

Supplementary Material for:
“A Linearly Convergent Method for Non-smooth Non-convex
Optimization on the Grassmannian with Applications to
Robust Subspace and Dictionary Learning”

May 23, 2019

Contents

1	Basic Notations and Definitions	2
1.1	$(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian Regularity Condition (RRC)	3
2	Proof of Proposition 1	3
3	Proof of Theorem 1	5
4	Proof of Theorem 2	6
5	Proof of Corollary 2	9
6	Initialization	10
7	Random Spherical Model	10
8	Comparison with the Regularity Condition for Smooth Function	12
9	Additional Experiments	13
9.1	DPCP for Robust Subspace Learning	13
9.2	Orthogonal Dictionary Learning	13
10	Proof of Lemma 5	14
10.1	Preliminaries	14
10.2	Proof of Lemma 5	16

1 Basic Notations and Definitions

We consider the following problem [1]

$$\underset{\mathbf{B} \in \mathbb{O}(c, D)}{\text{minimize}} f(\mathbf{B}) \quad (1)$$

where $f : \mathbb{R}^{D \times c} \rightarrow \mathbb{R}$ is lower semi-continuous, possibly non-convex and non-smooth, and homogeneous that $f(\mathbf{B}) = f(\mathbf{B}\mathbf{Q})$ for any $\mathbf{Q} \in \mathbb{O}(c)$. Note that the global minimum of (1) is not unique as if \mathbf{B}^* is a global minimum, then any point in $[\mathbf{B}^*]$ is also a global minimum.

For any $\mathbf{A}, \mathbf{B} \in \mathbb{O}(c, D)$, the principal angles between $\text{Span}(\mathbf{A})$ and $\text{Span}(\mathbf{B})$ are defined as [2] $\phi_i(\mathbf{A}, \mathbf{B}) = \arccos(\sigma_i(\mathbf{A}^\top \mathbf{B}))$, $i = 1, \dots, c$, where $\sigma_i(\cdot)$ denotes the i -th singular value. We then define the distance between \mathbf{A} and \mathbf{B} as

$$\text{dist}(\mathbf{A}, \mathbf{B}) := \sqrt{2 \sum_{i=1}^c (1 - \cos(\phi_i(\mathbf{A}, \mathbf{B})))} = \inf_{\mathbf{Q} \in \mathbb{O}(c)} \|\mathbf{B} - \mathbf{A}\mathbf{Q}\|_F, \quad (2)$$

where the last term is also known as the *orthogonal Procrustes problem* and the last equality follows from the result [3] that the optimal rotation matrix \mathbf{Q} minimizing $\|\mathbf{B} - \mathbf{A}\mathbf{Q}\|_F$ is equal to $\mathbf{Q} = \mathbf{U}\mathbf{V}^\top$, where $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$ is the SVD of $\mathbf{A}^\top \mathbf{B}$. We also define the projection of \mathbf{B} onto $[\mathbf{A}]$ as

$$\mathcal{P}_{\mathbf{A}}(\mathbf{B}) = \mathbf{A}\mathbf{Q}^*, \quad \text{where } \mathbf{Q}^* = \arg \min_{\mathbf{Q} \in \mathbb{O}(c)} \|\mathbf{B} - \mathbf{A}\mathbf{Q}\|_F. \quad (3)$$

Since f can be non-smooth and non-convex, we utilize the Fréchet subdifferential, which generalizes the gradient for smooth functions and the subdifferential in convex analysis. The subdifferential of a lower semi-continuous function f at \mathbf{B} is defined as

$$\partial f(\mathbf{B}) := \left\{ \mathbf{D} \in \mathbb{R}^{D \times c} : \liminf_{\mathbf{A} \rightarrow \mathbf{B}} \frac{f(\mathbf{A}) - f(\mathbf{B}) - \langle \mathbf{D}, \mathbf{A} - \mathbf{B} \rangle}{\|\mathbf{B} - \mathbf{A}\|_F} \geq 0 \right\}.$$

Roughly speaking, for each subgradient of f at \mathbf{B} (i.e., for each $\mathbf{D} \in \partial f(\mathbf{B})$), the graph of $\mathbf{A} \mapsto f(\mathbf{A}) + \langle \mathbf{D}, \mathbf{A} - \mathbf{B} \rangle$ constructs a local supporting hyperplane to the graph of f at \mathbf{B} . If f is a convex function, then a local supporting hyperplane turns out to be a global one, and $\partial f(\mathbf{B})$ reduces to $\{\mathbf{D} \in \mathbb{R}^{D \times c} : f(\mathbf{A}) - f(\mathbf{B}) \geq \langle \mathbf{D}, \mathbf{A} - \mathbf{B} \rangle \text{ for all } \mathbf{A} \in \mathbb{R}^{D \times c}\}$ [4]. If f is a smooth function, then the subdifferential $\partial f(\mathbf{B})$ is simply $\{\nabla f(\mathbf{B})\}$.

Since we consider problems on the Grassmannian, we use tools from Riemannian geometry to state optimality conditions. From [1], the tangent space of the Grassmannian at $[\mathbf{B}]$ is defined as $\{\mathbf{W} \in \mathbb{R}^{D \times c} : \mathbf{W}^\top \mathbf{B} = \mathbf{0}\}$, and the orthogonal projector onto the tangent space is $\mathbf{I} - \mathbf{B}\mathbf{B}^\top$, which is well-defined and does not depend on the class representative as $\mathbf{A}\mathbf{A}^\top = \mathbf{B}\mathbf{B}^\top$ for any $\mathbf{A} \in [\mathbf{B}]$. If $\mathbf{D} \in \partial f(\mathbf{B})$, then we call $(\mathbf{I} - \mathbf{B}\mathbf{B}^\top)\mathbf{D}$ a Riemannian subgradient of f at \mathbf{B} ; we define the collection of all such Riemannian subgradients of f at \mathbf{B} as

$$\tilde{\partial} f(\mathbf{B}) := \left\{ (\mathbf{I} - \mathbf{B}\mathbf{B}^\top)\mathbf{D} : \mathbf{D} \in \partial f(\mathbf{B}) \right\}. \quad (4)$$

We say that \mathbf{B} is a critical point of (1) if $\mathbf{0} \in \tilde{\partial} f(\mathbf{B})$, which is a necessary condition for being a minimizer to (1) as shown in [5].

Armed with the Riemannian subgradient, a simple yet popular method for solving (1) is the projected Riemannian subgradient method, which is stated in Algorithm 1.

Algorithm 1 Projected Riemannian Subgradient Method

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

- 1: **for** $k = 0, 1, \dots$ **do**
- 2: obtain $\mathcal{G}(\mathbf{B}_k) \in \tilde{\partial}f(\mathbf{B}_k)$ satisfying (6) with $\mathbf{B} = \mathbf{B}_k$;
- 3: compute a step size μ_k according to a certain rule;
- 4: update the iterate:

$$\widehat{\mathbf{B}}_{k+1} \leftarrow \mathbf{B}_k - \mu_k \mathcal{G}(\mathbf{B}_k), \quad \mathbf{B}_{k+1} \leftarrow \text{orth}(\widehat{\mathbf{B}}_{k+1}); \quad (5)$$

5: **end for**

1.1 $(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian Regularity Condition (RRC)

Definition 1. We say that $f : \mathbb{R}^{D \times c} \rightarrow \mathbb{R}$ satisfies the $(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian regularity condition (RRC) for parameters $\{\alpha, \epsilon\} > 0$ and $\mathbf{B}^* \in \mathbb{O}(c, D)$, if for every $\mathbf{B} \in \mathbb{O}(c, D)$ satisfying $\text{dist}(\mathbf{B}, \mathbf{B}^*) \leq \epsilon$, there exists a Riemannian subgradient $\mathcal{G}(\mathbf{B}) \in \tilde{\partial}f(\mathbf{B})$ such that

$$\langle \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}) - \mathbf{B}, -\mathcal{G}(\mathbf{B}) \rangle \geq \alpha \text{dist}(\mathbf{B}, \mathbf{B}^*). \quad (6)$$

Let

$$\xi := \sup \{ \|\mathcal{G}(\mathbf{B})\|_F : \text{dist}(\mathbf{B}, \mathbf{B}^*) \leq \epsilon \} \quad (7)$$

denote an upper bound on the size of the Riemannian subgradients in a neighborhood of \mathbf{B}^* . Assume $\xi < \infty$. To compare α and ξ , we plug the Cauchy-Schwarz inequality $\langle \mathbf{B} - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}), \mathcal{G}(\mathbf{B}) \rangle \leq \|\mathcal{G}(\mathbf{B})\|_F \text{dist}(\mathbf{B}, \mathbf{B}^*)$ into (6), giving

$$\|\mathcal{G}(\mathbf{B})\|_F \geq \alpha, \quad \forall \mathbf{B} \notin [\mathbf{B}^*], \text{dist}(\mathbf{B}, \mathbf{B}^*) \leq \epsilon, \quad (8)$$

which implies that

$$\xi \geq \alpha. \quad (9)$$

2 Proof of Proposition 1

We first repeat Proposition 1.

Proposition 1. Suppose that for some $(\alpha, \epsilon, \mathbf{B}^*)$ the function f satisfies the $(\alpha, \epsilon, \mathbf{B}^*)$ -RRC in Definition 1. Let $\{\mathbf{B}_k\}$ be generated by Algorithm 1 with step size $\mu_k \equiv \mu \leq \alpha\epsilon/\xi^2$ and initial iterate \mathbf{B}_0 satisfying $\text{dist}_0 \leq \epsilon$. Then, for all $k \geq 0$, it holds that

$$\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \leq \max \left\{ \text{dist}(\mathbf{B}_0, \mathbf{B}^*) - \mu\alpha k/2, \mu\xi^2/\alpha \right\}. \quad (10)$$

