

1 **WEAKLY CONVEX OPTIMIZATION OVER STIEFEL MANIFOLD**  
 2 **USING RIEMANNIAN SUBGRADIENT-TYPE METHODS\***

3 XIAO LI<sup>†</sup>, SHIXIANG CHEN<sup>‡</sup>, ZENGDE DENG<sup>§</sup>, QING QU<sup>¶</sup>, ZHIHUI ZHU<sup>||</sup>, AND  
 4 ANTHONY MAN-CHO SO<sup>#</sup>

5 **Abstract.** We consider a class of nonsmooth optimization problems over the Stiefel manifold,  
 6 in which the objective function is weakly convex in the ambient Euclidean space. Such problems  
 7 are ubiquitous in engineering applications but still largely unexplored. We present a family of Rie-  
 8 mannian subgradient-type methods—namely Riemannian subgradient, incremental subgradient, and  
 9 stochastic subgradient methods—to solve these problems and show that they all have an iteration  
 10 complexity of  $\mathcal{O}(\varepsilon^{-4})$  for driving a natural stationarity measure below  $\varepsilon$ . In addition, we establish  
 11 the local linear convergence of the Riemannian subgradient and incremental subgradient methods  
 12 when the problem at hand further satisfies a sharpness property and the algorithms are properly  
 13 initialized and use geometrically diminishing stepsizes. To the best of our knowledge, these are the  
 14 first convergence guarantees for using Riemannian subgradient-type methods to optimize a class of  
 15 nonconvex nonsmooth functions over the Stiefel manifold. The fundamental ingredient in the proof  
 16 of the aforementioned convergence results is a new *Riemannian subgradient inequality* for restrictions  
 17 of weakly convex functions on the Stiefel manifold, which could be of independent interest. We also  
 18 show that our convergence results can be extended to handle a class of compact embedded subman-  
 19 ifolds of the Euclidean space. Finally, we discuss the sharpness properties of various formulations  
 20 of the robust subspace recovery and orthogonal dictionary learning problems and demonstrate the  
 21 convergence performance of the algorithms on both problems via numerical simulations.

22 **Key words.** manifold optimization, nonconvex optimization, orthogonality constraint, iteration  
 23 complexity, linear convergence, robust subspace recovery, dictionary learning

24 **AMS subject classifications.** 68Q25, 65K10, 90C90, 90C26, 90C06.

25 **1. Introduction.** In this paper, we consider the problem of optimizing a func-  
 26 tion with finite-sum structure over the Stiefel manifold—i.e.,

27 (1.1) 
$$\begin{aligned} \text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) &:= \frac{1}{m} \sum_{i=1}^m f_i(\mathbf{X}) \\ \text{subject to } \mathbf{X} &\in \text{St}(n, r) \end{aligned}$$

28 with  $\text{St}(n, r) := \{\mathbf{X} \in \mathbb{R}^{n \times r} : \mathbf{X}^\top \mathbf{X} = \mathbf{I}_r\}$  and  $\mathbf{I}_r$  being the  $r \times r$  identity matrix—  
 29 where each component  $f_i : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$  ( $i = 1, \dots, m$ ) is assumed to be weakly convex

---

\*Submitted to the editors March 24, 2021. The first and second authors contributed equally to this paper. Most of the work of the first author was done when he was affiliated with the Department of Electronic Engineering, The Chinese University of Hong Kong.

**Funding:** X. Li was partially supported by the University Development Fund UDF01001808 of CUHK (SZ). Q. Qu was partially supported by the Moore-Sloan fellowship. Z. Zhu was partially supported by NSF Grant 1704458 and NSF Grant CCF-2008460. A. M.-C. So was partially supported by the Hong Kong Research Grants Council (RGC) General Research Fund (GRF) Project CUHK 14208117 and the CUHK Research Sustainability of Major RGC Funding Schemes Project 3133236.

<sup>†</sup>Corresponding author. School of Data Science, The Chinese University of Hong Kong, Shenzhen. (lixiao@cuhk.edu.cn, <https://sites.google.com/view/xli>).

<sup>‡</sup>Department of Industrial and Systems Engineering, Texas A&M University. (sxchen@tamu.edu).

<sup>§</sup>Cainiao Network, Hangzhou, China. (dengzengde@gmail.com).

<sup>¶</sup>Department of Electrical Engineering and Computer Science, University of Michigan. (qingqu@umich.edu, <https://qingqu.engin.umich.edu/>).

<sup>||</sup>Department of Electrical and Computer Engineering, University of Denver. (zihui.zhu@du.edu, <http://mysite.du.edu/~zzhu61/>).

<sup>#</sup>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>).

30 in the ambient Euclidean space  $\mathbb{R}^{n \times r}$ . Recall that a function  $h$  is said to be *weakly*  
 31 *convex* if  $h(\cdot) + \frac{\tau}{2} \|\cdot\|_2^2$  is convex for some constant  $\tau \geq 0$  [55]. In particular, the  
 32 objective function in (1.1) can be nonconvex and nonsmooth. Our interest in (1.1)  
 33 stems from the fact that it arises in many applications from different engineering fields  
 34 such as representation learning and imaging science. As an illustration, let us present  
 35 two motivating applications, in which nonsmooth formulations have clear advantages  
 36 over smooth ones.

### 37 1.1. Motivating applications.

38 **Application 1: Robust subspace recovery.** Fitting a linear subspace to a  
 39 dataset corrupted by outliers is a fundamental problem in machine learning and statis-  
 40 tics, primarily known as robust principal component analysis (RPCA) [56] or robust  
 41 subspace recovery (RSR) [33]. In this problem, one is given measurements  $\tilde{\mathbf{Y}}$  of the  
 42 form  $\tilde{\mathbf{Y}} = [\mathbf{Y} \ \mathbf{O}] \mathbf{\Gamma} \in \mathbb{R}^{n \times m}$ , where the columns of  $\mathbf{Y} \in \mathbb{R}^{n \times m_1}$  form inlier points  
 43 spanning a  $d$ -dimensional subspace  $\mathcal{S}$ ; the columns of  $\mathbf{O} \in \mathbb{R}^{n \times m_2}$  form outlier points  
 44 with no linear structure;  $\mathbf{\Gamma} \in \mathbb{R}^{m \times m}$  is an unknown permutation, and the goal is to  
 45 recover the subspace  $\mathcal{S}$ . It is well-known that the presence of outliers can severely  
 46 affect the quality of the solutions obtained by the classic PCA approach, which in-  
 47 volves minimizing a smooth least-squares loss [56]. In order to obtain solutions that  
 48 are more robust against outliers, the recent works [33, 34, 40] propose to minimize the  
 49 nonsmooth least absolute deviation (LAD) loss. This leads to the formulation

$$50 \quad (1.2) \quad \begin{aligned} & \underset{\mathbf{X} \in \mathbb{R}^{n \times d}}{\text{minimize}} \quad f(\mathbf{X}) := \frac{1}{m} \sum_{i=1}^m \|(\mathbf{I}_n - \mathbf{X} \mathbf{X}^\top) \tilde{\mathbf{y}}_i\|_2 \\ & \text{subject to} \quad \mathbf{X} \in \text{St}(n, d), \end{aligned}$$

51 where  $\tilde{\mathbf{y}}_i \in \mathbb{R}^n$  ( $i = 1, \dots, m$ ) denotes the  $i$ -th column of  $\tilde{\mathbf{Y}}$  and the columns of a global  
 52 minimizer of (1.2) are expected to form an orthonormal basis of the subspace  $\mathcal{S}$ . The  
 53 weak convexity of the components of the objective function in (1.2) can be verified by  
 54 following the arguments in the proof of [37, Proposition 6]. Thus, the formulation (1.2)  
 55 is an instance of problem (1.1). On another front, the works [54, 69, 70] consider a dual  
 56 form of the problem, which leads to the so-called dual principal component pursuit  
 57 (DPCP) formulation:

$$58 \quad (1.3) \quad \begin{aligned} & \underset{\mathbf{X} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}) := \frac{1}{m} \sum_{i=1}^m \|\tilde{\mathbf{y}}_i^\top \mathbf{X}\|_2 \\ & \text{subject to} \quad \mathbf{X} \in \text{St}(n, r). \end{aligned}$$

59 In contrast to the primal formulation (1.2), the dual formulation (1.3) aims to find an  
 60 orthogonal basis of  $\mathcal{S}^\perp$  (the orthogonal complement to  $\mathcal{S}$ ) with dimension  $r = n - d$ .  
 61 It is clear that the components of the objective function in (1.3) are convex, thus  
 62 showing that the formulation (1.3) is also an instance of problem (1.1).

63 **Application 2: Learning sparsely-used dictionaries.** A problem that arises  
 64 in many machine learning and computer vision applications is dictionary learning  
 65 (DL), whose goal is to find a suitable compact representation of certain input data  
 66  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_m] \in \mathbb{R}^{n \times m}$  [39, 47, 61]. Informally, this entails factorizing the data  
 67  $\mathbf{Y}$  into a dictionary  $\mathbf{A}$  and a sparse code matrix  $\mathbf{S} = [\mathbf{s}_1, \dots, \mathbf{s}_m]$ ; i.e.,  $\mathbf{Y} \approx \mathbf{A} \mathbf{S}$ .  
 68 When the dictionary  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is orthogonal and the code matrix  $\mathbf{S} \in \mathbb{R}^{n \times m}$  is

69 sufficiently sparse, the product  $\mathbf{A}^\top \mathbf{Y} \approx \mathbf{S}$  should be sparse. Thus, one may approach  
 70 the problem by finding the sparsest vectors in the row space of  $\mathbf{Y}$  [45, 50, 52]. This  
 71 motivates the following formulation [3]:

$$72 \quad (1.4) \quad \begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) := \frac{1}{m} \|\mathbf{Y}^\top \mathbf{x}\|_1 = \frac{1}{m} \sum_{i=1}^m |\mathbf{y}_i^\top \mathbf{x}| \\ & \text{subject to} \quad \mathbf{x} \in \text{St}(n, 1). \end{aligned}$$

73 Note that the solution to (1.4) only returns one column of  $\mathbf{A}$ . Thus, some extra  
 74 refinement technique, such as deflation [53] or repetitive independent trials [3], is  
 75 needed to fully solve the DL problem. It has been shown in [3] that under a suitable  
 76 statistical model, the formulation (1.4) requires fewer samples for exact recovery of the  
 77 dictionary  $\mathbf{A}$  than the smooth variant considered in [52, 53]. Still, since the approach  
 78 based on (1.4) recovers the columns of  $\mathbf{A}$  one at a time, it can be rather sensitive to  
 79 noise. To circumvent this difficulty, one possibility is to directly recover the orthogonal  
 80 dictionary  $\mathbf{A}$  by

$$81 \quad (1.5) \quad \begin{aligned} & \underset{\mathbf{X} \in \mathbb{R}^{n \times n}}{\text{minimize}} \quad f(\mathbf{X}) := \frac{1}{m} \|\mathbf{Y}^\top \mathbf{X}\|_1 = \frac{1}{m} \sum_{i=1}^m \|\mathbf{y}_i^\top \mathbf{X}\|_1 \\ & \text{subject to} \quad \mathbf{X} \in \text{St}(n, n); \end{aligned}$$

82 cf. [59, 65]. This approach can be easily extended to handle any complete (i.e., square  
 83 and invertible) dictionaries via preconditioning [52, 65]. Clearly, both (1.4) and (1.5)  
 84 are instances of (1.1).

85 **1.2. Main contributions.** We study three *Riemannian subgradient-type meth-*  
 86 *ods* for solving problem (1.1), namely Riemannian subgradient method, Riemannian  
 87 incremental subgradient method, and Riemannian stochastic subgradient method (see  
 88 Subsection 2.2). To analyze the convergence behavior of these methods, we first ex-  
 89 tend the surrogate stationarity measure developed in [12, 16] for weakly convex min-  
 90 imization in the Euclidean space to one for weakly convex minimization over the  
 91 Stiefel manifold (see Subsection 4.1). Then, we show that the iterates generated by  
 92 the aforementioned Riemannian subgradient-type methods will drive the surrogate  
 93 stationarity measure to zero at a rate of  $\mathcal{O}(k^{-\frac{1}{4}})$ , where  $k$  is the iteration index  
 94 (see Subsections 4.2 and 4.3). Such a complexity guarantee matches that established  
 95 in [12] for a host of algorithms that solve weakly convex minimization problems in the  
 96 *Euclidean space*. Next, we show that if problem (1.1) further satisfies the *sharpness*  
 97 property (see Definition 1), then the Riemannian subgradient and incremental sub-  
 98 gradient methods with properly designed geometrically diminishing stepsizes and a  
 99 good initialization will converge to the set of local minima associated with the sharp-  
 100 ness property at a *linear* rate (see Section 5). To the best of our knowledge, our  
 101 work is the *first* to establish the iteration complexities and convergence rates of Rie-  
 102 mannian subgradient-type methods for optimizing a class of nonconvex nonsmooth  
 103 functions over the Stiefel manifold. We also extend the above convergence results to  
 104 the setting where the constraint is a compact embedded submanifold of the Euclidean  
 105 space (see Section 6). Lastly, we show that under certain conditions on the inlier  
 106 and outlier distributions, the LAD (1.2) and DPCP (1.3) formulations of the RSR  
 107 problem satisfy the sharpness property (see Subsection 7.1). Consequently, we are  
 108 able to obtain recovery guarantees for the so-called Haystack model of the input data  
 109 that are competitive with state-of-the-art results.

110 The key to establishing the aforementioned convergence results is an algorithm-  
 111 independent property that we discovered for restrictions of weakly convex functions  
 112 on the Stiefel manifold, which we term the *Riemannian subgradient inequality* (see  
 113 [Section 3](#)). This is one of the main contributions of this work and could be of inde-  
 114 pendent interest for other Riemannian optimization problems. We believe that our  
 115 results will have broad implications on understanding the convergence behavior of  
 116 algorithms for solving more general manifold optimization problems with nonsmooth  
 117 objectives.

### 118 1.3. Connections with prior arts.

119 **Nonsmooth optimization in Euclidean space.** The problem of minimizing  
 120 a weakly convex function over a convex constraint set is well studied in the literature.  
 121 The main algorithms for this task include subgradient-type methods [\[12, 13, 36\]](#) and  
 122 proximal point-type methods [\[15\]](#). The convergence analyses of these algorithms rely  
 123 on a certain *weakly convex inequality*. We extend this line of work by considering  
 124 a nonconvex constraint set—i.e., the Stiefel manifold—and develop an analog of the  
 125 weakly convex inequality on the Stiefel manifold called the *Riemannian subgradient*  
 126 *inequality*. Such an inequality allows us to resort to the analysis techniques for weakly  
 127 convex minimization in the Euclidean space and prove new convergence results for  
 128 our Riemannian subgradient-type methods when solving the problem of weakly convex  
 129 minimization over the Stiefel manifold [\(1.1\)](#).

130 **Smooth optimization over Riemannian manifold.** Riemannian smooth op-  
 131 timization has been extensively studied over the years; see, e.g., [\[2, 7, 24, 26, 38\]](#) and  
 132 the references therein. Recently, global sublinear convergence results for Riemannian  
 133 gradient descent and Riemannian trust region have been presented in [\[7\]](#). The analysis  
 134 relies on the assumption that the pullback of the objective function  $f$  to the tangent  
 135 spaces of the manifold has a Lipschitz continuous gradient, which allows one to follow  
 136 the analyses of the corresponding methods for unconstrained smooth optimization.  
 137 However, such an approach breaks down when  $f$  is nonsmooth, as the gradient of the  
 138 pullback of  $f$  may not exist.

139 **Nonsmooth optimization over Riemannian manifold.** In contrast to Rie-  
 140 mannian smooth optimization, Riemannian nonsmooth optimization is relatively less  
 141 explored [\[1\]](#). In the following, we briefly review some state-of-the-art results in this  
 142 area and explain their limitations and connections to our results.

143 *Riemannian nonsmooth optimization with geodesic convexity.* Recently, the works  
 144 [\[4, 18, 19, 66\]](#) study the convergence behavior of Riemannian subgradient-type meth-  
 145 ods when the objective function is geodesically convex over a Riemannian manifold.  
 146 Thanks to the availability of a geodesic version of the convex subgradient inequality,  
 147 the conventional analysis for convex optimization in the Euclidean space can be carried  
 148 over to geodesically convex optimization over a Riemannian manifold. In particular,  
 149 an asymptotic convergence result is first established in [\[19\]](#), while a global convergence  
 150 rate of  $\mathcal{O}(k^{-\frac{1}{2}})$  is established in [\[4, 18\]](#), for the Riemannian subgradient method. The  
 151 work [\[66\]](#) considers the setting where the objective function is geodesically strongly  
 152 convex over the Riemannian manifold and shows that the rate can be improved to  
 153  $\mathcal{O}(k^{-1})$  for Riemannian projected subgradient methods. Unfortunately, these results  
 154 are not useful for understanding problem [\(1.1\)](#). This is because the constraint in [\(1.1\)](#)  
 155 is a *compact* manifold, and every continuous function that is geodesically convex on  
 156 a compact Riemannian manifold can only be a constant; see, e.g., [\[5, Proposition 2.2\]](#)  
 157 and [\[64\]](#).

158 *Riemannian gradient sampling algorithms.* For general Riemannian nonsmooth  
 159 optimization, the recent works [22, 23] propose Riemannian gradient sampling algo-  
 160 rithms, which are motivated by the gradient sampling algorithms for nonconvex non-  
 161 smooth optimization in the Euclidean space [9]. As introduced in [22, 23], given the  
 162 current iterate  $\mathbf{X}_k$ , a typical Riemannian gradient sampling algorithm first samples  
 163 some points  $\{\mathbf{X}_k^j\}_{j=1}^J$  in the neighborhood of  $\mathbf{X}_k$  at which the objective function  $f$  is  
 164 differentiable, where the number of sampled points  $J$  usually needs to be larger than  
 165 the dimension of the manifold  $\mathcal{M}$ . Then, to obtain a descent direction, it solves the  
 166 quadratic program

$$167 \quad (1.6) \quad \boldsymbol{\xi}_k = - \underset{\mathbf{G} \in \text{conv}(\mathcal{W})}{\text{argmin}} \|\mathbf{G}\|^2,$$

168 where  $\text{conv}(\mathcal{W})$  denotes the convex hull of  $\mathcal{W} := \{\text{grad } f(\mathbf{X}_k^1), \dots, \text{grad } f(\mathbf{X}_k^J)\}$  and  
 169  $\text{grad } f$  is the Riemannian gradient of  $f$  on  $\mathcal{M}$ . The update can then be performed via  
 170 classical retractions on  $\mathcal{M}$  using the descent direction  $\boldsymbol{\xi}_k$ . This type of algorithms can  
 171 potentially be utilized to solve a large class of Riemannian nonsmooth optimization  
 172 problems. However, they are only known to converge asymptotically without any rate  
 173 guarantee [22, 23]. Moreover, in order to tackle problem (1.1) with large  $n$  and  $r$  using  
 174 a Riemannian gradient sampling algorithm, one has to sample a large number of Rie-  
 175 mannian gradients in each iteration, which makes the subproblem (1.6) very expensive  
 176 to solve. By contrast, although we assume that the objective function in (1.1) has  
 177 weakly convex components, we can establish the convergence of various Riemannian  
 178 subgradient-type methods with explicit rate guarantees. In addition, each iteration  
 179 of those methods involves only the computation of a Riemannian subgradient, which  
 180 can potentially be much cheaper.

181 *Two types of proximal point methods.* Another classic approach to tackling Rie-  
 182 mannian nonsmooth optimization is to apply proximal point-type methods. The idea  
 183 is to iteratively compute the proximal mapping of the objective function over the  
 184 Riemannian manifold [14, 20]. These methods are shown to converge globally at a  
 185 sublinear rate, based on the so-called sufficient decrease property. However, the main  
 186 issue with this type of methods is that each subproblem is as difficult as the original  
 187 problem, which renders them not practical. When specialized to the Stiefel manifold,  
 188 such a difficulty has been alleviated by some recent advances [10, 11, 25]. Specifically,  
 189 they propose to compute the proximal mapping over the tangent space instead of  
 190 over the Stiefel manifold, which results in a linearly constrained convex subproblem  
 191 that is much easier to solve than the original problem. They also prove that the new  
 192 algorithms converge globally at a sublinear rate. Nonetheless, the subproblem still  
 193 needs to be solved by an iterative algorithm. By contrast, the methods considered  
 194 in this paper do not need to solve expensive subproblems except for the computation  
 195 of one Riemannian subgradient. As such, our overall computational complexities are  
 196 much lower.

