

# Finding a sparse vector in a subspace: Linear sparsity using alternating directions

Qing Qu, Ju Sun, and John Wright  
{qq2105, js4038, jw2966}@columbia.edu

December 15, 2014

## Abstract

We consider the problem of recovering the sparsest vector in a subspace  $\mathcal{S} \subseteq \mathbb{R}^p$  with  $\dim(\mathcal{S}) = n < p$ . This problem can be considered a homogeneous variant of the sparse recovery problem, and finds applications in sparse dictionary learning, sparse PCA, and other problems in signal processing and machine learning. Simple convex heuristics for this problem provably break down when the fraction of nonzero entries in the target sparse vector substantially exceeds  $1/\sqrt{n}$ . In contrast, we exhibit a relatively simple nonconvex approach based on alternating directions, which provably succeeds even when the fraction of nonzero entries is  $\Omega(1)$ . To our knowledge, this is the first practical algorithm to achieve this linear scaling. This result assumes a planted sparse model, in which the target sparse vector is embedded in an otherwise random subspace. Empirically, our proposed algorithm also succeeds in more challenging data models arising, e.g., from sparse dictionary learning.

## 1 Introduction

Suppose we are given a linear subspace  $\mathcal{S}$  of a high-dimensional space  $\mathbb{R}^p$ , which contains a sparse vector  $\mathbf{x}_0 \neq \mathbf{0}$ . Given arbitrary basis of  $\mathcal{S}$ , can we efficiently recover  $\mathbf{x}_0$ ? Equivalently, provided a matrix  $\mathbf{A} \in \mathbb{R}^{(p-n) \times p}$ , can we efficiently find a nonzero sparse vector  $\mathbf{x}$  such that  $\mathbf{A}\mathbf{x} = \mathbf{0}$ ? In the language of sparse approximation, can we solve

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \mathbf{A}\mathbf{x} = \mathbf{0}, \mathbf{x} \neq \mathbf{0} \quad ? \quad (1.1)$$

Variants of this problem have been studied in the context of applications to numerical linear algebra [CP86], system control and optimizations [ZF13], nonrigid structure from motion [DLH12], spectral estimation and Prony's problem [BM05], sparse PCA [ZHT06], blind source separation [ZP01], dictionary learning [SWW12], graphical model learning [AHJK13], and sparse coding on manifolds [HXV13].

However, in contrast to the standard sparse regression problem ( $\mathbf{A}\mathbf{x} = \mathbf{b}$ ,  $\mathbf{b} \neq \mathbf{0}$ ), for which convex relaxations perform nearly optimally for broad classes of designs  $\mathbf{A}$  [CT05, Don06], the computational properties of problem (1.1) are not nearly as well understood. It has been known for several decades that the basic formulation

$$\min_{\mathbf{x}} \|\mathbf{x}\|_0, \quad \text{s.t.} \quad \mathbf{x} \in \mathcal{S} \setminus \{\mathbf{0}\}, \quad (1.2)$$

is NP-hard [CP86]. However, it is only recently that efficient computational surrogates with nontrivial recovery guarantees have been discovered for some practical cases of interest. In the context of sparse dictionary learning, Spielman et al. [SWW12] introduced a relaxation which replaces the nonconvex problem (1.2) with a sequence of linear programs:

$$\ell^1/\ell^\infty \text{ Relaxation:} \quad \min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad x_i = 1, \mathbf{x} \in \mathcal{S}, 1 \leq i \leq p, \quad (1.3)$$

**Table 1:** Comparison of existing methods for recovering a planted sparse vector in a subspace

Method	Recovery Condition	Total Complexity
$\ell^1/\ell^\infty$ Relaxation[HD13]	$\theta \in O(1/\sqrt{n})$	$O(np^3)$
SDP Relaxation	$\theta \in O(1/\sqrt{n})$	$O(p^{3.5})$
SOS Relaxation [BKS13]	$p \geq \Omega(n^2), \theta \in O(1)$	high order $poly(p)$
This work	$p \geq \Omega(n^4 \log n), \theta \in O(1)$	$O(n^5 p^2 \log n)$

and proved that when  $\mathcal{S}$  is generated as a span of  $n$  random sparse vectors, with high probability the relaxation recovers these vectors, provided the probability of an entry being nonzero is at most  $\theta \in O(1/\sqrt{n})$ . In a *planted sparse model*, in which  $\mathcal{S}$  consists of a single sparse vector  $\mathbf{x}_0$  embedded in a “generic” subspace, Hand and Demanant proved that (1.3) also correctly recovers  $\mathbf{x}_0$ , provided the fraction of nonzeros in  $\mathbf{x}_0$  scales as  $\theta \in O(1/\sqrt{n})$  [HD13]. Unfortunately, the results of [SWW12, HD13] are essentially sharp: *when  $\theta$  substantially exceeds  $1/\sqrt{n}$ , in both models the relaxation (1.3) provably breaks down*. Moreover, the most natural semidefinite programming (SDP) relaxation of (1.1),

$$\min_{\mathbf{X}} \|\mathbf{X}\|_1, \quad \text{s.t.} \quad \langle \mathbf{A}^\top \mathbf{A}, \mathbf{X} \rangle = 0, \text{ trace}[\mathbf{X}] = 1, \mathbf{X} \succeq \mathbf{0}. \quad (1.4)$$

also breaks down at exactly the same threshold of  $\theta \sim 1/\sqrt{n}$ .<sup>1</sup>

One might naturally conjecture that this  $1/\sqrt{n}$  threshold is simply an intrinsic price we must pay for having an efficient algorithm, even in these random models. Some evidence towards this conjecture might be borrowed from the superficial similarity of (1.2)-(1.4) and *sparse PCA* [ZHT06]. In sparse PCA, there is a substantial gap between what can be achieved with efficient algorithms and the information theoretic optimum [BR13]. Is this also the case for recovering a sparse vector in a subspace? *Is  $\theta \in O(1/\sqrt{n})$  simply the best we can do with efficient, guaranteed algorithms?*

Remarkably, this is not the case. Recently, Barak et al. introduced a new rounding technique for sum-of-squares relaxations, and showed that the sparse vector  $\mathbf{x}_0$  in the planted sparse model can be recovered when  $p \geq \Omega(n^2)$  and  $\theta = \Omega(1)$  [BKS13]. It is perhaps surprising that this is possible at all with a polynomial time algorithm. Unfortunately, the runtime of this approach is a high-degree polynomial in  $p$ , and so for machine learning problems in which  $p$  is either a feature dimension or sample size, this algorithm is of theoretical interest only. However, it raises an interesting algorithmic question: *Is there a practical algorithm that provably recovers a sparse vector with  $\theta \gg 1/\sqrt{n}$  portion of nonzeros from a generic subspace  $\mathcal{S}$ ?*

In this paper, we address this problem, under the following hypotheses: we assume the planted sparse model, in which a target sparse vector  $\mathbf{x}_0$  is embedded in an otherwise random  $n$ -dimensional subspace of  $\mathbb{R}^p$ . We allow  $\mathbf{x}_0$  to have up to  $\theta_0 p$  nonzero entries, where  $\theta_0 \in (0, 1)$  is a constant. We provide a relatively simple algorithm which, with very high probability, exactly recovers  $\mathbf{x}_0$ , provided that  $p \geq \Omega(n^4 \log n)$ , where the comparison with existing results is shown in Table 1.

