

# SUBGRADIENT DESCENT LEARNS ORTHOGONAL DICTIONARIES

Yu Bai, Qijia Jiang & Ju Sun

Stanford University

{yub, qjiang2, sunju}@stanford.edu

## ABSTRACT

This paper concerns dictionary learning, i.e., sparse coding, a fundamental representation learning problem. We show that a subgradient descent algorithm, with random initialization, can recover orthogonal dictionaries on a natural nonsmooth, nonconvex  $\ell_1$  minimization formulation of the problem, under mild statistical assumption on the data. This is in contrast to previous provable methods that require either expensive computation or delicate initialization schemes. Our analysis develops several tools for characterizing landscapes of nonsmooth functions, which might be of independent interest for provable training of deep networks with nonsmooth activations (e.g., ReLU), among other applications. Preliminary synthetic and real experiments corroborate our analysis and show that our algorithm works well empirically in recovering orthogonal dictionaries.

## 1 INTRODUCTION

Dictionary learning (DL), i.e., sparse coding, concerns the problem of learning compact representations, i.e., given data  $\mathbf{Y}$ , one tries to find a representation basis  $\mathbf{A}$  and coefficients  $\mathbf{X}$ , so that  $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$  where  $\mathbf{X}$  is most sparse. DL has numerous applications especially in image processing and computer vision (Mairal et al., 2014). When posed in analytical form, DL seeks a transformation  $\mathbf{Q}$  such that  $\mathbf{Q}\mathbf{Y}$  is sparse; in this sense DL can be considered as an (extremely!) primitive “deep” network (Ravishanker & Bresler, 2013).

Many heuristic algorithms have been proposed to solve DL since the seminal work of Olshausen & Field (1996), most of them surprisingly effective in practice (Mairal et al., 2014; Sun et al., 2015). However, understandings on when and how DL is solvable have only recently started to emerge. Under appropriate generating models on  $\mathbf{A}$  and  $\mathbf{X}$ , Spielman et al. (2012) showed that complete (i.e., square, invertible)  $\mathbf{A}$  can be recovered from  $\mathbf{Y}$ , provided that  $\mathbf{X}$  is ultra-sparse. Subsequent works (Agarwal et al., 2017; Arora et al., 2014; 2015; Chatterji & Bartlett, 2017; Awasthi & Vijayaraghavan, 2018) provided similar guarantees for overcomplete (i.e. fat)  $\mathbf{A}$ , again in the ultra-sparse regime. The latter methods are invariably based on nonconvex optimization with model-dependent initialization, rendering their practicality on real data questionable.

The ensuing developments have focused on breaking the sparsity barrier and addressing the practicality issue. Convex relaxations based on the sum-of-squares (SOS) SDP hierarchy can recover overcomplete  $\mathbf{A}$  when  $\mathbf{X}$  has linear sparsity (Barak et al., 2015; Ma et al., 2016; Schramm & Steurer, 2017), while incurring expensive computation (solving large-scale SDP’s or large-scale tensor decomposition). By contrast, Sun et al. (2015) showed that complete  $\mathbf{A}$  can be recovered in the linear sparsity regime by solving a certain nonconvex problem with arbitrary initialization. However, the second-order optimization method proposed there is still expensive. This problem is partially addressed by (Gilboa et al., 2018) which proved that the first-order gradient descent with random initialization enjoys a similar performance guarantee.

A standing barrier toward practicality is dealing with nonsmooth functions. To promote sparsity in the coefficients, the  $\ell_1$  norm is the function of choice in practical DL, as is common in modern signal processing and machine learning (Candès, 2014): despite its nonsmoothness, this choice often admits highly scalable numerical methods, such as proximal gradient method and alternating direction

---

The reader is welcome to refer to our [arXiv version](#) for future updates.

method (Mairal et al., 2014). The analyses in Sun et al. (2015); Gilboa et al. (2018), however, focused on characterizing the algorithm-independent function landscape of a certain nonconvex formulation of DL, which takes a smooth surrogate to  $\ell_1$  to get around the nonsmoothness. The tactic smoothing there introduced substantial analysis difficulty, and broke the practical advantage of computing with the simple  $\ell_1$  function.

In this paper, we show that working directly with a natural  $\ell_1$  norm formulation results in neat analysis and a practical algorithm. We focus on the problem of learning orthogonal dictionaries: given data  $\{\mathbf{y}_i\}_{i \in [m]}$  generated as  $\mathbf{y}_i = \mathbf{A}\mathbf{x}_i$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is a fixed unknown orthogonal matrix and each  $\mathbf{x}_i \in \mathbb{R}^n$  is an iid Bernoulli-Gaussian random vector with parameter  $\theta \in (0, 1)$ , recover  $\mathbf{A}$ . This statistical model is the same as in previous works (Spielman et al., 2012; Sun et al., 2015).

Write  $\mathbf{Y} \doteq [\mathbf{y}_1, \dots, \mathbf{y}_m]$  and similarly  $\mathbf{X} \doteq [\mathbf{x}_1, \dots, \mathbf{x}_m]$ . We propose to recover  $\mathbf{A}$  by solving the following nonconvex (due to the constraint), nonsmooth (due to the objective) optimization problem:

$$\text{minimize}_{\mathbf{q} \in \mathbb{R}^n} f(\mathbf{q}) \doteq \frac{1}{m} \|\mathbf{q}^\top \mathbf{Y}\|_1 = \frac{1}{m} \sum_{i=1}^m |\mathbf{q}^\top \mathbf{y}_i| \quad \text{subject to } \|\mathbf{q}\|_2 = 1. \quad (1.1)$$

Based on the statistical model,  $\mathbf{q}^\top \mathbf{Y} = \mathbf{q}^\top \mathbf{A}\mathbf{X}$  has the highest sparsity when  $\mathbf{q}$  is a column of  $\mathbf{A}$  (up to sign) so that  $\mathbf{q}^\top \mathbf{A}$  is 1-sparse. Spielman et al. (2012) formalized this intuition and optimized the same objective as Eq. (1.1) with a  $\|\mathbf{q}\|_\infty = 1$  constraint, which only works when  $\theta \sim O(1/\sqrt{n})$ . Sun et al. (2015) worked with the sphere constraint but replaced the  $\ell_1$  objective with a smooth surrogate, introducing substantial analytical and computational deficiencies as alluded above.

In contrast, we show that with sufficiently many samples, the optimization landscape of formulation (1.1) is benign with high probability (over the randomness of  $\mathbf{X}$ ), and a simple Riemannian subgradient descent algorithm can provably recover  $\mathbf{A}$  in polynomial time.

**Theorem 1.1** (Main result, informal version of Theorem 3.1). *Assume  $\theta \in [1/n, 1/2]$ . For  $m \geq \Omega(\theta^{-2} n^4 \log^4 n)$ , the following holds with high probability: there exists a  $\text{poly}(m, \epsilon^{-1})$ -time algorithm, which runs Riemannian subgradient descent on formulation (1.1) from at most  $O(n \log n)$  independent, uniformly random initial points, and outputs a set of vectors  $\{\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_n\}$  such that up to permutation and sign change,  $\|\hat{\mathbf{a}}_i - \mathbf{a}_i\|_2 \leq \epsilon$  for all  $i \in [n]$ .*

In words, our algorithm works also in the linear sparsity regime, the same as established in Sun et al. (2015); Gilboa et al. (2018), at a lower sample complexity  $O(n^4)$  in contrast to the existing  $O(n^{5.5})$  in Sun et al. (2015).<sup>1</sup> As for the landscape, we show that (Theorems 3.4 and 3.6) each of the desired solutions  $\{\pm \mathbf{a}_i\}_{i \in [n]}$  is a local minimizer of formulation (1.1) with a sufficiently large basin of attraction so that a random initialization will land into one of the basins with at least constant probability. To obtain the result, we integrate and develop elements from nonsmooth analysis (on Riemannian manifolds), set-valued analysis, and random set theory, which might be valuable to studying other nonconvex, nonsmooth optimization problems.

## 1.1 RELATED WORK

**Dictionary learning** Besides the many results sampled above, we highlight similarities of our result to Gilboa et al. (2018). Both propose first-order optimization methods with random initialization, and several quantities we work with in the proofs are the same. A defining difference is we work with the nonsmooth  $\ell_1$  objective directly, while Gilboa et al. (2018) built on the smoothed objective from Sun et al. (2015). We put considerable emphasis on practicality: the subgradient of the nonsmooth objective is considerably cheaper to evaluate than that of the smooth objective in Sun et al. (2015), and in the algorithm we use Euclidean projection rather than exponential mapping to remain feasible—again, the former is much lighter for computation.

**General nonsmooth analysis** While nonsmooth analytic tools such as subdifferential for convex functions are now well received in machine learning and relevant communities, that for general functions are much less so. The Clarke subdifferential and relevant calculus developed for the family of locally Lipschitz functions seem to be particularly relevant, and cover several families of functions of interest, such as convex functions, differentiable functions, and many forms of composition

<sup>1</sup>The sample complexity in Gilboa et al. (2018) is not explicitly stated.

(Clarke, 1990; Aubin, 1998; Bagirov et al., 2014). Remarkably, majority of the tools and results can be generalized to locally Lipschitz functions on Riemannian manifolds (Ledyaev & Zhu, 2007; Hosseini & Pouryayevali, 2011). Our formulation (1.1) is exactly optimization of a locally Lipschitz function (as it is convex) on a Riemannian manifold (the sphere). For simplicity, we try to avoid the full manifold language, nonetheless.

**Nonsmooth optimization on Riemannian manifolds or with constraints** Equally remarkable is many of the smooth optimization techniques and convergence results can be naturally adapted to optimization of locally Lipschitz functions on Riemannian manifolds (Grohs & Hosseini, 2015; Hosseini, 2015; Hosseini & Uschmajew, 2017; Grohs & Hosseini, 2016). New optimization methods such as gradient sampling and variants have been invented to solve general nonsmooth problems (Burke et al., 2005; 2018; Bagirov et al., 2014; Curtis & Que, 2015; Curtis et al., 2017). Almost all available convergence results pertain to only global convergence, which is too weak for our purpose. Our specific convergence analysis gives us a local convergence result (Theorem 3.8).

**Nonsmooth landscape characterization** Nonsmoothness is not a big optimization barrier if the problem is convex; here we review some recent work on analyzing nonconvex nonsmooth problems. Loh & Wainwright (2015) study the regularized empirical risk minimization problem with nonsmooth regularizers and show results of the type “all stationary points are within statistical error of ground truth” under certain restricted strong convexity of the smooth risk. Duchi & Ruan (2017); Davis et al. (2017) study the phase retrieval problem with  $\ell_1$  loss, characterizing its nonconvex nonsmooth landscape and providing efficient algorithms.

There is a recent surge of work on analyzing one-hidden-layer ReLU networks, which are nonconvex and nonsmooth. Algorithm-independent characterizations of the landscape are mostly local and require strong initialization procedures (Zhong et al., 2017), whereas stronger global results can be established via designing new loss functions (Ge et al., 2017), relating to PDEs (Mei et al., 2018), or problem-dependent analysis of the SGD (Li & Yuan, 2017; Li & Liang, 2018). Our result provides an algorithm-independent characterization of the landscape of non-smooth dictionary learning, and is “almost global” in the sense that the initialization condition is satisfied by random initialization with high probability.

**Other nonsmooth problems in application** Prevalence of nonsmooth problems in optimal control and economics is evident from all monographs on nonsmooth analysis (Clarke, 1990; Aubin, 1998; Bagirov et al., 2014). In modern machine learning and data analysis, nonsmooth functions are often taken to encode structural information (e.g., sparsity, low-rankness, quantization), or whenever robust estimation is desired. In deep learning, the optimization problem is nonsmooth when nonsmooth activations are in use, e.g., the popular ReLU. The technical ideas around nonsmooth analysis, set-valued analysis, and random set theory that we gather and develop here are particularly relevant to these applications.

## 2 PRELIMINARIES

**Problem setup** Given an unknown orthogonal dictionary  $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{n \times n}$ , we wish to recover  $\mathbf{A}$  through  $m$  observations of the form

$$\mathbf{y}_i = \mathbf{A}\mathbf{x}_i, \tag{2.1}$$

or  $\mathbf{Y} = \mathbf{A}\mathbf{X}$  in matrix form, where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$  and  $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_m]$ .

The coefficient vectors  $\mathbf{x}_i$  are sampled from the Bernoulli-Gaussian distribution with parameter  $\theta \in (0, 1)$ , denoted as  $\text{BG}(\theta)$ : each entry  $x_{ij}$  is independently drawn from a standard Gaussian with probability  $\theta$  and zero otherwise. The Bernoulli-Gaussian is a good prototype distribution for sparse vectors, as  $\mathbf{x}_i$  will be on average  $\theta$ -sparse. For any  $\mathbf{z} \sim_{iid} \text{Ber}(\theta)$ , we let  $\Omega$  denote the set of non-zero indices, which is a random set itself.

We assume that  $n \geq 3$  and  $\theta \in [1/n, 1/2]$ . In particular,  $\theta \geq 1/n$  is to require that each  $\mathbf{x}_i$  has at least one non-zero entry on average.

**First-order geometry** We will focus on the first-order geometry of the non-smooth objective Eq. (1.1):  $f(\mathbf{q}) = \frac{1}{m} \sum_{i=1}^m |\mathbf{q}^\top \mathbf{y}_i|$ . In the whole Euclidean space  $\mathbb{R}^n$ ,  $f$  is convex with

sub-differential set

$$\partial f(\mathbf{q}) = \frac{1}{m} \sum_{i=1}^m \text{sign}(\mathbf{q}^\top \mathbf{y}_i) \mathbf{y}_i, \quad (2.2)$$

where  $\text{sign}(\cdot)$  is the set-valued sign function (i.e.  $\text{sign}(0) = [-1, 1]$ ). As we minimize  $f$  subject to the constraint  $\|\mathbf{q}\|_2 = 1$ , our problem is no longer convex. The Riemannian sub-differential of  $f$  on  $\mathbb{S}^{n-1}$  is defined as (Hosseini & Uschmajew, 2017):

$$\partial_R f(\mathbf{q}) \doteq (\mathbf{I} - \mathbf{q}\mathbf{q}^\top) \partial f(\mathbf{q}). \quad (2.3)$$

A point  $\mathbf{q}$  is stationary for problem Eq. (1.1) if  $\mathbf{0} \in \partial_R f(\mathbf{q})$ . We will not distinguish between local maxima and saddle points—we call a stationary point  $\mathbf{q}$  a saddle point if there is a descent direction (i.e. direction along which the function is locally maximized at  $\mathbf{q}$ ).

**Set-valued analysis** As the subdifferential is a set-valued mapping, analyzing it requires some set-valued analysis, which we briefly present here. The addition of two sets is defined as the Minkowski summation:  $X + Y = \{x + y : x \in X, y \in Y\}$ . The expectation of random sets is a straightforward extension of the Minkowski sum allowing any measurable “selection” procedure; for the concrete definition see (Molchanov, 2013). The Hausdorff distance between two sets is defined as

$$d_H(X_1, X_2) \doteq \sup \left\{ \sup_{\mathbf{x}_1 \in X_1} d(\mathbf{x}_1, X_2), \sup_{\mathbf{x}_2 \in X_2} d(\mathbf{x}_2, X_1) \right\}. \quad (2.4)$$

Basic properties about the Hausdorff distance are provided in Appendix A.1.

**Notations** Bold small letters (e.g.,  $\mathbf{x}$ ) are vectors and bold capitals are matrices (e.g.,  $\mathbf{X}$ ). The dotted equality  $\doteq$  is for definition. For any positive integer  $k$ ,  $[k] \doteq \{1, \dots, k\}$ . By default,  $\|\cdot\|$  is the  $\ell_2$  norm if applied to a vector, and the operator norm if applied to a matrix.  $C$  and  $c$  or any indexed versions are reserved for universal constants that may change from place to place.

### 3 MAIN RESULT

We now state our main result, the recovery guarantee for learning orthogonal dictionary by solving formulation (1.1).

**Theorem 3.1** (Recovering orthogonal dictionary via subgradient descent). *Suppose we observe*

$$m \geq Cn^4 \theta^{-2} \log^4 n \quad (3.1)$$

*samples in the dictionary learning problem and we desire an accuracy  $\epsilon \in (0, 1)$  for recovering the dictionary. With probability at least  $1 - \exp(-cm\theta^3 n^{-3} \log^{-3} m) - \exp(-c'R/n)$ , an algorithm which runs Riemannian subgradient descent  $R = C'n \log n$  times with independent random initializations on  $\mathbb{S}^{n-1}$  outputs a set of vectors  $\{\hat{\mathbf{a}}_1, \dots, \hat{\mathbf{a}}_n\}$  such that up to permutation and sign change,  $\|\hat{\mathbf{a}}_i - \mathbf{a}_i\|_2 \leq \epsilon$  for all  $i \in [n]$ . The total number of subgradient descent iterations is bounded by*

$$C'' R \theta^{-16/3} \epsilon^{-8/3} n^4 \log^{8/3} n. \quad (3.2)$$

Here  $C, C', C'', c, c' > 0$  are universal constants.

At a high level, the proof of Theorem 3.1 consists of the following steps, which we elaborate throughout the rest of this section.

1. Partition the sphere into  $2n$  symmetric “good sets” and show certain directional gradient is strong on population objective  $\mathbb{E}[f]$  inside the good sets (Section 3.1).
2. Show that the same geometric properties carry over to the empirical objective  $f$  with high probability. This involves proving the uniform convergence of the subdifferential set  $\partial f$  to  $\mathbb{E}[\partial f]$  (Section 3.2).
3. Under the benign geometry, establish the convergence of Riemannian subgradient descent to one of  $\{\pm \mathbf{a}_i : i \in [n]\}$  when initialized in the corresponding “good set” (Section 3.3).
4. Calling the randomly initialized optimization procedure  $O(n \log n)$  times will recover all of  $\{\mathbf{a}_1, \dots, \mathbf{a}_n\}$  with high probability, by a coupon collector’s argument (Section 3.4).

**Scaling and rotating to identity** Throughout the rest of this paper, we are going to assume WLOG that the dictionary is the identity matrix, i.e.  $\mathbf{A} = \mathbf{I}_n$ , so that  $\mathbf{Y} = \mathbf{X}$ ,  $f(\mathbf{q}) = \|\mathbf{q}^\top \mathbf{X}\|_1$ , and the goal is to find the standard basis vectors  $\{\pm \mathbf{e}_1, \dots, \pm \mathbf{e}_n\}$ . The case of a general orthogonal  $\mathbf{A}$  can be reduced to this special case via rotating by  $\mathbf{A}^\top$ :  $\mathbf{q}^\top \mathbf{Y} = \mathbf{q}^\top \mathbf{A} \mathbf{X} = (\mathbf{q}')^\top \mathbf{X}$  where  $\mathbf{q}' = \mathbf{A}^\top \mathbf{q}$  and applying the result on  $\mathbf{q}'$ . We also scale the objective by  $\sqrt{\pi/2}$  for convenience of later analysis.

### 3.1 PROPERTIES OF THE POPULATION OBJECTIVE

We begin by characterizing the geometry of the expected objective  $\mathbb{E}[f]$ . Recall that we have rotated  $\mathbf{A}$  to be identity, so that we have

$$f(\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \frac{1}{m} \|\mathbf{q}^\top \mathbf{X}\|_1 = \sqrt{\frac{\pi}{2}} \cdot \frac{1}{m} \sum_{i=1}^m |\mathbf{q}^\top \mathbf{x}_i|, \quad \partial f(\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \frac{1}{m} \sum_{i=1}^m \text{sign}(\mathbf{q}^\top \mathbf{x}_i) \mathbf{x}_i. \quad (3.3)$$

**Minimizers and saddles of the population objective** We begin by computing the function value and subdifferential set of the population objective and giving a complete characterization of its stationary points, i.e. local minimizers and saddles.

**Proposition 3.2** (Population objective value and gradient). *We have*

$$\mathbb{E}[f](\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}[|\mathbf{q}^\top \mathbf{x}|] = \mathbb{E}_\Omega \|\mathbf{q}_\Omega\| \quad (3.4)$$

$$\partial \mathbb{E}[f](\mathbf{q}) = \mathbb{E}[\partial f](\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}[\text{sign}(\mathbf{q}^\top \mathbf{x}) \mathbf{x}] = \mathbb{E}_\Omega \left\{ \begin{array}{l} \mathbf{q}_\Omega / \|\mathbf{q}_\Omega\|, \quad \mathbf{q}_\Omega \neq 0, \\ \{\mathbf{v}_\Omega : \|\mathbf{v}_\Omega\| \leq 1\}, \quad \mathbf{q}_\Omega = 0. \end{array} \right. \quad (3.5)$$

**Proposition 3.3** (Stationary points). *The stationary points of  $\mathbb{E}[f]$  on the sphere are*

$$\mathcal{S} = \left\{ \frac{1}{\sqrt{k}} \mathbf{q} : \mathbf{q} \in \{-1, 0, 1\}^n, \|\mathbf{q}\|_0 = k, k \in [n] \right\}. \quad (3.6)$$

The case  $k = 1$  corresponds to the  $2n$  global minimizers  $\mathbf{q} = \pm \mathbf{e}_i$ , and all other values of  $k$  correspond to saddle points.

A consequence of Proposition 3.3 is that the population objective has no ‘‘spurious local minima’’: each stationary point is either a global minimizer or a saddle point, though the problem itself is non-convex due to the constraint.

**Identifying  $2n$  ‘‘good’’ subsets** We now define  $2n$  subsets on the sphere, each containing one of the global minimizers  $\{\pm \mathbf{e}_i\}$  and possessing benign geometry for both the population and empirical objective, following (Gilboa et al., 2018). For any  $\zeta \in [0, \infty)$  and  $i \in [n]$  define

$$\mathcal{S}_\zeta^{(i+)} \doteq \left\{ \mathbf{q} : q_i > 0, \frac{q_i^2}{\|\mathbf{q}_{-i}\|_\infty^2} \geq 1 + \zeta \right\}, \quad \mathcal{S}_\zeta^{(i-)} \doteq \left\{ \mathbf{q} : q_i < 0, \frac{q_i^2}{\|\mathbf{q}_{-i}\|_\infty^2} \geq 1 + \zeta \right\}. \quad (3.7)$$

For points in  $\mathcal{S}_\zeta^{(i+)} \cup \mathcal{S}_\zeta^{(i-)}$ , the  $i$ -th index is larger than all other indices (in absolute value) by a multiplicative factor of  $\zeta$ . In particular, for any point in these subsets, the largest index is unique, so by Proposition 3.3 all population saddle points are excluded from these  $2n$  subsets.

Intuitively, this partition can serve as a ‘‘tiebreaker’’: points in  $\mathcal{S}_{\zeta_0}^{(i+)}$  is closer to  $\mathbf{e}_i$  than all the other  $2n - 1$  signed basis vectors. Therefore, we hope that optimization algorithms initialized in this region could favor  $\mathbf{e}_i$  over the other standard basis vectors, which we are going to show is indeed the case. For simplicity, we are going to state our geometry results in  $\mathcal{S}_\zeta^{(n+)}$ ; by symmetry the results will automatically carry over to all the other  $2n - 1$  subsets.