197 *Splitting-type methods.* There are also splitting-type methods for solving Rie-  
 198 mannian nonsmooth optimization problems, such as the manifold ADMM-type algo-  
 199 rithms in [30, 31]. In this approach, the problem at hand is typically split into two  
 200 subproblems—one involves optimizing a smooth function over the Riemannian man-  
 201 ifold, the other involves optimizing a nonsmooth function without any constraint.  
 202 These subproblems are then solved in an alternating manner. Despite their simplic-  
 203 ity, these methods often do not have any convergence guarantee.

### 204 Nonsmooth optimization over Stiefel manifold for specific problems.

205 Finally, we close this subsection by mentioning several problem-specific results. The  
 206 recent works [3] and [69, 70] propose to use the Riemannian subgradient method to  
 207 solve the orthogonal DL problem (1.4) and RSR problem (1.3), respectively, and  
 208 establish its local linear convergence when solving these problems. The proofs are  
 209 based on a certain regularity condition instead of the sharpness property studied in  
 210 this work. We will give a detailed comparison between the said regularity condition  
 211 and the sharpness property in Section 5. For now, it is worth noting that the analyses  
 212 in [3, 69, 70] critically depend on the specific model structure of the problem at hand  
 213 and cannot be easily generalized. By contrast, we develop a more general frame-  
 214 work for analyzing Riemannian subgradient-type methods when applied to a family  
 215 of nonsmooth nonconvex optimization problems over certain compact Riemannian  
 216 submanifolds, which can yield both global and local convergence guarantees.

217 **1.4. Notation.** We use  $T_{\mathbf{X}} \text{St} := \{\boldsymbol{\xi} \in \mathbb{R}^{n \times r} : \boldsymbol{\xi}^\top \mathbf{X} + \mathbf{X}^\top \boldsymbol{\xi} = 0\}$  to denote  
 218 the tangent space to the Stiefel manifold  $\text{St}(n, r)$  at the point  $\mathbf{X} \in \text{St}(n, r)$ . Let  
 219  $\langle \mathbf{A}, \mathbf{B} \rangle = \text{trace}(\mathbf{A}^\top \mathbf{B})$  denote the Euclidean inner product of two matrices  $\mathbf{A}, \mathbf{B}$   
 220 of the same dimensions and  $\|\mathbf{A}\|_F = \sqrt{\langle \mathbf{A}, \mathbf{A} \rangle}$  denote the Frobenius norm of  $\mathbf{A}$ .  
 221 We endow the Stiefel manifold  $\text{St}(n, r)$  with the Riemannian metric inherited from  
 222 the Euclidean inner product; i.e.,  $\langle \mathbf{X}, \mathbf{Y} \rangle = \text{trace}(\mathbf{X}^\top \mathbf{Y})$  for any  $\mathbf{X}, \mathbf{Y} \in T_{\mathbf{Z}} \text{St}$  and  
 223  $\mathbf{Z} \in \text{St}(n, r)$ . For a closed set  $\mathcal{C} \subseteq \mathbb{R}^{n \times r}$ , we use  $\mathcal{P}_{\mathcal{C}}$  to denote the orthogonal projector  
 224 onto  $\mathcal{C}$  and  $\text{dist}(\mathbf{X}, \mathcal{C}) := \inf_{\mathbf{Y} \in \mathcal{C}} \|\mathbf{X} - \mathbf{Y}\|_F$  to denote the distance between  $\mathbf{X}$  and  
 225  $\mathcal{C}$ . We use  $x \lesssim y$  and  $x \gtrsim y$  to denote  $x \leq cy$  and  $x \geq cy$  for some universal constant  
 226  $c$ , respectively.

227 **2. Preliminaries.** In this section, we first review some basic notions in Riemannian  
 228 optimization and then present the Riemannian subgradient-type algorithms for  
 229 solving problem (1.1).

#### 230 2.1. Optimization over Stiefel manifold.

231 **Riemannian subgradient and first-order optimality condition.** By our  
 232 assumption, the objective function  $f$  in (1.1) is  $\tau$ -weakly convex for some  $\tau \geq 0$ ; i.e.,  
 233 there exists a convex function  $g : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$  such that  $f(\mathbf{X}) = g(\mathbf{X}) - \frac{\tau}{2} \|\mathbf{X}\|_F^2$  for  
 234 any  $\mathbf{X} \in \mathbb{R}^{n \times r}$  [55, Proposition 4.3]. Although  $f$  may not be convex, we may define  
 235 its (Euclidean) subdifferential  $\partial f$  via

$$236 \quad (2.1) \quad \partial f(\mathbf{X}) = \partial g(\mathbf{X}) - \tau \mathbf{X}, \quad \forall \mathbf{X} \in \mathbb{R}^{n \times r};$$

237 see [55, Proposition 4.6]. Note that since  $g$  is convex,  $\partial g$  is simply its usual convex  
 238 subdifferential. Hence, the subdifferential  $\partial f$  in (2.1) is well defined.

239 Using the properties of weakly convex functions in [55, Proposition 4.5] and the  
 240 result in [63, Theorem 5.1], the Riemannian subdifferential  $\partial_{\mathcal{R}} f$  of  $f$  on the Stiefel  
 241 manifold  $\text{St}(n, r)$  is given by

$$242 \quad (2.2) \quad \partial_{\mathcal{R}} f(\mathbf{X}) = \mathcal{P}_{T_{\mathbf{X}} \text{St}}(\partial f(\mathbf{X})), \quad \forall \mathbf{X} \in \text{St}(n, r).$$

243 In particular, given an Euclidean subgradient  $\tilde{\nabla} f(\mathbf{X}) \in \partial f(\mathbf{X})$  of  $f$  at  $\mathbf{X} \in \text{St}(n, r)$ ,  
 244 we obtain a corresponding Riemannian subgradient  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}) \in \partial_{\mathcal{R}} f(\mathbf{X})$  through  
 245  $\tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}) = \mathcal{P}_{T_{\mathbf{X}} \text{St}}(\tilde{\nabla} f(\mathbf{X}))$ . Recall that for any  $\mathbf{B} \in \mathbb{R}^{n \times r}$ , the projection of  $\mathbf{B}$  onto  
 246  $T_{\mathbf{X}} \text{St}$  is given by  $\mathcal{P}_{T_{\mathbf{X}} \text{St}}(\mathbf{B}) = \mathbf{B} - \frac{1}{2} \mathbf{X} (\mathbf{B}^\top \mathbf{X} + \mathbf{X}^\top \mathbf{B})$  [2, Example 3.6.2].

247 Using (2.2), we call  $\mathbf{X} \in \text{St}(n, r)$  a *stationary point* of problem (1.1) if it satisfies  
 248 the following first-order optimality condition:

$$249 \quad (2.3) \quad \mathbf{0} \in \partial_{\mathcal{R}} f(\mathbf{X}).$$

250 **Retractions on Stiefel manifold.** To enable search along curves on the Stiefel  
 251 manifold, we need the notion of a retraction (see [2, Definition 4.1.1] for the definition).  
 252 There are four commonly used retractions on the Stiefel manifold. These include the  
 253 exponential map [17] and those based on the  $QR$  decomposition, Cayley transforma-  
 254 tion [60], and polar decomposition. It is mentioned in [11] that among the above  
 255 four retractions, the polar decomposition-based one is the most efficient in terms of  
 256 computational complexity. Therefore, we shall focus on polar decomposition-based  
 257 retraction, which is given by

$$258 \quad (2.4) \quad \text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) = (\mathbf{X} + \boldsymbol{\xi})(\mathbf{I}_r + \boldsymbol{\xi}^\top \boldsymbol{\xi})^{-\frac{1}{2}}.$$

259 However, we remark that our results also apply to the other three retractions; see  
 260 Section 6 for a detailed discussion.

261 As the following lemma shows, given any  $\mathbf{X} \in \text{St}(n, r)$  and  $\boldsymbol{\xi} \in \text{T}_{\mathbf{X}} \text{St}$ , the polar  
 262 decomposition-based retraction at  $\mathbf{X}$  essentially computes the projection of  $\mathbf{X} + \boldsymbol{\xi}$   
 263 onto  $\text{St}(n, r)$ . Moreover, this projection has a Lipschitz-like behavior, even though  
 264  $\text{St}(n, r)$  is nonconvex.

265 **LEMMA 1.** *Let  $\mathbf{X} \in \text{St}(n, r)$  and  $\boldsymbol{\xi} \in \text{T}_{\mathbf{X}} \text{St}$  be given. Consider the point  $\mathbf{X}^+ =$   
 266  $\mathbf{X} + \boldsymbol{\xi}$ . Then, the polar decomposition-based retraction (2.4) satisfies  $\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) =$   
 267  $\mathbf{X}^+ (\mathbf{X}^{+\top} \mathbf{X}^+)^{-\frac{1}{2}} = \mathcal{P}_{\text{St}}(\mathbf{X}^+)$  and*

$$268 \quad \|\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) - \overline{\mathbf{X}}\|_F \leq \|\mathbf{X}^+ - \overline{\mathbf{X}}\|_F = \|\mathbf{X} + \boldsymbol{\xi} - \overline{\mathbf{X}}\|_F, \quad \forall \overline{\mathbf{X}} \in \text{St}(n, r).$$

269 *Proof.* It is well known that the convex hull of the Stiefel manifold  $\text{St}(n, r)$  is  
 270 given by  $H \equiv H(n, r) := \{\mathbf{Y} \in \mathbb{R}^{n \times r} : \|\mathbf{Y}\|_2 \leq 1\}$ , where  $\|\mathbf{Y}\|_2$  denotes the spectral  
 271 norm (i.e. the largest singular value) of  $\mathbf{Y}$ ; see, e.g., [27]. Let us first show that  
 272  $\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) = \mathcal{P}_{\text{St}}(\mathbf{X}^+) = \mathcal{P}_H(\mathbf{X}^+)$ . Let  $\mathbf{X}^+ = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^\top$  be an SVD of  $\mathbf{X}^+$ . Since  
 273  $\boldsymbol{\xi} \in \text{T}_{\mathbf{X}} \text{St}$ , we have  $\mathbf{X}^{+\top} \mathbf{X}^+ = \mathbf{I}_r + \boldsymbol{\xi}^\top \boldsymbol{\xi}$ , which implies that all the singular values of  
 274  $\mathbf{X}^+$  are at least 1. This, together with the Hoffman-Wielandt Theorem for singular  
 275 values (see, e.g., [51]), implies that  $\mathcal{P}_{\text{St}}(\mathbf{X}^+) = \mathcal{P}_H(\mathbf{X}^+) = \mathbf{U}\mathbf{V}^\top$ , as desired.

276 Now, observe that  $\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) = \mathbf{X}^+ (\mathbf{X}^{+\top} \mathbf{X}^+)^{-\frac{1}{2}} = \mathbf{U}\mathbf{V}^\top$  and  $\overline{\mathbf{X}} \in H(n, r)$ .  
 277 Hence, we have  $\|\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) - \overline{\mathbf{X}}\|_F = \|\mathcal{P}_H(\mathbf{X}^+) - \mathcal{P}_H(\overline{\mathbf{X}})\|_F$ . Upon noting that  
 278 projections onto closed convex sets are 1-Lipschitz, the proof is complete.  $\square$

## 279 2.2. A family of Riemannian subgradient-type methods.

280 **Riemannian subgradient method.** We begin by revisiting the Riemannian  
 281 gradient method for smooth optimization over the Stiefel manifold. Let  $h : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$   
 282 be a smooth function and consider

$$283 \quad \begin{aligned} & \text{minimize } h(\mathbf{X}) \\ & \mathbf{X} \in \mathbb{R}^{n \times r} \\ & \text{subject to } \mathbf{X} \in \text{St}(n, r). \end{aligned}$$

284 A generic Riemannian gradient method for solving the above problem is given by

$$285 \quad \mathbf{X}_{k+1} = \text{Retr}_{\mathbf{X}_k}(\boldsymbol{\xi}_k) \quad \text{with} \quad \boldsymbol{\xi}_k = -\gamma_k \text{grad } h(\mathbf{X}_k),$$

286 where  $\text{grad } h(\mathbf{X}_k)$  is the Riemannian gradient of  $h$  at  $\mathbf{X}_k$ ,  $\gamma_k > 0$  is the stepsize,  
 287 and  $\text{Retr}$  is any retraction on the Stiefel manifold; see, e.g., [2, Section 4.2]. Since  
 288 problem (1.1) involves a possibly nonsmooth objective function, one approach to

289 tackling it is to apply a natural generalization of the Riemannian gradient method,  
290 namely the Riemannian subgradient method:

$$291 \quad (2.5) \quad \boxed{\mathbf{X}_{k+1} = \text{Retr}_{\mathbf{X}_k}(\boldsymbol{\xi}_k) \quad \text{with} \quad \boldsymbol{\xi}_k = -\gamma_k \widetilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_k).}$$

292 Here, recall that  $\widetilde{\nabla}_{\mathcal{R}} f(\mathbf{X}_k) \in \partial_{\mathcal{R}} f(\mathbf{X}_k)$  is a Riemannian subgradient of  $f$  at  $\mathbf{X}_k \in$   
293  $\text{St}(n, r)$ , which can be obtained by taking  $\widetilde{\nabla} f(\mathbf{X}) \in \partial f(\mathbf{X})$  and setting  $\widetilde{\nabla}_{\mathcal{R}} f(\mathbf{X}) =$   
294  $\mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}(\widetilde{\nabla} f(\mathbf{X}))$ ; see [Subsection 2.1](#).

295 **Riemannian incremental and stochastic subgradient methods.** Recall  
296 that the objective function in (1.1) has the *finite-sum structure*  $f = \frac{1}{m} \sum_{i=1}^m f_i$ . In  
297 many modern machine learning tasks, the number of components  $m$  can be very  
298 large. Thus, it is not desirable to evaluate the full Riemannian subgradient of  $f$ . This  
299 motivates us to introduce two variants of the Riemannian subgradient method (2.5),  
300 namely the Riemannian incremental subgradient method and Riemannian stochastic  
301 subgradient method, to better exploit the finite-sum structure in (1.1). Let us now  
302 give a high-level description of these two methods.

303 Starting with the current iterate  $\mathbf{X}_k$ , both methods generate a sequence of  $m$   
304 inner iterates  $\mathbf{X}_{k,1}, \dots, \mathbf{X}_{k,i}, \dots, \mathbf{X}_{k,m}$  via

$$305 \quad (2.6) \quad \boxed{\mathbf{X}_{k,i} = \text{Retr}_{\mathbf{X}_{k,i-1}}(\boldsymbol{\xi}_{k,i-1}) \quad \text{with} \quad \boldsymbol{\xi}_{k,i-1} = -\gamma_k \widetilde{\nabla}_{\mathcal{R}} f_{\ell_i}(\mathbf{X}_{k,i-1})}$$

306 with  $\mathbf{X}_{k,0} = \mathbf{X}_k$ , where  $f_{\ell_i}$  is selected from  $\{f_1, \dots, f_m\}$  according to a certain rule.  
307 The next iterate  $\mathbf{X}_{k+1}$  is then obtained by setting  $\mathbf{X}_{k+1} = \mathbf{X}_{k,m}$ . The difference  
308 between the incremental and stochastic methods lies in the rule for selecting the  
309 component function  $f_{\ell_i}$ . In particular,

- 310 • *Riemannian incremental subgradient method* picks the component function  
311  $f_{\ell_i}$  *sequentially* from  $f_1$  to  $f_m$ —i.e.,  $\boldsymbol{\xi}_{k,i-1} = -\gamma_k \widetilde{\nabla}_{\mathcal{R}} f_i(\mathbf{X}_{k,i-1})$ ;
- 312 • *Riemannian stochastic subgradient method* picks the component function  $f_{\ell_i}$   
313 *independently* and *uniformly* from  $\{f_1, \dots, f_m\}$  in each inner iteration (2.6)—  
314 i.e.,  $\boldsymbol{\xi}_{k,i-1} = -\gamma_k \widetilde{\nabla}_{\mathcal{R}} f_{\ell_i}(\mathbf{X}_{k,i-1})$  with  $\ell_i \sim_{i.i.d.} \text{Uniform}(\{1, \dots, m\})$ .

315 **3. Riemannian Subgradient Inequality over Stiefel Manifold.** Naturally,  
316 we are interested in the convergence behavior of the Riemannian subgradient-type  
317 methods introduced in [Subsection 2.2](#) when applied to problem (1.1). Towards that  
318 end, let us derive a useful inequality, which we call the *Riemannian subgradient in-*  
319 *equality*, for restrictions of weakly convex functions on the Stiefel manifold. The main  
320 motivation for deriving such an inequality is that an analogous one for weakly con-  
321 vex functions in the Euclidean space, known as the *weakly convex inequality*, plays a  
322 fundamental role in the convergence analysis of subgradient-type methods for solv-  
323 ing weakly convex minimization problems [12, 13, 36, 37]. To begin, recall that for a  
324  $\tau$ -weakly convex function  $h : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$ , the weakly convex inequality states that

$$325 \quad (3.1) \quad h(\mathbf{Y}) \geq h(\mathbf{X}) + \left\langle \widetilde{\nabla} h(\mathbf{X}), \mathbf{Y} - \mathbf{X} \right\rangle - \frac{\tau}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2,$$

$$\forall \widetilde{\nabla} h(\mathbf{X}) \in \partial h(\mathbf{X}); \quad \mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times r}$$

326 [55, Proposition 4.8]. The following is our extension of the above inequality to one  
327 for weakly convex functions that are restricted on the Stiefel manifold.

328 **THEOREM 1 (Riemannian Subgradient Inequality).** *Suppose that  $h : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$*   
329 *is  $\tau$ -weakly convex for some  $\tau \geq 0$ . Then, for any bounded open convex set  $\mathcal{U}$  that*



330 contains  $\text{St}(n, r)$ , there exists a constant  $L > 0$  such that  $h$  is  $L$ -Lipschitz continuous  
 331 on  $\mathcal{U}$  and satisfies

$$332 \quad (3.2) \quad h(\mathbf{Y}) \geq h(\mathbf{X}) + \left\langle \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}), \mathbf{Y} - \mathbf{X} \right\rangle - \frac{\tau + L}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2, \\ \forall \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}) \in \partial_{\mathcal{R}} h(\mathbf{X}); \mathbf{X}, \mathbf{Y} \in \text{St}(n, r).$$

333 Before we proceed to prove [Theorem 1](#), let us highlight the differences between  
 334 the weakly convex inequality (3.1) and the Riemannian subgradient inequality (3.2).  
 335 First, the former involves elements in the Euclidean subdifferential  $\partial h$ , while the latter  
 336 involves elements in the Riemannian subdifferential  $\partial_{\mathcal{R}} h$ . Second, the former holds  
 337 for all pairs of points in the Euclidean space  $\mathbb{R}^{n \times r}$ , while the latter only holds for all  
 338 pairs of points on the Stiefel manifold  $\text{St}(n, r)$ . Third, the latter involves the extra  
 339 compensation term  $-\frac{L}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2$ , which accounts for the behavior of the restriction  
 340 of  $h$  on the Stiefel manifold  $\text{St}(n, r)$ .