Our algorithm is based on alternating directions, with two special twists. First, we introduce a special data driven initialization, which seems to be important for achieving  $\theta = \Omega(1)$ . Second, our theoretical results require a second, linear programming based rounding phase, which is similar to [SWW12]. Our core algorithm has very simple iterations, of linear complexity in the size of the data, and hence should be scalable to moderate-to-large scale problems.

In addition to enjoying theoretical guarantees in a regime ( $\theta = \Omega(1)$ ) that is out of the reach of previous practical algorithms, it performs well in simulations – succeeding with  $p \geq \Omega(n \log n)$ . It also performs well empirically on more challenging data models, such as the dictionary learning model, in which the subspace of interest contains not one, but  $n$  target sparse vectors. Breaking the  $O(1/\sqrt{n})$  sparsity barrier with a practical algorithm is an important open problem in the nascent literature on algorithmic guarantees for dictionary learning [AGM13, ABGM14, AAN13, AAJ<sup>+</sup>13]. We are optimistic that the techniques introduced here will be applicable in this direction.<sup>2</sup>

<sup>1</sup>This breakdown behavior is again in sharp contrast to the standard sparse approximation problem (with  $\mathbf{b} \neq \mathbf{0}$ ), in which it is possible to handle very large fractions of nonzeros (say,  $\theta = \Omega(1/\log n)$ ), or even  $\theta = \Omega(1)$  using a very simple  $\ell^1$  relaxation [CT05, Don06]

<sup>2</sup>In work currently in preparation [SQW14], we show that in the dictionary learning problem, efficient algorithms based on nonconvex

## 2 Problem Formulation and Global Optimality

We study the problem of recovering a sparse vector  $\mathbf{x}_0 \neq \mathbf{0}$  (up to scale), which is an element of a known subspace  $\mathcal{S} \subset \mathbb{R}^p$  of dimension  $n$ , provided an arbitrary orthonormal basis  $\mathbf{Y} \in \mathbb{R}^{p \times n}$  for  $\mathcal{S}$ . Our starting point is the nonconvex formulation (1.2). Both the objective and constraint are nonconvex, and hence not easy to optimize over. We relax (1.2) by replacing the  $\ell^0$  norm with the  $\ell^1$  norm. For the constraint  $\mathbf{x} \neq \mathbf{0}$ , which is necessary to avoid a trivial solution, we force  $\mathbf{x}$  to live on the unit sphere  $\|\mathbf{x}\|_2 = 1$ , giving

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \mathbf{x} \in \mathcal{S}, \|\mathbf{x}\|_2 = 1. \quad (2.1)$$

This formulation is still nonconvex, and so we should not expect to obtain an efficient algorithm that can solve it globally for general inputs  $\mathcal{S}$ . Nevertheless, the geometry of the sphere is benign enough that for well-structured inputs it actually *will* be possible to give algorithms that find the global optimum.

The formulation (2.1) can be contrasted with (1.3), in which effectively we optimize the  $\ell^1$  norm subject to the constraint  $\|\mathbf{x}\|_\infty = 1$ . Because  $\|\cdot\|_\infty$  is polyhedral, that formulation immediately yields a sequence of linear programs. This is very convenient for computation and analysis, but suffers from the aforementioned breakdown behavior around  $\|\mathbf{x}_0\|_0 \sim p/\sqrt{n}$ . In contrast, the sphere  $\|\mathbf{x}\|_2 = 1$  is a more complicated geometric constraint, but will allow much larger numbers of nonzeros in  $\mathbf{x}_0$ . Indeed, if we consider the global optimizer of a reformulation of (2.1):

$$\min_{\mathbf{q} \in \mathbb{R}^n} \|\mathbf{Y}\mathbf{q}\|_1, \quad \text{s.t.} \quad \|\mathbf{q}\|_2 = 1, \quad (2.2)$$

where  $\mathbf{Y}$  is any orthonormal basis for  $\mathcal{S}$ , we have strong recovery guarantees as follows:

**Theorem 2.1** ( $\ell^1/\ell^2$  recovery, planted sparse model). *There exists a constant  $\theta_0 > 0$  such that if the subspace  $\mathcal{S}$  follows the planted sparse model*

$$\mathcal{S} = \text{span}(\mathbf{x}_0, \mathbf{g}_1, \dots, \mathbf{g}_{n-1}) \subset \mathbb{R}^p, \quad (2.3)$$

where  $\mathbf{g}_i \sim_{i.i.d.} \mathcal{N}(\mathbf{0}, \frac{1}{p}\mathbf{I})$ , and  $\mathbf{x}_0 \sim_{i.i.d.} \frac{1}{\sqrt{\theta p}} \text{Ber}(\theta)$  are all mutually independent and  $1/\sqrt{n} < \theta < \theta_0$ , then optimum  $\mathbf{q}^*$  to (2.2), for any orthonormal basis  $\mathbf{Y}$  of  $\mathcal{S}$ , produces  $\mathbf{Y}\mathbf{q}^* = \xi\mathbf{x}_0$  for some  $\xi \neq 0$  with high probability, provided  $p \geq \Omega(n \log n)$ .<sup>3</sup>

Hence, if we could find the global optimizer of (2.2), we would be able to recover  $\mathbf{x}_0$  whose number of nonzero entries is quite large – even linear in the dimension  $p$  ( $\theta = \Omega(1)$ ). On the other hand, it is not obvious that this should be possible: (2.2) is nonconvex. In the next section, we will describe a simple heuristic algorithm for (a near approximation of) the  $\ell^1/\ell^2$  problem (2.2), which guarantees to find a stationary point. More surprisingly, we will then prove that for a class of random problem instances, this algorithm, plus an auxiliary rounding technique, actually recovers the global optimum – the target sparse vector  $\mathbf{x}_0$ . The proof requires a detailed probabilistic analysis, which is sketched in Section 4.2.

Before continuing, it is worth noting that the formulation (2.1) is in no way novel – see, e.g., the work of [ZP01] in blind source separation for precedent. However, our algorithms and subsequent analysis are novel.

## 3 Algorithm based on Alternating Direction Method (ADM)

To develop an algorithm for solving (2.2), it is useful to consider a slight relaxation of (2.2), in which we introduce an auxiliary variable  $\mathbf{x} \approx \mathbf{Y}\mathbf{q}$ :

$$\min_{\mathbf{q}, \mathbf{x}} \frac{1}{2} \|\mathbf{Y}\mathbf{q} - \mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad \text{s.t.} \quad \|\mathbf{q}\|_2 = 1, \quad (3.1)$$

optimization also produce global solutions, even when  $\theta = \Omega(1)$ .

<sup>3</sup>Note that this version is much stronger and more practical than that appearing in the conference version [QSW14].