**Theorem 3.4** (Lower bound on directional subgradients). *Fix any  $\zeta_0 \in (0, 1)$ . We have*

(a) For all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$  and all indices  $j \neq n$  such that  $q_j \neq 0$ ,

$$\inf \left\langle \mathbb{E}[\partial_R f](\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \geq \frac{1}{2n} \theta (1 - \theta) \frac{\zeta_0}{1 + \zeta_0}. \quad (3.8)$$

(b) For all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$ , we have that

$$\inf \langle \mathbb{E} [\partial_R f] (\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \geq \frac{1}{8} \theta (1 - \theta) \zeta_0 n^{-3/2} \|\mathbf{q}_{-n}\|. \quad (3.9)$$

These lower bounds verify our intuition: points inside  $\mathcal{S}_{\zeta_0}^{(n+)}$  have subgradients pointing towards  $\mathbf{e}_n$ , both in a coordinate-wise sense and a combined sense: the direction  $\mathbf{e}_n - q_n \mathbf{q}$  is exactly the tangent direction of the sphere at  $\mathbf{q}$  that points towards  $\mathbf{e}_n$ .

### 3.2 BENIGN GEOMETRY OF THE EMPIRICAL OBJECTIVE

We now show that the benign geometry in [Theorem 3.4](#) is carried onto the empirical objective  $f$  given sufficiently many samples, using a concentration argument. The key result behind is the concentration of the empirical subdifferential set to the population subdifferential, where concentration is measured in the Hausdorff distance between sets.

**Proposition 3.5** (Uniform convergence of subdifferential). *For any  $t \in (0, 1]$ , when*

$$m \geq Ct^{-2} n \log^2(n/t), \quad (3.10)$$

with probability at least  $1 - \exp(-cm\theta t^2/\log m)$ , we have

$$d_H(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})) \leq t \quad \text{for all } \mathbf{q} \in \mathbb{S}^{n-1}. \quad (3.11)$$

Here  $C, c \geq 0$  are universal constants.

The concentration result guarantees that the sub-differential set is close to its expectation given sufficiently many samples with high probability. Choosing an appropriate concentration level  $t$ , the lower bounds on the directional subgradients carry over to the empirical objective  $f$ , which we state in the following theorem.

**Theorem 3.6** (Directional subgradient lower bound, empirical objective). *There exist universal constants  $C, c \geq 0$  so that the following holds: for all  $\zeta_0 \in (0, 1)$ , when  $m \geq Cn^4\theta^{-2}\zeta_0^{-2}\log^2(n/\zeta_0)$ , with probability at least  $1 - \exp(-cm\theta^3\zeta_0^2n^{-3}\log^{-1}m)$ , the following properties hold simultaneously for all the  $2n$  subsets  $\{\mathcal{S}_{\zeta_0}^{(i+)}, \mathcal{S}_{\zeta_0}^{(i-)} : i \in [n]\}$ : (stated only for  $\mathcal{S}_{\zeta_0}^{(n+)}$ )*

(a) For all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$  and all  $j \in [n]$  with  $q_j \neq 0$  and  $q_n^2/q_j^2 \leq 3$ ,

$$\inf \left\langle \partial_R f(\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \geq \frac{1}{4n} \theta (1 - \theta) \frac{\zeta_0}{1 + \zeta_0}. \quad (3.12)$$

(b) For all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$ ,

$$\inf \langle \partial_R f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \geq \frac{\sqrt{2}}{16} \theta (1 - \theta) n^{-\frac{3}{2}} \zeta_0 \|\mathbf{q}_{-n}\| \geq \frac{1}{16} \theta (1 - \theta) n^{-\frac{3}{2}} \zeta_0 \|\mathbf{q} - \mathbf{e}_n\|. \quad (3.13)$$

The consequence of [Theorem 3.6](#) is two-fold. First, it guarantees that the only possible stationary point of  $f$  in  $\mathcal{S}_{\zeta_0}^{(n+)}$  is  $\mathbf{e}_n$ : for every other point  $\mathbf{q} \neq \mathbf{e}_n$ , property (b) guarantees that  $0 \notin \partial_R f(\mathbf{q})$ , therefore  $\mathbf{q}$  is non-stationary. Second, the directional subgradient lower bounds allow us to establish convergence of the Riemannian subgradient descent algorithm, in a way similar to showing convergence of unconstrained gradient descent on star strongly convex functions.

We now present an upper bound on the norm of the subdifferential sets, which is needed for the convergence analysis.

**Proposition 3.7.** *There exist universal constants  $C, c \geq 0$  such that*

$$\sup \|\partial f(\mathbf{q})\| \leq 2 \quad \forall \mathbf{q} \in \mathbb{S}^{n-1} \quad (3.14)$$

with probability at least  $1 - \exp(-cm\theta \log^{-1}m)$ , provided that  $m \geq Cn \log n$ . This particularly implies that

$$\sup \|\partial_R f(\mathbf{q})\| \leq 2 \quad \forall \mathbf{q} \in \mathbb{S}^{n-1}. \quad (3.15)$$

### 3.3 FINDING ONE BASIS VIA RIEMANNIAN SUBGRADIENT DESCENT

The benign geometry of the empirical objective allows a simple Riemannian subgradient descent algorithm to find one basis vector a time. The Riemannian subgradient descent algorithm with initialization  $\mathbf{q}^{(0)}$  and step size  $\{\eta^{(k)}\}_{k \geq 0}$  is as follows. For an arbitrary  $\mathbf{v} \in \partial_R f(\mathbf{q}^{(k)})$ ,

$$\mathbf{q}^{(k+1)} = \frac{\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}}{\|\mathbf{q}^{(k)} - \eta^{(k)} \mathbf{v}\|}, \quad \text{for } k = 0, 1, 2, \dots \quad (3.16)$$

Each iteration moves in an arbitrary Riemannian subgradient direction followed by a projection back onto the sphere. We show that the algorithm is guaranteed to find one basis as long as the initialization is in the ‘‘right’’ region. To give a concrete result, we set  $\zeta_0 = 1/(5 \log n)$ .<sup>2</sup>

**Theorem 3.8** (One run of subgradient descent recovers one basis). *Let  $m \geq C\theta^{-2}n^4 \log^4 n$  and  $\epsilon \in (0, 2\theta/25]$ . With probability at least  $1 - \exp(-cm\theta^3 n^{-3} \log^{-3} m)$  the following happens. If the initialization  $\mathbf{q}^{(0)} \in \mathcal{S}_{1/(5 \log n)}^{(n+)}$ , and we run the projected Riemannian subgradient descent with step size  $\eta^{(k)} = k^{-\alpha}/(100\sqrt{n})$  with  $\alpha \in (0, 1/2)$ , and keep track of the best function value so far until after iterate  $K$  is performed, producing  $\mathbf{q}^{\text{best}}$ . Then,  $\mathbf{q}^{\text{best}}$  obeys*

$$f(\mathbf{q}^{\text{best}}) - f(\mathbf{e}_n) \leq \epsilon, \quad \text{and} \quad \|\mathbf{q}^{\text{best}} - \mathbf{e}_n\| \leq \frac{16}{\theta(1-\theta)}\epsilon, \quad (3.17)$$

provided that

$$K \geq \max \left\{ \left( \frac{32000n^{5/2} \log n (1-\alpha)}{\theta(1-\theta)\epsilon} \right)^{1/(1-\alpha)}, \left( \frac{64 \frac{1-\alpha}{1-2\alpha} n^{3/2} \log n}{5\theta(1-\theta)\epsilon} \right)^{1/\alpha} \right\}. \quad (3.18)$$

In particular, choosing  $\alpha = 3/8 < 1/2$ , it suffices to let

$$K \geq K_{3/8} \doteq C'\theta^{-8/3}\epsilon^{-8/3}n^4 \log^{8/3} n. \quad (3.19)$$

Here  $C, C', c \geq 0$  are universal constants.

The above optimization result (Theorem 3.8) shows that Riemannian subgradient descent is able to find the basis vector  $\mathbf{e}_n$  when initialized in the associated region  $\mathcal{S}_{1/(5 \log n)}^{(n+)}$ . We now show that a simple uniformly random initialization on the sphere is guaranteed to be in one of these  $2n$  regions with at least probability  $1/2$ .

**Lemma 3.9** (Random initialization falls in ‘‘good set’’). *Let  $\mathbf{q}^{(0)} \sim \text{Uniform}(\mathbb{S}^{n-1})$ , then with probability at least  $1/2$ ,  $\mathbf{q}^{(0)}$  belongs to one of the  $2n$  sets  $\left\{ \mathcal{S}_{1/(5 \log n)}^{(i+)}, \mathcal{S}_{1/(5 \log n)}^{(i-)} : i \in [n] \right\}$ .*

### 3.4 RECOVERING ALL BASES FROM MULTIPLE RUNS

As long as the initialization belongs to  $\mathcal{S}_{1/(5 \log n)}^{(i+)}$  or  $\mathcal{S}_{1/(5 \log n)}^{(i-)}$ , our finding-one-basis result in Theorem 3.8 guarantees that Riemannian subgradient descent will converge to  $\mathbf{e}_i$  or  $-\mathbf{e}_i$  respectively. Therefore if we run the algorithm with independent, uniformly random initializations on the sphere multiple times, by a coupon collector’s argument, we will recover all the basis vectors. This is formalized in the following theorem.

**Theorem 3.10** (Recovering the identity dictionary from multiple random initializations). *Let  $m \geq Cn^4\theta^{-2} \log^4 n$  and  $\epsilon \in (0, 1)$ , with probability at least  $1 - \exp(-cm\theta^3 n^{-3} \log^{-3} m)$  the following happens. Suppose we run the Riemannian subgradient descent algorithm independently for  $R$  times, each with a uniformly random initialization on  $\mathbb{S}^{n-1}$ , and choose the step size as  $\eta^{(k)} = k^{-3/8}/(100\sqrt{n})$ . Then, provided that  $R \geq C'n \log n$ , all standard basis vectors will be recovered up to  $\epsilon$  accuracy with probability at least  $1 - \exp(-cR/n)$  in  $C'R\theta^{-16/3}\epsilon^{-8/3}n^4 \log^{8/3} n$  iterations. Here  $C, C', c \geq 0$  are universal constants.*

When the dictionary  $\mathbf{A}$  is not the identity matrix, we can apply the rotation argument sketched in the beginning of this section to get the same result, which leads to our main result in Theorem 3.1.

<sup>2</sup>It is possible to set  $\zeta_0$  to other values, inducing different combinations of the final sample complexity, iteration complexity, and repetition complexity in Theorem 3.10.

## 4 PROOF HIGHLIGHTS

A key technical challenge is establishing the uniform convergence of subdifferential sets in [Proposition 3.5](#), which we now elaborate. Recall that the population and empirical subdifferentials are

$$\partial f(\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \frac{1}{m} \sum_{i=1}^m \text{sign}(\mathbf{q}^\top \mathbf{x}_i) \mathbf{x}_i, \quad \mathbb{E}[\partial f](\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}_{\mathbf{x} \sim \text{BG}(\theta)} [\text{sign}(\mathbf{q}^\top \mathbf{x}) \mathbf{x}], \quad (4.1)$$

and we wish to show that the difference between  $\partial f(\mathbf{q})$  and  $\mathbb{E}[\partial f](\mathbf{q})$  is small uniformly over  $\mathbf{q} \in \mathcal{Q} = \mathbb{S}^{n-1}$ . Two challenges stand out in showing such a uniform convergence:

1. The subdifferential is set-valued and random, and it is unclear a-priori how one could formulate and analyze the concentration of random sets.
2. The usual covering argument won't work here, as the Lipschitz gradient property does not hold:  $\partial f(\mathbf{q})$  and  $\mathbb{E}[\partial f](\mathbf{q})$  are not Lipschitz in  $\mathbf{q}$ . Therefore, no matter how fine we cover the sphere in Euclidean distance, points not in this covering can have radically different subdifferential sets.

### 4.1 CONCENTRATION OF RANDOM SETS

We state and analyze concentration of random sets in the Hausdorff distance (defined in [Section 2](#)). We now illustrate how the Hausdorff distance is the “right” distance to consider for concentration of subdifferentials—the reason is that the Hausdorff distance is closely related to the *support function* of sets, which for any set  $S \in \mathbb{R}^n$  is defined as

$$h_S(\mathbf{u}) \doteq \sup_{\mathbf{x} \in S} \langle \mathbf{x}, \mathbf{u} \rangle. \quad (4.2)$$

For convex compact sets, the sup difference between their support functions is exactly the Hausdorff distance.

**Lemma 4.1** (Section 1.3.2, [Molchanov \(2013\)](#)). *For convex compact sets  $X, Y \subset \mathbb{R}^n$ , we have*

$$d_H(X, Y) = \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_X(\mathbf{u}) - h_Y(\mathbf{u})|. \quad (4.3)$$

[Lemma 4.1](#) is convenient for us in the following sense. Suppose we wish to upper bound the difference of  $\partial f(\mathbf{q})$  and  $\mathbb{E}[\partial f](\mathbf{q})$  along some direction  $\mathbf{u} \in \mathbb{S}^{n-1}$  (as we need in proving the key empirical geometry result [Theorem 3.6](#)). As both subdifferential sets are convex and compact, by [Lemma 4.1](#) we immediately have

$$\left| \inf_{\mathbf{g} \in \partial f(\mathbf{q})} \langle \mathbf{g}, \mathbf{u} \rangle - \inf_{\mathbf{g} \in \mathbb{E}[\partial f](\mathbf{q})} \langle \mathbf{g}, \mathbf{u} \rangle \right| = \left| -h_{\partial f(\mathbf{q})}(-\mathbf{u}) + h_{\mathbb{E}[\partial f](\mathbf{q})}(-\mathbf{u}) \right| \leq d_H(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})). \quad (4.4)$$

Therefore, as long as we are able to bound the Hausdorff distance, all directional differences between the subdifferentials are simultaneously bounded, which is exactly what we want to show to carry the benign geometry from the population to the empirical objective.

### 4.2 COVERING IN THE $d_{\mathbb{E}}$ METRIC

We argue that the absence of gradient Lipschitzness is because the Euclidean distance is not the “right” metric in this problem. Think of the toy example  $f(x) = |x|$ , whose subdifferential set  $\partial f(x) = \text{sign}(x)$  is not Lipschitz across  $x = 0$ . However, once we partition  $\mathbb{R}$  into  $\mathbb{R}_{>0}$ ,  $\mathbb{R}_{<0}$  and  $\{0\}$  (i.e. according to the sign pattern), the subdifferential set is Lipschitz on each subset.

The situation with the dictionary learning objective is quite similar: we resolve the gradient non-Lipschitzness by proposing a stronger metric  $d_{\mathbb{E}}$  on the sphere which is sign-pattern aware and averages all “subset angles” between two points. Formally, we define  $d_{\mathbb{E}}$  as

$$d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \doteq \mathbb{P}_{\mathbf{x} \sim \text{BG}(\theta)} [\text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x})] = \frac{1}{\pi} \mathbb{E}_{\Omega} \angle(\mathbf{p}_{\Omega}, \mathbf{q}_{\Omega}), \quad (4.5)$$



(the second equality shown in [Lemma C.1](#).) Our plan is to perform the covering argument in  $d_{\mathbb{E}}$ , which requires showing gradient Lipschitzness in  $d_{\mathbb{E}}$  and bounding the covering number.

**Lipschitzness of  $\partial f$  and  $\mathbb{E}[\partial f]$  in  $d_{\mathbb{E}}$**  For the population subdifferential  $\mathbb{E}[\partial f]$ , note that  $\mathbb{E}[\partial f](\mathbf{q}) = \mathbb{E}_{\mathbf{x} \sim \text{BG}(\theta)}[\text{sign}(\mathbf{q}^{\top} \mathbf{x})]$  (modulo rescaling). Therefore, to bound  $d_{\mathbb{H}}(\mathbb{E}[\partial f](\mathbf{p}), \mathbb{E}[\partial f](\mathbf{q}))$  by [Lemma 4.1](#), we have the bound for all  $\mathbf{u} \in \mathbb{S}^{n-1}$

$$|h_{\mathbb{E}[\partial f](\mathbf{p})}(\mathbf{u}) - h_{\mathbb{E}[\partial f](\mathbf{q})}(\mathbf{u})| = \mathbb{E}[\sup |\text{sign}(\mathbf{p}^{\top} \mathbf{x}) - \text{sign}(\mathbf{q}^{\top} \mathbf{x})| \cdot |\mathbf{x}^{\top} \mathbf{u}|] \quad (4.6)$$

$$\leq 2\mathbb{E}[\mathbb{1}\{\text{sign}(\mathbf{p}^{\top} \mathbf{x}) \neq \text{sign}(\mathbf{q}^{\top} \mathbf{x})\} |\mathbf{x}^{\top} \mathbf{u}|]. \quad (4.7)$$

As long as  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$ , the indicator is non-zero with probability at most  $\epsilon$ , and thus the above expectation should also be small – we bound it by  $O(\epsilon\sqrt{\log(1/\epsilon)})$  in [Lemma F.5](#).

To show the same for the empirical subdifferential  $\partial f$ , one only needs to bound the observed proportion of sign differences for all  $\mathbf{p}, \mathbf{q}$  such that  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$ , which by a VC dimension argument is uniformly bounded by  $2\epsilon$  with high probability ([Lemma C.5](#)).

**Bounding the covering number in  $d_{\mathbb{E}}$**  Our first step is to reduce  $d_{\mathbb{E}}$  to the **maximum length-2 angle** (the  $d_2$  metric) over any consistent support pattern. This is achieved through the following *vector angle inequality* ([Lemma C.2](#)): for any  $\mathbf{p}, \mathbf{q} \in \mathbb{R}^d$  ( $d \geq 3$ ), we have

$$\angle(\mathbf{p}, \mathbf{q}) \leq \sum_{\Omega \subset [d], |\Omega|=2} \angle(\mathbf{p}_{\Omega}, \mathbf{q}_{\Omega}) \quad \text{provided } \angle(\mathbf{p}, \mathbf{q}) \leq \pi/2. \quad (4.8)$$

Therefore, as long as  $\text{sign}(\mathbf{p}) = \text{sign}(\mathbf{q})$  (coordinate-wise) and  $\max_{|\Omega|=2} \angle(\mathbf{p}_{\Omega}, \mathbf{q}_{\Omega}) \leq \epsilon/n^2$ , we would have for all  $|\Omega| \geq 3$  that

$$\angle(\mathbf{p}_{\Omega}, \mathbf{q}_{\Omega}) \leq \pi/2 \quad \text{and} \quad \angle(\mathbf{p}_{\Omega}, \mathbf{q}_{\Omega}) \leq \sum_{\Omega' \subset \Omega, |\Omega'|=2} \angle(\mathbf{p}_{\Omega'}, \mathbf{q}_{\Omega'}) \leq \binom{|\Omega|}{2} \cdot \frac{\epsilon}{n^2} \leq \epsilon. \quad (4.9)$$

By [Eq. \(4.5\)](#), the above implies that  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon/\pi$ , the desired result. Hence the task reduces to constructing an  $\eta = \epsilon/n^2$  covering in  $d_2$  over any consistent sign pattern.

Our second step is a **tight bound on this covering number**: the  $\eta$ -covering number in  $d_2$  is bounded by  $\exp(Cn \log(n/\eta))$  ([Lemma C.3](#)). For bounding this, a first thought would be to take the covering in all size-2 angles (there are  $\binom{n}{2}$  of them) and take the common refinement of all their partitions, which gives covering number  $(C/\eta)^{O(n^2)} = \exp(Cn^2 \log(1/\eta))$ . We improve upon this strategy by *sorting* the coordinates in  $\mathbf{p}$  and restricting attentions in the consecutive size-2 angles after the sorting (there are  $n-1$  of them). We show that a proper covering in these consecutive size-2 angles by  $\eta/n$  will yield a covering for all size-2 angles by  $\eta$ . The corresponding covering number in this case is thus  $(Cn/\eta)^{O(n)} = \exp(Cn \log(n/\eta))$ , which modulo the  $\log n$  factor is the tightest we can get.

## 5 EXPERIMENTS

**Setup** We set the true dictionary  $\mathbf{A}$  to be the identity and random orthogonal matrices, respectively. For each choice, we sweep the combinations of  $(m, n)$  with  $n \in \{30, 50, 70, 100\}$  and  $m = 10n^{\{0.5, 1, 1.5, 2, 2.5\}}$ , and fix the sparsity level at  $\theta = 0.1, 0.3, 0.5$ , respectively. For each  $(m, n)$  pair, we generate 10 problem instances, corresponding to re-sampling the coefficient matrix  $\mathbf{X}$  for 10 times. Note that our theoretical guarantee applies for  $m = \Omega(n^4)$ , and the sample complexity we experiment with here is lower than what our theory requires. To recover the dictionary, we run the Riemannian subgradient descent algorithm [Eq. \(3.16\)](#) with decaying step size  $\eta^{(k)} = 1/\sqrt{k}$ , corresponding to the boundary case  $\alpha = 1/2$  in [Theorem 3.8](#) with a much better base size.

**Metric** As [Theorem 3.1](#) guarantees recovering the entire dictionary with  $R \geq Cn \log n$  independent runs, we perform  $R = \text{round}(5n \log n)$  runs on each instance. For each run, a true dictionary element  $\mathbf{a}_i$  is considered to be found if  $\|\mathbf{a}_i - \mathbf{q}_{\text{best}}\| \leq 10^{-3}$ . For each instance, we regard it a successful recovery if the  $R = \text{round}(5n \log n)$  runs have found all the dictionary elements, and we report the empirical success rate over the 10 instances.

**Result** From our simulations, Riemannian subgradient descent succeeds in recovering the dictionary as long as  $m \geq Cn^2$  (Fig. 2), across different sparsity level  $\theta$ . The dependency on  $n$  is consistent with our theory and suggests that the actual sample complexity requirement for guaranteed recovery might be even lower than  $\tilde{O}(n^4)$  we established.<sup>3</sup> The  $\tilde{O}(n^2)$  rate we observe also matches the results based on the SOS method (Barak et al., 2015; Ma et al., 2016; Schramm & Steurer, 2017). Moreover, the problem seems to become harder when  $\theta$  grows, evident from the observation that the success transition threshold being pushed to the right.

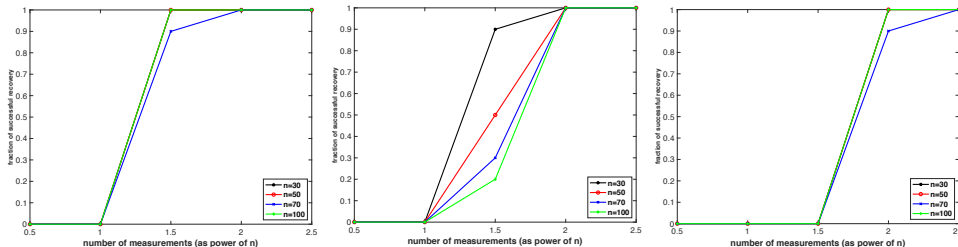


Figure 1: Empirical success rates of recovery of the Riemannian subgradient descent with  $R = 5n \log n$  runs, averaged over 10 instances. Left to right: identity dictionaries with  $\theta = 0.1, 0.3, 0.5$ . See Appendix G for the results on orthogonal dictionaries, which have qualitatively the same behaviors.

**Additional experiments** A faster alternative algorithm for large-scale instances is tested in Appendix H. A complementary experiment on real images is included as Appendix I.

## 6 CONCLUSION AND FUTURE DIRECTIONS

This paper presents the first theoretical guarantee for orthogonal dictionary learning using subgradient descent on a natural  $\ell_1$  minimization formulation. Along the way, we develop tools for analyzing the optimization landscape of nonconvex nonsmooth functions, which could be of broader interest.

For future work, there is an  $O(n^2)$  sample complexity gap between what we established in Theorem 3.1, and what we observed in the simulations alongside previous results based on the SOS method (Barak et al., 2015; Ma et al., 2016; Schramm & Steurer, 2017). As our main geometric result Theorem 3.6 already achieved tight bounds on the directional derivatives, further sample complexity improvement could potentially come out of utilizing second-order information such as the strong negative curvature (Lemma B.2), or careful algorithm-dependent analysis.