Proof of Proposition 1. We first prove that $\widehat{\mathbf{B}}_{k+1}$ always has full column rank since $\mathcal{G}(\mathbf{B}_k)$ is orthogonal to \mathbf{B}_k . Let $\widehat{\mathbf{B}}_{k+1} = \mathbf{P}\mathbf{\Omega}\mathbf{Q}^\top$ be a reduced SVD of $\widehat{\mathbf{B}}_{k+1}$, where $\mathbf{\Omega}$ is an $c \times c$ diagonal matrix with singular values w_1, \dots, w_c along the diagonals. Since $\widehat{\mathbf{B}}_{k+1} = \mathbf{B}_k - \mu_k \mathcal{G}(\mathbf{B}_k)$, we have

$$\widehat{\mathbf{B}}_{k+1}^\top \widehat{\mathbf{B}}_{k+1} = \mathbf{I} + \mu_k^2 (\mathcal{G}(\mathbf{B}_k))^\top \mathcal{G}(\mathbf{B}_k),$$

where the equality follows because $\mathbf{B}_k \in \mathbb{O}(c, D)$ is orthogonal to $\mathcal{G}(\mathbf{B}_k)$. Thus, the eigenvalues of $\widehat{\mathbf{B}}_{k+1}^\top \widehat{\mathbf{B}}_{k+1}$ is always greater than or equal to 1, which implies that $w_1, \dots, w_c \geq 1$. Therefore, all singular values of $\widehat{\mathbf{B}}_{k+1}$ are non-vanishing, which means $\widehat{\mathbf{B}}_{k+1}$ has full column rank. Additionally, for any $\mathbf{U} \in \mathbb{O}(c, D)$, it follows that

$$\begin{aligned} & \|\widehat{\mathbf{B}}_{k+1} - \mathbf{U}\|_F^2 - \|\mathbf{B}_{k+1} - \mathbf{U}\|_F^2 \\ &= \|\mathbf{P}\boldsymbol{\Omega}\mathbf{Q}^\top\|_F^2 - \|\mathbf{P}\mathbf{Q}^\top\|_F^2 - 2\text{trace}((\boldsymbol{\Omega} - \mathbf{I})\mathbf{P}^\top\mathbf{U}\mathbf{Q}) \\ &\geq \sum_{i=1}^c \omega_i^2 - 1 - 2(\omega_i - 1) = \sum_{i=1}^c (\omega_i - 1)^2 \geq 0, \end{aligned} \tag{11}$$

where we have chosen \mathbf{B}_{k+1} to be $\mathbf{P}\mathbf{Q}^\top$, and the last line directly follows Von Neumann's inequality, i.e., $\text{trace}(\mathbf{F}^\top\mathbf{G}) \leq \sum_i \sigma_i(\mathbf{F})\sigma_i(\mathbf{G})$ where $\sigma_1(\cdot) \geq \sigma_2(\cdot) \geq \dots \geq 0$.

We now prove (10) by induction. It is clear that (10) holds when $k = 0$. Now assume that (10) holds at the k -th iteration, which implies that $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \leq \epsilon$. Then,

$$\begin{aligned} \text{dist}^2(\mathbf{B}_{k+1}, \mathbf{B}^*) &\leq \|\mathbf{B}_{k+1} - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_k)\|_F^2 \\ &\leq \|\widehat{\mathbf{B}}_{k+1} - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_k)\|_F^2 \\ &= \|\mathbf{B}_k - \mu\mathcal{G}(\mathbf{B}_k) - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_k)\|_F^2 \\ &= \|\mathbf{B}_k - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_k)\|_F^2 - 2\mu\langle \mathbf{B}_k - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_k), \mathcal{G}(\mathbf{B}_k) \rangle + \mu^2\|\mathcal{G}(\mathbf{B}_k)\|_F^2 \\ &\leq \text{dist}^2(\mathbf{B}_k, \mathbf{B}^*) - 2\alpha\mu \text{dist}(\mathbf{B}_k, \mathbf{B}^*) + \mu^2\xi^2, \end{aligned} \tag{12}$$

where the second line utilizes (11), and the last line utilizes the Riemannian regularity condition (6).

It is clear from (12) that $\text{dist}^2(\mathbf{B}_{k+1}, \mathbf{B}^*) \leq \text{dist}^2(\mathbf{B}_k, \mathbf{B}^*)$ if $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \geq \frac{\mu\xi^2}{2\alpha}$. In particular, when $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \geq \frac{\mu\xi^2}{\alpha}$, we have

$$\begin{aligned} \text{dist}^2(\mathbf{B}_{k+1}, \mathbf{B}^*) &\leq \text{dist}^2(\mathbf{B}_k, \mathbf{B}^*) - \alpha\mu \text{dist}(\mathbf{B}_k, \mathbf{B}^*) + \mu^2\xi^2 - \alpha\mu \text{dist}(\mathbf{B}_k, \mathbf{B}^*) \\ &\leq \left(\text{dist}(\mathbf{B}_k, \mathbf{B}^*) - \frac{\mu\alpha}{2} \right)^2, \end{aligned}$$

which implies that

$$\text{dist}(\mathbf{B}_{k+1}, \mathbf{B}^*) \leq \text{dist}(\mathbf{B}_k, \mathbf{B}^*) - \frac{\mu\alpha}{2}$$

since $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \geq \frac{\mu\xi^2}{\alpha} \geq \mu\alpha$.

On the other hand, when $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \leq \frac{\mu\xi^2}{\alpha}$, it also follows from (12) that

$$\begin{aligned} \text{dist}^2(\mathbf{B}_{k+1}, \mathbf{B}^*) &\leq \max \left\{ \left(\frac{\mu\xi^2}{\alpha} \right)^2 - 2\mu\alpha \frac{\mu\xi^2}{\alpha} + \mu^2\xi^2, \mu^2\xi^2 \right\} \\ &= \max \left\{ \left(\frac{\mu\xi^2}{\alpha} \right)^2 - \mu^2\xi^2, \mu^2\xi^2 \right\} \\ &\leq \max \left\{ \left(\frac{\mu\xi^2}{\alpha} \right)^2 - \mu^2\xi^2, \mu^2\xi^2 \frac{\xi^2}{\alpha^2} \right\} \\ &= \left(\frac{\mu\xi^2}{\alpha} \right)^2, \end{aligned}$$

where the first inequality follows that $h(t) := t^2 - 2\alpha\mu t$ is increasing in $[t', \infty]$ for any t' such that $h(t') \geq 0$, and the second inequality utilizes (9). Thus by induction, (10) holds for all $k \geq 0$. \square

3 Proof of Theorem 1

We first repeat Theorem 1 (which is a lightly stronger result than the one presented in the paper).

Theorem 1. *Suppose that f satisfies the $(\alpha, \epsilon, \mathbf{B}^*)$ -RRC in Definition 1. Let $\{\mathbf{B}_k\}$ be the sequence generated by Algorithm 1 with step size*

$$\mu_k = \mu_0 \beta^k \quad (13)$$

and initialization \mathbf{B}_0 satisfying $\text{dist}(\mathbf{B}_0, \mathbf{B}^*) \leq \epsilon$. Assume

$$\mu_0 \leq \frac{\alpha \text{dist}_0}{2\xi^2} \quad \text{and} \quad \sqrt{1 - 2\frac{\alpha\mu_0}{\text{dist}_0} + \frac{\mu_0^2 \xi^2}{\text{dist}_0^2}} =: \underline{\beta} \leq \beta < 1, \quad (14)$$

for some $\text{dist}(\mathbf{B}_0, \mathbf{B}^*) \leq \text{dist}_0 \leq \epsilon$, where ξ is defined in (7). Then, the sequence $\{\mathbf{B}_k\}$ satisfies

$$\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \leq \text{dist}_0 \beta^k \quad \text{for all } k \geq 0. \quad (15)$$

Proof of Theorem 1. We first show that $\underline{\beta}$ in (14) is well-defined and satisfies $0 < \underline{\beta} < 1$. To see this, on one hand, $\mu_0 \leq \alpha \text{dist}_0 / 2\xi^2$ and (9) together imply $1 - 2\alpha\mu_0 / \text{dist}_0 \geq 0$. On the other hand, $-2\alpha\mu_0 / \text{dist}_0 + \mu_0^2 \xi^2 / \text{dist}_0^2 < 0$ is a decreasing function of μ_0 when $\mu_0 \in (0, \alpha \text{dist}_0 / 2\xi^2]$. In particular, when $\mu_0 = \alpha \text{dist}_0 / 2\xi^2$, we have $\underline{\beta} = \sqrt{1 - 3\alpha^2 / 4\xi^2}$, giving the fastest decaying rate by setting $\beta = \underline{\beta}$.

We now prove (15) by induction. It is clear that (15) holds when $k = 0$. Now assume that (15) holds at the k -th iteration, which implies that $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \leq \text{dist}_0 \beta^k$. Since \mathbf{B}_k satisfies the Riemannian regularity condition (6), we know that (12) holds:

$$\begin{aligned} \text{dist}^2(\mathbf{B}_{k+1}, \mathbf{B}^*) &\leq \text{dist}^2(\mathbf{B}_k, \mathbf{B}^*) - 2\alpha\mu_k \text{dist}(\mathbf{B}_k, \mathbf{B}^*) + \mu_k^2 \xi^2 \\ &= (\text{dist}(\mathbf{B}_k, \mathbf{B}^*) - \alpha\mu_k)^2 + \mu_k^2 (\xi^2 - \alpha^2). \end{aligned} \quad (16)$$

From $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) \leq \text{dist}_0 \beta^k$ and

$$\text{dist}_0 \beta^k \geq 2\frac{\mu_0 \xi^2}{\alpha} \beta^k \geq 2\alpha\mu_0 \beta^k = 2\alpha\mu_k \geq \alpha\mu_k,$$

where the first inequality is due to the assumption (14) and the second inequality follows $\xi \geq \alpha$ in (9). Therefore, (16) achieves its maximum at $\text{dist}(\mathbf{B}_k, \mathbf{B}^*) = \text{dist}_0 \beta^k$. Plugging this observation into (16) gives

$$\begin{aligned} \text{dist}^2(\mathbf{B}_{k+1}, \mathbf{B}^*) &\leq \text{dist}_0^2 \beta^{2k} - 2\alpha\mu_k \text{dist}_0 \beta^k + \mu_k^2 \xi^2 \\ &= \text{dist}_0^2 \beta^{2k} - 2\alpha\mu_0 \text{dist}_0 \beta^{2k} + \mu_0^2 \beta^{2k} \xi^2 \\ &= \text{dist}_0^2 \beta^{2k} \left(1 - 2\frac{\alpha\mu_0}{\text{dist}_0} + \frac{\mu_0^2 \xi^2}{\text{dist}_0^2} \right) \\ &\leq \text{dist}_0^2 \beta^{2(k+1)}, \end{aligned} \quad (17)$$

where the last line holds because $\beta \geq \underline{\beta} = \sqrt{1 - 2 \frac{\alpha \mu_0}{\text{dist}_0} + \frac{\mu_0^2 \xi^2}{\text{dist}_0^2}}$. Hence, the proof is completed by induction. \square

4 Proof of Theorem 2

We first repeat the DPCP problem. Given a dataset $\tilde{\mathcal{X}} = [\mathcal{X} \ \mathcal{O}] \mathbf{\Gamma} \in \mathbb{R}^{D \times L}$, where $\mathcal{X} \in \mathbb{R}^{D \times N}$ are inlier points spanning a d -dimensional subspace \mathcal{S} of \mathbb{R}^D , \mathcal{O} are outlier points, and $\mathbf{\Gamma}$ is an unknown permutation, the DPCP problem is

$$\underset{\mathbf{B} \in \mathbb{O}(c, D)}{\text{minimize}} f(\mathbf{B}) := \|\tilde{\mathcal{X}}^\top \mathbf{B}\|_{1,2} = \sum_{i=1}^L \|\tilde{\mathbf{x}}_i^\top \mathbf{B}\|_2, \quad (18)$$