Here,  $\lambda > 0$  is a penalty parameter. It is not difficult to see that this problem is equivalent to minimizing the *Huber M-estimator* over  $\mathbf{Y}\mathbf{q}$ . This relaxation makes it possible to apply the alternating direction method to this problem. This method starts from some initial point  $\mathbf{q}^{(0)}$ , alternates between optimizing with respect to  $\mathbf{x}$  and optimizing with respect to  $\mathbf{q}$ :

$$\mathbf{x}^{(k+1)} = \arg \min_{\mathbf{x}} \frac{1}{2} \left\| \mathbf{Y}\mathbf{q}^{(k)} - \mathbf{x} \right\|_2^2 + \lambda \|\mathbf{x}\|_1, \quad (3.2)$$

$$\mathbf{q}^{(k+1)} = \arg \min_{\mathbf{q}} \frac{1}{2} \left\| \mathbf{Y}\mathbf{q} - \mathbf{x}^{(k+1)} \right\|_2^2 \text{ s.t. } \|\mathbf{q}\|_2 = 1. \quad (3.3)$$

Both (3.2) and (3.3) have simple closed form solutions:

$$\mathbf{x}^{(k+1)} = S_\lambda[\mathbf{Y}\mathbf{q}^{(k)}], \quad \mathbf{q}^{(k+1)} = \frac{\mathbf{Y}^\top \mathbf{x}^{(k+1)}}{\left\| \mathbf{Y}^\top \mathbf{x}^{(k+1)} \right\|_2}, \quad (3.4)$$

where  $S_\lambda[x] = \text{sign}(x) \max\{|x| - \lambda, 0\}$  is the soft-thresholding operator. The proposed ADM algorithm is summarized in Algorithm 1.

---

**Algorithm 1** Nonconvex ADM for solving (3.1)

---

**Input:** A matrix  $\mathbf{Y} \in \mathbb{R}^{p \times n}$  with  $\mathbf{Y}^\top \mathbf{Y} = \mathbf{I}$ , initialization  $\mathbf{q}^{(0)}$ , threshold parameter  $\lambda > 0$ .

**Output:** The recovered sparse vector  $\hat{\mathbf{x}}_0 = \mathbf{Y}\mathbf{q}^{(k)}$

1: **for**  $k = 0, \dots, O(n^4 \log n)$  **do**

2:  $\mathbf{x}^{(k+1)} = S_\lambda[\mathbf{Y}\mathbf{q}^{(k)}],$

3:  $\mathbf{q}^{(k+1)} = \frac{\mathbf{Y}^\top \mathbf{x}^{(k+1)}}{\left\| \mathbf{Y}^\top \mathbf{x}^{(k+1)} \right\|_2},$

4: **end for**

---

The algorithm is simple to state and easy to implement. However, if our goal is to recover the *sparsest* vector  $\mathbf{x}_0$ , some additional tricks are needed.

**Initialization.** Because the problem (2.2) is nonconvex, an arbitrary or random initialization is unlikely to produce a global minimizer.<sup>4</sup> Therefore, good initializations are critical for the proposed ADM algorithm to succeed. For this purpose, we suggest to use every normalized row of  $\mathbf{Y}$  as initializations for  $\mathbf{q}$ , and solve a sequence of  $p$  nonconvex programs (2.2) by the ADM algorithm.

To get an intuition of why our initialization works, recall the planted sparse model:  $\mathcal{S} = \text{span}(\mathbf{x}_0, \mathbf{g}_1, \dots, \mathbf{g}_{n-1})$ . Write  $\mathbf{Z} = [\mathbf{x}_0 \mid \mathbf{g}_1 \mid \dots \mid \mathbf{g}_{n-1}] \in \mathbb{R}^{p \times n}$ . Suppose we take a row  $\mathbf{z}^i$  of  $\mathbf{Z}$ , in which  $\mathbf{x}_0(i)$  is nonzero, then  $\mathbf{x}_0(i) = \Theta(1/\sqrt{\theta p})$ . Meanwhile, the entries of  $\mathbf{g}_1(i), \dots, \mathbf{g}_{n-1}(i)$  are all  $\mathcal{N}(0, 1/p)$ , and so have size about  $1/\sqrt{p}$ . Hence, when  $\theta$  is not too large,  $\mathbf{x}_0(i)$  will be somewhat bigger than most of the other entries in  $\mathbf{z}_i$ . Put another way,  $\mathbf{z}_i$  is *biased towards the first standard basis vector*  $\mathbf{e}_1$ .

Now, under our probabilistic assumptions,  $\mathbf{Z}$  is very well conditioned:  $\mathbf{Z}^\top \mathbf{Z} \approx \mathbf{I}$ .<sup>5</sup> Using, e.g., Gram-Schmidt, we can find a basis  $\bar{\mathbf{Y}}$  for  $\mathcal{S}$  of the form

$$\bar{\mathbf{Y}} = \mathbf{Z}\mathbf{R}, \quad (3.5)$$

where  $\mathbf{R}$  is upper triangular, and  $\mathbf{R}$  is itself well-conditioned:  $\mathbf{R} \approx \mathbf{I}$ . Since the  $i$ -th row of  $\mathbf{Z}$  is biased in the direction of  $\mathbf{e}_1$  and  $\mathbf{R}$  is well-conditioned, the  $i$ -th row  $\bar{\mathbf{y}}^i$  is also biased in the direction of  $\mathbf{e}_1$ . Moreover, we know that the global optimizer  $\mathbf{q}_*$  should satisfy  $\bar{\mathbf{Y}}\mathbf{q}_* = \mathbf{x}_0$ . Since  $\mathbf{Z}\mathbf{e}_1 = \mathbf{x}_0$ , we have  $\mathbf{q}_* = \mathbf{R}^{-1}\mathbf{e}_1 \approx \mathbf{e}_1$ . Here, the approximation comes from  $\mathbf{R} \approx \mathbf{I}$ . Hence, for this particular choice of  $\mathbf{Y}$ , described in (3.5), *the  $i$ -th row is biased in the direction of the global optimizer*.

<sup>4</sup>More precisely, in our models, random initialization *does* work, but only when the subspace dimension  $n$  is extremely low compared to the ambient dimension  $p$ .

<sup>5</sup>This is the common heuristic that “tall random matrices are well conditioned” [Ver10].

What if we are handed some other basis  $\mathbf{Y} = \bar{\mathbf{Y}}\mathbf{U}$ , where  $\mathbf{U}$  is an orthogonal matrix? Suppose  $\mathbf{q}_*$  is a global optimizer to (2.2) with input matrix  $\bar{\mathbf{Y}}$ , then it is easy to check that, with input matrix  $\mathbf{Y}$ ,  $\mathbf{U}^\top \mathbf{q}_*$  is also a global optimizer to (2.2), which implies that our initialization is *invariant* to any rotation of the basis. Hence, *even if we are handed an arbitrary basis for  $\mathcal{S}$ , the  $i$ -th row is still biased in the direction of the global optimizer.*

**Rounding.** Let  $\bar{\mathbf{q}}$  denote the output of Algorithm 1. We will prove that with our particular initialization and an appropriate choice of  $\lambda$ , the solution of our ADM algorithm falls within a certain radius of the globally optimal solution  $\mathbf{q}_*$  to (2.2). To recover  $\mathbf{q}_*$ , or equivalently to recover the sparse vector  $\mathbf{x}_0 = \xi \mathbf{Y} \mathbf{q}_*$  for some  $\xi \neq 0$ , we solve the linear program

$$\min_{\mathbf{q}} \|\mathbf{Y}\mathbf{q}\|_1 \quad \text{s.t.} \quad \langle \mathbf{r}, \mathbf{q} \rangle = 1 \quad (3.6)$$

with  $\mathbf{r} = \bar{\mathbf{q}}$ . We will prove that if  $\bar{\mathbf{q}}$  is close enough to  $\mathbf{q}^*$ , then (3.6) exactly recovers  $\mathbf{q}^*$ , and hence  $\mathbf{x}_0$ .