While our result applies only to (complete) orthogonal dictionaries, a natural question is whether we can generalize to overcomplete dictionaries. To date the only known provable algorithms for learning overcomplete dictionaries in the linear sparsity regime are based on the SOS method (Barak et al., 2015; Ma et al., 2016; Schramm & Steurer, 2017). We believe that our nonsmooth analysis has the potential of handling over-complete dictionaries, as for reasonably well-conditioned overcomplete dictionaries  $\mathbf{A}$ , each  $\mathbf{a}_i$  (columns of  $\mathbf{A}$ ) makes  $\mathbf{a}_i^\top \mathbf{A}$  approximately 1-sparse and so  $\mathbf{a}_i^\top \mathbf{A} \mathbf{X}$  gives noisy estimate of a certain row of  $\mathbf{X}$ . So the same formulation as Eq. (1.1) intuitively still works. We would like to leave that to future work.

Nonsmooth phase retrieval and deep networks with ReLU mentioned in Section 1.1 are examples of many nonsmooth, nonconvex problems encountered in practice. Most existing theoretical results on these problems tend to be technically vague about handling the nonsmooth points: they either prescribe a rule for choosing a subgradient element, which effectively disconnects theory and practice because numerical testing of nonsmooth points is often not reliable, or ignore the nonsmooth points altogether, assuming that practically numerical methods would never touch these points—this sounds intuitive but no formalism on this appears in the relevant literature yet. Besides our work, (Laurent & von Brecht, 2017; Kakade & Lee, 2018) also warns about potential problems of ignoring nonsmooth points when studying optimization of nonsmooth functions in machine learning.

<sup>3</sup>The  $\tilde{O}(\cdot)$  notation ignores the dependency on logarithmic terms and other factors.

## REFERENCES

- Alekh Agarwal, Animashree Anandkumar, and Praneeth Netrapalli. A clustering approach to learning sparsely used overcomplete dictionaries. *IEEE Trans. Information Theory*, 63(1):575–592, 2017.
- Sanjeev Arora, Rong Ge, and Ankur Moitra. New algorithms for learning incoherent and overcomplete dictionaries. In *Conference on Learning Theory*, pp. 779–806, 2014.
- Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms for sparse coding. 2015.
- Jean-Pierre Aubin. *Optima and equilibria: an introduction to nonlinear analysis*, volume 140. Springer Science & Business Media, 1998.
- Pranjal Awasthi and Aravindan Vijayaraghavan. Towards learning sparsely used dictionaries with arbitrary supports. *arXiv preprint arXiv:1804.08603*, 2018.
- Adil Bagirov, Napsu Karmitsa, and Marko M Mäkelä. *Introduction to Nonsmooth Optimization: theory, practice and software*. Springer, 2014.
- Boaz Barak, Jonathan A Kelner, and David Steurer. Dictionary learning and tensor decomposition via the sum-of-squares method. In *Proceedings of the forty-seventh annual ACM symposium on Theory of computing*, pp. 143–151. ACM, 2015.
- Stéphane Boucheron, Olivier Bousquet, and Gábor Lugosi. Theory of classification: A survey of some recent advances. *ESAIM: probability and statistics*, 9:323–375, 2005.
- James V Burke, Adrian S Lewis, and Michael L Overton. A robust gradient sampling algorithm for nonsmooth, nonconvex optimization. *SIAM Journal on Optimization*, 15(3):751–779, 2005. doi: 10.1137/030601296.
- James V. Burke, Frank E. Curtis, Adrian S. Lewis, Michael L. Overton, and Lucas E. A. Simes. Gradient sampling methods for nonsmooth optimization. *arXiv:1804.11003*, 2018.
- Emmanuel Candès. Mathematics of sparsity (and a few other things). In *Proceedings of the International Congress of Mathematicians*, volume 123, 2014.
- Niladri Chatterji and Peter L Bartlett. Alternating minimization for dictionary learning with random initialization. In *Advances in Neural Information Processing Systems*, pp. 1997–2006, 2017.
- Frank H Clarke. *Optimization and nonsmooth analysis*. SIAM, 1990. doi: 10.1137/1.9781611971309.
- Frank E Curtis and Xiaocun Que. A quasi-newton algorithm for nonconvex, nonsmooth optimization with global convergence guarantees. *Mathematical Programming Computation*, 7(4):399–428, 2015. doi: 10.1007/s12532-015-0086-2.
- Frank E Curtis, Tim Mitchell, and Michael L Overton. A bfgs-sqp method for nonsmooth, nonconvex, constrained optimization and its evaluation using relative minimization profiles. *Optimization Methods and Software*, 32(1):148–181, 2017.
- Damek Davis, Dmitriy Drusvyatskiy, and Courtney Paquette. The nonsmooth landscape of phase retrieval. *arxiv:1711.03247*, 2017. URL <http://arxiv.org/abs/1711.03247>.
- John C. Duchi and Feng Ruan. Solving (most) of a set of quadratic equalities: Composite optimization for robust phase retrieval. *arxiv:1705.02356*, 2017. URL <http://arxiv.org/abs/1705.02356>.
- Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. *arxiv:1711.00501*, 2017. URL <http://arxiv.org/abs/1711.00501>.
- Dar Gilboa, Sam Buchanan, and John Wright. Efficient dictionary learning with gradient descent. In *ICML 2018 Workshop on Modern Trends in Nonconvex Optimization for Machine Learning*, 2018. URL <https://arxiv.org/abs/1809.10313>.

- P Grohs and S Hosseini. Nonsmooth trust region algorithms for locally Lipschitz functions on Riemannian manifolds. *IMA Journal of Numerical Analysis*, 36(3):1167–1192, 2015. doi: 10.1093/imanum/drv043.
- Philipp Grohs and Seyedehsomyeh Hosseini.  $\varepsilon$ -subgradient algorithms for locally Lipschitz functions on Riemannian manifolds. *Advances in Computational Mathematics*, 42(2):333–360, 2016. doi: 10.1007/s10444-015-9426-z.
- Jean-Baptiste Hiriart-Urruty and Claude Lemarchal. *Fundamentals of Convex Analysis*. Springer-Verlag Berlin Heidelberg, 2001. doi: 10.1007/978-3-642-56468-0.
- S Hosseini and MR Pouryayevali. Generalized gradients and characterization of epi-lipschitz sets in Riemannian manifolds. *Nonlinear Analysis: Theory, Methods & Applications*, 74(12):3884–3895, 2011. doi: 10.1016/j.na.2011.02.023.
- Seyedehsomyeh Hosseini. Optimality conditions for global minima of nonconvex functions on Riemannian manifolds. *Pacific Journal of Optimization*, 2015. URL <http://uschmajew.ins.uni-bonn.de/research/pub/hosseini/3.pdf>.
- Seyedehsomyeh Hosseini and André Uschmajew. A Riemannian gradient sampling algorithm for nonsmooth optimization on manifolds. *SIAM Journal on Optimization*, 27(1):173–189, 2017. doi: 10.1137/16M1069298.
- Sham Kakade and Jason D. Lee. Provably correct automatic subdifferentiation for qualified programs. *arXiv:1809.08530*, 2018.
- Thomas Laurent and James von Brecht. The multilinear structure of relu networks. *arxiv:1712.10132*, 2017. URL <http://arxiv.org/abs/1712.10132>.
- Yu Ledyae and Qiji Zhu. Nonsmooth analysis on smooth manifolds. *Transactions of the American Mathematical Society*, 359(8):3687–3732, 2007. doi: 10.1090/S0002-9947-07-04075-5.
- Yuanzhi Li and Yingyu Liang. Learning overparameterized neural networks via stochastic gradient descent on structured data. *arXiv:1808.01204*, 2018.
- Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with relu activation. *arxiv:1705.09886*, 2017. URL <http://arxiv.org/abs/1705.09886>.
- Po-Ling Loh and Martin J Wainwright. Regularized m-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *Journal of Machine Learning Research*, 16:559–616, 2015.
- Tengyu Ma, Jonathan Shi, and David Steurer. Polynomial-time tensor decompositions with sum-of-squares. 2016.
- Julien Mairal, Francis Bach, and Jean Ponce. Sparse modeling for image and vision processing. *Foundations and Trends® in Computer Graphics and Vision*, 8(2-3):85–283, 2014.
- Song Mei, Andrea Montanari, and Phan-Minh Nguyen. A mean field view of the landscape of two-layers neural networks. *arXiv:1804.06561*, 2018.
- Ilya Molchanov. Foundations of stochastic geometry and theory of random sets. In *Stochastic Geometry, Spatial Statistics and Random Fields*, pp. 1–20. Springer, 2013.
- Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607, 1996.
- Barrett O’neill. *Semi-Riemannian geometry with applications to relativity*, volume 103. Academic press, 1983.
- Saiprasad Ravishankar and Yoram Bresler. Learning sparsifying transforms. *IEEE Transactions on Signal Processing*, 61(5):1072–1086, 2013.
- Tselil Schramm and David Steurer. Fast and robust tensor decomposition with applications to dictionary learning. *arXiv:1706.08672*, 2017.

Daniel A Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *Conference on Learning Theory*, pp. 37–1, 2012.

Shlomo Sternberg. *Dynamical Systems*. Dover Publications, Inc, 2013.

Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere. *arxiv:1504.06785*, 2015. URL <http://arxiv.org/abs/1504.06785>.

Aad Van Der Vaart and Jon A Wellner. A note on bounds for VC dimensions. *Institute of Mathematical Statistics collections*, 5:103, 2009.

Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge University Press, 2018. doi: 10.1017/9781108231596.

Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. *arxiv:1706.03175*, 2017. URL <http://arxiv.org/abs/1706.03175>.

## A TECHNICAL TOOLS

### A.1 HAUSDORFF DISTANCE

We need the Hausdorff metric to measure differences between nonempty sets. For any set  $X$  and a point  $\mathbf{p}$  in  $\mathbb{R}^n$ , the point-to-set distance is defined as

$$d(\mathbf{q}, X) \doteq \inf_{\mathbf{x} \in X} \|\mathbf{x} - \mathbf{p}\|. \quad (\text{A.1})$$

For any two sets  $X_1, X_2 \in \mathbb{R}^n$ , the Hausdorff distance is defined as

$$d_H(X_1, X_2) \doteq \sup \left\{ \sup_{\mathbf{x}_1 \in X_1} d(\mathbf{x}_1, X_2), \sup_{\mathbf{x}_2 \in X_2} d(\mathbf{x}_2, X_1) \right\}. \quad (\text{A.2})$$

When  $X_1$  is a singleton, say  $X_1 = \{\mathbf{p}\}$ . Then

$$d_H(\{\mathbf{p}\}, X_2) = \sup_{\mathbf{x}_2 \in X_2} \|\mathbf{x}_2 - \mathbf{p}\|. \quad (\text{A.3})$$

Moreover, for any sets  $X_1, X_2, Y_1, Y_2 \subset \mathbb{R}^n$ ,

$$d_H(X_1 + Y_1, X_2 + Y_2) \leq d_H(X_1, X_2) + d_H(Y_1, Y_2). \quad (\text{A.4})$$

On the sets of nonempty, compact subsets of  $\mathbb{R}^n$ , the Hausdorff metric is a valid metric; particularly, it obeys the triangular inequality: for nonempty, compact subsets  $X, Y, Z \subset \mathbb{R}^n$ ,

$$d_H(X, Z) \leq d_H(X, Y) + d_H(Y, Z). \quad (\text{A.5})$$

See, e.g., Sec. 7.1 of [Sternberg \(2013\)](#) for a proof.

**Lemma A.1** (Restatement of [Lemma A.1](#)). *For convex compact sets  $X, Y \subset \mathbb{R}^n$ , we have*

$$d_H(X, Y) = \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_X(\mathbf{u}) - h_Y(\mathbf{u})|, \quad (\text{A.6})$$

where  $h_S(\mathbf{u}) \doteq \sup_{\mathbf{x} \in S} \langle \mathbf{x}, \mathbf{u} \rangle$  is the support function associated with the set  $S$ .

### A.2 SUB-GAUSSIAN RANDOM MATRICES AND PROCESSES

**Proposition A.2** (Talagrand’s comparison inequality, Corollary 8.6.3 and Exercise 8.6.5 of [Vershynin \(2018\)](#)). *Let  $\{X_{\mathbf{x}}\}_{\mathbf{x} \in T}$  be a zero-mean random process on a subset  $T \subset \mathbb{R}^n$ . Assume that for all  $\mathbf{x}, \mathbf{y} \in T$  we have*

$$\|X_{\mathbf{x}} - X_{\mathbf{y}}\|_{\psi_2} \leq K \|\mathbf{x} - \mathbf{y}\|. \quad (\text{A.7})$$

Then, for any  $t > 0$

$$\sup_{\mathbf{x} \in T} |X_{\mathbf{x}}| \leq CK [w(T) + t \cdot \text{rad}(T)] \quad (\text{A.8})$$

with probability at least  $1 - 2 \exp(-t^2)$ . Here  $w(T) \doteq \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \sup_{\mathbf{x} \in T} \langle \mathbf{x}, \mathbf{g} \rangle$  is the Gaussian width of  $T$  and  $\text{rad}(T) = \sup_{\mathbf{x} \in T} \|\mathbf{x}\|$  is the radius of  $T$ .

**Proposition A.3** (Deviation inequality for sub-Gaussian matrices, Theorem 9.1.1 and Exercise 9.1.8 of Vershynin (2018)). *Let  $\mathbf{A}$  be an  $n \times m$  matrix whose rows  $\mathbf{A}^i$ 's are independent, isotropic, and sub-Gaussian random vectors in  $\mathbb{R}^m$ . Then for any subset  $T \subset \mathbb{R}^m$ , we have*

$$\mathbb{P} \left[ \sup_{\mathbf{x} \in T} \left| \|\mathbf{A}\mathbf{x}\| - \sqrt{n} \|\mathbf{x}\| \right| > CK^2 [w(T) + t \cdot \text{rad}(T)] \right] \leq 2 \exp(-t^2). \quad (\text{A.9})$$

Here  $K = \max_i \|\mathbf{A}^i\|_{\psi_2}$ .

## B PROOFS FOR SECTION 3.1

### B.1 PROOF OF PROPOSITION 3.2

We have

$$\mathbb{E}[f](\mathbf{q}) = \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}[\|\mathbf{q}^\top \mathbf{x}\|] = \sqrt{\frac{\pi}{2}} \cdot \mathbb{E}[\|\mathbf{q}_\Omega^\top \mathbf{x}_\Omega\|] = \mathbb{E}[\|\mathbf{q}_\Omega\|_2], \quad (\text{B.1})$$

where the last equality is obtained by conditioning on  $\Omega$  and the fact that  $\mathbb{E}_{Z \sim N(0, \sigma^2)}[|Z|] = \sqrt{2/\pi} \sigma$ .

The subdifferential expression comes from

$$\partial \|\mathbf{q}_\Omega\|_2 = \begin{cases} \frac{\mathbf{q}_\Omega}{\|\mathbf{q}_\Omega\|_2}, & \mathbf{q}_\Omega \neq 0; \\ \{\mathbf{v}_\Omega : \|\mathbf{v}_\Omega\|_2 \leq 1\}, & \mathbf{q}_\Omega = 0, \end{cases} \quad (\text{B.2})$$

and the fact that  $\partial \mathbb{E}[f](\mathbf{q}) = \partial \mathbb{E}[\|\mathbf{q}_\Omega\|_2] = \mathbb{E}[\partial \|\mathbf{q}_\Omega\|_2]$  as the sub-differential and expectation can be exchanged for convex functions (Hiriart-Urruty & Lemarchal, 2001). By the same exchangeability result, we also have  $\mathbb{E}[\partial f(\mathbf{q})] = \partial \mathbb{E}[f](\mathbf{q})$ .

### B.2 PROOF OF PROPOSITION 3.3

We first show that points in the claimed set are indeed stationary points by taking the choice  $\mathbf{v}_\Omega = \mathbf{0}$  in Eq. (3.5), giving the subgradient choice  $\mathbb{E}[\partial f](\mathbf{q}) = \mathbb{E}[\mathbf{q}_\Omega / \|\mathbf{q}_\Omega\|_2 \mathbb{1}\{\mathbf{q}_\Omega \neq 0\}]$ . Let  $\mathbf{q} \in \mathcal{S}$  and such that  $\|\mathbf{q}\|_0 = k$ . For all  $j \in \text{supp}(\mathbf{q})$ , we have

$$\mathbf{e}_j^\top \mathbb{E}[\partial f](\mathbf{q}) = \theta q_j \cdot \mathbb{E}_\Omega \frac{1}{\|\mathbf{q}_\Omega\|} \mathbb{1}\{j \in \Omega\} \quad (\text{B.3})$$

$$= \theta q_j \cdot \left[ \sum_{i=1}^k \theta^{i-1} (1-\theta)^{k-i} \cdot \sqrt{\frac{k}{i}} \right] \quad (\text{B.4})$$

$$= q_j \cdot \sum_{i=1}^k \theta^i (1-\theta)^{k-i} \cdot \frac{1}{\sqrt{i}} \doteq c(\theta, k) q_j \quad (\text{B.5})$$

On the other hand, for all  $j \notin \text{supp}(\mathbf{q})$ , we always have  $[\mathbf{q}_\Omega]_j = 0$ , so  $\mathbf{e}_j^\top \mathbb{E}[\partial f](\mathbf{q}) = 0$ . Therefore, we have that  $\mathbb{E}[\partial f](\mathbf{q}) = c(\theta, k) \mathbf{q}$ , and so

$$(\mathbf{I} - \mathbf{q}\mathbf{q}^\top) \mathbb{E}[\partial f](\mathbf{q}) = c(\theta, k) \mathbf{q} - c(\theta, k) \mathbf{q} = \mathbf{0}. \quad (\text{B.6})$$

Therefore  $\mathbf{q} \in \mathcal{S}$  is stationary. To see that  $\{\pm \mathbf{e}_i : i \in [n]\}$  are the global minima, note that for all  $\mathbf{q} \in \mathbb{S}^{n-1}$ , we have

$$\mathbb{E}[f](\mathbf{q}) = \mathbb{E}[\|\mathbf{q}_\Omega\|_2] \geq \mathbb{E}[\|\mathbf{q}_\Omega\|_2^2] = \theta. \quad (\text{B.7})$$

Equality holds if and only if  $\|\mathbf{q}_\Omega\|_2 \in \{0, 1\}$  almost surely, which is only satisfied at  $\mathbf{q} \in \{\pm \mathbf{e}_i : i \in [n]\}$ .

To see that the other  $\mathbf{q}$ 's are saddles, we only need to show that there exists a tangent direction along which  $\mathbf{q}$  is local max. Indeed, for any other  $\mathbf{q}$ , there exists at least two non-zero entries (with equal absolute value): WLOG assume that  $q_1 = q_n > 0$ . Using the reparametrization in Appendix B.3 and applying Lemma B.2, we get that  $\mathbb{E}[f](\mathbf{q})$  is directionally differentiable along  $[-\mathbf{q}_{-n}; \frac{1-q_n^2}{q_n}]$ , with derivative zero (necessarily, because  $\mathbf{0} \in \mathbb{E}[\partial_R f](\mathbf{q})$ ) and strictly negative second derivative.

Therefore  $\mathbb{E}[f](\mathbf{q})$  is locally maximized at  $\mathbf{q}$  along this tangent direction, which shows that  $\mathbf{q}$  is a saddle point.

The other direction (all other points are not stationary) is implied by [Theorem 3.4](#), which guarantees that  $\mathbf{0} \notin \mathbb{E}[\partial_R f](\mathbf{q})$  whenever  $\mathbf{q} \notin \mathcal{S}$ . Indeed, as long as  $\mathbf{q} \notin \mathcal{S}$ ,  $\mathbf{q}$  has a max absolute value coordinate (say  $n$ ) and another non-zero coordinate with strictly smaller absolute value (say  $j$ ). For this pair of indices, the proof of [Theorem 3.4\(a\)](#) goes through for index  $j$  (even if  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$  does not necessarily hold because the max index might not be unique), which implies that  $\mathbf{0} \notin \mathbb{E}[\partial_R f](\mathbf{q})$ .

### B.3 REPARAMETRIZATION

For analysis purposes, we introduce the reparametrization  $\mathbf{w} = \mathbf{q}_{1:(n-1)}$  in the region  $\mathcal{S}_0^{(n+)}$ , following ([Sun et al., 2015](#)). With this reparametrization, the problem becomes

$$\text{minimize}_{\mathbf{w} \in \mathbb{R}^{n-1}} g(\mathbf{w}) \doteq \sqrt{\frac{\pi}{2}} \cdot \frac{1}{m} \left\| \left[ \mathbf{w}; \sqrt{1 - \|\mathbf{w}\|^2} \right]^\top \mathbf{X} \right\|_1 \quad \text{subject to } \|\mathbf{w}\| \leq \sqrt{\frac{n-1}{n}}. \quad (\text{B.8})$$

The constraint comes from the fact that  $q_n \geq 1/\sqrt{n}$  and thus  $\|\mathbf{w}\| \leq \sqrt{(n-1)/n}$ .