One choice for its Riemannian subgradient is

$$\mathcal{G}(\mathbf{B}) = (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{i=1}^L \tilde{\mathbf{x}}_i \text{sign}(\tilde{\mathbf{x}}_i^\top \mathbf{B}), \quad \text{sign}(\mathbf{a}) := \begin{cases} \mathbf{a}/\|\mathbf{a}\|_2 & \text{if } \mathbf{a} \neq \mathbf{0}, \\ \mathbf{0} & \text{if } \mathbf{a} = \mathbf{0}. \end{cases} \quad (19)$$

Let us recall several quantities: the first one, related to the outliers, characterizes the maximum Riemannian subgradient of $\frac{1}{M} \sum_{i=1}^M \|\mathbf{o}_i^\top \mathbf{B}\|_2$:

$$\eta_{\mathcal{O}} := \frac{1}{M} \max_{\mathbf{B} \in \mathbb{O}(c, D)} \left\| (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{i=1}^M \mathbf{o}_i \text{sign}(\mathbf{o}_i^\top \mathbf{B}) \right\|_F, \quad (20)$$

which appears in [6] when \mathbf{B} is on \mathbb{S}^{D-1} . The second one is related to the inliers and is given by

$$c_{\mathcal{X}, \min} := \frac{1}{N} \min_{\mathbf{b} \in \mathcal{S} \cap \mathbb{S}^{D-1}} \|\mathcal{X}^\top \mathbf{b}\|_1, \quad (21)$$

which is referred to as the *permeance statistic* in [7]. These quantities reflect how well distributed the inliers and outliers are, with larger values of $c_{\mathcal{X}, \min}$ (respectively, smaller values of $\eta_{\mathcal{O}}$) corresponding to more uniform distributions of inliers (respectively, outliers).

We require one more result concerning the *principal angles* between two subspaces as follows.

Definition 2. [8] Suppose $\mathbf{U} \in \mathbb{R}^{D \times p}$ and $\mathbf{V} \in \mathbb{R}^{D \times q}$ are two orthonormal bases. Suppose $p \geq q$. Then the principal angles between $\text{Span}(\mathbf{U})$ and $\text{Span}(\mathbf{V})$, $\phi_1(\mathbf{U}, \mathbf{V}) \leq \phi_2(\mathbf{U}, \mathbf{V}) \leq \dots \leq \phi_q(\mathbf{U}, \mathbf{V})$, are defined as

$$\phi_i(\mathbf{U}, \mathbf{V}) = \arccos(\sigma_i(\mathbf{U}^\top \mathbf{V}))$$

for all $i \in \{1, 2, \dots, q\}$, where $\sigma_i(\cdot)$ denotes the i -th largest singular value. The largest principal angle $\phi_q(\mathbf{U}, \mathbf{V})$ is referred to as the *subspace angle* between $\text{Span}(\mathbf{U})$ and $\text{Span}(\mathbf{V})$.

Lemma 1. [8] Suppose $\mathbf{U} \in \mathbb{R}^{D \times p}$ and $\mathbf{V} \in \mathbb{R}^{D \times q}$ are two orthonormal bases and $[\mathbf{U} \ \mathbf{U}^\perp]$ is an orthonormal basis of \mathbb{R}^D . Suppose $p \geq q$. Then the principal angles between $\text{Span}(\mathbf{U})$ and $\text{Span}(\mathbf{V})$ and the principal angles between $\text{Span}(\mathbf{U}^\perp)$ and $\text{Span}(\mathbf{V})$ satisfy the following relationship

$$\left[\frac{\pi}{2}, \dots, \frac{\pi}{2}, \phi_q(\mathbf{U}, \mathbf{V}), \dots, \phi_1(\mathbf{U}, \mathbf{V}) \right] = \left[\frac{\pi}{2} - \phi_1(\mathbf{U}^\perp, \mathbf{V}), \dots, \frac{\pi}{2} - \phi_q(\mathbf{U}^\perp, \mathbf{V}), 0, \dots, 0 \right],$$

where extra $\frac{\pi}{2}$'s and 0's are added on either side to match the sizes.

We now review Theorem 2 and then prove it.

Theorem 2 ($(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian regularity condition for the DPCP problem (18)). *For any $\epsilon < \sqrt{2(1 - M\eta_{\mathcal{O}}/Nc\boldsymbol{\chi}_{\min})}$, the DPCP problem (18) satisfies $(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian regularity condition with $\mathbf{B}^* = [\mathbf{S}^\perp]$ and $\alpha = ((1 - \epsilon^2/2)Nc\boldsymbol{\chi}_{\min} - M\eta_{\mathcal{O}})/\sqrt{2c}$, where \mathbf{S}^\perp is an orthonormal basis for \mathcal{S}^\perp . Also,*

$$\|\mathcal{G}(\mathbf{B})\|_F \leq \sqrt{N} \|\boldsymbol{\chi}\|_2 + M\eta_{\mathcal{O}}, \quad \forall \mathbf{B} \in \mathbb{O}(c, D). \quad (22)$$

Proof of Theorem 2. We first establish the following result which is key to the Riemannian regularity condition for the DPCP problem.

Lemma 2. *Let $\mathbf{B}^* = [\mathbf{S}^\perp]$. Then, for any $\mathbf{B} \in \mathbb{O}(c, D)$, we have*

$$\langle -\mathcal{G}(\mathbf{B}), \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}) - \mathbf{B} \rangle \geq \sin(\phi_{\max})(\cos(\phi_{\max})Nc\boldsymbol{\chi}_{\min} - M\eta_{\mathcal{O}}), \quad (23)$$

where $\phi_{\max} := \phi_{\max}(\mathbf{B}, \mathbf{S}^\perp)$ is the largest principal angle between \mathbf{B} and \mathbf{S}^\perp .

Proof. Let $\mathbf{S} \in \mathbb{R}^{D \times d}$ be an orthonormal basis of the subspace \mathcal{S} and let $\mathbf{S}^\perp \in \mathbb{R}^{D \times c}$ be an orthonormal basis of the orthogonal complement \mathcal{S}^\perp . We rewrite \mathbf{B} as

$$\mathbf{B} = \mathbf{S}\mathbf{S}^\top \mathbf{B} + \mathbf{S}^\perp(\mathbf{S}^\perp)^\top \mathbf{B},$$

where $\mathbf{S}\mathbf{S}^\top \mathbf{B}$ represents the projection of \mathbf{B} onto the subspace \mathcal{S} , while the other term $\mathbf{S}^\perp(\mathbf{S}^\perp)^\top \mathbf{B}$ represents the projection of \mathbf{B} onto the complement \mathcal{S}^\perp . Let $(\mathbf{S}^\perp)^\top \mathbf{B} = \mathbf{U} \cos(\boldsymbol{\Phi}) \mathbf{R}^\top$ be the canonical SVD of $(\mathbf{S}^\perp)^\top \mathbf{B}$, where $\cos(\boldsymbol{\Phi})$ is a diagonal matrix with $\cos(\phi_1), \dots, \cos(\phi_c)$ along its diagonal, $\mathbf{U} \in \mathbb{R}^{c \times c}$, $\mathbf{R} \in \mathbb{R}^{c \times c}$ are orthonormal matrices. Here ϕ_i is the i -th principal angles between \mathbf{B} and \mathbf{S}^\perp . When $\phi_1 = \dots = \phi_c = 0$, it implies that $\mathbf{B} \in [\mathbf{S}^\perp]$, i.e., \mathbf{B} is equivalent to \mathbf{S}^\perp .

Without loss of generality, we assume $c \leq d$.¹ In this case, we can then rewrite $\mathbf{S}^\top \mathbf{B} = \mathbf{V} \sin(\boldsymbol{\Phi}) \mathbf{R}^\top$, where $\mathbf{V} \in \mathbb{R}^{d \times c}$ is an orthonormal matrix. Thus, we have

$$\mathbf{B} = \mathbf{S}\mathbf{V} \sin(\boldsymbol{\Phi}) \mathbf{R}^\top + \mathbf{S}^\perp \mathbf{U} \cos(\boldsymbol{\Phi}) \mathbf{R}^\top. \quad (24)$$

After defining

$$\mathbf{G} = \mathbf{S}\mathbf{V} \cos(\boldsymbol{\Phi}) \sin(\boldsymbol{\Phi}) \mathbf{R}^\top - \mathbf{S}^\perp \mathbf{U} \sin^2(\boldsymbol{\Phi}) \mathbf{R}^\top, \quad (25)$$

we have

$$\begin{aligned} \langle -\mathcal{G}(\mathbf{B}), \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}) - \mathbf{B} \rangle &= \langle -\mathcal{G}(\mathbf{B}), \mathbf{S}^\perp \mathbf{U} \mathbf{R}^\top \rangle \\ &= - \left\langle (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \left(\sum_{j=1}^L \tilde{\mathbf{x}}_j \text{sign}(\tilde{\mathbf{x}}_j^\top \mathbf{B}) \right), \mathbf{S}^\perp \mathbf{U} \mathbf{R}^\top \right\rangle \\ &= \left\langle \sum_{j=1}^N \mathbf{x}_j \text{sign}(\mathbf{x}_j^\top \mathbf{B}) + \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}), \mathbf{G} \right\rangle \\ &= \sum_{j=1}^N \langle \mathbf{x}_j^\top \mathbf{S}\mathbf{V} \cos(\boldsymbol{\Phi}) \sin(\boldsymbol{\Phi}), \text{sign}(\mathbf{x}_j^\top \mathbf{S}\mathbf{V} \sin(\boldsymbol{\Phi})) \rangle - \left\langle \mathbf{G}, (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}) \right\rangle \end{aligned} \quad (26)$$

¹For the case $c > d$, we have at least $\phi_1 = \dots = \phi_{c-d} = 0$. Similar to the case $c \leq d$, we can also rewrite $\mathbf{S}^\top \mathbf{B} = \mathbf{V} \sin(\boldsymbol{\Phi}) \mathbf{R}^\top$, where $\mathbf{V} = [\mathbf{0} \quad \bar{\mathbf{V}}]$ with $\bar{\mathbf{V}} \in \mathbb{R}^{d \times d}$ an orthonormal matrix. Thus, we also have (24) and the following proofs are the same.

where the very last line utilizes $\text{sign}(\mathbf{a}\mathbf{R}^\top) = \text{sign}(\mathbf{a})\mathbf{R}^\top$ ($\mathbf{R} \in \mathbb{R}^{c \times c}$ is an orthonormal matrix) and the fact that $\mathbf{G} \in \text{Span}(\mathbf{B}^\perp)$.