## 4 Analysis

### 4.1 Main Results

In this section, we describe our main theoretical result, which shows that with high probability, the algorithm described in the previous section succeeds.

**Theorem 4.1.** *Suppose that  $\mathcal{S}$  satisfies the planted sparse model, and let the columns of  $\mathbf{Y}$  be an arbitrary orthonormal basis for the subspace  $\mathcal{S}$ . Let  $\mathbf{y}^1, \dots, \mathbf{y}^p \in \mathbb{R}^n$  denote the (transposes of) the rows of  $\mathbf{Y}$ . Apply Algorithm 1 with  $\lambda = 1/\sqrt{p}$ , using initializations  $\mathbf{q}^{(0)} = \mathbf{y}^1, \dots, \mathbf{y}^p$ , to produce outputs  $\bar{\mathbf{q}}_1, \dots, \bar{\mathbf{q}}_p$ . Solve the linear program (3.6) with  $\mathbf{r} = \bar{\mathbf{q}}_1, \dots, \bar{\mathbf{q}}_p$ , to produce  $\hat{\mathbf{q}}_1, \dots, \hat{\mathbf{q}}_p$ . Set  $i^* \in \arg \min_i \|\mathbf{Y}\hat{\mathbf{q}}_i\|_0$ . Then*

$$\mathbf{Y}\hat{\mathbf{q}}_{i^*} = \gamma \mathbf{x}_0, \quad (4.1)$$

for some  $\gamma \neq 0$  with overwhelming probability, provided

$$\exp(n/2)/2 \geq p \geq Cn^4 \log n, \quad \text{and} \quad \frac{1}{\sqrt{n}} \leq \theta \leq \theta_0. \quad (4.2)$$

Here,  $C$  and  $\theta_0 > 0$  are universal constants.

We can see that the result in Theorem 4.1 is suboptimal compared to the global optimality result in Theorem 2.1 and Barak et al.'s result [BKS13] in sampling complexity. For successful recovery, we require  $p \geq \Omega(n^4 \log n)$ , while the global optimality and Barak et al. demand  $p \geq \Omega(n \log n)$  and  $p \geq \Omega(n^2)$ , respectively. Aside from possible deficiencies in our current analysis, compared to Barak et al., we believe this is still the first practical and efficient method which is guaranteed to achieve  $\theta \sim O(1)$  rate. The lower bound on  $\theta$  in Theorem 4.1 is mostly for convenience in the proof; in fact, the LP rounding stage of our algorithm already succeeds with high probability when  $\theta \in O(1/\sqrt{n})$ .

### 4.2 A Sketch of Analysis

The proof of our main result requires rather detailed technical analysis of the iteration-by-iteration properties of Algorithm 1. In this section, as illustrated in Fig. 1, we briefly sketch the main ideas. Detailed proofs are deferred to the appendices.

As noted in Section 3, the ADM algorithm is invariant to change of basis. So, we can assume without loss of generality that we are working with the particular basis  $\bar{\mathbf{Y}} = \mathbf{Z}\mathbf{R}$  defined in that section. In order to further streamline the presentation, we are going to sketch the proof under the assumption that

$$\mathbf{Y} = [\mathbf{x}_0 \mid \mathbf{g}_1 \mid \dots \mid \mathbf{g}_{n-1}], \quad (4.3)$$

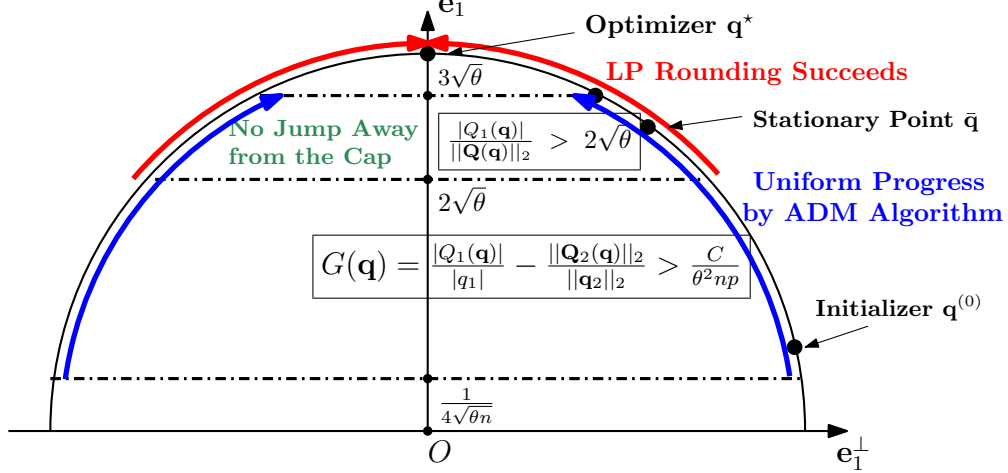


Figure 1: An illustration of the proof sketch for our ADM algorithm.

rather than the orthogonalized version  $\bar{\mathbf{Y}}$ . When  $p$  is large  $\mathbf{Y}$  is already nearly orthogonal, and hence  $\mathbf{Y}$  is very close to  $\bar{\mathbf{Y}}$ . In fact, in our proof, we simply carry through the argument for  $\mathbf{Y}$ , and then note that  $\mathbf{Y}$  and  $\bar{\mathbf{Y}}$  are close enough that all steps of the proof still hold with  $\mathbf{Y}$  replaced by  $\bar{\mathbf{Y}}$ . With that noted, let  $\mathbf{y}^1, \dots, \mathbf{y}^p \in \mathbb{R}^n$  denote the transposes of the rows of  $\mathbf{Y}$ , and note that these are independent random vectors. From (3.4), we can see one step of the ADM algorithm takes the form:

$$\mathbf{q}^{(k+1)} = \frac{\frac{1}{p} \sum_{i=1}^p \mathbf{y}^i S_\lambda \left[ (\mathbf{q}^{(k)})^\top \mathbf{y}^i \right]}{\left\| \frac{1}{p} \sum_{i=1}^p \mathbf{y}^i S_\lambda \left[ (\mathbf{q}^{(k)})^\top \mathbf{y}^i \right] \right\|_2}. \quad (4.4)$$

This is a very favorable form for analysis: if  $\mathbf{q}^{(k)}$  is viewed as fixed, the term in the numerator is a sum of  $p$  independent random vectors. To this end, we define a vector valued random process  $\mathbf{Q}(\mathbf{q})$  on  $\mathbf{q} \in \mathbb{S}^{n-1}$ , via

$$\mathbf{Q}(\mathbf{q}) = \frac{1}{p} \sum_{i=1}^p \mathbf{y}^i S_\lambda [\mathbf{q}^\top \mathbf{y}^i]. \quad (4.5)$$

