

# 373 Appendices

374 The appendices are organized as follows. In Appendix A we introduce the basic notations and  
 375 problem reductions that are used throughout the main draft and the appendix. In Appendix B we  
 376 describe and prove the main geometric properties of the optimization landscape for Huber loss. In  
 377 Appendix C, we provide global convergence analysis for the propose Riemannian gradient descent  
 378 methods for optimizing the Huber loss, and the subgradient methods for solving LP rounding. We  
 379 list the basic technical tools and results in Appendix D. All the technical geometric analysis are  
 380 postponed to Appendix E, Appendix F, Appendix G, and Appendix H. Finally, in Appendix I we  
 381 describe the proposed optimization algorithms in full details for all  $\ell^1$ , Huber, and  $\ell^4$  losses.

## 382 A Basic Notations and Problem Reductions

### 383 A.1 Notations

384 Throughout this paper, all vectors/matrices are written in bold font  $a/A$ ; indexed values are written  
 385 as  $a_i, A_{ij}$ . We use  $v_{-i}$  to denote a subvector of  $v$  without the  $i$ -th entry. Zeros or ones vectors are  
 386 defined as  $\mathbf{0}_m$  or  $\mathbf{1}_m$  with  $m$  denoting its length, and  $i$ -th canonical basis vector defined as  $e_i$ . We  
 387 use  $\mathbb{S}^{n-1}$  to denote an  $n$ -dimensional unit sphere in the Euclidean space  $\mathbb{R}^n$ . We use  $z^{(k)}$  to denote  
 388 an optimization variable  $z$  at  $k$ -th iteration. We let  $[m] = \{1, 2, \dots, m\}$ . Let  $F_n \in \mathbb{C}^{n \times n}$  denote a  
 389 unnormalized  $n \times n$  Fourier matrix, with  $\|F_n\| = \sqrt{n}$ , and  $F_n^{-1} = n^{-1}F_n^*$ .

390 We define sign( $\cdot$ ) as

$$\text{sign}(z) = \begin{cases} z/|z|, & z \neq 0 \\ 0, & z = 0 \end{cases}$$

391 **Some basic operators.** We use  $\iota_{n \rightarrow m}$  to denote the zero-padding operator  $\iota_{n \rightarrow m}v = \begin{bmatrix} v \\ \mathbf{0}_{n-m} \end{bmatrix}$ ,  
 392 which zero-pads a length  $n$  vector  $v \in \mathbb{R}^n$  to length  $m$  ( $n \leq m$ ). Correspondingly, its adjoint operator  
 393  $\iota_{n \rightarrow m}^*$  denotes the restriction of a vector to its first  $n$  coordinate (and  $\iota_{n \rightarrow m}^* = \iota_{m \rightarrow n}$ ). Similarly,  
 394 given a subset  $\mathcal{I} \subseteq [m]$  and a vector  $v \in \mathbb{R}^{|\mathcal{I}|}$ , we use  $\iota_{\mathcal{I} \rightarrow m} : \mathbb{R}^{|\mathcal{I}|} \mapsto \mathbb{R}^m$  to denote an operator that  
 395 maps  $v$  to a zero-padded vector whose entries in  $\mathcal{I}$  corresponding to those of  $v$ .

396 We use  $\mathcal{P}_v$  and  $\mathcal{P}_{v^\perp}$  to denote projections onto  $v$  and its orthogonal complement, respectively. We let  
 397  $\mathcal{P}_{\mathbb{S}^{n-1}}$  to be the  $\ell^2$ -normalization operator. To sum up, we have

$$\mathcal{P}_{v^\perp}u = u - \frac{vv^\top}{\|v\|^2}v, \quad \mathcal{P}_v u = \frac{vv^\top}{\|v\|^2}u, \quad \mathcal{P}_{\mathbb{S}^{n-1}}v = \frac{v}{\|v\|}.$$

398 **Circular convolution and circulant matrices.** The convolution operator  $\circledast$  is *circular* with  
 399 modulo- $m$ :  $(a \circledast x)_i = \sum_{j=0}^{m-1} a_j x_{i-j}$ , and we use  $\boxtimes$  to specify the *circular* convolution in 2D. For  
 400 a vector  $v \in \mathbb{R}^m$ , let  $s_\ell[v]$  denote the cyclic shift of  $v$  with length  $\ell$ . Thus, we can introduce the  
 401 circulant matrix  $C_v \in \mathbb{R}^{m \times m}$  generated through  $v \in \mathbb{R}^m$ ,

$$C_v = \begin{bmatrix} v_1 & v_m & \cdots & v_3 & v_2 \\ v_2 & v_1 & v_m & & v_3 \\ \vdots & v_2 & v_1 & \ddots & \vdots \\ v_{m-1} & & \ddots & \ddots & v_m \\ v_m & v_{m-1} & \cdots & v_2 & v_1 \end{bmatrix} = [s_0[v] \ s_1[v] \ \cdots \ s_{m-1}[v]]. \quad (16)$$

402 Now the circulant convolution can also be written in a simpler matrix-vector product form. For  
 403 instance, for any  $u \in \mathbb{R}^m$  and  $v \in \mathbb{R}^n$  ( $n \leq m$ ),

$$u \circledast v = C_u \cdot \iota_{n \rightarrow m}v = C_{\iota_{n \rightarrow m}v} \cdot u.$$

404 In addition, the correlation between  $u$  and  $v$  can be also written in a similar form of convolution  
 405 operator which reverses one vector before convolution. Let  $\tilde{v}$  denote a *cyclic reversal* of  $v \in \mathbb{R}^m$ ,

406 i.e.,  $\check{\mathbf{v}} = [v_1, v_m, v_{m-1}, \dots, v_2]^\top$ , and define two correlation matrices  $\mathbf{C}_\mathbf{v}^* \mathbf{e}_j = s_j[\mathbf{v}]$  and  $\check{\mathbf{C}}_\mathbf{v} \mathbf{e}_j =$   
407  $s_{-j}[\mathbf{v}]$ . The two operators satisfy

$$\mathbf{C}_{\iota_{n \rightarrow m} \mathbf{v}}^* \mathbf{u} = \check{\mathbf{v}} \circledast \mathbf{u}, \quad \check{\mathbf{C}}_{\iota_{n \rightarrow m} \mathbf{v}} \mathbf{u} = \mathbf{v} \circledast \check{\mathbf{u}}.$$

408 **Notation for several distributions.** We use *i.i.d.* to denote *identically* and *independently distributed* random variables:  
409

- 410 • we use  $\mathcal{N}(\mu, \sigma^2)$  to denote Gaussian distribution with mean  $\mu$  and variance  $\sigma^2$ ;
- 411 • we use  $\mathcal{U}(\mathbb{S}^{n-1})$  to denote the uniform distribution over the sphere  $\mathbb{S}^{n-1}$ ;
- 412 • we use  $\mathcal{B}(\theta)$  to denote the Bernoulli distribution with parameter  $\theta$  controlling the nonzero  
413 probability;
- 414 • we use  $\mathcal{BG}(\theta)$  to denote Bernoulli-Gaussian distribution, i.e., if  $u \sim \mathcal{BG}(\theta)$ , then  $u = b \cdot g$   
415 with  $b \sim \mathcal{B}(\theta)$  and  $g \sim \mathcal{N}(0, 1)$ ;
- 416 • we use  $\mathcal{BR}(\theta)$  to denote Bernoulli-Rademacher distribution, i.e., if  $u \sim \mathcal{BR}(\theta)$ , then  
417  $u = b \cdot r$  with  $b \sim \mathcal{B}(\theta)$  and  $r$  follows Rademacher distribution.

## 418 A.2 Problem Reduction

419 In the following sections of the appendices, we study the optimization of

$$\min_{\mathbf{q}} \varphi_h(\mathbf{q}) := \frac{1}{np} \sum_{i=1}^p H_\mu (\mathbf{C}_{\mathbf{y}_i} \mathbf{P} \mathbf{q}), \quad \text{s.t. } \mathbf{q} \in \mathbb{S}^{n-1},$$

420 where

$$\mathbf{P} = \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2}.$$

421 We simplify the problem by a change of variable  $\bar{\mathbf{q}} = \mathbf{Q} \mathbf{q}$ , which rotates the space by the orthogonal  
422 matrix  $\mathbf{Q}$  in (7). Since the rotation  $\mathbf{Q}$  does not change the optimization landscape, by an abuse of  
423 notation of  $\mathbf{q}$  and  $\bar{\mathbf{q}}$ , we can rewrite the problem (5) as

$$\min_{\mathbf{q}} f(\mathbf{q}) := \frac{1}{np} \sum_{i=1}^p H_\mu (\mathbf{C}_{\mathbf{x}_i} \mathbf{R} \mathbf{Q}^{-1} \mathbf{q}), \quad \text{s.t. } \|\mathbf{q}\| = 1, \quad (17)$$

424 where

$$\mathbf{R} = \mathbf{C}_{\mathbf{a}} \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2}, \quad \mathbf{Q} = \mathbf{C}_{\mathbf{a}} (\mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}})^{-1/2},$$

425 and

$$\mathbf{R} \mathbf{Q}^{-1} = \mathbf{C}_{\mathbf{a}} \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2} (\mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}})^{1/2} \mathbf{C}_{\mathbf{a}}^{-1}.$$

426 For the reduction from the original problem to (17), we used the fact that  $\mathbf{C}_{\mathbf{y}_i} \mathbf{P} = \mathbf{C}_{\mathbf{x}_i} \mathbf{R}$  in (7).  
427 Moreover, since  $\mathbf{R} \approx \mathbf{Q}$  is near orthogonal, by assuming  $\mathbf{R} \mathbf{Q}^{-1} = \mathbf{I}$  we can further reduce (17) to

$$\min_{\mathbf{q}} \tilde{f}(\mathbf{q}) = \frac{1}{np} \sum_{i=1}^p H_\mu (\mathbf{C}_{\mathbf{x}_i} \mathbf{q}), \quad \text{s.t. } \|\mathbf{q}\| = 1. \quad (18)$$

428 The objective (10) is simpler and much easier for analysis. By a similar analysis as [22, 32], it can be  
429 shown that asymptotically the landscape is highly symmetric and the standard basis vectors  $\{\pm \mathbf{e}_i\}_{i=1}^n$   
430 are the only global minimizers.

431 For the following sections of the appendices, without loss of generality, we study the optimization  
432 landscape of  $f(\mathbf{q})$  over the sphere, and proving global convergence of vanilla Riemannian gradient  
433 descent methods. We will show that  $\mathbf{R} \mathbf{Q}^{-1} \approx \mathbf{I}$ , so that we can study the landscape of  $f(\mathbf{q})$  via  
434 studying the landscape of  $\tilde{f}(\mathbf{q})$  followed by a perturbation analysis.

435 **B Geometry: Main Results**

436 In this part of appendix, we study the optimization landscape of  $f(\mathbf{q})$  over regions

$$\mathcal{S}_\xi^{i\pm} := \left\{ \mathbf{q} \in \mathbb{S}^{n-1} \mid \frac{|q_i|}{\|\mathbf{q}_{-i}\|_\infty} \geq \sqrt{1+\xi}, q_i \geq 0 \right\}, \quad \xi \in [0, \infty).$$

437 We will show that the function  $f(\mathbf{q})$  over each one of  $\mathcal{S}_\xi^{i\pm}$  has benign first-order geometric structure,  
438 which enables efficient optimization via vanilla Riemannian gradient descent methods.

439 **Proposition B.1 (Regularity condition)** Suppose  $\theta \geq \frac{1}{n}$  and  $\mu \leq c_0 \min \left\{ \theta, \frac{1}{\sqrt{n}} \right\}$ . There exists  
440 some numerical constant  $\gamma \in (0, 1)$ , when the sample complexity

$$p \geq C \max \left\{ n, \frac{\kappa^8}{\theta \mu^2 \sigma_{\min}^2} \log^4 n \right\} \xi^{-2} \theta^{-2} n^4 \log \left( \frac{\theta n}{\mu} \right),$$

441 with probability at least  $1 - n^{-c_1} - c_2 np^{-c_3 n \theta}$  over the randomness of  $\{\mathbf{x}_i\}_{i=1}^p$ , we have

$$\langle \text{grad } f(\mathbf{q}), q_i \mathbf{q} - \mathbf{e}_i \rangle \geq c_4 \theta (1-\theta) q_i \|\mathbf{q} - \mathbf{e}_i\|, \quad \sqrt{1-q_i^2} \in [\mu, \gamma], \quad (19)$$

$$\langle \text{grad } f(\mathbf{q}), q_i \mathbf{q} - \mathbf{e}_i \rangle \geq c_4 \theta (1-\theta) q_i n^{-1} \|\mathbf{q} - \mathbf{e}_i\|, \quad \sqrt{1-q_i^2} \in \left[ \gamma, \sqrt{\frac{n-1}{n}} \right], \quad (20)$$

442 holds for any  $\mathbf{q} \in \mathcal{S}_\xi^{i+}$  and each index  $i \in [n]$ . Here,  $c_0, c_1, c_2, c_3, c_4$ , and  $C$  are positive numerical  
443 constants.

444 **Proof** Without loss of generality, it is enough to consider the case  $i = n$ . For all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ , we have

$$\begin{aligned} & \langle \text{grad } f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \\ &= \langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}) + \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})] + \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle \\ &\geq \langle \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle - |\langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle| \\ &\quad - |\langle \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle|. \end{aligned}$$

445 From Proposition E.1, when  $\theta \geq \frac{1}{n}$  and  $\mu \leq c_0 \min \left\{ \theta, \frac{1}{\sqrt{n}} \right\}$ , we know that in the worst case  
446 scenario,

$$\langle \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle \geq c_1 \theta (1-\theta) \xi n^{-3/2} \|\mathbf{q}_{-n}\|$$

447 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ . On the other hand, by Corollary G.2, when  $p \geq C_1 \theta^{-2} \xi^{-2} n^5 \log \left( \frac{\theta n}{\mu} \right)$ , we  
448 have

$$\begin{aligned} |\langle \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle| &\leq \|\text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})]\| \|q_n \mathbf{q} - \mathbf{e}_n\| \\ &\leq \frac{c_1}{3} \theta (1-\theta) \xi n^{-3/2} \|q_n \mathbf{q} - \mathbf{e}_n\| \end{aligned}$$

449 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$  with probability at least  $1 - np^{-c_2 \theta n} - n \exp(-c_3 n^2)$ . Moreover, from  
450 Proposition H.1, we know that when  $p \geq C \frac{\kappa^8 n^4}{\mu^2 \theta^3 \sigma_{\min}^2 \xi^2} \log^4 n \log \left( \frac{\theta n}{\mu} \right)$

$$\begin{aligned} |\langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle| &\leq \|q_n \mathbf{q} - \mathbf{e}_n\| \cdot \|\text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q})\| \\ &\leq \frac{c_1}{3} \theta (1-\theta) \xi n^{-3/2} \|q_n \mathbf{q} - \mathbf{e}_n\| \end{aligned}$$

451 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$  with probability at least  $1 - c_4 p^{-c_5 n \theta} - n^{-c_6} - n e^{-c_7 \theta n p}$ . By combining all  
452 the bounds above, we obtain the desired result.  $\blacksquare$

453 **Proposition B.2 (Negative curvature on the gradient)** Suppose  $\theta \geq \frac{1}{n}$  and  $\mu \leq \frac{c_0}{\sqrt{n}}$ . For any  
454 index  $i \in [n]$ , when the sample

$$p \geq C \max \left\{ n, \frac{\kappa^8}{\theta \mu^2 \sigma_{\min}^2} \log^4 n \right\} \xi^{-2} \theta^{-2} n^4 \log \left( \frac{\theta n}{\mu} \right),$$

455 with probability at least  $1 - n^{-c_1} - c_2 np^{-c_3 n \theta}$  over the randomness of  $\{\mathbf{x}_i\}_{i=1}^p$ , we have

$$\left\langle \text{grad } f(\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \geq c_4 \frac{\theta(1-\theta)}{n} \frac{\xi}{1+\xi}, \quad (21)$$

456 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{i+}$  and any  $q_j$  such that  $j \neq i$  and  $q_j^2 \geq \frac{1}{3} q_i^2$ . Here,  $c_0, c_1, c_2, c_3, c_4$ , and  $C$  are  
457 positive numerical constants.

458 **Proof** Without loss of generality, it is enough to consider the case  $i = n$ . For all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ , we have

$$\begin{aligned} & \left\langle \text{grad } f(\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \\ &= \left\langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}) + \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})] + \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \\ &\geq \left\langle \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle - \left| \left\langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \right| \\ &\quad - \left| \left\langle \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \right|. \end{aligned}$$

459 From Proposition F.1, when  $\theta \geq \frac{1}{n}$  and  $\mu \leq \frac{c_0}{\sqrt{n}}$ , we know that

$$\left\langle \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \geq \frac{\theta(1-\theta)}{4n} \frac{\xi}{1+\xi}$$

460 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$  and any  $q_j$  such that  $q_j^2 \geq \frac{1}{3} q_i^2$ . On the other hand, by Corollary G.2, when  
461  $p \geq C_1 \theta^{-2} \xi^{-2} n^5 \log \left( \frac{\theta n}{\mu} \right)$ , we have

$$\begin{aligned} \left| \left\langle \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \right| &\leq \left\| \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})] \right\| \cdot \left\| \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\| \\ &\leq \frac{\theta(1-\theta)}{12n} \frac{\xi}{1+\xi} \end{aligned}$$

462 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$  with probability at least  $1 - np^{-c_2 \theta n} - n \exp(-c_3 n^2)$ . For the last inequality,  
463 we used the fact that

$$\left\| \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\| = \sqrt{\frac{1}{q_j^2} + \frac{1}{q_n^2}} \leq 2\sqrt{n}.$$

464 Moreover, from Proposition H.1, we know that when  $p \geq C \frac{\kappa^8 n^4}{\mu^2 \theta^3 \sigma_{\min}^2 \xi^2} \log^4 n \log \left( \frac{\theta n}{\mu} \right)$

$$\begin{aligned} \left| \left\langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \right\rangle \right| &\leq \left\| \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}) \right\| \cdot \left\| \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\| \\ &\leq \frac{\theta(1-\theta)}{12n} \frac{\xi}{1+\xi} \end{aligned}$$

465 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$  with probability at least  $1 - c_4 p^{-c_5 n \theta} - n^{-c_6} - n e^{-c_7 \theta n p}$ . By combining all  
466 the bounds above, we obtain the desired result.  $\blacksquare$

467 **Proposition B.3 (Bounded gradient)** Suppose  $\theta \geq \frac{1}{n}$  and  $\mu \leq \frac{c_0}{\sqrt{n}}$ . For any index  $i \in [n]$ , when  
468 the sample

$$p \geq C \max \left\{ n, \frac{\kappa^8}{\theta \mu^2 \sigma_{\min}^2} \log^4 n \right\} \theta^{-2} n \log \left( \frac{\theta n}{\mu} \right),$$

469 with probability at least  $1 - n^{-c_1} - c_2 np^{-c_3 n \theta}$  over the randomness of  $\{\mathbf{x}_i\}_{i=1}^p$ , we have

$$|\langle \text{grad } f(\mathbf{q}), \mathbf{e}_i \rangle| \leq 2, \quad (22)$$

$$\|\text{grad } f(\mathbf{q})\| \leq 2\sqrt{\theta n}. \quad (23)$$

470 holds for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  and any index  $i \in [n]$ . Here,  $c_0, c_1, c_2, c_3$  and  $C$  are positive numerical  
471 constants.

472 **Proof** For any index  $i \in [n]$ , we have

$$\begin{aligned} \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } f(\mathbf{q}), \mathbf{e}_i \rangle| &\leq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| + \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| \\ &\leq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| + \|\text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q})\|. \end{aligned}$$

473 By Corollary G.3, when  $p \geq C_1 n \log \left( \frac{\theta n}{\mu} \right)$ , we have

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| \leq \frac{3}{2}$$

474 holds for any index  $i \in [n]$  with probability at least  $1 - np^{-c_1 \theta n} - n \exp(-c_2 p)$ . On the other hand,  
475 Proposition H.1 implies that, when  $p \geq C_2 \frac{\kappa^8 n}{\mu^2 \theta \sigma_{\min}^2} \log^4 n \log \left( \frac{\theta n}{\mu} \right)$ , we have

$$\|\text{grad } f(\mathbf{q}) - \text{grad } \tilde{f}(\mathbf{q})\| \leq \frac{1}{2},$$

476 holds with probability at least  $1 - c_3 p^{-c_4 n \theta} - n^{-c_5} - n e^{-c_6 \theta n p}$ . Combining the bounds above gives  
477 (22). The bound (23) can be proved in a similar fashion.  $\blacksquare$

## 478 C Convergence Analysis

479 In this section, we show linear convergence of *vanilla* gradient descent to target solutions. Firstly,  
480 for Huber loss, we prove that the gradient descent method converges to an approximate solution in  
481 polynomial steps. Second, we show linear convergence of subgradient method to the target solution,  
482 which solves phase-2 LP rounding problem.

483 Our analysis is based on the geometric properties of the optimization landscape showed in  
484 Appendix B. Namely, our following proofs are based on the results in Proposition B.1, Proposition  
485 B.2, and Proposition B.3 (i.e., the equations (19), (20), (22), and (23)) holding for the rest of this  
486 section.

### 487 C.1 Proof of linear convergence for Algorithm 1

488 First, assuming the geometric properties in Appendix B hold, we show that starting from a random  
489 initialization, the gradient descent method optimizing

$$\min_{\mathbf{q}} f(\mathbf{q}) = \frac{1}{np} \sum_{i=1}^p H_\mu (\mathbf{C}_{\mathbf{x}_i} \mathbf{R} \mathbf{Q}^{-1} \mathbf{q}), \quad \text{s.t. } \mathbf{q} \in \mathbb{S}^{n-1} \quad (24)$$

490 recovers an approximate solution in polynomial steps.