We now bound the first term in the last line of (26) by

$$\begin{aligned}
& \sum_{j=1}^N \langle \mathbf{x}_j^\top \mathbf{S}\mathbf{V} \sin(\Phi) \cos(\Phi), \text{sign}(\mathbf{x}_j^\top \mathbf{S}\mathbf{V} \sin(\Phi)) \rangle \\
& \geq \cos(\phi_c) \sum_{j=1}^N \langle \mathbf{x}_j^\top \mathbf{S}\mathbf{V} \sin(\Phi), \text{sign}(\mathbf{x}_j^\top \mathbf{S}\mathbf{V} \sin(\Phi)) \rangle \\
& = \cos(\phi_c) \sum_{j=1}^N \|\mathbf{x}_j^\top \mathbf{S}\mathbf{V} \sin(\Phi)\|_2 \geq \cos(\phi_c) \sin(\phi_c) \sum_{j=1}^N |\mathbf{x}_j^\top \mathbf{S}\mathbf{v}_c| \\
& \geq \cos(\phi_c) \sin(\phi_c) N c \boldsymbol{\chi}_{\min},
\end{aligned}$$

where the first inequality follows because $0 \leq \phi_1 \leq \phi_2 \leq \dots \leq \phi_c \leq \frac{\pi}{2}$, and the last inequality utilizes (21) since $\mathbf{S}\mathbf{v}_c \in \mathcal{S} \cap \mathbb{S}^{D-1}$. On the other hand, the second term in the last line of (26) can be bounded by

$$\begin{aligned}
& \left| \left\langle \mathbf{G}, (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}) \right\rangle \right| \\
& = \left| \left\langle \mathbf{S}\mathbf{V} \cos(\Phi) \sin(\Phi) \mathbf{R}^\top - \mathbf{S}^\perp \mathbf{U} \sin^2(\Phi) \mathbf{R}^\top, (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}) \right\rangle \right| \\
& \leq \sin(\phi_c) \left| \left\langle \mathbf{S}\mathbf{V} \cos(\Phi) \mathbf{R}^\top - \mathbf{S}^\perp \mathbf{U} \sin(\Phi) \mathbf{R}^\top, (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}) \right\rangle \right| \\
& \leq \sin(\phi_c) \left\| (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}) \right\|_F \leq \sin(\phi_c) M \eta \boldsymbol{\sigma}
\end{aligned}$$

where the first inequality follows because $0 \leq \phi_1 \leq \phi_2 \leq \dots \leq \phi_c \leq \frac{\pi}{2}$, the second inequality utilizes the fact that $\mathbf{S}\mathbf{V} \cos(\Phi) \mathbf{R}^\top - \mathbf{S}^\perp \mathbf{U} \sin(\Phi) \mathbf{R}^\top$ is an orthonormal matrix, and the last inequality follows from (20). This completes the proof. \square

We now turn to prove the $(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian regularity condition. First note that

$$\begin{aligned}
\|\mathcal{P}_{\mathbf{B}^*}(\mathbf{B}) - \mathbf{B}\|_F^2 &= 2 \sum_{i=1}^c (1 - \cos(\phi_i)) \leq 2c(1 - \cos(\phi_c)) \\
&= 4c \sin^2\left(\frac{\phi_c}{2}\right) \leq 2c \sin^2(\phi_c),
\end{aligned} \tag{27}$$

where the last inequality utilizes $\sin(\alpha) \geq \sqrt{2} \sin(\alpha/2)$ for any $\alpha \in [0, \frac{\pi}{2}]$. Combining (27) together with (23) give

$$\langle -\mathcal{G}(\mathbf{B}), \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}) - \mathbf{B} \rangle \geq \frac{\cos(\phi_c) N c \boldsymbol{\chi}_{\min} - M \eta \boldsymbol{\sigma}}{\sqrt{2c}} \|\mathcal{P}_{\mathbf{B}^*}(\mathbf{B}) - \mathbf{B}\|_F. \tag{28}$$

On the other hand, we have

$$\|\mathbf{B} - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B})\|_F^2 = 2 \sum_{i=1}^c (1 - \cos(\phi_i)) \geq 2(1 - \cos(\phi_c)),$$

which implies that $\cos(\phi_c) \geq 1 - \frac{\|\mathbf{B} - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B})\|_F^2}{2}$. This together with (28) and $\|\mathbf{B} - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B})\|_F \leq \epsilon$ complete the proof of the $(\alpha, \epsilon, \mathbf{B}^*)$ -Riemannian regularity condition. The rest is to prove (22).

Proof of (22) For convenience, denote by $\text{sign}(\mathcal{X}^\top \mathbf{B}) = [\text{sign}(\mathbf{x}_1^\top \mathbf{B}) \ \cdots \ \text{sign}(\mathbf{x}_N^\top \mathbf{B})]^\top$. Using (20) and the fact that $\|\text{sign}(\mathcal{X}^\top \mathbf{B})\|_F \leq \sqrt{N}$ allows us to bound the Riemannian subgradient in (19) as

$$\begin{aligned} \|\mathcal{G}(\mathbf{B})\|_F &\leq \left\| (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \mathcal{X} \text{sign}(\mathcal{X}^\top \mathbf{B}) \right\|_F + \left\| (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \sum_{j=1}^M \mathbf{o}_j \text{sign}(\mathbf{o}_j^\top \mathbf{B}) \right\|_F \\ &\leq \left\| (\mathbf{I} - \mathbf{B}\mathbf{B}^\top) \mathcal{X} \right\|_2 \left\| \text{sign}(\mathcal{X}^\top \mathbf{B}) \right\|_F + M\eta\mathfrak{o} \\ &\leq \sqrt{N} \|\mathcal{X}\|_2 + M\eta\mathfrak{o}, \end{aligned}$$

where the second inequality follows from $\|\mathbf{A}\mathbf{B}\|_F \leq \|\mathbf{A}\|_2 \|\mathbf{B}\|_F$. □

5 Proof of Corollary 2

Definition 3 (Random model for ODL [9]). *Assume $\mathbf{A} \in \mathbb{R}^{D \times D}$ is a fixed but unknown orthonormal matrix. The data is generated as $\tilde{\mathcal{X}} = \mathbf{A}\mathbf{S}$, where each column of $\mathbf{S} \in \mathbb{R}^{D \times D}$ is an i.i.d. Bernoulli-Gaussian random vector with parameter $\rho \in (0, 1)$ that controls the sparsity.*

We first repeat the Riemannian regularity condition for the ODL problem.

Theorem 3. [9, Theorem 3.6] *Assume $\rho \in [1/D, 1/2]$ in the random model of Definition 3. There exist universal constants $C, c > 0$ such that if $N \geq CD^4 \zeta^{-2} \rho^{-2} \log(D/\zeta)$, $\forall \zeta \in (0, 1)$, then with probability at least $1 - \exp(-cN\rho^3 \zeta^2 D^{-3} / \log N)$ the ODL problem satisfies (6) for any $\mathbf{b} \in \mathcal{I}_\zeta^i$ with $\mathcal{G}(\mathbf{b})$ and $\mathbf{B}^* = \mathbf{e}_i$ for any i , and*

$$\alpha = \frac{1}{16} \rho (1 - \rho) \zeta D^{-\frac{3}{2}}. \quad (29)$$

Now we repeat Corollary 2 and then prove it.

Corollary 2. *Let $\{\mathbf{b}_k\}$ be the sequence generated by Algorithm 1 for the ODL problem with $\mathbf{b}_0 \in \mathcal{I}_\zeta^i$ ($\zeta \leq \frac{55}{64}$) and step size $\mu_k = \mu_0 \beta^k$, where μ_0 and β satisfy the conditions in Theorem 1 with $\xi = 2$ and $\epsilon = \sqrt{2}$, and $\alpha = \frac{1}{16} \rho (1 - \rho) \zeta D^{-\frac{3}{2}}$. Under the same setup as in Theorem 3, with probability at least $1 - \exp(-cN\rho^3 \zeta^2 D^{-3} / \log N)$, $\{\mathbf{b}_k\}$ converges to \mathbf{e}_i at an R -linear rate, i.e.,*

$$\text{dist}(\mathbf{b}_k, \mathbf{e}_i) \leq \beta^k \text{dist}(\mathbf{b}_0, \mathbf{e}_i). \quad (30)$$

Proof. To apply Theorem 1, we require an upper bound on the norm of the Riemannian subgradients. It follows from [9, Proposition 3.7] that if $N \geq CD \log D$, then $\sup_{\mathbf{B} \in \mathbb{S}^{D-1}} \|\mathcal{G}(\mathbf{B})\|_2 \leq 2$ with probability at least $1 - \exp(-cN\rho/\log N)$. We also require $\mathbf{b}_k \in \mathcal{I}_\zeta^i$ for all k so that the Riemannian regularity condition (6) holds at all the iterates. A sufficient condition to guarantee $\mathbf{b}_k \in \mathcal{I}_\zeta^i$ is that $\mu_k \leq \min\{\frac{1}{100}, \frac{1-\zeta}{9}\} \frac{1}{D^{1/2}}$ [9, Proposition D.2]. Plugging (29), $\epsilon = \sqrt{2}$, and $\xi = 2$ into (14) gives $\mu_k \leq \mu_0 \leq \frac{\zeta}{64D^{3/2}} \leq \min\{\frac{1}{100}, \frac{1-\zeta}{9}\} \frac{1}{D^{1/2}}$ since $\zeta \leq \frac{55}{64}$. \square

6 Initialization

Lemma 3. *Consider a spectral initialization \mathbf{B}_0 by taking the bottom c eigenvectors of $\tilde{\mathcal{X}}\tilde{\mathcal{X}}^\top$. Then, it satisfies*

$$\|\mathbf{B}_0 - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_0)\|_F^2 \leq \frac{\sum_{j=1}^c \sigma_j^2(\mathcal{O}) - \sum_{j=D-c+1}^D \sigma_j^2(\mathcal{O})}{\sigma_d^2(\mathcal{X})}, \quad (31)$$

where σ_ℓ denotes the ℓ -th largest singular value.