We study the behavior of the iteration (4.4) through the random process  $\mathbf{Q}(\mathbf{q})$ . We wish to show that with overwhelming probability in our choice of  $\mathbf{Y}$ ,  $\mathbf{q}^{(k)}$  converges to some small neighborhood of  $\pm \mathbf{e}_1$ , so that the ADM algorithm plus the LP rounding (described in Section 3) successfully retrieves the sparse vector  $\mathbf{x}_0 = \mathbf{Y} \mathbf{e}_1$ . Thus, we hope that in general,  $\mathbf{Q}(\mathbf{q})$  is more concentrated on the first coordinate than  $\mathbf{q}$ . Let us partition the vector  $\mathbf{q}$  as  $\mathbf{q} = \begin{bmatrix} q_1 \\ \mathbf{q}_2 \end{bmatrix}$ , with  $q_1 \in \mathbb{R}$  and  $\mathbf{q}_2 \in \mathbb{R}^{n-1}$ , and correspondingly partition

$\mathbf{Q}(\mathbf{q}) = \begin{bmatrix} Q_1(\mathbf{q}) \\ \mathbf{Q}_2(\mathbf{q}) \end{bmatrix}$ , where

$$Q_1(\mathbf{q}) = \frac{1}{p} \sum_{i=1}^p x_{0i} S_\lambda [\mathbf{q}^\top \mathbf{y}^i] \quad \text{and} \quad \mathbf{Q}_2(\mathbf{q}) = \frac{1}{p} \sum_{i=1}^p \mathbf{g}_i S_\lambda [\mathbf{q}^\top \mathbf{y}^i]. \quad (4.6)$$

The inner product of  $\mathbf{Q}(\mathbf{q}) / \|\mathbf{Q}(\mathbf{q})\|_2$  and  $\mathbf{e}_1$  is strictly larger than the inner product of  $\mathbf{q}$  and  $\mathbf{e}_1$  if and only if

$$\frac{|Q_1(\mathbf{q})|}{|q_1|} > \frac{\|\mathbf{Q}_2(\mathbf{q})\|_2}{\|\mathbf{q}_2\|_2}. \quad (4.7)$$

In the appendix, we show that with overwhelming probability, this inequality holds uniformly over a significant portion of the sphere, so the algorithm moves in the correct direction. To complete the proof of Theorem 4.1, we combine the following observations, provided  $\exp(n/2)/2 \geq p \geq \Omega(n^4 \log n)$ :

1. **Good initializers.** With overwhelming probability, at least one of the initializers  $\mathbf{q}^{(0)}$  satisfies  $|q_1^{(0)}| > \frac{1}{4\sqrt{\theta n}}$ .
2. **Uniform progress away from the equator.** With overwhelming probability, for every  $\mathbf{q} \in \mathbb{S}^{n-1}$  such that  $\frac{1}{4\sqrt{\theta n}} \leq |q_1| \leq 3\sqrt{\theta}$ , the bound

$$\frac{|Q_1(\mathbf{q})|}{|q_1|} - \frac{\|\mathbf{Q}_2(\mathbf{q})\|_2}{\|\mathbf{q}\|_2} > \frac{C}{\theta^2 np} \quad (4.8)$$

holds, for some numerical constant  $C > 0$ . This implies that if at any iteration  $k$  of the algorithm,  $|q_1^{(k)}| > \frac{1}{4\sqrt{\theta n}}$ , the algorithm will eventually obtain a point  $\mathbf{q}^{(k')}$ ,  $k' > k$ , for which  $|q_1^{(k')}| > 3\sqrt{\theta}$ , if sufficiently many iterations are allowed.

3. **No jumps away from the caps.** With overwhelming probability,

$$\frac{|Q_1(\mathbf{q})|}{\sqrt{|Q_1(\mathbf{q})|^2 + \|\mathbf{Q}_2(\mathbf{q})\|_2^2}} \geq 2\sqrt{\theta} \quad (4.9)$$

for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  with  $|q_1| > 3\sqrt{\theta}$ .

4. **Location of stationary points.** The above steps imply that, with overwhelming probability, Algorithm 1 fed with the proposed initialization scheme produces at least one stopping point  $\bar{\mathbf{q}} \in \mathbb{S}^{n-1}$  satisfying  $|\bar{q}_1| \geq 2\sqrt{\theta}$ .
5. **Rounding succeeds when  $|r_1| > 2\sqrt{\theta}$ .** With overwhelming probability, the linear programming based rounding (3.6) will produce  $\pm \mathbf{x}_0$ , up to scale, whenever it is provided with an input  $\mathbf{r}$  whose first coordinate has magnitude at least  $2\sqrt{\theta}$ .

Taken together, these claims imply that from at least one of the initializers  $\mathbf{q}^{(0)}$ , the ADM algorithm will produce an output  $\bar{\mathbf{q}}$  which is accurate enough for LP rounding to exactly return  $\mathbf{x}_0$ , up to scale. As  $\mathbf{x}_0$  is the sparsest nonzero vector in the subspace  $\mathcal{S}$  with overwhelming probability, it will be selected as  $\mathbf{Y}\mathbf{q}_{i^*}$ , and hence produced by the algorithm.

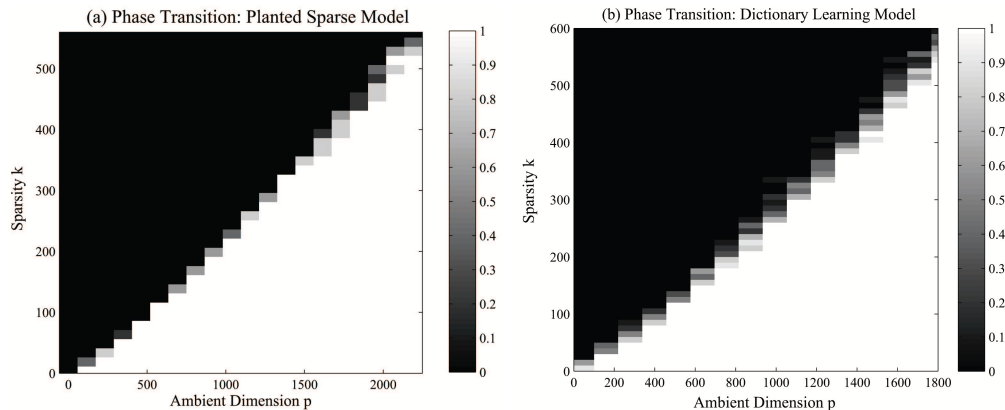
## 5 Experimental Results

In this section, we show the performance of the proposed ADM algorithm on both synthetic and real datasets. On the synthetic dataset, we show the phase transition of our algorithm on both the planted sparse vector and dictionary learning models; for the real dataset, we demonstrate how seeking sparse vectors can help discover interesting patterns.