491 **Theorem C.1 (Linear convergence of Algorithm 1)** *Given an initialization  $\mathbf{q}^{(0)} \sim \mathcal{U}(\mathbb{S}^{n-1})$  uniform random drawn from the sphere, choose a stepsize*

$$\tau = c \min \left\{ \frac{1}{n^{5/2}}, \frac{\mu}{n} \right\},$$

493 *then the vanilla gradient descent method for (5) produces a solution*

$$\|\mathbf{q}^{(k)} - \mathbf{e}_i\| \leq 2\mu$$

494 *for some  $i \in [n]$ , whenever*

$$k \geq K := \frac{C}{\theta} \max \left\{ n^4, \frac{n^{5/2}}{\mu} \right\} \log \left( \frac{1}{\mu} \right).$$

495 **Proof** [Proof of Theorem C.1]

496 **Initialization and iterate stays within the region.** First, from Lemma C.3, we know that when  
497  $\xi = \frac{1}{5 \log n}$ , with probability at least 1/2, our random initialization  $\mathbf{q}^{(0)}$  falls into one of the sets  
498  $\{\mathcal{S}_\xi^{1+}, \mathcal{S}_\xi^{1-}, \dots, \mathcal{S}_\xi^{n+}, \mathcal{S}_\xi^{n-}\}$ . Without loss of generality, we assume that  $\mathbf{q}^{(0)} \in \mathcal{S}_\xi^{n+}$ .

499 Once  $\mathbf{q}^{(0)}$  initialized within the region  $\mathcal{S}_\xi^{n+}$ , from Lemma C.4, whenever the stepsize  $\tau \leq c_0/\sqrt{n}$ ,  
500 we know that our gradient descent stays within the region  $\mathcal{S}_\xi^{n+}$  when the stepsize  $\tau \leq c_1/\sqrt{n}$  for  
501 some  $c_1 > 0$ . Based on this, to complete the proof, we now proceed by proving the following results.

502 **Linear convergence until reaching  $\|\mathbf{q} - \mathbf{e}_n\| \leq \mu$ .** From Proposition B.1, there exists some  
503 numerical constant  $\gamma \in (\mu, 1)$ , such that the regularity condition

$$\langle \text{grad } f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \geq \underbrace{c_2 \theta (1 - \theta) n^{-3/2}}_{\alpha_1} \cdot \|\mathbf{q} - \mathbf{e}_n\|, \quad \sqrt{1 - q_n^2} \in \left[ \gamma, \sqrt{\frac{n-1}{n}} \right], \quad (25)$$

$$\langle \text{grad } f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \geq \underbrace{c'_2 \theta (1 - \theta)}_{\alpha_2} \cdot \|\mathbf{q} - \mathbf{e}_n\|, \quad \sqrt{1 - q_n^2} \in [\mu, \gamma], \quad (26)$$

504 holds w.h.p. for all  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ . As  $\alpha_2 \geq \alpha_1$ , the regularity condition holds for all  $\mathbf{q}$  with  $\alpha = \alpha_1$ .

505 Select a stepsize  $\tau$  such that  $\tau \leq \gamma \frac{\alpha_1}{2\sqrt{2\theta n}}$ . By Lemma C.5 and the regularity condition (25), we have

$$\|\mathbf{q}^{(k)} - \mathbf{e}_n\|^2 - \frac{\gamma^2}{2} \leq (1 - \tau \alpha_1)^k \left[ \|\mathbf{q}^{(0)} - \mathbf{e}_n\|^2 - \frac{\gamma^2}{2} \right] \leq 2(1 - \tau \alpha_1)^k,$$

506 where the last inequality utilizes the fact that  $\|\mathbf{q}^{(0)} - \mathbf{e}_n\|^2 \leq 2$ . This further implies that

$$1 - q_n^2 \leq \|\mathbf{q}^{(k)} - \mathbf{e}_n\|^2 \leq \frac{\gamma^2}{2} + 2(1 - \tau \alpha_1)^k \leq \gamma^2,$$

507 when

$$2(1 - \tau \alpha_1)^k \leq \frac{\gamma^2}{2} \implies k \geq K_1 := \frac{\log(\gamma^2/4)}{\log(1 - \tau \alpha_1)}.$$

508 This implies that  $\sqrt{1 - q_n^2} \leq \gamma$  for  $\forall k \geq K_1$ . Thus, from (26), we know that the regularity condition  
509 holds with  $\alpha = \alpha_2$ . Choose stepsize  $\tau \leq \frac{\mu \alpha_2}{2\sqrt{2\theta n}}$ , apply Lemma C.5 again with  $\alpha = \alpha_2$ , for all  $k \geq 1$ ,  
510 we have

$$\|\mathbf{q}^{(K_1+k)} - \mathbf{e}_n\|^2 - \frac{\mu^2}{2} \leq (1 - \tau \alpha_2)^k \left( \|\mathbf{q}^{(0)} - \mathbf{e}_n\|^2 - \frac{\mu^2}{2} \right) \leq (\gamma^2 - \mu^2)(1 - \tau \alpha_2)^k.$$

511 This further implies that

$$\|\mathbf{q}^{(K_1+k)} - \mathbf{e}_n\|^2 \leq \frac{\mu^2}{2} + \left( \gamma^2 - \frac{\mu^2}{2} \right) (1 - \tau \alpha_2)^k \leq \mu^2$$

512 whenever

$$\left(\gamma^2 - \frac{\mu^2}{2}\right)(1 - \tau\alpha_2)^k \leq \frac{\mu^2}{2} \implies k \geq K_2 := \frac{\log(\mu^2/(2\gamma^2 - \mu^2))}{\log(1 - \tau\alpha_2)}.$$

513 Therefore, combining the results above, by using the fact that  $\alpha_1 = c_2\theta(1 - \theta)n^{-3/2}$  and  $\alpha_2 =$   
514  $c'_2\theta(1 - \theta)$ , we have  $\|\mathbf{q}^{(k)} - \mathbf{e}_n\| \leq \mu$  whenever

$$\tau \leq \min\left\{\frac{\gamma\alpha_1}{2\sqrt{2\theta}n}, \frac{\mu\alpha_2}{2\sqrt{2\theta}n}\right\} = C \min\left\{\frac{1}{n^{5/2}}, \frac{\mu}{n}\right\}$$

515 and  $k \geq K := K_1 + K_2$  with

$$\begin{aligned} K &= \frac{\log(4/\gamma^2)}{\log((1 - \tau\alpha_1)^{-1})} + \frac{\log((2\gamma^2 - \mu^2)/\mu^2)}{\log((1 - \tau\alpha_2)^{-1})} \\ &\leq \frac{c_3}{\tau\alpha_1} + \frac{c_4}{\tau\alpha_2} \log\left(\frac{1}{\mu}\right) \leq \frac{c_5}{\theta} \max\left\{n^4, \frac{n^{5/2}}{\mu}\right\} \log\left(\frac{1}{\mu}\right), \end{aligned}$$

516 where we used the fact that  $\log^{-1}((1 - x)^{-1}) \leq 2/x$  for small  $x$ .

517 **No jump away from an approximate solution  $\mathbf{e}_n$ .** Finally, we show that once our iterate reaches  
518 the region

$$\mathcal{S} := \{\mathbf{q} \in \mathbb{S}^{n-1} \mid \|\mathbf{q} - \mathbf{e}_n\| \leq 2\mu\},$$

519 it will stay within the region  $\mathcal{S}$ , such that our final iterates will always stay close to an approximate  
520 solution  $\mathbf{e}_n$ . Towards this end, suppose  $\mathbf{q}^{(k)} \in \mathcal{S}$ . Therefore two possibilities: (i)  $\mu \leq \|\mathbf{q}^{(k)} - \mathbf{e}_n\| \leq$   
521  $2\mu$  (ii)  $\|\mathbf{q}^{(k)} - \mathbf{e}_n\| \leq \mu$ . If the case (i) holds, then our argument above implies that  $\|\mathbf{q}^{(k+1)} - \mathbf{e}_n\| \leq$   
522  $\|\mathbf{q}^{(k)} - \mathbf{e}_n\| \leq 2\mu$ . Otherwise  $\|\mathbf{q}^{(k)} - \mathbf{e}_n\| \leq \mu$ , for which we have

$$\begin{aligned} \|\mathbf{q}^{(k+1)} - \mathbf{e}_n\| &\leq \|\mathbf{q}^{(k)} - \tau \operatorname{grad} f(\mathbf{q}) - \mathbf{e}_n\| \\ &\leq \|\mathbf{q}^{(k)} - \mathbf{e}_n\| + \tau \|\operatorname{grad} f(\mathbf{q})\| \leq \mu + 2\tau\sqrt{\theta n} \leq 2\mu, \end{aligned}$$

523 where we used the fact that  $\tau \leq \frac{\mu}{\sqrt{\theta n}}$ . Thus, by induction, we have  $\mathbf{q}^{(k')} \in \mathcal{S}$  for all future iterates  
524  $k' = k + 1, k + 2, \dots$ . This completes the proof. ■

525 **Lemma C.2** For any  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ , we have

$$1 - q_n^2 \leq \|\mathbf{q} - \mathbf{e}_n\|^2 \leq 2(1 - q_n^2) \leq 2.$$

526 **Proof** We have

$$1 - q_n^2 \leq \|\mathbf{q} - \mathbf{e}_n\|^2 = \|\mathbf{q}_{-n}\|^2 + (1 - q_n)^2 \|\mathbf{e}_n\|^2 = 2(1 - q_n) = 2 \frac{1 - q_n^2}{1 + q_n^2} \leq 2(1 - q_n^2)$$

527 as desired. ■

528 **Lemma C.3 (Random initialization falls into good region)** Let  $\mathbf{q}^{(0)} \sim \mathcal{U}(\mathbb{S}^{n-1})$  be uniformly  
529 random generated from the unit sphere  $\mathbb{S}^{n-1}$ . When  $\xi = \frac{1}{5\log n}$ , then with probability at least  
530  $1/2$ ,  $\mathbf{q}^{(0)}$  belongs to one of the  $2n$  sets  $\{\mathcal{S}_\xi^{1+}, \mathcal{S}_\xi^{1-}, \dots, \mathcal{S}_\xi^{n+}, \mathcal{S}_\xi^{n-}\}$ . The set  $\mathbf{q}^{(0)}$  belongs to is  
531 uniformly at random.

532 **Proof** We refer the readers to Lemma 3.9 of [23] and Theorem 1 of [22] for detailed proofs. ■

533 **Lemma C.4 (Stay within the region  $\mathcal{S}_\xi^{n+}$ )** Suppose  $\mathbf{q}^{(0)} \in \mathcal{S}_\xi^{n+}$  with  $\xi \leq 1$ . There exists some  
534 constant  $c > 0$ , such that when the stepsize satisfies  $\tau \leq \frac{c}{\sqrt{n}}$ , our Riemannian gradient iterate  
535  $\mathbf{q}^{(k)} = \mathcal{P}_{\mathbb{S}^{n-1}}(\mathbf{q}^{(k-1)} - \tau \cdot \operatorname{grad} f(\mathbf{q}^{(k-1)}))$  satisfies  $\mathbf{q}^{(k)} \in \mathcal{S}_\xi^{n+}$  for all  $k \geq 1$ .

536 **Proof** We prove this by induction. For any  $k \geq 1$ , suppose  $\mathbf{q}^{(k)} \in \mathcal{S}_\xi^{n+}$ . For convenience, let  
537  $\mathbf{g}^{(k)} = \text{grad } f(\mathbf{q}^{(k)})$ . Then, for any  $j \neq k$ , we have

$$\left( \frac{q_n^{(k+1)}}{q_j^{(k+1)}} \right)^2 = \left( \frac{q_n^{(k)} - \tau g_n^{(k)}}{q_j^{(k)} - \tau g_j^{(k)}} \right)^2.$$

538 We proceed by considering the following two cases.

539 **Case (i):**  $|q_n^{(k)}/q_j^{(k)}| \geq \sqrt{3}$ . In this case, we have

$$\left( \frac{q_n^{(k+1)}}{q_j^{(k+1)}} \right)^2 = \left( \frac{q_n^{(k)} - \tau g_n^{(k)}}{q_j^{(k)} - \tau g_j^{(k)}} \right)^2 \geq \left( \frac{1 - \tau \cdot g_n^{(k)}/q_n^{(k)}}{q_j^{(k)}/q_n^{(k)} - \tau g_j^{(k)}/q_n^{(k)}} \right)^2 \geq \left( \frac{1 - 2\tau\sqrt{n}}{1/\sqrt{3} + 2\tau\sqrt{n}} \right)^2 \geq 2,$$

540 where the second inequality utilizes (22) and the fact  $q_n^{(k)} \geq \frac{1}{\sqrt{n}}$ , and the last inequality follows when

$$541 \tau \leq \frac{\sqrt{3}-\sqrt{2}}{2(\sqrt{6}+\sqrt{3})} \frac{1}{\sqrt{n}}.$$

542 **Case (ii):**  $|q_n^{(k)}/q_j^{(k)}| \leq \sqrt{3}$ . Proposition B.1 and Proposition B.2 implies that

$$\frac{g_j^{(k)}}{q_j^{(k)}} \geq 0, \quad \frac{g_j^{(k)}}{q_j^{(k)}} - \frac{g_n^{(k)}}{q_n^{(k)}} \geq 0. \quad (27)$$

543 By noting that  $|q_j^{(k)}| \geq |q_n^{(k)}|/\sqrt{3} \geq 1/\sqrt{3n}$  and  $|g_j^{(k)}| \leq 2$ , we have

$$\tau \leq \frac{1}{2\sqrt{3n}} \leq \frac{g_j^{(k)}}{q_j^{(k)}} \implies \tau \cdot \frac{g_j^{(k)}}{q_j^{(k)}} \leq 1. \quad (28)$$

544 Thus, we have

$$\begin{aligned} \left( \frac{q_n^{(k+1)}}{q_j^{(k+1)}} \right)^2 &= \left( \frac{q_n^{(k)}}{q_j^{(k)}} \right)^2 \left( 1 + \tau \cdot \frac{g_j^{(k)}/q_j^{(k)} - g_n^{(k)}/q_n^{(k)}}{1 - \tau g_j^{(k)}/q_j^{(k)}} \right)^2 \\ &\geq \left( \frac{q_n^{(k)}}{q_j^{(k)}} \right)^2 \left( 1 + \tau \cdot \left( \frac{g_j^{(k)}}{q_j^{(k)}} - \frac{g_n^{(k)}}{q_n^{(k)}} \right) \right)^2 \geq \left( \frac{q_n^{(k)}}{q_j^{(k)}} \right)^2 \left( 1 + \tau \cdot \frac{\theta(1-\theta)}{4n} \frac{\xi}{1+\xi} \right)^2. \end{aligned}$$

545 The first inequality follows from (27) and (28), and the second inequality directly follows from  
546 Proposition B.2. Therefore, when  $\xi \leq 1$ , this implies that  $\mathbf{q}^{(k+1)} \in \mathcal{S}_\xi^{n+}$ . By induction, this holds for  
547 all  $k \geq 1$ .  $\blacksquare$

548 In the following, we show that the iterates get closer to  $\mathbf{e}_n$ .

549 **Lemma C.5 (Iterate contraction)** For any  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ , assuming the following regularity condition

$$\langle \text{grad } f(\mathbf{q}), q_i \mathbf{q} - \mathbf{e}_n \rangle \geq \alpha \|\mathbf{q} - \mathbf{e}_n\| \quad (29)$$

550 holds for a parameter  $\alpha > 0$ . Then if  $\mathbf{q}^{(k)} \in \mathcal{S}_\xi^{n+}$  and the stepsize  $\tau \leq c \frac{\alpha}{\theta n}$ , the iterate  $\mathbf{q}^{(k+1)} =$   
551  $\mathcal{P}_{\mathbb{S}^{n-1}}(\mathbf{q} - \tau \cdot \text{grad } f(\mathbf{q}))$  satisfies

$$\|\mathbf{q}^{(k+1)} - \mathbf{e}_n\|^2 - \left( \frac{2\tau\theta n}{\alpha} \right)^2 \leq (1 - \tau\alpha) \left[ \|\mathbf{q}^{(k)} - \mathbf{e}_n\|^2 - \left( \frac{2\tau\theta n}{\alpha} \right)^2 \right].$$

552 **Proof** First, note that

$$\begin{aligned}
\|\mathbf{q}^{(k+1)} - \mathbf{e}_n\|^2 &= \left\| \mathcal{P}_{\mathbb{S}^{n-1}} \left( \mathbf{q}^{(k)} - \tau \cdot \text{grad } f(\mathbf{q}^{(k)}) \right) - \mathcal{P}_{\mathbb{S}^{n-1}}(\mathbf{e}_n) \right\|^2 \\
&\leq \left\| \mathbf{q}^{(k)} - \tau \cdot \text{grad } f(\mathbf{q}^{(k)}) - \mathbf{e}_n \right\|^2 \\
&= \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\|^2 - 2\tau \cdot \langle \text{grad } f(\mathbf{q}^{(k)}), \mathbf{q}^{(k)} - \mathbf{e}_n \rangle + \tau^2 \left\| \text{grad } f(\mathbf{q}^{(k)}) \right\|^2 \\
&\leq \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\|^2 - 2\tau\alpha \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\| + 4\tau^2\theta n,
\end{aligned}$$

553 where the first inequality utilizes the fact that  $\mathcal{P}_{\mathbb{S}^{n-1}}(\cdot)$  is 1-Lipschitz continuous, and the last line  
554 follows from (29) and (23) in Proposition B.3. We now subtract both sides by  $(\frac{2\tau\theta n}{\alpha})^2$ ,

$$\begin{aligned}
\left\| \mathbf{q}^{(k+1)} - \mathbf{e}_n \right\|^2 - \left( \frac{2\tau\theta n}{\alpha} \right)^2 &\leq \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\|^2 - \left( \frac{2\tau\theta n}{\alpha} \right)^2 - 2\tau\alpha \left( \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\| - \frac{2\tau\theta n}{\alpha} \right) \\
&= \left[ 1 - 2\tau\alpha \left( \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\| + \frac{2\tau\theta n}{\alpha} \right)^{-1} \right] \left[ \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\|^2 - \left( \frac{2\tau\theta n}{\alpha} \right)^2 \right] \\
&\leq (1 - \tau\alpha) \left[ \left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\|^2 - \left( \frac{2\tau\theta n}{\alpha} \right)^2 \right],
\end{aligned}$$

555 where the last inequality follows because

$$\left\| \mathbf{q}^{(k)} - \mathbf{e}_n \right\|^2 \leq 2, \quad \tau \leq \left( 1 - \frac{1}{\sqrt{2}} \right) \frac{\alpha}{\theta n},$$

556 such that

$$\left\| \mathbf{q} - \mathbf{e}_n \right\| + \frac{2\tau\theta n}{\alpha} \leq 2.$$

557 This completes the proof. ■

## 558 C.2 Exact solution via LP rounding

559 To obtain exact solutions, we use the approximate solution  $\mathbf{q}_*$  from phase-1 gradient descent method  
560 as a warm start  $\mathbf{r} = \mathbf{q}_*$ , and consider solving a *convex* phase-2 LP rounding problem,

$$\min_{\mathbf{q}} \zeta(\mathbf{q}) := \frac{1}{np} \sum_{i=1}^p \|C_{\mathbf{x}_i} \mathbf{RQ}^{-1} \mathbf{q}\|_1, \quad \text{s.t.} \quad \langle \mathbf{r}, \mathbf{q} \rangle = 1.$$

561 In the following, we show the function is sharp around the target solution, so that projected subgradient  
562 descent methods converge linearly to the truth.

### 563 C.2.1 Sharpness of the objective function.

564 **Proposition C.6** Suppose  $\theta \in (\frac{1}{n}, \frac{1}{3})$  and  $\mathbf{r}$  satisfies

$$\frac{\|\mathbf{r}_{-n}\|}{r_n} \leq \frac{1}{20}. \tag{30}$$

565 Whenever  $p \geq C_{\theta\sigma_{\min}^2(C_\alpha)} \frac{\kappa^8}{\log^3 n}$ , with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ , the function  $\zeta(\mathbf{q})$   
566 is sharp in a sense that

$$\zeta(\mathbf{q}) - \zeta \left( (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right) \geq \frac{1}{50} \sqrt{\frac{2}{\pi}} \theta \left\| \mathbf{q} - (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| \tag{31}$$

567 for any feasible  $\mathbf{q}$  with  $\langle \mathbf{r}, \mathbf{q} \rangle = 1$ . Here,  $\tilde{\mathbf{r}} = (\mathbf{RQ}^{-1})^{-\top} \mathbf{r}$ .