**Lemma B.1.** *We have*

$$\mathbb{E}_{\mathbf{x} \sim \text{iid BG}(\theta)} g(\mathbf{w}) = (1 - \theta) \mathbb{E}_\Omega \|\mathbf{w}_\Omega\| + \theta \mathbb{E}_\Omega \sqrt{1 - \|\mathbf{w}_{\Omega^c}\|^2}. \quad (\text{B.9})$$

**Proof.** Direct calculation gives

$$\mathbb{E}_{\mathbf{x} \sim \text{iid BG}(\theta)} g(\mathbf{w}) \quad (\text{B.10})$$

$$= \mathbb{E}_{\Omega \sim \text{iid Ber}(\theta), \omega \sim \text{Ber}(\theta)} \mathbb{E}_{\mathbf{z} \sim \text{iid } \mathcal{N}(0,1), z \sim \mathcal{N}(0,1)} \left| \left[ \mathbf{w}; \sqrt{1 - \|\mathbf{w}\|^2} \right]^\top ([\Omega; \omega] \odot [\mathbf{z}; z]) \right| \quad (\text{B.11})$$

$$= (1 - \theta) \mathbb{E}_{\Omega \sim \text{iid Ber}(\theta)} \mathbb{E}_{\mathbf{z} \sim \text{iid } \mathcal{N}(0,1)} |\mathbf{w}_\Omega^\top \mathbf{z}| \\ + \theta \mathbb{E}_{\Omega \sim \text{iid Ber}(\theta)} \mathbb{E}_{\mathbf{z} \sim \text{iid } \mathcal{N}(0,1), z \sim \mathcal{N}(0,1)} \left| \mathbf{w}_\Omega^\top \mathbf{z} + \sqrt{1 - \|\mathbf{w}\|^2} z \right| \quad (\text{B.12})$$

$$= (1 - \theta) \mathbb{E}_{\Omega \sim \text{iid Ber}(\theta)} \|\mathbf{w}_\Omega\| + \theta \mathbb{E}_{\Omega \sim \text{iid Ber}(\theta)} \sqrt{1 - \|\mathbf{w}_{\Omega^c}\|^2}, \quad (\text{B.13})$$

as claimed.  $\blacksquare$

**Lemma B.2** (Negative-curvature region). *For all unit vector  $\mathbf{v} \in \mathbb{S}^{n-1}$  and all  $s \in (0, 1)$ , let*

$$h_{\mathbf{v}}(s) \doteq \mathbb{E}[g](s\mathbf{v}) \quad (\text{B.14})$$

*it holds that*

$$\nabla^2 h_{\mathbf{v}}(s) \leq -\theta(1 - \theta). \quad (\text{B.15})$$

*In other words, for all  $\mathbf{w} \neq \mathbf{0}$ ,  $\pm \mathbf{w}/\|\mathbf{w}\|$  is a direction of negative curvature.*

**Proof.** By [Lemma B.1](#),

$$h_{\mathbf{v}}(s) = (1 - \theta) s \mathbb{E}_\Omega \|\mathbf{v}_\Omega\| + \theta \mathbb{E}_\Omega \sqrt{1 - s^2 \|\mathbf{v}_{\Omega^c}\|^2}. \quad (\text{B.16})$$

For  $s \in (0, 1)$ ,  $h_{\mathbf{v}}(s)$  is twice differentiable, and we have

$$\nabla^2 h_{\mathbf{v}}(s) = -\theta \mathbb{E}_\Omega \frac{\|\mathbf{v}_{\Omega^c}\|^2}{(1 - s^2 \|\mathbf{v}_{\Omega^c}\|^2)^{3/2}} \quad (\text{B.17})$$

$$\leq -\theta \mathbb{E}_\Omega \|\mathbf{v}_{\Omega^c}\|^2 = -\theta(1 - \theta), \quad (\text{B.18})$$

completing the proof.  $\blacksquare$

**Lemma B.3** (Inward gradient). *For any  $\mathbf{w}$  with  $\|\mathbf{w}\|^2 + \|\mathbf{w}\|_\infty^2 \leq 1$ ,*

$$D_{\mathbf{w}/\|\mathbf{w}\|}^c \mathbb{E}[g](\mathbf{w}) \geq \theta(1-\theta) \left( 1/\sqrt{1 + \|\mathbf{w}\|_\infty^2 / \|\mathbf{w}\|^2} - \|\mathbf{w}\| \right). \quad (\text{B.19})$$

**Proof.** For any unit vector  $\mathbf{v} \in \mathbb{R}^{n-1}$ , define  $h_{\mathbf{v}}(t) \doteq \mathbb{E}[g](t\mathbf{v})$  for  $t \in (0, 1)$ . We have from Lemma B.1

$$h_{\mathbf{v}}(t) = (1-\theta)t\mathbb{E}_\Omega \|\mathbf{v}_\Omega\| + \theta\mathbb{E}_\Omega \sqrt{1 - t^2 \|\mathbf{v}_\Omega^c\|^2}. \quad (\text{B.20})$$

Moreover,

$$\nabla_i h_{\mathbf{v}}(t) = (1-\theta)\mathbb{E}_\Omega \|\mathbf{v}_\Omega\| - \theta\mathbb{E}_\Omega \frac{t \|\mathbf{v}_\Omega^c\|^2}{\sqrt{1 - t^2 \|\mathbf{v}_\Omega^c\|^2}} \quad (\text{B.21})$$

$$= (1-\theta)\mathbb{E}_\Omega \frac{\|\mathbf{v}_\Omega\|^2}{\|\mathbf{v}_\Omega\|} - \theta\mathbb{E}_\Omega \frac{t \|\mathbf{v}_\Omega^c\|^2}{\sqrt{1 - t^2 \|\mathbf{v}_\Omega^c\|^2}} \quad (\text{assuming } \frac{0}{0} \doteq 0) \quad (\text{B.22})$$

$$= (1-\theta) \sum_{i=1}^{n-1} \mathbb{E}_\Omega \frac{v_i^2 \mathbb{1}\{i \in \Omega\}}{\sqrt{v_i^2 \mathbb{1}\{i \in \Omega\} + \|\mathbf{v}_{\Omega \setminus i}\|^2}} - \theta \sum_{i=1}^{n-1} \mathbb{E}_\Omega \frac{tv_i^2 \mathbb{1}\{i \notin \Omega\}}{\sqrt{1 - t^2 v_i^2 \mathbb{1}\{i \notin \Omega\} - t^2 \|\mathbf{v}_{\Omega^c \setminus i}\|^2}} \quad (\text{B.23})$$

$$= \theta(1-\theta) \sum_{i=1}^{n-1} \mathbb{E}_\Omega \left[ \frac{v_i^2}{\sqrt{v_i^2 + \|\mathbf{v}_{\Omega \setminus i}\|^2}} - \frac{tv_i^2}{\sqrt{1 - t^2 v_i^2 - t^2 \|\mathbf{v}_{\Omega^c \setminus i}\|^2}} \right] \quad (\text{B.24})$$

$$= \theta(1-\theta)t \sum_{i=1}^{n-1} v_i^2 \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{t^2 v_i^2 + t^2 \|\mathbf{v}_{\Omega \setminus i}\|^2}} - \frac{1}{\sqrt{1 - t^2 v_i^2 - t^2 \|\mathbf{v}_{\Omega^c \setminus i}\|^2}} \right] \quad (\text{B.25})$$

$$= \theta(1-\theta)t \sum_{i=1}^{n-1} v_i^2 \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{t^2 v_i^2 + t^2 \|\mathbf{v}_{\Omega \setminus i}\|^2}} - \frac{1}{\sqrt{1 - t^2 \|\mathbf{v}\|^2 + t^2 \|\mathbf{v}_{\Omega \setminus i}\|^2}} \right]. \quad (\text{B.26})$$

We are interested in the regime of  $t$  so that

$$1 - t^2 \|\mathbf{v}\|^2 \geq t^2 \|\mathbf{v}\|_\infty^2 \implies t \leq 1/\sqrt{1 + \|\mathbf{v}\|_\infty^2}. \quad (\text{B.27})$$

So  $\nabla_t h_{\mathbf{v}}(t) \geq 0$  holds always for  $t \leq 1/\sqrt{1 + \|\mathbf{v}\|_\infty^2}$ .

By Lemma B.2,  $\nabla^2 h_{\mathbf{v}}(t) \leq -\theta(1-\theta)$  over  $t \in (0, 1)$ , which implies

$$\langle \nabla_t h_{\mathbf{v}}(t_1) - \nabla_t h_{\mathbf{v}}(t_2), t_1 - t_2 \rangle \leq -\theta(1-\theta)(t_1 - t_2)^2. \quad (\text{B.28})$$

Taking  $t_1 = 1/\sqrt{1 + \|\mathbf{v}\|_\infty^2}$  and considering  $t_2 \in [0, t_1]$ , we have

$$\nabla_t h_{\mathbf{v}}(t_2) \geq \nabla_t h_{\mathbf{v}}(t_1) + \theta(1-\theta)(t_1 - t_2) \geq \theta(1-\theta) \left( 1/\sqrt{1 + \|\mathbf{v}\|_\infty^2} - t_2 \right). \quad (\text{B.29})$$

For any  $\mathbf{w}$ , applying the above result to the unit vector  $\mathbf{w}/\|\mathbf{w}\|$  and recognizing that  $\nabla_t h_{\mathbf{w}/\|\mathbf{w}\|}(t) = D_{\mathbf{w}/\|\mathbf{w}\|}^c g(\mathbf{w}) = D_{\mathbf{w}/\|\mathbf{w}\|}^c g(\mathbf{w})$ , we complete the proof. ■

#### B.4 PROOF OF THEOREM 3.4

We first show Eq. (3.9) using the reparametrization in Appendix B.3. We have

$$\langle \partial_R f(\mathbf{q}), \mathbf{q} - \mathbf{e}_n \rangle = \langle \partial f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle = q_n \langle \partial g(\mathbf{w}), \mathbf{w} \rangle, \quad (\text{B.30})$$



where the second equality follows by differentiating  $g$  via the chain rule. Now, by [Lemma B.3](#),

$$q_n \langle \mathbb{E} [\partial g] (\mathbf{w}), \mathbf{w} \rangle \geq \|\mathbf{w}\| \theta (1 - \theta) \cdot q_n \left( \frac{\|\mathbf{w}\|}{\sqrt{\|\mathbf{w}\|^2 + \|\mathbf{w}\|_\infty^2}} - \|\mathbf{w}\| \right). \quad (\text{B.31})$$

For each radial direction  $\mathbf{v} \doteq \mathbf{w}/\|\mathbf{w}\|$ , consider points of the form  $t\mathbf{v}$  with  $t \leq 1/\sqrt{1 + \|\mathbf{v}\|_\infty^2}$ . Obviously, the function

$$\tilde{h}(t) \doteq q_n(t\mathbf{v}) \left( \frac{\|t\mathbf{v}\|}{\sqrt{\|t\mathbf{v}\|^2 + \|t\mathbf{v}\|_\infty^2}} - \|t\mathbf{v}\| \right) = q_n(t\mathbf{v}) \left( \frac{1}{\sqrt{1 + \|\mathbf{v}\|_\infty^2}} - t \right) \quad (\text{B.32})$$

is monotonically decreasing wrt  $t$ . Thus, to derive a lower bound, it is enough to consider the largest  $t$  allowed. In  $\mathcal{S}_{\zeta_0}^{(n+)}$ , the limit amounts to requiring  $q_n^2/\|\mathbf{w}\|_\infty^2 = 1 + \zeta_0$ ,

$$1 - t_0^2 = t_0^2 \|\mathbf{v}\|_\infty^2 (1 + \zeta_0) \implies t_0 = \frac{1}{\sqrt{1 + (1 + \zeta_0) \|\mathbf{v}\|_\infty^2}}. \quad (\text{B.33})$$

So for any fixed  $\mathbf{v}$  and all allowed  $t$  for points in  $\mathcal{S}_{\zeta_0}^{(n+)}$ , a uniform lower bound is

$$q_n(t_0\mathbf{v}) \left( \frac{1}{\sqrt{1 + \|\mathbf{v}\|_\infty^2}} - t_0 \right) \quad (\text{B.34})$$

$$\geq \frac{1}{\sqrt{n}} \left( \frac{1}{\sqrt{1 + \|\mathbf{v}\|_\infty^2}} - \frac{1}{\sqrt{1 + (1 + \zeta_0) \|\mathbf{v}\|_\infty^2}} \right) \geq \frac{1}{8\sqrt{n}} \zeta_0 \|\mathbf{v}\|_\infty^2 \geq \frac{1}{8} \zeta_0 n^{-3/2}. \quad (\text{B.35})$$

So we conclude that for all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$ ,

$$\langle \mathbb{E} [\partial f] (\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \geq \frac{1}{8} \theta (1 - \theta) \zeta_0 n^{-3/2} \|\mathbf{w}\| = \frac{1}{8} \theta (1 - \theta) \zeta_0 n^{-3/2} \|\mathbf{q}_{-n}\|. \quad (\text{B.36})$$

We now turn to showing [Eq. \(3.8\)](#). For  $\mathbf{e}_j$  with  $q_j \neq 0$ ,

$$\begin{aligned} \frac{1}{q_j} \mathbf{e}_j^\top \mathbb{E} [\partial f] (\mathbf{q}) &= \frac{1}{q_j} \mathbf{e}_j^\top \mathbb{E}_\Omega \left[ \frac{\mathbf{q}_\Omega}{\|\mathbf{q}_\Omega\|} \mathbb{1} \{ \mathbf{q}_\Omega \neq \mathbf{0} \} + \{ \mathbf{v}_\Omega : \|\mathbf{v}_\Omega\| \leq 1 \} \mathbb{1} \{ \mathbf{q}_\Omega = \mathbf{0} \} \right] \\ &= \frac{1}{q_j} \mathbb{E}_\Omega \left[ \frac{\langle \mathbf{q}_\Omega, \mathbf{e}_j \rangle}{\|\mathbf{q}_\Omega\|} \mathbb{1} \{ \mathbf{q}_\Omega \neq \mathbf{0} \} \right] = \frac{1}{q_j} \theta q_j \mathbb{E}_\Omega \left[ \frac{1}{\|\mathbf{q}_\Omega\|} \mathbb{1} \{ j \in \Omega \} \right] = \theta \mathbb{E} \left[ \frac{1}{\sqrt{q_j^2 + \|\mathbf{q}_{\Omega \setminus \{j\}}\|^2}} \right] \geq \theta. \end{aligned} \quad (\text{B.37})$$

So for all  $j$  with  $q_j \neq 0$ , we have

$$\begin{aligned} &\left\langle \mathbb{E} [\partial_R f] (\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \\ &= \left\langle (\mathbf{I} - \mathbf{q}\mathbf{q}^\top) \mathbb{E} [\partial f] (\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \end{aligned} \quad (\text{B.38})$$

$$= \left\langle \mathbb{E} [\partial f] (\mathbf{q}), \frac{1}{q_j} \mathbf{e}_j - \frac{1}{q_n} \mathbf{e}_n \right\rangle \quad (\text{B.39})$$

$$= \theta \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_j^2 + \|\mathbf{q}_{\Omega \setminus \{j\}}\|^2}} \right] - \theta \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_n^2 + \|\mathbf{q}_{\Omega \setminus \{n\}}\|^2}} \right] \quad (\text{B.40})$$

$$= \theta^2 \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_j^2 + q_n^2 + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2}} \right] + \theta (1 - \theta) \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_j^2 + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2}} \right]$$

$$-\theta^2 \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_j^2 + q_n^2 + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2}} \right] - \theta(1-\theta) \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_n^2 + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2}} \right] \quad (\text{B.41})$$

$$= \theta(1-\theta) \mathbb{E}_\Omega \left[ \frac{1}{\sqrt{q_j^2 + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2}} - \frac{1}{\sqrt{q_n^2 + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2}} \right] \quad (\text{B.42})$$

$$= \theta(1-\theta) \mathbb{E}_\Omega \int_{q_j^2}^{q_n^2} \frac{1}{2(t + \|\mathbf{q}_{\Omega \setminus \{j,n\}}\|^2)^{3/2}} dt \quad (\text{B.43})$$

$$\geq \theta(1-\theta) \frac{1}{2} (q_n^2 - q_j^2) \geq \frac{1}{2} \theta(1-\theta) (q_n^2 - \|\mathbf{q}_{-n}\|_\infty^2) \geq \frac{1}{2} \theta(1-\theta) \frac{\zeta_0}{1+\zeta_0} q_n^2 \quad (\text{B.44})$$

$$\geq \frac{1}{2n} \theta(1-\theta) \frac{\zeta_0}{1+\zeta_0} \quad (\text{B.45})$$

completing the proof.

## C PROOFS FOR SECTION 3.2

### C.1 COVERING IN THE $d_{\mathbb{E}}$ METRIC

For any  $\theta \in (0, 1)$ , define

$$d_{\mathbb{E},\theta}(\mathbf{p}, \mathbf{q}) \doteq \mathbb{E} [\mathbb{1} \{ \text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}) \}] \quad \text{with } \mathbf{x} \sim_{iid} \text{BG}(\theta). \quad (\text{C.1})$$

We stress that this notion always depend on  $\theta$ , and we will omit the subscript  $\theta$  when no confusion is expected. This indeed defines a metric on subsets of  $\mathbb{S}^{n-1}$ .

**Lemma C.1.** *Over any subset of  $\mathbb{S}^{n-1}$  with a consistent support pattern,  $d_{\mathbb{E}}$  is a valid metric.*

**Proof.** Recall that  $\angle(\mathbf{x}, \mathbf{y}) \doteq \arccos \langle \mathbf{x}, \mathbf{y} \rangle$  defines a valid metric on  $\mathbb{S}^{n-1}$ .<sup>4</sup> In particular, the triangular inequality holds. For  $d_{\mathbb{E}}$  and  $\mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}$  with the same support pattern, we have

$$d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) = \mathbb{E} [\mathbb{1} \{ \text{sign}(\mathbf{p}^\top \mathbf{z}) \neq \text{sign}(\mathbf{q}^\top \mathbf{z}) \}] \quad (\text{C.2})$$

$$= \mathbb{E}_\Omega \mathbb{E}_{\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \mathbb{1} \{ \text{sign}(\mathbf{p}_\Omega^\top \mathbf{z}) \neq \text{sign}(\mathbf{q}_\Omega^\top \mathbf{z}) \} \quad (\text{C.3})$$

$$= \mathbb{E}_\Omega (\mathbb{E}_{\mathbf{z}} \mathbb{1} \{ \mathbf{p}_\Omega^\top \mathbf{z} \mathbf{q}_\Omega^\top \mathbf{z} < 0 \} + \mathbb{E}_{\mathbf{z}} \mathbb{1} \{ \mathbf{p}_\Omega^\top \mathbf{z} = 0 \text{ or } \mathbf{q}_\Omega^\top \mathbf{z} = 0, \text{ not both} \}) \quad (\text{C.4})$$

$$= \mathbb{E}_\Omega (\mathbb{E}_{\mathbf{z}} \mathbb{1} \{ \mathbf{p}_\Omega^\top \mathbf{z} \mathbf{q}_\Omega^\top \mathbf{z} < 0 \}) \quad (\text{C.5})$$

$$= \frac{1}{\pi} \mathbb{E}_\Omega \angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega), \quad (\text{C.6})$$

where we have adopted the convention that  $\angle(\mathbf{0}, \mathbf{v}) \doteq 0$  for any  $\mathbf{v}$ . It is easy to verify that  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) = 0 \iff \mathbf{p} = \mathbf{q}$ , and  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) = d_{\mathbb{E}}(\mathbf{q}, \mathbf{p})$ . To show the triangular inequality, note that for any  $\mathbf{p}, \mathbf{q}$  and  $\mathbf{r}$  with the same support pattern,  $\mathbf{p}_\Omega, \mathbf{q}_\Omega$ , and  $\mathbf{r}_\Omega$  are either identically zero, or all nonzero. For the former case,

$$\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \angle(\mathbf{p}_\Omega, \mathbf{r}_\Omega) + \angle(\mathbf{q}_\Omega, \mathbf{r}_\Omega) \quad (\text{C.7})$$

holds trivially. For the latter, since  $\angle(\cdot, \cdot)$  obeys the triangular inequality uniformly over the sphere,

$$\angle\left(\frac{\mathbf{p}_\Omega}{\|\mathbf{p}_\Omega\|}, \frac{\mathbf{q}_\Omega}{\|\mathbf{q}_\Omega\|}\right) \leq \angle\left(\frac{\mathbf{p}_\Omega}{\|\mathbf{p}_\Omega\|}, \frac{\mathbf{r}_\Omega}{\|\mathbf{r}_\Omega\|}\right) + \angle\left(\frac{\mathbf{q}_\Omega}{\|\mathbf{q}_\Omega\|}, \frac{\mathbf{r}_\Omega}{\|\mathbf{r}_\Omega\|}\right), \quad (\text{C.8})$$

which implies

$$\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \angle(\mathbf{p}_\Omega, \mathbf{r}_\Omega) + \angle(\mathbf{q}_\Omega, \mathbf{r}_\Omega). \quad (\text{C.9})$$

<sup>4</sup>This fact can be proved either directly, see, e.g., page 12 of this online notes: <http://www.math.mcgill.ca/drury/notes354.pdf>, or by realizing that the angle equal to the geodesic length, which is the Riemannian distance over the sphere; see, e.g., Riemannian Distance of Chapter 5 of the book O'neill (1983).

So

$$\mathbb{E}_\Omega \angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \mathbb{E}_\Omega \angle(\mathbf{p}_\Omega, \mathbf{r}_\Omega) + \mathbb{E}_\Omega \angle(\mathbf{q}_\Omega, \mathbf{r}_\Omega), \quad (\text{C.10})$$

completing the proof.  $\blacksquare$

**Lemma C.2** (Vector angle inequality). *For  $n \geq 2$ , consider  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^n$  so that  $\angle(\mathbf{u}, \mathbf{v}) \leq \pi/2$ . It holds that*

$$\angle(\mathbf{u}, \mathbf{v}) \leq \sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega). \quad (\text{C.11})$$

**Proof.** The inequality holds trivially when either of  $\mathbf{u}, \mathbf{v}$  is zero. Suppose they are both nonzero and wlog assume both are normalized, i.e.,  $\|\mathbf{u}\| = \|\mathbf{v}\| = 1$ . Then,

$$\sin^2 \angle(\mathbf{u}, \mathbf{v}) = 1 - \cos^2 \angle(\mathbf{u}, \mathbf{v}) \quad (\text{C.12})$$

$$= \|\mathbf{u}\|^2 \|\mathbf{v}\|^2 - \langle \mathbf{u}, \mathbf{v} \rangle^2 \quad (\text{C.13})$$

$$= \sum_{i,j>i} (u_i v_j - u_j v_i)^2 \quad (\text{Lagrange's identity}) \quad (\text{C.14})$$

$$= \sum_{\Omega \in \binom{[n]}{2}} \|\mathbf{u}_\Omega\|^2 \|\mathbf{v}_\Omega\|^2 - \langle \mathbf{u}_\Omega, \mathbf{v}_\Omega \rangle^2 \quad (\text{C.15})$$

$$= \sum_{\Omega \in \binom{[n]}{2}} \|\mathbf{u}_\Omega\|^2 \|\mathbf{v}_\Omega\|^2 \sin^2 \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) \quad (\text{C.16})$$

$$\leq \sum_{\Omega \in \binom{[n]}{2}} \sin^2 \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega). \quad (\text{C.17})$$

If  $\sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) > \pi/2$ , the claimed inequality holds trivially, as  $\angle(\mathbf{u}, \mathbf{v}) \leq \pi/2$  by our assumption. Suppose  $\sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) \leq \pi/2$ . Then,

$$\sum_{\Omega \in \binom{[n]}{2}} \sin^2 \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) \leq \sin^2 \sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) \quad (\text{C.18})$$

by recursive application of the following inequality:  $\forall \theta_1, \theta_2 \in [0, \pi/2]$  with  $\theta_1 + \theta_2 \leq \pi/2$ ,

$$\sin^2(\theta_1 + \theta_2) = \sin^2 \theta_1 + \sin^2 \theta_2 + 2 \sin \theta_1 \sin \theta_2 \cos(\theta_1 + \theta_2) \geq \sin^2 \theta_1 + \sin^2 \theta_2. \quad (\text{C.19})$$

So we have that when  $\sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) \leq \pi/2$ ,

$$\sin^2 \angle(\mathbf{u}, \mathbf{v}) \leq \sin^2 \sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega) \implies \angle(\mathbf{u}, \mathbf{v}) \leq \sum_{\Omega \in \binom{[n]}{2}} \angle(\mathbf{u}_\Omega, \mathbf{v}_\Omega), \quad (\text{C.20})$$

as claimed.  $\blacksquare$

**Lemma C.3** (Covering in maximum length-2 angles). *For any  $\eta \in (0, 1/3)$ , there exists a subset  $\mathcal{Q} \subset \mathbb{S}^{n-1}$  of size at most  $(5n \log(1/\eta)/\eta)^{2n-1}$  satisfying the following: for any  $\mathbf{p} \in \mathbb{S}^{n-1}$ , there exists some  $\mathbf{q} \in \mathcal{Q}$  such that  $\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \eta$  for all  $\Omega \subset [n]$  with  $|\Omega| \leq 2$ .*

**Proof.** Define

$$d_2(\mathbf{p}, \mathbf{q}) = \max_{|\Omega| \leq 2} \angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega), \quad (\text{C.21})$$

our goal is to give an  $\eta$ -covering of  $\mathbb{S}^{n-1}$  in the  $d_2$  metric.

**Step 1** We partition  $\mathbb{S}^{n-1}$  according to the support, the sign pattern, and the ordering of the non-zero elements. For each configuration, we are going to construct a covering with the same configuration of support, sign pattern, and ordering. There are no more than  $3^n \cdot n!$  such configurations. Note that we only need to construct one such covering for each support size, and for each support size we can ignore the zero entries – the angle  $\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega)$  is always zero when  $\mathbf{p}, \mathbf{q}$  have matching support and  $\Omega$  contains at least one zero index.