341 *Proof of Theorem 1.* The Lipschitz continuity of  $h$  on  $\mathcal{U}$  follows directly from [[55](#),  
 342 Proposition 4.4] and the boundedness of  $\mathcal{U}$ . Since  $h$  is  $\tau$ -weakly convex on  $\mathbb{R}^{n \times r}$ , for  
 343 any  $\mathbf{X}, \mathbf{Y} \in \text{St}(n, r) \subseteq \mathbb{R}^{n \times r}$ , the inequality (3.1) implies that

$$344 \quad (3.3) \quad h(\mathbf{Y}) \geq h(\mathbf{X}) + \left\langle \tilde{\nabla} h(\mathbf{X}), \mathbf{Y} - \mathbf{X} \right\rangle - \frac{\tau}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2 \\ = h(\mathbf{X}) + \left\langle \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}(\tilde{\nabla} h(\mathbf{X})) + \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\tilde{\nabla} h(\mathbf{X})), \mathbf{Y} - \mathbf{X} \right\rangle - \frac{\tau}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2,$$

345 where

$$346 \quad (3.4) \quad \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\mathbf{B}) = \frac{1}{2} \mathbf{X} (\mathbf{B}^{\top} \mathbf{X} + \mathbf{X}^{\top} \mathbf{B}), \quad \forall \mathbf{B} \in \mathbb{R}^{n \times r}$$

347 [[2](#), Example 3.6.2]. Now, we compute

$$348 \quad (3.5) \quad \left\langle \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\tilde{\nabla} h(\mathbf{X})), \mathbf{Y} - \mathbf{X} \right\rangle = \left\langle \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\tilde{\nabla} h(\mathbf{X})), \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\mathbf{Y} - \mathbf{X}) + \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}(\mathbf{Y} - \mathbf{X}) \right\rangle \\ = \left\langle \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\tilde{\nabla} h(\mathbf{X})), \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\mathbf{Y} - \mathbf{X}) \right\rangle \\ = \left\langle \tilde{\nabla} h(\mathbf{X}), \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}^{\perp}(\mathbf{Y} - \mathbf{X}) \right\rangle \\ \stackrel{(i)}{=} \frac{1}{2} \left\langle \tilde{\nabla} h(\mathbf{X}), \mathbf{X} (\mathbf{Y}^{\top} \mathbf{X} + \mathbf{X}^{\top} \mathbf{Y} - 2\mathbf{I}_r) \right\rangle \\ \stackrel{(ii)}{\geq} -\frac{1}{2} \|\tilde{\nabla} h(\mathbf{X})\|_F \|\mathbf{Y}^{\top} \mathbf{X} + \mathbf{X}^{\top} \mathbf{Y} - 2\mathbf{I}_r\|_F \\ \stackrel{(iii)}{\geq} -\frac{L}{2} \|\mathbf{Y}^{\top} \mathbf{X} + \mathbf{X}^{\top} \mathbf{Y} - 2\mathbf{I}_r\|_F,$$

349 where (i) comes from (3.4), (ii) is due to the fact that  $\mathbf{X} \in \text{St}(n, r)$ , and (iii) follows  
 350 from [[46](#), Theorem 9.13] and the  $L$ -Lipschitz continuity of  $h$  on  $\mathcal{U}$ . Note that

$$351 \quad (3.6) \quad \|\mathbf{Y}^{\top} \mathbf{X} + \mathbf{X}^{\top} \mathbf{Y} - 2\mathbf{I}_r\|_F = \|(\mathbf{X} - \mathbf{Y})^{\top} (\mathbf{X} - \mathbf{Y})\|_F \leq \|\mathbf{X} - \mathbf{Y}\|_F^2$$

352 since  $\mathbf{X}, \mathbf{Y} \in \text{St}(n, r)$ . Combining (3.5) and (3.6) and recalling (3.3), we get

$$353 \quad h(\mathbf{Y}) \geq h(\mathbf{X}) + \left\langle \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}(\tilde{\nabla} h(\mathbf{X})), \mathbf{Y} - \mathbf{X} \right\rangle - \frac{\tau + L}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2.$$

354 Since  $\mathbf{X}, \mathbf{Y} \in \text{St}(n, r)$ ,  $\tilde{\nabla} h(\mathbf{X}) \in \partial h(\mathbf{X})$  are arbitrary and  $\partial_{\mathcal{R}} h(\mathbf{X}) = \mathcal{P}_{\text{T}_{\mathbf{X}} \text{St}}(\partial h(\mathbf{X}))$   
 355 (see (2.2)), the proof is complete.  $\square$

356 As we shall see in subsequent sections, the Riemannian subgradient inequality  
 357 plays a similar role to the weakly convex inequality and allows us to connect the  
 358 analysis of Riemannian subgradient-type methods with that of their Euclidean coun-  
 359 terparts. In particular, equipped with [Theorem 1](#), we can obtain the iteration com-  
 360 plexities of the Riemannian subgradient-type methods introduced in [Subsection 2.2](#)  
 361 for addressing problem [\(1.1\)](#). Moreover, if problem [\(1.1\)](#) further possesses certain  
 362 sharpness property (see [Definition 1](#)), then the aforementioned methods with geo-  
 363 metrically diminishing stepsizes and a proper initialization will achieve local linear  
 364 convergence to the set of so-called weak sharp minima (again, see [Definition 1](#)).

365 Although the Riemannian subgradient inequality in [Theorem 1](#) focuses on the  
 366 Stiefel manifold, it can be extended to a class of compact embedded submanifolds of  
 367 the Euclidean space. We shall present such an extension in [Section 6](#).

368 **4. Global Convergence.** In this section, we study the iteration complexities of  
 369 Riemannian subgradient-type methods for solving problem [\(1.1\)](#). Our analysis relies  
 370 on the Riemannian subgradient inequality in [Theorem 1](#).

371 **4.1. Surrogate stationarity measure.** In classical Euclidean nonsmooth con-  
 372 vex optimization, the iteration complexities of subgradient-type methods are typically  
 373 presented in terms of the suboptimality gap  $f(\mathbf{X}_k) - \min f$ ; see, e.g., [\[44, Theorem](#)  
 374 [3.2.2\]](#), [\[42, Proposition 2.3\]](#). On the other hand, in Riemannian smooth optimiza-  
 375 tion, which typically involves nonconvex constraints, the iteration complexities of  
 376 various methods can be expressed in terms of the continuous stationarity measure  
 377  $\|\text{grad}f(\mathbf{X}_k)\|_F$  [\[7\]](#). However, for the Riemannian nonsmooth optimization problem  
 378 [\(1.1\)](#), neither the suboptimality gap  $f(\mathbf{X}_k) - \min f$  (due to nonconvexity) nor the  
 379 minimum-norm Riemannian subgradient  $\text{dist}(\mathbf{0}, \partial_{\mathcal{R}}f(\mathbf{X}_k))$  (due to nonsmoothness)  
 380 is an appropriate stationarity measure. Therefore, in order to establish the iteration  
 381 complexities of Riemannian subgradient-type methods, we need to find a surrogate  
 382 stationarity measure that can track the progress of those methods.

383 Towards that end, we borrow ideas from the recent works [\[12, 16\]](#) on weakly convex  
 384 minimization in the Euclidean space, which propose to use the gradient of the Moreau  
 385 envelope of the weakly convex function at hand as a surrogate stationarity measure.  
 386 To begin, let us define, for any  $\lambda > 0$ , the following analogs of the Moreau envelope  
 387 and proximal mapping for problem [\(1.1\)](#), which take into account the effect of the  
 388 Stiefel manifold constraint on the problem:

$$389 \quad (4.1) \quad \begin{cases} f_\lambda(\mathbf{X}) = \min_{\mathbf{Y} \in \text{St}(n, r)} \left\{ f(\mathbf{Y}) + \frac{1}{2\lambda} \|\mathbf{Y} - \mathbf{X}\|_F^2 \right\}, & \mathbf{X} \in \text{St}(n, r), \\ P_{\lambda f}(\mathbf{X}) \in \underset{\mathbf{Y} \in \text{St}(n, r)}{\text{argmin}} \left\{ f(\mathbf{Y}) + \frac{1}{2\lambda} \|\mathbf{Y} - \mathbf{X}\|_F^2 \right\}, & \mathbf{X} \in \text{St}(n, r). \end{cases}$$

390 We remark that the Moreau envelope and proximal mapping defined above differ from  
 391 those in [\[20\]](#) in that the proximal term  $\mathbf{Y} \mapsto \frac{1}{2\lambda} \|\mathbf{Y} - \mathbf{X}\|_F^2$  is based on the Euclidean  
 392 distance rather than the geodesic distance. This will facilitate our later analysis.

393 By [\(2.2\)](#) and [\(2.3\)](#), the point  $P_{\lambda f}(\mathbf{X})$  satisfies the first-order optimality condition  
 394  $\mathbf{0} \in \partial_{\mathcal{R}}f(P_{\lambda f}(\mathbf{X})) + \frac{1}{\lambda} \mathcal{P}_{\text{T}_{P_{\lambda f}(\mathbf{X})} \text{St}}(P_{\lambda f}(\mathbf{X}) - \mathbf{X})$ . It follows that

$$395 \quad (4.2) \quad \boxed{\begin{aligned} \text{dist}(\mathbf{0}, \partial_{\mathcal{R}}f(P_{\lambda f}(\mathbf{X}))) &\leq \lambda^{-1} \cdot \left\| \mathcal{P}_{\text{T}_{P_{\lambda f}(\mathbf{X})} \text{St}}(P_{\lambda f}(\mathbf{X}) - \mathbf{X}) \right\|_F \\ &\leq \lambda^{-1} \cdot \|P_{\lambda f}(\mathbf{X}) - \mathbf{X}\|_F =: \Theta(\mathbf{X}). \end{aligned}}$$

396 In particular, we see from (2.3) that  $\mathbf{X} \in \text{St}(n, r)$  is a stationary point of problem (1.1)  
 397 when  $\Theta(\mathbf{X}) = 0$ . This motivates us to use  $\mathbf{X} \mapsto \Theta(\mathbf{X})$  as a surrogate stationarity  
 398 measure of problem (1.1) and call  $\mathbf{X} \in \text{St}(n, r)$  an  $\varepsilon$ -nearly stationary point of problem  
 399 (1.1) if it satisfies  $\Theta(\mathbf{X}) \leq \varepsilon$ .

400 The careful reader may note that the proximal mapping  $P_{\lambda f}$  in (4.1) needs not  
 401 yield a unique point at a given  $\mathbf{X} \in \text{St}(n, r)$ . Nevertheless, for the purpose of defin-  
 402 ing the surrogate stationarity measure, we can choose any point returned by  $P_{\lambda f}$  at  
 403  $\mathbf{X}$ , as each of them plays exactly the same role in our analysis and will satisfy the  
 404 convergence rate bounds in Theorem 2.

#### 405 4.2. Riemannian subgradient and incremental subgradient methods.

406 Using the surrogate stationarity measure  $\Theta$ , we are now ready to establish the iteration  
 407 complexities of the Riemannian subgradient and incremental subgradient methods.  
 408 We will focus on analyzing the Riemannian incremental subgradient method, as the  
 409 Riemannian subgradient method can be regarded as its special case where there is  
 410 only one (i.e.,  $m = 1$ ) component function.

411 To begin, let us establish a relationship between the surrogate stationarity mea-  
 412 sure  $\Theta$  and the sufficient decrease of the Moreau envelope  $f_\lambda$ .

413 PROPOSITION 1. *Suppose that each component function  $f_i : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$  ( $i =$   
 414  $1, \dots, m$ ) in problem (1.1) is  $\tau$ -weakly convex on  $\mathbb{R}^{n \times r}$  for some  $\tau \geq 0$ . Let  $\mathcal{U}$  be  
 415 any bounded open convex set that contains  $\text{St}(n, r)$ . Furthermore, let  $\{\mathbf{X}_k\}$  be the  
 416 sequence generated by the Riemannian incremental subgradient method (2.6) with an  
 417 arbitrary initialization for solving problem (1.1). Then, for any  $\lambda < \frac{1}{2(L+\tau)}$  in (4.1),  
 418 we have for any  $k \geq 0$*

$$419 \quad m\gamma_k \Theta^2(\mathbf{X}_k) \leq \frac{2(f_\lambda(\mathbf{X}_k) - f_\lambda(\mathbf{X}_{k+1})) + \frac{\gamma_k^2 L^2}{\lambda} m^2 + \frac{\gamma_k^3 L^2(L+\tau)}{\lambda} C(m)}{2\lambda \left(\frac{1}{2\lambda} - (L+\tau)\right)},$$

420 where  $L > 0$  is an upper bound on the Lipschitz constants of  $f_1, \dots, f_m$  on  $\mathcal{U}$  and  
 421  $C(m) = \frac{1}{3}m(m-1)(2m-1)$ .

422 *Proof.* According to (4.1), we have

$$423 \quad (4.3) \quad \begin{aligned} f_\lambda(\mathbf{X}_{k+1}) &= f(P_{\lambda f}(\mathbf{X}_{k+1})) + \frac{1}{2\lambda} \|P_{\lambda f}(\mathbf{X}_{k+1}) - \mathbf{X}_{k+1}\|_F^2 \\ &\leq f(P_{\lambda f}(\mathbf{X}_k)) + \frac{1}{2\lambda} \|P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_{k+1}\|_F^2, \end{aligned}$$

424 where the last inequality follows from the optimality of  $P_{\lambda f}(\mathbf{X}_{k+1})$  and the fact that  
 425  $P_{\lambda f}(\mathbf{X}_k) \in \text{St}(n, r)$ . We claim that for  $l = 1, \dots, m$ ,

$$426 \quad (4.4) \quad \begin{aligned} \|P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_{k,l}\|_F^2 &\leq \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 - 2\gamma_k \sum_{i=1}^l (f_i(\mathbf{X}_{k,i-1}) - f_i(P_{\lambda f}(\mathbf{X}_k))) \\ &\quad + \gamma_k(L+\tau) \sum_{i=1}^l \|\mathbf{X}_{k,i-1} - P_{\lambda f}(\mathbf{X}_k)\|_F^2 + l\gamma_k^2 L^2. \end{aligned}$$

427 The proof is by induction on  $l$ . For  $l = 1$ , recalling that  $\mathbf{X}_{k,0} = \mathbf{X}_k$ , we compute

$$428 \quad (4.5) \quad \begin{aligned} \|P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_{k,1}\|_F^2 &\leq \|\mathbf{X}_k + \boldsymbol{\xi}_{k,0} - P_{\lambda f}(\mathbf{X}_k)\|_F^2 \\ &\leq \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 - 2\gamma_k (f_1(\mathbf{X}_k) - f_1(P_{\lambda f}(\mathbf{X}_k))) \\ &\quad + \gamma_k(L+\tau) \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 + \gamma_k^2 L^2, \end{aligned}$$

429 where we used (2.6) and Lemma 1 in the first inequality and Theorem 1 and the  
 430 fact that  $\left\| \tilde{\nabla}_{\mathcal{R}} f_i(\mathbf{X}_{k,i-1}) \right\|_F \leq \left\| \tilde{\nabla} f_i(\mathbf{X}_{k,i-1}) \right\|_F \leq L$  in the second inequality. The  
 431 inductive step can be completed by following the same derivations as in (4.5). Thus,  
 432 the claim (4.4) is established. Setting  $l = m$  in (4.4) and plugging it into (4.3), we  
 433 obtain

$$434 \quad (4.6) \quad f_{\lambda}(\mathbf{X}_{k+1}) \leq f_{\lambda}(\mathbf{X}_k) + \frac{\gamma_k}{\lambda} \sum_{i=1}^m (f_i(P_{\lambda f}(\mathbf{X}_k)) - f_i(\mathbf{X}_{k,i-1})) \\ + \frac{\gamma_k(L + \tau)}{2\lambda} \sum_{i=1}^m \|\mathbf{X}_{k,i-1} - P_{\lambda f}(\mathbf{X}_k)\|_F^2 + \frac{m\gamma_k^2 L^2}{2\lambda},$$

435 where we used the relation  $f_{\lambda}(\mathbf{X}_k) = f(P_{\lambda f}(\mathbf{X}_k)) + \frac{1}{2\lambda} \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2$  (since  
 436  $\mathbf{X}_k \in \text{St}(n, r)$ ).

437 Next, we claim that for  $i = 1, \dots, m$ ,

$$438 \quad (4.7) \quad \|\mathbf{X}_{k,i-1} - \mathbf{X}_k\|_F \leq (i-1)\gamma_k L.$$

439 The proof is again by induction on  $i$ . The claim trivially holds when  $i = 1$ . Suppose  
 440 that (4.7) holds for  $i = j$ . For  $i = j + 1$ , we compute  $\|\mathbf{X}_{k,j} - \mathbf{X}_k\|_F \leq \|\mathbf{X}_{k,j-1} +$   
 441  $\boldsymbol{\xi}_{k,j-1} - \mathbf{X}_k\|_F \leq j\gamma_k L$ , where we used (2.6) and Lemma 1 in the first inequality. This  
 442 completes the inductive step and the proof of the claim.

443 With (4.7), we have

$$444 \quad (4.8) \quad f_i(P_{\lambda f}(\mathbf{X}_k)) - f_i(\mathbf{X}_{k,i-1}) = f_i(P_{\lambda f}(\mathbf{X}_k)) - f_i(\mathbf{X}_k) + f_i(\mathbf{X}_k) - f_i(\mathbf{X}_{k,i-1}) \\ \leq (i-1)\gamma_k L^2 + f_i(P_{\lambda f}(\mathbf{X}_k)) - f_i(\mathbf{X}_k)$$

445 and

$$446 \quad (4.9) \quad \|\mathbf{X}_{k,i-1} - P_{\lambda f}(\mathbf{X}_k)\|_F^2 = \|\mathbf{X}_{k,i-1} - \mathbf{X}_k + \mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 \\ \leq 2(i-1)^2\gamma_k^2 L^2 + 2\|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2.$$

447 Plugging (4.8) and (4.9) into (4.6) yields

$$448 \quad (4.10) \quad f_{\lambda}(\mathbf{X}_{k+1}) \leq f_{\lambda}(\mathbf{X}_k) + \frac{m^2\gamma_k^2 L^2 + \frac{1}{3}m(m-1)(2m-1)\gamma_k^3 L^2(L + \tau)}{2\lambda} \\ + \frac{m\gamma_k}{\lambda} (f(P_{\lambda f}(\mathbf{X}_k)) - f(\mathbf{X}_k)) + \frac{m\gamma_k(L + \tau)}{\lambda} \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2.$$

449 By definition of the Moreau envelope and proximal mapping in (4.1), we have

$$450 \quad (4.11) \quad - \left[ f(\mathbf{X}_k) - f(P_{\lambda f}(\mathbf{X}_k)) - (L + \tau) \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 \right] \\ = - \left[ f(\mathbf{X}_k) - \left( f(P_{\lambda f}(\mathbf{X}_k)) + \frac{1}{2\lambda} \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 \right) \right. \\ \left. + \left( \frac{1}{2\lambda} - (L + \tau) \right) \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2 \right] \\ \leq - \left( \frac{1}{2\lambda} - (L + \tau) \right) \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2,$$

451 where the last inequality is due to  $f_\lambda(\mathbf{X}_k) = f(P_{\lambda f}(\mathbf{X}_k)) + \frac{1}{2\lambda} \|\mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k)\|_F^2$   
 452 (since  $\mathbf{X}_k \in \text{St}(n, r)$ ) and  $f_\lambda(\mathbf{X}_k) \leq f(\mathbf{X}_k)$ . Since  $\lambda < \frac{1}{2(L+\tau)}$  by assumption, the  
 453 desired result then follows by substituting (4.11) into (4.10) and recognizing that  
 454  $\Theta(\mathbf{X}_k) = \lambda^{-1} \|P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_k\|_F$  (see (4.2)).  $\square$

455 Using Proposition 1, we obtain our iteration complexity result for the Riemannian  
 456 subgradient and incremental subgradient methods.