568 **Proof** Let us denote  $\tilde{\mathbf{q}} = \mathbf{R}\mathbf{Q}^{-1}\mathbf{q}$ . Then we can rewrite our original problem as

$$\min_{\tilde{\mathbf{q}}} \tilde{\zeta}(\tilde{\mathbf{q}}) = \frac{1}{np} \sum_{i=1}^p \|\mathbf{C}_{\mathbf{x}_i} \tilde{\mathbf{q}}\|_1 \quad \text{s.t.} \quad \langle \tilde{\mathbf{r}}, \tilde{\mathbf{q}} \rangle = 1,$$

569 which is reduced to the orthogonal problem in (32) of Lemma C.7. To utilize the result in Lemma C.7,  
570 we first prove that  $\tilde{\mathbf{r}}$  satisfies (33) if  $\mathbf{r}$  satisfies (30). Towards that end, note that

$$\tilde{\mathbf{r}} = (\mathbf{R}\mathbf{Q}^{-1})^{-\top} \mathbf{r} = \mathbf{r} + ((\mathbf{R}\mathbf{Q}^{-1})^{-\top} - \mathbf{I}) \mathbf{r}.$$

571 By Lemma H.4, we know that, for any  $\delta \in (0, 1)$ , whenever  $p \geq C \frac{\kappa^8}{\theta \delta^2 \sigma_{\min}^2(\mathbf{C}_a)} \log^3 n$ ,

$$\|((\mathbf{R}\mathbf{Q}^{-1})^{-\top} - \mathbf{I}) \mathbf{r}\| \leq \|(\mathbf{R}\mathbf{Q}^{-1})^{-1} - \mathbf{I}\| \|\mathbf{r}\| \leq 2\delta \|\mathbf{r}\|$$

572 holds with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ . This further implies that

$$\tilde{r}_n \geq r_n - 2\delta \|\mathbf{r}\|, \quad \|\tilde{\mathbf{r}}_{-n}\| \leq \|\mathbf{r}_{-n}\| + 2\delta \|\mathbf{r}\|.$$

573 Therefore, by choose  $\delta$  sufficiently small, we have

$$\frac{\|\tilde{\mathbf{r}}_{-n}\|}{\tilde{r}_n} \leq \frac{\|\mathbf{r}_{-n}\| + 2\delta \|\mathbf{r}\|}{r_n - 2\delta \|\mathbf{r}\|} = \frac{\|\mathbf{r}_{-n}\| / r_n + 2\delta \sqrt{1 + (\|\mathbf{r}_{-n}\| / r_n)^2}}{1 - 2\delta \sqrt{1 + (\|\mathbf{r}_{-n}\| / r_n)^2}} \leq \frac{1}{10},$$

574 where the last inequality follows from (30). Therefore, by Lemma C.7, we obtain

$$\begin{aligned} \zeta(\mathbf{q}) - \zeta\left((\mathbf{R}\mathbf{Q}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n}\right) &= \tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{\tilde{r}_n}\right) \\ &\geq \frac{1}{25} \sqrt{\frac{2}{\pi}} \theta \left\| \tilde{\mathbf{q}} - \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| \\ &= \frac{1}{25} \sqrt{\frac{2}{\pi}} \theta \left\| (\mathbf{R}\mathbf{Q}^{-1}) \cdot \left( \mathbf{q} - (\mathbf{R}\mathbf{Q}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right) \right\| \\ &\geq \frac{1}{25} \sqrt{\frac{2}{\pi}} \theta \cdot \sigma_{\min}(\mathbf{R}\mathbf{Q}^{-1}) \cdot \left\| \mathbf{q} - (\mathbf{R}\mathbf{Q}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| \end{aligned}$$

575 By Lemma H.4, we know that  $\|(\mathbf{R}\mathbf{Q}^{-1})^{-1}\| \leq 1 + 2\delta$ , so that

$$\sigma_{\min}(\mathbf{R}\mathbf{Q}^{-1}) = \|(\mathbf{R}\mathbf{Q}^{-1})^{-1}\|^{-1} \geq \frac{1}{1 + 2\delta}.$$

576 Thus, this further implies that

$$\zeta(\mathbf{q}) - \zeta\left((\mathbf{R}\mathbf{Q}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n}\right) \geq \frac{1}{25} \sqrt{\frac{2}{\pi}} \frac{\theta}{1 + 2\delta} \cdot \left\| \mathbf{q} - (\mathbf{R}\mathbf{Q}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\|,$$

577 as desired. ■

578 **Lemma C.7 (Sharpness for the orthogonal case)** Consider the following problem

$$\min_{\mathbf{q}} \tilde{\zeta}(\mathbf{q}) := \frac{1}{np} \sum_{i=1}^p \|\mathbf{C}_{\mathbf{x}_i} \mathbf{q}\|_1 \quad \text{s.t.} \quad \langle \mathbf{r}, \mathbf{q} \rangle = 1, \tag{32}$$

579 with  $\mathbf{r} \in \mathbb{S}^{n-1}$  satisfying

$$\frac{\|\mathbf{r}_{-n}\|}{r_n} \leq \frac{1}{10}, \quad r_n > 0. \tag{33}$$

580 Whenever  $p \geq \frac{C}{\theta^2} n \log\left(\frac{n}{\theta}\right)$ , with probability at least  $1 - c_1 np^{-6} - c_2 n e^{-c_3 \theta^2 p}$ , the function  $\tilde{\zeta}(\mathbf{q})$  is  
581 sharp in a sense that

$$\tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{r_n}\right) \geq \frac{1}{25} \sqrt{\frac{2}{\pi}} \theta \left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\|$$

582 for any feasible  $\mathbf{q}$  with  $\langle \mathbf{r}, \mathbf{q} \rangle = 1$ .

583 **Proof** Observing that  $\langle \mathbf{r}, \mathbf{q} \rangle = \mathbf{r}_{-n}^\top \mathbf{q}_{-n} + r_n q_n = 1$ , we have

$$\|\mathbf{r}_{-n}\| \|\mathbf{q}_{-n}\| \geq \mathbf{r}_{-n}^\top \mathbf{q}_{-n} = r_n \left( \frac{1}{r_n} - q_n \right) \geq r_n \left( \frac{1}{r_n} - |q_n| \right).$$

584 This further implies that

$$\frac{1}{r_n} - |q_n| \leq \frac{\|\mathbf{r}_{-n}\|}{r_n} \|\mathbf{q}_{-n}\|. \quad (34)$$

585 Second, we have

$$\left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\| = \sqrt{\left( \frac{1}{r_n} - q_n \right)^2 + \|\mathbf{q}_{-n}\|^2} \leq \sqrt{1 + \left( \frac{\|\mathbf{r}_{-n}\|}{r_n} \right)^2} \|\mathbf{q}_{-n}\|,$$

586 which implies that

$$\left( 1 + \left( \frac{\|\mathbf{r}_{-n}\|}{r_n} \right)^2 \right)^{-1/2} \left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\| \leq \|\mathbf{q}_{-n}\|. \quad (35)$$

587 We now proceed by considering the following two cases.

588 **Case i:**  $|q_n| \geq \frac{1}{r_n}$ . In this case, we have

$$\begin{aligned} \tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{r_n}\right) &\geq \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\| \geq \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \left( 1 + \left( \frac{\|\mathbf{r}_{-n}\|}{r_n} \right)^2 \right)^{-1/2} \left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\| \\ &\geq \frac{5}{33} \sqrt{\frac{2}{\pi}} \theta \left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\|, \end{aligned}$$

589 where the first inequality follows by (36), the second inequality follows by (35), and the last inequality  
590 follows because  $\frac{\|\mathbf{r}_{-n}\|}{r_n} \leq \frac{1}{10}$ .

591 **Case ii:**  $|q_n| \leq \frac{1}{r_n}$ . In this case, we have

$$\begin{aligned} \tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{r_n}\right) &\geq \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\| - \frac{5}{4} \sqrt{\frac{2}{\pi}} \theta \left( \frac{1}{r_n} - |q_n| \right) \\ &\geq \theta \left( \frac{1}{6} \sqrt{\frac{2}{\pi}} - \frac{5}{4} \sqrt{\frac{2}{\pi}} \frac{\|\mathbf{r}_{-n}\|}{r_n} \right) \|\mathbf{q}_{-n}\| \\ &\geq \theta \left( \frac{1}{6} \sqrt{\frac{2}{\pi}} - \frac{5}{4} \sqrt{\frac{2}{\pi}} \frac{\|\mathbf{r}_{-n}\|}{r_n} \right) \left( 1 + \left( \frac{\|\mathbf{r}_{-n}\|}{r_n} \right)^2 \right)^{-1/2} \left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\| \\ &\geq \frac{\theta}{25} \sqrt{\frac{2}{\pi}} \left\| \mathbf{q} - \frac{\mathbf{e}_n}{r_n} \right\|, \end{aligned}$$

592 where the first inequality follows by (36), the second inequality follows from (34), the third inequality  
593 follows from (35), and the last one follows because  $\frac{\|\mathbf{r}_{-n}\|}{r_n} \leq \frac{1}{10}$ .

594 Combining the results in both cases, we obtain the desired result. ■

595 **Lemma C.8** Suppose  $\theta \in (\frac{1}{n}, \frac{1}{3})$ . Whenever  $p \geq \frac{C}{\theta^2} n \log(\frac{n}{\theta})$ , we have

$$\tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{r_n}\right) \geq \begin{cases} \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\|, & \text{if } |q_n| - \frac{1}{r_n} \geq 0, \\ \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\| - \frac{5}{4} \sqrt{\frac{2}{\pi}} \theta \left( \frac{1}{r_n} - |q_n| \right), & \text{if } |q_n| - \frac{1}{r_n} < 0, \end{cases} \quad (36)$$

596 holds with probability at least  $1 - c_1 np^{-6} - c_2 ne^{-c_3 \theta^2 p}$ .

597 **Proof** For each  $j \in [n]$ , let us define an index set  $\mathcal{I}_j := \{i \in [p] : (s_j [\tilde{\mathbf{x}}_i])_n \neq 0\}$ , and let us define  
598 events

$$\mathcal{E} := \bigcap_{j=0}^{n-1} \mathcal{E}_j, \quad \mathcal{E}_j := \left\{ |\mathcal{I}_j| \leq \frac{9}{8} \theta p \right\}, \quad (0 \leq j \leq n-1).$$

599 By Hoeffding's inequality and a union bound, we know that

$$\mathbb{P}(\mathcal{E}^c) \leq \sum_{j=0}^{n-1} \mathbb{P}(\mathcal{E}_j^c) \leq n \exp(-p\theta^2/2).$$

600 Based on this, we have

$$\begin{aligned} & \tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{r_n}\right) \\ &= \frac{1}{np} \sum_{i=1}^p \|\mathbf{C}_{\mathbf{x}_i} \mathbf{q}\|_1 - \frac{1}{np} \frac{1}{r_n} \sum_{i=1}^p \|\mathbf{x}_i\|_1 \\ &= \frac{1}{np} \sum_{i=1}^p \sum_{j=0}^{n-1} |\langle s_j [\tilde{\mathbf{x}}_i], \mathbf{q} \rangle| - \frac{1}{np} \frac{1}{r_n} \sum_{i=1}^p \|\mathbf{x}_i\|_1 \\ &\geq \frac{1}{np} \left( |q_n| - \frac{1}{r_n} \right) \sum_{i=1}^p \|\mathbf{x}_i\|_1 + \frac{1}{np} \sum_{j=0}^{n-1} \left( \sum_{i \in \mathcal{I}_j^c} |\langle (s_j [\tilde{\mathbf{x}}_i])_{-n}, \mathbf{q}_{-n} \rangle| - \sum_{i \in \mathcal{I}_j} |\langle (s_j [\tilde{\mathbf{x}}_i])_{-n}, \mathbf{q}_{-n} \rangle| \right) \\ &= \frac{1}{np} \left( |q_n| - \frac{1}{r_n} \right) \sum_{i=1}^p \|\mathbf{x}_i\|_1 + \frac{1}{np} \sum_{j=0}^{n-1} \left( \|\mathbf{q}_{-n}^\top \mathbf{M}_{\mathcal{I}_j^c}^j\|_1 - \|\mathbf{q}_{-n}^\top \mathbf{M}_{\mathcal{I}_j}^j\|_1 \right), \end{aligned}$$

601 where we denote  $\mathbf{M}^j = [(s_j [\tilde{\mathbf{x}}_1])_{-n} \quad (s_j [\tilde{\mathbf{x}}_2])_{-n} \quad \cdots \quad (s_j [\tilde{\mathbf{x}}_p])_{-n}]$ , and  $\mathbf{M}_{\mathcal{I}}^j$  denote a submatrix of  $\mathbf{M}^j$  with columns indexed by  $\mathcal{I}$ . Conditioned on the event  $\mathcal{E}$ , by Lemma D.5 and a union  
602 bound, whenever  $p \geq \frac{C}{\theta^2} n \log(\frac{n}{\theta})$ , we have

$$\|\mathbf{q}_{-n}^\top \mathbf{M}_{\mathcal{I}_j^c}^j\|_1 - \|\mathbf{q}_{-n}^\top \mathbf{M}_{\mathcal{I}_j}^j\|_1 \geq \frac{p}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\|, \quad \forall \mathbf{q}_{-n} \in \mathbb{R}^{n-1}, \quad (0 \leq j \leq n-1)$$

604 with probability at least  $1 - cnp^{-6}$ . On the other hand, by Gaussian concentration inequality, we have

$$\mathbb{P}\left(\frac{1}{np} \sum_{i=1}^p \|\mathbf{x}_i\|_1 \geq \frac{5}{4} \sqrt{\frac{2}{\pi}} \theta\right) \leq \exp\left(-\frac{\theta^2 p}{64\pi}\right).$$

605 Therefore, combining all the results above, we have

$$\tilde{\zeta}(\mathbf{q}) - \tilde{\zeta}\left(\frac{\mathbf{e}_n}{r_n}\right) \geq \begin{cases} \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\|, & \text{if } |q_n| - \frac{1}{r_n} \geq 0, \\ \frac{1}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{q}_{-n}\| - \frac{5}{4} \sqrt{\frac{2}{\pi}} \theta \left(\frac{1}{r_n} - |q_n|\right), & \text{if } |q_n| - \frac{1}{r_n} < 0, \end{cases}$$

606 as desired. ■

### 607 C.3 Linear convergence for projection subgradient descent in Algorithm 3

608 Now based on the sharpness condition, we are ready to show that the projected subgradient descent  
609 method

$$\mathbf{q}^{(k+1)} = \mathbf{q}^{(k)} - \tau^{(k)} \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)}, \quad \mathbf{g}^{(k)} = \sum_{i=1}^p (\mathbf{RQ}^{-1})^\top \mathbf{C}_{\mathbf{x}_i}^\top \text{sign}(\mathbf{C}_{\mathbf{x}_i} \mathbf{RQ}^{-1} \mathbf{q}^{(k)}).$$

610 on  $\zeta(\mathbf{q})$  converges linearly to the target solution up to a scaling factor. For convenience, let us first  
611 define the distance between the iterate and the target solution

$$d^{(k)} := \|\mathbf{s}^{(k)}\|, \quad \mathbf{s}^{(k)} := \mathbf{q}^{(k)} - (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n},$$

612 and several parameters

$$\alpha := \frac{1}{50} \sqrt{\frac{2}{\pi}} \theta, \quad \beta := 36 \log(np).$$

613 We show the following result.

614 **Proposition C.9** Suppose  $\theta \in (\frac{1}{n}, \frac{1}{3})$  and  $\mathbf{r}$  satisfies

$$\frac{\|\mathbf{r}_{-n}\|}{r_n} \leq \frac{1}{20}, \quad r_n > 0, \quad \|\mathbf{r}\| = 1. \quad (37)$$

615 Let  $\mathbf{q}^{(k)}$  be the sequence generated by the projected subgradient method (cf. Algorithm 3) with  
616 initialization  $\mathbf{q}^{(0)} = \mathbf{r}$  and geometrically decreasing step size

$$\tau^{(k)} = \eta^k \tau^{(0)}, \quad \tau^{(0)} = \frac{16}{25} \frac{\alpha}{\beta^2}, \quad \sqrt{1 - \frac{\alpha^2}{2\beta^2}} \leq \eta < 1 \quad (38)$$

617 Whenever  $p \geq C \frac{\kappa^8}{\theta \sigma_{\min}^2(C_a)} \log^3 n$ , with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ , the sequence  
618  $\{\mathbf{q}^{(k)}\}_{k \geq 0}$  satisfies

$$\left\| \mathbf{q}^{(k)} - (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| \leq \frac{2}{5} \eta^k, \quad (39)$$

619 for all iteration  $k = 0, 1, 2, \dots$ .

620 **Proof** Given the initialization  $\mathbf{q}^{(0)} = \mathbf{r}$ , we have

$$\begin{aligned} d^{(0)} &= \left\| \mathbf{r} - (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| \leq \left\| (\mathbf{RQ}^{-1})^{-1} \right\| \left\| \tilde{\mathbf{r}} - \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| \\ &\leq \frac{10}{9} \cdot \left( \|\tilde{\mathbf{r}}_{-n}\|^2 + \left( \tilde{r}_n - \frac{1}{\tilde{r}_n} \right)^2 \right)^{1/2}, \end{aligned}$$

621 where the last inequality we used Lemma H.4. From the argument in Proposition C.6, we know that  
622 (37) implies  $\|\tilde{\mathbf{r}}_{-n}\| / \tilde{r}_n \leq 1/10$ . By the fact that  $\|\tilde{\mathbf{r}}\| \leq 10/9$ , we have

$$\|\tilde{\mathbf{r}}_{-n}\| \leq \frac{1}{9}, \quad \left| \tilde{r}_n - \frac{1}{\tilde{r}_n} \right| \leq \left| \frac{8}{9} - \frac{9}{8} \right|^2 \leq \frac{1}{4} \implies d^{(0)} \leq \frac{2}{5}. \quad (40)$$

623 On the other hand, notice that

$$\begin{aligned} \left( d^{(k+1)} \right)^2 &= \left\| \mathbf{q}^{(k)} - \tau^{(k)} \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)} - (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\|^2 \\ &= \left( d^{(k)} \right)^2 - 2\tau^{(k)} \langle \mathbf{s}^{(k)}, \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)} \rangle + \left( \tau^{(k)} \right)^2 \left\| \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)} \right\|^2 \end{aligned}$$

624 By Lemma C.10, we know that when  $p \geq C \frac{\kappa^8}{\theta \sigma_{\min}^2(C_a)} \log^3 n$ , for any  $k = 1, 2, \dots$ ,

$$\left\| \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)} \right\|^2 \leq 36 \log(np) = \beta$$

625 holds with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ . On the other hand, by the sharpness property of  
626 the function in Proposition C.6, for any  $k = 1, 2, \dots$ ,

$$\begin{aligned} \langle \mathbf{s}^{(k)}, \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)} \rangle &= \langle \mathbf{s}^{(k)}, \mathbf{g}^{(k)} \rangle \geq \zeta(\mathbf{q}^{(k)}) - \zeta\left((\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n}\right) \\ &\geq \frac{1}{50} \sqrt{\frac{2}{\pi}} \theta \left\| \mathbf{q}^{(k)} - (\mathbf{RQ}^{-1})^{-1} \frac{\mathbf{e}_n}{\tilde{r}_n} \right\| = \alpha \cdot d^{(k)}, \end{aligned}$$

627 where the first equality follows from the fact that  $\langle \mathbf{r}, \mathbf{s}^{(k)} \rangle = 0$  so that  $\mathcal{P}_{\mathbf{r}^\perp} \mathbf{s}^{(k)} = \mathbf{s}^{(k)}$ , the first  
628 inequality follows from the fact that  $\zeta(\mathbf{q})$  is convex, and the second inequality utilizes the sharpness  
629 of the function in Proposition C.6 given the condition (37). Thus, we have

$$\left( d^{(k+1)} \right)^2 \leq \left( d^{(k)} \right)^2 - 2\alpha \cdot \tau^{(k)} \cdot d^{(k)} + \beta^2 \cdot \left( \tau^{(k)} \right)^2.$$

630 Now we proceed to prove (39) by induction. It is clear that (39) holds for  $\mathbf{q}^{(0)}$ . Suppose  $\mathbf{q}^{(k)}$  satisfies  
631 (39), i.e.,  $d^{(k)} \leq \eta^k d^{(0)}$  for some  $k \geq 1$ . The quadratic term of  $d^{(k)}$  on the right hand side of the  
632 inequality above will obtain its maximum at  $\frac{2}{5}\eta^k$  due to the definition of  $\tau^{(0)}$  and  $d^{(0)} \leq \frac{2}{5}$  as shown  
633 in (40). This, together with  $\tau^{(k)} = \eta\tau^{(k-1)}$ , it gives

$$\begin{aligned} \left( d^{(k+1)} \right)^2 &\leq \frac{4}{25}\eta^{2k} - \frac{4}{5}\alpha \cdot \eta^{2k}\tau^{(0)} + \beta^2 \cdot \eta^{2k} \left( \tau^{(0)} \right)^2 \\ &= \frac{4}{25}\eta^{2k} \cdot \left[ 1 - 5\alpha\tau^{(0)} + \frac{25}{4}\beta^2 \left( \tau^{(0)} \right)^2 \right] \leq \eta^{2k+2} \cdot \left( d^{(0)} \right)^2 \end{aligned}$$

634 where the last inequality follows from (38), where

$$1 - 5\alpha\tau^{(0)} + \frac{25}{4}\beta^2 \left( \tau^{(0)} \right)^2 \leq 1 - \alpha\tau^{(0)} \leq 1 - \frac{\alpha^2}{2\beta^2} \leq \eta^2 < 1.$$

635 This completes the proof. ■

636 **Lemma C.10** Suppose  $\theta \in (\frac{1}{n}, \frac{1}{3})$ . Whenever  $p \geq C \frac{\kappa^8}{\theta\sigma_{\min}^2(C_\alpha)} \log^3 n$ , we have

$$\rho := \sup_{\mathbf{q}: \mathbf{q}^\top \mathbf{r}=1} \frac{1}{np} \left\| \mathcal{P}_{\mathbf{r}^\perp} \sum_{i=1}^p (\mathbf{RQ}^{-1})^\top \mathbf{C}_{\mathbf{x}_i}^\top \text{sign}(\mathbf{C}_{\mathbf{x}_i} \mathbf{RQ}^{-1} \mathbf{q}) \right\| \leq 6\sqrt{\log(np)} \quad (41)$$

637 holds with probability at least  $1 - p^{-c_1 n^\theta} - n^{-c_2}$ .