Therefore, the task reduces to bounding the covering number of

$$A_n = \{\mathbf{p} \in \mathbb{S}^{n-1} : \mathbf{p} \geq 0, 0 < p_1 \leq p_2 \leq \dots \leq p_n\} \quad (\text{C.22})$$

in  $d_2$  for all  $n$ .

**Step 2** We bound the covering number of  $A_n$  by induction. Suppose that

$$N(A_{n'}) \leq \left( \frac{5 \log(1/\eta)}{\eta} \cdot n' \right)^{n'-1} = (C_\eta n')^{n'-1} \quad (\text{C.23})$$

holds for all  $n' \leq n-1$ . (The base case  $m=2$  clearly holds.) Let  $\mathcal{C}_{n'} \subset \mathbb{S}^{n'-1}$  be the corresponding covering sets.

We now construct a covering for  $A_n$ . Let  $R \doteq 1/\eta = r^k$  for some  $r \geq 1$  and  $k$  to be determined. Consider the set

$$\mathcal{Q}_{r,k} = \left\{ \mathbf{q} \in \mathbb{S}^{n-1} : \frac{q_{i+1}}{q_i} \in \{1, r, r^2, \dots, r^{k-1}\} \text{ for all } 1 \leq i \leq n-1 \right\}. \quad (\text{C.24})$$

We claim that  $\mathcal{Q}_{r,k}$  with properly chosen  $(r, k)$  gives a covering of

$$A_{n,R} = \left\{ \mathbf{p} \in A_n : \frac{p_{i+1}}{p_i} \leq R, \forall i \right\} \subset A_n. \quad (\text{C.25})$$

Indeed, we can decompose  $[1, R]$  into  $[1, r), [r, r^2), \dots, [r^{k-1}, R]$ . Each consecutive ratio  $p_{i+1}/p_i$  falls in one of these intervals, and we choose  $\mathbf{q}$  so that  $q_{i+1}/q_i$  is the left endpoint of this interval. Such a  $\mathbf{q}$  satisfies  $\mathbf{q} \in \mathcal{Q}_{r,k}$  and

$$\frac{p_{i+1}/p_i}{q_{i+1}/q_i} \in [1, r) \text{ for all } i \in [n-1]. \quad (\text{C.26})$$

By multiplying these bounds, we obtain that for all  $1 \leq i < j \leq n$ ,

$$\frac{p_j/p_i}{q_j/q_i} \in [1, r^{n-1}). \quad (\text{C.27})$$

Take  $r = 1 + \eta/2n$ , we have  $r^{n-1} = (1 + \eta/2n)^{n-1} \leq \exp(\eta/2) \leq 1 + \eta$ . Therefore, for all  $i, j$ , we have  $\frac{p_j/p_i}{q_j/q_i} \in [1, 1 + \eta)$ , which further implies that  $\angle((p_i, p_j), (q_i, q_j)) \leq \eta$  by [Lemma F.4](#). Thus we have for all  $|\Omega| \leq 2$  that  $\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \eta$ . (The size-1 angles are all zero as we have sign match.)

For this choice of  $r$ , we have  $k = \log R / \log r$  and thus

$$|\mathcal{Q}_{r,k}| = k^{n-1} = \left( \frac{\log R}{\log r} \right)^{n-1} = \left( \frac{\log(1/\eta)}{\log(1 + \eta/(2n))} \right)^{n-1} \leq \left( \frac{4n \log(1/\eta)}{\eta} \right)^{n-1} \doteq \tilde{N}_n, \quad (\text{C.28})$$

and we have  $N(A_{n,R}) \leq \tilde{N}_n$ .

**Step 3** We now construct the covering of  $A_n \setminus A_{n,R}$ . For any  $\mathbf{p} \in A_n \setminus A_{n,R}$ , there exists some  $i$  such that  $p_{i+1}/p_i \in [R, \infty)$ , which means that the angle of the ray  $(p_i, p_{i+1})$  is in between  $[\arctan(R), \pi/2) = [\pi/2 - \eta, \pi/2)$ . As  $\mathbf{p}$  is sorted, we have that

$$\frac{p_{i+j}}{p_{i-l}} \geq R \text{ for all } j \geq 1, l \geq 0, \quad (\text{C.29})$$

So if we take  $\mathbf{q}$  such that  $q_{i+1}/q_i \in [R, \infty)$ ,  $\mathbf{q}$  also has the above property, which gives that

$$\angle((p_{i-l}, p_{i+j}), (q_{i-l}, q_{i+j})) \leq \pi/2 - (\pi/2 - \eta) = \eta \text{ for all } j \geq 1, l \geq 0. \quad (\text{C.30})$$

Therefore to obtain the cover in  $d_2$ , we only need to consider the angles for  $\Omega \subset \{1, \dots, i\}$  and  $\Omega \subset \{i+1, \dots, n\}$ , which can be done by taking the product of the covers in  $A_i$  and  $A_{n-i}$ .

By considering all  $i \in \{1, \dots, n-1\}$ , we obtain the bound

$$N(A_n \setminus A_{n,R}) \leq \sum_{i=1}^{n-1} N(A_i) N(A_{n-i}). \quad (\text{C.31})$$

**Step 4** Putting together Step 2 and Step 3 and using the inductive assumption, we get that

$$N(A_n) \leq N(A_{n,R}) + N(A_n \setminus A_{n,R}) \leq \tilde{N}_n + \sum_{i=1}^{n-1} N(A_i)N(A_{n-i}) \quad (\text{C.32})$$

$$\leq \left( \frac{4n \log(1/\eta)}{\eta} \right)^{n-1} + \sum_{i=1}^{n-1} (C_\eta i)^{i-1} (C_\eta (n-i))^{n-i-1} \quad (\text{C.33})$$

$$\leq \left( \frac{4}{5} \right)^{n-1} (C_\eta n)^{n-1} + (n-1) \cdot C_\eta^{n-2} n^{n-2} \quad (\text{C.34})$$

$$\leq \left( \left( \frac{4}{5} \right)^{n-1} + \frac{1}{C_\eta} \right) (C_\eta n)^{n-1} \leq (C_\eta n)^{n-1}. \quad (\text{C.35})$$

This shows the case for  $m = n$  and completes the induction.

**Step 5** Considering all configurations of {support, sign pattern, ordering}, we have

$$N(\mathbb{S}^{n-1}) \leq 3^n \cdot n! \cdot N(A_n) \leq (3n)^n \left( \frac{5 \log(1/\eta)}{\eta} n \right)^{n-1} \leq \left( \frac{5 \log(1/\eta)}{\eta} n \right)^{2n-1}. \quad (\text{C.36})$$

■

**Lemma C.4** (Covering number in the  $d_{\mathbb{E}}$  metric). *Assume  $n \geq 3$ . There exists a numerical constant  $C > 0$  such that for any  $\epsilon \in (0, 1)$ ,  $\mathbb{S}^{n-1}$  admits an  $\epsilon$ -net of size  $\exp(Cn \log \frac{n}{\epsilon})$  w.r.t.  $d_{\mathbb{E}}$  defined in Eq. (C.1): for any  $\mathbf{p} \in \mathbb{S}^{n-1}$ , there exists a  $\mathbf{q}$  in the net with  $\text{supp}(\mathbf{q}) = \text{supp}(\mathbf{p})$  and  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$ . We say such  $\epsilon$  nets are admissible for  $\mathbb{S}^{n-1}$  wrt  $d_{\mathbb{E}}$ .*

**Proof.** Let  $\eta = \epsilon/n^2$ . By Lemma C.3, there exists a subset  $\mathcal{Q} \subset \mathbb{S}^{n-1}$  of size at most

$$\left( \frac{5n \log(1/\eta)}{\eta} \right)^{2n-1} = \left( \frac{5n^3 \log(n^2/\epsilon)}{\epsilon} \right)^{2n-1} \leq \exp\left(Cn \log \frac{n}{\epsilon}\right) \quad (\text{C.37})$$

such that for any  $\mathbf{p} \in \mathbb{S}^{n-1}$ , there exists  $\mathbf{q} \in \mathcal{Q}$  such that  $\text{supp}(\mathbf{p}) = \text{supp}(\mathbf{q})$  and  $\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \eta$  for all  $|\Omega| \leq 2$ . In particular, the  $|\Omega| = 1$  case says that  $\text{sign}(\mathbf{p}) = \text{sign}(\mathbf{q})$ , which implies that

$$\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \pi/2 \quad \forall \Omega \in \{0, 1\}^n. \quad (\text{C.38})$$

Thus, applying the vector angle inequality (Lemma C.2), for any  $\mathbf{p} \in \mathbb{S}^{n-1}$  and the corresponding  $\mathbf{q} \in \mathcal{Q}$ , we have

$$\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \sum_{|\Omega'|=2, \Omega' \subset \Omega} \angle(\mathbf{p}_{\Omega'}, \mathbf{q}_{\Omega'}) \leq 2 \binom{|\Omega|}{2} \eta \leq |\Omega|^2 \eta \quad \forall \Omega \text{ with } 3 \leq |\Omega| \leq n. \quad (\text{C.39})$$

Summing up, we get

$$\angle(\mathbf{p}_\Omega, \mathbf{q}_\Omega) \leq \max(2, |\Omega|^2) \eta \leq n^2 \eta = \epsilon \quad \forall \Omega. \quad (\text{C.40})$$

Therefore  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$ . ■

Below we establish the "Lipschitz" property in terms of  $d_{\mathbb{E}}$  distance.

**Lemma C.5.** *Fix  $\theta \in (0, 1)$ . For any  $\epsilon \in (0, 1)$ , let  $N_\epsilon$  be an admissible  $\epsilon$ -net for  $\mathbb{S}^{n-1}$  wrt  $d_{\mathbb{E}}$ . Let  $\mathbf{x}_1, \dots, \mathbf{x}_m$  be iid copies of  $\mathbf{x} \sim_{iid} \text{BG}(\theta)$  in  $\mathbb{R}^n$ . When  $m \geq C\epsilon^{-2}n$ , the inequality*

$$\sup_{\substack{\mathbf{p} \in \mathbb{S}^{n-1}, \mathbf{q} \in N_\epsilon \\ \text{supp}(\mathbf{p}) = \text{supp}(\mathbf{q}), d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon}} R(\mathbf{p}, \mathbf{q}) \doteq \frac{1}{m} \sum_{i=1}^m \mathbb{1} \{ \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}_i) \} \leq 2\epsilon \quad (\text{C.41})$$

holds with probability at least  $1 - \exp(-c\epsilon^2 m)$ . Here  $C$  and  $c$  are universal constants independent of  $\epsilon$ .

**Proof.** We call any pair of  $\mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}$  with  $\mathbf{q} \in N_\epsilon$ ,  $\text{supp}(\mathbf{p}) = \text{supp}(\mathbf{q})$ , and  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$  an admissible pair. Over any admissible pair  $(\mathbf{p}, \mathbf{q})$ ,  $\mathbb{E}[R] = d_{\mathbb{E}}(\mathbf{p}, \mathbf{q})$ . We next bound the deviation  $R - \mathbb{E}[R]$  uniformly over all admissible  $(\mathbf{p}, \mathbf{q})$  pairs. Observe that the process  $R$  is the sample average of  $m$  indicator functions. Define the hypothesis class

$$\mathcal{H} = \{\mathbf{x} \mapsto \mathbb{1}\{\text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x})\} : (\mathbf{p}, \mathbf{q}) \text{ is an admissible pair}\}. \quad (\text{C.42})$$

and let  $d_{\text{vc}}(\mathcal{H})$  be the VC-dimension of  $\mathcal{H}$ . From concentration results for VC-classes (see, e.g., Eq (3) and Theorem 3.4 of [Boucheron et al. \(2005\)](#)), we have

$$\mathbb{P}\left[\sup_{(\mathbf{p}, \mathbf{q}) \text{ admissible}} \{R(\mathbf{p}, \mathbf{q}) - \mathbb{E}[R](\mathbf{p}, \mathbf{q})\} \geq C_0 \sqrt{\frac{d_{\text{vc}}(\mathcal{H})}{m}} + t\right] \leq \exp(-mt^2) \quad (\text{C.43})$$

for any  $t > 0$ . It remains to bound the VC-dimension  $d_{\text{vc}}(\mathcal{H})$ . First, we have

$$d_{\text{vc}}(\mathcal{H}) \leq d_{\text{vc}}\{\mathbf{x} \mapsto \mathbb{1}\{\text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x})\} : \mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}\}. \quad (\text{C.44})$$

Observe that each set in the latter hypothesis class can be written as

$$\begin{aligned} & \{\mathbf{x} \mapsto \mathbb{1}\{\text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x})\} : \mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}\} \\ &= \{\mathbf{x} \mapsto \mathbb{1}\{\mathbf{p}^\top \mathbf{x} > 0, \mathbf{q}^\top \mathbf{x} \leq 0\} : \mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}\} \cup \{\mathbf{x} \mapsto \mathbb{1}\{\mathbf{p}^\top \mathbf{x} \geq 0, \mathbf{q}^\top \mathbf{x} < 0\} : \mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}\} \end{aligned} \quad (\text{C.45})$$

$$\cup \{\mathbf{x} \mapsto \mathbb{1}\{\mathbf{p}^\top \mathbf{x} < 0, \mathbf{q}^\top \mathbf{x} \geq 0\} : \mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}\} \cup \{\mathbf{x} \mapsto \mathbb{1}\{\mathbf{p}^\top \mathbf{x} \leq 0, \mathbf{q}^\top \mathbf{x} > 0\} : \mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}\}. \quad (\text{C.46})$$

the union of intersections of two halfspaces. Thus, letting

$$\mathcal{H}_0 = \{\mathbf{x} \mapsto \mathbb{1}\{\mathbf{x}^\top \mathbf{z} \geq 0\} : \mathbf{z} \in \mathbb{R}^n\} \quad (\text{C.47})$$

be the class of halfspaces, we have

$$\mathcal{H} \subset (\mathcal{H}_0 \cap \mathcal{H}_0) \sqcup (\mathcal{H}_0 \cap \mathcal{H}_0) \sqcup (\mathcal{H}_0 \cap \mathcal{H}_0) \sqcup (\mathcal{H}_0 \cap \mathcal{H}_0). \quad (\text{C.48})$$

Note that  $\mathcal{H}_0$  has VC-dimension  $n + 1$ . Applying bounds on the VC-dimension of unions and intersections (Theorem 1.1, [Van Der Vaart & Wellner \(2009\)](#)), we get that

$$d_{\text{vc}}(\mathcal{H}) \leq C_1 d_{\text{vc}}(\mathcal{H}_0 \cap \mathcal{H}_0) \leq C_2 d_{\text{vc}}(\mathcal{H}_0) \leq C_3 n. \quad (\text{C.49})$$

Plugging this bound into [Eq. \(C.43\)](#), we can set  $t = \epsilon/2$  and make  $m$  large enough so that  $C_0 \sqrt{C_3} \sqrt{n/m} \leq \epsilon/2$ , completing the proof.  $\blacksquare$

## C.2 POINTWISE CONVERGENCE OF SUB-DIFFERENTIAL

**Proposition C.6** (Pointwise convergence). *For any fixed  $\mathbf{q} \in \mathbb{S}^{n-1}$ ,*

$$\mathbb{P}\left[d_{\text{H}}(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})) > C_a \sqrt{n/m} + C_b t / \sqrt{m}\right] \leq 2 \exp(-t^2) \quad \forall t > 0. \quad (\text{C.50})$$

Here  $C_a, C_b \geq 0$  are universal constants.

**Proof.** Recall that

$$d_{\text{H}}(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})) = \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_{\partial f(\mathbf{q})}(\mathbf{u}) - h_{\mathbb{E}[\partial f](\mathbf{q})}(\mathbf{u})| = \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_{\partial f(\mathbf{q})}(\mathbf{u}) - \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{u})|. \quad (\text{C.51})$$

Write  $X_{\mathbf{u}} \doteq h_{\partial f(\mathbf{q})}(\mathbf{u}) - \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{u})$  and consider the zero-mean random process  $\{X_{\mathbf{u}}\}$  defined on  $\mathbb{S}^{n-1}$ . For any  $\mathbf{u}, \mathbf{v} \in \mathbb{S}^{n-1}$ , we have

$$\|X_{\mathbf{u}} - X_{\mathbf{v}}\|_{\psi_2} = \|h_{\partial f(\mathbf{q})}(\mathbf{u}) - \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{u}) - h_{\partial f(\mathbf{q})}(\mathbf{v}) + \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{v})\|_{\psi_2} \quad (\text{C.52})$$

$$= C_0 \left\| \frac{1}{m} \sum_{i \in [m]} (h_{Q_i}(\mathbf{u}) - \mathbb{E}h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v}) + \mathbb{E}h_{Q_i}(\mathbf{v})) \right\|_{\psi_2} \quad (\text{C.53})$$

$$\leq C_1 \frac{1}{m} \left( \sum_{i \in [m]} \|h_{Q_i}(\mathbf{u}) - \mathbb{E}h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v}) + \mathbb{E}h_{Q_i}(\mathbf{v})\|_{\psi_2}^2 \right)^{1/2} \quad (\text{C.54})$$

$$\leq C_2 \frac{1}{m} \left( \sum_{i \in [m]} \|h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v})\|_{\psi_2}^2 \right)^{1/2} \quad (\text{centering}), \quad (\text{C.55})$$

where we write  $Q_i \doteq \text{sign}(\mathbf{q}^\top \mathbf{x}_i)$  for all  $i \in [m]$ . Next we estimate  $\|h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v})\|_{\psi_2}$ . By definition,

$$h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v}) = \sup_{\mathbf{z} \in Q_i} \langle \mathbf{z}, \mathbf{u} \rangle - \sup_{\mathbf{z}' \in Q_i} \langle \mathbf{z}', \mathbf{v} \rangle. \quad (\text{C.56})$$

If  $h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v}) \geq 0$  and let  $\mathbf{z}_* \doteq \arg \max_{\mathbf{z} \in Q_i} \langle \mathbf{z}, \mathbf{u} \rangle$ , we have

$$h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v}) \leq \langle \mathbf{z}_*, \mathbf{u} \rangle - \langle \mathbf{z}_*, \mathbf{v} \rangle = \langle \mathbf{z}_*, \mathbf{u} - \mathbf{v} \rangle, \quad (\text{C.57})$$

and

$$\|h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v})\|_{\psi_2} \leq \|\langle \mathbf{z}_*, \mathbf{u} - \mathbf{v} \rangle\|_{\psi_2} \leq \|\mathbf{x}_i^\top (\mathbf{u} - \mathbf{v})\|_{\psi_2} \leq C_3 \|\mathbf{u} - \mathbf{v}\|, \quad (\text{C.58})$$

where we have used [Lemma F.1](#) to obtain the last upper bound. If  $h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v}) \leq 0$ ,  $h_{Q_i}(\mathbf{v}) - h_{Q_i}(\mathbf{u}) \geq 0$  and we can use similar argument to conclude that

$$\|h_{Q_i}(\mathbf{u}) - h_{Q_i}(\mathbf{v})\|_{\psi_2} \leq C_3 \|\mathbf{u} - \mathbf{v}\|. \quad (\text{C.59})$$

So

$$\|X_{\mathbf{u}} - X_{\mathbf{v}}\|_{\psi_2} \leq \frac{C_4}{\sqrt{m}} \|\mathbf{u} - \mathbf{v}\|. \quad (\text{C.60})$$

Thus,  $\{X_{\mathbf{u}}\}$  is a centered random process with sub-Gaussian increments with a parameter  $C_4/\sqrt{m}$ . We can apply [Proposition A.2](#) to conclude that

$$\mathbb{P} \left[ \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_{\partial f(\mathbf{q})}(\mathbf{u}) - \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{u})| > C_5 \sqrt{n/m} + C_6 t / \sqrt{m} \right] \leq 2 \exp(-t^2) \quad \forall t > 0, \quad (\text{C.61})$$

which implies the claimed result.  $\blacksquare$

### C.3 PROOF OF [PROPOSITION 3.5](#) (UNIFORM CONVERGENCE)

Throughout the proof, we let  $c, C$  denote universal constants that could change from step to step.

Fix an  $\epsilon \in (0, 1/2)$  to be decided later. Let  $N_\epsilon$  be an admissible  $\epsilon$  net for  $\mathbb{S}^{n-1}$  wrt  $d_{\mathbb{E}}$ , with  $|N_\epsilon| \leq \exp(Cn \log(n/\epsilon))$  ([Lemma C.4](#)). By [Proposition C.6](#) and the union bound,

$$\mathbb{P}[\exists \mathbf{q} \in N_\epsilon, d_{\mathbb{H}}(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})) > t/3] \leq \exp\left(-cmt^2 + Cn \log \frac{n}{\epsilon}\right) \quad (\text{C.62})$$

provided that  $m \geq Ct^{-2}n$ .