457 **THEOREM 2.** *Under the setting of Proposition 1, the following hold:*

458 (a) *If we choose the constant stepsize  $\gamma_k = \frac{1}{m\sqrt{T+1}}$ ,  $k = 0, 1, \dots$  with  $T$  being the  
 459 total number of iterations, then*

$$460 \quad \min_{0 \leq k \leq T} \Theta^2(\mathbf{X}_k) \leq \frac{2(f_\lambda(\mathbf{X}_0) - \min f_\lambda) + \frac{L^2}{\lambda} + \frac{L^2(L+\tau)}{\lambda m^3} C(m)}{2\lambda \left(\frac{1}{2\lambda} - (L+\tau)\right) \sqrt{T+1}}.$$

461 (b) *If we choose the diminishing stepsizes  $\gamma_k = \frac{1}{m\sqrt{k+1}}$ ,  $k = 0, 1, \dots$ , then*

$$462 \quad \min_{0 \leq k \leq T} \Theta^2(\mathbf{X}_k) \leq \frac{2(f_\lambda(\mathbf{X}_0) - \min f_\lambda) + \left(\frac{L^2}{\lambda} + \frac{L^2(L+\tau)}{\lambda m^3} C(m)\right) (\ln(T+1) + 1)}{2\lambda \left(\frac{1}{2\lambda} - (L+\tau)\right) \sqrt{T+1}}.$$

463 *Proof.* By summing both sides of the relation in Proposition 1 over  $k = 0, 1, \dots, T$ ,  
 464 we deduce that

$$465 \quad \min_{0 \leq k \leq T} \Theta^2(\mathbf{X}_k) \leq \frac{2(f_\lambda(\mathbf{X}_0) - \min f_\lambda) + \frac{L^2}{\lambda} m^2 \sum_{k=0}^T \gamma_k^2 + \frac{L^2(L+\tau)}{\lambda} C(m) \sum_{k=0}^T \gamma_k^3}{2\lambda \left(\frac{1}{2\lambda} - (L+\tau)\right) m \sum_{k=0}^T \gamma_k}.$$

466 The result in (a) follows immediately by substituting  $\gamma_k = \frac{1}{m\sqrt{T+1}}$  into the above  
 467 inequality, while that in (b) follows by substituting  $\gamma_k = \frac{1}{m\sqrt{k+1}}$  into the above  
 468 inequality and noting that  $\sum_{k=0}^T \frac{1}{\sqrt{k+1}} > \sqrt{T+1}$  and  $\sum_{k=0}^T \frac{1}{k+1} < \ln(T+1) + 1$ .  $\square$

469 By taking  $\lambda = \frac{1}{4(L+\tau)}$  and using the constant stepsize  $\gamma_k = \frac{1}{m\sqrt{T+1}}$ ,  $k = 0, 1, \dots$ ,  
 470 we see from Theorem 2 that

$$471 \quad \min_{0 \leq k \leq T} \Theta(\mathbf{X}_k) \leq \frac{2\sqrt{(f_\lambda(\mathbf{X}_0) - \min f_\lambda) + 2L^2(L+\tau)(1+L+\tau)}}{(T+1)^{1/4}}.$$

472 In particular, the iteration complexity of the Riemannian (incremental) subgradient  
 473 method for computing an  $\varepsilon$ -nearly stationary point of problem (1.1) is  $\mathcal{O}(\varepsilon^{-4})$ . It  
 474 is worth noting that this matches the iteration complexity of a host of methods for  
 475 solving weakly convex minimization problems in the Euclidean space [12].

476 **4.3. Riemannian stochastic subgradient method.** Now, let us turn to ana-  
 477 lyze the Riemannian stochastic subgradient method. Instead of focusing on objective  
 478 functions with a finite-sum structure as in (1.1), we consider the following more gen-  
 479 eral stochastic optimization problem over the Stiefel manifold:

$$480 \quad (4.12) \quad \begin{aligned} & \text{minimize } f(\mathbf{X}) := \mathbb{E}_{\zeta \sim D}[g(\mathbf{X}, \zeta)] \\ & \text{subject to } \mathbf{X} \in \text{St}(n, r). \end{aligned}$$

481 Here, we assume that the function  $\mathbf{X} \mapsto g(\mathbf{X}, \zeta)$  is  $\tau$ -weakly convex ( $\tau \geq 0$ ) for each  
 482 realization  $\zeta$  and the function  $f$  is finite-valued on  $\mathbb{R}^{n \times r}$ . Furthermore, we assume the

483 existence of a bounded open convex set  $\mathcal{U}$  containing  $\text{St}(n, r)$  such that  $\mathbf{X} \mapsto g(\mathbf{X}, \zeta)$  is  
 484 Lipschitz continuous on  $\mathcal{U}$  with some constant  $L(\zeta) > 0$  and  $L^2 = \mathbb{E}_{\zeta \sim D}[L(\zeta)^2] < +\infty$ .  
 485 This would then imply that for any  $\mathbf{X} \in \text{St}(n, r)$ , we have

$$486 \quad (4.13) \quad \mathbb{E}_{\zeta \sim D} \left[ \left\| \tilde{\nabla} g(\mathbf{X}, \zeta) \right\|_F^2 \right] \leq L^2,$$

487 where  $\tilde{\nabla} g(\mathbf{X}, \zeta) \in \partial g(\mathbf{X}, \zeta)$ . Moreover, the function  $f$  is  $L$ -Lipschitz continuous on  
 488  $\mathcal{U}$ . When  $D$  is the empirical distribution on  $m$  data samples, problem (4.12) reduces  
 489 to our original finite-sum optimization problem (1.1). If all the component functions  
 490 are finite-valued and weakly convex, then the above two assumptions hold.

491 Now, suppose that the Riemannian stochastic subgradient method is equipped  
 492 with a *Riemannian stochastic subgradient oracle*, which has the following properties:

- 493 (a) The oracle can generate i.i.d. samples according to the distribution  $D$ .  
 494 (b) Given a point  $\mathbf{X} \in \text{St}(n, r)$ , the oracle generates a sample  $\zeta \sim D$  and returns  
 495 a stochastic subgradient  $\tilde{\nabla} g(\mathbf{X}, \zeta) \in \partial g(\mathbf{X}, \zeta)$  with  $\mathbb{E}_{\zeta \sim D}[\tilde{\nabla} g(\mathbf{X}, \zeta)] \in \partial f(\mathbf{X})$ ,  
 496 from which one can obtain a Riemannian stochastic subgradient  $\tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}, \zeta) \in$   
 497  $\partial_{\mathcal{R}} g(\mathbf{X}, \zeta)$  with  $\mathbb{E}_{\zeta \sim D}[\tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}, \zeta)] \in \partial_{\mathcal{R}} f(\mathbf{X})$ .

498 We remark that the above properties mirror those of the stochastic subgradient oracle  
 499 for stochastic optimization in the Euclidean space; see, e.g., Assumptions (A1) and  
 500 (A2) in [43].

501 At the current iterate  $\mathbf{X}_k$ , the Riemannian stochastic subgradient oracle generates  
 502 a sample  $\zeta_k \sim D$  that is independent of  $\{\zeta_0, \dots, \zeta_{k-1}\}$  and returns a Riemannian  
 503 stochastic subgradient  $\tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_k, \zeta_k)$ . Then, the Riemannian stochastic subgradient  
 504 method generates the next iterate  $\mathbf{X}_{k+1}$  via

$$505 \quad (4.14) \quad \mathbf{X}_{k+1} = \text{Retr}_{\mathbf{X}_k}(\boldsymbol{\xi}_k) \quad \text{with} \quad \boldsymbol{\xi}_k = -\gamma_k \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_k, \zeta_k).$$

506 This generalizes the update (2.6) introduced in Subsection 2.2 for the case where  $D$   
 507 is the empirical distribution on  $m$  data samples.

508 Similar to the analysis of the Riemannian subgradient and incremental subgradi-  
 509 ent methods, we begin by establishing the following result; cf. Proposition 1:

510 **PROPOSITION 2.** *Suppose that the aforementioned assumptions on problem (4.12)*  
 511 *hold, and that a Riemannian stochastic subgradient oracle having properties (a)–(b)*  
 512 *above is available. Let  $\{\mathbf{X}_k\}$  be the sequence generated by the Riemannian stochas-*  
 513 *tic subgradient method (4.14) with arbitrary initialization for solving problem (4.12).*  
 514 *Then, for any  $\lambda < \frac{1}{L+\tau}$  in (4.1), we have*

$$515 \quad \gamma_k \mathbb{E}[\Theta^2(\mathbf{X}_k)] \leq \frac{2(\mathbb{E}[f_\lambda(\mathbf{X}_k)] - \mathbb{E}[f_\lambda(\mathbf{X}_{k+1})]) + \frac{\gamma_k^2 L^2}{\lambda}}{\lambda \left( \frac{1}{\lambda} - (L + \tau) \right)}, \quad \forall k \geq 0.$$

516 *Proof.* Using (4.1), the optimality of  $P_{\lambda f}(\mathbf{X}_{k+1})$ , Lemma 1, and the fact that  
 517  $P_{\lambda f}(\mathbf{X}_k) \in \text{St}(n, r)$ , we obtain

$$518 \quad \begin{aligned} \mathbb{E}_{\zeta_k \sim D}[f_\lambda(\mathbf{X}_{k+1})] &\leq f(P_{\lambda f}(\mathbf{X}_k)) + \frac{1}{2\lambda} \mathbb{E}_{\zeta_k \sim D} \left[ \|P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_{k+1}\|_F^2 \right] \\ 519 \quad &\leq f(P_{\lambda f}(\mathbf{X}_k)) + \frac{1}{2\lambda} \mathbb{E}_{\zeta_k \sim D} \left[ \left\| \mathbf{X}_k - \gamma_k \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_k, \zeta_k) - P_{\lambda f}(\mathbf{X}_k) \right\|_F^2 \right] \\ 520 \quad &\leq f_\lambda(\mathbf{X}_k) + \frac{\gamma_k}{\lambda} \mathbb{E}_{\zeta_k \sim D} \left[ \left\langle \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_k, \zeta_k), P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_k \right\rangle \right] + \frac{\gamma_k^2 L^2}{2\lambda}, \end{aligned}$$

522 where the first inequality is due to the optimality of  $P_{\lambda f}(\mathbf{X}_{k+1})$ , the second inequality  
 523 comes from [Lemma 1](#) and the fact that  $P_{\lambda f}(\mathbf{X}_k) \in \text{St}(n, r)$ , and the third inequality  
 524 is due to [\(4.13\)](#) and the fact that  $\left\| \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}, \zeta) \right\|_F \leq \left\| \tilde{\nabla} g(\mathbf{X}, \zeta) \right\|_F$ . Since we have  
 525  $\mathbb{E}_{\zeta_k \sim D} \left[ \tilde{\nabla}_{\mathcal{R}} g(\mathbf{X}_k, \zeta_k) \right] \in \partial_{\mathcal{R}} f(\mathbf{X}_k)$ , the  $L$ -Lipschitz continuity of  $f$  on  $\mathcal{U}$ , [Theorem 1](#),  
 526 and [\(4.11\)](#) imply that

$$527 \quad \mathbb{E}_{\zeta_k \sim D} [f_{\lambda}(\mathbf{X}_{k+1})] \leq f_{\lambda}(\mathbf{X}_k) - \frac{\gamma_k}{2\lambda} \left( \frac{1}{\lambda} - (L + \tau) \right) \left\| \mathbf{X}_k - P_{\lambda f}(\mathbf{X}_k) \right\|_F^2 + \frac{\gamma_k^2 L^2}{2\lambda}.$$

529 Upon taking expectation with respect to all the previous realizations  $\zeta_0, \dots, \zeta_{k-1}$  on  
 530 both sides, we get

$$531 \quad \mathbb{E} [f_{\lambda}(\mathbf{X}_{k+1})] \leq \mathbb{E} [f_{\lambda}(\mathbf{X}_k)] - \frac{\gamma_k}{2\lambda} \left( \frac{1}{\lambda} - (L + \tau) \right) \mathbb{E} \left[ \left\| P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_k \right\|_F^2 \right] + \frac{\gamma_k^2 L^2}{2\lambda}.$$

533 The desired result then follows by rearranging the above inequality and recognizing  
 534 that  $\Theta(\mathbf{X}_k) = \lambda^{-1} \left\| P_{\lambda f}(\mathbf{X}_k) - \mathbf{X}_k \right\|_F$  (see [\(4.2\)](#)).  $\square$

535 Now, we can bound the iteration complexity of the Riemannian stochastic sub-  
 536 gradient method using [Proposition 2](#).

537 **THEOREM 3.** *Under the setting of [Proposition 2](#), suppose that we choose the con-*  
 538 *stant stepsize  $\gamma_k = \frac{1}{\sqrt{T+1}}$ ,  $k = 0, 1, \dots$  with  $T$  being the total number of iterations*  
 539 *and the algorithm returns  $\mathbf{X}_{\bar{k}}$  with  $\bar{k}$  sampled from  $\{1, \dots, T\}$  uniformly at random.*  
 540 *Then, we have*

$$541 \quad \mathbb{E} [\Theta^2(\mathbf{X}_{\bar{k}})] \leq \frac{1}{\lambda \left( \frac{1}{\lambda} - (L + \tau) \right)} \frac{2(f_{\lambda}(\mathbf{X}_0) - \min f_{\lambda}) + \frac{L^2}{\lambda}}{\sqrt{T+1}},$$

542 where the expectation is taken over all random choices by the algorithm.

543 *Proof.* By summing both sides of the relation in [Proposition 2](#) over  $k = 0, 1, \dots, T$ ,  
 544 we have

$$545 \quad \sum_{k=0}^T \gamma_k \mathbb{E} [\Theta^2(\mathbf{X}_k)] \leq \frac{2(f_{\lambda}(\mathbf{X}_0) - \min f_{\lambda}) + \frac{L^2}{\lambda} \sum_{k=0}^T \gamma_k^2}{\lambda \left( \frac{1}{\lambda} - (L + \tau) \right)}.$$

546 It follows that

$$547 \quad \sum_{k=0}^T \frac{\gamma_k}{\sum_{k=0}^T \gamma_k} \mathbb{E} [\Theta^2(\mathbf{X}_k)] \leq \frac{1}{\lambda \left( \frac{1}{\lambda} - (L + \tau) \right)} \frac{2(f_{\lambda}(\mathbf{X}_0) - \min f_{\lambda}) + \frac{L^2}{\lambda} \sum_{k=0}^T \gamma_k^2}{\sum_{k=0}^T \gamma_k}.$$

548 To complete the proof, it remains to substitute  $\gamma_k = \frac{1}{\sqrt{T+1}}$  into the above inequality  
 549 and note that the resulting LHS is exactly  $\mathbb{E} [\Theta^2(\mathbf{X}_{\bar{k}})]$  with the expectation being  
 550 taken with respect to  $\zeta_0, \dots, \zeta_{T-1}, \bar{k}$ .  $\square$

551 **5. Local Linear Convergence for Sharp Instances.** So far our discussion on  
 552 problem [\(1.1\)](#) does not assume any structure on the objective function  $f$  besides weak  
 553 convexity. However, many applications, such as those discussed in [Subsection 1.1](#),  
 554 give rise to weakly convex objective functions that are not arbitrary but have rather  
 555 concrete structure. It is thus natural to ask whether the methods we considered can  
 556 exploit this structure and provably achieve faster convergence rates than those estab-  
 557 lished in [Section 4](#). In this section, we introduce a regularity property of problem [\(1.1\)](#)

558 called *sharpness* and show that the Riemannian subgradient and incremental subgra-  
 559 dient methods will achieve a local linear convergence rate when applied to instances  
 560 of (1.1) that possess the sharpness property. Then, we will discuss in Section 7 how the  
 561 notion of sharpness captures, in a unified manner, the structure of both the dual prin-  
 562 cipal component pursuit (DPCP) formulation (1.3) of the robust subspace recovery  
 563 (RSR) problem and the single-column formulation (1.4) of the orthogonal dictionary  
 564 learning (DL) problem.

565 **5.1. Sharpness: Weak sharp minima.** To begin, let us introduce the notion  
 566 of a weak sharp minima set.

567 DEFINITION 1 (Sharpness; cf. [8, 28, 35]). *We say that  $\mathcal{X} \subseteq \text{St}(n, r)$  is a set of*  
 568 *weak sharp minima for the function  $h : \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$  with parameter  $\alpha > 0$  if there exists*  
 569 *a constant  $\rho > 0$  such that for any  $\mathbf{X} \in \mathcal{B} := \{\mathbf{X} \in \mathbb{R}^{n \times r} : \text{dist}(\mathbf{X}, \mathcal{X}) < \rho\} \cap \text{St}(n, r)$ ,*  
 570 *we have*

$$571 \quad h(\mathbf{X}) - h(\mathbf{Y}) \geq \alpha \text{dist}(\mathbf{X}, \mathcal{X})$$

572 *for all  $\mathbf{Y} \in \mathcal{X}$ , where  $\text{dist}(\mathbf{X}, \mathcal{X}) := \inf_{\mathbf{Y} \in \mathcal{X}} \|\mathbf{Y} - \mathbf{X}\|_F$ .*

573 From the definition, it is immediate that if  $\mathcal{X}$  is a set of weak sharp minima for  
 574  $h$ , then it is the set of minimizers of  $h$  over  $\mathcal{B}$ , and the function value grows linearly  
 575 with the distance to  $\mathcal{X}$ . Moreover, if  $h$  is continuous (e.g., when  $h$  is weakly convex),  
 576 then  $\mathcal{X}$  can be taken as closed.

577 Similar notions of sharpness play a fundamental role in establishing the linear  
 578 convergence of a host of methods for weakly convex minimization in the Euclidean  
 579 space. For instance, it is shown in [21] that the subgradient method with geometrically  
 580 diminishing stepsizes will converge linearly to the optimal solution set when applied  
 581 to minimize a sharp convex function. Later, the work [13] establishes a similar linear  
 582 convergence result for sharp weakly convex minimization. In the recent work [36], it  
 583 is shown that the incremental subgradient, proximal point, and prox-linear methods  
 584 will converge linearly when applied to minimize a sharp weakly convex function. In  
 585 this paper, we extend, for the first time, the above results to the manifold setting  
 586 by establishing the linear convergence of Riemannian subgradient-type methods for  
 587 minimizing a weakly convex function over the Stiefel manifold under the sharpness  
 588 property in Definition 1.

589 **5.2. Riemannian subgradient and incremental subgradient methods.**  
 590 Again, we will focus on analyzing the Riemannian incremental subgradient method.  
 591 The analysis of the Riemannian subgradient method will follow as a special case. We  
 592 first present the following result, which is crucial for our subsequent development.