638 **Proof** We have

$$\rho \leq \frac{1}{np} \|\mathbf{RQ}^{-1}\| \sum_{i=1}^p \left( \|\mathbf{C}_{\mathbf{x}_i}\| \sup_{\mathbf{q}: \mathbf{q}^\top \mathbf{r}=1} \|\text{sign}(\mathbf{C}_{\mathbf{x}_i} \mathbf{RQ}^{-1} \mathbf{q})\| \right).$$

639 Since the  $\text{sign}(\cdot)$  function is bounded by 1, we have

$$\rho \leq \frac{1}{np} \|\mathbf{RQ}^{-1}\| \cdot \left( \sum_{i=1}^p \|\mathbf{F}\mathbf{x}_i\|_\infty \right) \cdot \sqrt{n},$$

640 where we used the fact that  $\|\mathbf{C}_{\mathbf{x}_i}\| = \|\mathbf{F}\mathbf{x}_i\|_\infty$ . As  $\mathbf{x}_i \sim i.i.d. \mathcal{B}\mathcal{G}(\theta)$ , let  $\mathbf{x}_i = \mathbf{b}_i \odot \mathbf{g}_i$  with  $\mathbf{b}_i \sim \mathcal{B}(\theta)$   
641 and  $\mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Then we have

$$\|\mathbf{C}_{\mathbf{x}_i}\| = \|\mathbf{F}\mathbf{x}_i\|_\infty = \max_{1 \leq j \leq n} |(\mathbf{f}_j \odot \mathbf{b}_i)^* \mathbf{g}_i|.$$

642 By Gaussian concentration inequality in Lemma D.4 and a union bound, we have

$$\mathbb{P} \left( \max_{1 \leq i \leq p} \|\mathbf{F}\mathbf{x}_i\| \geq t \right) \leq (np) \cdot \exp \left( -\frac{t^2}{2n} \right).$$

643 Choose  $t = 4\sqrt{n \log(np)}$ , then we have

$$\max_{1 \leq i \leq p} \|\mathbf{F}\mathbf{x}_i\| \leq 4\sqrt{n \log(np)},$$

644 with probability at least  $1 - (np)^{-7}$ . On the other hand, by Lemma H.4, we know that whenever  
645  $p \geq C \frac{\kappa^8}{\theta\sigma_{\min}^2(C_\alpha)} \log^3 n$ , we have

$$\|\mathbf{RQ}^{-1}\| \leq \frac{3}{2},$$

646 holds with probability at least  $1 - p^{-c_1 n^\theta} - n^{-c_2}$ . Combining all the results above, we obtain

$$\rho \leq \frac{1}{np} \cdot \frac{3}{2} \cdot \left( 4p\sqrt{n \log(np)} \right) \cdot \sqrt{n} = 6\sqrt{\log(np)},$$

647 as desired. ■

648 **D Basics**

649 **Lemma D.1 (Moments of the Gaussian Random Variable)** *If  $X \sim \mathcal{N}(0, \sigma_X^2)$ , then it holds for  
650 all integer  $m \geq 1$  that*

$$\mathbb{E}[|X|^m] = \sigma_X^m (m-1)!! \left[ \sqrt{\frac{2}{\pi}} \mathbb{1}_{m=2k+1} + \mathbb{1}_{m=2k} \right] \leq \sigma_X^m (m-1)!! , k = \lfloor m/2 \rfloor.$$

651 **Lemma D.2 (sub-Gaussian Random Variables)** *Let  $X$  be a centered  $\sigma^2$  sub-Gaussian random  
652 variable, such that*

$$\mathbb{P}(|X| \geq t) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right),$$

653 *then for any integer  $p \geq 1$ , we have*

$$\mathbb{E}[|X|^p] \leq (2\sigma^2)^{p/2} p\Gamma(p/2).$$

654 *In particular, we have*

$$\|X\|_{L^p} = (\mathbb{E}[|X|^p])^{1/p} \leq \sigma e^{1/e} \sqrt{p}, \quad p \geq 2,$$

655 *and  $\mathbb{E}[|X|] \leq \sigma\sqrt{2\pi}$ .*

656 **Lemma D.3 (Moment-Control Bernstein's Inequality for Random Variables [42])** *Let  
657  $X_1, \dots, X_N$  be i.i.d. real-valued random variables. Suppose that there exist some positive  
658 numbers  $R$  and  $\sigma_X^2$  such that*

$$\mathbb{E}[|X_k|^m] \leq \frac{m!}{2} \sigma_X^2 R^{m-2}, \text{ for all integers } m \geq 2.$$

659 *Let  $S \doteq \frac{1}{N} \sum_{k=1}^N X_k$ , then for all  $t > 0$ , it holds that*

$$\mathbb{P}[|S - \mathbb{E}[S]| \geq t] \leq 2 \exp\left(-\frac{Nt^2}{2\sigma_X^2 + 2Rt}\right).$$

660 **Lemma D.4 (Gaussian Concentration Inequality)** *Let  $\mathbf{g} \in \mathbb{R}^n$  be a standard Gaussian random  
661 variable  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and let  $f : \mathbb{R}^n \mapsto \mathbb{R}$  denote an  $L$ -Lipschitz function. Then for all  $t > 0$ ,*

$$\mathbb{P}(|f(\mathbf{g}) - \mathbb{E}[f(\mathbf{g})]| \geq t) \leq 2 \exp\left(-\frac{t^2}{2L^2}\right).$$

662 **Lemma D.5 (Lemma VII.1, [34])** *Let  $\mathbf{M} \in \mathbb{R}^{n_1 \times n_2}$  with  $\mathbf{M} \sim \mathcal{BG}(\theta)$  and  $\theta \in (0, 1/3)$ . For a  
663 given set  $\mathcal{I} \subseteq [n_2]$  with  $|\mathcal{I}| \leq \frac{9}{8}\theta n_2$ , whenever  $n_2 \geq \frac{C}{\theta^2} n_1 \log(\frac{n_1}{\theta})$ , it holds*

$$\|\mathbf{v}^\top \mathbf{M}_{\mathcal{I}^c}\|_1 - \|\mathbf{v}^\top \mathbf{M}_{\mathcal{I}}\|_1 \geq \frac{n_2}{6} \sqrt{\frac{2}{\pi}} \theta \|\mathbf{v}\|$$

664 *for all  $\mathbf{v} \in \mathbb{R}^{n_1}$ , with probability at least  $1 - cn_2^{-6}$ .*

665 **Lemma D.6 (Derivates of  $h_\mu(z)$ )** *The first two derivatives of  $h_\mu(z)$  are*

$$\nabla h_\mu(z) = \begin{cases} \text{sign}(z) & |z| \geq \mu \\ z/\mu & |z| < \mu \end{cases}, \quad \nabla^2 h_\mu(z) = \begin{cases} 0 & |z| > \mu \\ 1/\mu & |z| < \mu \end{cases}. \quad (42)$$

666 *Whenever necessary, we define  $\nabla^2 h_\mu(\mu) = 0$ , and write the “second derivative” as  $\nabla^2 \bar{h}_\mu(\mu)$   
667 instead. Moreover for all  $z, z'$ ,*

$$|\nabla h_\mu(z) - \nabla h_\mu(z')| \leq \frac{1}{\mu} |z - z'|. \quad (43)$$

668 **Lemma D.7** Let  $X \sim \mathcal{N}(0, \sigma_x^2)$  and  $Y \sim \mathcal{N}(0, \sigma_y^2)$  and  $Z \sim \mathcal{N}(0, \sigma_z^2)$  be independent random  
669 variables. Then we have

$$\mathbb{E}[X \mathbb{1}_{X+Y \geq \mu}] = \frac{\sigma_x^2}{\sqrt{2\pi} \sqrt{\sigma_x^2 + \sigma_y^2}} \exp\left(-\frac{\mu^2}{2(\sigma_x^2 + \sigma_y^2)}\right), \quad (44)$$

$$\mathbb{E}[XY \mathbb{1}_{|X+Y| \leq \mu}] = -\sqrt{\frac{2}{\pi}} \frac{\mu \sigma_x^2 \sigma_y^2}{(\sigma_x^2 + \sigma_y^2)^{3/2}} \exp\left(-\frac{\mu^2}{2(\sigma_x^2 + \sigma_y^2)}\right), \quad (45)$$

$$\mathbb{E}[|X| \mathbb{1}_{|X| > \mu}] = \sqrt{\frac{2}{\pi}} \sigma_x \exp\left(-\frac{\mu^2}{2\sigma_x^2}\right), \quad (46)$$

$$\mathbb{E}[XY \mathbb{1}_{|X+Y+Z| < \mu}] = -\sqrt{\frac{2}{\pi}} \mu \exp\left(-\frac{\mu^2}{2(\sigma_x^2 + \sigma_y^2 + \sigma_z^2)}\right) \frac{\sigma_x^2 \sigma_y^2}{(\sigma_x^2 + \sigma_y^2 + \sigma_z^2)^{3/2}}, \quad (47)$$

$$\mathbb{E}[X^2 \mathbb{1}_{|X| < \mu}] = -\sqrt{\frac{2}{\pi}} \sigma_x \mu \exp\left(-\frac{\mu^2}{2\sigma_x^2}\right) + \sigma_x^2 \mathbb{P}[|X| < \mu], \quad (48)$$

$$\mathbb{E}[X^2 \mathbb{1}_{|X+Y| < \mu}] = -\sqrt{\frac{2}{\pi}} \mu \frac{\sigma_x^4}{(\sigma_x^2 + \sigma_y^2)^{3/2}} \exp\left(-\frac{\mu^2}{2(\sigma_x^2 + \sigma_y^2)}\right) + \sigma_x^2 \mathbb{P}[|X+Y| < \mu]. \quad (49)$$

670 **Proof** Direct calculations. ■

671 **Lemma D.8 (Calculus for Function of Matrices, Chapter X of [43])** Let  $\mathcal{S}^{n \times n}$  be the set of symmetric  
672 matrices of size  $n \times n$ . We define a map  $f : \mathcal{S}^{n \times n} \mapsto \mathcal{S}^{n \times n}$  as

$$f(\mathbf{A}) = \mathbf{U} f(\boldsymbol{\Lambda}) \mathbf{U}^*,$$

673 where  $\mathbf{A} \in \mathcal{S}^{n \times n}$  has the eigen-decomposition  $\mathbf{A} = \mathbf{U} \boldsymbol{\Lambda} \mathbf{U}^*$ . The map  $f$  is called (Fréchet)  
674 differentiable at  $\mathbf{A}$  if there exists a linear transformation on  $\mathcal{S}^{n \times n}$  such that for all  $\boldsymbol{\Delta}$

$$\|f(\mathbf{A} + \boldsymbol{\Delta}) - f(\mathbf{A}) - \mathrm{D}f(\mathbf{A})[\boldsymbol{\Delta}]\| = o(\|\boldsymbol{\Delta}\|).$$

675 The linear operator  $\mathrm{D}f(\mathbf{A})$  is called the derivative of  $f$  at  $\mathbf{A}$ , and  $\mathrm{D}f(\mathbf{A})[\boldsymbol{\Delta}]$  is the directional  
676 derivative of  $f$  along  $\boldsymbol{\Delta}$ . If  $f$  is differentiable at  $\mathbf{A}$ , then

$$\mathrm{D}f(\mathbf{A})[\boldsymbol{\Delta}] = \frac{d}{dt} f(\mathbf{A} + t\boldsymbol{\Delta}) \Big|_{t=0}.$$

677 We denote the operator norm of the derivative  $\mathrm{D}f(\mathbf{A})$  as

$$\|\mathrm{D}f(\mathbf{A})\| \doteq \sup_{\|\boldsymbol{\Delta}\|=1} \|\mathrm{D}f(\mathbf{A})[\boldsymbol{\Delta}]\|.$$

678 **Lemma D.9 (Mean Value Theorem for Function of Matrices)** Let  $f$  be a differentiable map from  
679 a convex subset  $\mathcal{U}$  of a Banach space  $\mathcal{X}$  into the Banach space  $\mathcal{Y}$ . Let  $\mathbf{A}, \mathbf{B} \in \mathcal{U}$ , and let  $\mathcal{L}$  be the  
680 line segment joining them. Then

$$\|f(\mathbf{B}) - f(\mathbf{A})\| \leq \|\mathbf{B} - \mathbf{A}\| \sup_{\mathbf{U} \in \mathcal{L}} \|\mathrm{D}f(\mathbf{U})\|.$$

681 **Lemma D.10 (Theorem VII.2.3 of [43])** Let  $\mathbf{A}$  and  $\mathbf{B}$  be operators whose spectra are contained  
682 in the open right half-plane and open left half-plane, respectively. Then the solution of the equation  
683  $\mathbf{AX} - \mathbf{XB} = \mathbf{Y}$  can be expressed as

$$\mathbf{X} = \int_0^\infty e^{-t\mathbf{A}} \mathbf{Y} e^{t\mathbf{B}} dt$$

684 **Lemma D.11** Let  $f(\mathbf{A}) = \mathbf{A}^{-1/2}$ , defined the set of all  $n \times n$  positive definite matrices  $\mathcal{S}_+^{n \times n}$ , then  
685 we have

$$\|\mathrm{D}f(\mathbf{A})\| \leq \frac{1}{\sigma_{\min}^2(\mathbf{A})},$$

686 where  $\sigma_{\min}(\mathbf{A})$  is the smallest singular value of  $\mathbf{A}$ .

687 **Proof** To bound the operator norm  $\|Df(\mathbf{A})\|$ , we introduce an auxiliary function

$$g(\mathbf{A}) = \mathbf{A}^{-2}, \quad f(\mathbf{A}) = g^{-1}(\mathbf{A}),$$

688 such that  $f$  and  $g$  are the inverse function to each other. Whenever  $g \circ f(\mathbf{A}) \neq 0$  (which is true for  
689 our case  $\mathbf{A} > \mathbf{0}$ ), this gives

$$Df(\mathbf{A}) = [D(g \circ f)(\mathbf{A})]^{-1} = [Dg(\mathbf{A}^{-1/2})]^{-1}. \quad (50)$$

690 This suggests that we can estimate  $Df(\mathbf{A})$  via estimating  $Dg(\mathbf{A})$  of its inverse function  $g$ . Let

$$g = h \circ w(\mathbf{A}), \quad h(\mathbf{A}) = \mathbf{A}^{-1}, \quad w(\mathbf{A}) = \mathbf{A}^2,$$

691 such that their directional derivatives have simple form

$$Dh(\mathbf{A})[\Delta] = -\mathbf{A}^{-1}\Delta\mathbf{A}^{-1}, \quad Dw(\mathbf{A})[\Delta] = \Delta\mathbf{A} + \mathbf{A}\Delta.$$

692 By using chain rule, simple calculation gives

$$\begin{aligned} Dg(\mathbf{A})[\Delta] &= Dh(w(\mathbf{A})) [Dw(\mathbf{A})[\Delta]], \\ &= -(A^{-2}\Delta A^{-1} + A^{-1}\Delta A^{-2}). \end{aligned}$$

693 Now by (50), the directional derivative

$$Z \doteq Df(\mathbf{A})[\Delta]$$

694 satisfies

$$\mathbf{A}ZA^{1/2} + A^{1/2}ZA = -\Delta.$$

695 Since  $\mathbf{A} > \mathbf{0}$ , we write the eigen decomposition as  $\mathbf{A} = \mathbf{U}\Lambda\mathbf{U}^*$ , with  $\mathbf{U}$  orthogonal and  $\Lambda > 0$   
696 diagonal. Let  $\tilde{\mathbf{Z}} = \mathbf{U}^*Z\mathbf{U}$  and  $\tilde{\Delta} = \mathbf{U}^*\Delta\mathbf{U}$ , then the equation above gives

$$\Lambda^{1/2}\tilde{\mathbf{Z}} - \tilde{\mathbf{Z}}(-\Lambda^{1/2}) = -\Lambda^{-1/2}\tilde{\Delta}\Lambda^{-1/2},$$

697 which is the Sylvester equation []. Since  $\Lambda^{1/2}$  and  $-\Lambda^{1/2}$  do not have common eigenvalues, Lemma  
698 D.10 gives

$$Df(\mathbf{A})[\Delta] = \mathbf{U} \left[ \int_0^\infty e^{-\Lambda^{1/2}\tau} (-\Lambda^{-1/2}\tilde{\Delta}\Lambda^{-1/2}) e^{-\Lambda^{1/2}\tau} d\tau \right] \mathbf{U}^*.$$

699 Thus, by Lemma D.8 we know that

$$\begin{aligned} \|Df(\mathbf{A})\| &= \sup_{\|\Delta\|=1} \|Df(\mathbf{A})[\Delta]\| \\ &\leq \int_0^\infty \left\| e^{-\Lambda^{1/2}\tau} (-\Lambda^{-1/2}\tilde{\Delta}\Lambda^{-1/2}) e^{-\Lambda^{1/2}\tau} \right\| d\tau \\ &\leq \left\| \Lambda^{-1/2}\tilde{\Delta}\Lambda^{-1/2} \right\| \int_0^\infty e^{-\sigma_{\min}\tau} d\tau \leq \frac{1}{\sigma_{\min}^2(\mathbf{A})}. \end{aligned}$$

700 ■

701 **Lemma D.12 (Matrix Perturbation Bound)** Suppose  $\mathbf{A} > \mathbf{0}$ . Then for any symmetric perturba-  
702 tion matrix  $\Delta$  with  $\|\Delta\| \leq \frac{1}{2}\sigma_{\min}(\mathbf{A})$ , it holds that

$$\|(\mathbf{A} + \Delta)^{-1/2} - \mathbf{A}^{-1/2}\| \leq \frac{4\|\Delta\|}{\sigma_{\min}^2(\mathbf{A})},$$

703 where  $\sigma_{\min}(\mathbf{A})$  denotes the minimum singular value of  $\mathbf{A}$ .

704 **Proof** Let us denote  $f(\mathbf{A}) = \mathbf{A}^{-1/2}$ . Given a symmetric perturbation matrix  $\Delta$ , by mean value  
705 theorem, we have

$$\begin{aligned} \|(\mathbf{A} + \Delta)^{-1/2} - \mathbf{A}^{-1/2}\| &= \left\| \int_0^1 Df(\mathbf{A} + t\Delta)[\Delta] dt \right\| \\ &\leq \left( \sup_{t \in [0,1]} \|Df(\mathbf{A} + t\Delta)\| \right) \cdot \|\Delta\|. \end{aligned}$$

706 Thus, by Lemma D.11 and by using the fact that  $\|\Delta\| \leq \frac{1}{2}\sigma_{\min}(\mathbf{A})$ , we have

$$\|(\mathbf{A} + \Delta)^{-1/2} - \mathbf{A}^{-1/2}\| \leq \left( \sup_{t \in [0,1]} \frac{1}{\sigma_{\min}^2(\mathbf{A} + t\Delta)} \right) \|\Delta\| \leq \frac{4\|\Delta\|}{\sigma_{\min}^2(\mathbf{A})}.$$

707 ■

## 708 E Regularity Condition in Population

709 **Proposition E.1** For every  $i \in [n]$ , define a set

$$\mathbf{q} \in \mathcal{S}_\xi^{i+} \doteq \left\{ \mathbf{q} \in \mathbb{R}^n \mid q_i > 0, \frac{q_i}{\|\mathbf{q}_{-i}\|_\infty} \geq \sqrt{1 + \xi} \right\}.$$

710 Whenever  $\theta \in (\frac{1}{n}, c_0)$  and  $\mu \leq c_1 \min\{\theta, \frac{1}{\sqrt{n}}\}$ , we have

$$\langle \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})], q_i \mathbf{q} - \mathbf{e}_i \rangle \geq c_2 \theta (1 - \theta) q_i \|\mathbf{q}_{-i}\|, \quad \sqrt{1 - q_i^2} \in [\mu, c_3] \quad (51)$$

$$\langle \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})], q_i \mathbf{q} - \mathbf{e}_i \rangle \geq c_2 \theta (1 - \theta) q_i n^{-1} \|\mathbf{q}_{-i}\|, \quad \sqrt{1 - q_i^2} \in \left[ c_3, \sqrt{\frac{n-1}{n}} \right], \quad (52)$$

711 hold for any  $\mathbf{q} \in \mathcal{S}_\xi^{i+}$  and each  $i \in [n]$ .

712 **Remarks.** For proving this result, we first introduce some basic notations. We use  $\mathcal{I}$  to denote the  
713 generic support set of  $\mathbf{q} \in \mathbb{S}^{n-1}$  of i.i.d.  $\mathcal{B}(\theta)$  law. Since the landscape is symmetric for each  $i \in [n]$ ,  
714 without loss of generality, it is enough to consider the case when  $i = n$ . We reparameterize  $\mathbf{q} \in \mathbb{S}^{n-1}$   
715 by

$$\mathbf{q}(\mathbf{w}) : \mathbf{w} \mapsto \begin{bmatrix} \mathbf{w} \\ \sqrt{1 - \|\mathbf{w}\|^2} \end{bmatrix}, \quad (53)$$

716 where  $\mathbf{w} \in \mathbb{R}^{n-1}$  with  $\|\mathbf{w}\| \leq \sqrt{\frac{n-1}{n}}$ . We write

$$\mathbf{q}_{\mathcal{I}} = \begin{bmatrix} \mathbf{w}_{\mathcal{J}} \\ q_n \mathbb{1}_{n \in \mathcal{I}} \end{bmatrix},$$

717 where we use  $\mathcal{J}$  to denote the support set of  $\mathbf{w}$  of i.i.d.  $\mathcal{B}(\theta)$  law.