For any  $\mathbf{p} \in \mathbb{S}^{n-1}$ , let  $\mathbf{q} \in N_\epsilon$  satisfy  $\text{supp}(\mathbf{q}) = \text{supp}(\mathbf{p})$  and  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$ . Then we have

$$d_{\mathbb{H}}(\partial f(\mathbf{p}), \mathbb{E}[\partial f](\mathbf{p})) \leq \underbrace{d_{\mathbb{H}}(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q}))}_{\text{I}} + \underbrace{d_{\mathbb{H}}(\mathbb{E}[\partial f](\mathbf{p}), \mathbb{E}[\partial f](\mathbf{q}))}_{\text{II}} + \underbrace{d_{\mathbb{H}}(\partial f(\mathbf{p}), \partial f(\mathbf{q}))}_{\text{III}} \quad (\text{C.63})$$

by the triangular inequality for the Hausdorff metric. By the preceding union bound, term I is bounded by  $t/3$  as long as the bad event does not happen. For term II, we have

$$\begin{aligned} & d_{\mathbb{H}}(\mathbb{E}[\partial f](\mathbf{p}), \mathbb{E}[\partial f](\mathbf{q})) \\ &= \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_{\mathbb{E}[\partial f](\mathbf{p})}(\mathbf{u}) - h_{\mathbb{E}[\partial f](\mathbf{q})}(\mathbf{u})| \end{aligned} \quad (\text{C.64})$$

$$= \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |\mathbb{E} [h_{\partial f(\mathbf{p})}(\mathbf{u}) - h_{\partial f(\mathbf{q})}(\mathbf{u})]| \quad (\text{C.65})$$

$$= \sqrt{\frac{\pi}{2}} \cdot \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |\mathbb{E} [\sup \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \mathbf{x}_i^\top \mathbf{u} - \sup \text{sign}(\mathbf{q}^\top \mathbf{x}_i) \mathbf{x}_i^\top \mathbf{u}]| \quad (\text{C.66})$$

$$\leq \sqrt{\frac{\pi}{2}} \cdot \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |\mathbb{E} [|\mathbf{x}_i^\top \mathbf{u}| \mathbb{1} \{\text{sign}(\mathbf{p}^\top \mathbf{x}_i) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}_i)\}]| \quad (\text{C.67})$$

$$\leq \sqrt{\frac{\pi}{2}} \cdot 3\epsilon \sqrt{\log \frac{1}{\epsilon}}. \quad (\text{C.68})$$

where the last line follows from [Lemma F.5](#). As long as  $\epsilon \leq ct/\sqrt{\log(1/t)}$ , the above term is upper bounded by  $t/3$ . For term III, we have

$$\begin{aligned} & d_{\text{H}}(\partial f(\mathbf{p}), \partial f(\mathbf{q})) \\ &= \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |h_{\partial f(\mathbf{p})}(\mathbf{u}) - h_{\partial f(\mathbf{q})}(\mathbf{u})| \end{aligned} \quad (\text{C.69})$$

$$= \sqrt{\frac{\pi}{2}} \cdot \frac{1}{m} \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} \left| \sum_{i \in [m]: \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}_i)} \sup \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \mathbf{x}_i^\top \mathbf{u} - \sup \text{sign}(\mathbf{q}^\top \mathbf{x}_i) \mathbf{x}_i^\top \mathbf{u} \right| \quad (\text{C.70})$$

$$= \sqrt{\frac{\pi}{2}} \cdot \frac{2}{m} \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} \left| \sum_{i \in [m]: \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}_i)} s_i \mathbf{x}_i^\top \mathbf{u} \right| \quad (s_i \in \{+1, -1\} \text{ dependent on } \mathbf{p}, \mathbf{q}, \mathbf{x}_i \text{ and } \mathbf{u}) \quad (\text{C.71})$$

$$= \sqrt{\frac{\pi}{2}} \cdot \frac{2}{m} \left\| \sum_{i \in [m]: \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}_i)} s_i \mathbf{x}_i \right\|. \quad (\text{C.72})$$

By [Lemma C.5](#), with probability at least  $1 - \exp(-c\epsilon^2 m)$ , the number of different signs is upper bounded by  $2m\epsilon$  for all  $\mathbf{p}, \mathbf{q}$  such that  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon$ . On this good event, the above quantity can be upper bounded as follows. Define a set  $T \doteq \{\mathbf{s} \in \mathbb{R}^m : s_i \in \{+1, -1, 0\}, \|\mathbf{s}\|_0 \leq 2m\epsilon\}$  and consider the quantity  $\sup_{\mathbf{s} \in T} \|\mathbf{X}\mathbf{s}\|$ , where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_m]$ . Then,

$$\left\| \sum_{i \in [m]: \text{sign}(\mathbf{p}^\top \mathbf{x}_i) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}_i)} s_i \mathbf{x}_i \right\| \leq \sup_{\mathbf{s} \in T} \|\mathbf{X}\mathbf{s}\| \quad (\text{C.73})$$

uniformly (i.e., independent of  $\mathbf{p}, \mathbf{q}$  and  $\mathbf{u}$ ). We have

$$w(T) = \mathbb{E} \sup_{\mathbf{s} \in T} \mathbf{s}^\top \mathbf{g} = \mathbb{E} \sup_{K \subset [m], |K| \leq 2m\epsilon} \sum_{i \in K} |g_i| \quad (\text{C.74})$$

$$\leq 2m\epsilon \mathbb{E} \|\mathbf{g}\|_\infty \leq 4m\sqrt{\log m}\epsilon, \quad (\text{here } \mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_m)) \quad (\text{C.75})$$

$$\text{rad}(T) = \sqrt{2m\epsilon}. \quad (\text{C.76})$$

Noting that  $1/\sqrt{\theta} \cdot \mathbf{X}$  has independent, isotropic, and sub-Gaussian rows with a parameter  $C/\sqrt{\theta}$ , we apply [Proposition A.3](#) and obtain that

$$\sup_{\mathbf{s} \in T} \|\mathbf{X}\mathbf{s}\| \leq \sqrt{\theta n} \sqrt{2m\epsilon} + \frac{C}{\sqrt{\theta}} \left( 4m\sqrt{\log m}\epsilon + t_0 \sqrt{2m\epsilon} \right) \quad (\text{C.77})$$

with probability at least  $1 - 2\exp(-t_0^2)$ . So we have over all admissible  $(\mathbf{p}, \mathbf{q})$  pairs,

$$d_{\text{H}}(\partial f(\mathbf{p}), \partial f(\mathbf{q})) \leq \sqrt{\frac{\pi}{2}} \cdot \frac{2}{m} \left[ \sqrt{\theta n} \sqrt{2m\epsilon} + \frac{C}{\sqrt{\theta}} \left( 4m\sqrt{\log m}\epsilon + t_0 \sqrt{2m\epsilon} \right) \right] \quad (\text{C.78})$$

$$= \sqrt{\frac{\pi}{2}} \cdot \left( \sqrt{\frac{8\theta n\epsilon}{m}} + C\epsilon \sqrt{\frac{\log m}{\theta}} + Ct_0 \sqrt{\frac{8\epsilon}{m}} \right). \quad (\text{C.79})$$



Setting  $t_0 = ct\sqrt{m}$  and  $\epsilon = ct\sqrt{\theta/\log m}$ , we have that

$$d_H(\partial f(\mathbf{p}), \partial f(\mathbf{q})) \leq \frac{t}{3}, \quad (\text{C.80})$$

provided that  $m \geq Cct^{-2}n = Ct^{-1}n\sqrt{\theta/\log m}$ , which is subsumed by the earlier requirement  $m \geq Ct^{-2}n$ .

Putting together the three bounds Eq. (C.62), Eq. (C.67), Eq. (C.80), we can choose

$$\epsilon = ct\sqrt{\frac{\theta}{\log(m/t)}} \leq ct \cdot \min \left\{ \sqrt{\frac{\theta}{\log m}}, \frac{1}{\sqrt{\log(1/t)}} \right\} \quad (\text{C.81})$$

and get that  $d_H(\partial f(\mathbf{p}), \mathbb{E}[\partial f](\mathbf{p})) \leq t$  with probability at least

$$1 - 2 \exp(-cmt^2) - \exp(-cm\epsilon^2) - \exp(-cmt^2 + Cn \log \frac{n}{\epsilon}) \quad (\text{C.82})$$

$$\geq 1 - 2 \exp(-cmt^2) - \exp\left(-\frac{cm\theta t^2}{\log(m/t)}\right) - \exp\left(-cmt^2 + Cn \log \frac{n \log(m/t)}{\theta t}\right) \quad (\text{C.83})$$

$$\geq 1 - \exp\left(-\frac{cm\theta t^2}{\log(m/t)}\right) \quad (\text{C.84})$$

provided that  $m \geq Cnt^{-2} \log \frac{n \log(m/t)}{\theta t}$ . A sufficient condition is that  $m \geq Cnt^{-2} \log^2(n/t)$  for sufficiently large  $C$ . When this is satisfied, the probability is further lower bounded by  $1 - \exp(-cm\theta t^2/\log m)$ .

#### C.4 PROOF OF THEOREM 3.6

Define

$$t = \frac{1}{32n^{3/2}}\theta(1-\theta)\zeta_0 \leq \min \left\{ \frac{1}{8n^{3/2}}\theta(1-\theta)\frac{\zeta_0}{1+\zeta_0}, \frac{2-\sqrt{2}}{16n^{3/2}}\theta(1-\theta)\zeta_0 \right\}. \quad (\text{C.85})$$

By Proposition 3.5, with probability at least  $1 - \exp(-cm\theta^3\zeta_0^2n^{-3}\log^{-1}m)$  we have

$$d_H(\mathbb{E}[\partial f](\mathbf{q}), \partial f(\mathbf{q})) \leq t, \quad (\text{C.86})$$

provided that  $m \geq Cn^4\theta^{-2}\zeta_0^{-2}\log(n/\zeta_0)$ . We now show the properties Eq. (3.12) and Eq. (3.13) on this good event, focusing on  $\mathcal{S}_{\zeta_0}^{(n+)}$  but obtaining the same results for all other  $2n-1$  subsets by the same arguments.

For Eq. (3.12), we have

$$\langle \partial_{Rf}(\mathbf{q}), \mathbf{e}_j/q_j - \mathbf{e}_n/q_n \rangle = \langle \partial f(\mathbf{q}), (\mathbf{I} - \mathbf{q}\mathbf{q}^\top)(\mathbf{e}_j/q_j - \mathbf{e}_n/q_n) \rangle = \langle \partial f(\mathbf{q}), \mathbf{e}_j/q_j - \mathbf{e}_n/q_n \rangle. \quad (\text{C.87})$$

Now

$$\sup \langle \partial f(\mathbf{q}), \mathbf{e}_n/q_n - \mathbf{e}_j/q_j \rangle = h_{\partial f(\mathbf{q})}(\mathbf{e}_n/q_n - \mathbf{e}_j/q_j) \quad (\text{C.88})$$

$$= \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{e}_n/q_n - \mathbf{e}_j/q_j) - \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{e}_n/q_n - \mathbf{e}_j/q_j) + h_{\partial f(\mathbf{q})}(\mathbf{e}_n/q_n - \mathbf{e}_j/q_j) \quad (\text{C.89})$$

$$\leq \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{e}_n/q_n - \mathbf{e}_j/q_j) + \|\mathbf{e}_n/q_n - \mathbf{e}_j/q_j\| \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |\mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{u}) - h_{\partial f(\mathbf{q})}(\mathbf{u})| \quad (\text{C.90})$$

$$= \sup \langle \mathbb{E}[\partial f](\mathbf{q}), \mathbf{e}_n/q_n - \mathbf{e}_j/q_j \rangle + \|\mathbf{e}_n/q_n - \mathbf{e}_j/q_j\| d_H(\mathbb{E}[\partial f](\mathbf{q}), \partial f(\mathbf{q})). \quad (\text{C.91})$$

By Theorem 3.4(a),

$$\sup \langle \mathbb{E}[\partial f](\mathbf{q}), \mathbf{e}_n - q_n \mathbf{q} \rangle \leq -\frac{1}{2n}\theta(1-\theta)\frac{\zeta_0}{1+\zeta_0}. \quad (\text{C.92})$$

Moreover,  $\|\mathbf{e}_n/q_n - \mathbf{e}_j/q_j\| = \sqrt{1/q_n^2 + 1/q_j^2} \leq \sqrt{1/q_n^2 + 3/q_n^2} \leq 2\sqrt{n}$ . Meanwhile, we have

$$d_H(\mathbb{E}[\partial f](\mathbf{q}), \partial f(\mathbf{q})) \leq t \leq \frac{1}{8n^{3/2}}\theta(1-\theta)\frac{\zeta_0}{1+\zeta_0}. \quad (\text{C.93})$$

We conclude that

$$\inf \langle \partial f(\mathbf{q}), \mathbf{e}_j/q_j - \mathbf{e}_n/q_n \rangle = -\sup \langle \partial f(\mathbf{q}), \mathbf{e}_n/q_n - \mathbf{e}_j/q_j \rangle \quad (\text{C.94})$$

$$\geq \frac{1}{2n}\theta(1-\theta)\frac{\zeta_0}{1+\zeta_0} - 2\sqrt{n} \cdot \frac{1}{4n}\theta(1-\theta)\frac{\zeta_0}{1+\zeta_0} \cdot \frac{1}{2\sqrt{n}} \quad (\text{C.95})$$

$$\geq \frac{1}{4n}\theta(1-\theta)\frac{\zeta_0}{1+\zeta_0}, \quad (\text{C.96})$$

as claimed.

For Eq. (3.13), we have by Theorem 3.4(b) that

$$\sup \langle \partial f(\mathbf{q}), \mathbf{e}_n - q_n \mathbf{q} \rangle = h_{\partial f(\mathbf{q})}(\mathbf{e}_n - q_n \mathbf{q}) \quad (\text{C.97})$$

$$= \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{e}_n - q_n \mathbf{q}) - \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{e}_n - q_n \mathbf{q}) + h_{\partial f(\mathbf{q})}(\mathbf{e}_n - q_n \mathbf{q}) \quad (\text{C.98})$$

$$\leq \mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{e}_n - q_n \mathbf{q}) + \|\mathbf{e}_n - q_n \mathbf{q}\| \sup_{\mathbf{u} \in \mathbb{S}^{n-1}} |\mathbb{E}h_{\partial f(\mathbf{q})}(\mathbf{u}) - h_{\partial f(\mathbf{q})}(\mathbf{u})| \quad (\text{C.99})$$

$$= \sup \langle \mathbb{E}[\partial f](\mathbf{q}), \mathbf{e}_n - q_n \mathbf{q} \rangle + \|\mathbf{q}_{-n}\| d_H(\mathbb{E}[\partial f](\mathbf{q}), \partial f(\mathbf{q})). \quad (\text{C.100})$$

As we are on the good event

$$d_H(\mathbb{E}[\partial f](\mathbf{q}), \partial f(\mathbf{q})) \leq t \leq \frac{2-\sqrt{2}}{16n^{3/2}} \cdot \theta(1-\theta)\zeta_0, \quad (\text{C.101})$$

we have

$$\inf \langle \partial f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle = -\sup \langle \partial f(\mathbf{q}), \mathbf{e}_n - q_n \mathbf{q} \rangle \quad (\text{C.102})$$

$$\geq \frac{1}{8}\theta(1-\theta)\zeta_0 n^{-3/2} \|\mathbf{w}\| - \|\mathbf{q}_{-n}\| \frac{2-\sqrt{2}}{16} \cdot \theta(1-\theta)\zeta_0 n^{-3/2} \quad (\text{C.103})$$

$$\geq \frac{\sqrt{2}}{16}\theta(1-\theta)\zeta_0 n^{-3/2} \|\mathbf{q}_{-n}\|. \quad (\text{C.104})$$

Noting that  $\|\mathbf{q}_{-n}\| \geq \frac{1}{\sqrt{2}}\|\mathbf{q} - \mathbf{e}_n\|$  for all  $\mathbf{q}$  with  $q_n \geq 0$  completes the proof.

### C.5 PROOF OF PROPOSITION 3.7

For any  $\mathbf{q} \in \mathbb{S}^{n-1}$ ,

$$\sup \|\partial f(\mathbf{q})\| = d_H(\{0\}, \partial f(\mathbf{q})) \leq d_H(\{0\}, \mathbb{E}[\partial f](\mathbf{q})) + d_H(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})) \quad (\text{C.105})$$

by the metric property of the Hausdorff metric. On one hand, we have

$$\sup \|\mathbb{E}[\partial f](\mathbf{q})\| = \sup \left\| \mathbb{E}_\Omega \left[ \frac{\mathbf{q}_\Omega}{\|\mathbf{q}_\Omega\|} \mathbb{1}\{\mathbf{q}_\Omega \neq \mathbf{0}\} + \{\mathbf{v}_\Omega : \|\mathbf{v}_\Omega\| \leq 1\} \mathbb{1}\{\mathbf{q}_\Omega = \mathbf{0}\} \right] \right\| \leq 1. \quad (\text{C.106})$$

On the other hand, by Proposition 3.5,

$$d_H(\partial f(\mathbf{q}), \mathbb{E}[\partial f](\mathbf{q})) \leq 1 \quad \forall \mathbf{q} \in \mathbb{S}^{n-1} \quad (\text{C.107})$$

with probability at least  $1 - \exp(-c_1 m \theta \log^{-1} m)$ , provided that  $m \geq C_2 n^2 \log n$  (simplified using  $\theta \geq 1/n$ ). Combining the two results complete the proof.

## C.6 ADDITIONAL GEOMETRIES ON THE EMPIRICAL OBJECTIVE

**Proposition C.7.** *On the good event in Proposition 3.7, for all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$ , we have*

$$f(\mathbf{q}) - f(\mathbf{e}_n) \leq 2\sqrt{n} \|\mathbf{q} - \mathbf{e}_n\|. \quad (\text{C.108})$$

**Proof.** We use Lebourg’s mean value theorem for locally Lipschitz functions<sup>5</sup>, i.e., Theorem 2.3.7 of (Clarke, 1990). It is convenient to work in the  $\mathbf{w}$  space here. By subdifferential chain rules,  $g(\mathbf{w})$  is locally Lipschitz over  $\{\mathbf{w} : \|\mathbf{w}\| < \sqrt{\frac{n-1}{n}}\}$ . Thus, we have

$$f(\mathbf{q}) - f(\mathbf{e}_n) = g(\mathbf{w}) - g(\mathbf{0}) = \langle \mathbf{v}, \mathbf{w} \rangle \quad (\text{C.109})$$

for a certain  $t_0 \in (0, 1)$  and a certain  $\mathbf{v} \in \partial g(t_0\mathbf{w})$ . Now for any  $\mathbf{q}$  and the corresponding  $\mathbf{w}$ ,

$$\langle \partial g(\mathbf{w}), \mathbf{w} \rangle = \frac{1}{q_n} \langle \partial_R f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle. \quad (\text{C.110})$$

It follows

$$\begin{aligned} \langle \mathbf{v}, \mathbf{w} \rangle &\leq \frac{1}{t_0} \sup \langle \partial g(t_0\mathbf{w}), t_0\mathbf{w} \rangle = \frac{1}{t_0} \frac{1}{q_n(t_0\mathbf{w})} \sup \langle \partial_R f(\mathbf{q}(t_0\mathbf{w})), q_n(t_0\mathbf{w}) \mathbf{q}(t_0\mathbf{w}) - \mathbf{e}_n \rangle \\ &\leq \sup \|\partial_R f(\mathbf{q}(t_0\mathbf{w}))\| \cdot \frac{\|q_n(t_0\mathbf{w}) \mathbf{q}(t_0\mathbf{w}) - \mathbf{e}_n\|}{t_0 q_n(t_0\mathbf{w})} \leq 2 \frac{\|q_n(t_0\mathbf{w}) \mathbf{q}(t_0\mathbf{w}) - \mathbf{e}_n\|}{t_0 q_n(t_0\mathbf{w})}, \end{aligned} \quad (\text{C.111})$$

where at the last inequality we have used Proposition 3.7. Continuing the calculation, we further have

$$\frac{\|q_n(t_0\mathbf{w}) \mathbf{q}(t_0\mathbf{w}) - \mathbf{e}_n\|}{t_0 q_n(t_0\mathbf{w})} = \frac{t_0 \|\mathbf{w}\|}{t_0 q_n(t_0\mathbf{w})} \leq \sqrt{n} \|\mathbf{w}\| \leq \sqrt{n} \|\mathbf{q} - \mathbf{e}_n\|, \quad (\text{C.112})$$

completing the proof.  $\blacksquare$

**Proposition C.8.** *Assume  $\theta \in [1/n, 1/2]$ . When  $m \geq C\theta^{-2}n \log n$ , with probability at least  $1 - \exp(-cm\theta^3 \log^{-1} m)$ , the following holds: for all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$  satisfying  $f(\mathbf{q}) - f(\mathbf{e}_n) \leq \frac{2}{25}\theta$ ,*

$$f(\mathbf{q}) - f(\mathbf{e}_n) \geq \frac{\sqrt{2}}{16} \theta (1 - \theta) \|\mathbf{q}_{-n}\| \geq \frac{1}{16} \theta (1 - \theta) \|\mathbf{q} - \mathbf{e}_n\|. \quad (\text{C.113})$$

Here  $C, c > 0$  are universal constants.

**Proof.** We first establish uniform convergence of  $f(\mathbf{p})$  to  $\mathbb{E}[f](\mathbf{p})$ . Consider the zero-centered random process  $X_{\mathbf{p}} \doteq f(\mathbf{p}) - \mathbb{E}[f](\mathbf{p})$  on  $\mathbb{S}^{n-1}$ . Similar to proof of Proposition C.6, we can show that for all  $\mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}$

$$\|X_{\mathbf{p}} - X_{\mathbf{q}}\|_{\psi_2} \leq \frac{C}{\sqrt{m}} \|\mathbf{p} - \mathbf{q}\|. \quad (\text{C.114})$$

Applying Proposition A.2 gives that

$$\|f(\mathbf{p}) - \mathbb{E}[f](\mathbf{q})\| \leq \frac{1}{100} \theta \quad \forall \mathbf{q} \in \mathbb{S}^{n-1} \quad (\text{C.115})$$

with probability at least  $1 - \exp(-cm\theta^2)$ , provided that  $m \geq C\theta^{-2}n$ .

Now we consider  $\mathbb{E}[f](\mathbf{q}) - \mathbb{E}[f](\mathbf{e}_n)$ . For convenience, we first work in the  $\mathbf{w}$  space and note that  $\mathbb{E}[f](\mathbf{q}) - \mathbb{E}[f](\mathbf{e}_n) = \mathbb{E}[g](\mathbf{w}(\mathbf{q})) - \mathbb{E}[g](\mathbf{0})$ . By Lemma B.3,  $\mathbb{E}[g]$  is monotonically increasing in every radial direction  $\mathbf{v}$  until  $\|\mathbf{w}\|^2 + \|\mathbf{w}\|_{\infty}^2 \leq 1$ , which implies that

$$\inf_{\|\mathbf{w}\| \geq 1/2} \mathbb{E}[g](\mathbf{w}(\mathbf{q})) - \mathbb{E}[g](\mathbf{0}) = \inf_{\|\mathbf{w}\|=1/2} \mathbb{E}[g](\mathbf{w}(\mathbf{q})) - \mathbb{E}[g](\mathbf{0}). \quad (\text{C.116})$$

<sup>5</sup>It is possible to directly apply the manifold version of Lebourg’s mean value theorem, i.e., Theorem 3.3 of Hosseini & Pouryayevali (2011). We avoid this technicality by working with the Euclidean version in  $\mathbf{w}$  space.

For  $\mathbf{w}$  with  $\|\mathbf{w}\| = 1/2$ ,

$$\mathbb{E}[g](\mathbf{w}) - \mathbb{E}[g](\mathbf{0}) = (1 - \theta) \mathbb{E}_\Omega \|\mathbf{w}_\Omega\| + \theta \mathbb{E}_\Omega \sqrt{1 - \|\mathbf{w}_{\Omega^c}\|^2} - \theta \quad (\text{Lemma B.1}) \quad (\text{C.117})$$

$$\geq (1 - \theta) \theta \|\mathbf{w}\| + \theta \mathbb{E}_\Omega \sqrt{1 - \|\mathbf{w}\|^2} - \theta \quad (\text{C.118})$$

$$\geq \frac{1}{4}\theta + \frac{\sqrt{3}}{2}\theta - \theta \quad (\text{using } \theta \leq 1/2 \text{ and } \|\mathbf{w}\| = 1/2) \quad (\text{C.119})$$

$$\geq \frac{1}{10}\theta. \quad (\text{C.120})$$

So, back to the  $\mathbf{q}$  space,

$$\inf_{\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}: \|\mathbf{q}_{-n}\| \geq 1/2} \mathbb{E}[f](\mathbf{q}) - \mathbb{E}[f](\mathbf{0}) \geq \frac{1}{10}\theta. \quad (\text{C.121})$$

Combining the results in Eq. (C.115) and Eq. (C.121), we conclude that with high probability

$$\inf_{\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}: \|\mathbf{q}_{-n}\| \geq 1/2} f(\mathbf{q}) - f(\mathbf{0}) \geq \frac{2}{25}\theta. \quad (\text{C.122})$$

So when  $f(\mathbf{q}) - f(\mathbf{0}) \leq 2/25 \cdot \theta$ ,  $\|\mathbf{q}_{-n}\| \leq 1/2$ , which is equivalent to  $\|\mathbf{w}\| \leq 1/2$  in the  $\mathbf{w}$  space. Under this constraint, by Lemma B.3,

$$D_{-\mathbf{w}/\|\mathbf{w}\|}^c \mathbb{E}[g](\mathbf{w}) \leq -\theta(1 - \theta) \left( 1/\sqrt{1 + \|\mathbf{w}\|_\infty^2 / \|\mathbf{w}\|^2} - \|\mathbf{w}\| \right) \quad (\text{C.123})$$

$$\leq -\theta(1 - \theta) \left( \frac{1}{\sqrt{2}} - \frac{1}{2} \right) \leq -\frac{1}{5}\theta(1 - \theta). \quad (\text{C.124})$$

So, emulating the proof of Eq. (3.9) in Theorem 3.4, we have that for  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)}$  with  $\|\mathbf{q}_{-n}\| \leq 1/2$ ,

$$\langle \mathbb{E}[\partial_R f](\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle = q_n \langle \mathbb{E}[\partial g](\mathbf{w}), \mathbf{w} \rangle \geq q_n \|\mathbf{w}\| \cdot \frac{1}{5}\theta(1 - \theta) \geq \frac{\sqrt{3}}{10}\theta(1 - \theta) \|\mathbf{w}\|, \quad (\text{C.125})$$

where at the last inequality we use  $q_n = \sqrt{1 - \|\mathbf{w}\|^2} \geq \sqrt{3}/2$  when  $\|\mathbf{w}\| \leq 1/2$ . Moreover, we emulate the proof of Eq. (3.13) in Theorem 3.6 to obtain that

$$\inf \langle \partial_R f(\mathbf{q}), q_n \mathbf{q} - \mathbf{e}_n \rangle \geq \frac{\sqrt{2}}{16}\theta(1 - \theta) \|\mathbf{q}_{-n}\| \geq \frac{1}{16}\theta(1 - \theta) \|\mathbf{q} - \mathbf{e}_n\| \quad (\text{C.126})$$

with probability at least  $1 - \exp(-cm\theta^3 \log^{-1} m)$ , provided that  $m \geq C\theta^{-2}n \log n$ .