Proof. Note that for any $\mathbf{B} \perp \mathcal{S}$, $\|\tilde{\mathcal{X}}^\top \mathbf{B}\|_F^2 = \|\mathcal{O}^\top \mathbf{B}\|_F^2 = \text{trace}(\mathbf{B}^\top \mathcal{O} \mathcal{O}^\top \mathbf{B}) = \sum_{j=1}^c \mathbf{b}_j^\top \mathcal{O} \mathcal{O}^\top \mathbf{b}_j \leq \sum_{j=1}^c \sigma_j^2(\mathcal{O})$. Thus, since \mathbf{B}_0 is the optimal solution to $\arg \min_{\mathbf{B} \in \mathcal{O}(c,D)} \|\tilde{\mathcal{X}}^\top \mathbf{B}\|_F^2$, we have

$$\|\tilde{\mathcal{X}}^\top \mathbf{B}_0\|_F^2 \leq \sum_{j=1}^c \sigma_j^2(\mathcal{O}).$$

On the other hand, let \mathbf{S} be an orthonormal basis for \mathcal{S} and let Θ be the coefficients of \mathcal{X} in \mathbf{S} , i.e., $\mathcal{X} = \mathbf{S}\Theta$, we have

$$\begin{aligned} \|\tilde{\mathcal{X}}^\top \mathbf{B}_0\|_F^2 &= \|\mathcal{X}^\top \mathbf{B}_0\|_F^2 + \|\mathcal{O}^\top \mathbf{B}_0\|_F^2 = \|\mathcal{X}^\top \mathbf{S} \mathbf{S}^\top \mathbf{B}_0\|_F^2 + \|\mathcal{O}^\top \mathbf{B}_0\|_F^2 \\ &= \|\Theta^\top \mathbf{S}^\top \mathbf{B}_0\|_F^2 + \|\mathcal{O}^\top \mathbf{B}_0\|_F^2 \geq \sigma_{\min}^2(\Theta) \|\mathbf{S}^\top \mathbf{B}_0\|_F^2 + \|\mathcal{O}^\top \mathbf{B}_0\|_F^2 \\ &\geq \sigma_d^2(\mathcal{X}) \|\mathbf{B}_0 - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_0)\|_F^2 + \sum_{j=D-c+1}^D \sigma_j^2(\mathcal{O}), \end{aligned}$$

where we first utilize the fact that \mathcal{X} lies in \mathcal{S} such that $\mathcal{X} = \mathbf{S} \mathbf{S}^\top \mathcal{X}$, the inequality follows because $\|\mathbf{A} \mathbf{B}\|_F^2 = \text{trace}(\mathbf{A}^\top \mathbf{A} \mathbf{B} \mathbf{B}^\top) \geq \sigma_{\min}(\mathbf{A}^\top \mathbf{A}) \|\mathbf{B} \mathbf{B}^\top\|_F$ for any \mathbf{A}, \mathbf{B} , and the last line follows from $\|\mathbf{S}^\top \mathbf{B}_0\|_F^2 = \|\mathbf{S} \mathbf{S}^\top \mathbf{B}_0\|_F^2 = \|\mathbf{B}_0 - \mathbf{S}^\perp (\mathbf{S}^\perp)^\top \mathbf{B}_0\|_F^2 = \|\mathbf{B}_0 - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_0)\|_F^2$. Combining the above two equations gives

$$\|\mathbf{B}_0 - \mathcal{P}_{\mathbf{B}^*}(\mathbf{B}_0)\|_F^2 \leq \frac{\sum_{j=1}^c \sigma_j^2(\mathcal{O}) - \sum_{j=D-c+1}^D \sigma_j^2(\mathcal{O})}{\sigma_d^2(\mathcal{X})}.$$

\square

7 Random Spherical Model

In this section, we consider the following random spherical model.

Definition 1. For any given subspace \mathcal{S} of dimension $d < D$, a random spherical model refers to that the columns of \mathbf{O} drawn independently and uniformly at random from the unit sphere \mathbb{S}^{D-1} , and the columns of \mathbf{X} are drawn independently and uniformly at random from the intersection of the unit sphere with the subspace \mathcal{S} .

We require the following result from [10, Lemma 4] concerning $c_{\mathbf{X},\min}$.

Lemma 4. [10, Lemma 4] Under the random spherical model in Definition 1, we have

$$\mathbb{P}\left(c_{\mathbf{X},\min} \geq c_d - (2 + \frac{t}{2})/\sqrt{N}\right) \geq 1 - 2e^{-\frac{t^2}{2}},$$

where

$$c_D := \frac{(D-2)!!}{(D-1)!!} \cdot \begin{cases} \frac{2}{\pi}, & D \text{ is even,} \\ 1, & D \text{ is odd,} \end{cases} \text{ where } k!! := \begin{cases} k(k-2)(k-4)\cdots 4 \cdot 2, & k \text{ is even,} \\ k(k-2)(k-4)\cdots 3 \cdot 1, & k \text{ is odd.} \end{cases} \quad (32)$$

Lemma 5. Let $\mathbf{o}_1, \dots, \mathbf{o}_M$ be uniformly distributed on \mathbb{S}^{D-1} . Then for any $t > 0$

$$\mathbb{P}\left[\eta_{\mathbf{O}} \gtrsim \frac{\sqrt{c_D D} \log(c_D D) + t}{\sqrt{M}}\right] \leq 2 \exp(-t^2/2), \quad (33)$$

where $c = D - d$ is the co-dimensions.

Its proof is given in Section 10.

The following results provide concentration inequalities for the singular values appeared in (31) when the inliners and outliers are generated from a random spherical model.

Lemma 6. [11, Theorem 5.39] Under the random spherical model in Definition 1, then for every $t > 0$, there exist constants C_1, C_2 such that

$$\begin{aligned} \mathbb{P}\left(\sigma_1(\mathbf{O}) \geq \frac{\sqrt{M} + C_2\sqrt{D} + t}{\sqrt{D}}\right) &\leq e^{-C_1 t^2}, \\ \mathbb{P}\left(\sigma_D(\mathbf{O}) \leq \frac{\sqrt{M} - C_2\sqrt{D} - t}{\sqrt{D}}\right) &\leq e^{-C_1 t^2}, \\ \mathbb{P}\left(\sigma_d(\mathbf{X}) \leq \frac{\sqrt{N} - C_2\sqrt{d} - t}{\sqrt{d}}\right) &\leq e^{-C_1 t^2}. \end{aligned} \quad (34)$$

The following result establishes that the spectral initialization satisfies the condition $\text{dist}^2(\mathbf{B}_0, \mathbf{S}^\perp) < 2(1 - M\eta_{\mathbf{O}}/Nc_{\mathbf{X},\min})$ with high probability when the data are generated from a random spherical model.

Corollary 3. Consider the same random spherical model as in Definition 1. Then for any positive number $t < \min\{\sqrt{N} - C_2\sqrt{d}, 2c_d\sqrt{N} - 4\}$, with probability at least $1 - 4e^{-t^2/2} - 3e^{-C_1 t^2}$, the spectral initialization \mathbf{B}_0 in Lemma 3 satisfies the condition $\text{dist}^2(\mathbf{B}_0, \mathbf{S}^\perp) < 2(1 - M\eta_{\mathbf{O}}/Nc_{\mathbf{X},\min})$,

provided that

$$\begin{aligned} \frac{2cd\sqrt{M}(C_2\sqrt{D}+t)}{D(\sqrt{N}-C_2\sqrt{d}-t)^2} + \frac{C_0(\sqrt{cMD}\log(c_D D) + \sqrt{Mt})}{c_d N - (2 + \frac{t}{2})\sqrt{N}} < 1, \\ c_d N - (2 + t/2)\sqrt{N} > C_0 \left(\sqrt{cD}\log(c_D D) + t \right) \sqrt{M}, \end{aligned} \quad (35)$$

where c_d and c_D are defined in (32), and C_0, C_1 and C_2 are universal constants independent of N, M, D, d and t .

Proof. This follows by combining Lemma 3, Lemma 4, Lemma 5, and Lemma 6. \square

Note that the first line in (35) suggests $O\left(\frac{cd\sqrt{M}}{\sqrt{DN}} + \frac{\sqrt{cD}\log D\sqrt{M}}{N}\right) < 1$, while the second line of (35) suggests that $O\left(\frac{\sqrt{cD}\log D\sqrt{M}}{N}\right) < 1$. The combination of both implies that the projected Riemannian subgradient method with a spectral initialization can converge to \mathbf{S}^\perp in a linear rate when there are $M = O\left(\frac{D}{c^2(d+D\log D)^2}N^2\right)$ outliers.

8 Comparison with the Regularity Condition for Smooth Function

Aside from the weak convexity and sharpness, another regularity condition related to Definition 1 is the one proposed in [12]: we say a continuously differentiable function g satisfies the $(\alpha, \gamma, \epsilon)$ -regularity condition if for all $\mathbf{x} \in \mathbb{R}^D$ such that $\text{dist}(\mathbf{x}, \mathcal{X}) \leq \epsilon$, we have

$$\begin{aligned} \langle \mathcal{P}_{\mathcal{X}}(\mathbf{x}) - \mathbf{x}, -\nabla g(\mathbf{x}) \rangle \\ \geq \alpha \text{dist}^2(\mathbf{x}, \mathcal{X}) + \gamma \|\nabla g(\mathbf{x})\|^2. \end{aligned} \quad (36)$$

We now compare (6) with (36). On one hand, (36) has similar form to (6) as both attempt to provide lower bounds for the inner product between the gradient (or Riemannian subgradient) and the vector $\mathbf{x} - \mathcal{P}_{\mathcal{X}}(\mathbf{x})$ for any \mathbf{x} that is close to \mathcal{X} . On the other hand, (36) mainly differs from (6) in two aspects. First note that unlike (36), there is no $\|\mathcal{G}(\mathbf{B})\|$ term (which is $\|\nabla g(\mathbf{x})\|^2$ in (36)) in (6). This is mainly because as we illustrated before, the Riemannian subgradient does not vanish even when \mathbf{B} approaching \mathbf{B}^* . Thus, it is impossible to include the $\|\mathcal{G}(\mathbf{B})\|$ term into (6) as its left hand side (LHS) goes to 0 when \mathbf{B} tends to \mathbf{B}^* . Besides, (6) involves the term $\text{dist}(\mathbf{B}, \mathbf{B}^*)$, while (36) has the term $\text{dist}^2(\mathbf{x}, \mathcal{X})$. If we apply the Cauchy-Schwarz inequality to the LHS of (36), we obtain

$$\gamma \|\nabla g(\mathbf{x})\| \leq \text{dist}(\mathbf{x}, \mathcal{X}) - \frac{\text{dist}^2(\mathbf{x}, \mathcal{X})}{\|\nabla g(\mathbf{x})\|},$$

which implies $\nabla g(\mathbf{x}) \rightarrow \mathbf{0}$ when $\text{dist}(\mathbf{x}, \mathcal{X}) \rightarrow 0$. This is in sharp contrast to (8). Informally speaking, the regularity condition in (36) describes certain geometric property of smooth functions while the Riemannian regularity condition in Definition 1 characterizes certain geometric property of non-smooth functions.