718 **Proof** We denote

$$g(\mathbf{w}) = h_\mu \left( \mathbf{w}^\top \mathbf{x}_{-n} + x_n \sqrt{1 - \|\mathbf{w}\|^2} \right) \quad (54)$$

719 Note that if  $\mathbf{e}_n$  is a local minimizer of  $\mathbb{E}[\tilde{f}(\mathbf{q})]$ , then  $\mathbb{E}[g(\mathbf{w})]$  has a corresponding local minimum  
720 at  $\mathbf{0}$ . Since  $g(\cdot)$  satisfies chain rule when computing its gradient, we have

$$\begin{aligned} \langle \mathbb{E}[\nabla g(\mathbf{w})], \mathbf{w} - \mathbf{0} \rangle &= \left\langle \begin{bmatrix} \mathbf{I}_{n-1} & \frac{-\mathbf{w}}{\sqrt{1 - \|\mathbf{w}\|^2}} \end{bmatrix} \nabla \mathbb{E}[\tilde{f}(\mathbf{q})], \mathbf{w} \right\rangle \\ &= \left\langle \mathbb{E}[\nabla \tilde{f}(\mathbf{q})], \mathbf{q} - \frac{1}{q_n} \mathbf{e}_n \right\rangle = \frac{1}{q_n} \langle \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle, \end{aligned}$$

721 which gives

$$\langle \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})], q_n \mathbf{q} - \mathbf{e}_n \rangle = q_n \langle \mathbb{E}[\nabla g(\mathbf{w})], \mathbf{w} \rangle. \quad (55)$$

722 Thus, the above relationship implies that we can work on the “unconstrained” function  $g(\mathbf{w})$  and  
723 establish the following: for any  $\mathbf{q}(\mathbf{w}) \in \mathcal{S}_\xi^{n+}$  with  $\xi > 0$ , or equivalently,

$$\|\mathbf{w}\|^2 + (1 + \xi) \|\mathbf{w}\|_\infty^2 \leq 1,$$

<sup>724</sup> the following holds

$$\langle \nabla \mathbb{E}[g(\mathbf{w})], \mathbf{w} - \mathbf{0} \rangle \gtrsim \|\mathbf{w}\|.$$

<sup>725</sup> When  $\|\mathbf{w}\| \in [c_0\mu, c_1]$ , Lemma E.4 implies that

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \geq c_2\theta(1-\theta)\|\mathbf{w}\|.$$

<sup>726</sup> By Lemma E.5, we know that when  $c_1 \leq \|\mathbf{w}\| \leq \sqrt{\frac{n-1}{n}}$ ,

$$\mathbf{w}^\top \nabla^2 \mathbb{E}[g(\mathbf{w})] \mathbf{w} \leq -c_3\theta(1-\theta)\|\mathbf{w}\|^2,$$

<sup>727</sup> which implies concavity of  $g(\mathbf{w})$  along the  $\mathbf{w}$  direction. Let us denote  $\mathbf{v} = \mathbf{w}/\|\mathbf{w}\|$ , then the  
<sup>728</sup> directional concavity implies that

$$\mathbf{t}\mathbf{v}^\top \nabla \mathbb{E}[g(t\mathbf{v})] \geq (\mathbf{t}'\mathbf{v})^\top \nabla \mathbb{E}[g(\mathbf{t}'\mathbf{v})] + c_4\theta(1-\theta)(\mathbf{t}' - \mathbf{t}),$$

<sup>729</sup> for any  $t, t' \in [c_1, \sqrt{\frac{n-1}{n}}]$ . Choose  $t' = \frac{\|\mathbf{w}\|}{\sqrt{\|\mathbf{w}\|^2 + \|\mathbf{w}\|_\infty^2}}$  and  $t = \|\mathbf{w}\|$ , by Lemma E.3, we know that

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \geq c_4\theta(1-\theta)\|\mathbf{w}\| \left( \frac{1}{\sqrt{\|\mathbf{w}\|^2 + \|\mathbf{w}\|_\infty^2}} - 1 \right).$$

<sup>730</sup> The function

$$h_{\mathbf{v}}(t) \doteq \frac{\|t\mathbf{v}\|}{\sqrt{\|t\mathbf{v}\|^2 + \|t\mathbf{v}\|_\infty^2}} - \|t\mathbf{v}\| = \frac{1}{\sqrt{1 + \|\mathbf{v}\|_\infty^2}} - t$$

<sup>731</sup> is obviously monotonically decreasing w.r.t.  $t$ . Since  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ , we have

$$\|t\mathbf{v}\|^2 + (1+\xi)\|t\mathbf{v}\|_\infty^2 \leq 1 \implies t \leq \frac{1}{\sqrt{1 + (1+\xi)\|\mathbf{v}\|_\infty^2}}.$$

<sup>732</sup> Therefore, we can uniformly lower bound  $h_{\mathbf{v}}(t)$  by

$$h_{\mathbf{v}}(t) \geq \frac{1}{\sqrt{1 + \|\mathbf{v}\|_\infty^2}} - \frac{1}{\sqrt{1 + (1+\xi)\|\mathbf{v}\|_\infty^2}} \geq \xi \|\mathbf{v}\|_\infty^2 \geq \xi n^{-1}$$

<sup>733</sup> Therefore, we have

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \geq c_4\xi\theta(1-\theta)n^{-1}\|\mathbf{w}\|,$$

<sup>734</sup> when  $\|\mathbf{w}\| \in [c_1, \sqrt{\frac{n-1}{n}}]$ . Combining the bounds above, we obtain the desired results. ■

<sup>735</sup> **Lemma E.2** Suppose  $\mathbf{g} \in \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ , we have

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] = \frac{1}{\mu} \mathbb{E}_{\mathcal{I}} \left[ \left( \|\mathbf{q}_{\mathcal{I}}\|^2 - \mathbb{1}_{n \in \mathcal{I}} \right) \mathbb{P}(|\mathbf{q}_{\mathcal{I}}^\top \mathbf{g}| \leq \mu) \right]. \quad (56)$$

<sup>736</sup> **Proof** In particular, exchange of gradient and expectation operator can again be justified. By simple  
<sup>737</sup> calculation, we obtain that

$$\nabla g(\mathbf{w}) = h'_\mu(\mathbf{q}^\top \mathbf{x}) \left( \mathbf{x}_{-n} - \frac{x_n}{q_n} \mathbf{w} \right) = \begin{cases} \frac{\mathbf{q}^\top \mathbf{x}}{\mu} \left( \mathbf{x}_{-n} - \frac{x_n}{q_n} \mathbf{w} \right), & |\mathbf{q}^\top \mathbf{x}| \leq \mu \\ \text{sign}(\mathbf{q}^\top \mathbf{x}) \left( \mathbf{x}_{-n} - \frac{x_n}{q_n} \mathbf{w} \right), & |\mathbf{q}^\top \mathbf{x}| > \mu. \end{cases} \quad (57)$$

<sup>738</sup> Thus, we obtain

$$\begin{aligned} & \mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \\ &= \mathbb{E} \left[ \text{sign}(\mathbf{q}^\top \mathbf{x}) \left( \mathbf{w}^\top \mathbf{x}_{-n} - \frac{x_n}{q_n} \|\mathbf{w}\|^2 \right) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \geq \mu} \right] + \mathbb{E} \left[ \frac{\mathbf{q}^\top \mathbf{x}}{\mu} \left( \mathbf{w}^\top \mathbf{x}_{-n} - \frac{x_n}{q_n} \|\mathbf{w}\|^2 \right) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu} \right] \\ &= \mathbb{E} \left[ \text{sign}(\mathbf{q}^\top \mathbf{x}) \left( \mathbf{q}^\top \mathbf{x} - \frac{x_n}{q_n} \right) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \geq \mu} \right] + \frac{1}{\mu} \mathbb{E} \left[ \left( \mathbf{q}^\top \mathbf{x} - \frac{x_n}{q_n} \right) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu} \right], \end{aligned}$$

<sup>739</sup> where we used the fact that

$$\mathbf{w}^\top \mathbf{x}_{-n} - \frac{x_n}{q_n} \|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{x}_{-n} + q_n x_n - x_n \frac{\|\mathbf{w}\|^2 + q_n^2}{q_n} = \mathbf{q}^\top \mathbf{x} - \frac{x_n}{q_n}.$$

<sup>740</sup> Let  $Z = X + Y$ , with

$$X = \mathbf{w}^\top \mathbf{x}_{-n} \sim \mathcal{N}(\mathbf{0}, \|\mathbf{w}_\mathcal{J}\|^2), \quad Y = q_n x_n \sim \mathcal{N}(0, q_n^2 \mathbb{1}_{n \in \mathcal{I}}), \quad Z \sim \mathcal{N}(\mathbf{0}, \|\mathbf{q}_\mathcal{I}\|^2). \quad (58)$$

<sup>741</sup> This gives

$$\begin{aligned} \mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] &= \mathbb{E}[|\mathbf{q}^\top \mathbf{x}| \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \geq \mu}] - \frac{1}{q_n} \mathbb{E}[\text{sign}(\mathbf{q}^\top \mathbf{x}) x_n \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \geq \mu}] \\ &\quad + \frac{1}{\mu} \mathbb{E}[(\mathbf{q}^\top \mathbf{x})^2 \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}] - \frac{1}{q_n \mu} \mathbb{E}[x_n (\mathbf{w}^\top \mathbf{x}_{-n}) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}] - \frac{1}{\mu} \mathbb{E}[x_n^2 \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}] \\ &= \mathbb{E}[|Z| \mathbb{1}_{|Z| \geq \mu}] - \frac{1}{q_n^2} \mathbb{E}[\text{sign}(X + Y) Y \mathbb{1}_{|X + Y| \geq \mu}] + \frac{1}{\mu} \mathbb{E}[Z^2 \mathbb{1}_{|Z| \leq \mu}] \\ &\quad - \frac{1}{\mu q_n^2} \mathbb{E}[XY \mathbb{1}_{|X + Y| \leq \mu}] - \frac{1}{\mu q_n^2} \mathbb{E}[Y^2 \mathbb{1}_{|X + Y| \leq \mu}]. \end{aligned}$$

<sup>742</sup> Now by Lemma D.7, we have

$$\begin{aligned} \mathbb{E}[|Z| \mathbb{1}_{|Z| \geq \mu}] &= \sqrt{\frac{2}{\pi}} \mathbb{E}_\mathcal{I} \left[ \|\mathbf{q}_\mathcal{I}\| \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_\mathcal{I}\|^2} \right) \right] \\ \mathbb{E}[\text{sign}(X + Y) Y \mathbb{1}_{|X + Y| \geq \mu}] &= q_n^2 \sqrt{\frac{2}{\pi}} \mathbb{E} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_\mathcal{I}\|} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_\mathcal{I}\|^2} \right) \right] \\ \mathbb{E}[Z^2 \mathbb{1}_{|Z| \leq \mu}] &= -\mu \sqrt{\frac{2}{\pi}} \mathbb{E}_\mathcal{I} \left[ \|\mathbf{q}_\mathcal{I}\| \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_\mathcal{I}\|^2} \right) \right] + \mathbb{E}_\mathcal{I} \left[ \|\mathbf{q}_\mathcal{I}\|^2 \mathbb{P}(|\mathbf{q}_\mathcal{I}^\top \mathbf{g}| \leq \mu) \right] \\ \mathbb{E}[XY \mathbb{1}_{|X + Y| \leq \mu}] &= -\mu q_n^2 \sqrt{\frac{2}{\pi}} \mathbb{E}_\mathcal{I} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}} \|\mathbf{w}_\mathcal{J}\|^2}{\|\mathbf{q}_\mathcal{I}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_\mathcal{I}\|^2} \right) \right] \\ \mathbb{E}[Y^2 \mathbb{1}_{|X + Y| \leq \mu}] &= -\mu q_n^4 \sqrt{\frac{2}{\pi}} \mathbb{E}_\mathcal{I} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_\mathcal{I}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_\mathcal{I}\|^2} \right) \right] + q_n^2 \mathbb{E}_\mathcal{I} [\mathbb{1}_{n \in \mathcal{I}} \mathbb{P}(|\mathbf{q}_\mathcal{I}^\top \mathbf{g}| \leq \mu)] \end{aligned}$$

<sup>743</sup> Putting the above calculations together and simplify, we obtain the desired result in (56). ■

<sup>744</sup>

<sup>745</sup> **Lemma E.3** When for any  $\mathbf{w} \in \mathbb{R}^{n-1}$  satisfies  $\|\mathbf{w}\|^2 + \|\mathbf{w}\|_\infty^2 \leq 1$ , we have

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \geq 0.$$

<sup>746</sup> **Proof** From Lemma E.2, we know that

$$\begin{aligned}
& \mu \cdot \mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \\
&= \mathbb{E}_{\mathcal{I}} \left[ \left( \|\mathbf{q}_{\mathcal{I}}\|^2 - \mathbb{1}_{n \in \mathcal{I}} \right) \mathbb{P}(|\mathbf{q}_{\mathcal{I}}^\top \mathbf{g}| \leq \mu) \right] \\
&= \mathbb{E}_{\mathcal{J}} \left[ (1-\theta) \|\mathbf{w}_{\mathcal{J}}\|^2 \mathbb{P}(|\mathbf{g}_{-\mathcal{J}}^\top \mathbf{w}_{\mathcal{J}}| \leq \mu) - \theta \|\mathbf{w}_{\mathcal{J}^c}\|^2 \mathbb{P}(|\mathbf{g}_{-\mathcal{J}}^\top \mathbf{w}_{\mathcal{J}} + q_n g_n| \leq \mu) \right] \\
&= \mathbb{E}_{\mathcal{J}} \left[ \int_{-\mu}^{\mu} \left( \frac{1-\theta}{\sqrt{2\pi}} \frac{\|\mathbf{w}_{\mathcal{J}}\|^2}{\|\mathbf{w}_{\mathcal{J}}\|} \exp\left(-\frac{t^2}{2\|\mathbf{w}_{\mathcal{J}}\|^2}\right) - \frac{\theta}{\sqrt{2\pi}} \frac{\|\mathbf{w}_{\mathcal{J}^c}\|^2}{\sqrt{1-\|\mathbf{w}_{\mathcal{J}^c}\|^2}} \exp\left(\frac{-t^2}{2-2\|\mathbf{w}_{\mathcal{J}^c}\|^2}\right) \right) dt \right] \\
&= \frac{1-\theta}{\sqrt{2\pi}} \sum_{i=1}^{n-1} \int_{-\mu}^{\mu} \mathbb{E}_{\mathcal{J}} \left[ \frac{w_i^2 \mathbb{1}_{i \in \mathcal{J}}}{\sqrt{w_i^2 \mathbb{1}_{i \in \mathcal{J}} + \|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2}} \exp\left(-\frac{t^2}{2w_i^2 \mathbb{1}_{i \in \mathcal{J}} + 2\|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2}\right) \right] dt \\
&\quad - \frac{\theta}{\sqrt{2\pi}} \sum_{i=1}^{n-1} \int_{-\mu}^{\mu} \mathbb{E}_{\mathcal{J}} \left[ \frac{w_i^2 \mathbb{1}_{i \notin \mathcal{J}}}{\sqrt{1-w_i^2 \mathbb{1}_{i \notin \mathcal{J}} - \|\mathbf{w}_{\mathcal{J}^c \setminus \{i\}}\|^2}} \exp\left(-\frac{t^2}{2-2w_i^2 \mathbb{1}_{i \notin \mathcal{J}} - 2\|\mathbf{w}_{\mathcal{J}^c \setminus \{i\}}\|^2}\right) \right] dt \\
&= \frac{(1-\theta)\theta}{\sqrt{2\pi}} \sum_{i=1}^{n-1} \int_{-\mu}^{\mu} \mathbb{E}_{\mathcal{J}} \left[ \frac{w_i^2}{\sqrt{w_i^2 + \|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2}} \exp\left(-\frac{t^2}{2w_i^2 + 2\|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2}\right) \right] dt \\
&\quad - \frac{(1-\theta)\theta}{\sqrt{2\pi}} \sum_{i=1}^{n-1} \int_{-\mu}^{\mu} \mathbb{E}_{\mathcal{J}} \left[ \frac{w_i^2}{\sqrt{1-\|\mathbf{w}\|^2 + \|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2}} \exp\left(-\frac{t^2}{2-2\|\mathbf{w}\|^2 + 2\|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2}\right) \right] dt \\
&= (1-\theta)\theta \sum_{i=1}^{n-1} w_i^2 \mathbb{E}_{\mathcal{J}} [\mathbb{P}(|Z_{i1}| \leq \mu) - \mathbb{P}(|Z_{i2}| \leq \mu)], \tag{59}
\end{aligned}$$

<sup>747</sup> where

$$Z_{i1} \sim \mathcal{N}(0, w_i^2 + \|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2), \quad Z_{i2} \sim \mathcal{N}(0, 1 - \|\mathbf{w}\|^2 + \|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2). \tag{60}$$

<sup>748</sup> Since we have  $1 - \|\mathbf{w}\|^2 \geq \|\mathbf{w}\|_\infty^2 \geq w_i^2$ , the variance of  $Z_i^2$  is larger than that of  $Z_i^1$ . Therefore, we <sup>749</sup> have  $\mathbb{P}(|Z_{i1}| \leq \mu) \geq \mathbb{P}(|Z_{i2}| \leq \mu)$  for each  $i = 1, \dots, n-1$ . Hence, we obtain

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] = \frac{1}{\mu} \theta(1-\theta) \sum_{i=1}^{n-1} w_i^2 \mathbb{E}_{\mathcal{J}} [\mathbb{P}(|Z_{i1}| \leq \mu) - \mathbb{P}(|Z_{i2}| \leq \mu)] \geq 0.$$

<sup>750</sup>

■

<sup>751</sup> **Lemma E.4** For any  $\mathbf{w}$  with  $c_0\mu \leq \|\mathbf{w}\| \leq c_1$ , we have

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] \geq c\theta(1-\theta) \|\mathbf{w}\|$$

<sup>752</sup> **Proof** Recall from (59), we have

$$\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] = \frac{1}{\mu} (1-\theta)\theta \sum_{i=1}^{n-1} w_i^2 \mathbb{E}_{\mathcal{J}} [\mathbb{P}(|Z_{i1}| \leq \mu) - \mathbb{P}(|Z_{i2}| \leq \mu)],$$

<sup>753</sup> where  $Z_{i1}$  and  $Z_{i2}$  are defined the same as (60). Let us denote

$$Z_1 \sim \mathcal{N}(0, \|\mathbf{w}\|^2), \quad Z_2 \sim \mathcal{N}(0, 1 - \|\mathbf{w}\|^2).$$

<sup>754</sup> Since we have  $\|\mathbf{w}\|^2 \geq w_i^2 + \|\mathbf{w}_{\mathcal{J} \setminus \{i\}}\|^2$ , the variance of  $Z_1$  is larger than that of  $Z_{i1}$ . Therefore, <sup>755</sup> we have  $\mathbb{P}(|Z_{i1}| \leq \mu) \geq \mathbb{P}(|Z_1| \leq \mu)$  for each  $i = 1, \dots, n-1$ . By a similar argument, we have

756  $\mathbb{P}(|Z_{i2}| \leq \mu) \leq \mathbb{P}(|Z_2| \leq \mu)$  for each  $i = 1, \dots, n-1$ . Thus, we obtain

$$\begin{aligned}
& \mathbb{P}(|Z_{i1}| \leq \mu) - \mathbb{P}(|Z_{i2}| \leq \mu) \\
& \geq \mathbb{P}(|Z_1| \leq \mu) - \mathbb{P}(|Z_2| \leq \mu) \\
& = \sqrt{\frac{2}{\pi}} \frac{1}{\|\mathbf{w}\|} \int_0^\mu \exp\left(-\frac{t^2}{2\|\mathbf{w}\|^2}\right) dt - \sqrt{\frac{2}{\pi}} \frac{1}{\sqrt{1-\|\mathbf{w}\|^2}} \int_0^\mu \exp\left(-\frac{t^2}{2-2\|\mathbf{w}\|^2}\right) dt \\
& \geq \sqrt{\frac{2}{\pi}} \left[ \frac{1}{\|\mathbf{w}\|} \int_0^\mu \left(1 - \frac{t^2}{2\|\mathbf{w}\|^2}\right) dt - \frac{\mu}{\sqrt{1-\|\mathbf{w}\|^2}} \right] \\
& = \sqrt{\frac{2}{\pi}} \left[ \frac{1}{\|\mathbf{w}\|} \left(\mu - \frac{1}{6} \frac{\mu^3}{\|\mathbf{w}\|^2}\right) - \frac{\mu}{\sqrt{1-\|\mathbf{w}\|^2}} \right] \\
& \geq \mu \sqrt{\frac{2}{\pi}} \left( \frac{1}{\|\mathbf{w}\|} - 2 \frac{1}{\sqrt{1-\|\mathbf{w}\|^2}} \right) \geq \frac{\mu}{2\sqrt{2\pi}} \frac{1}{\|\mathbf{w}\|} \tag{61}
\end{aligned}$$

757 where we used the fact that  $\mu/\sqrt{3} \leq \|\mathbf{w}\| \leq 1/\sqrt{17}$  for the last two inequalities. Plugging (61) back  
758 into (59) gives

$$\begin{aligned}
\mathbf{w}^\top \nabla \mathbb{E}[g(\mathbf{w})] &= \frac{1}{\mu} (1-\theta)\theta \sum_{i=1}^{n-1} w_i^2 \mathbb{E}_{\mathcal{J}}[\mathbb{P}(|Z_{i1}| \leq \mu) - \mathbb{P}(|Z_{i2}| \leq \mu)] \\
&\geq \frac{(1-\theta)\theta}{2\sqrt{2\pi}\|\mathbf{w}\|} \sum_{i=1}^{n-1} w_i^2 = \frac{1}{2\sqrt{2\pi}} (1-\theta)\theta \|\mathbf{w}\|,
\end{aligned}$$

759 as desired. ■

760 **Lemma E.5** When  $\mu \leq c_0 \min\left\{\frac{1}{\sqrt{n}}, \theta\right\}$  and  $\theta \in (\frac{1}{n}, c_1)$ , we have

$$\mathbf{w}^\top \nabla^2 \mathbb{E}[g(\mathbf{w})] \mathbf{w} \leq -c_2 \theta (1-\theta) \|\mathbf{w}\|^2$$

761 for all  $\mathbf{w}$  with  $c_3 \leq \|\mathbf{w}\| \leq \sqrt{\frac{n-1}{n}}$ . Here,  $c_0$ ,  $c_1$ ,  $c_2$ , and  $c_3$  are some numerical constants.