The last step of our proof is invoking the mean value theorem, similar to the proof of Proposition C.7. For any  $\mathbf{q}$ , we have

$$f(\mathbf{q}) - f(\mathbf{e}_n) = g(\mathbf{w}) - g(\mathbf{0}) = \langle \mathbf{v}, \mathbf{w} \rangle \quad (\text{C.127})$$

for a certain  $t \in (0, 1)$  and a certain  $\mathbf{v} \in \partial g(t\mathbf{w})$ . We have

$$\langle \mathbf{v}, \mathbf{w} \rangle \geq \frac{1}{t_0} \inf \langle \partial g(t_0 \mathbf{w}), t_0 \mathbf{w} \rangle = \frac{1}{t_0} \frac{1}{q_n(t_0 \mathbf{w})} \inf \langle \partial_R f(\mathbf{q}(t_0 \mathbf{w})), q_n(t_0 \mathbf{w}) \mathbf{q}(t_0 \mathbf{w}) - \mathbf{e}_n \rangle \quad (\text{C.128})$$

$$\geq \frac{1}{t_0} \frac{1}{q_n(t_0 \mathbf{w})} \frac{\sqrt{2}}{16} \theta(1 - \theta) \|t_0 \mathbf{w}\| \quad (\text{C.129})$$

$$\geq \frac{\sqrt{2}}{16} \theta(1 - \theta) \|\mathbf{w}\| \quad (\text{C.130})$$

$$\geq \frac{1}{16} \theta(1 - \theta) \|\mathbf{q} - \mathbf{e}_n\|, \quad (\text{C.131})$$

completing the proof.  $\blacksquare$

## D PROOFS FOR SECTION 3.3

### D.1 STAYING IN THE REGION $\mathcal{S}_{\zeta_0}^{(n+)}$

**Lemma D.1** (Progress in  $\mathcal{S}_{\zeta_0}^{(n+)} \setminus \mathcal{S}_1^{(n+)}$ ). *Set  $\eta = t_0/(100\sqrt{n})$  for  $t_0 \in (0, 1)$ . For any  $\zeta_0 \in (0, 1)$ , on the good events stated in [Proposition 3.7](#) and [Theorem 3.6](#), we have for all  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)} \setminus \mathcal{S}_1^{(n+)}$  and  $\mathbf{q}_+$  being the next step of Riemannian subgradient descent that*

$$\frac{q_{+,n}^2}{\|\mathbf{q}_{+,-n}\|_\infty^2} \geq \frac{q_n^2}{\|\mathbf{q}_{-n}\|_\infty^2} \left(1 + t \frac{\theta(1-\theta)\zeta_0}{400n^{3/2}(1+\zeta_0)}\right)^2. \quad (\text{D.1})$$

In particular, we have  $\mathbf{q}_+ \in \mathcal{S}_{\zeta_0}^{(n+)}$ .

**Proof.** We divide the index set  $[n-1]$  into three sets

$$\mathcal{I}_0 \doteq \{j \in [n-1] : q_j = 0\}, \quad (\text{D.2})$$

$$\mathcal{I}_1 \doteq \{j \in [n-1] : q_n^2/q_j^2 > 1 + 2\zeta_1 = 3, q_j \neq 0\} \quad (\text{D.3})$$

$$\mathcal{I}_2 \doteq \{j \in [n-1] : q_n^2/q_j^2 \leq 1 + 2\zeta_1 = 3\}. \quad (\text{D.4})$$

We perform different arguments on different sets. We let  $\mathbf{g}(\mathbf{q}) \in \partial_R f(\mathbf{q})$  be the subgradient taken at  $\mathbf{q}$  and note by [Proposition 3.7](#) that  $\|\mathbf{g}\| \leq 2$ , and so  $|g_i| \leq 2$  for all  $i \in [n]$ . We have

$$\frac{q_{+,n}^2}{q_{+,j}^2} = \frac{(q_n - \eta g_n)^2 / \|\mathbf{q} - \eta \mathbf{g}\|^2}{(q_j - \eta g_j)^2 / \|\mathbf{q} - \eta \mathbf{g}\|^2} = \frac{(q_n - \eta g_n)^2}{(q_j - \eta g_j)^2}. \quad (\text{D.5})$$

For any  $j \in \mathcal{I}_0$ ,

$$\frac{q_{+,n}^2}{q_{+,j}^2} = \frac{(q_n - \eta g_n)^2}{\eta^2 g_j^2} = q_n^2 \frac{(1 - \eta g_n/q_n)^2}{\eta^2 g_j^2} \geq \frac{(1 - 2\eta\sqrt{n})^2}{4n\eta^2}. \quad (\text{D.6})$$

Provided that  $\eta \leq 1/(4\sqrt{n})$ ,  $1 - 2\eta\sqrt{n} \geq 1/2$ , and so

$$\frac{(1 - 2\eta\sqrt{n})^2}{4n\eta^2} \geq \frac{1}{16n\eta^2} \geq \frac{5}{2}, \quad (\text{D.7})$$

where the last inequality holds when  $\eta \leq 1/\sqrt{40n}$ .

For any  $j \in \mathcal{I}_1$ ,

$$\frac{q_{+,n}^2}{q_{+,j}^2} \geq \frac{q_n^2 (1 - \eta g_n/q_n)^2}{q_j^2 + \eta^2 g_j^2} \geq \frac{q_n^2 (1 - \eta g_n/q_n)^2}{q_n^2/3 + 4\eta^2} = \frac{3(1 - \eta g_n/q_n)^2}{1 + 12\eta^2/q_n^2} \geq \frac{3(1 - 2\eta\sqrt{n})^2}{1 + 12n\eta^2} \geq \frac{5}{2}, \quad (\text{D.8})$$

where the very last inequality holds when  $\eta \leq 1/(26\sqrt{n})$ .

Since  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)} \setminus \mathcal{S}_1^{(n+)}$ ,  $\mathcal{I}_2$  is nonempty. For any  $j \in \mathcal{I}_2$ ,

$$\frac{q_{+,n}^2}{q_{+,j}^2} = \frac{q_n^2}{q_j^2} \left(1 + \eta \frac{g_j/q_j - g_n/q_n}{1 - \eta g_j/q_j}\right)^2. \quad (\text{D.9})$$

Since  $g_j/q_j \leq 2\sqrt{3n}$ ,  $1 - \eta g_j/q_j \geq 1/2$  when  $\eta \leq 1/(4\sqrt{3n})$ . Conditioned on this and due to that  $g_j/q_j - g_n/q_n \geq 0$ , it follows

$$\left(1 + \eta \frac{g_j/q_j - g_n/q_n}{1 - \eta g_j/q_j}\right)^2 \leq [1 + 2\eta(g_j/q_j - g_n/q_n)]^2 \leq [1 + 2\eta(2\sqrt{3n} + 2\sqrt{n})]^2 \leq (1 + 11\eta\sqrt{n})^2 \quad (\text{D.10})$$

If  $q_n^2/q_j^2 \leq 2$ ,  $q_{+,n}^2/q_{+,j}^2 \leq 5/2$  provided that

$$(1 + 11\eta\sqrt{n})^2 \leq \frac{5/2}{2} = \frac{5}{4} \iff \eta \leq \frac{1}{100\sqrt{n}}. \quad (\text{D.11})$$

As  $\mathbf{q} \notin S_1^{(n+)}$ , we have  $q_n^2/\|\mathbf{q}_{-n}\|_\infty^2 \leq 2$ , so there must be a certain  $j \in \mathcal{I}_2$  satisfying  $q_n^2/q_j^2 \leq 2$ . We conclude that when

$$\eta \leq \min \left\{ \frac{1}{\sqrt{40n}}, \frac{1}{26\sqrt{n}}, \frac{1}{100\sqrt{n}} \right\} = \frac{1}{100\sqrt{n}}, \quad (\text{D.12})$$

the index of largest entries of  $\mathbf{q}_{+,-n}$  remains in  $\mathcal{I}_2$ .

On the other hand, when  $\eta \leq 1/(100\sqrt{n})$ , for all  $j \in \mathcal{I}_2$ ,

$$\left(1 + \eta \frac{g_j/q_j - g_n/q_n}{1 - \eta g_j/q_j}\right)^2 \geq [1 + \eta(g_j/q_j - g_n/q_n)]^2 \geq \left(1 + \frac{\eta}{4n} \theta(1 - \theta) \frac{\zeta_0}{1 + \zeta_0}\right)^2. \quad (\text{D.13})$$

So when  $\eta = t/(100\sqrt{n})$  for any  $t \in (0, 1)$ ,

$$\frac{q_{+,n}^2}{\|\mathbf{q}_{+,-n}\|_\infty^2} \geq \frac{q_n^2}{\|\mathbf{q}_{-n}\|_\infty^2} \left(1 + t \frac{\theta(1 - \theta)\zeta_0}{400n^{3/2}(1 + \zeta_0)}\right)^2, \quad (\text{D.14})$$

completing the proof.  $\blacksquare$

**Proposition D.2.** For any  $\zeta_0 \in (0, 1)$ , on the good events stated in [Proposition 3.7](#) and [Theorem 3.6](#), if the step sizes satisfy

$$\eta^{(k)} \leq \min \left\{ \frac{1}{100\sqrt{n}}, \frac{1 - \zeta_0}{9\sqrt{n}} \right\} \text{ for all } k, \quad (\text{D.15})$$

the iteration sequence will stay in  $\mathcal{S}_{\zeta_0}^{(n+)}$  provided that our initialization  $\mathbf{q}^{(0)} \in \mathcal{S}_{\zeta_0}^{(n+)}$ .

**Proof.** By [Lemma D.1](#), if the current iterate  $\mathbf{q} \in \mathcal{S}_{\zeta_0}^{(n+)} \setminus \mathcal{S}_1^{(n+)}$ , the next iterate  $\mathbf{q}_+ \in \mathcal{S}_{\zeta_0}^{(n+)}$ , provided that  $\eta \leq 1/(100\sqrt{n})$ . Now if the current  $\mathbf{q} \in \mathcal{S}_1^{(n+)}$ , i.e.,  $q_n^2/q_j^2 \geq 2$  for all  $j \in [n - 1]$ , we can emulate the analysis of the set  $\mathcal{I}_1$  in proof of [Lemma D.1](#). Indeed, for any  $j \in [n - 1]$ ,

$$\frac{q_{+,n}^2}{q_{+,j}^2} \geq \frac{q_n^2(1 - \eta g_n/q_n)^2}{q_j^2 + \eta^2 g_j^2} \geq \frac{q_n^2(1 - 2\eta\sqrt{n})^2}{q_n^2/2 + 4\eta^2} \geq \frac{2(1 - 2\eta\sqrt{n})^2}{1 + 8n\eta^2} \geq 1 + \zeta_0, \quad (\text{D.16})$$

where the last inequality holds provided that  $\eta \leq (1 - \zeta_0)/(9\sqrt{n})$ . Combining the two cases finishes the proof.  $\blacksquare$

## D.2 PROOF OF [THEOREM 3.8](#)

As we have  $\eta^{(k)} \leq \frac{1}{100\sqrt{n}}$  and  $\mathbf{q}^{(0)} \in \mathcal{S}_{\zeta_0}^{(n+)}$ , the entire sequence  $\{\mathbf{q}^{(k)}\}_{k \geq 0}$  will stay in  $\mathcal{S}_{\zeta_0}^{(n+)}$  by [Proposition D.2](#).

For any  $\mathbf{q}$  and any  $\mathbf{v} \in \partial_R f(\mathbf{q})$ , we have  $\langle \mathbf{v}, \mathbf{q} \rangle = 0$  and therefore

$$\|\mathbf{q} - \eta\mathbf{v}\|^2 = \|\mathbf{q}\|^2 + \eta^2 \|\mathbf{v}\|^2 \geq 1. \quad (\text{D.17})$$

So  $\mathbf{q} - \eta\mathbf{v}$  is not inside  $\mathbb{B}^n$ . Since projection onto  $\mathbb{B}^n$  is a contraction, we have

$$\begin{aligned} \|\mathbf{q}_+ - \mathbf{e}_n\|^2 &= \left\| \frac{\mathbf{q} - \eta\mathbf{v}}{\|\mathbf{q} - \eta\mathbf{v}\|} - \mathbf{e}_n \right\|^2 \leq \|\mathbf{q} - \eta\mathbf{v} - \mathbf{e}_n\|^2 \\ &\leq \|\mathbf{q} - \mathbf{e}_n\|^2 + \eta^2 \|\mathbf{v}\|^2 - 2\eta \langle \mathbf{v}, \mathbf{q} - \mathbf{e}_n \rangle \leq \|\mathbf{q} - \mathbf{e}_n\|^2 + 4\eta^2 - \frac{1}{8}\eta\theta(1 - \theta)n^{-3/2}\zeta_0 \|\mathbf{q} - \mathbf{e}_n\|, \end{aligned} \quad (\text{D.18})$$

where we have used the bounds in [Proposition 3.7](#) and [Theorem 3.6](#) to obtain the last inequality. Further applying [Proposition C.7](#), we have

$$\|\mathbf{q}_+ - \mathbf{e}_n\|^2 \leq \|\mathbf{q} - \mathbf{e}_n\|^2 + 4\eta^2 - \frac{1}{16}\eta\theta(1-\theta)n^{-2}\zeta_0(f(\mathbf{q}) - f(\mathbf{e}_n)). \quad (\text{D.19})$$

Summing up the inequalities until step  $K$  (assumed  $\geq 5$ ), we have

$$0 \leq \|\mathbf{q}^{(K)} - \mathbf{e}_n\|^2 + 4 \sum_{j=0}^K (\eta^{(j)})^2 - \frac{1}{16}\theta(1-\theta)n^{-2}\zeta_0 \sum_{j=0}^K \eta^{(j)} (f(\mathbf{q}^{(j)}) - f(\mathbf{e}_n)) \quad (\text{D.20})$$

$$\implies \sum_{j=0}^K \eta^{(j)} (f(\mathbf{q}^{(j)}) - f(\mathbf{e}_n)) \leq \frac{16\|\mathbf{q}^{(K)} - \mathbf{e}_n\|^2 + 64 \sum_{j=0}^K (\eta^{(j)})^2}{\theta(1-\theta)n^{-2}\zeta_0} \quad (\text{D.21})$$

$$\implies f(\mathbf{q}^{\text{best}}) - f(\mathbf{e}_n) \leq \frac{16\|\mathbf{q}^{(K)} - \mathbf{e}_n\|^2 + 64 \sum_{j=0}^K (\eta^{(j)})^2}{\theta(1-\theta)n^{-2}\zeta_0 \sum_{j=0}^K \eta^{(j)}}. \quad (\text{D.22})$$

Substituting the following estimates

$$\sum_{j=0}^K (\eta^{(j)})^2 \leq \frac{1}{10^4 n} \left(1 + \int_0^K t^{-2\alpha} dt\right) \leq \frac{1}{10^4 n} \frac{1}{1-2\alpha} (K^{1-2\alpha} + 1), \quad (\text{D.23})$$

$$\sum_{j=0}^{K'} \eta^{(j)} \geq \frac{1}{10^2 \sqrt{n}} \int_0^K t^{-\alpha} dt \geq \frac{1}{10^2 \sqrt{n}} \frac{K^{1-\alpha}}{1-\alpha}, \quad (\text{D.24})$$

and noting  $16\|\mathbf{q}^{(K)} - \mathbf{e}_n\|^2 \leq 32$ , we have

$$f(\mathbf{q}^{\text{best}}) - f(\mathbf{e}_n) \leq \frac{3200n^{5/2}(1-\alpha) + 16/25 \cdot n^{3/2} \left(\frac{1-\alpha}{1-2\alpha} K^{1-2\alpha} + 1 - \alpha\right)}{\theta(1-\theta)\zeta_0 K^{1-\alpha}}. \quad (\text{D.25})$$

Noting that

$$K \geq \left(\frac{6400n^{5/2}(1-\alpha)}{\theta(1-\theta)\zeta_0\epsilon}\right)^{\frac{1}{1-\alpha}} \implies \frac{3200n^{5/2}(1-\alpha)}{\theta(1-\theta)\zeta_0 K^{1-\alpha}} \leq \frac{\epsilon}{2}, \quad (\text{D.26})$$

and when  $K \geq 1$ ,  $K^{1-2\alpha} \geq 1$ , yielding that

$$K \geq \left(\frac{64n^{3/2} \frac{1-\alpha}{1-2\alpha}}{25\epsilon\theta(1-\theta)\zeta_0}\right)^{\frac{1}{\alpha}} \implies \frac{32n^{3/2} \frac{1-\alpha}{1-2\alpha} K^{-\alpha}}{25\theta(1-\theta)\zeta_0} \leq \frac{\epsilon}{2} \implies \frac{16n^{3/2} \cdot (1-\alpha) \left(\frac{1}{1-2\alpha} K^{1-2\alpha} + 1\right)}{25\theta(1-\theta)\zeta_0 K^{1-\alpha}} \leq \frac{\epsilon}{2}. \quad (\text{D.27})$$

So we conclude that when

$$K \geq \max \left\{ \left(\frac{6400n^{5/2}(1-\alpha)}{\theta(1-\theta)\zeta_0\epsilon}\right)^{\frac{1}{1-\alpha}}, \left(\frac{64n^{3/2} \frac{1-\alpha}{1-2\alpha}}{25\epsilon\theta(1-\theta)\zeta_0}\right)^{\frac{1}{\alpha}} \right\}, \quad (\text{D.28})$$

$f(\mathbf{q}^{\text{best}}) - f(\mathbf{e}_n) \leq \epsilon$ . When this happens, by [Proposition C.8](#),

$$\|\mathbf{q}^{\text{best}} - \mathbf{e}_n\| \leq \frac{16}{\theta(1-\theta)}\epsilon. \quad (\text{D.29})$$

Plugging in the choice  $\zeta_0 = 1/(5 \log n)$  in [Eq. \(D.28\)](#) gives the desired bound on the number of iterations.

## E PROOFS FOR [SECTION 3.4](#)

### E.1 PROOF OF [LEMMA 3.9](#)

**Lemma E.1.** For all  $n \geq 3$  and  $\zeta \geq 0$ , it holds that

$$\frac{\text{vol}(\mathcal{S}_\zeta^{(n+)})}{\text{vol}(\mathbb{S}^{n-1})} \geq \frac{1}{2n} - \frac{9 \log n}{8n} \zeta. \quad (\text{E.1})$$

We note that a similar result appears in (Gilboa et al., 2018) but our definitions of the region  $\mathcal{S}_\zeta$  are slightly different. For completeness we provide a proof in Lemma F.3.

We now prove Lemma 3.9. Taking  $\zeta = 1/(5 \log n)$  in Lemma E.1, we obtain

$$\frac{\text{vol}(\mathcal{S}_{1/(5 \log n)}^{(n+)})}{\text{vol}(\mathbb{S}^{n-1})} \geq \frac{1}{2n} - \frac{9 \log n}{8n} \cdot \frac{1}{5 \log n} \geq \frac{1}{4n}. \quad (\text{E.2})$$

By symmetry, all the  $2n$  sets  $\{\mathcal{S}_{1/(5 \log n)}^{(i+)}, \mathcal{S}_{1/(5 \log n)}^{(i-)} : i \in [n]\}$  have the same volume which is at least  $1/(4n)$ . As  $\mathbf{q}^{(0)} \sim \text{Uniform}(\mathbb{S}^{n-1})$ , it falls into their union with probability at least  $2n \cdot 1/(4n) = 1/2$ , on which it belongs to a uniformly random one of these  $2n$  sets.

## E.2 PROOF OF THEOREM 3.10

Assume that the good event in Proposition 3.7 happens and that in Theorem 3.6 happens to all the  $2n$  sets  $\{\mathcal{S}_{1/(5 \log n)}^{(i+)}, \mathcal{S}_{1/(5 \log n)}^{(i-)} : i \in [n]\}$ , which by setting  $\zeta_0 = 1/(5 \log n)$  has probability at least

$$1 - \exp(-cm\theta^3 \zeta_0^2 n^{-3} \log^{-1} m) - \exp(-cm\theta \log^{-1} m) = 1 - \exp(-c'm\theta^3 n^{-3} \log m^{-3}). \quad (\text{E.3})$$

By Lemma 3.9, random initialization will fall these  $2n$  sets with probability at least  $1/2$ . When it falls in one of these  $2n$  sets, by Theorem 3.8, one run of the algorithm will find a signed standard basis vector up to  $\epsilon$  accuracy. With  $R$  independent runs, at least  $S \doteq \frac{1}{4}R$  of them are effective with probability at least  $1 - \exp(-(R/4)^2/(R/4 \cdot 2)) = 1 - \exp(-R/8)$ , due to Bernstein's inequality. After these effective runs, the probability any standard basis vector is missed (up to sign) is bounded by

$$n \left(1 - \frac{1}{n}\right)^S \leq \exp\left(-\frac{S}{n} + \log n\right) \leq \exp\left(-\frac{S}{2n}\right), \quad (\text{E.4})$$

where the second inequality holds whenever  $S \geq 2n \log n$ .

## F AUXILIARY CALCULATIONS

**Lemma F.1.** For  $x \sim \text{BG}(\theta)$ ,  $\|x\|_{\psi_2} \leq C_a$ . For any vector  $\mathbf{u} \in \mathbb{R}^n$  and  $\mathbf{x} \sim_{iid} \text{BG}(\theta)$ ,  $\|\mathbf{x}^\top \mathbf{u}\|_{\psi_2} \leq C_b \|\mathbf{u}\|$ . Here  $C_a, C_b \geq 0$  are universal constants.

**Proof.** For any  $\lambda \in \mathbb{R}$ ,

$$\exp(\lambda x) = \theta \exp(\lambda x) \leq \exp(\lambda x). \quad (\text{F.1})$$

So  $\|x\|_{\psi_2}$  is bounded by a universal constant. Moreover,

$$\|\mathbf{u}^\top \mathbf{x}\|_{\psi_2} = \left\| \sum_i u_i x_i \right\|_{\psi_2} \leq C_1 \left( \sum_i u_i^2 \|x_i\|_{\psi_2}^2 \right)^{1/2} \leq C_2 \|\mathbf{u}\|, \quad (\text{F.2})$$

as claimed. ■

**Lemma F.2.** Let  $\mathbf{a}_1, \dots, \mathbf{a}_m$  be iid copies of  $\mathbf{a} \sim_{iid} \text{BG}(\theta)$ . Then,

$$\mathbb{P} \left[ \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left| \sum_{i \in [m]} |\mathbf{q}^\top \mathbf{x}_i| - \mathbb{E} [|\mathbf{q}^\top \mathbf{x}|] \right| > C_a \sqrt{mn} + C_b \sqrt{mt} \right] \leq 2 \exp(-t^2). \quad (\text{F.3})$$

for any  $t \geq 0$ . Here  $C_a, C_b \geq 0$  are universal constants.