593 PROPOSITION 3. *Under the setting of Proposition 1, for any  $\bar{\mathbf{X}} \in \text{St}(n, r)$ , we*  
 594 *have*

$$595 \quad \begin{aligned} \|\mathbf{X}_{k+1} - \bar{\mathbf{X}}\|_F^2 &\leq (1 + 2m\gamma_k(L + \tau)) \|\mathbf{X}_k - \bar{\mathbf{X}}\|_F^2 - 2m\gamma_k (f(\mathbf{X}_k) - f(\bar{\mathbf{X}})) \\ &\quad + m^2\gamma_k^2 L^2 + C(m)\gamma_k^3 L^2(L + \tau), \quad \forall k \geq 0, \end{aligned}$$

596 *where  $C(m) = \frac{1}{3}m(m-1)(2m-1)$ .*



597 *Proof.* According to [Lemma 1](#), for any  $\bar{\mathbf{X}} \in \text{St}(n, r)$ , we have

$$\begin{aligned}
& \|\mathbf{X}_{k,i} - \bar{\mathbf{X}}\|_F^2 \leq \|\mathbf{X}_{k,i-1} + \boldsymbol{\xi}_{k,i-1} - \bar{\mathbf{X}}\|_F^2 \\
& \stackrel{(i)}{\leq} \|\mathbf{X}_{k,i-1} - \bar{\mathbf{X}}\|_F^2 - 2\gamma_k \langle \tilde{\nabla}_{\mathcal{R}} f_i(\mathbf{X}_{k,i-1}), \mathbf{X}_{k,i-1} - \bar{\mathbf{X}} \rangle + \gamma_k^2 L^2 \\
& \stackrel{(ii)}{\leq} \|\mathbf{X}_{k,i-1} - \bar{\mathbf{X}}\|_F^2 - 2\gamma_k (f_i(\mathbf{X}_{k,i-1}) - f_i(\bar{\mathbf{X}})) \\
& \quad + \gamma_k(L + \tau) \|\mathbf{X}_{k,i-1} - \bar{\mathbf{X}}\|_F^2 + \gamma_k^2 L^2,
\end{aligned}
\tag{5.1}$$

599 where (i) follows from the fact that  $\|\tilde{\nabla}_{\mathcal{R}} f_i(\mathbf{X}_{k,i-1})\| \leq \|\tilde{\nabla} f_i(\mathbf{X}_{k,i-1})\|_F \leq L$  and (ii)  
600 is from [Theorem 1](#). Following the derivations of (4.7)–(4.9), we get

$$\begin{aligned}
& f_i(\bar{\mathbf{X}}) - f_i(\mathbf{X}_{k,i-1}) \leq (i-1)\gamma_k L^2 - (f_i(\mathbf{X}_k) - f_i(\bar{\mathbf{X}})), \\
& \|\mathbf{X}_{k,i-1} - \bar{\mathbf{X}}\|_F^2 \leq 2(i-1)^2 \gamma_k^2 L^2 + 2\|\mathbf{X}_k - \bar{\mathbf{X}}\|_F^2.
\end{aligned}$$

604 Substituting the above two upper bounds into (5.1) gives

$$\begin{aligned}
& \|\mathbf{X}_{k,i} - \bar{\mathbf{X}}\|_F^2 \leq \|\mathbf{X}_{k,i-1} - \bar{\mathbf{X}}\|_F^2 - 2\gamma_k (f_i(\mathbf{X}_k) - f_i(\bar{\mathbf{X}})) + 2\gamma_k(L + \tau) \|\mathbf{X}_k - \bar{\mathbf{X}}\|_F^2 \\
& \quad + (2i-1)\gamma_k^2 L^2 + 2(i-1)^2 \gamma_k^3 L^2 (L + \tau).
\end{aligned}$$

608 Upon summing both sides of the above inequality over  $i = 1, \dots, m$ , we obtain

$$\begin{aligned}
& \|\mathbf{X}_{k+1} - \bar{\mathbf{X}}\|_F^2 \leq (1 + 2m\gamma_k(L + \tau)) \|\mathbf{X}_k - \bar{\mathbf{X}}\|_F^2 - 2m\gamma_k (f(\mathbf{X}_k) - f(\bar{\mathbf{X}})) \\
& \quad + m^2 \gamma_k^2 L^2 + \frac{1}{3} m(m-1)(2m-1) \gamma_k^3 L^2 (L + \tau),
\end{aligned}$$

612 which completes the proof.  $\square$

613 In order for Riemannian subgradient-type methods to achieve linear convergence  
614 when solving sharp instances of problem (1.1), we need to choose the stepsizes ap-  
615 propriately. Motivated by previous works [13, 21, 36, 42, 49] on sharp weakly convex  
616 minimization in the Euclidean space, let us consider using geometrically diminishing  
617 stepsizes of the form  $\gamma_k = \beta^k \gamma_0$ ,  $k = 0, 1, \dots$ . Then, by applying [Proposition 3](#), we  
618 can establish the following local linear convergence result:

619 **THEOREM 4.** *Consider the setting of [Proposition 1](#). Suppose further that  $\mathcal{X}$  is a*  
620 *set of weak sharp minima for the objective function  $f$  in (1.1) with parameter  $\alpha >$*   
621 *0 over the set  $\mathcal{B}$  defined in [Definition 1](#). Let  $\{\mathbf{X}_k\}$  be the sequence generated by*  
622 *Riemannian incremental subgradient method (2.6) for solving problem (1.1), in which*  
623 *the initial point  $\mathbf{X}_0$  satisfies  $\text{dist}(\mathbf{X}_0, \mathcal{X}) < \min\left\{\frac{\alpha}{L+\tau}, \rho\right\}$  (so that  $\mathbf{X}_0 \in \mathcal{B}$ ) and the*  
624 *stepsizes satisfy  $\gamma_k = \beta^k \gamma_0$ ,  $k = 0, 1, \dots$ , where*

$$\gamma_0 < \min \left\{ \frac{2me_0(\alpha - (L + \tau)e_0)}{d(m)L^2}, \frac{e_0}{2m(\alpha - (L + \tau)e_0)} \right\},$$

$$\beta \in [\beta_{\min}, 1) \quad \text{with} \quad \beta_{\min} := \sqrt{1 + 2m \left( L + \tau - \frac{\alpha}{e_0} \right) \gamma_0 + \frac{d(m)L^2}{e_0^2} \gamma_0^2},$$

$$d(m) = \frac{5}{3}m^2 - m + \frac{1}{3}, \quad \text{and} \quad e_0 = \min \left\{ \max \left\{ \text{dist}(\mathbf{X}_0, \mathcal{X}), \frac{\alpha}{2(L + \tau)} \right\}, \rho \right\}.$$

630 Then, we have

$$\text{dist}(\mathbf{X}_k, \mathcal{X}) \leq \beta^k \cdot e_0, \quad \forall k \geq 0.$$

632 *Proof.* We first show that  $\beta_{\min} \in (0, 1)$  and  $\gamma_0 > 0$  are well defined. Towards  
 633 that end, note that  $\beta_{\min} = \sqrt{1 + v(\gamma_0)}$  with  $v(\gamma) = 2m \left( L + \tau - \frac{\alpha}{e_0} \right) \gamma + \frac{d(m)L^2}{e_0^2} \gamma^2$   
 634 being quadratic in  $\gamma$ . By definition of  $\gamma_0$ , we immediately have  $v(\gamma_0) < 0$ . Moreover,  
 635 the function  $\gamma \mapsto v(\gamma)$  attains its minimum at  $\bar{\gamma} = \frac{me_0(\alpha - (L + \tau)e_0)}{d(m)L^2}$  with value  $v(\bar{\gamma}) =$   
 636  $-\frac{m^2(\alpha - (L + \tau)e_0)^2}{d(m)L^2} > -\frac{\alpha^2}{L^2} \geq -1$ , where the first inequality is due to  $\frac{m^2}{d(m)} \leq 1$  for  $m \geq 1$   
 637 and  $e_0 < \frac{\alpha}{L + \tau}$ , and the second inequality is implied by the sharpness assumption  
 638 because  $\alpha \|\mathbf{X} - \bar{\mathbf{X}}\|_F \leq f(\mathbf{X}) - f(\bar{\mathbf{X}}) \leq L \|\mathbf{X} - \bar{\mathbf{X}}\|_F$  for any  $\mathbf{X} \in \mathcal{B}$  and  $\bar{\mathbf{X}} \in \mathcal{P}_{\mathcal{X}}(\mathbf{X})$ .  
 639 Hence, we have  $v(\gamma_0) \in (-1, 0)$ , which implies that  $\beta_{\min} \in (0, 1)$ . On the other hand,  
 640 since  $e_0 < \frac{\alpha}{L + \tau}$ , the upper bound on the initial stepsize  $\gamma_0$  is positive. It follows that  
 641  $\gamma_0$  is well defined.

642 We now prove the theorem by induction on  $k$ . The base case  $k = 0$  follows directly  
 643 from the definition of  $e_0$ . For the inductive step, suppose that  $\text{dist}(\mathbf{X}_k, \mathcal{X}) \leq \beta^k \cdot e_0$   
 644 for some  $k \geq 0$ . Note that this implies  $\mathbf{X}_k \in \mathcal{B}$ . Let  $\bar{\mathbf{X}} \in \mathcal{P}_{\mathcal{X}}(\mathbf{X}_k)$ . Clearly, we  
 645 have  $\text{dist}(\mathbf{X}_k, \mathcal{X}) = \|\mathbf{X}_k - \bar{\mathbf{X}}\|_F$  and  $\text{dist}(\mathbf{X}_{k+1}, \mathcal{X}) \leq \|\mathbf{X}_{k+1} - \bar{\mathbf{X}}\|_F$ . Hence, by  
 646 **Proposition 3**, the sharpness assumption, and the fact that  $\gamma_k \leq \gamma_0$  for  $k = 0, 1, \dots$ ,  
 647 we get

$$648 \quad (5.2) \quad \text{dist}^2(\mathbf{X}_{k+1}, \mathcal{X}) \leq (1 + 2m\gamma_0(L + \tau)) \text{dist}^2(\mathbf{X}_k, \mathcal{X}) - 2m\gamma_k\alpha \text{dist}(\mathbf{X}_k, \mathcal{X}) \\ + m^2\gamma_k^2L^2 + C(m)\gamma_k^3L^2(L + \tau).$$

649 Observe that the RHS of the above recursion is quadratic in  $\text{dist}(\mathbf{X}_k, \mathcal{X})$ . By definition  
 650 of  $\gamma_0$ , we have  $\gamma_0 < \frac{e_0}{2m(\alpha - (L + \tau)e_0)}$  and hence  $\frac{2m\gamma_0\alpha}{1 + 2m\gamma_0(L + \tau)} < e_0$ . This implies that the  
 651 RHS of (5.2) achieves its maximum when  $\text{dist}(\mathbf{X}_k, \mathcal{X}) = \beta^k \cdot e_0$ . Since  $\text{dist}(\mathbf{X}_k, \mathcal{X}) \leq$   
 652  $\beta^k \cdot e_0$  by the inductive hypothesis, plugging  $\gamma_k = \beta^k\gamma_0$  and  $\text{dist}(\mathbf{X}_k, \mathcal{X}) = \beta^k \cdot e_0$  into  
 653 (5.2) yields

$$654 \quad (5.3) \quad \text{dist}^2(\mathbf{X}_{k+1}, \mathcal{X}) \\ \leq \beta^{2k}e_0^2 \left[ 1 + 2m \left( L + \tau - \frac{\alpha}{e_0} \right) \gamma_0 + L^2 \left( \frac{m^2 + C(m)\gamma_0(L + \tau)}{e_0^2} \right) \gamma_0^2 \right].$$

655 Note that  $\gamma_0 < \frac{2m(\alpha e_0 - (L + \tau)e_0^2)}{d(m)L^2} \leq \frac{m\alpha^2}{2d(m)L^2(L + \tau)} < \frac{1}{m(L + \tau)}$ . It then follows from (5.3)  
 656 that

$$657 \quad \text{dist}^2(\mathbf{X}_{k+1}, \mathcal{X}) \leq \beta^{2k}e_0^2 \left[ 1 + 2m \left( L + \tau - \frac{\alpha}{e_0} \right) \gamma_0 + \frac{d(m)L^2}{e_0^2} \gamma_0^2 \right] \leq \beta^{2(k+1)}e_0^2.$$

658 This completes the inductive step and hence the proof of **Theorem 4**.  $\square$

659 From **Theorem 4**, we see that in order to achieve a fast linear convergence rate, one  
 660 should choose an appropriate  $\gamma_0$  so that the minimum decay factor  $\beta_{\min}$  is as small as  
 661 possible. By minimizing  $\beta_{\min}$  with respect to  $\gamma_0$ , we see that the theoretical minimum  
 662 value of  $\beta_{\min}$  is  $\sqrt{1 - \frac{m^2(\alpha - (L + \tau)e_0)^2}{d(m)L^2}}$ , which is attained at  $\gamma_0 = \bar{\gamma}_0 = \frac{me_0(\alpha - (L + \tau)e_0)}{d(m)L^2}$ .  
 663 This suggests that subject to the requirement in **Theorem 4**, the initial stepsize  $\gamma_0$   
 664 should be set as close to  $\bar{\gamma}_0$  as possible. As an illustration, consider the case where  
 665 the sharpness property holds globally over the Stiefel manifold (i.e.,  $\mathcal{B} = \text{St}(n, r)$  in  
 666 **Definition 1**). Then, the parameter  $\rho$  can be set as large as possible. In this case, we  
 667 have  $e_0 = \max \left\{ \text{dist}(\mathbf{X}_0, \mathcal{X}), \frac{\alpha}{2(L + \tau)} \right\}$ , and the condition on  $\gamma_0$  in **Theorem 4** becomes  
 668  $\gamma_0 < \frac{2me_0(\alpha - (L + \tau)e_0)}{d(m)L^2}$ . This implies that we can choose  $\gamma_0 = \bar{\gamma}_0$  to obtain the smallest

669 possible  $\beta_{\min}$ . Note, however, that the larger the initialization error  $\text{dist}(\mathbf{X}_0, \mathcal{X})$ , the  
 670 larger the minimum decay factor  $\beta_{\min}$ . In particular, from the expression for  $\beta_{\min}$   
 671 above, we see that  $\beta_{\min}$  approaches 1 as  $\text{dist}(\mathbf{X}_0, \mathcal{X})$  approaches its maximum  $\frac{\alpha}{L+\tau}$ .

672 We end this section by comparing the sharpness property with the *Riemannian*  
 673 *regularity condition* used in [3] and [69] for orthogonal DL and RSR, respectively. For  
 674 a target solution set  $\mathcal{X}$ , the Riemannian regularity condition stipulates the existence  
 675 of a constant  $\kappa > 0$  such that  $\langle \tilde{\nabla}_{\mathcal{R}} f(\mathbf{X}), \mathbf{X} - \mathbf{Y} \rangle \geq \kappa \text{dist}(\mathbf{X}, \mathcal{X})$  for all  $\mathbf{X}$  in a  
 676 *small neighborhood* of  $\mathcal{X}$  and  $\mathbf{Y} \in \mathcal{P}_{\mathcal{X}}(\mathbf{X})$ . This condition is motivated by the need  
 677 to bound the inner product term on the LHS in the convergence analysis of the  
 678 Riemannian subgradient method; see (5.1) with  $f_i = f$  and  $\mathbf{X}_{k,i-1} = \mathbf{X}_k$ . Informally,  
 679 the Riemannian regularity condition is a combination of the Riemannian subgradient  
 680 inequality in Theorem 1 and the sharpness property in Definition 1. However, the  
 681 tangling of these two elements potentially restricts the applicability of the Riemannian  
 682 regularity condition. In particular, since the Riemannian regularity condition can only  
 683 hold locally, it cannot be used to establish global convergence and iteration complexity  
 684 results for the Riemannian subgradient method.

## 685 6. Extension to Optimization over a Compact Embedded Submanifold.

686 There is of course no conceptual difficulty in adapting the Riemannian subgradient-  
 687 type methods in Subsection 2.2 to minimize weakly convex functions over more general  
 688 manifolds. All that is needed is an efficiently computable retraction on the manifold of  
 689 interest. In this section, let us briefly demonstrate how the machinery developed in the  
 690 previous sections can be extended to study the convergence behavior of Riemannian  
 691 subgradient-type methods when the manifold in question is compact and defined by  
 692 a certain smooth mapping.

693 **Riemannian subgradient inequality.** Our starting point is the following gen-  
 694 eralization of the Riemannian subgradient inequality in Theorem 1, which applies  
 695 to restrictions of weakly convex functions on a class of compact embedded subman-  
 696 ifolds of the Euclidean space. Some examples of manifolds in this class include the  
 697 generalized Stiefel manifold, oblique manifold, and symplectic manifold; see, e.g., [2].

698 **COROLLARY 1.** *Let  $\mathcal{M}$  be a compact submanifold of  $\mathbb{R}^p$  given by  $\mathcal{M} = \{\mathbf{X} \in \mathbb{R}^p :$   
 699  $F(\mathbf{X}) = \mathbf{0}\}$ , where  $F : \mathbb{R}^p \rightarrow \mathbb{R}^q$  is a smooth mapping whose derivative  $DF(\mathbf{X})$  at  $\mathbf{X}$   
 700 has full row rank for all  $\mathbf{X} \in \mathcal{M}$ . Then, for any weakly convex function  $h : \mathbb{R}^p \rightarrow \mathbb{R}$ ,  
 701 there exists a constant  $c > 0$  such that*

$$702 \quad h(\mathbf{Y}) \geq h(\mathbf{X}) + \langle \tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}), \mathbf{Y} - \mathbf{X} \rangle - c \|\mathbf{Y} - \mathbf{X}\|_F^2$$

703 for all  $\mathbf{X}, \mathbf{Y} \in \mathcal{M}$  and  $\tilde{\nabla}_{\mathcal{R}} h(\mathbf{X}) \in \partial_{\mathcal{R}} h(\mathbf{X})$ .

704 *Proof.* By our assumptions on  $F$  and [2, Equation (3.19)], we have  $T_{\mathbf{X}}\mathcal{M} =$   
 705  $\ker(DF(\mathbf{X}))$ , where  $\ker(T)$  denotes the kernel of the operator  $T$ . Thus, the projector  
 706  $\mathcal{P}_{T_{\mathbf{X}}\mathcal{M}}^{\perp}$  is given by  $DF(\mathbf{X})^{\top}(DF(\mathbf{X})DF(\mathbf{X})^{\top})^{-1}DF(\mathbf{X})$ . Following the proof of  
 707 Theorem 1, we need to bound

$$708 \quad \langle \tilde{\nabla} h(\mathbf{X}), \mathcal{P}_{T_{\mathbf{X}}\mathcal{M}}^{\perp}(\mathbf{Y} - \mathbf{X}) \rangle \geq -\|\tilde{\nabla} h(\mathbf{X})\|_F \cdot \|\mathcal{P}_{T_{\mathbf{X}}\mathcal{M}}^{\perp}(\mathbf{Y} - \mathbf{X})\|_F$$

$$709 \quad = -\|\tilde{\nabla} h(\mathbf{X})\|_F \cdot \|DF(\mathbf{X})^{\top}(DF(\mathbf{X})DF(\mathbf{X})^{\top})^{-1}DF(\mathbf{X})(\mathbf{Y} - \mathbf{X})\|_F$$

$$710 \quad \geq -\|\tilde{\nabla} h(\mathbf{X})\|_F \cdot \max_{\mathbf{X} \in \mathcal{M}} \|DF(\mathbf{X})^{\top}(DF(\mathbf{X})DF(\mathbf{X})^{\top})^{-1}\|_F \cdot \|DF(\mathbf{X})(\mathbf{Y} - \mathbf{X})\|_F.$$

711  
 712 Since  $h$  is weakly convex on  $\mathbb{R}^p$ , it is Lipschitz continuous on any bounded open convex  
 713 set  $\mathcal{U}$  that contains  $\mathcal{M}$ . Thus, the term  $\|\tilde{\nabla} h(\mathbf{X})\|_F$  is bounded above. Moreover, the

compactness of  $\mathcal{M}$  implies that the term  $\max_{\mathbf{X} \in \mathcal{M}} \|DF(\mathbf{X})^\top (DF(\mathbf{X})DF(\mathbf{X})^\top)^{-1}\|_F$  is also bounded above. Lastly, observe that  $F(\mathbf{Y}) = F(\mathbf{X}) + DF(\mathbf{X})(\mathbf{Y} - \mathbf{X}) + \mathcal{O}(\|\mathbf{Y} - \mathbf{X}\|_F^2)$  by Taylor's theorem and  $F(\mathbf{X}) = F(\mathbf{Y}) = \mathbf{0}$  whenever  $\mathbf{X}, \mathbf{Y} \in \mathcal{M}$ . Hence, we have  $\|DF(\mathbf{X})(\mathbf{Y} - \mathbf{X})\|_F = \mathcal{O}(\|\mathbf{Y} - \mathbf{X}\|_F^2)$ . Putting these together, we conclude that  $\left| \left\langle \tilde{\nabla} h(\mathbf{X}), \mathcal{P}_{\mathbf{T}_{\mathbf{X}}\mathcal{M}}^\perp(\mathbf{Y} - \mathbf{X}) \right\rangle \right| = \mathcal{O}(\|\mathbf{Y} - \mathbf{X}\|_F^2)$ . The rest of the argument is similar to that in the proof of [Theorem 1](#).  $\square$

**General retractions.** The notion of retraction introduced in [Subsection 2.1](#) for the Stiefel manifold can be easily adapted to that for general manifolds. Specifically, a retraction on the manifold  $\mathcal{M}$  is a smooth map  $\text{Retr} : \text{T}\mathcal{M} \rightarrow \mathcal{M}$  from the tangent bundle  $\text{T}\mathcal{M}$  onto the manifold  $\mathcal{M}$  that satisfies  $\text{Retr}_{\mathbf{X}}(\mathbf{0}) = \mathbf{X}$  and  $D\text{Retr}_{\mathbf{X}}(\mathbf{0}) = \text{Id}$  for all  $\mathbf{X} \in \mathcal{M}$ . Unlike the polar decomposition-based retraction on the Stiefel manifold, a general retraction may not have the Lipschitz-like property in [Lemma 1](#). Nevertheless, a retraction on a compact submanifold  $\mathcal{M}$  satisfies a *second-order boundedness* property [[7](#)]; i.e., there exists a constant  $b \geq 0$  such that for all  $\mathbf{X} \in \mathcal{M}$  and  $\boldsymbol{\xi} \in \text{T}_{\mathbf{X}}\mathcal{M}$ ,

$$\|\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) - \mathbf{X} - \boldsymbol{\xi}\|_F \leq b\|\boldsymbol{\xi}\|_F^2.$$

This allows us to replace the result in [Lemma 1](#) by

$$\begin{aligned} \|\text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) - \bar{\mathbf{X}}\|_F &= \|(\mathbf{X} + \boldsymbol{\xi}) - \bar{\mathbf{X}} + \text{Retr}_{\mathbf{X}}(\boldsymbol{\xi}) - (\mathbf{X} + \boldsymbol{\xi})\|_F \\ &\leq \|\mathbf{X} + \boldsymbol{\xi} - \bar{\mathbf{X}}\|_F + b\|\boldsymbol{\xi}\|_F^2, \end{aligned}$$

which holds for any  $\mathbf{X}, \bar{\mathbf{X}} \in \mathcal{M}$  and  $\boldsymbol{\xi} \in \text{T}_{\mathbf{X}}\mathcal{M}$ . Although the above inequality has the extra term  $b\|\boldsymbol{\xi}\|_F^2$ , it can still be used to establish convergence guarantees (with slightly worse constants) for the Riemannian subgradient-type methods considered in [Subsection 2.2](#). Specifically, by following the analyses in [Sections 4](#) and [5](#), we can show that for problem [\(1.1\)](#) with the Stiefel manifold  $\text{St}(n, r)$  being replaced by a manifold of the type considered in [Corollary 1](#), the iteration complexity of Riemannian subgradient-type methods for computing an  $\varepsilon$ -nearly stationary point is  $\mathcal{O}(\varepsilon^{-4})$ , and the Riemannian subgradient and incremental subgradient methods will achieve a local linear convergence rate if the instance satisfies the sharpness property in [Definition 1](#).