762 **Proof** Since the expectation and derivative are exchangeable, we have

$$\mathbf{w}^\top \nabla^2 \mathbb{E}[g(\mathbf{w})] \mathbf{w} = \mathbf{w}^\top \mathbb{E}[\nabla^2 g(\mathbf{w})] \mathbf{w}.$$

763 From (57), we obtain

$$\mathbf{w}^\top \nabla^2 g(\mathbf{w}) \mathbf{w} = \begin{cases} \frac{1}{\mu} \left[ (\mathbf{q}^\top \mathbf{x})^2 - \frac{x_n}{q_n} (\mathbf{q}^\top \mathbf{x}) - \frac{x_n}{q_n^3} (\mathbf{x}_{-n}^\top \mathbf{w}) \right], & |\mathbf{q}^\top \mathbf{x}| \leq \mu \\ -\frac{x_n}{q_n^3} \|\mathbf{w}\|^2 \text{sign}(\mathbf{q}^\top \mathbf{x}), & |\mathbf{q}^\top \mathbf{x}| \geq \mu. \end{cases}$$

764 Thus, we have

$$\begin{aligned}
\mathbb{E}[\mathbf{w}^\top \nabla^2 g(\mathbf{w}) \mathbf{w} \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \geq \mu}] &= -\frac{\|\mathbf{w}\|^2}{q_n^4} \mathbb{E}[q_n x_n \text{sign}(\mathbf{q}^\top \mathbf{x}) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \geq \mu}] \\
&= -\sqrt{\frac{2}{\pi}} \frac{\|\mathbf{w}\|^2}{q_n^2} \mathbb{E}_{\mathcal{I}} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|} \exp\left(-\frac{\mu^2}{2\|\mathbf{q}_{\mathcal{I}}\|^2}\right) \right]
\end{aligned}$$

765 and

$$\begin{aligned}
& \mathbb{E}[\mathbf{w}^\top \nabla^2 g(\mathbf{w}) \mathbf{w} \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}] \\
&= \frac{1}{\mu} \mathbb{E}\left[(\mathbf{q}^\top \mathbf{x})^2 \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}\right] - \frac{1}{\mu} \mathbb{E}\left[\frac{x_n}{q_n} (\mathbf{q}^\top \mathbf{x}) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}\right] - \frac{1}{\mu} \mathbb{E}\left[\frac{x_n}{q_n^3} (\mathbf{x}_{-n}^\top \mathbf{w}) \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}\right] \\
&= \frac{1}{\mu} \mathbb{E}[Z^2 \mathbb{1}_{|Z| \leq \mu}] - \frac{1}{\mu q_n^2} \mathbb{E}[Y^2 \mathbb{1}_{|X+Y| \leq \mu}] - \frac{1}{\mu} \left(\frac{1}{q_n^2} + \frac{1}{q_n^4}\right) \mathbb{E}[XY \mathbb{1}_{|X+Y| \leq \mu}],
\end{aligned}$$

<sup>766</sup> where  $X, Y$  and  $Z = X + Y$  are defined the same as (58). Similar to Lemma E.2, by using  
<sup>767</sup> Lemma D.7, we obtain

$$\begin{aligned} & \mathbb{E} [\mathbf{w}^\top \nabla^2 g(\mathbf{w}) \mathbf{w} \mathbb{1}_{|\mathbf{q}^\top \mathbf{x}| \leq \mu}] \\ &= -\sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{q}_{\mathcal{I}}\| \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] + \frac{1}{\mu} \mathbb{E} \left[ (\|\mathbf{q}_{\mathcal{I}}\|^2 - \mathbb{1}_{n \in \mathcal{I}}) \mathbb{P}(|\mathbf{q}_{\mathcal{I}}^\top \mathbf{g}| \leq \mu) \right] \\ & \quad + \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{q_n^2 \mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] + \sqrt{\frac{2}{\pi}} \left( 1 + \frac{1}{q_n^2} \right) \mathbb{E}_{\mathcal{I}} \left[ \frac{\|\mathbf{w}_{\mathcal{J}^c}\|^2 \mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right]. \end{aligned}$$

<sup>768</sup> Combining the results above and using integral by parts, we obtain

$$\begin{aligned} & \mathbf{w}^\top \nabla^2 \mathbb{E}[g(\mathbf{w})] \mathbf{w} \\ &= -\sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] + 2\sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] \\ & \quad - \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{q}_{\mathcal{I}}\| \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] + \frac{1}{\mu} \mathbb{E} \left[ (\|\mathbf{q}_{\mathcal{I}}\|^2 - \mathbb{1}_{n \in \mathcal{I}}) \mathbb{P}(|\mathbf{q}_{\mathcal{I}}^\top \mathbf{g}| \leq \mu) \right] \\ &= -\sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{\|\mathbf{w}_{\mathcal{J}^c}\|^2 \mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] \\ & \quad + \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|} \left( \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) - \frac{\|\mathbf{q}_{\mathcal{I}}\|}{\mu} \int_0^{\mu/\|\mathbf{q}_{\mathcal{I}}\|} t^2 \exp(-t^2/2) dt \right) \right] \\ & \quad - \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{q}_{\mathcal{I}}\| \left( \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) - \frac{\|\mathbf{q}_{\mathcal{I}}\|}{\mu} \int_0^{\mu/\|\mathbf{q}_{\mathcal{I}}\|} t^2 \exp(-t^2/2) dt \right) \right] \\ &= -\sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{w}_{\mathcal{J}^c}\|^2 \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] - \frac{1}{\mu} \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \mathbb{1}_{n \in \mathcal{I}} \int_0^{\mu/\|\mathbf{q}_{\mathcal{I}}\|} t^2 \exp(-t^2/2) dt \right] \\ & \quad + \frac{1}{\mu} \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{q}_{\mathcal{I}}\|^2 \int_0^{\mu/\|\mathbf{q}_{\mathcal{I}}\|} t^2 \exp(-t^2/2) dt \right] \\ &\leq -\sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{w}_{\mathcal{J}^c}\|^2 \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] + \frac{1}{\mu} \sqrt{\frac{2}{\pi}} \int_0^\mu t^2 \mathbb{E}_{\mathcal{I}} \left[ \frac{1}{\|\mathbf{q}_{\mathcal{I}}\|} \exp \left( -\frac{t^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] dt. \end{aligned}$$

<sup>769</sup> First, when  $\sqrt{\frac{n-1}{n}} \geq \|\mathbf{w}\| \geq c_0$ , we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{I}} \left[ \|\mathbf{w}_{\mathcal{J}^c}\|^2 \frac{\mathbb{1}_{n \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] \\ &= \theta \mathbb{E}_{\mathcal{J}} \left[ \|\mathbf{w}_{\mathcal{J}^c}\|^2 \frac{1}{(q_n^2 + \|\mathbf{w}_{\mathcal{J}}\|^2)^{3/2}} \exp \left( -\frac{\mu^2}{2(q_n^2 + \|\mathbf{w}_{\mathcal{J}}\|^2)} \right) \right] \\ &\geq \theta \mathbb{E}_{\mathcal{J}} \left[ \|\mathbf{w}_{\mathcal{J}^c}\|^2 \exp \left( -\frac{\mu^2}{2q_n^2 + 2\|\mathbf{w}_{\mathcal{J}}\|^2} \right) \right] \\ &\geq \theta \mathbb{E}_{\mathcal{J}} \left[ \|\mathbf{w}_{\mathcal{J}^c}\|^2 \exp \left( -\frac{\mu^2}{2q_n^2} \right) \right] \geq c_1 \theta (1-\theta) \|\mathbf{w}\|^2. \end{aligned}$$

<sup>770</sup> Second, notice that the function

$$h(x) = x^{-1} \exp \left( -\frac{t^2}{2x^2} \right), \quad x \in [0, 1]$$

771 reaches the maximum when  $x = t$ . Thus, we have

$$\frac{1}{\mu} \sqrt{\frac{2}{\pi}} \int_0^\mu t^2 \mathbb{E}_{\mathcal{I}} \left[ \frac{1}{\|q_{\mathcal{I}}\|} \exp \left( -\frac{t^2}{2 \|q_{\mathcal{I}}\|^2} \right) \right] dt \leqslant \frac{1}{\mu} \sqrt{\frac{2}{\pi}} \int_0^\mu t \exp \left( -\frac{1}{2} \right) dt \leqslant \frac{1}{\sqrt{2\pi}} e^{-1/2} \mu.$$

772 Therefore, when  $\mu \leqslant \frac{1}{n} \leqslant \theta$ , we have

$$\mathbf{w}^\top \nabla^2 \mathbb{E}[g(\mathbf{w})] \mathbf{w} \leqslant -c_2 \theta(1-\theta) \|\mathbf{w}\|^2$$

773 for any  $\sqrt{\frac{n-1}{n}} \geqslant \|\mathbf{w}\| \geqslant c_0$ .

774  $\blacksquare$

## 775 F Negative Curvature on Gradient in Population

776 **Proposition F.1** Suppose  $\theta \geqslant \frac{1}{n}$ . Given any index  $i \in [n]$ , when  $\mu \leqslant \frac{1}{\sqrt{3n}}$ , we have

$$\left\langle \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_i} \mathbf{e}_i \right\rangle \geqslant \frac{\theta(1-\theta)}{4n} \frac{\xi}{1+\xi},$$

777 holds for all  $\mathbf{q} \in \mathcal{S}_\xi^{i+}$  and any  $q_j$  such that  $j \neq i$  and  $q_j^2 \geqslant \frac{1}{3} q_i^2$

778 **Proof** Without loss of generality, let us consider the case  $i = n$ . For any  $j \neq n$ , we have

$$\begin{aligned} & \left\langle \text{grad } \mathbb{E}[\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \\ &= \left( \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right)^\top \mathcal{P}_{\mathbf{q}^\perp} \mathbb{E}[\mathbf{x} \cdot h'_\mu(\mathbf{x}^\top \mathbf{q})] \\ &= \left( \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right)^\top \mathbb{E}[\mathbf{x} \cdot h'_\mu(\mathbf{x}^\top \mathbf{q})]. \end{aligned}$$

779 Let

$$Z = Z_1 + Z_2, \quad Z_1 = q_i x_i \sim \mathcal{N}(0, (b_i q_i)^2), \quad Z_2 = \mathbf{q}_{-i}^\top \mathbf{x}_{-i} \sim \mathcal{N}(0, \|\mathbf{q}_{-i} \odot \mathbf{b}_{-i}\|^2).$$

780 Notice that for every  $i \in [n]$ , we have

$$\begin{aligned} & \frac{1}{q_i} \mathbf{e}_i^\top \mathbb{E}[\mathbf{x} \cdot h'_\mu(\mathbf{x}^\top \mathbf{q})] \\ &= \frac{1}{q_i^2} \frac{1}{\mu} \mathbb{E}[Z_1^2 \mathbb{1}_{|Z_1+Z_2| \leqslant \mu}] + \frac{1}{q_i^2} \frac{1}{\mu} \mathbb{E}[Z_1 Z_2 \mathbb{1}_{|Z_1+Z_2| \leqslant \mu}] + \frac{1}{q_i^2} \mathbb{E}[Z_1 \text{sign}(Z_1+Z_2) \mathbb{1}_{|Z_1+Z_2| \geqslant \mu}]. \end{aligned}$$

781 By Lemma D.7, we have

$$\begin{aligned} \mathbb{E}[Z_1^2 \mathbb{1}_{|Z_1+Z_2| \leqslant \mu}] &= -\sqrt{\frac{2}{\pi}} \mu \mathbb{E}_{\mathcal{I}} \left[ \frac{q_i^4 \mathbb{1}_{i \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] \\ &\quad + \mathbb{E}[q_i^2 \mathbb{1}_{i \in \mathcal{I}} \mathbb{P}(|Z| \leqslant \mu)], \\ \mathbb{E}[Z_1 Z_2 \mathbb{1}_{|Z_1+Z_2| \leqslant \mu}] &= -\sqrt{\frac{2}{\pi}} \mu \mathbb{E}_{\mathcal{I}} \left[ \frac{q_i^2 \mathbb{1}_{i \in \mathcal{I}} \|\mathbf{q}_{-i}\|^2}{\|\mathbf{q}_{\mathcal{I}}\|^3} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right] \\ \mathbb{E}[Z_1 \text{sign}(Z_1+Z_2) \mathbb{1}_{|Z_1+Z_2| \geqslant \mu}] &= \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{q_i^2 \mathbb{1}_{i \in \mathcal{I}}}{\|\mathbf{q}_{\mathcal{I}}\|^2} \exp \left( -\frac{\mu^2}{2 \|\mathbf{q}_{\mathcal{I}}\|^2} \right) \right]. \end{aligned}$$

782 Combining the results above, we obtain

$$\frac{1}{q_i} \mathbf{e}_i^\top \mathbb{E}[\mathbf{x} \cdot h'_\mu(\mathbf{x}^\top \mathbf{q})] = \frac{1}{\mu} \mathbb{E}[\mathbb{1}_{i \in \mathcal{I}} \mathbb{P}(|Z| \leqslant \mu)].$$

<sup>783</sup> Therefore, we have

$$\begin{aligned}
& \left\langle \text{grad } \mathbb{E} [\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \\
&= \frac{1}{\mu} (\mathbb{E} [\mathbb{1}_{j \in \mathcal{I}} \mathbb{P} (|Z| \leq \mu)] - \mathbb{E} [\mathbb{1}_{n \in \mathcal{I}} \mathbb{P} (|Z| \leq \mu)]) \\
&= \frac{\theta}{\mu} \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{1}{\sqrt{q_j^2 + \|\mathbf{q}_{\mathcal{I} \setminus j}\|^2}} \int_0^\mu \exp \left( -\frac{t^2}{q_j^2 + \|\mathbf{q}_{\mathcal{I} \setminus j}\|^2} \right) dt \right] \\
&\quad - \frac{\theta}{\mu} \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{1}{\sqrt{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus n}\|^2}} \int_0^\mu \exp \left( -\frac{t^2}{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus n}\|^2} \right) dt \right] \\
&= \frac{\theta(1-\theta)}{\mu} \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{1}{\sqrt{q_j^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \int_0^\mu \exp \left( -\frac{t^2}{q_j^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2} \right) dt \right] \\
&\quad - \frac{\theta(1-\theta)}{\mu} \sqrt{\frac{2}{\pi}} \mathbb{E}_{\mathcal{I}} \left[ \frac{1}{\sqrt{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \int_0^\mu \exp \left( -\frac{t^2}{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2} \right) dt \right] \\
&= \frac{\theta(1-\theta)}{\mu} \mathbb{E}_{\mathcal{I}} \left[ \text{erf} \left( \frac{\mu}{\sqrt{q_i^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \right) - \text{erf} \left( \frac{\mu}{\sqrt{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \right) \right]
\end{aligned}$$

<sup>784</sup> where  $\text{erf}(x)$  is the Gaussian error function

$$\text{erf}(x) = \frac{1}{\sqrt{2\pi}} \int_{-x}^x \exp(-t^2/2) dt = \sqrt{\frac{2}{2\pi}} \int_0^x \exp(-t^2/2) dt, \quad x \geq 0.$$

<sup>785</sup> When  $\mu \leq \frac{1}{\sqrt{3n}}$  such that  $\frac{\mu}{\sqrt{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \leq 1$  for  $\mathbf{q} \in \mathcal{S}_\xi^{n+}$ , by Taylor approximation we have

$$\begin{aligned}
& \text{erf} \left( \frac{\mu}{\sqrt{q_i^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \right) - \text{erf} \left( \frac{\mu}{\sqrt{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \right) \\
& \geq \frac{\mu}{2} \left[ \frac{1}{\sqrt{q_i^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} - \frac{1}{\sqrt{q_n^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2}} \right] = \frac{\mu}{4} \int_{q_i^2}^{q_n^2} \frac{1}{(t^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2)^{3/2}} dt.
\end{aligned}$$

<sup>786</sup> Therefore, we have

$$\begin{aligned}
& \left\langle \text{grad } \mathbb{E} [\tilde{f}(\mathbf{q})], \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \\
& \geq \frac{\theta(1-\theta)}{4} \int_{q_i^2}^{q_n^2} \frac{1}{(t^2 + \|\mathbf{q}_{\mathcal{I} \setminus \{j,n\}}\|^2)^{3/2}} dt \\
& \geq \frac{\theta(1-\theta)}{4} (q_n^2 - \|\mathbf{q}_{-n}\|_\infty^2) \geq \frac{\theta(1-\theta)}{4} \frac{\xi}{1+\xi} q_n^2 \geq \frac{\theta(1-\theta)}{4n} \frac{\xi}{1+\xi}.
\end{aligned}$$

<sup>787</sup> This gives the desired result. ■

## <sup>788</sup> G Gradient Concentration

<sup>789</sup> In this section, we uniformly bound the deviation between the empirical process  $\text{grad } \tilde{f}(\mathbf{q})$  and its  
<sup>790</sup> mean  $\mathbb{E} [\text{grad } \tilde{f}(\mathbf{q})]$  over the sphere. Namely, we show the following

791 **Proposition G.1** For every  $i \in [n]$  and any  $\delta \in (0, 1)$ , when

$$p \geq C\delta^{-2}n \log\left(\frac{\theta n}{\mu\delta}\right), \quad (62)$$

792 we have

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left| \left\langle \text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E} \left[ \text{grad } \tilde{f}(\mathbf{q}) \right], \mathbf{e}_i \right\rangle \right| \leq \delta$$

793 holds with probability at least  $1 - np^{-c_1\theta n} - n \exp(-c_2 p \delta^2)$ , for any  $\mathbf{e}_i$ . Here,  $c_1$ ,  $c_2$ , and  $C$  are  
794 some universal positive numerical constants.

795 **Remarks.** Here, our bound is loose by roughly a factor of  $n$  because of the looseness in handling  
796 the probabilistic dependency due to the convolution measurement. We believe this bound can be  
797 improved by an order of  $\mathcal{O}(n)$  using more advanced probability tools, such as decoupling and  
798 chaining [44–46].

799 **Proof** First, note that

$$\tilde{f}(\mathbf{q}) = \frac{1}{np} \sum_{i=1}^p H_\mu(C_{\mathbf{x}_i} \mathbf{q}), \quad \text{grad } \tilde{f}(\mathbf{q}) = \frac{1}{np} \mathcal{P}_{\mathbf{q}^\perp} \sum_{i=1}^p C_{\mathbf{x}_i}^\top h'_\mu(C_{\mathbf{x}_i} \mathbf{q}). \quad (63)$$

800 Thus, we have

$$\begin{aligned} & \left\langle \text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E} \left[ \text{grad } \tilde{f}(\mathbf{q}) \right], \mathbf{e}_n \right\rangle \\ &= \frac{1}{np} \sum_{i=1}^p \sum_{j=0}^{n-1} \left[ \left\langle \mathcal{P}_{\mathbf{q}^\perp} s_j [\tilde{\mathbf{x}}_i], \mathbf{e}_n \right\rangle h'_\mu(s_j [\tilde{\mathbf{x}}_i]^\top \mathbf{q}) - \mathbb{E} \left[ (\mathbf{e}_n^\top \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}) h'_\mu(\mathbf{x}^\top \mathbf{q}) \right] \right]. \end{aligned}$$

801 This is a summation of dependent random variables, which is very difficult to show measurement  
802 concentration in general. We alleviate this difficulty by only considering a partial summation of  
803 independent random variables, namely,

$$\mathcal{L}(\mathbf{q}) = \frac{1}{p} \frac{1}{\|\mathcal{P}_{\mathbf{q}^\perp} \mathbf{e}_n\|} \sum_{i=1}^p \left[ \left\langle \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}_i, \mathbf{e}_n \right\rangle h'_\mu(\mathbf{x}_i^\top \mathbf{q}) - \mathbb{E} \left[ (\mathbf{e}_n^\top \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}) h'_\mu(\mathbf{x}^\top \mathbf{q}) \right] \right],$$

804 where  $\mathbf{x}_i \sim_{i.i.d.} \mathcal{BG}(\theta)$ . Note that the bound of  $\mathcal{L}(\mathbf{q})$  automatically gives an upper bound of  
805  $\left\langle \text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E} \left[ \text{grad } \tilde{f}(\mathbf{q}) \right], \mathbf{e}_n \right\rangle$  in distribution. To uniformly control  $\mathcal{L}(\mathbf{q})$  over the sphere, we  
806 first consider controlling  $\mathcal{L}(\mathbf{q})$  for a fixed  $\mathbf{q} \in \mathbb{S}^{n-1}$ . For each  $\ell = 1, 2, \dots$ , we have the moments

$$\mathbb{E} \left[ \left| \left\langle \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}_i, \mathbf{e}_n \right\rangle h'_\mu(\mathbf{x}_i^\top \mathbf{q}) \right|^\ell \right] \leq \mathbb{E} \left[ |\mathbf{e}_n^\top \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}_i|^\ell \right] = \mathbb{E} \left[ |Z_i|^\ell \right],$$

807 where conditioned on the Bernoulli distribution, we have  $Z_i \sim \mathcal{N} \left( 0, \left\| (\mathcal{P}_{\mathbf{q}^\perp} \mathbf{e}_n)_{\mathcal{J}} \right\|^2 \right)$ . By Lemma  
808 D.1, we have

$$\mathbb{E} \left[ \left| \left\langle \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}_i, \mathbf{e}_n \right\rangle h'_\mu(\mathbf{x}_i^\top \mathbf{q}) \right|^\ell \right] \leq \mathbb{E}_{\mathcal{J}} \left[ (\ell-1)!! \left\| (\mathcal{P}_{\mathbf{q}^\perp} \mathbf{e}_n)_{\mathcal{J}} \right\|^\ell \right] \leq \frac{\ell!}{2} \|\mathcal{P}_{\mathbf{q}^\perp} \mathbf{e}_n\|^\ell,$$

809 where we used the fact that  $|h'_\mu(z)| \leq 1$  for any  $z$ . Thus, we are controlling the concentration of  
810 summation of sub-Gaussian r.v., for which we have

$$\mathbb{P}(|\mathcal{L}(\mathbf{q})| \geq t) \leq \exp \left( -C \frac{pt^2}{2} \right).$$

811 Next, we turn this point-wise concentration into a uniform bound for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  via a standard  
812 covering argument. Let  $\mathcal{N}(\varepsilon)$  be an  $\varepsilon$ -net of the sphere, whose cardinality can be controlled by

$$|\mathcal{N}(\varepsilon)| \leq \left( \frac{3}{\varepsilon} \right)^{n-1}.$$

813 Thus, we have

$$\mathbb{P} \left( \sup_{\mathbf{q} \in \mathcal{N}(\varepsilon)} |\mathcal{L}(\mathbf{q})| \geq t \right) \leq \left( \frac{3}{\varepsilon} \right)^{n-1} \exp \left( -\frac{pt^2}{2+2t} \right).$$

814 For any point  $\mathbf{q} \in \mathbb{S}^{n-1}$ , it can written as  $\mathbf{q} = \mathbf{q}' + \mathbf{e}$ , where  $\mathbf{q}' \in \mathcal{N}(\varepsilon)$  and  $\|\mathbf{e}\| \leq \varepsilon$ . Now we control  
815 the all points over the sphere through the Lipschitz property of  $\mathcal{L}$ .

$$\begin{aligned} & \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\mathcal{L}(\mathbf{q})| \\ &= \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} |\mathcal{L}(\mathbf{q}' + \mathbf{e})| \\ &\leq \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon)} |\mathcal{L}(\mathbf{q}')| + \underbrace{\sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} |\mathbb{E} [(\mathbf{e}_n^\top \mathcal{P}_{(\mathbf{q}'+\mathbf{e})^\perp} \mathbf{x} - \mathbf{e}_n^\top \mathcal{P}_{(\mathbf{q}')^\perp} \mathbf{x}) h'_\mu(\mathbf{x}^\top \mathbf{q}')]|}_{\mathcal{L}_1} \\ &\quad + \underbrace{\sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} |\mathbb{E} [(\mathbf{e}_n^\top \mathcal{P}_{(\mathbf{q}'+\mathbf{e})^\perp} \mathbf{x}) (h'_\mu(\mathbf{x}^\top (\mathbf{q}' + \mathbf{e})) - h'_\mu(\mathbf{x}^\top \mathbf{q}'))]|}_{\mathcal{L}_2} \\ &\quad + \underbrace{\sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} \left| \frac{1}{p} \sum_{i=1}^p [\mathbf{e}_n^\top \mathcal{P}_{(\mathbf{q}'+\mathbf{e})^\perp} \mathbf{x}_i - \mathbf{e}_n^\top \mathcal{P}_{(\mathbf{q}')^\perp} \mathbf{x}_i] h'_\mu(\mathbf{x}_i^\top \mathbf{q}') \right|}_{\mathcal{L}_3} \\ &\quad + \underbrace{\sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} \left| \frac{1}{p} \sum_{i=1}^p (\mathbf{e}_n^\top \mathcal{P}_{(\mathbf{q}'+\mathbf{e})^\perp} \mathbf{x}_i) [h'_\mu(\mathbf{x}_i^\top (\mathbf{q}' + \mathbf{e})) - h'_\mu(\mathbf{x}_i^\top \mathbf{q}')] \right|}_{\mathcal{L}_4}. \end{aligned}$$