### 5.1 Phase Transition on Synthetic Data

For the planted sparse model, for each pair of  $(k, p)$ , we generate the  $n$  dimensional subspace  $\mathcal{S} \in \mathbb{R}^p$  by a  $k$  sparse vector  $\mathbf{x}_0$  with nonzero entries equal to 1 and a random Gaussian matrix  $\mathbf{G} \in \mathbb{R}^{p \times (n-1)}$  with  $G_{ij} \sim_{i.i.d.} \mathcal{N}(0, 1/p)$ , so that one basis  $\mathbf{Y}$  of the subspace  $\mathcal{S}$  can be constructed by  $\mathbf{Y} = \text{GS}([\mathbf{x}_0, \mathbf{G}])\mathbf{U}$ , where  $\text{GS}(\cdot)$  denotes the Gram-Schmidt orthonormalization operator and  $\mathbf{U} \in \mathbb{R}^{n \times n}$  is an arbitrary orthogonal matrix. We fix the relationship between  $n$  and  $p$  as  $p = 5n \log n$ , and set the regularization parameter in (3.1) as  $\lambda = 1/\sqrt{p}$ . We use all the normalized rows of  $\mathbf{Y}$  as initializations of  $\mathbf{q}$  for the proposed ADM

algorithm, and run every program for 5000 iterations. We determine the recovery to be successful whenever  $\|\mathbf{x}_0/\|\mathbf{x}_0\|_2 - \mathbf{Y}\mathbf{q}\|_2 \leq \varepsilon$  for at least one of the  $p$  programs, where  $\varepsilon = 10^{-3}$ . For each pair of  $(k, p)$ , we repeat the simulation for five times.



**Figure 2:** Phase transition for the planted sparse model (left) and dictionary learning model (right) using the ADM algorithm, with fixed relationship between  $p$  and  $n$ :  $p = 5n \log n$ . White indicates success and black indicates failure.

Second, we consider the same dictionary learning model as in [SWW12]. Specifically, the observation is assumed to be  $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$ , where  $\mathbf{A}_0$  is a square, invertible matrix, and  $\mathbf{X}_0$  a  $n \times p$  sparse matrix. Since  $\mathbf{A}_0$  is invertible, the row space of  $\mathbf{Y}$  is the same as that of  $\mathbf{X}_0$ . For each pair of  $(k, n)$ , we generate  $\mathbf{X}_0 = [\mathbf{x}_1, \dots, \mathbf{x}_n]^\top$ , where each vector  $\mathbf{x}_i \in \mathbb{R}^p$  is  $k$ -sparse with every nonzero entry following i.i.d. Gaussian distribution, and construct the observation by  $\mathbf{Y}^\top = \mathbf{G}\mathbf{S}(\mathbf{X}_0^\top) \mathbf{U}^\top$ . We repeat the same experiment as for the planted sparse model presented above. The only difference is that here we determine the recovery to be successful as long as one sparse row of  $\mathbf{X}_0$  is recovered by one of those  $p$  programs.

Figure 2 shows the phase transition between the sparsity level  $k$  and  $p$  for both models. It seems clear for both problems our algorithm can work well into (even beyond) the linear sparsity regime whenever  $p \sim n \log n$ . Hence for the planted sparse model, to close the gap between our algorithm and practice is one future direction. Also, how to extend our analysis for dictionary learning is another interesting direction.

## 5.2 Exploratory Experiments on Faces

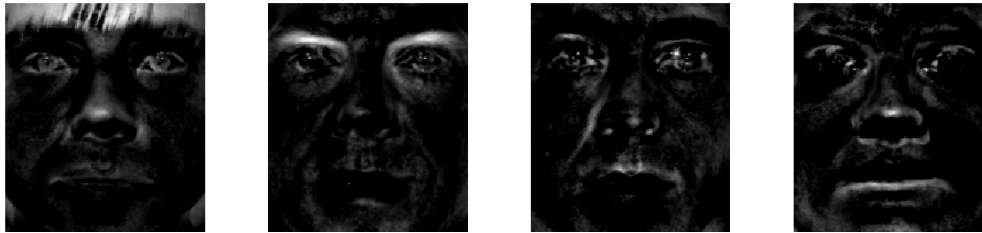
It is well known in computer vision that appearance of convex objects only subject to illumination changes leads to image collection that can be well approximated by low-dimensional space in raw-pixel space [BJ03]. We will play with face subspaces here. First, we extract face images of one person (65 images) under different illumination conditions. Then we apply *robust principal component analysis* [CLMW11] to the data and get a low dimensional subspace of dimension 10, i.e., the basis  $\mathbf{Y} \in \mathbb{R}^{32256 \times 10}$ . We apply the ADM algorithm to find the sparsest element in such a subspace, by randomly selecting 10% rows as initializations for  $\mathbf{q}$ . We judge the sparsity in a  $\ell^1/\ell^2$  sense, that is, the sparsest vector  $\hat{\mathbf{x}}_0 = \mathbf{Y}\mathbf{q}^*$  should produce the smallest  $\|\mathbf{Y}\mathbf{q}\|_1 / \|\mathbf{Y}\mathbf{q}\|_2$  among all results. Once some sparse vectors are found, we project the subspace onto orthogonal complement of the sparse vectors already found, and continue the seeking process in the projected subspace. Figure 3 shows the first four sparse vectors we get from the data. We can see they correspond well to different extreme illumination conditions.

Second, we manually select ten different persons' faces under the normal lighting condition. Again, the dimension of the subspace is 10 and  $\mathbf{Y} \in \mathbb{R}^{32256 \times 10}$ . We repeat the same experiment as stated above. Figure 4 shows four sparse vectors we get from the data. Interestingly, the sparse vectors roughly correspond to differences of face images concentrated around facial parts that different people tend to differ from each other, e.g., eye bows, forehead hair, nose, etc.





**Figure 3:** Four sparse vectors extracted by the ADM algorithm for one person in the Yale B database under different illuminations.



**Figure 4:** Four sparse vectors extracted by the ADM algorithm for 10 persons in the Yale B database under normal illuminations.

In sum, our algorithm seems to find useful sparse vectors for potential applications, like peculiarity discovery in first setting, and locating differences in second setting. Nevertheless, the main goal of this experiment is to invite readers to think about similar pattern discovery problems that might be cast as the problem of seeking sparse vectors in a subspace. The experiment also demonstrates in a concrete way the practicality of our algorithm, both in handling data sets of realistic size and in producing meaningful results even beyond the (idealized) planted sparse model that we adopt for analysis.

## 6 Discussion

The random models we assume for the subspace can be easily extended to other random models, particularly for dictionary learning. Moreover we believe the algorithm paradigm works far beyond the idealized models, as our preliminary experiments on face data have clearly shown. For the particular planted sparse model, the performance gap in terms of  $(p, n, \theta)$  between the empirical simulation and our result is likely due to analysis itself. Advanced techniques to bound the empirical process, such as decoupling [DIPG99] techniques, can be deployed in place of our crude union bound to cover all iterates. On the application side, the potential of seeking sparse/structured element in a subspace seems largely unexplored, despite the cases we mentioned at the start. We hope this work can invite more application ideas.

This paper is part of a recent surge of research into provable and practical nonconvex approaches to estimating various types of low-dimensional structures, often in large-scale settings [CLS14, JNS13, Har13, NJS13, YCS13]. The dominant approach is to start with a clever, problem-specific initialization, and then perform a local analysis of the subsequent iterates. Our forthcoming work [SQW14] on dictionary learning takes a more geometric approach, and proves global recovery via efficient algorithms, with arbitrary initialization. The approach developed there may be applicable to the planted sparse model studied here, as well as to many other interesting nonconvex problems.

## Acknowledgement

JS thanks the Wei Family Private Foundation for their generous support. We thank Cun Mu, IEOR Department of Columbia University, for helpful discussion and input regarding this work. This work was partially supported by grants ONR N00014-13-1-0492, NSF 1343282, and funding from the Moore and Sloan Foundations.