**Proof.** Consider the zero-centered random process defined on  $\mathbb{S}^{n-1}$ :  $X_{\mathbf{q}} \doteq \sum_{i \in [m]} (|\mathbf{q}^\top \mathbf{x}_i| - \mathbb{E} |\mathbf{q}^\top \mathbf{x}|)$ . Then, for any  $\mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}$ ,

$$\|X_{\mathbf{p}} - X_{\mathbf{q}}\|_{\psi_2} = \left\| \sum_{i \in [m]} (|\mathbf{p}^\top \mathbf{x}_i| - |\mathbf{q}^\top \mathbf{x}_i| - \mathbb{E} |\mathbf{p}^\top \mathbf{x}| + \mathbb{E} |\mathbf{q}^\top \mathbf{x}|) \right\|_{\psi_2} \quad (\text{F.4})$$



$$\leq C_1 \left( \sum_{i \in [m]} \left| |\mathbf{p}^\top \mathbf{x}_i| - |\mathbf{q}^\top \mathbf{x}_i| - \mathbb{E} |\mathbf{p}^\top \mathbf{x}| + \mathbb{E} |\mathbf{q}^\top \mathbf{x}| \right|_{\psi_2}^2 \right)^{1/2} \quad (\text{F.5})$$

$$\leq C_2 \left( \sum_{i \in [m]} \left| |\mathbf{p}^\top \mathbf{x}_i| - |\mathbf{q}^\top \mathbf{x}_i| \right|_{\psi_2}^2 \right)^{1/2} \quad (\text{centering}) \quad (\text{F.6})$$

$$\leq C_2 \left( \sum_{i \in [m]} \left\| (\mathbf{p} - \mathbf{q})^\top \mathbf{x}_i \right\|_{\psi_2}^2 \right)^{1/2} \quad (\text{F.7})$$

$$= C_3 \sqrt{m} \|\mathbf{p} - \mathbf{q}\|, \quad (\text{F.8})$$

where we use the estimate in [Lemma F.1](#) to obtain the last inequality. Note that  $X_{\mathbf{q}}$  is a mean-zero random process, and we can invoke [Proposition A.2](#) with  $w(\mathbb{S}^{n-1}) = C_4 \sqrt{n}$  and  $\text{rad}(\mathbb{S}^{n-1}) = 2$  to get the claimed result.  $\blacksquare$

**Lemma F.3.** *For all  $n \geq 3$  and  $\zeta \geq 0$ , it holds that*

$$\frac{\text{vol}(\mathcal{S}_\zeta^{(n+)})}{\text{vol}(\mathbb{S}^{n-1})} \geq \frac{1}{2n} - \frac{9 \log n}{8n} \zeta. \quad (\text{F.9})$$

**Proof.** We have

$$\frac{\text{vol}(\mathcal{S}_\zeta^{(n+)})}{\text{vol}(\mathbb{S}^{n-1})} = \mathbb{P}_{\mathbf{q} \sim \text{uniform}(\mathbb{S}^{n-1})} \left[ q_n^2 \geq (1 + \zeta) \|\mathbf{q}_{-n}\|_\infty^2, q_n \geq 0 \right] \quad (\text{F.10})$$

$$= \mathbb{P}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)} \left[ \mathbf{x}_n \geq 0, x_n^2 \geq (1 + \zeta) \mathbf{x}_i^2 \forall i \neq n \right] \quad (\text{F.11})$$

$$= (2\pi)^{n/2} \int_0^\infty e^{-x_n^2/2} \left( \prod_{j=1}^{n-1} \int_{-x_n/\sqrt{1+\zeta}}^{x_n/\sqrt{1+\zeta}} e^{-x_j^2/2} dx_j \right) dx_n \quad (\text{F.12})$$

$$= (2\pi)^{1/2} \int_0^\infty e^{-x_n^2/2} \psi^{n-1} \left( x_n / \sqrt{1 + \zeta} \right) dx_n \quad (\text{F.13})$$

$$= \frac{\sqrt{1 + \zeta}}{\sqrt{2\pi}} \int_0^\infty e^{-(1+\zeta)x^2/2} \psi^{n-1}(x) dx \doteq \tilde{h}(\zeta) > 0, \quad (\text{F.14})$$

where we write  $\psi(t) \doteq \frac{1}{\sqrt{2\pi}} \int_{-t}^t \exp(-s^2/2) ds$ . Now we derive a lower bound of the volume ratio by considering a first-order Taylor expansion of the last equation around  $\zeta = 0$  (as we are mostly interested in small  $\zeta$ ). By symmetry,  $\tilde{h}(0) = 1/(2n)$ . Moreover, we have

$$\left. \frac{\partial \tilde{h}(\zeta)}{\partial \zeta} \right|_{\zeta=0} = \frac{1}{2} \frac{1}{\sqrt{2\pi}} \int_0^\infty e^{-x^2/2} \psi^{n-1}(x) dx - \frac{1}{\sqrt{2\pi}} \int_0^\infty e^{-x^2/2} x^2 \psi^{n-1}(x) dx \quad (\text{F.15})$$

$$= \frac{1}{4n} - \frac{1}{2\sqrt{2\pi}} \int_0^\infty e^{-x^2/2} x^2 \psi^{n-1}(x) dx. \quad (\text{F.16})$$

Now we provide an upper bound for the second term of the last equation. Note that

$$\frac{1}{\sqrt{2\pi}} \int_0^\infty e^{-x^2/2} x^2 \psi^{n-1}(x) dx = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)} \left[ x_n^2 \mathbb{1} \left\{ x_n^2 \geq \|\mathbf{x}_{-n}\|_\infty^2 \right\} \mathbb{1} \{x_n \geq 0\} \right] \quad (\text{F.17})$$

$$= \frac{1}{2n} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n)} \|\mathbf{x}\|_\infty^2. \quad (\text{F.18})$$

Now for any  $\lambda \in (0, 1/2)$ ,

$$\exp\left(\lambda \mathbb{E} \|\mathbf{x}\|_\infty^2\right) \leq \mathbb{E} \exp\left(\lambda \|\mathbf{x}\|_\infty^2\right) \leq \sum_{j=1}^n \mathbb{E} \exp\left(\lambda x_j^2\right) = n \mathbb{E}_{x \sim \mathcal{N}(0,1)} \exp\left(\lambda x^2\right) \leq \frac{n}{\sqrt{1-2\lambda}}. \quad (\text{F.19})$$

Taking logarithm on both sides, rearranging the terms, and setting  $\lambda = 1/4$ , we obtain

$$\mathbb{E} \|\mathbf{x}\|_\infty^2 \leq \inf_{\lambda \in (0, 1/2)} \frac{\log n + \frac{1}{2} \log(1 - 2\lambda)^{-1}}{\lambda} \leq 4 \log n + 2 \log 2. \quad (\text{F.20})$$

So

$$\left. \frac{\partial \bar{h}(\zeta)}{\partial \zeta} \right|_{\zeta=0} \geq \frac{1}{4n} - \frac{1}{4n} (4 \log n + 2 \log 2) \geq -\frac{9 \log n}{8n}, \quad (\text{F.21})$$

provided that  $n \geq 3$ . Now we show that  $\bar{h}(\zeta) \geq \bar{h}(0) + \bar{h}'(0)\zeta$  by showing that  $\bar{h}''(\zeta) \geq 0$ . We have

$$\frac{\partial^2 \bar{h}(\zeta)}{\partial \zeta^2} = \frac{\sqrt{1+\zeta}}{4\sqrt{2\pi}} \int_0^\infty \left[ x^4 - \frac{2x^2}{1+\zeta} - \frac{1}{(1+\zeta)^2} \right] e^{-\frac{1+\zeta}{2}x^2} \psi^{n-1}(x) dx. \quad (\text{F.22})$$

Using integration by part, we have

$$\int_0^\infty \left[ x^4 - \frac{3x^2}{1+\zeta} \right] e^{-\frac{1+\zeta}{2}x^2} \psi^{n-1}(x) dx \quad (\text{F.23})$$

$$= -\frac{1}{1+\zeta} e^{-\frac{1+\zeta}{2}x^2} x^3 \cdot \psi^{n-1}(x) \Big|_0^\infty + \int_0^\infty \frac{1}{1+\zeta} e^{-\frac{1+\zeta}{2}x^2} x^3 (n-1) \psi^{n-2}(x) \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}} dx \quad (\text{F.24})$$

$$= \int_0^\infty \frac{1}{1+\zeta} e^{-\frac{1+\zeta}{2}x^2} x^3 (n-1) \psi^{n-2}(x) \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}} dx \geq 0, \quad (\text{F.25})$$

and similarly

$$\int_0^\infty \left[ x^2 - \frac{1}{1+\zeta} \right] e^{-\frac{1+\zeta}{2}x^2} \psi^{n-1}(x) dx \quad (\text{F.26})$$

$$= -\frac{1}{1+\zeta} e^{-\frac{1+\zeta}{2}x^2} x \cdot \psi^{n-1}(x) \Big|_0^\infty + \int_0^\infty \frac{1}{1+\zeta} e^{-\frac{1+\zeta}{2}x^2} x (n-1) \psi^{n-2}(x) \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}} dx \quad (\text{F.27})$$

$$= \int_0^\infty \frac{1}{1+\zeta} e^{-\frac{1+\zeta}{2}x^2} x (n-1) \psi^{n-2}(x) \sqrt{\frac{2}{\pi}} e^{-\frac{x^2}{2}} dx \geq 0. \quad (\text{F.28})$$

Noting that

$$x^4 - \frac{2x^2}{1+\zeta} - \frac{1}{(1+\zeta)^2} = x^4 - \frac{3x^2}{1+\zeta} + \frac{1}{1+\zeta} \left( x^2 - \frac{1}{1+\zeta} \right) \quad (\text{F.29})$$

and combining the above integral results, we conclude that  $\bar{h}''(\zeta) \geq 0$  and complete the proof.  $\blacksquare$

**Lemma F.4.** Let  $(x_1, y_1), (x_2, y_2) \in \mathbb{R}_{>0}^2$  be two points in the first quadrant satisfying  $y_1 \geq x_1$  and  $y_2 \geq x_2$ , and  $\frac{y_2/x_2}{y_1/x_1} \in [1, 1 + \eta]$  for some  $\eta \leq 1$ , then we have  $\angle((x_1, y_1), (x_2, y_2)) \leq \eta$ .

**Proof.** For  $i = 1, 2$ , let  $\theta_i$  be the angle between the ray  $(x_i, y_i)$  and the  $x$ -axis. Our assumption implies that  $\theta_i \in [\pi/4, \pi/2)$  and  $\theta_2 \geq \theta_1$ , thus  $\angle((x_1, y_1), (x_2, y_2)) = \theta_2 - \theta_1$ , so we have

$$\begin{aligned} & \tan \angle((x_1, y_1), (x_2, y_2)) \\ &= \frac{\tan \theta_2 - \tan \theta_1}{1 + \tan \theta_2 \tan \theta_1} \\ &= \frac{y_2/x_2 - y_1/x_1}{1 + y_2 y_1 / (x_2 x_1)} \\ &= \frac{\frac{y_2/x_2}{y_1/x_1} - 1}{y_2/x_2 + x_1/y_1} \\ &\leq \frac{y_2/x_2}{y_1/x_1} - 1 \\ &\leq \eta. \end{aligned} \quad (\text{F.30})$$

Therefore  $\angle((x_1, y_1), (x_2, y_2)) \leq \arctan(\eta) \leq \eta$ .  $\blacksquare$

**Lemma F.5.** For any  $\mathbf{p}, \mathbf{q} \in \mathbb{S}^{n-1}$  with the same support pattern such that  $d_{\mathbb{E}}(\mathbf{p}, \mathbf{q}) \leq \epsilon \leq \frac{1}{2}$ , we have for all  $\mathbf{u} \in \mathbb{S}^{n-1}$  that

$$\mathbb{E}_{\mathbf{x} \sim \text{BG}(\theta)} [\mathbf{u}^\top \mathbf{x} \mathbb{1} \{ \text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}) \}] \leq 3\epsilon \sqrt{\log \frac{1}{\epsilon}}. \quad (\text{F.31})$$

**Proof.** Fix some threshold  $t > 0$  to be determined. We have

$$\mathbb{E} [\mathbf{u}^\top \mathbf{x} \mathbb{1} \{ \text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}) \}] \quad (\text{F.32})$$

$$\leq \mathbb{E} [|\mathbf{u}^\top \mathbf{x}| \mathbb{1} \{ |\mathbf{u}^\top \mathbf{x}| > t \}] + \mathbb{E} [|\mathbf{u}^\top \mathbf{x}| \mathbb{1} \{ |\mathbf{u}^\top \mathbf{x}| \leq t, \text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}) \}] \quad (\text{F.33})$$

$$\leq (\mathbb{E} [(\mathbf{u}^\top \mathbf{x})^2] \cdot \mathbb{P} [|\mathbf{u}^\top \mathbf{x}| > t])^{1/2} + t \mathbb{E} [\mathbb{1} \{ \text{sign}(\mathbf{p}^\top \mathbf{x}) \neq \text{sign}(\mathbf{q}^\top \mathbf{x}) \}] \quad (\text{F.34})$$

$$\leq (\theta \cdot 2 \exp(-t^2/2))^{1/2} + \epsilon t. \quad (\text{F.35})$$

The second to last inequality uses Cauchy-Schwarz, and the last inequality uses the fact that  $\mathbf{u}^\top \mathbf{x} = \mathbf{u}_\Omega^\top \mathbf{x}_\Omega$  is  $\|\mathbf{u}_\Omega\|_2^2$ -sub-Gaussian conditioned on  $\Omega$  and thus 1-sub-Gaussian marginally. Taking  $t = \sqrt{2 \log \frac{1}{\epsilon^2}}$ , the above bound simplifies to

$$\sqrt{2\theta \exp\left(-\log \frac{1}{\epsilon^2}\right)} + \epsilon \sqrt{2 \log \frac{1}{\epsilon^2}} = \sqrt{2\theta} \epsilon \left(1 + \log \frac{1}{\epsilon^2}\right) = \epsilon \left(\sqrt{2\theta} + \sqrt{2 \log \frac{1}{\epsilon^2}}\right) \leq 3\epsilon \sqrt{\log \frac{1}{\epsilon}} \quad (\text{F.36})$$

where we have used  $\theta \leq 1/2$  and  $\epsilon \leq 1/2$ .  $\blacksquare$

## G RESULTS ON ORTHOGONAL DICTIONARIES

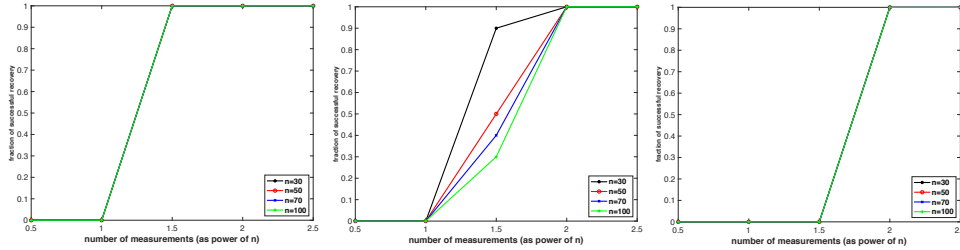


Figure 2: Empirical success rates of recovery of the Riemannian subgradient descent with  $R = 5n \log n$  runs, averaged over 10 instances. Left to right: orthogonal dictionaries with  $\theta = 0.1, 0.3, 0.5$ .

## H FASTER ALTERNATIVE ALGORITHM FOR LARGE-SCALE INSTANCES

The Riemannian subgradient descent is cheap per iteration but slow in overall convergence, similar to many other first-order methods. We also test a faster quasi-Newton type method, GRANSO,<sup>6</sup> that employs BFGS for solving constrained nonsmooth problems based on sequential quadratic optimization (Curtis et al., 2017). For a large dictionary of dimension  $n = 400$  and sample complexity  $m = 10n^2$  (i.e.,  $1.6 \times 10^6$ ), GRANSO successfully identifies a basis after 1500 iterations with CPU time 4 hours on a two-socket Intel Xeon E5-2640v4 processor (10-core Broadwell, 2.40 GHz)—this is approximately  $10\times$  faster than the Riemannian subgradient descent method, showing the potential of quasi-Newton type methods for solving large-scale problems.

## I EXPERIMENT WITH IMAGES

To experiment with images, we follow a typical setup for dictionary learning as used in image processing (Mairal et al., 2014). We focus on testing if complete (i.e., square and invertible) dictionaries are reasonable sparsification bases for real images, instead on any particular image processing or vision tasks.

<sup>6</sup>Available online: <http://www.timmitchell.com/software/GRANSO/>.

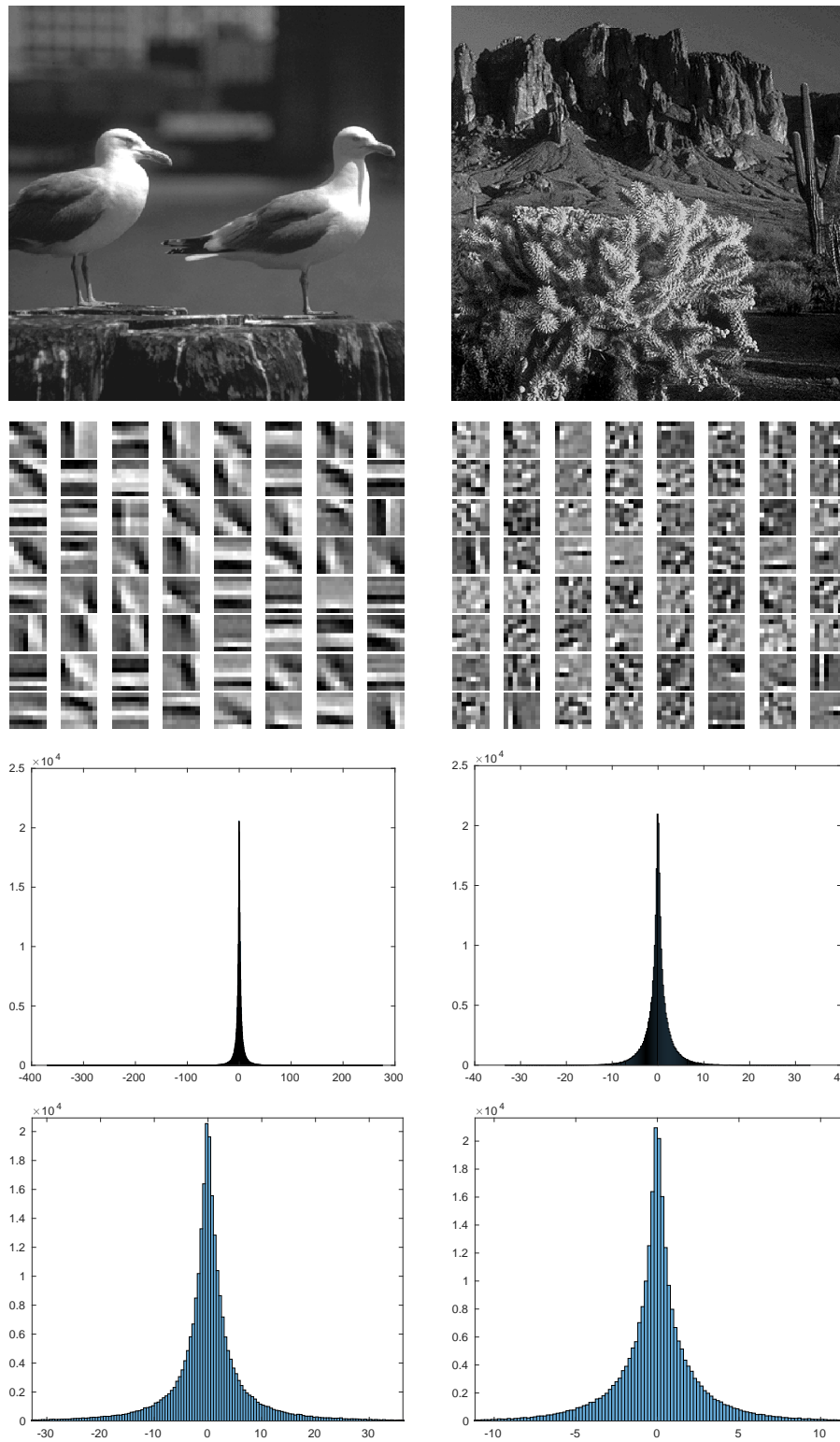


Figure 3: Results on two images. First row: the images; Second row: learned dictionaries; Third row: histograms of the representation coefficients; Fourth row: zoomed-in versions of the histograms around zero.

**Setup** Two natural images are picked for this experiment, as shown in the first row of Fig. 3, each of resolution  $512 \times 512$ . Each image is divided into  $8 \times 8$  non-overlapping blocks, resulting in  $64 \times 64 = 4096$  blocks. The blocks are then vectorized, and stacked columnwise into a data matrix  $\mathbf{Y} \in \mathbb{R}^{64 \times 4096}$ . We precondition the data to obtain

$$\bar{\mathbf{Y}} = (\mathbf{Y}\mathbf{Y}^\top)^{-1/2} \mathbf{Y}, \quad (\text{I.1})$$

so that nonvanishing singular values of  $\bar{\mathbf{Y}}$  are identically one. We then solve formulation (I.1) round  $(5n \log n)$  times with  $n = 64$  using the BFGS solver based on GRANSO, obtaining round  $(5n \log n)$  vectors. Negative equivalent copies are pruned and vectors with large correlations with other remaining vectors are sequentially removed until only 64 vectors are left. This forms the final complete dictionary.

**Results** The learned complete dictionaries for the two test images are displayed in the second row of Fig. 3. Visually, the dictionaries seem reasonably adaptive to the image contents: for the left image with prevalent sharp edges, the learned dictionary consists of almost exclusively oriented sharp corners and edges, while for the right image with blurred textures and occasional sharp features, the learned dictionary does seem to be composed of the two kinds of elements. Let the learned dictionary be  $\mathbf{A}$ . We estimate the representation coefficients as  $\mathbf{A}^{-1}\bar{\mathbf{Y}}$ . The third row of Fig. 3 contains the histograms of the coefficients. For both images, the coefficients are sharply concentrated around zero (see also the fourth row for zoomed versions of the portions around zero), and the distribution resembles a typical zero-centered Laplace distribution—which is a good indication of sparsity. Quantitatively, we calculate the mean sparsity level of the coefficient vectors (i.e., columns of  $\mathbf{A}^{-1}\bar{\mathbf{Y}}$ ) by the metric  $\|\cdot\|_1 / \|\cdot\|_2$ : for a vector  $\mathbf{v} \in \mathbb{R}^n$ ,  $\|\mathbf{v}\|_1 / \|\mathbf{v}\|_2$  ranges from 1 (when  $\mathbf{v}$  is one-sparse) to  $\sqrt{n}$  (when  $\mathbf{v}$  is fully dense with elements of equal magnitudes), which serves as a good measure of sparsity level for  $\mathbf{v}$ . For our two images, the sparsity levels by the norm-ratio metric are 5.9135 and 6.4339, respectively, while the fully dense extreme would have a value  $\sqrt{64} = 8$ , suggesting the complete dictionaries we learned are reasonable sparsification bases for the two natural images, respectively.