816 By Lipschitz continuity and the fact that  $h'_\mu(z) \leq 1$  for any  $z$ , we obtain

$$\begin{aligned} \mathcal{L}_1 &\leq \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} \sqrt{\theta} \|(\mathcal{P}_{(\mathbf{q}'+\mathbf{e})^\perp} - \mathcal{P}_{(\mathbf{q}')^\perp}) \mathbf{e}_n\| \leq 3\sqrt{\theta}\varepsilon \\ \mathcal{L}_2 &\leq \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} \frac{1}{\mu} \mathbb{E} [\|\mathbf{x}\| \|\mathbf{x}^\top \mathbf{e}\|] \leq \frac{\theta n}{\mu} \varepsilon. \end{aligned}$$

817 For each  $\mathbf{x}_i$ , we know that  $\mathbf{x}_i = \mathbf{g}_i \odot \mathbf{b}_i$  with  $\mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $\mathbf{b}_i \sim i.i.d. \mathcal{B}(\theta)$ . By Gaussian  
818 concentration inequality, we know that for each  $\mathbf{x}_i$ ,

$$\mathbb{P} (\|\mathbf{x}_i\| - \sqrt{\theta n} \geq t) \leq \mathbb{P} (\|\mathbf{x}_i\| - \mathbb{E} [\|\mathbf{x}_i\|] \geq t) \leq \exp \left( -\frac{t^2}{2 \|\mathbf{b}_i\|_\infty} \right) \leq \exp \left( -\frac{t^2}{2} \right).$$

819 Therefore, by a union bound, we have

$$\max_{1 \leq i \leq p} \|\mathbf{x}_i\| \leq 5\sqrt{\theta n \log p}$$

820 holds with probability at least  $1 - p^{-8\theta n}$ . Therefore, w.h.p we have

$$\begin{aligned} \mathcal{L}_3 &\leq \left( \max_{1 \leq i \leq p} \|\mathbf{x}_i\| \right) \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} \|\mathcal{P}_{(\mathbf{q}'+\mathbf{e})^\perp} - \mathcal{P}_{(\mathbf{q}')^\perp}\| \leq 15\sqrt{\theta n \log p} \varepsilon, \\ \mathcal{L}_4 &\leq \frac{1}{\mu} \left( \max_{1 \leq i \leq p} \|\mathbf{x}_i\|^2 \right) \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon), \|\mathbf{e}\| \leq \varepsilon} \|\mathbf{e}\| \leq 25 \frac{\theta n \log p}{\mu} \varepsilon. \end{aligned}$$

821 Combining the bounds above, choose  $\varepsilon = \frac{\mu t}{c\theta n \log p}$ , we have

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\mathcal{L}(\mathbf{q})| \leq \sup_{\mathbf{q}' \in \mathcal{N}(\varepsilon)} |\mathcal{L}(\mathbf{q}')| + c \frac{\theta n \log p}{\mu} \varepsilon \leq 2t$$

822 holds with probability at least

$$1 - p^{-8\theta n} - \exp \left( -C \frac{pt^2}{2} + c'n \log \left( \frac{\theta n}{\mu t} \right) \right).$$

823 Thus, applying a union bound, we obtain the desired result holding for every  $i \in [n]$ . ■

824 Similarly, we also show the following result.

825 **Corollary G.2** *For any  $\delta \in (0, 1)$ , when*

$$p \geq C\delta^{-2}n^2 \log\left(\frac{\theta n}{\mu\delta}\right), \quad (64)$$

826 we have

$$\begin{aligned} \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})]\| &\leq \delta, \\ \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\nabla \tilde{f}(\mathbf{q}) - \mathbb{E}[\nabla \tilde{f}(\mathbf{q})]\| &\leq \delta, \end{aligned}$$

827 hold with probability at least  $1 - p^{-c_1\theta n} - n \exp(-c_2 p \delta^2)$ . Here,  $c_1$ ,  $c_2$ , and  $C$  are some universal  
828 positive numerical constants.

829 **Proof** From Proposition G.1, we know that when  $p \geq C_0 \varepsilon^{-2} n \log\left(\frac{\theta n}{\mu\varepsilon}\right)$ ,

$$\begin{aligned} &\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})]\|^2 \\ &\leq \sum_{i=1}^n \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})], \mathbf{e}_i \rangle|^2 \leq n\varepsilon^2. \end{aligned}$$

830 holds with probability at least  $1 - p^{-c_1\theta n} - n \exp(-c_2 p \delta^2)$ . Therefore, by letting  $\delta = \sqrt{n}\varepsilon$ , w.h.p.  
831 we have

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})]\| \leq \delta,$$

832 whenever  $p \geq C\delta^{-2}n^2 \log\left(\frac{\theta n}{\mu\delta}\right)$ . By a similar argument, we can also provide the same bound for  
833  $\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\nabla \tilde{f}(\mathbf{q}) - \mathbb{E}[\nabla \tilde{f}(\mathbf{q})]\|$ . ■

834 **Corollary G.3** *For each  $i \in [n]$  and any  $\delta \in (0, 1)$ , when  $p \geq C\delta^{-2}n \log\left(\frac{\theta n}{\mu\delta}\right)$ , we have*

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| \leq 1 + \delta,$$

835 hold with probability at least  $1 - np^{-c_1\theta n} - n \exp(-c_2 p \delta^2)$ . Here,  $c_1$ ,  $c_2$ , and  $C$  are some universal  
836 positive numerical constants.

837 **Proof** For any  $\mathbf{q} \in \mathbb{S}^{n-1}$  and every  $i \in [n]$ , we have

$$\mathbb{E}[|\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle|] = \mathbb{E}[|(\mathbf{e}_i^\top \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}) \cdot h'_\mu(\mathbf{x}^\top \mathbf{q})|] \leq \mathbb{E}[\|\mathbf{e}_i^\top \mathcal{P}_{\mathbf{q}^\perp} \mathbf{x}\|] \leq 1.$$

838 Thus, we have

$$\begin{aligned} &\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}) - \mathbb{E}[\text{grad } \tilde{f}(\mathbf{q})], \mathbf{e}_i \rangle| \\ &\geq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} (|\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| - \mathbb{E}[|\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle|]) \\ &\geq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle| - \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \mathbb{E}[|\langle \text{grad } \tilde{f}(\mathbf{q}), \mathbf{e}_i \rangle|]. \end{aligned}$$

839 Therefore, by using the result in Proposition G.1, we obtain the desired result. ■

840 **Corollary G.4** For any  $\delta \in (0, 1)$ , when  $p$  satisfies (64), we have

$$\sup_{q \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(q)\| \leq \sqrt{\theta n} + \delta,$$

841 hold with probability at least  $1 - p^{-c_1\theta n} - n \exp(-c_2 p \delta^2)$ . Here,  $c_1$ ,  $c_2$ , and  $C$  are some universal  
842 positive numerical constants.

843 **Proof** For any  $q \in \mathbb{S}^{n-1}$ , we have

$$\mathbb{E} [\|\text{grad } \tilde{f}(q)\|] = \mathbb{E} [\|\mathcal{P}_{q^\perp} \mathbf{x} h'_\mu(\mathbf{x}^\top q)\|] \leq \mathbb{E} [\|\mathbf{x}\|] \leq \sqrt{\theta n}.$$

844 Note that

$$\begin{aligned} \sup_{q \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(q) - \mathbb{E} [\text{grad } \tilde{f}(q)]\| &\geq \sup_{q \in \mathbb{S}^{n-1}} (\|\text{grad } \tilde{f}(q)\| - \mathbb{E} [\|\text{grad } \tilde{f}(q)\|]) \\ &\geq \sup_{q \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(q)\| - \sup_{q \in \mathbb{S}^{n-1}} \mathbb{E} [\|\text{grad } \tilde{f}(q)\|]. \end{aligned}$$

845 Thus, by using the result in Corollary G.2, we obtain the desired result.  $\blacksquare$

## 846 H Preconditioning

847 In this section, given the Riemannian gradient of  $\tilde{f}(q)$  and its preconditioned variant

$$\begin{aligned} \text{grad } \tilde{f}(q) &= \frac{1}{np} \mathcal{P}_{q^\perp} \sum_{i=1}^p \mathbf{C}_{x_i}^\top h'_\mu(\mathbf{C}_{x_i} q), \\ \text{grad } f(q) &= \frac{1}{np} \mathcal{P}_{q^\perp} \sum_{i=1}^p (\mathbf{R} \mathbf{Q}^{-1})^\top \mathbf{C}_{x_i}^\top h'_\mu(\mathbf{C}_{x_i} (\mathbf{R} \mathbf{Q}^{-1}) q), \end{aligned}$$

848 we prove the following result.

849 **Proposition H.1** Suppose  $\theta \geq \frac{1}{n}$ . For any  $\delta \in (0, 1)$ , whenever

$$p \geq C \frac{\kappa^8 n}{\mu^2 \theta \delta^2 \sigma_{\min}^2} \log^4 n \log \left( \frac{\theta n}{\mu} \right),$$

850 we have

$$\sup_{q \in \mathbb{S}^{n-1}} \|\text{grad } \tilde{f}(q) - \text{grad } f(q)\| \leq \delta$$

851 holds with probability at least  $1 - c_1 p^{-c_2 n \theta} - n^{-c_3} - n e^{-c_4 \theta np}$ . Here,  $\kappa$  and  $\sigma_{\min}$  denote the  
852 condition number and minimum singular value of  $\mathbf{C}_a$ , and  $c_1$ ,  $c_2$ ,  $c_3$ ,  $c_4$  and  $C$  are some positive  
853 numerical constants.

854 **Proof** Notice that

$$\mathbf{R} = \mathbf{C}_a \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{y_i}^\top \mathbf{C}_{y_i} \right)^{-1/2}, \quad \mathbf{Q} = \mathbf{C}_a (\mathbf{C}_a^\top \mathbf{C}_a)^{-1/2}$$

855 so that

$$\mathbf{R} \mathbf{Q}^{-1} = \mathbf{C}_a \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{y_i}^\top \mathbf{C}_{y_i} \right)^{-1/2} (\mathbf{C}_a^\top \mathbf{C}_a)^{1/2} \mathbf{C}_a^{-1}.$$

856 Thus, we have

$$\begin{aligned}
& \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left\| \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } f(\mathbf{q}) \right\| \\
& \leq \frac{1}{np} \left\| \mathcal{P}_{\mathbf{q}^\perp} (\mathbf{I} - (\mathbf{RQ}^{-1}))^\top \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top h'_\mu(\mathbf{C}_{\mathbf{x}_i} \mathbf{q}) \right\| \\
& \quad + \frac{1}{np} \left\| \mathcal{P}_{\mathbf{q}^\perp} (\mathbf{RQ}^{-1})^\top \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top [h'_\mu(\mathbf{C}_{\mathbf{x}_i} \mathbf{q}) - h'_\mu(\mathbf{C}_{\mathbf{x}_i} (\mathbf{RQ}^{-1}) \mathbf{q})] \right\| \\
& \leq \|\mathbf{I} - \mathbf{RQ}^{-1}\| \left\| \nabla \tilde{f}(\mathbf{q}) \right\| + \|\mathbf{RQ}^{-1}\| \left\| \frac{1}{np} \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top [h'_\mu(\mathbf{C}_{\mathbf{x}_i} \mathbf{q}) - h'_\mu(\mathbf{C}_{\mathbf{x}_i} (\mathbf{RQ}^{-1}) \mathbf{q})] \right\| \\
& \leq \|\mathbf{I} - \mathbf{RQ}^{-1}\| \left\| \nabla \tilde{f}(\mathbf{q}) \right\| + \frac{1}{\mu\sqrt{n}} \|\mathbf{RQ}^{-1}\| \left( \max_{1 \leq i \leq p} \|\mathbf{x}_i\| \|\mathbf{F}\mathbf{x}_i\|_\infty \right) \|\mathbf{I} - \mathbf{RQ}^{-1}\|. \tag{65}
\end{aligned}$$

857 Here, by Lemma H.4, for any given  $\varepsilon \in (0, 1)$ , when  $p \geq C \frac{\kappa^8}{\theta\varepsilon^2\sigma_{\min}^2(\mathbf{C}_\alpha)} \log^3 n$ , we have

$$\|\mathbf{RQ}^{-1} - \mathbf{I}\| \leq \varepsilon, \quad \|\mathbf{RQ}^{-1}\| \leq 1 + \varepsilon, \tag{66}$$

858 holding with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ . On the other hand, by Gaussian concentration  
859 inequality and a union bound, we have

$$\max_{1 \leq i \leq p} \|\mathbf{x}_i\| \leq 4\sqrt{n \log p}, \quad \max_{1 \leq i \leq p} \|\mathbf{F}\mathbf{x}_i\|_\infty \leq 4\sqrt{n \log p}, \tag{67}$$

860 hold with probability at least  $1 - p^{-c_3 n}$ . By Corollary G.4, when  $p \geq C_2 \theta^{-1} n \log \left( \frac{\theta n}{\mu} \right)$ , we have

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left\| \text{grad } \tilde{f}(\mathbf{q}) \right\| \leq 2\sqrt{\theta n} \tag{68}$$

861 holds with probability at least  $1 - p^{-c_4 \theta n} - n e^{-c_5 \theta np}$ . Plugging the bounds in (66) and (67) into  
862 (65), we obtain

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left\| \text{grad } \tilde{f}(\mathbf{q}) - \text{grad } f(\mathbf{q}) \right\| \leq \varepsilon \left[ 2\sqrt{\theta n} + \frac{16\sqrt{n} \log p}{\mu} \cdot (1 + \varepsilon) \right].$$

863 By a change of variable, we obtain the desired result. ■

864 **Lemma H.2** When  $\theta \geq 1/n$ ,

$$\left\| \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top \mathbf{C}_{\mathbf{x}_i} - \mathbf{I} \right\| \leq t \tag{69}$$

865 holds with probability at least  $1 - p^{-c_1 n \theta} - n \exp \left( -c_2 \min \left\{ \frac{pt^2}{\theta \log p}, \frac{pt}{\sqrt{\theta \log p}} \right\} \right)$  for some numerical  
866 constants  $c_1, c_2 > 0$ .

867 **Proof**

868 Notice that

$$\mathbf{C}_{\mathbf{x}_i}^\top \mathbf{C}_{\mathbf{x}_i} = \mathbf{F}^* \text{diag} \left( |\mathbf{F}\mathbf{x}_i|^{\odot 2} \right) \mathbf{F}.$$

869 Then

$$\begin{aligned}
\left\| \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top \mathbf{C}_{\mathbf{x}_i} - \mathbf{I} \right\| &= \left\| \mathbf{F}^* \left( \text{diag} \left( \frac{1}{\theta np} \sum_{i=1}^p |\mathbf{F}\mathbf{x}_i|^{\odot 2} \right) - \mathbf{F}^{-1} (\mathbf{F}^*)^{-1} \right) \mathbf{F} \right\| \\
&= \left\| \frac{1}{\theta np} \sum_{i=1}^p |\mathbf{F}\mathbf{x}_i|^{\odot 2} - \mathbf{1} \right\|_\infty. \tag{70}
\end{aligned}$$

870 Let  $\mathbf{x}_i = \mathbf{b}_i \odot \mathbf{g}_i$  with  $\mathbf{b}_i \sim_{i.i.d.} \mathcal{B}(\theta)$  and  $\mathbf{g}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ , and let us define events

$$\mathcal{E}_{i,j} \doteq \left\{ \|\mathbf{b}_i \odot \mathbf{f}_j\|^2 \leq 5n\sqrt{\theta \log p} \right\}, \quad 1 \leq i \leq p, 1 \leq j \leq n.$$

871 We use  $\mathcal{E}_j = \bigcap_{i=1}^p \mathcal{E}_{i,j}$ . For each individual  $i$  and  $j$ , by the Hoeffding's inequality, we have

$$\mathbb{P}(\mathcal{E}_{i,j}^c) \leq \exp(-8n\theta \log p)$$

872 For each  $j = 1, \dots, n$ , by conditional probability and union bound, we have

$$\begin{aligned} \mathbb{P}\left(\left|\frac{1}{\theta np} \sum_{i=1}^p |\mathbf{f}_j^* \mathbf{x}_i|^2 - 1\right| \geq t\right) &\leq \mathbb{P}\left(\bigcup_{i=1}^p \mathcal{E}_{i,j}^c\right) + \mathbb{P}\left(\left|\frac{1}{\theta np} \sum_{i=1}^p |\mathbf{f}_j^* \mathbf{x}_i|^2 - 1\right| \geq t \mid \mathcal{E}_j\right) \\ &\leq \sum_{i=1}^p \mathbb{P}(\mathcal{E}_{i,j}^c) + \mathbb{P}\left(\left|\frac{1}{\theta np} \sum_{i=1}^p |\mathbf{f}_j^* \mathbf{x}_i|^2 - 1\right| \geq t \mid \mathcal{E}_j\right) \\ &\leq pe^{-8n\theta \log p} + \mathbb{P}\left(\left|\frac{1}{\theta np} \sum_{i=1}^p |\mathbf{f}_j^* \mathbf{x}_i|^2 - 1\right| \geq t \mid \mathcal{E}_j\right). \end{aligned} \quad (71)$$

873 For the second term, since  $\mathbf{x}_i \sim \mathcal{BG}(\theta)$ , we have

$$\mathbf{f}_j^* \mathbf{x}_i = \sum_{k=1}^n f_{ji} b_{ik} g_{ik} \sim \mathcal{N}(0, \|\mathbf{b}_i \odot \mathbf{f}_j\|^2)$$

874 for all  $\ell \geq 1$ , by Lemma D.1, we have

$$\begin{aligned} \mathbb{E}\left[(\theta n)^{-\ell} |\mathbf{f}_j^* \mathbf{x}_i|^{2\ell} \mid \mathcal{E}_{i,j}\right] &= \frac{(2\ell-1)!!}{(\theta n)^\ell} \mathbb{E}\left[\|\mathbf{b}_i \odot \mathbf{f}_j\|^{2\ell} \mid \mathcal{E}_{i,j}\right] \\ &\leq \frac{\ell!}{2} 10^\ell \theta^{-\ell/2} \log^{\ell/2} p. \end{aligned}$$

875 Thus, by Bernstein inequality in Lemma D.3, we have

$$\begin{aligned} \mathbb{P}\left(\left|\frac{1}{\theta np} \sum_{i=1}^p |\mathbf{f}_j^* \mathbf{x}_i|^2 - 1\right| \geq t \mid \mathcal{E}_j\right) &\leq \exp\left(-\frac{pt^2}{200 \log p + 20\sqrt{\log pt}}\right) \\ &\leq \exp\left(-\min\left\{\frac{pt^2}{400\theta \log p}, \frac{pt}{40\sqrt{\theta \log p}}\right\}\right). \end{aligned} \quad (72)$$

876 Plugging (72) into (71), we obtain

$$\left|\frac{1}{\theta np} \sum_{i=1}^p |\mathbf{f}_j^* \mathbf{x}_i|^2 - 1\right| \leq t$$

877 holds with high probability for each  $j = 1, \dots, n$ . We apply a union bound to control the  $\ell_\infty$ -norm  
878 in (70), and hence get the desired result.  $\blacksquare$

879 **Lemma H.3** For any  $\varepsilon \in (0, 1)$ , when  $p \geq C\theta^{-1}\varepsilon^{-2} \log^3 n$ , we have

$$\begin{aligned} \left\| \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right\| &\leq (1 + \varepsilon) \|\mathbf{C}_a\|^2 \\ \left\| \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2} - (\mathbf{C}_a^\top \mathbf{C}_a)^{-1/2} \right\| &\leq \frac{4\kappa^2 \varepsilon}{\sigma_{\min}^2(\mathbf{C}_a)} \end{aligned}$$

880 holds with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ . Here,  $\kappa$  is the condition number of  $\mathbf{C}_a$ , and  
881  $\sigma_{\min}(\mathbf{C}_a)$  is the smallest singular value of  $\mathbf{C}_a$ .

882 **Proof** For any  $\varepsilon \in (0, 1)$ , from Lemma H.2, when  $p \geq C\theta^{-1}\varepsilon^{-2} \log^3 n$  we know that the event

$$\mathcal{E}(\varepsilon) \doteq \left\{ \left\| \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top \mathbf{C}_{\mathbf{x}_i} - \mathbf{I} \right\| \leq \varepsilon \right\}$$

883 holds with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ . Conditioned on the event  $\mathcal{E}(\varepsilon)$ , let us denote

$$\mathbf{A} = \mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}} > 0,$$

884 and let  $\sigma_{\max}(\mathbf{A})$ ,  $\sigma_{\min}(\mathbf{A})$  be the largest and smallest singular values of  $\mathbf{A}$ , respectively. Then we  
885 observe,

$$\begin{aligned} \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} &= \mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}} + \underbrace{\mathbf{C}_{\mathbf{a}}^\top \left[ \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{x}_i}^\top \mathbf{C}_{\mathbf{x}_i} - \mathbf{I} \right] \mathbf{C}_{\mathbf{a}} }_{\Delta} \\ &= \mathbf{A} + \Delta, \quad \|\Delta\| \leq \varepsilon \cdot \sigma_{\max}(\mathbf{A}). \end{aligned}$$

886 Therefore, we have

$$\left\| \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right\| \leq \|\mathbf{A}\| + \|\Delta\| \leq (1 + \varepsilon) \|\mathbf{C}_{\mathbf{a}}\|^2.$$

887 By Lemma D.12, whenever

$$\|\Delta\| \leq \frac{1}{2} \sigma_{\min}(\mathbf{A}) \implies \varepsilon \leq \frac{1}{2} \frac{\sigma_{\min}(\mathbf{A})}{\sigma_{\max}(\mathbf{A})} = \frac{1}{2\kappa^2},$$

888 we know that

$$\begin{aligned} \left\| \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2} - (\mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}})^{-1/2} \right\| &= \|(\mathbf{A} + \Delta)^{-1/2} - \mathbf{A}^{-1/2}\| \\ &\leq \frac{4\|\Delta\|}{\sigma_{\min}^2(\mathbf{A})} \leq \frac{4\varepsilon\sigma_{\max}(\mathbf{A})}{\sigma_{\min}^2(\mathbf{A})} = \frac{4\kappa^2\varepsilon}{\sigma_{\min}^2(\mathbf{C}_{\mathbf{a}})}. \end{aligned}$$

889 ■

890 **Lemma H.4** Let  $\theta \in (1/n, 1/3)$ , and given a  $\delta \in (0, 1)$ . Whenever

$$p \geq C \frac{\kappa^8}{\theta \delta^2 \sigma_{\min}^2(\mathbf{C}_{\mathbf{a}})} \log^3 n,$$

891 we have

$$\begin{aligned} \|\mathbf{RQ}^{-1} - \mathbf{I}\| &\leq \delta, & \|\mathbf{RQ}^{-1}\| &\leq 1 + \delta, \\ \|(\mathbf{RQ}^{-1})^{-1} - \mathbf{I}\| &\leq 2\delta, & \|(\mathbf{RQ}^{-1})^{-1}\| &\leq 1 + 2\delta \end{aligned}$$

892 hold with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ .

893 **Proof** First, by Lemma H.3, for a given  $\varepsilon \in (0, 1)$ , when  $p \geq C_1 \theta^{-1} \varepsilon^{-2} \log^3 n$ , we have

$$\begin{aligned} \|\mathbf{RQ}^{-1} - \mathbf{I}\| &= \left\| \mathbf{I} - \mathbf{C}_{\mathbf{a}} \left( \frac{1}{\sqrt{\theta np}} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2} (\mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}})^{1/2} \mathbf{C}_{\mathbf{a}}^{-1} \right\| \\ &\leq \kappa \cdot \|\mathbf{C}_{\mathbf{a}}\| \cdot \left\| \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2} - (\mathbf{C}_{\mathbf{a}}^\top \mathbf{C}_{\mathbf{a}})^{-1/2} \right\| \\ &\leq \kappa \|\mathbf{C}_{\mathbf{a}}\| \frac{4\kappa^2\varepsilon}{\sigma_{\min}^2(\mathbf{C}_{\mathbf{a}})} \leq \frac{4\kappa^4\varepsilon}{\sigma_{\min}(\mathbf{C}_{\mathbf{a}})}, \end{aligned}$$

Table 2: Gradient for each different loss function

Loss function	$\nabla \varphi(\mathbf{q})$ for 1D problem (73)	$\nabla \varphi(\mathbf{Z})$ for 2D problem <sup>13</sup> (74)
$\ell^1$ -loss	$\frac{1}{np} \sum_{i=1}^p \check{\bar{\mathbf{y}}}_i \circledast \text{sign}(\bar{\mathbf{y}}_i \circledast \mathbf{q})$	$\frac{1}{n^2 p} \sum_{i=1}^p \check{\bar{\mathbf{Y}}}_i \boxtimes \text{sign}(\bar{\mathbf{Y}}_i \boxtimes \mathbf{Z})$
Huber-loss	$\frac{1}{np} \sum_{i=1}^p \check{\bar{\mathbf{y}}}_i \circledast h'_\mu(\bar{\mathbf{y}}_i \circledast \mathbf{q})$	$\frac{1}{n^2 p} \sum_{i=1}^p \check{\bar{\mathbf{Y}}}_i \boxtimes h'_\mu(\bar{\mathbf{Y}}_i \boxtimes \mathbf{Z})$
$\ell^4$ -loss	$-\frac{1}{np} \sum_{i=1}^p \check{\bar{\mathbf{y}}}_i \circledast (\bar{\mathbf{y}}_i \circledast \mathbf{q})^{\odot 3}$	$-\frac{1}{n^2 p} \sum_{i=1}^p \check{\bar{\mathbf{Y}}}_i \boxtimes (\bar{\mathbf{Y}}_i \boxtimes \mathbf{Z})^{\odot 3}$

894 and

$$\|\mathbf{R}\mathbf{Q}^{-1}\| \leq 1 + \|\mathbf{I} - \mathbf{R}\mathbf{Q}^{-1}\| \leq 1 + \frac{4\kappa^4 \varepsilon}{\sigma_{\min}(\mathbf{C}_a)}$$

895 hold with probability at least  $1 - p^{-c_1 n \theta} - n^{-c_2}$ . Similarly, by Lemma H.3,

$$\begin{aligned} \|\mathbf{I} - (\mathbf{R}\mathbf{Q}^{-1})^{-1}\| &= \left\| \mathbf{I} - \mathbf{C}_a (\mathbf{C}_a^\top \mathbf{C}_a)^{-1/2} \left( \frac{1}{\sqrt{\theta np}} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{1/2} \mathbf{C}_a^{-1} \right\| \\ &\leq \kappa \cdot \left\| \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right\|^{1/2} \cdot \left\| \left( \frac{1}{\theta np} \sum_{i=1}^p \mathbf{C}_{\mathbf{y}_i}^\top \mathbf{C}_{\mathbf{y}_i} \right)^{-1/2} - (\mathbf{C}_a^\top \mathbf{C}_a)^{-1/2} \right\| \\ &\leq \kappa \cdot \frac{4\kappa^2 \varepsilon}{\sigma_{\min}^2(\mathbf{C}_a)} \cdot (1 + \varepsilon)^{1/2} \|\mathbf{C}_a\| \leq \frac{8\kappa^4 \varepsilon}{\sigma_{\min}(\mathbf{C}_a)}, \end{aligned}$$

896 and

$$\|(\mathbf{R}\mathbf{Q}^{-1})^{-1}\| \leq 1 + \|\mathbf{I} - (\mathbf{R}\mathbf{Q}^{-1})^{-1}\| \leq 1 + \frac{8\kappa^4 \varepsilon}{\sigma_{\min}(\mathbf{C}_a)}$$

897 Thus, replace  $\delta = \frac{4\kappa^4 \varepsilon}{\sigma_{\min}(\mathbf{C}_a)}$ , we obtain the desired result.  $\blacksquare$

## 898 I Algorithms and Implementation Details

899

900 In this section, we provide detailed descriptions of our algorithms. First, we introduce the details  
 901 Riemannian (sub)gradient descent method for 1D problem. Second, we discuss about subgradient  
 902 methods for solving the LP rounding problem. Finally, we provide more details about how to solve  
 903 problems in 2D.

904 For the purpose of implementation efficiency, we describe the problem and algorithms based on  
 905 circulant convolution, which is slightly different from the main sections. Because our gradient descent  
 906 method works for any sparse promoting loss function (other than Huber loss), in the following we  
 907 describe the problem and the algorithm in a more general form. However, it should be noted that our  
 908 analysis is only specified for Huber loss in the following sections.

### 909 I.1 Riemannian (sub)gradient descent methods

910 Here, we consider (sub)gradient descent for optimizing a more general problem

$$\min_{\mathbf{q}} \varphi(\mathbf{q}) := \frac{1}{np} \sum_{i=1}^p \psi(\mathbf{C}_{\mathbf{y}_i} \mathbf{P} \mathbf{q}), \quad \text{s.t. } \|\mathbf{q}\| = 1,$$

---

<sup>13</sup>Here, for 2D problem,  $\check{\mathbf{Z}}$  denotes a flip operator that flips a matrix  $\mathbf{Z} \in \mathbb{R}^{n_1 \times n_2}$  both vertically and horizontally, i.e.,  $\check{\mathbf{Z}}_{i,j} = Z_{n_1-i+1, n_2-j+1}$ .

---

**Algorithm 1** Riemannian (sub)gradient descent algorithm

---

**Input:** observation  $\{\mathbf{y}_i\}_{i=1}^m$   
**Output:** the vector  $\mathbf{q}_*$ ,

Precondition the data by  $\bar{\mathbf{y}}_i = \mathbf{y}_i \circledast \mathbf{v}$ , with  $\mathbf{v} = \left( \frac{1}{\theta np} \sum_{i=1}^p |\mathbf{y}_i|^{\odot 2} \right)^{\odot -1/2}$ .

Initialize the iterate  $\mathbf{q}^{(0)}$  and stepsize  $\tau^{(0)}$ .

**while** not converged **do**

    Update the iterate by

$$\mathbf{q}^{(k+1)} = \mathcal{P}_{\mathbb{S}^{n-1}} \left( \mathbf{q}^{(k)} - \tau^{(k)} \operatorname{grad} \varphi(\mathbf{q}^{(k)}) \right).$$

    Choose a new stepsize  $\tau^{(k+1)}$ , and set  $k \leftarrow k + 1$ .

**end while**

---

911 where  $\psi(\mathbf{z})$  can be  $\ell^1$ -loss ( $\psi(\mathbf{z}) = \|\mathbf{z}\|_1$ ), Huber-loss ( $\psi(\mathbf{z}) = H_\mu(\mathbf{z})$ ), and  $\ell^4$ -loss ( $\psi(\mathbf{z}) =$   
912  $-\|\mathbf{z}\|_4^4$ ). The preconditioning matrix  $\mathbf{P}$  can be written as

$$\mathbf{P} = \mathbf{C}_v, \quad \mathbf{v} = \mathbf{F}^{-1} \left( \left( \frac{1}{\theta np} \sum_{i=1}^p |\hat{\mathbf{y}}_i|^{\odot 2} \right)^{\odot -1/2} \right),$$

913 where  $\hat{\mathbf{y}}_i = \mathbf{F}\mathbf{y}_i$ , so that

$$\mathbf{C}_{\mathbf{y}_i} \mathbf{P} = \mathbf{C}_{\mathbf{y}_i} \mathbf{C}_v = \mathbf{C}_{\mathbf{y}_i \circledast \mathbf{v}} = \mathbf{C}_{\bar{\mathbf{y}}_i}, \quad \bar{\mathbf{y}}_i = \mathbf{y}_i \circledast \mathbf{v}.$$

914 Therefore, our problem can be rewritten as

$$\min_{\mathbf{q}} \varphi(\mathbf{q}) := \frac{1}{np} \sum_{i=1}^p \psi(\bar{\mathbf{y}}_i \circledast \mathbf{q}), \quad \text{s.t. } \|\mathbf{q}\| = 1. \quad (73)$$

915 Starting from an initialization, we solve the problem via Riemannian (sub)gradient descent,

$$\mathbf{q}^{(k+1)} = \mathcal{P}_{\mathbb{S}^{n-1}} \left( \mathbf{q}^{(k)} - \tau^{(k)} \cdot \operatorname{grad} \varphi(\mathbf{q}^{(k)}) \right),$$

916 where  $\tau^{(k)}$  is the stepsize, and the Riemannian (sub)gradient is

$$\operatorname{grad} \varphi(\mathbf{q}) = \mathcal{P}_{\mathbf{q}^\perp} \nabla \varphi(\mathbf{q}),$$

917 which is defined on the *tangent space*<sup>14</sup>  $T_q \mathbb{S}^{n-1}$  at a point  $\mathbf{q} \in \mathbb{S}^{n-1}$ . Table 2 lists the calculation of (sub)gradients  $\nabla \varphi(\mathbf{q})$  for different loss functions. For each iteration, the projection operator  $\mathcal{P}_{\mathbb{S}^{n-1}}(\mathbf{z}) = \mathbf{z}/\|\mathbf{z}\|$  retracts the iterate back to the sphere. Let  $\odot$  denotes entry-wise  
918 power/multiplication, the overall algorithm is summarized in Algorithm 1.  
920

921 **Initialization.** In our theory, we showed that starting from a random initialization drawn uniformly  
922 over the sphere,

$$\mathbf{q}^{(0)} = \mathbf{d}, \quad \mathbf{d} \sim \mathcal{U}(\mathbb{S}^{n-1}),$$

923 for Huber-loss, Riemannian gradient descent method provably recovers the target solution. On the  
924 other hand, we could also cook up a data-driven initialization by choosing a row of  $\mathbf{C}_{\bar{\mathbf{y}}_i}$ ,

$$\mathbf{q}^{(0)} = \mathcal{P}_{\mathbb{S}^{n-1}} \left( \mathbf{C}_{\bar{\mathbf{y}}_i}^\top \mathbf{e}_j \right)$$

925 for some randomly chosen  $1 \leq i \leq p$  and  $1 \leq j \leq n$ . By observing

$$\mathbf{C}_{\bar{\mathbf{y}}_i} \approx \mathbf{C}_{\mathbf{x}_i} \mathbf{C}_a (\mathbf{C}_a^\top \mathbf{C}_a)^{-1/2}, \quad \mathbf{q}^{(0)} \approx \mathcal{P}_{\mathbb{S}^{n-1}} \left( (\mathbf{C}_a^\top \mathbf{C}_a)^{-1/2} \mathbf{C}_a^\top s_j [\tilde{\mathbf{x}}_i] \right),$$

926 we have

$$\mathbf{C}_{\bar{\mathbf{y}}_i} \mathbf{q}^{(0)} \approx \alpha \mathbf{C}_{\mathbf{x}_i} \mathbf{C}_a (\mathbf{C}_a^\top \mathbf{C}_a)^{-1} \mathbf{C}_a^\top s_\ell [\tilde{\mathbf{x}}_i] = \alpha \mathbf{C}_{\mathbf{x}_i} s_\ell [\tilde{\mathbf{x}}_i].$$

927 This suggests that our particular initialization  $\mathbf{q}^{(0)}$  is acting like  $s_\ell [\tilde{\mathbf{x}}_i]$  in the rotated domain. It is  
928 sparse and possesses several large spiky entries more biased towards the target solutions. Empirically,  
929 we find this data-driven initialization often works better than random initializations.

<sup>14</sup> We refer the readers to Chapter 3 of [47] for more details.

930 **Choice of stepsizes.** For Huber and  $\ell^4$  losses, we can choose a fixed stepsize  $\tau^{(k)}$  for all iterates to  
 931 guarantee linear convergence. For subgradient descent of  $\ell^1$ -loss, it often achieves linear convergence  
 932 when we choose a geometrically decreasing sequence of stepsize  $\tau^{(k)}$  [48]. Empirically, we find that  
 933 the algorithm converges much faster when Riemannian linesearch is deployed (see Algorithm 2).

---

**Algorithm 2** Riemannian linesearch for stepsize  $\tau$ 


---

**Input:**  $a, x, \tau_0, \eta \in (0.5, 1), \beta \in (0, 1)$ ,  
**Output:**  $\tau, \mathcal{R}_a^{\mathcal{M}}(-\tau P_{T_{\mathcal{M}}} \nabla \psi_x(a))$   
 Initialize  $\tau \leftarrow \tau_0$ ,  
 Set  $\tilde{q} = \mathcal{P}_{\mathbb{S}^{n-1}}(q - \tau \text{grad } \varphi(q))$   
**while**  $\varphi(\tilde{q}) \geq \varphi(q) - \tau \cdot \eta \cdot \|\text{grad } \varphi(q)\|^2$  **do**  
 $\tau \leftarrow \beta \tau$ ,  
 Update  $\tilde{q} = \mathcal{P}_{\mathbb{S}^{n-1}}(q - \tau \text{grad } \varphi(q))$ .  
**end while**

---

934 **I.2 LP rounding**

935 Due to preconditioning or smoothing effects of our choice of loss functions, the Riemannian  
 936 (sub)gradient descent methods can only produce an approximate solution. To obtain the exact  
 937 solution, we use the solution  $r = q_*$  produced by gradient methods as a warm start, and solve another  
 938 phase-two LP rounding problem,

$$\min_{q} \zeta(q) := \frac{1}{np} \sum_{i=1}^p \|\bar{y}_i \circledast q\|_1 \quad \text{s.t.} \quad \langle r, q \rangle = 1.$$

939 Since the feasible set  $\langle r, q \rangle = 1$  is essentially the tangent space of the sphere  $\mathbb{S}^{n-1}$  at  $q_*$ , whenever  $q_*$   
 940 is close enough to one of the target solutions, one should expect that the optimizer  $q_r$  of LP rounding  
 941 exactly recovers the inverse of the kernel  $a$  up to a scaled-shift. To address this computational issue,  
 942 we utilize a *projected subgradient method* for solving the LP rounding problem. Namely, we take

$$\begin{aligned} q^{(k+1)} &= r + (I - rr^\top) (q^{(k)} - \tau^{(k)} g^{(k)}) \\ &= q^{(k)} - \tau^{(k)} \mathcal{P}_{r^\perp} g^{(k)}, \end{aligned}$$

943 where  $g^{(k)}$  is the subgradient at  $q^{(k)}$  with

$$g^{(k)} = \frac{1}{np} \sum_{i=1}^p \check{y}_i \circledast \text{sign}(\bar{y}_i \circledast q^{(k)}).$$

944 By choosing a geometrically shrinking stepsizes

$$\tau^{(k+1)} = \beta \tau^{(k)}, \quad \beta \in (0, 1).$$

945 we show that the subgradient descent linearly converges to the target solution. The overall method is  
 946 summarized in Algorithm 3.

947 **I.3 Solving problems in 2D**

948 Finally, we briefly discuss about technical details about solving the MCS-BD problem in 2D, which  
 949 appears broadly in imaging applications such as image deblurring [13–15] and microscopy imaging  
 950 [3, 16, 17].

951 **Problem formulation.** Given the measurements

$$Y_i = A \boxtimes X_i, \quad 1 \leq i \leq p,$$

952 where  $\boxtimes$  denotes 2D convolution,  $A \in \mathbb{R}^{n \times n}$  is a 2D kernel, and  $X_i \in \mathbb{R}^{n \times n}$  is a sparse activation  
 953 map, we want to recover  $A$  and  $\{X_i\}_{i=1}^p$  simultaneously. We first precondition the data via

$$\bar{Y}_i = Y_i \boxtimes V, \quad V = \mathcal{F}^{-1} \left( \left( \frac{1}{\theta n^2 p} \sum_{i=1}^p |\mathcal{F}(Y_i)|^{\odot 2} \right)^{\odot -1/2} \right),$$

---

**Algorithm 3** Projected subgradient method for solving the LP rounding problem

---

**Input:** observation  $\{\mathbf{y}_i\}_{i=1}^m$ , vector  $\mathbf{r}$ , stepsize  $\tau_0$ , and  $\beta \in (0, 1)$ .  
**Output:** the solution  $\mathbf{q}_*$ ,

Precondition the data by  $\bar{\mathbf{y}}_i = \mathbf{y}_i \odot \mathbf{v}$ , with  $\mathbf{v} = \left( \frac{1}{\theta np} \sum_{i=1}^p |\mathbf{y}_i|^{\odot 2} \right)^{\odot -1/2}$ .

Initialize  $\mathbf{q}^{(0)} = \mathbf{r}$ ,  $\tau^{(0)} = \tau_0$

**while** not converged **do**

Update the iterate

$$\mathbf{q}^{(k+1)} = \mathbf{q}^{(k)} - \tau^{(k)} \mathcal{P}_{\mathbf{r}^\perp} \mathbf{g}^{(k)}.$$

Set  $\tau^{(k+1)} = \beta \tau^{(k)}$ , and  $k \leftarrow k + 1$ .

**end while**

---

954 where  $\mathcal{F}(\cdot)$  denote the 2D DFT operator. By using the preconditioned data, we solve the following  
955 optimization problem

$$\min_{\mathbf{Z}} \varphi(\mathbf{Z}) := \frac{1}{n^2 p} \sum_{i=1}^p \psi(\bar{\mathbf{Y}}_i \boxtimes \mathbf{Z}), \quad \text{s.t. } \|\mathbf{Z}\|_F = 1, \quad (74)$$

956 where  $\varphi(\cdot)$  is the loss function (e.g.,  $\ell^1$ , Huber,  $\ell^4$ -loss), and  $\|\cdot\|_F$  denotes the Frobenius norm. If the  
957 problem (74) can be solved to the target solution  $\mathbf{Z}_*$ , then we can recover the kernel and the sparse  
958 activation map up to a signed-shift by

$$\mathbf{A}_* = \mathcal{F}^{-1} \left( \mathcal{F}(\mathbf{V} \boxtimes \mathbf{Z}_*)^{\odot -1} \right), \quad \mathbf{X}_i^* = (\mathbf{Y}_i \boxtimes \mathbf{V}) \boxtimes \mathbf{Z}_*, \quad 1 \leq i \leq p.$$

959 **Riemannian (sub)gradient descent.** Similar to the 1D case, we can optimize the problem (74) via  
960 Riemannian (sub)gradient descent,

$$\mathbf{Z}^{(k+1)} = \mathcal{P}_F \left( \mathbf{Z}^{(k)} - \tau^{(k)} \cdot \text{grad } \varphi(\mathbf{Z}^{(k)}) \right),$$

961 where the Riemannian (sub)gradient

$$\text{grad } \varphi(\mathbf{Z}) = \mathcal{P}_{\mathbf{Z}^\perp} \nabla \varphi(\mathbf{Z}).$$

962 The gradient  $\nabla \varphi(\mathbf{Z})$  for different loss functions are recorded in Table 2. For any  $\mathbf{W} \in \mathbb{R}^{n \times n}$ , the  
963 normalization operator  $\mathcal{P}_F(\cdot)$  and projection operator  $\mathcal{P}_{\mathbf{Z}^\perp}(\cdot)$  are defined as

$$\mathcal{P}_F(\mathbf{W}) := \mathbf{W} / \|\mathbf{W}\|_F, \quad \mathcal{P}_{\mathbf{Z}^\perp}(\mathbf{W}) := \mathbf{W} - \|\mathbf{Z}\|_F^{-2} \langle \mathbf{Z}, \mathbf{W} \rangle \mathbf{Z}.$$

964 The initialization and stepsize  $\tau^{(k)}$  can be chosen similarly as the 1D case.

965 **LP rounding.** Similar to 1D case, we solve a phase-two linear program to obtain exact solution. By  
966 using the solution  $\mathbf{Z}_*$  produced by Riemannian gradient descent as a warm start  $\mathbf{U} = \mathbf{Z}_*$ , we solve

$$\min_{\mathbf{Z}} \frac{1}{n^2 p} \sum_{i=1}^p \|\bar{\mathbf{Y}}_i \boxtimes \mathbf{Z}\|_1, \quad \text{s.t. } \langle \mathbf{U}, \mathbf{Z} \rangle = 1.$$

967 We optimize the LP rounding problem via subgradient descent,

$$\mathbf{Z}^{(k+1)} = \mathbf{Z}^{(k)} - \tau^{(k)} \mathcal{P}_{\mathbf{U}^\perp} \mathbf{G}^{(k)},$$

968 where we choose a geometrically decreasing stepsize  $\tau^{(k)}$  and set the subgradient

$$\mathbf{G}^{(k)} = \frac{1}{n^2 p} \sum_{i=1}^p \check{\bar{\mathbf{Y}}}_i \boxtimes \text{sign} \left( \bar{\mathbf{Y}}_i \boxtimes \mathbf{Z}^{(k)} \right).$$