## References

- [AAJ<sup>+</sup>13] Alekh Agarwal, Animashree Anandkumar, Prateek Jain, Praneeth Netrapalli, and Rashish Tandon. Learning sparsely used overcomplete dictionaries via alternating minimization. *arXiv preprint arXiv:1310.7991*, 2013.
- [AAN13] Alekh Agarwal, Animashree Anandkumar, and Praneeth Netrapalli. Exact recovery of sparsely used overcomplete dictionaries. *arXiv preprint arXiv:1309.1952*, 2013.
- [ABGM14] Sanjeev Arora, Aditya Bhaskara, Rong Ge, and Tengyu Ma. More algorithms for provable dictionary learning. *arXiv preprint arXiv:1401.0579*, 2014.
- [AGM13] Sanjeev Arora, Rong Ge, and Ankur Moitra. New algorithms for learning incoherent and overcomplete dictionaries. *arXiv preprint arXiv:1308.6273*, 2013.
- [AHJK13] Anima Anandkumar, Daniel Hsu, Majid Janzamin, and Sham M Kakade. When are overcomplete topic models identifiable? uniqueness of tensor tucker decompositions with structured sparsity. In *Advances in Neural Information Processing Systems*, pages 1986–1994, 2013.
- [BJ03] Ronen Basri and David W Jacobs. Lambertian reflectance and linear subspaces. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 25(2):218–233, 2003.
- [BKS13] Boaz Barak, Jonathan Kelner, and David Steurer. Rounding sum-of-squares relaxations. *arXiv preprint arXiv:1312.6652*, 2013.
- [BM05] Gregory Beylkin and Lucas Monzón. On approximation of functions by exponential sums. *Applied and Computational Harmonic Analysis*, 19(1):17–48, 2005.
- [BR13] Quentin Berthet and Philippe Rigollet. Complexity theoretic lower bounds for sparse principal component detection. In *Conference on Learning Theory*, pages 1046–1066, 2013.
- [CLMW11] Emmanuel Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM*, 58(3), May 2011.
- [CLS14] Emmanuel J. Candès, Xiaodong Li, and Mahdi Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *arXiv preprint arXiv:1407.1065*, 2014.
- [CP86] Thomas F Coleman and Alex Pothén. The null space problem i. complexity. *SIAM Journal on Algebraic Discrete Methods*, 7(4):527–537, 1986.
- [CT05] Emmanuel J Candès and Terence Tao. Decoding by linear programming. *Information Theory, IEEE Transactions on*, 51(12):4203–4215, 2005.
- [DLH12] Yuchao Dai, Hongdong Li, and Mingyi He. A simple prior-free method for non-rigid structure-from-motion factorization. In *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*, pages 2018–2025. IEEE, 2012.
- [DIPG99] Victor De la Pena and Evarist Giné. *Decoupling: from dependence to independence*. Springer, 1999.

- [Don06] David L Donoho. For most large underdetermined systems of linear equations the minimal  $\ell^1$ -norm solution is also the sparsest solution. *Communications on pure and applied mathematics*, 59(6):797–829, 2006.
- [FLM77] Tadeusz Figiel, Joram Lindenstrauss, and Vitali D Milman. The dimension of almost spherical sections of convex bodies. *Acta Mathematica*, 139(1):53–94, 1977.
- [GG84] Andrej Y Garnaev and Efim D Gluskin. The widths of a euclidean ball. In *Dokl. Akad. Nauk SSSR*, volume 277, pages 1048–1052, 1984.
- [GM03] E Gluskin and V Milman. Note on the geometric-arithmetic mean inequality. In *Geometric aspects of Functional analysis*, pages 131–135. Springer, 2003.
- [Har13] Moritz Hardt. On the provable convergence of alternating minimization for matrix completion. *arXiv preprint arXiv:1312.0925*, 2013.
- [HD13] Paul Hand and Laurent Demanet. Recovering the sparsest element in a subspace. *arXiv preprint arXiv:1310.1654*, 2013.
- [HXV13] Jeffrey Ho, Yuchen Xie, and Baba Vemuri. On a nonlinear generalization of sparse coding and dictionary learning. In *Proceedings of The 30th International Conference on Machine Learning*, pages 1480–1488, 2013.
- [JNS13] Prateek Jain, Praneeth Netrapalli, and Sujay Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the 45th annual ACM symposium on Symposium on theory of computing*, pages 665–674. ACM, 2013.
- [NJS13] Praneeth Netrapalli, Prateek Jain, and Sujay Sanghavi. Phase retrieval using alternating minimization. In *Advances in Neural Information Processing Systems*, pages 2796–2804, 2013.
- [Pis99] Gilles Pisier. *The volume of convex bodies and Banach space geometry*, volume 94. Cambridge University Press, 1999.
- [QSW14] Qing Qu, Ju Sun, and John Wright. Finding a sparse vector in a subspace: Linear sparsity using alternating directions. In *Advances in Neural Information Processing Systems*, 2014.
- [SQW14] Ju Sun, Qing Qu, and John Wright. Complete dictionary learning over the sphere. *In preparation*, 2014.
- [SWW12] Daniel A Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *Proceedings of the 25th Annual Conference on Learning Theory*, 2012.
- [Ver10] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [YCS13] Xinyang Yi, Constantine Caramanis, and Sujay Sanghavi. Alternating minimization for mixed linear regression. *arXiv preprint arXiv:1310.3745*, 2013.
- [ZF13] Yun-Bin Zhao and Masao Fukushima. Rank-one solutions for homogeneous linear matrix equations over the positive semidefinite cone. *Applied Mathematics and Computation*, 219(10):5569–5583, 2013.
- [ZHT06] Hui Zou, Trevor Hastie, and Robert Tibshirani. Sparse principal component analysis. *Journal of computational and graphical statistics*, 15(2):265–286, 2006.
- [ZP01] Michael Zibulevsky and Barak A Pearlmutter. Blind source separation by sparse decomposition in a signal dictionary. *Neural computation*, 13(4):863–882, 2001.

# Appendices

**Notes on notations.** For a matrix  $\mathbf{X}$ ,  $\mathbf{x}_i$  denotes its  $i$ -th column, and  $\mathbf{x}^j$  denotes its  $j$ -th row, all in column vector form. So  $(\mathbf{x}^j)^\top$  is the  $j$ -th row in row vector form. We will use the compact notation  $[k] \doteq \{1, \dots, k\}$  for any positive integer  $k$ . We will use  $C$  or  $c$ , and their indexed versions to denote constants. The scope of these constants are always local, namely within a particular lemma, proposition, or proof, such that the apparently same constant in different contexts may carry different values. For probable events, sometimes we will just say the event holds with “high probability” if the probability of failure is dominated by some polynomial  $\text{poly}(n, p)$  which diminished to zero whenever  $n$  or  $p$  is large, with “overwhelming probability” if the failure probability is dominated some exponential function  $\exp(\text{poly}(n, p))$  which diminishes to zero whenever  $n$  or  $p$  is large.