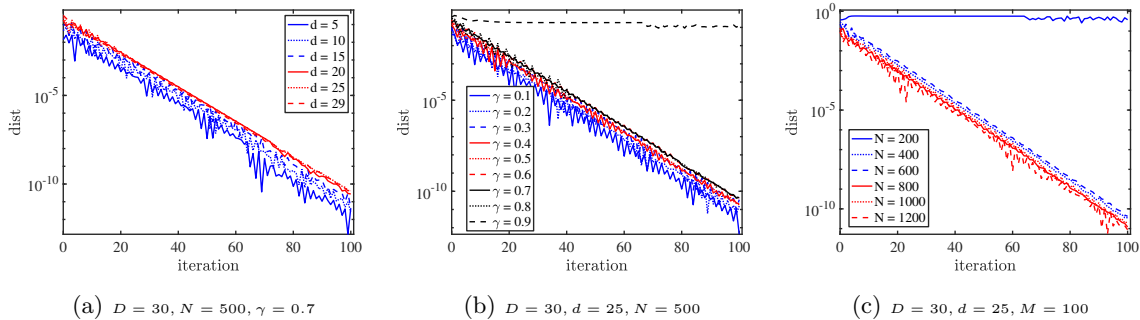


Figure 1: The convergence of Algorithm 1 (the Riemannian SubGM with geometrically diminishing step size) for the DPCP problems with $\beta = 0.8$ and μ_0 determined by line search method. Here $\gamma = \frac{M}{M+N}$ denotes the outlier ratio.

9 Additional Experiments

9.1 DPCP for Robust Subspace Learning

We use synthetic experiments under different settings to further verify the projected Riemannian subgradient method with geometrically diminishing step sizes for the DPCP problem. We randomly sample a subspace \mathcal{S} of dimension $d < D - 1$, and uniformly at random sample N inliers and M outliers with unit ℓ_2 -norm. Denote by $\gamma = \frac{M}{M+N}$ the outlier ratio. We set $\beta = 0.8$ for geometrically diminishing step size with initial step size obtained by one iteration of a backtracking line search. We define \mathbf{B}_0 to be the bottom eigenvector of $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^\top$. Figure 1 displays the convergence of θ (to 0) with different d , N and outlier ratio γ . In particular, Figure 1a shows the convergence of θ with $D = 30, N = 500, \gamma = 0.7$ and different subspace dimension d . We observe linear convergence in this case, irrespectively the subspace dimension d . In Figure 1b, we set $D = 30, d = 25, N = 500$ and vary the outlier ratio γ from 0.1 to 0.9. We also observe linear convergence except for the case $\gamma = 0.9$, in which we have much more outliers than inliers. Finally we display experiments with varied N in Figure 1c. We also observe linear convergence for Algorithm 1 given sufficient number of inliers.

9.2 Orthogonal Dictionary Learning

As illustrated in the paper, we first generate a random orthogonal dictionary $\mathbf{A} \in \mathbb{R}^{D \times D}$ with $D = 70$. Set the sparsity level $\rho = 0.3$ and the number of data points $N \approx 10D^{1.5} = 5857$. The initialization \mathbf{B}_0 is randomly generated from the unit sphere \mathbb{S}^{D-1} . Figure 2 displays the effect of the initial step size μ_0 and the decaying factor β for Algorithm 1 with geometrically diminishing step size $\mu_k = \mu_0\beta^k$. We observe similar phenomena as for the DPCP problem. First observe from Figure 2a that, as expected, β controls the convergence speed: a value of β too small ($\beta = 0.7$) may result in no convergence, in agreement with (14) and (15); whereas when $\beta \geq 0.8$, the algorithm converges in a linear rate, with a larger value of β resulting in a slower convergence speed (comparing $\beta = 0.8, 0.9, 0.95$). Figure 2b displays the effect of μ_0 when β is fixed. We observe that a value of μ_0 too small ($\mu_0 = 1$) results in no convergence. This can be explained following the discussion after Theorem 1: when μ_0 is small, then the smallest allowable decaying factor $\underline{\beta}$ in (14) increases

when μ_0 decreases and particularly $\underline{\beta} \rightarrow 1$ when $\mu_0 \rightarrow 0$, thus contradicting the requirement $\beta \geq \underline{\beta}$ in (14) when we fix $\beta = 0.9$ and set $\underline{\mu}_0$ too small ($\mu_0 = 1$).

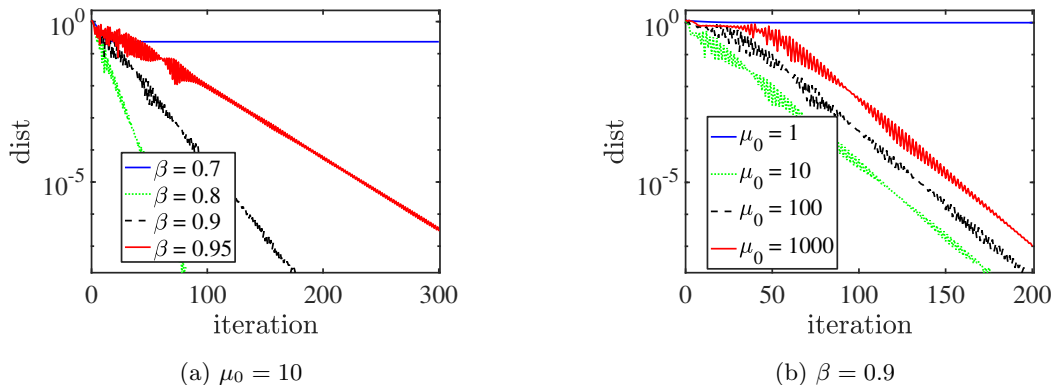


Figure 2: Convergence of Algorithm 1 with different initial step size μ_0 and decaying factor β for dictionary learning.

10 Proof of Lemma 5

10.1 Preliminaries

Suppose X_1, \dots, X_n are n independent and identically distributed (i.i.d.) random observations from a probability measure P on a measurable space $(\mathcal{X}, \mathcal{A})$. Given a measurable function $f : \mathcal{X} \rightarrow \mathbb{R}$, the *empirical process* evaluated at f is defined as

$$\mathbb{G}_n f := \sqrt{n} \left(\frac{1}{n} \sum_{i=1}^n f(X_i) - \int f dP \right), \quad (37)$$

where $\int f dP$ is the expectation of f under P and $\frac{1}{n} \sum_{i=1}^n f(X_i)$ is called the *empirical distribution*. There are several results concerning the supreme of $\mathbb{G}_n f$ over a given class \mathcal{F} of measurable functions.

Define an *envelope function* $F : \mathcal{X} \rightarrow \mathbb{R}$ such that $|f| \leq F$ for every $f \in \mathcal{F}$. The $L_r(P)$ -norm is defined as $\|f\|_{L_r(P)} = (\int |f|^r dP)^{1/r}$. We need one more definition for the so-called *bracket number* which (informally speaking) measures the size of a class functions \mathcal{F} . Given two functions l and u , the *bracket* $[l, u]$ is the set of all functions f with $l \leq f \leq u$. An ϵ -bracket in $L_r(P)$ is a bracket $[l, u]$ with $\int (u - l)^r dP \leq \epsilon^r$ (since $l \leq u$, it is equivalent to say $\|u - l\|_{L_r(P)} \leq \epsilon$). The bracket number $N_{[]}(\epsilon, \mathcal{F}, L_2(P))$ is the minimum number of ϵ -brackets needed to cover \mathcal{F} .

Lemma 7 ([13], Cor. 19.35). *For any class \mathcal{F} of measurable functions with envelope function F ,*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} |\mathbb{G}_n f| \right] \lesssim J_{[]}(\|F\|_{P,2}, \mathcal{F}, L_2(P)), \quad (38)$$

where $J_{[]}(\|F\|_{P,2}, \mathcal{F}, L_2(P))$ is called the bracketing integral:

$$J_{[]}(\|F\|_{L_2(P)}, \mathcal{F}, L_2(P)) = \int_0^{\|F\|_{L_2(P)}} \sqrt{\log(N_{[]}(\epsilon, \mathcal{F}, L_2(P)))} d\epsilon. \quad (39)$$

Lemma 8 (McDiarmid's Inequality, [14]). *Let Z_1, \dots, Z_n be real-valued independent random variables. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a function that satisfies*

$$\sup_{z_1, \dots, z_n, z'_i} \left| f(z_1, \dots, z_{i-1}, z_i, z_{i+1}, \dots, z_n) - f(z_1, \dots, z_{i-1}, z'_i, z_{i+1}, \dots, z_n) \right| \leq c_i,$$

for every $i = 1, \dots, n$. Then

$$\mathbb{P} \left[\left| f(Z_1, \dots, Z_n) - \mathbb{E} \left[f(Z_1, \dots, Z_n) \right] \right| \geq \epsilon \right] \leq 2 \exp \left(- \frac{2\epsilon^2}{\sum_{i=1}^n c_i^2} \right).$$

Lemma 9 (Vector-valued Comparison Inequality for Rademacher Process, [15], Corollary 2). *Let \mathcal{F} be a class of functions $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^c$ and let $h_i : \mathbb{R}^c \rightarrow \mathbb{R}$ be 1-Lipschitz functions. Then, for any $\mathbf{v}_1, \dots, \mathbf{v}_N \in \mathbb{R}^D$, we have*

$$\mathbb{E} \left[\sup_{\mathbf{f} \in \mathcal{F}} \sum_{i=1}^N \varepsilon_i h_i(\mathbf{f}(\mathbf{v}_i)) \right] \leq \sqrt{2} \mathbb{E} \left[\sup_{\mathbf{f} \in \mathcal{F}} \sum_{i=1}^N \varepsilon_i^\top \mathbf{f}(\mathbf{v}_i) \right], \quad (40)$$

where ε_i are independent Rademacher random variables, and each $\varepsilon_i \in \mathbb{R}^c$ is independent and it contains independent Rademacher random variables.

Lemma 10 (Rademacher Symmetrization, [16], Thm. 1.1). *Let \mathcal{F} be a class of functions $f : \mathbb{R}^D \rightarrow \mathbb{R}$ such that $0 \leq f(\mathbf{v}) \leq 1$. Let ε_i be Rademacher random variables. Then for independent and identically distributed random variables $\mathbf{v}_1, \dots, \mathbf{v}_n$, we have*

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left(\frac{1}{N} \sum_{i=1}^N f(\mathbf{v}_i) - \mathbb{E}[f(\mathbf{v})] \right) \right] \leq 2 \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N \varepsilon_i f(\mathbf{v}_i) \right], \quad (41)$$

$$\mathbb{E} \left[\sup_{f \in \mathcal{F}} \left(\mathbb{E}[f(\mathbf{v})] - \frac{1}{N} \sum_{i=1}^N f(\mathbf{v}_i) \right) \right] \leq 2 \mathbb{E} \left[\sup_{f \in \mathcal{F}} \frac{1}{N} \sum_{i=1}^N \varepsilon_i f(\mathbf{v}_i) \right]. \quad (42)$$

We also require the covering number of $\mathbb{S}(D, c)$, which could be easily derived from the standard result for the sphere. Denote by \mathcal{N}_ϵ an ϵ -net of $\mathbb{S}(D, c)$ if every point $\mathbf{B} \in \mathbb{S}(D, c)$ can be approximated to within ϵ by some point $\mathbf{B}' \in \mathcal{N}_\epsilon$. The minimal cardinality of an ϵ -net, denoted by $\mathcal{N}(\mathbb{S}^{D-1}, \epsilon)$, is called the covering number of $\mathbb{S}(D, c)$.