**7. Applications and Numerical Results.** In this section, we apply the Riemannian subgradient-type methods in [Subsection 2.2](#) to solve the RSR and orthogonal DL problems. As described in [Section 1](#), the objective functions of both problems are weakly convex. Thus, [Theorem 2](#) and [Theorem 3](#) ensure that the Riemannian subgradient-type methods with arbitrary initialization will have a global convergence rate of  $\mathcal{O}(k^{-1/4})$  when utilized to solve those problems. We also discuss the sharpness properties of the RSR and orthogonal DL problems. For reproducible research, our code for generating the numerical results can be found at

<https://github.com/lixiao0982/Riemannian-subgradient-methods>

**7.1. Robust subspace recovery (RSR).** We begin with the DPCP formulation [\(1.3\)](#) of the RSR problem, which has a relatively simpler form than the least absolute deviation (LAD) formulation [\(1.2\)](#). Recall that the objective function in [\(1.3\)](#) takes the form  $\text{St}(n, r) \ni \mathbf{X} \mapsto f(\mathbf{X}) = \frac{1}{m} \sum_{i=1}^m \|\tilde{\mathbf{y}}_i^\top \mathbf{X}\|_2$ , where  $\tilde{\mathbf{y}}_i \in \mathbb{R}^n$  ( $i = 1, \dots, m$ ) denotes the  $i$ -th column of  $\tilde{\mathbf{Y}} = [\mathbf{Y} \ \mathbf{O}] \boldsymbol{\Gamma} \in \mathbb{R}^{n \times m}$ , the columns  $\mathbf{y}_i$  of  $\mathbf{Y} \in \mathbb{R}^{n \times m_1}$  form inlier points spanning a  $d$ -dimensional subspace  $\mathcal{S}$  with  $r = n - d$ , the columns  $\mathbf{o}_i$  of  $\mathbf{O} \in \mathbb{R}^{n \times m_2}$  form outlier points, and  $\boldsymbol{\Gamma} \in \mathbb{R}^{m \times m}$  is an unknown

759 permutation. Note that  $f$  is rotationally invariant; i.e.,  $f(\mathbf{X}) = f(\mathbf{X}\mathbf{R})$  for any  
760  $\mathbf{X} \in \text{St}(n, r)$  and  $\mathbf{R} \in \text{St}(r, r)$ .

761 **Sharpness.** Let  $\mathbf{S}^\perp \in \text{St}(n, r)$  be an orthonormal basis of  $\mathcal{S}^\perp$ . Since the goal  
762 of DPCP is to find an orthonormal basis (but not necessary  $\mathbf{S}^\perp$ ) for  $\mathcal{S}^\perp$ , we are  
763 interested in the elements in the set  $\mathcal{X} = \{\mathbf{S}^\perp \mathbf{R} \in \mathbb{R}^{n \times r} : \mathbf{R} \in \text{St}(r, r)\}$ . Due to  
764 rotation invariance,  $f$  is constant on  $\mathcal{X}$ . To study the sharpness property of problem  
765 (1.3), let us introduce two quantities that reflect how well distributed the inliers and  
766 outliers are:

$$767 \quad (7.1) \quad c_{\mathbf{Y}, \min} := \frac{1}{m_1} \inf_{\substack{\mathbf{D} \in \mathbb{R}^{n \times \ell} \\ \|\mathbf{D}\|_F=1, \text{col}(\mathbf{D}) \subseteq \mathcal{S}}} \sum_{i=1}^{m_1} \|\mathbf{y}_i^\top \mathbf{D}\|_2,$$

$$768 \quad (7.2) \quad c_{\mathbf{O}, \max} := \frac{1}{m_2} \sup_{\substack{\mathbf{B} \in \mathbb{R}^{n \times r} \\ \|\mathbf{B}\|_F=1}} \sum_{i=1}^{m_2} \|\mathbf{o}_i^\top \mathbf{B}\|_2.$$

770 Here,  $\ell = \min\{d, r\}$  and  $\text{col}(\mathbf{D})$  denotes the column space of  $\mathbf{D}$ . In a nutshell, larger  
771 values of  $c_{\mathbf{Y}, \min}$  (respectively, smaller values of  $c_{\mathbf{O}, \max}$ ) correspond to a more uniform  
772 distribution of inliers (respectively, outliers). As the following proposition shows, the  
773 quantities  $c_{\mathbf{Y}, \min}$  and  $c_{\mathbf{O}, \max}$  can be used to capture the sharpness property of the  
774 DPCP formulation (1.3).

775 **PROPOSITION 4.** *Suppose that  $m_2 c_{\mathbf{O}, \max} \leq \frac{1}{2} m_1 c_{\mathbf{Y}, \min}$ . Then, the DPCP formu-*  
776 *lation (1.3) has  $\mathcal{X}$  as a set of weak sharp minima with parameter  $\alpha = \frac{1}{m} (\frac{1}{2} m_1 c_{\mathbf{Y}, \min} -$   
777  $m_2 c_{\mathbf{O}, \max}) > 0$  over the set  $\mathcal{B} = \text{St}(n, r)$ ; i.e.,*

$$778 \quad f(\mathbf{X}) - f(\mathbf{S}^\perp) \geq \alpha \text{dist}(\mathbf{X}, \mathcal{X}), \quad \forall \mathbf{X} \in \text{St}(n, r).$$

779 *Proof.* Let  $\mathbf{X} \in \text{St}(n, r)$  be arbitrary. For any  $\mathbf{X}^* \in \mathcal{P}_{\mathcal{X}}(\mathbf{X})$ , we have  $f(\mathbf{X}^*) =$   
780  $f(\mathbf{S}^\perp)$  and

$$781 \quad (7.3) \quad \begin{aligned} f(\mathbf{X}) - f(\mathbf{S}^\perp) &= \frac{1}{m} \sum_{i=1}^m \|\tilde{\mathbf{y}}_i^\top \mathbf{X}\|_2 - \frac{1}{m} \sum_{i=1}^m \|\tilde{\mathbf{y}}_i^\top \mathbf{S}^\perp\|_2 \\ &= \frac{1}{m} \sum_{i=1}^{m_1} \|\mathbf{y}_i^\top \mathbf{X}\|_2 + \frac{1}{m} \left( \sum_{i=1}^{m_2} \|\mathbf{o}_i^\top \mathbf{X}\|_2 - \sum_{i=1}^{m_2} \|\mathbf{o}_i^\top \mathbf{S}^\perp\|_2 \right), \end{aligned}$$

782 where  $\mathbf{y}_i$  (respectively,  $\mathbf{o}_i$ ) is the  $i$ -th column of  $\mathbf{Y}$  (respectively,  $\mathbf{O}$ ), and the second  
783 line follows because the inliers  $\{\mathbf{y}_i\}_{i=1}^{m_1}$  are orthogonal to  $\mathbf{S}^\perp$ . Now, let us derive lower  
784 bounds for the two terms on the right-hand side separately.

785 For the first term, let  $\mathcal{S} \in \mathbb{R}^{n \times d}$  be an orthonormal basis of the subspace  $\mathcal{S}$ . By  
786 projecting  $\mathbf{X}$  onto the orthogonal subspaces  $\mathcal{S}$  and  $\mathcal{S}^\perp$ , we have

$$787 \quad (7.4) \quad \mathbf{X} = \mathbf{S}\mathbf{S}^\top \mathbf{X} + \mathbf{S}^\perp (\mathbf{S}^\perp)^\top \mathbf{X}.$$

788 For  $i = 1, \dots, r$ , let  $\phi_i = \arccos(\sigma_i((\mathbf{S}^\perp)^\top \mathbf{X}))$  be the  $i$ -th smallest principal angle  
789 between the subspaces spanned by  $\mathbf{X}$  and  $\mathbf{S}^\perp$ , where  $\sigma_i(\cdot)$  denotes the  $i$ -th largest  
790 singular value [51]. Then, we can write  $(\mathbf{S}^\perp)^\top \mathbf{X} = \mathbf{U} \cos(\Phi) \mathbf{W}^\top$ , where  $\cos(\Phi) \in$   
791  $\mathbb{R}^{r \times r}$  is the diagonal matrix with  $\cos(\phi_1) \geq \dots \geq \cos(\phi_r)$  on its diagonal and  $\mathbf{U} \in$   
792  $\mathbb{R}^{r \times r}$ ,  $\mathbf{W} \in \mathbb{R}^{r \times r}$  are orthogonal matrices. On the other hand, according to [29,  
793 Theorem 2.7], the  $i$ -th smallest principal angle between the subspaces spanned by  
794  $\mathbf{X}$  and  $\mathcal{S}$  is  $\tilde{\phi}_i = \frac{\pi}{2} - \phi_{r-i+1}$ , where  $i = 1, \dots, \ell$  with  $\ell = \min\{d, r\}$ . Hence, we

795 can write  $\mathbf{S}^\top \mathbf{X} = \mathbf{V} \sin(\tilde{\Phi}) \mathbf{H}^\top$ , where  $\sin(\tilde{\Phi}) \in \mathbb{R}^{\ell \times \ell}$  is the diagonal matrix with  
 796  $\sin(\phi_r) \geq \dots \geq \sin(\phi_{r-\ell+1})$  on its diagonal and  $\mathbf{V} \in \mathbb{R}^{d \times \ell}$ ,  $\mathbf{H} \in \mathbb{R}^{r \times \ell}$  are orthogonal  
 797 matrices. These, together with (7.4), yield  $\mathbf{X} = \mathbf{S} \mathbf{V} \sin(\tilde{\Phi}) \mathbf{H}^\top + \mathbf{S}^\perp \mathbf{U} \cos(\Phi) \mathbf{W}^\top$ .  
 798 Hence, we can bound

$$\begin{aligned}
 & \sum_{i=1}^{m_1} \|\mathbf{y}_i^\top \mathbf{X}\|_2 = \sum_{i=1}^{m_1} \|\mathbf{y}_i^\top \mathbf{S} \mathbf{V} \sin(\tilde{\Phi})\|_2 \\
 799 \quad (7.5) \quad & = \|\mathbf{S} \mathbf{V} \sin(\tilde{\Phi})\|_F \sum_{i=1}^{m_1} \left\| \mathbf{y}_i^\top \frac{\mathbf{S} \mathbf{V} \sin(\tilde{\Phi})}{\|\mathbf{S} \mathbf{V} \sin(\tilde{\Phi})\|_F} \right\|_2 \geq m_1 c_{\mathbf{Y}, \min} \|\sin(\tilde{\Phi})\|_F,
 \end{aligned}$$

800 where  $c_{\mathbf{Y}, \min}$  is defined in (7.1). On the other hand, observe that

$$\begin{aligned}
 & \text{dist}^2(\mathbf{X}, \mathcal{X}) = \underset{\mathbf{R} \in \text{St}(r, r)}{\text{minimize}} \|\mathbf{X} - \mathbf{S}^\perp \mathbf{R}\|_F^2 = \|\mathbf{X} - \mathbf{S}^\perp \mathbf{U} \mathbf{W}^\top\|_F^2 \\
 801 \quad (7.6) \quad & = 2r - 2 \text{trace}(\cos(\Phi)) = 2 \sum_{i=1}^{\ell} (1 - \cos(\phi_i)) = 4 \sum_{i=1}^{\ell} \sin^2(\phi_i/2) \leq 4 \|\sin(\tilde{\Phi})\|_F^2,
 \end{aligned}$$

802 where the second equality follows from the solution to the orthogonal Procrustes  
 803 problem [48] and the fourth equality utilizes the fact that the number of nonzero  
 804 principal angles in  $\Phi$  is at most  $\ell = \min\{d, r\}$  [29, Theorem 2.7]. Combining (7.5)  
 805 and (7.6) gives

$$806 \quad (7.7) \quad \sum_{i=1}^{m_1} \|\mathbf{y}_i^\top \mathbf{X}\|_2 \geq \frac{1}{2} m_1 c_{\mathbf{Y}, \min} \text{dist}(\mathbf{X}, \mathcal{X}).$$

807 Now, let us consider the second term on the right-hand side of (7.3). Let  $\mathbf{R}^* =$   
 808  $\text{argmin}_{\mathbf{R} \in \text{St}(r, r)} \|\mathbf{X} - \mathbf{S}^\perp \mathbf{R}\|_F$ . Then, we have

$$\begin{aligned}
 & \left| \sum_{i=1}^{m_2} \|\mathbf{o}_i^\top \mathbf{X}\|_2 - \|\mathbf{o}_i^\top \mathbf{S}^\perp \mathbf{R}\|_2 \right| \leq \sum_{i=1}^{m_2} \left| \|\mathbf{o}_i^\top \mathbf{X}\|_2 - \|\mathbf{o}_i^\top \mathbf{S}^\perp \mathbf{R}^*\|_2 \right| \\
 809 \quad (7.8) \quad & \leq \sum_{i=1}^{m_2} \|\mathbf{o}_i^\top (\mathbf{X} - \mathbf{S}^\perp \mathbf{R}^*)\|_2 = \|\mathbf{X} - \mathbf{S}^\perp \mathbf{R}^*\|_F \sum_{i=1}^{m_2} \left\| \mathbf{o}_i^\top \frac{\mathbf{X} - \mathbf{S}^\perp \mathbf{R}^*}{\|\mathbf{X} - \mathbf{S}^\perp \mathbf{R}^*\|_F} \right\|_2 \\
 & \leq m_2 c_{\mathbf{O}, \max} \text{dist}(\mathbf{X}, \mathcal{X}),
 \end{aligned}$$

810 where  $c_{\mathbf{O}, \max}$  is defined in (7.2).

811 By plugging (7.7) and (7.8) into (7.3), the desired result follows.  $\square$

812 The requirement  $m_2 c_{\mathbf{O}, \max} \leq \frac{1}{2} m_1 c_{\mathbf{Y}, \min}$  in Proposition 4 determines the number  
 813 of outliers that can be tolerated. Now, let us give probabilistic estimates of the  
 814 quantities  $c_{\mathbf{Y}, \min}$  and  $c_{\mathbf{O}, \max}$  under the popular *Haystack model* (see, e.g., [34, 40, 68])  
 815 of the input data. The model stipulates that the inliers  $\{\mathbf{y}_i\}_{i=1}^{m_1}$  are i.i.d. according  
 816 to the Gaussian distribution  $\mathcal{N}(\mathbf{0}, \frac{1}{d} \mathcal{P}_{\mathcal{S}})$  with  $\mathcal{P}_{\mathcal{S}}$  being the orthogonal projector onto  
 817 the  $d$ -dimensional subspace  $\mathcal{S}$ , while the outliers  $\{\mathbf{o}_i\}_{i=1}^{m_2}$  are i.i.d. according to the  
 818 Gaussian distribution  $\mathcal{N}(\mathbf{0}, \frac{1}{n} \mathbf{I}_n)$ .

819 LEMMA 2. *Under the Haystack model, the event*

$$820 \quad c_{\mathbf{Y}, \min} \geq \sqrt{\frac{2}{d\pi}} - \sqrt{\frac{8\ell}{m_1}} - \frac{c_1}{\sqrt{m_1}}$$

821 will hold with probability at least  $1 - 2\exp(-\frac{c_1^2 d}{2})$  for some constant  $c_1 > 0$ , where  
 822  $\ell = \min\{d, r\}$ . Moreover, the event

$$823 \quad c_{\mathcal{O}, \max} \leq \frac{1}{\sqrt{n}} + \sqrt{\frac{8r}{m_2}} + \frac{c_2}{\sqrt{m_2}}$$

824 will hold with probability at least  $1 - 2\exp(-\frac{c_2^2 n}{2})$  for some constant  $c_2 > 0$ .

825 The proof of Lemma 2 can be found in Appendix A. Lemma 2 implies that under  
 826 the Haystack model, if the numbers of inliers  $m_1$  and outliers  $m_2$  satisfy  $m_1 \gtrsim d\ell$  and  
 827  $m_2 \gtrsim nr$ , then we will have  $c_{\mathcal{Y}, \min} \gtrsim \frac{1}{\sqrt{d}}$  and  $c_{\mathcal{O}, \max} \lesssim \frac{1}{\sqrt{n}}$  with high probability.  
 828 Combining Theorem 4, Proposition 4, and Lemma 2, we see that as long as

$$829 \quad (7.9) \quad m_2 \lesssim \sqrt{\frac{n}{d}} m_1$$

830 so that  $m_2 c_{\mathcal{O}, \max} \leq \frac{1}{2} m_1 c_{\mathcal{Y}, \min}$ , the Riemannian subgradient and incremental sub-  
 831 gradient methods with geometrically diminishing stepsizes and a proper initialization  
 832 will converge linearly to an orthonormal basis of  $\mathcal{S}^\perp$ . One initialization strategy is to  
 833 take the bottom eigenvectors of  $\tilde{\mathbf{Y}}\tilde{\mathbf{Y}}^\top$  [40, 70].

834 It is instructive to compare the bound (7.9) with those in the literature. When  
 835  $d = \mathcal{O}(1)$  or when both  $d$  and  $r$  are on the order of  $n$ , our bound (7.9) holds in the  
 836 regime  $m \gtrsim n^2$ . In this regime, the algorithms proposed in [40, 68] can recover  $\mathcal{S}$  as  
 837 long as  $m_2 \lesssim \frac{\sqrt{n(n-d)}}{d} m_1$ ; see [40, Section 5.5.2]. Such a bound is superior to ours  
 838 when  $d = \mathcal{O}(1)$  but is comparable when both  $d$  and  $r$  are on the order of  $n$ . When  
 839  $r = \mathcal{O}(1)$ , our bound (7.9) holds in the regime  $m \gtrsim n$ , which is superior to the bound  
 840  $m_2 \lesssim \frac{n-d}{d} m_1$  established in [34, 67] for the same regime. We remark that there are  
 841 other works [32, 62, 70] studying the RSR problem. However, they differ from our work  
 842 in that they either assume different data models, require additional data structures,  
 843 or consider the asymptotic setting  $m \rightarrow \infty$ .

