psi_2} \leq c \|\mathbf{X}\|_F$; i.e., Z_i is a sub-Gaussian random variable, as desired.

929 Now, applying the Hoeffding-type inequality in Lemma 2 with $t = \delta \|\mathbf{X}\|_F$ and
 930 $K = c \|\mathbf{X}\|_F$ gives

$$931 \quad \Pr\left[\frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 - \mathbb{E}[\|\mathcal{A}(\mathbf{X})\|_1] > \delta \|\mathbf{X}\|_F\right] \leq 2 \exp(-c_1 m \delta^2)$$

932 for some constant $c_1 > 0$. Using (A.4), we conclude that (A.3) holds with probability
 933 at least $1 - 2 \exp(-c_1 m \delta^2)$. This completes the proof. \square

934 **A.3. Proof of Proposition 1.** We now utilize an ϵ -net argument to show that
 935 (A.3) holds for all rank- r matrices with high probability as long as $m \gtrsim nr$. Since
 936 the inequality (A.3) is scale invariant, without loss of generality, we may assume that
 937 $\|\mathbf{X}\|_F = 1$ and focus on the set \mathbb{S}_r defined in Lemma 3.

938 *Proof of Proposition 1.* We begin by showing that (A.3) holds for all $\mathbf{X} \in \overline{\mathbb{S}}_{r, \epsilon}$
 939 with high probability. Indeed, upon setting $\epsilon = \frac{\delta \sqrt{\pi}}{16}$ in Lemma 3 and utilizing a union

940 bound together with Lemma 4, we have

$$\begin{aligned}
 941 \quad (\text{A.5}) \quad & \Pr \left[\max_{\bar{\mathbf{X}} \in \bar{\mathbb{S}}_{r,\epsilon}} \frac{1}{m} \left| \|\mathcal{A}(\bar{\mathbf{X}})\|_1 - m \sqrt{\frac{2}{\pi}} \|\bar{\mathbf{X}}\|_F \right| \geq \frac{\delta}{2} \right] \leq 2|\bar{\mathbb{S}}_{r,\epsilon}| \exp(-c_1 m \delta^2) \\
 & \leq 2 \left(\frac{9}{\epsilon} \right)^{(2n+1)r} \exp(-c_1 m \delta^2) \leq \exp(-c_2 m \delta^2)
 \end{aligned}$$

942 whenever $m \gtrsim nr$.

943 Next, we show that (A.3) holds for all $\mathbf{X} \in \mathbb{S}_r$. Towards that end, set

$$944 \quad (\text{A.6}) \quad \kappa_r = \frac{1}{m} \sup_{\mathbf{X} \in \mathbb{S}_r} \|\mathcal{A}(\mathbf{X})\|_1$$

945 and let $\mathbf{X} \in \mathbb{S}_r$ be arbitrary. Then, there exists an $\bar{\mathbf{X}} \in \bar{\mathbb{S}}_{r,\epsilon}$ such that $\|\mathbf{X} - \bar{\mathbf{X}}\|_F \leq \epsilon$.

946 It follows from (A.5) that with high probability,

$$\begin{aligned}
 947 \quad (\text{A.7}) \quad & \frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 = \frac{1}{m} \|\mathcal{A}(\mathbf{X} - \bar{\mathbf{X}}) + \mathcal{A}(\bar{\mathbf{X}})\|_1 \leq \frac{1}{m} \|\mathcal{A}(\mathbf{X} - \bar{\mathbf{X}})\|_1 + \frac{1}{m} \|\mathcal{A}(\bar{\mathbf{X}})\|_1 \\
 & \leq \frac{1}{m} \|\mathcal{A}(\mathbf{X} - \bar{\mathbf{X}})\|_1 + \sqrt{\frac{2}{\pi}} + \frac{\delta}{2}.
 \end{aligned}$$

948 Noting that $\mathbf{X} - \bar{\mathbf{X}}$ has rank at most $2r$, we can decompose it as $\mathbf{X} - \bar{\mathbf{X}} = \mathbf{\Delta}_1 + \mathbf{\Delta}_2$,
 949 where $\langle \mathbf{\Delta}_1, \mathbf{\Delta}_2 \rangle = 0$ and $\text{rank}(\mathbf{\Delta}_1), \text{rank}(\mathbf{\Delta}_2) \leq r$ (this follows essentially from the
 950 SVD). Hence, we can compute

$$\begin{aligned}
 & \frac{1}{m} \|\mathcal{A}(\mathbf{X} - \bar{\mathbf{X}})\|_1 \leq \frac{1}{m} [\|\mathcal{A}(\mathbf{\Delta}_1)\|_1 + \|\mathcal{A}(\mathbf{\Delta}_2)\|_1] \\
 951 \quad & = \frac{1}{m} [\|\mathbf{\Delta}_1\|_F \|\mathcal{A}(\mathbf{\Delta}_1/\|\mathbf{\Delta}_1\|_F)\|_1 + \|\mathbf{\Delta}_2\|_F \|\mathcal{A}(\mathbf{\Delta}_2/\|\mathbf{\Delta}_2\|_F)\|_1] \\
 & \leq \kappa_r (\|\mathbf{\Delta}_1\|_F + \|\mathbf{\Delta}_2\|_F) \leq \sqrt{2} \kappa_r \epsilon,
 \end{aligned}$$

952 where the last inequality is due to $\|\mathbf{\Delta}_1\|_F^2 + \|\mathbf{\Delta}_2\|_F^2 = \|\mathbf{X} - \bar{\mathbf{X}}\|_F^2 \leq \epsilon^2$. This, together
 953 with (A.7), gives

$$954 \quad (\text{A.8}) \quad \frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 \leq \sqrt{\frac{2}{\pi}} + \frac{\delta}{2} + \sqrt{2} \kappa_r \epsilon.$$

955 In particular, using the definition of κ_r in (A.6), we obtain $\kappa_r \leq \sqrt{\frac{2}{\pi}} + \frac{\delta}{2} + \sqrt{2} \kappa_r \epsilon$,

956 or equivalently, $\kappa_r \leq \frac{\sqrt{2/\pi + \delta/2}}{1 - \sqrt{2}\epsilon}$. Plugging in our choice of ϵ yields $\sqrt{2} \kappa_r \epsilon \leq \frac{\delta}{2}$. This,
 957 together with (A.8) and the fact that $\|\mathbf{X}\|_F = 1$, implies

$$958 \quad \frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 \leq \left(\sqrt{\frac{2}{\pi}} + \delta \right) \|\mathbf{X}\|_F.$$

959 Similarly, using (A.5), we have

$$\begin{aligned}
 960 \quad & \frac{1}{m} \|\mathcal{A}(\mathbf{X})\|_1 \geq \frac{1}{m} \|\mathcal{A}(\bar{\mathbf{X}})\|_1 - \frac{1}{m} \|\mathcal{A}(\mathbf{X} - \bar{\mathbf{X}})\|_1 \\
 961 \quad & \geq \sqrt{\frac{2}{\pi}} - \frac{\delta}{2} - \frac{1}{m} \|\mathcal{A}(\mathbf{X} - \bar{\mathbf{X}})\|_1 \\
 962 \quad & \geq \sqrt{\frac{2}{\pi}} - \frac{\delta}{2} - \sqrt{2} \kappa_r \epsilon \geq \sqrt{\frac{2}{\pi}} - \delta = \left(\sqrt{\frac{2}{\pi}} - \delta \right) \|\mathbf{X}\|_F \\
 963 \quad &
 \end{aligned}$$

964 with high probability. This completes the proof. \square