Lemma 11. (Covering Number of $\mathbb{S}(D, c)$, [11, Lemma 5.2]) *For every $\epsilon > 0$, the covering number of the sphere \mathbb{S}^{D-1} satisfies*

$$\mathcal{N}(\mathbb{S}^{D-1}, \epsilon) \leq \left(1 + \frac{2}{\epsilon} \right)^{cD}. \quad (43)$$

We finally require one more result concerning the probability that $\|\text{sign}(\mathbf{o}^\top \mathbf{B}) - \text{sign}(\mathbf{o}^\top \mathbf{B}')\|$ is small when \mathbf{B} is very close to \mathbf{B}' .

Lemma 12. Denote by $\mathbb{B}(\mathbf{B}, \epsilon_1)$ the set of points that around \mathbf{B} :

$$\mathbb{B}(\mathbf{B}, \epsilon_1) := \{\mathbf{B}' \in \mathbb{O}(c, D) : \|\mathbf{B} - \mathbf{B}'\|_2 \leq \epsilon_1\}.$$

Let $\mathbf{o} \in \mathbb{S}^{D-1}$ be drawn independently and uniformly at random from the unit sphere \mathbb{S}^{D-1} . For any $\mathbf{B} \in \mathbb{S}^{D-1}$ and $\epsilon_2 > 0$, define

$$\bar{\mathbb{A}} := \{\mathbf{o} \in \mathbb{R}^D : \|\text{sign}(\mathbf{o}^\top \mathbf{B}) - \text{sign}(\mathbf{o}^\top \mathbf{B}')\| \leq \epsilon_2, \forall \mathbf{B}' \in \mathbb{B}(\mathbf{B}, \epsilon_1)\}. \quad (44)$$

Then

$$\mathbb{P}[\mathbf{o} \in \bar{\mathbb{A}}^c] \lesssim c_D D \frac{\epsilon_1^2}{\epsilon_2^2},$$

where \lesssim means smaller than up to a universal constant which is independent of D .

Proof. We first bound the difference between $\text{sign}(\mathbf{o}^\top \mathbf{B})$ and $\text{sign}(\mathbf{o}^\top \mathbf{B}')$ by

$$\begin{aligned} \|\text{sign}(\mathbf{o}^\top \mathbf{B}) - \text{sign}(\mathbf{o}^\top \mathbf{B}')\| &= \left\| \frac{\mathbf{o}^\top \mathbf{B}}{\|\mathbf{o}^\top \mathbf{B}\|} - \frac{\mathbf{o}^\top \mathbf{B}'}{\|\mathbf{o}^\top \mathbf{B}'\|} \right\| = \left\| \frac{\|\mathbf{o}^\top \mathbf{B}'\| \mathbf{o}^\top \mathbf{B} - \|\mathbf{o}^\top \mathbf{B}\| \mathbf{o}^\top \mathbf{B}'}{\|\mathbf{o}^\top \mathbf{B}\| \|\mathbf{o}^\top \mathbf{B}'\|} \right\| \\ &= \left\| \frac{\|\mathbf{o}^\top \mathbf{B}'\| \mathbf{o}^\top (\mathbf{B} - \mathbf{B}') - (\|\mathbf{o}^\top \mathbf{B}\| - \|\mathbf{o}^\top \mathbf{B}'\|) \mathbf{o}^\top \mathbf{B}'}{\|\mathbf{o}^\top \mathbf{B}\| \|\mathbf{o}^\top \mathbf{B}'\|} \right\| \\ &\leq \frac{\|\mathbf{B} - \mathbf{B}'\|}{\|\mathbf{o}^\top \mathbf{B}\|} + \frac{|\|\mathbf{o}^\top \mathbf{B}\| - \|\mathbf{o}^\top \mathbf{B}'\||}{\|\mathbf{o}^\top \mathbf{B}\|} \\ &\leq 2 \frac{\|\mathbf{B} - \mathbf{B}'\|}{\|\mathbf{o}^\top \mathbf{B}\|} \leq 2 \frac{\epsilon_1}{\|\mathbf{o}^\top \mathbf{B}\|}. \end{aligned}$$

Thus, as long as $\|\mathbf{o}^\top \mathbf{B}\| \geq \frac{\epsilon_1}{2\epsilon_2}$, we have $\|\text{sign}(\mathbf{o}^\top \mathbf{B}) - \text{sign}(\mathbf{o}^\top \mathbf{B}')\| \leq \epsilon_2$. Without loss of generality, suppose $\mathbf{B} = [\mathbf{e}_1 \ \cdots \ \mathbf{e}_c]$. Then, the probability that $\mathbf{o} \in \bar{\mathbb{A}}^c$ is bounded by the probability that $\|\mathbf{o}^\top \mathbf{B}\| \leq \frac{\epsilon_1}{2\epsilon_2}$:

$$\mathbb{P}[\mathbf{o} \in \bar{\mathbb{A}}^c] \leq \mathbb{P}\left[\|\mathbf{o}^\top \mathbf{B}\| \leq \frac{\epsilon_1}{2\epsilon_2}\right] \leq \mathbb{P}\left[o_1 \leq \frac{\epsilon_1}{2\epsilon_2}\right] \lesssim c_D D \frac{\epsilon_1^2}{\epsilon_2^2}.$$

where o_1 is the first element in \mathbf{o} , and the last inequality follows from [10, Lemma 12]. \square

10.2 Proof of Lemma 5

Before givin out the main proofs, we first preset the following useful result concerning the expectation of $\eta_{\mathbf{o}}$.

Lemma 13. Suppose $\mathbf{o}_1, \dots, \mathbf{o}_M$, are drawn independently and uniformly at random from the unit sphere \mathbb{S}^{D-1} . Then

$$\mathbb{E} \left[\sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j=1}^M \langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \rangle \right| \right] \lesssim \sqrt{cD} \log(\sqrt{cD}) \sqrt{M}, \quad (45)$$

where \lesssim means smaller than up to a universal constant which is independent of D and M .

Proof. The main idea for proving Lemma 13 is to view

$$\frac{1}{\sqrt{M}} \sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j=1}^M \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \right\rangle \right|$$

as an empirical process and then utilize Lemma 7. Towards that end, define the set

$$\mathbb{F} := \{(\mathbf{B}, \mathbf{G}) : \mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}\}.$$

We further define the parameterized function as

$$f_{\mathbf{B}, \mathbf{G}}(\mathbf{o}) := \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}), \mathbf{G}^\top \mathbf{o} \right\rangle.$$

The class of functions we are interested in is $\mathcal{F} := \{f_{\mathbf{B}, \mathbf{G}} : (\mathbf{B}, \mathbf{G}) \in \mathbb{F}\}$.

Note that for any $f_{\mathbf{B}, \mathbf{G}} \in \mathcal{F}$ (i.e., $(\mathbf{B}, \mathbf{G}) \in \mathbb{F}$), we have

$$\mathbb{E} [f_{\mathbf{B}, \mathbf{G}}(\mathbf{o})] = \mathbb{E} \left[\left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}), \mathbf{G}^\top \mathbf{o} \right\rangle \right] = 0,$$

which together with (37) indicates that

$$\sum_{j=1}^M \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \right\rangle = \sqrt{M} \mathbb{G}_M f_{\mathbf{B}, \mathbf{G}},$$

where $\mathbb{G}_M f_{\mathbf{B}, \mathbf{G}}$ is the empirical process of $f_{\mathbf{B}, \mathbf{G}}$.

To utilize Lemma 7, the rest of the proof is to show the corresponding bracketing integral is finite for our problem. Since $|f_{\mathbf{B}, \mathbf{G}}(\mathbf{o})| \leq \|\mathbf{o}\|_2$ for any $(\mathbf{B}, \mathbf{G}) \in \mathbb{F}$, we know $F(\mathbf{o}) = \|\mathbf{o}\|_2$ is the envelope function of \mathcal{F} and $\|F\|_{P, 2} = 1$. Thus, we only need to consider the bracket integral $J_{[]} (1, \mathcal{F}, L_2(P))$, where P is now a probability measure on $\mathbb{B}(\mathbf{B}, \mathbf{G})$. To that end, we first compute the bracket number $N_{[]}(\epsilon, \mathcal{F}, L_2(P))$.

Since our function $f_{\mathbf{B}, \mathbf{G}}$ is parameterized by (\mathbf{B}, \mathbf{G}) , covering the class of functions \mathcal{F} is related to covering the set \mathbb{F} . For any fixed $(\mathbf{B}, \mathbf{G}) \in \mathbb{F}$, define the set of points that around (\mathbf{B}, \mathbf{G}) :

$$\mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1) := \left\{ (\mathbf{B}', \mathbf{G}') \in \mathbb{F} : \sqrt{\|\mathbf{B} - \mathbf{B}'\|_F^2 + \|\mathbf{G} - \mathbf{G}'\|_F^2} \leq \epsilon_1 \right\}.$$

Then, denote by

$$\mathbb{A} := \left\{ \mathbf{o} \in \mathbb{R}^D : \|\text{sign}(\mathbf{o}^\top \mathbf{B}) - \text{sign}(\mathbf{o}^\top \mathbf{B}')\| \leq \epsilon_2, \forall (\mathbf{B}', \mathbf{G}') \in \mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1) \right\}.$$

When \mathbf{B} is close to \mathbf{B}' , then \mathbb{A} should cover most of \mathbf{o} . If $\mathbf{o} \in \mathbb{A}$, then for any $(\mathbf{B}', \mathbf{G}') \in \mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1)$ we have

$$\begin{aligned} |f_{\mathbf{B}, \mathbf{G}}(\mathbf{o}) - f_{\mathbf{B}', \mathbf{G}'}(\mathbf{o})| &= \left| \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}), \mathbf{G}^\top \mathbf{o} \right\rangle - \left\langle \text{sign}(\mathbf{B}'^\top \mathbf{o}), \mathbf{G}'^\top \mathbf{o} \right\rangle \right| \\ &= \left| \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}), (\mathbf{G} - \mathbf{G}')^\top \mathbf{o} \right\rangle - \left\langle (\text{sign}(\mathbf{B}'^\top \mathbf{o}) - \text{sign}(\mathbf{B}^\top \mathbf{o})), \mathbf{G}'^\top \mathbf{o} \right\rangle \right| \\ &\leq \|\mathbf{G} - \mathbf{G}'\| + \left\| \text{sign}(\mathbf{B}'^\top \mathbf{o}) - \text{sign}(\mathbf{B}^\top \mathbf{o}) \right\| \\ &\leq \epsilon_1 + \epsilon_2. \end{aligned}$$