844 To further demonstrate the power of Proposition 4, let us use it to establish the  
 845 sharpness property of the LAD formulation (1.2). To begin, let  $g$  be the objective  
 846 function in (1.2) and  $\mathbf{S} \in \text{St}(n, d)$  be an orthonormal basis of  $\mathcal{S}$ . We are interested in  
 847 the set  $\mathcal{D} = \{\mathbf{S}\mathbf{R} : \mathbf{R} \in \text{St}(d, d)\}$ , whose elements are different orthonormal bases of  $\mathcal{S}$ .  
 848 Now, observe that for any  $\mathbf{X} \in \text{St}(n, r)$ , we can find an orthonormal basis  $\mathbf{Z} \in \text{St}(n, d)$   
 849 of  $\text{col}(\mathbf{X})^\perp$ , and vice versa, such that  $f(\mathbf{X}) = g(\mathbf{Z})$  (recall that  $f$  is the objective  
 850 function in (1.3)). Hence, Proposition 4 asserts that  $g(\mathbf{Z}) - g(\mathbf{S}) \geq \alpha \text{dist}(\mathbf{X}, \mathcal{X})$ . By  
 851 invoking [29, Theorem 2.7], we obtain  $\text{dist}(\mathbf{X}, \mathcal{X}) = \text{dist}(\mathbf{Z}, \mathcal{D})$ , which shows that  $\mathcal{D}$   
 852 is a set of weak sharp minima with parameter  $\alpha$  over the set  $\mathcal{B} = \text{St}(n, d)$ .

853 **Experiments.** We first randomly sample a subspace  $\mathcal{S}$  with co-dimension  $r = 10$   
 854 in ambient dimension  $n = 100$ . We then generate  $m_1 = 1500$  inliers uniformly at  
 855 random from the unit sphere in  $\mathcal{S}$  and  $m_2 = 3500$  outliers uniformly at random  
 856 from the unit sphere in  $\mathbb{R}^n$ . We generate a standard Gaussian random vector and  
 857 use it to initialize all the algorithms, as such an initialization provides comparable  
 858 performance with the carefully designed initialization in [40, 70]. The numerical results  
 859 are displayed in Figure 1. Sublinear convergence can be observed from the log-log  
 860 plot in Figure 1a, where we use the diminishing stepsizes suggested in Theorem 2  
 861 and Theorem 3. In Figure 1b, we use geometrically diminishing stepsizes of the form  
 862  $\gamma_k = \beta^k \gamma_0$ . We fix  $\gamma_0 = 0.1$  and tune the best decay factor  $\beta$  for each algorithm. A  
 863 linear rate of convergence can be observed, which corroborates our theoretical results.

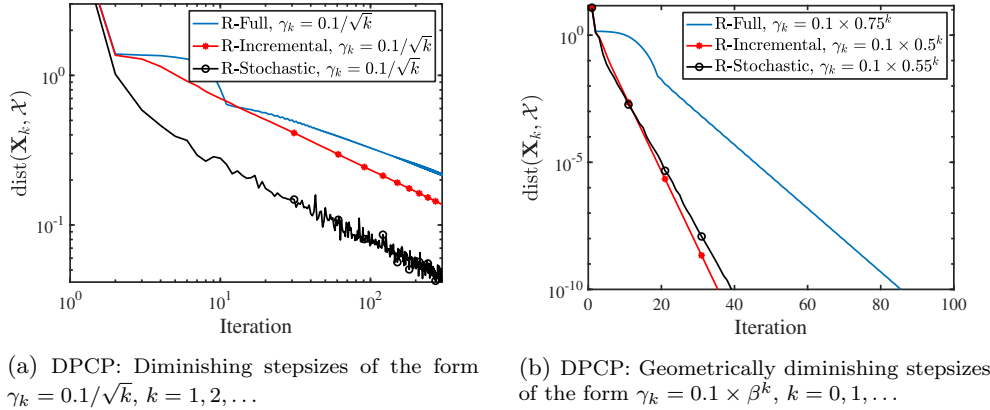


Fig. 1: Convergence performance of Riemannian subgradient-type methods for the DPCP formulation (1.3).

864 **7.2. Orthogonal dictionary learning (ODL).** We now turn to the orthog-  
 865 onal DL problem. Given  $\mathbf{Y} = \mathbf{A}\mathbf{S} \in \mathbb{R}^{n \times m}$ , where  $\mathbf{A} \in \text{St}(n, n)$  is an unknown  
 866 orthonormal dictionary and each column of  $\mathbf{S} \in \mathbb{R}^{n \times m}$  is sparse, we can try to re-  
 867 cover the columns of  $\mathbf{A}$  one at a time by considering the formulation (1.4), whose  
 868 objective function takes the form  $\mathbb{S}^{n-1} \ni \mathbf{x} \mapsto f(\mathbf{x}) = \frac{1}{m} \sum_{i=1}^m |\mathbf{y}_i^\top \mathbf{x}|$ , or to recover  
 869 the entire dictionary by considering the formulation (1.5), whose objective function  
 870 takes the form  $\text{St}(n, n) \ni \mathbf{X} \mapsto f(\mathbf{X}) = \frac{1}{m} \sum_{i=1}^m \|\mathbf{y}_i^\top \mathbf{X}\|_1$ .

871 **Sharpness.** The sharpness property of the formulation (1.4) has been studied  
 872 in [3], while that of (1.5) has been studied in [59] only in the *asymptotic* regime;  
 873 i.e., when the number of samples  $m$  tends to infinity. Although we do not yet know  
 874 how to establish the sharpness property of (1.5) in the *finite-sample* regime, the  
 875 following numerical results suggest that problem (1.5) likely possesses such a property,  
 876 as the Riemannian subgradient-type methods with geometrically diminishing stepsizes  
 877 exhibit linear convergence behavior, even with a random initialization. We leave the  
 878 study of the sharpness property of (1.5) in the finite-sample regime as a future work.

879 **Experiments.** For the orthogonal DL application, we generate synthetic data  
 880 in the same way as [3]. Specifically, we first generate the underlying orthogonal  
 881 dictionary  $\mathbf{A} \in \text{St}(n, n)$  with  $n = 30$  randomly and set the number of samples  $m$   
 882 to be  $m = 1643 \approx 10 \times n^{1.5}$ . We then generate a sparse coefficient matrix  $\mathbf{S} \in \mathbb{R}^{n \times m}$ ,  
 883 in which each entry follows the Bernoulli-Gaussian distribution with parameter 0.3  
 884 (sparsity)—i.e., each entry  $\mathbf{S}_{i,j}$  is drawn independently from the standard Gaussian  
 885 distribution with probability 0.3 and is set to zero otherwise. Lastly, we obtain the  
 886 observation  $\mathbf{Y} = \mathbf{A}\mathbf{S}$ . As before, we generate a standard Gaussian random vector and  
 887 use it to initialize all the algorithms. To evaluate the performance of the algorithms,  
 888 we define the error between  $\mathbf{X}$  and  $\mathbf{A}$  as  $\text{err}(\mathbf{X}, \mathbf{A}) = \sum_{i=1}^n |\max_{1 \leq j \leq n} |[\mathbf{x}_i^\top \mathbf{A}]_j| - 1|$ ,  
 889 where  $\mathbf{x}_i$  is the  $i$ -th column of  $\mathbf{X}$ . Clearly,  $\text{err}(\mathbf{X}, \mathbf{A}) = 0$  when  $\mathbf{X}$  and  $\mathbf{A}$  are  
 890 equal up to permutation and sign ambiguities. The numerical results are shown in  
 891 Figure 2. The log-log plot in Figure 2a shows the sublinear convergence of Riemannian  
 892 subgradient-type methods when the diminishing stepsizes suggested in Theorem 2 and



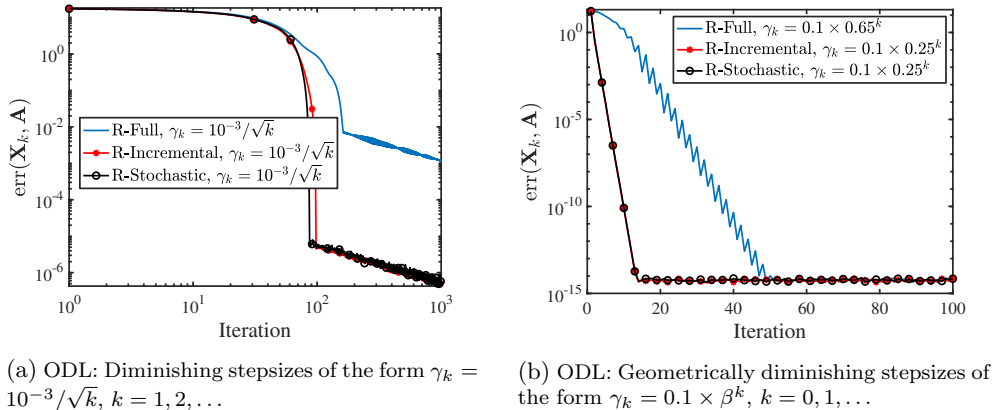


Fig. 2: Convergence performance of Riemannian subgradient-type methods for the orthogonal DL problem (1.5).

893 **Theorem 3** are used. **Figure 2b** shows the linear convergence of those methods when  
 894 geometrically diminishing stepsizes of the form  $\gamma_k = \beta^k \gamma_0$  are used. Here,  $\gamma_0 = 0.1$   
 895 and the best decay factor  $\beta$  is chosen for each algorithm.

896 **8. Conclusion.** In this work, we introduced a family of Riemannian subgradient-  
 897 type methods for minimizing weakly convex functions over the Stiefel manifold. We  
 898 proved, for the first time, iteration complexity and local convergence rate results for  
 899 these methods. Specifically, we showed that all these methods have a global sublinear  
 900 convergence rate, and that if the problem at hand further possesses the sharpness  
 901 property, then the Riemannian subgradient and incremental subgradient methods  
 902 with geometrically diminishing stepsizes and a proper initialization will converge lin-  
 903 early to the set of weak sharp minima of the problem. The key to establishing these  
 904 results is a new Riemannian subgradient inequality for restrictions of weakly convex  
 905 functions on the Stiefel manifold, which could be of independent interest. Our results  
 906 can be extended to cover weakly convex minimization over a class of compact embed-  
 907 ded submanifolds of the Euclidean space. Lastly, we showed that certain formulations  
 908 of the RSR and orthogonal DL problems possess the sharpness property and verified  
 909 the convergence performance of the Riemannian subgradient-type methods on these  
 910 problems via numerical simulations.

911 Our work has opened up several interesting directions for future investigation.  
 912 First, one can readily generalize our results to weakly convex minimization over a  
 913 Cartesian product of Stiefel manifolds, which has applications in  $\ell_1$ -PCA [33, 58]  
 914 and robust phase synchronization [57]. Next, since our results are specific to weakly  
 915 convex minimization over the Stiefel manifold, it would be interesting to see if they can  
 916 be extended to handle more general nonconvex nonsmooth functions over a broader  
 917 class of Riemannian manifolds. We believe that this should be possible based on the  
 918 analytic framework developed here. Finally, we suspect that the global convergence  
 919 rate  $\mathcal{O}(k^{-\frac{1}{4}})$  we established for the Riemannian subgradient-type methods is not tight.  
 920 This is because the Riemannian proximal point method for solving problem (1.1) has  
 921 a global convergence rate of  $\mathcal{O}(k^{-\frac{1}{2}})$  [10], and in smooth optimization the gradient

922 descent method has the same global convergence rate as the proximal point method.  
 923 Hence, it would be interesting to see if the global convergence rate established in this  
 924 paper can be improved.

925 **Acknowledgments.** We would like to thank Dr. Huikang Liu for fruitful dis-  
 926 cussions. We also thank the Associate Editor and two anonymous reviewers for their  
 927 detailed and helpful comments.

928

## REFERENCES

- 929 [1] P.-A. ABSIL AND S. HOSSEINI, *A collection of nonsmooth Riemannian optimization problems*,  
 930 in *Nonsmooth Optimization and Its Applications*, Springer, 2019, pp. 1–15.
- 931 [2] P.-A. ABSIL, R. MAHONY, AND R. SEPULCHRE, *Optimization Algorithms on Matrix Manifolds*,  
 932 Princeton University Press, 2009.
- 933 [3] Y. BAI, Q. JIANG, AND J. SUN, *Subgradient descent learns orthogonal dictionaries*, in *International  
 934 Conference on Learning Representations*, 2019.
- 935 [4] G. C. BENTO, O. P. FERREIRA, AND J. G. MELO, *Iteration-complexity of gradient, subgradient  
 936 and proximal point methods on Riemannian manifolds*, *Journal of Optimization Theory  
 937 and Applications*, 173 (2017), pp. 548–562.
- 938 [5] R. L. BISHOP AND B. O’NEILL, *Manifolds of negative curvature*, *Transactions of the American  
 939 Mathematical Society*, 145 (1969), pp. 1–49.
- 940 [6] S. BOUCHERON, G. LUGOSI, AND P. MASSART, *Concentration Inequalities: A Nonasymptotic  
 941 Theory of Independence*, Oxford University Press, 2013.
- 942 [7] N. BOUMAL, P.-A. ABSIL, AND C. CARTIS, *Global rates of convergence for nonconvex optimiza-  
 943 tion on manifolds*, *IMA Journal of Numerical Analysis*, 39 (2019), pp. 1–33.
- 944 [8] J. V. BURKE AND M. C. FERRIS, *Weak sharp minima in mathematical programming*, *SIAM  
 945 Journal on Control and Optimization*, 31 (1993), pp. 1340–1359.
- 946 [9] J. V. BURKE, A. S. LEWIS, AND M. L. OVERTON, *A robust gradient sampling algorithm for  
 947 nonsmooth, nonconvex optimization*, *SIAM Journal on Optimization*, 15 (2005), pp. 751–  
 948 779.
- 949 [10] S. CHEN, Z. DENG, S. MA, AND A. M.-C. SO, *Manifold proximal point algorithms for  
 950 dual principal component pursuit and orthogonal dictionary learning*, arXiv preprint  
 951 arXiv:2005.02356, (2020).
- 952 [11] S. CHEN, S. MA, A. M.-C. SO, AND T. ZHANG, *Proximal gradient method for nonsmooth  
 953 optimization over the Stiefel manifold*, *SIAM Journal on Optimization*, 30 (2020), pp. 210–  
 954 239.
- 955 [12] D. DAVIS AND D. DRUSVYATSKIY, *Stochastic model-based minimization of weakly convex func-  
 956 tions*, *SIAM Journal on Optimization*, 29 (2019), pp. 207–239.
- 957 [13] D. DAVIS, D. DRUSVYATSKIY, K. J. MACPHEE, AND C. PAQUETTE, *Subgradient methods for  
 958 sharp weakly convex functions*, *Journal of Optimization Theory and Applications*, 179  
 959 (2018), pp. 962–982.
- 960 [14] G. DE CARVALHO BENTO, J. X. DA CRUZ NETO, AND P. R. OLIVEIRA, *A new approach to  
 961 the proximal point method: Convergence on general Riemannian manifolds*, *Journal of  
 962 Optimization Theory and Applications*, 168 (2016), pp. 743–755.
- 963 [15] D. DRUSVYATSKIY, *The proximal point method revisited*, *SIAG/OPT Views and News*, 26  
 964 (2018), pp. 1–7.
- 965 [16] D. DRUSVYATSKIY AND C. PAQUETTE, *Efficiency of minimizing compositions of convex func-  
 966 tions and smooth maps*, *Mathematical Programming*, 178 (2019), pp. 503–558.
- 967 [17] A. EDELMAN, T. A. ARIAS, AND S. T. SMITH, *The geometry of algorithms with orthogonality  
 968 constraints*, *SIAM Journal on Matrix Analysis and Applications*, 20 (1998), pp. 303–353.
- 969 [18] O. FERREIRA, M. LOUZEIRO, AND L. PRUDENTE, *Iteration-complexity of the subgradient method  
 970 on Riemannian manifolds with lower bounded curvature*, *Optimization*, 68 (2019), pp. 713–  
 971 729.
- 972 [19] O. P. FERREIRA AND P. R. OLIVEIRA, *Subgradient algorithm on Riemannian manifolds*, *Journal  
 973 of Optimization Theory and Applications*, 97 (1998), pp. 93–104.
- 974 [20] O. P. FERREIRA AND P. R. OLIVEIRA, *Proximal point algorithm on Riemannian manifold*,  
 975 *Optimization*, 51 (2002), pp. 257–270.
- 976 [21] J.-L. GOFFIN, *On convergence rates of subgradient optimization methods*, *Mathematical Pro-  
 977 gramming*, 13 (1977), pp. 329–347.
- 978 [22] S. HOSSEINI, W. HUANG, AND R. YOUSEFPOUR, *Line search algorithms for locally Lipschitz  
 979 functions on Riemannian manifolds*, *SIAM Journal on Optimization*, 28 (2018), pp. 596–

- 980 619.
- 981 [23] S. HOSSEINI AND A. USCHMAJEV, *A Riemannian gradient sampling algorithm for nonsmooth*  
982 *optimization on manifolds*, SIAM Journal on Optimization, 27 (2017), pp. 173–189.
- 983 [24] J. HU, X. LIU, Z.-W. WEN, AND Y.-X. YUAN, *A brief introduction to manifold optimization*,  
984 Journal of the Operations Research Society of China, 8 (2020), pp. 199–248.
- 985 [25] W. HUANG AND K. WEI, *Riemannian proximal gradient methods*, arXiv preprint  
986 arXiv:1909.06065, (2019).
- 987 [26] B. JIANG, S. MA, A. M.-C. SO, AND S. ZHANG, *Vector transport-free SVRG with general re-*  
988 *traction for Riemannian optimization: Complexity analysis and practical implementation*,  
989 arXiv preprint arXiv:1705.09059, (2017).
- 990 [27] M. JOURNÉE, Y. NESTEROV, P. RICHTÁRIK, AND R. SEPULCHRE, *Generalized power method for*  
991 *sparse principal component analysis*, Journal of Machine Learning Research, 11 (2010),  
992 pp. 517–553.
- 993 [28] M. M. KARKHANEI AND N. MAHDAVI-AMIRI, *Nonconvex weak sharp minima on Riemannian*  
994 *manifolds*, Journal of Optimization Theory and Applications, 183 (2019), pp. 85–104.
- 995 [29] A. V. KNYAZEV AND M. E. ARGENTATI, *Majorization for changes in angles between subspaces,*  
996 *Ritz values, and graph Laplacian spectra*, SIAM Journal on Matrix Analysis and Applica-  
997 tions, 29 (2007), pp. 15–32.
- 998 [30] A. KOVNATSKY, K. GLASHOFF, AND M. M. BRONSTEIN, *MADMM: a generic algorithm for non-*  
999 *smooth optimization on manifolds*, in European Conference on Computer Vision, Springer,  
1000 2016, pp. 680–696.
- 1001 [31] R. LAI AND S. OSHER, *A splitting method for orthogonality constrained problems*, Journal of  
1002 Scientific Computing, 58 (2014), pp. 431–449.
- 1003 [32] G. LERMAN AND T. MAUNU, *Fast, robust and non-convex subspace recovery*, Information and  
1004 Inference: A Journal of the IMA, 7 (2018), pp. 277–336.
- 1005 [33] G. LERMAN AND T. MAUNU, *An overview of robust subspace recovery*, Proceedings of the IEEE,  
1006 106 (2018), pp. 1380–1410.
- 1007 [34] G. LERMAN, M. B. MCCOY, J. A. TROPP, AND T. ZHANG, *Robust computation of linear models*  
1008 *by convex relaxation*, Foundations of Computational Mathematics, 15 (2015), pp. 363–410.
- 1009 [35] C. LI, B. S. MORDUKHOVICH, J. WANG, AND J.-C. YAO, *Weak sharp minima on Riemannian*  
1010 *manifolds*, SIAM Journal on Optimization, 21 (2011), pp. 1523–1560.
- 1011 [36] X. LI, Z. ZHU, A. M.-C. SO, AND J. D. LEE, *Incremental methods for weakly convex optimiza-*  
1012 *tion*, arXiv preprint arXiv:1907.11687, (2019).
- 1013 [37] X. LI, Z. ZHU, A. M.-C. SO, AND R. VIDAL, *Nonconvex robust low-rank matrix recovery*, SIAM  
1014 Journal on Optimization, 30 (2020), pp. 660–686.
- 1015 [38] H. LIU, A. M.-C. SO, AND W. WU, *Quadratic optimization with orthogonality constraint:*  
1016 *Explicit Lojasiewicz exponent and linear convergence of retraction-based line-search and*  
1017 *stochastic variance-reduced gradient methods*, Mathematical Programming, 178 (2019),  
1018 pp. 215–262.
- 1019 [39] J. MAIRAL, F. BACH, AND J. PONCE, *Sparse modeling for image and vision processing*, Foun-  
1020 dations and Trends<sup>®</sup> in Computer Graphics and Vision, 8 (2014), pp. 85–283.
- 1021 [40] T. MAUNU, T. ZHANG, AND G. LERMAN, *A well-tempered landscape for non-convex robust*  
1022 *subspace recovery*, Journal of Machine Learning Research, 20 (2019), pp. 1–59.
- 1023 [41] A. MAURER, *A vector-contraction inequality for Rademacher complexities*, in Proceedings of  
1024 the 27th International Conference on Algorithmic Learning Theory (ALT 2016), R. Ortner,  
1025 H. U. Simon, and S. Zilles, eds., vol. 9925 of Lecture Notes in Artificial Intelligence, 2016,  
1026 pp. 3–17.
- 1027 [42] A. NEDIĆ AND D. BERTSEKAS, *Convergence rate of incremental subgradient algorithms*, in  
1028 Stochastic Optimization: Algorithms and Applications, S. Uryasev and P. M. Pardalos,  
1029 eds., vol. 54 of Applied Optimization, Springer Science+Business Media, Dordrecht, 2001,  
1030 pp. 223–264.
- 1031 [43] A. NEMIROVSKI, A. JUDITSKY, G. LAN, AND A. SHAPIRO, *Robust stochastic approximation*  
1032 *approach to stochastic programming*, SIAM Journal on Optimization, 19 (2009), pp. 1574–  
1033 1609.
- 1034 [44] YU. NESTEROV, *Introductory Lectures on Convex Optimization: A Basic Course*, Kluwer Aca-  
1035 demic Publishers, Boston, 2004.
- 1036 [45] Q. QU, J. SUN, AND J. WRIGHT, *Finding a sparse vector in a subspace: Linear sparsity using*  
1037 *alternating directions*, IEEE Transactions on Information Theory, 62 (2016), pp. 5855–  
1038 5880.
- 1039 [46] R. T. ROCKAFELLAR AND R. J.-B. WETS, *Variational Analysis*, vol. 317 of Grundlehren der  
1040 mathematischen Wissenschaften, Springer Science & Business Media, second ed., 2009.
- 1041 [47] R. RUBINSTEIN, A. M. BRUCKSTEIN, AND M. ELAD, *Dictionaries for sparse representation*