On the other hand, if $\mathbf{o} \in \mathbb{A}^c$, then for any $(\mathbf{B}', \mathbf{G}') \in \mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1)$ we have

$$|f_{\mathbf{B}, \mathbf{G}}(\mathbf{o}) - f_{\mathbf{B}', \mathbf{G}'}(\mathbf{o})| = \left| \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}), \mathbf{G}^\top \mathbf{o} \right\rangle \right| + \left| \left\langle \text{sign}(\mathbf{B}'^\top \mathbf{o}), \mathbf{G}'^\top \mathbf{o} \right\rangle \right| \leq 2.$$

To summary, we have

$$|f_{\mathbf{B}, \mathbf{G}}(\mathbf{o}) - f_{\mathbf{B}', \mathbf{G}'}(\mathbf{o})| \leq \epsilon_1 \delta_{\mathbb{A}}(\mathbf{o}) + 2\delta_{\mathbb{A}^c}(\mathbf{o}), \quad \forall (\mathbf{B}', \mathbf{G}') \in \mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1). \quad (46)$$

We now define a bracket $[l, u]$ by

$$\begin{aligned} l(\mathbf{o}) &= f_{\mathbf{B}, \mathbf{G}}(\mathbf{o}) - \epsilon_1 \delta_{\mathbb{A}}(\mathbf{o}) - 2\delta_{\mathbb{A}^c}(\mathbf{o}), \\ u(\mathbf{o}) &= f_{\mathbf{B}, \mathbf{G}}(\mathbf{o}) + \epsilon_1 \delta_{\mathbb{A}}(\mathbf{o}) + 2\delta_{\mathbb{A}^c}(\mathbf{o}), \end{aligned}$$

where the indicator function $\delta_{\mathbb{A}}(\mathbf{o})$ is defined as $\delta_{\mathbb{A}}(\mathbf{o}) = \begin{cases} 1, & \mathbf{o} \in \mathbb{A} \\ 0, & \mathbf{o} \in \mathbb{A}^c \end{cases}$. Due to (46), we have $f_{\mathbf{B}', \mathbf{G}'} \in [l, u]$ for all $(\mathbf{B}', \mathbf{G}') \in \mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1)$. Also,

$$\begin{aligned} \|u - l\|_{L_2(P)} &= \|2\epsilon_1 \delta_{\mathbb{A}}(\mathbf{o}) + 4\delta_{\mathbb{A}^c}(\mathbf{o})\|_{L_2(P)} = \sqrt{4(\epsilon_1 + \epsilon_2)^2 \mathbb{P}[\mathbf{o} \in \mathbb{A}] + 16\mathbb{P}[\mathbf{o} \in \mathbb{A}^c]} \\ &< 2(\epsilon_1 + \epsilon_2) + 4\sqrt{\mathbb{P}[\mathbf{o} \in \mathbb{A}^c]} \leq 2(\epsilon_1 + \epsilon_2) + 4\sqrt{c_1 c_D D} \frac{\epsilon_1}{\epsilon_2}, \end{aligned} \quad (47)$$

and the last inequality follows because $\mathbb{P}[\mathbf{o} \in \mathbb{A}^c] \leq c_1 c_D D \epsilon_1^2$ according to Lemma 12 with c_1 a universal constant

Finally, the number of brackets to cover \mathcal{F} is equal to the number of such balls $\mathbb{B}((\mathbf{B}, \mathbf{G}), \epsilon_1)$ that cover \mathbb{F} . Utilizing Lemma 11, the covering number for \mathbb{F} is

$$\mathcal{N}(\mathbb{F}, \epsilon_1) \leq \left(1 + \frac{2\sqrt{2}}{\epsilon_1} \right)^{2cD}. \quad (48)$$

Recall the definition that the bracket number $N_{[]}(\epsilon, \mathcal{F}, L_2(P))$ is the minimum number of ϵ -brackets needed to cover \mathcal{F} , where an ϵ -bracket in $L_2(P)$ is a bracket $[l, u]$ with $\|u - l\|_{L_2(P)} \leq \epsilon$. Thus, by letting $\epsilon_2 = \sqrt{\epsilon_1}$ and $2(\epsilon_1 + \sqrt{\epsilon_2}) + 4\sqrt{c_1 c_D D} \sqrt{\epsilon_1} = \epsilon$ and plugging this into (48), we obtain the bracket number

$$N_{[]}(\epsilon, \mathcal{F}, L_2(P)) \leq \left(1 + c_2 \frac{c_D D}{\epsilon^2} \right)^{2cD},$$

where c_2 is a universal constant. Now plug this into Lemma 7 gives

$$\begin{aligned} \frac{1}{\sqrt{M}} \mathbb{E} \left[\sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j=1}^M \left\langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \right\rangle \right| \right] &\lesssim \int_0^1 \sqrt{\left(1 + c_2 \frac{c_D D}{\epsilon^2} \right)^{2cD}} d\epsilon \\ &\lesssim \sqrt{cD} \log(c_D D). \end{aligned}$$

□

We are now ready to prove Lemma 5. Let \mathbf{o}'_k be any points of \mathbb{S}^{D-1} . Since the product of compact spaces is compact, there exist $\mathbf{B}^*, \mathbf{G}^* \in \mathbb{O}(c, D)$ for which the value $\sup_{\mathbf{B}, \mathbf{G} \in \mathbb{S}^{D-1}, \mathbf{B} \perp \mathbf{g}} \left| \sum_{j=1}^M \text{sign}(\mathbf{B}^\top \mathbf{o}_j) \mathbf{g}^\top \mathbf{o}_j \right|$ is achieved. Then, we have

$$\left| \sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j=1}^M \langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \rangle \right| \right| \quad (49)$$

$$- \sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j \neq k} \langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \rangle + \langle \text{sign}(\mathbf{B}^\top \mathbf{o}'_k), \mathbf{G}^\top \mathbf{o}'_k \rangle \right| \quad (50)$$

$$\leq \left| \sum_{j=1}^M \langle \text{sign}(\mathbf{B}^{*\top} \mathbf{o}_j), \mathbf{G}^{*\top} \mathbf{o}_j \rangle \right| - \left| \sum_{j \neq k} \langle \text{sign}(\mathbf{B}^{*\top} \mathbf{o}_j), \mathbf{G}^{*\top} \mathbf{o}_j \rangle + \langle \text{sign}(\mathbf{B}^{*\top} \mathbf{o}'_k), \mathbf{G}^{*\top} \mathbf{o}'_k \rangle \right| \quad (51)$$

$$\leq \left| \langle \text{sign}(\mathbf{B}^{*\top} \mathbf{o}_k), \mathbf{G}^{*\top} \mathbf{o}_k \rangle - \langle \text{sign}(\mathbf{B}^{*\top} \mathbf{o}'_k), \mathbf{G}^{*\top} \mathbf{o}'_k \rangle \right| \leq 2, \quad (52)$$

where the second inequality follows from the reverse triangle inequality. Applying Lemma 8 with $c_k = 2$ and using Lemma 13, we obtain

$$\mathbb{P} \left[\sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j=1}^M \langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \rangle \right| \geq \sqrt{cD} \log(c_D D) \sqrt{M} + \epsilon \right] \leq 2 \exp \left(-\frac{2\epsilon^2}{4M} \right). \quad (53)$$

Finally, set $\epsilon = t\sqrt{M}$ to get

$$\mathbb{P} \left[\sup_{\mathbf{B}, \mathbf{G} \in \mathbb{O}(c, D), \mathbf{G} \perp \mathbf{B}} \left| \sum_{j=1}^M \langle \text{sign}(\mathbf{B}^\top \mathbf{o}_j), \mathbf{G}^\top \mathbf{o}_j \rangle \right| \geq \left(\sqrt{cD} \log(c_D D) + t \right) \sqrt{M} \right] \leq 2 \exp(-t^2/2). \quad (54)$$

References

- [1] A. Edelman, T. Arias, and S. T. Smith, “The geometry of algorithms with orthogonality constraints,” *SIAM Journal of Matrix Analysis Applications*, vol. 20, no. 2, pp. 303–353, 1998.
- [2] G. W. Stewart and J.-g. Sun, *Matrix perturbation theory*. Academic press, 1990.
- [3] N. Higham and P. Papadimitriou, “Matrix procrustes problems,” *Rapport technique, University of Manchester*, 1995.
- [4] R. T. Rockafellar and R. J.-B. Wets, *Variational Analysis*, vol. 317. Springer, 1998.
- [5] O. Ferreira and P. Oliveira, “Subgradient algorithm on riemannian manifolds,” *Journal of Optimization Theory and Applications*, vol. 97, no. 1, pp. 93–104, 1998.
- [6] Z. Zhu, Y. Wang, D. P. Robinson, D. Naiman, R. Vidal, and M. C. Tsakiris, “Dual principal component pursuit: Improved analysis and efficient algorithms,” in *Neural Information Processing Systems*, 2018.

- [7] G. Lerman, M. B. McCoy, J. A. Tropp, and T. Zhang, “Robust computation of linear models by convex relaxation,” *Foundations of Computational Mathematics*, vol. 15, no. 2, pp. 363–410, 2015.
- [8] I. C. Ipsen and C. D. Meyer, “The angle between complementary subspaces,” *The American mathematical monthly*, vol. 102, no. 10, pp. 904–911, 1995.
- [9] Y. Bai, Q. Jiang, and J. Sun, “Subgradient Descent Learns Orthogonal Dictionaries,” in *International Conference on Learning Representations*, 2019.
- [10] Z. Zhu, Y. Wang, D. Robinson, D. Naiman, R. Vidal, and M. Tsakiris, “Dual principal component pursuit: Improved analysis and efficient algorithms,” *arXiv preprint arXiv:1812.09924*, 2018.
- [11] R. Vershynin, “Introduction to the non-asymptotic analysis of random matrices,” *arXiv preprint arXiv:1011.3027*, 2010.
- [12] E. J. Candès, X. Li, and M. Soltanolkotabi, “Phase retrieval via wirtinger flow: Theory and algorithms,” *IEEE Transactions on Information Theory*, vol. 61, no. 4, pp. 1985–2007, 2015.
- [13] A. W. Van der Vaart, *Asymptotic statistics*, vol. 3. Cambridge university press, 1998.
- [14] C. McDiarmid, “On the method of bounded differences,” *London Math. Soc. Lecture Note Ser*, vol. 141, pp. 148–188, 1989.
- [15] A. Maurer, “A vector-contraction inequality for rademacher complexities,” in *International Conference on Algorithmic Learning Theory*, pp. 3–17, Springer, 2016.
- [16] S. Kakade, “Symmetrization and rademacher averages,” *Lecture Notes on Statistical Learning Theory, (Lecture 11)*, 2011.