- 1042 *modeling*, Proceedings of the IEEE, 98 (2010), pp. 1045–1057.
- 1043 [48] P. H. SCHÖNEMANN, *A generalized solution of the orthogonal procrustes problem*, Psychome-  
1044 trika, 31 (1966), pp. 1–10.
- 1045 [49] N. Z. SHOR, *Minimization Methods for Non-Differentiable Functions*, vol. 3 of Springer Series  
1046 in Computational Mathematics, Springer-Verlag, Berlin Heidelberg, 1985.
- 1047 [50] D. A. SPIELMAN, H. WANG, AND J. WRIGHT, *Exact recovery of sparsely-used dictionaries*, in  
1048 Proceedings of the 25th Annual Conference on Learning Theory, 2012, pp. 37.1–37.18.
- 1049 [51] G. W. STEWART AND J. SUN, *Matrix Perturbation Theory*, Academic Press, Boston, 1990.
- 1050 [52] J. SUN, Q. QU, AND J. WRIGHT, *Complete dictionary recovery over the sphere I: Overview and*  
1051 *the geometric picture*, IEEE Transactions on Information Theory, 63 (2016), pp. 853–884.
- 1052 [53] J. SUN, Q. QU, AND J. WRIGHT, *Complete dictionary recovery over the sphere II: Recovery by*  
1053 *Riemannian trust-region method*, IEEE Transactions on Information Theory, 63 (2016),  
1054 pp. 885–914.
- 1055 [54] M. C. TSAKIRIS AND R. VIDAL, *Dual principal component pursuit*, Journal of Machine Learning  
1056 Research, 19 (2018), pp. 1–49.
- 1057 [55] J.-P. VIAL, *Strong and weak convexity of sets and functions*, Mathematics of Operations Re-  
1058 search, 8 (1983), pp. 231–259.
- 1059 [56] R. VIDAL, Y. MA, AND S. S. SASTRY, *Generalized Principal Component Analysis*, vol. 40 of  
1060 Interdisciplinary Applied Mathematics, Springer-Verlag, New York, 2016.
- 1061 [57] L. WANG AND A. SINGER, *Exact and stable recovery of rotations for robust synchronization*,  
1062 Information and Inference: A Journal of the IMA, 2 (2013), pp. 145–193.
- 1063 [58] P. WANG, H. LIU, AND A. M.-C. SO, *Globally convergent accelerated proximal alternating max-  
1064 imization method for  $L_1$ -principal component analysis*, in Proceedings of the 2019 IEEE  
1065 International Conference on Acoustics, Speech, and Signal Processing (ICASSP 2019),  
1066 2019, pp. 8147–8151.
- 1067 [59] Y. WANG, S. WU, AND B. YU, *Unique sharp local minimum in  $\ell_1$ -minimization complete  
1068 dictionary learning*, Journal of Machine Learning Research, 21 (2020), pp. 1–52.
- 1069 [60] Z. WEN AND W. YIN, *A feasible method for optimization with orthogonality constraints*, Math-  
1070 ematical Programming, 142 (2013), pp. 397–434.
- 1071 [61] J. WRIGHT, Y. MA, J. MAIRAL, G. SAPIRO, T. S. HUANG, AND S. YAN, *Sparse representation  
1072 for computer vision and pattern recognition*, Proceedings of the IEEE, 98 (2010), pp. 1031–  
1073 1044.
- 1074 [62] H. XU, C. CARAMANIS, AND S. MANNOR, *Outlier-robust PCA: The high-dimensional case*,  
1075 IEEE Transactions on Information Theory, 59 (2012), pp. 546–572.
- 1076 [63] W. H. YANG, L.-H. ZHANG, AND R. SONG, *Optimality conditions for the nonlinear programming  
1077 problems on Riemannian manifolds*, Pacific Journal of Optimization, 10 (2014), pp. 415–  
1078 434.
- 1079 [64] S.-T. YAU, *Non-existence of continuous convex functions on certain Riemannian manifolds*,  
1080 Mathematische Annalen, 207 (1974), pp. 269–270.
- 1081 [65] Y. ZHAI, Z. YANG, Z. LIAO, J. WRIGHT, AND Y. MA, *Complete dictionary learning via  $\ell^4$ -norm  
1082 maximization over the orthogonal group*, Journal of Machine Learning Research, 21 (2020),  
1083 pp. 1–68.
- 1084 [66] H. ZHANG AND S. SRA, *First-order methods for geodesically convex optimization*, in Proceedings  
1085 of the 29th Annual Conference on Learning Theory, 2016, pp. 1617–1638.
- 1086 [67] T. ZHANG, *Robust subspace recovery by Tyler’s  $M$ -estimator*, Information and Inference: A  
1087 Journal of the IMA, 5 (2016), pp. 1–21.
- 1088 [68] T. ZHANG AND G. LERMAN, *A novel  $M$ -estimator for robust PCA*, Journal of Machine Learning  
1089 Research, 15 (2014), pp. 749–808.
- 1090 [69] Z. ZHU, T. DING, M. TSAKIRIS, D. ROBINSON, AND R. VIDAL, *A linearly convergent method  
1091 for non-smooth non-convex optimization on Grassmannian with applications to robust  
1092 subspace and dictionary learning*, in Advances in Neural Information Processing Systems,  
1093 2019, pp. 9437–9447.
- 1094 [70] Z. ZHU, Y. WANG, D. ROBINSON, D. NAIMAN, R. VIDAL, AND M. TSAKIRIS, *Dual principal  
1095 component pursuit: Improved analysis and efficient algorithms*, in Advances in Neural  
1096 Information Processing Systems, 2018, pp. 2171–2181.

## 1097 Appendix A. Proof of Lemma 2.

1098 The proof follows the framework in [34, Section 8.1.1] with nontrivial modifica-  
1099 tions in order to handle our matrix-based definitions of  $c_{\mathbf{Y},\min}$  and  $c_{\mathbf{O},\max}$ .

1100 *Part I.* We first derive an upper bound on  $c_{\mathbf{O},\max}$ . Recall that under the Haystack  
1101 model, the outliers  $\mathbf{o}_1, \dots, \mathbf{o}_{m_2} \in \mathbb{R}^n$  are i.i.d. according to the Gaussian distribution

1102  $\mathcal{N}(\mathbf{0}, \frac{1}{n}\mathbf{I}_n)$ . Let  $\mathbf{o} \sim \mathcal{N}(\mathbf{0}, \frac{1}{n}\mathbf{I}_n)$  denote an i.i.d. copy of  $\mathbf{o}_i$ . Then, we have

$$1103 \quad (\text{A.1}) \quad \begin{aligned} & \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \|\mathbf{o}_i^\top \mathbf{B}\|_2 \\ & \leq \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} (\|\mathbf{o}_i^\top \mathbf{B}\|_2 - \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2]) + \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2]. \end{aligned}$$

1104 Using Jensen's inequality, we bound the second term as follows:

$$1105 \quad (\text{A.2}) \quad \begin{aligned} & \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2] \leq \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \sqrt{\mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2^2]} \\ & = \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \sqrt{\sum_{j=1}^r \sum_{k=1}^n \mathbb{E} [o_k^2 b_{kj}^2]} = \frac{m_2}{\sqrt{n}}. \end{aligned}$$

1106 To estimate the first term in (A.1), let  $\{\epsilon_i : i = 1, \dots, m_2\}$  be independent Rademacher  
1107 random variables (i.e.,  $\Pr[\epsilon_i = +1] = \Pr[\epsilon_i = -1] = 1/2$  for  $i = 1, \dots, m_2$ ) that are  
1108 independent of  $\{\mathbf{o}_i : i = 1, \dots, m_2\}$ . By a standard symmetrization argument (see,  
1109 e.g., [6, Lemma 11.4]), we have

$$1110 \quad (\text{A.3}) \quad \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} (\|\mathbf{o}_i^\top \mathbf{B}\|_2 - \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2]) \right] \leq 2 \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \epsilon_i \|\mathbf{o}_i^\top \mathbf{B}\|_2 \right].$$

1111 Furthermore, let  $\{\epsilon_{ij} : i = 1, \dots, m_2; j = 1, \dots, r\}$  be independent Rademacher  
1112 random variables that are independent of  $\{\mathbf{o}_i : i = 1, \dots, m_2\}$  and  $\mathbf{V} \in \mathbb{R}^{n \times r}$  be  
1113 the matrix whose  $j$ -th column ( $j = 1, \dots, r$ ) is  $\sum_{i=1}^{m_2} \epsilon_{ij} \mathbf{o}_i$ . Then, by the vector  
1114 contraction inequality in [41, Corollary 1] and Jensen's inequality, we have

$$1115 \quad \begin{aligned} & \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \epsilon_i \|\mathbf{o}_i^\top \mathbf{B}\|_2 \right] \leq \sqrt{2} \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} \sum_{j=1}^r \epsilon_{ij} \mathbf{o}_i^\top \mathbf{b}_j \right] \\ 1116 & = \sqrt{2} \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{j=1}^r \left( \sum_{i=1}^{m_2} \epsilon_{ij} \mathbf{o}_i \right)^\top \mathbf{b}_j \right] = \sqrt{2} \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \langle \mathbf{V}, \mathbf{B} \rangle \right] = \sqrt{2} \mathbb{E} [\|\mathbf{V}\|_F] \\ 1117 & \leq \sqrt{2} \sqrt{\sum_{j=1}^r \mathbb{E} \left[ \left\| \sum_{i=1}^{m_2} \epsilon_{ij} \mathbf{o}_i \right\|_2^2 \right]} \leq \sqrt{2m_2 r}. \end{aligned}$$

1118  
1119 This, together with (A.3), yields

$$1120 \quad (\text{A.4}) \quad \mathbb{E} \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} (\|\mathbf{o}_i^\top \mathbf{B}\|_2 - \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2]) \right] \leq 2\sqrt{2m_2 r}.$$

1121 Now, observe that the function

$$1122 \quad (\mathbf{o}_1, \dots, \mathbf{o}_{m_2}) \mapsto h(\mathbf{o}_1, \dots, \mathbf{o}_{m_2}) := \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} (\|\mathbf{o}_i^\top \mathbf{B}\|_2 - \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2])$$

1123 is Lipschitz continuous with constant at most  $\sqrt{m_2}$ . Hence, using the Gaussian con-  
 1124 centration inequality for Lipschitz functions [6, Theorem 5.6] and (A.4), we get  
 (A.5)

$$1125 \quad \Pr \left[ \sup_{\|\mathbf{B}\|_F=1} \sum_{i=1}^{m_2} (\|\mathbf{o}_i^\top \mathbf{B}\|_2 - \mathbb{E} [\|\mathbf{o}^\top \mathbf{B}\|_2]) \leq 2\sqrt{2m_2 r} + t \right] \geq 1 - 2 \exp \left( -\frac{nt^2}{2m_2} \right).$$

1126 Upon substituting (A.5) and (A.2) into (A.1) and letting  $t = c_2\sqrt{m_2}$ , the desired  
 1127 result follows.

1128 *Part II.* We now derive a lower bound on  $c_{\mathbf{Y}, \min}$ . Again, recall that under the  
 1129 Haystack model, the inliers  $\mathbf{y}_1, \dots, \mathbf{y}_{m_1} \in \mathbb{R}^n$  are i.i.d. according to the Gaussian  
 1130 distribution  $\mathcal{N}(\mathbf{0}, \frac{1}{d}\mathcal{P}_S)$ . Thus, for  $i = 1, \dots, m_1$ , we have  $\mathbb{E}[\mathbf{y}_i \mathbf{y}_i^\top] = \frac{1}{d}\mathcal{P}_S = \frac{1}{d}\mathbf{S}\mathbf{S}^\top$   
 1131 for some orthonormal basis  $\mathbf{S} \in \text{St}(n, d)$  of  $\mathcal{S}$  and  $\mathbf{y}_i = \mathbf{S}\tilde{\mathbf{y}}_i$  for some  $\tilde{\mathbf{y}}_i \in \mathbb{R}^d$ . Now,  
 1132 let  $\mathbf{D} \in \mathbb{R}^{n \times \ell}$  be such that  $\|\mathbf{D}\|_F = 1$  and  $\text{col}(\mathbf{D}) \subseteq \mathcal{S}$ ; see (7.1). Then, there exists a  
 1133  $\tilde{\mathbf{D}} \in \mathbb{R}^{d \times \ell}$  such that  $\mathbf{D} = \mathbf{S}\tilde{\mathbf{D}}$  and  $\|\tilde{\mathbf{D}}\|_F = 1$ . In particular, we have  $\mathbf{y}_i^\top \mathbf{D} = \tilde{\mathbf{y}}_i^\top \tilde{\mathbf{D}}$ ,  
 1134 and by the rotational invariance of the Gaussian distribution, the vector  $\tilde{\mathbf{y}}_i$  follows the  
 1135 Gaussian distribution  $\mathcal{N}(\mathbf{0}, \frac{1}{d}\mathbf{I}_d)$  in  $\mathbb{R}^d$ . Consequently, we may assume without loss  
 1136 of generality that  $\mathbf{y}_1, \dots, \mathbf{y}_{m_1} \in \mathbb{R}^d$  are i.i.d. according to the Gaussian distribution  
 1137  $\mathcal{N}(\mathbf{0}, \frac{1}{d}\mathbf{I}_d)$  and  $\mathbf{D} \in \mathbb{R}^{d \times \ell}$  satisfies  $\|\mathbf{D}\|_F = 1$ . The rest of the proof will be similar to  
 1138 that of Part I.

1139 Let  $\mathbf{y} \sim \mathcal{N}(\mathbf{0}, \frac{1}{d}\mathbf{I}_d)$  denote an i.i.d. copy of  $\mathbf{y}_i$ . Then, we have

$$1140 \quad (\text{A.6}) \quad \inf_{\|\mathbf{D}\|_F=1} \sum_{i=1}^{m_1} \|\mathbf{y}_i^\top \mathbf{D}\|_2$$

$$\geq \inf_{\|\mathbf{D}\|_F=1} \sum_{i=1}^{m_1} (\|\mathbf{y}_i^\top \mathbf{D}\|_2 - \mathbb{E} [\|\mathbf{y}^\top \mathbf{D}\|_2]) + \inf_{\|\mathbf{D}\|_F=1} \sum_{i=1}^{m_1} \mathbb{E} [\|\mathbf{y}^\top \mathbf{D}\|_2].$$

1141 The first term can be written as  $-\sup_{\|\mathbf{D}\|_F=1} \sum_{i=1}^{m_1} (\mathbb{E} [\|\mathbf{y}^\top \mathbf{D}\|_2] - \|\mathbf{y}_i^\top \mathbf{D}\|_2)$ . By  
 1142 following the same arguments as in Part I, we obtain  
 (A.7)

$$1143 \quad \Pr \left[ \sup_{\|\mathbf{D}\|_F=1} \sum_{i=1}^{m_1} (\mathbb{E} [\|\mathbf{y}^\top \mathbf{D}\|_2] - \|\mathbf{y}_i^\top \mathbf{D}\|_2) \leq 2\sqrt{2m_1 \ell} + t \right] \geq 1 - 2 \exp \left( -\frac{dt^2}{2m_1} \right).$$

1144 It remains to estimate the second term in (A.6). Let  $\mathbf{d}_i$  be the  $i$ -th column of  $\mathbf{D}$ ,  
 1145 where  $i = 1, \dots, \ell$ . By the Cauchy-Schwarz inequality and the fact that  $\|\mathbf{D}\|_F = 1$ ,  
 1146 we have

$$1147 \quad \underbrace{(\|\mathbf{d}_1\|_2^2 + \dots + \|\mathbf{d}_\ell\|_2^2)}_{=\|\mathbf{D}\|_F^2=1} \underbrace{[(\mathbf{y}^\top \mathbf{d}_1)^2 + \dots + (\mathbf{y}^\top \mathbf{d}_\ell)^2]}_{=\|\mathbf{y}^\top \mathbf{D}\|_2^2}$$

$$\geq [\|\mathbf{d}_1\|_2 |\mathbf{y}^\top \mathbf{d}_1| + \dots + \|\mathbf{d}_\ell\|_2 |\mathbf{y}^\top \mathbf{d}_\ell|]^2.$$

1148 Since  $\mathbf{y}^\top \mathbf{d}_i \sim \mathcal{N}(0, \frac{1}{d}\|\mathbf{d}_i\|_2^2)$ , we obtain

$$1149 \quad \mathbb{E} [\|\mathbf{y}^\top \mathbf{D}\|_2] \geq \mathbb{E} [\|\mathbf{d}_1\|_2 |\mathbf{y}^\top \mathbf{d}_1| + \dots + \|\mathbf{d}_\ell\|_2 |\mathbf{y}^\top \mathbf{d}_\ell|] = \sqrt{\frac{2}{d\pi}}.$$

1150 Note that the above inequality holds as equality when  $\mathbf{d}_1 = \dots = \mathbf{d}_\ell$ . This implies  
 1151 that

$$1152 \quad (\text{A.8}) \quad \inf_{\|\mathbf{D}\|_F=1} \sum_{i=1}^{m_1} \mathbb{E} [\|\mathbf{y}^\top \mathbf{D}\|_2] = m_1 \sqrt{\frac{2}{d\pi}}.$$

1153 By substituting (A.7) with  $t = c_1\sqrt{m_1}$  and (A.8) into (A.6), we complete the proof.