## A Technical Tools and Preliminaries

**Lemma A.1.** Let  $\psi(x)$  and  $\Psi(x)$  to denote the probability density function (pdf) and the cumulative distribution function (cdf) for the standard normal distribution:

$$\text{(Standard Normal pdf)} \quad \psi(x) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{x^2}{2}\right\} \quad (\text{A.1})$$

$$\text{(Standard Normal cdf)} \quad \Psi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x \exp\left\{-\frac{t^2}{2}\right\} dt, \quad (\text{A.2})$$

Suppose a random variable  $X \sim \mathcal{N}(0, \sigma^2)$ , with the pdf  $f_\sigma(x) = \frac{1}{\sigma} \psi\left(\frac{x}{\sigma}\right)$ , then for any  $t_2 > t_1$  we have

$$\int_{t_1}^{t_2} f_\sigma(x) dx = \Psi\left(\frac{t_2}{\sigma}\right) - \Psi\left(\frac{t_1}{\sigma}\right), \quad (\text{A.3})$$

$$\int_{t_1}^{t_2} x f_\sigma(x) dx = -\sigma \left[ \psi\left(\frac{t_2}{\sigma}\right) - \psi\left(\frac{t_1}{\sigma}\right) \right], \quad (\text{A.4})$$

$$\int_{t_1}^{t_2} x^2 f_\sigma(x) dx = \sigma^2 \left[ \Psi\left(\frac{t_2}{\sigma}\right) - \Psi\left(\frac{t_1}{\sigma}\right) \right] - \sigma \left[ t_2 \psi\left(\frac{t_2}{\sigma}\right) - t_1 \psi\left(\frac{t_1}{\sigma}\right) \right]. \quad (\text{A.5})$$

**Lemma A.2** (Taylor Expansion of Standard Gaussian cdf and pdf). Assume  $\psi(x)$  and  $\Psi(x)$  be defined as above. There exists some universal constant  $C_\psi > 0$  such that

$$|\psi(x) - [\psi(x_0) - x_0 \psi(x_0)(x - x_0)]| \leq C_\psi (x - x_0)^2, \quad (\text{A.6})$$

$$|\Psi(x) - [\Psi(x_0) + \psi(x_0)(x - x_0)]| \leq C_\psi (x - x_0)^2. \quad (\text{A.7})$$

**Lemma A.3** (Matrix Induced Norms). For any matrix  $\mathbf{A} \in \mathbb{R}^{p \times n}$ , the induced matrix norm from  $\ell^p \rightarrow \ell^q$  is defined as

$$\|\mathbf{A}\|_{\ell^p \rightarrow \ell^q} \doteq \sup_{\|\mathbf{x}\|_p=1} \|\mathbf{A}\mathbf{x}\|_q. \quad (\text{A.8})$$

In particular, we have

$$\|\mathbf{A}\|_{\ell^2 \rightarrow \ell^1} = \sup_{\|\mathbf{x}\|_2=1} \sum_{k=1}^p |\mathbf{a}_k^\top \mathbf{x}|, \quad \|\mathbf{A}\|_{\ell^2 \rightarrow \ell^\infty} = \max_{1 \leq k \leq p} \|\mathbf{a}^k\|_2, \quad (\text{A.9})$$

$$\|\mathbf{A}\mathbf{B}\|_{\ell^p \rightarrow \ell^r} \leq \|\mathbf{A}\|_{\ell^q \rightarrow \ell^r} \|\mathbf{B}\|_{\ell^p \rightarrow \ell^q}, \quad (\text{A.10})$$

and  $\mathbf{A}$  and  $\mathbf{B}$  are any matrices of compatible size.

**Lemma A.4** (Moments of the Gaussian Random Variable). *If  $X \sim \mathcal{N}(0, \sigma_X^2)$ , then it holds for all integer  $m \geq 1$  that*

$$\mathbb{E}[|X|^m] = \sigma_X^m (m-1)!! \left[ \sqrt{\frac{2}{\pi}} \mathbb{1}_{m=2k+1} + \mathbb{1}_{m=2k} \right] \leq \sigma_X^m (m-1)!!, \quad k = \lfloor m/2 \rfloor. \quad (\text{A.11})$$

**Lemma A.5** (Moments of the  $\chi$  Random Variable). *If  $X \sim \chi(n)$ , i.e.,  $X \equiv_d \|\mathbf{x}\|_2$ <sup>6</sup> for  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Then it holds for all integer  $m \geq 1$  that*

$$\mathbb{E}[X^m] = 2^{m/2} \frac{\Gamma(m/2 + n/2)}{\Gamma(n/2)} \leq m! n^{m/2} \quad (\text{A.12})$$

**Lemma A.6** (Moments of the  $\chi^2$  Random Variable). *If  $X \sim \chi^2(n)$ , i.e.,  $X \equiv_d \|\mathbf{x}\|_2^2$  for  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . Then it holds for all integer  $m \geq 1$  that*

$$\mathbb{E}[X^m] = 2^m \frac{\Gamma(m + n/2)}{\Gamma(n/2)} = \prod_{k=1}^m (n + 2k - 2) \leq \frac{m!}{2} (2n)^m. \quad (\text{A.13})$$

**Lemma A.7** (Moment-Control Bernstein's Inequality for Random Variables). *Let  $X_1, \dots, X_p$  be i.i.d. real-valued random variables. Suppose that there exist some positive number  $R$  and  $\sigma_X^2$  such that*

$$\mathbb{E}[|X_k|^m] \leq \frac{m!}{2} \sigma_X^2 R^{m-2}, \quad \text{for all integers } m \geq 2. \quad (\text{A.14})$$

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

$$\mathbb{P}[|S - \mathbb{E}[S]| \geq t] \leq 2 \exp\left(-\frac{pt^2}{2\sigma_X^2 + 2Rt}\right). \quad (\text{A.15})$$

**Lemma A.8** (Moment-Control Bernstein's Inequality for Random Vectors). *Let  $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^d$  be i.i.d. random vectors. Suppose there exist some positive number  $R$  and  $\sigma_X^2$  such that*

$$\mathbb{E}[\|\mathbf{x}_k\|_2^m] \leq \frac{m!}{2} \sigma_X^2 R^{m-2}, \quad \text{for all integers } m \geq 2. \quad (\text{A.16})$$

Let  $\mathbf{s} = \frac{1}{p} \sum_{k=1}^p \mathbf{s}_k$ , then for any  $t > 0$ , it holds that

$$\mathbb{P}[\|\mathbf{s} - \mathbb{E}[\mathbf{s}]\|_2 \geq t] \leq 2(d+1) \exp\left(-\frac{pt^2}{2\sigma_X^2 + 2Rt}\right). \quad (\text{A.17})$$

**Lemma A.9** (Hoeffding's Inequality). *Let  $X_1, \dots, X_p$  be independent random variables such that  $X_k$  takes its values in  $[a_k, b_k]$  almost surely for all  $1 \leq k \leq p$ . Let  $S = \sum_{k=1}^p (X_k - \mathbb{E}X_k)$ , then for every  $t > 0$ ,*

$$\mathbb{P}[S \geq t] \leq \exp\left(-\frac{2t^2}{\sum_{k=1}^p (b_k - a_k)^2}\right). \quad (\text{A.18})$$

**Lemma A.10** (Gaussian Concentration Inequality). *Let  $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_p)$ . Let  $f: \mathbb{R}^p \mapsto \mathbb{R}$  be an  $L$ -Lipschitz function. Then we have for all  $t > 0$  that*

$$\mathbb{P}[f(\mathbf{X}) - \mathbb{E}f(\mathbf{X}) \geq t] \leq \exp\left(-\frac{t^2}{2L^2}\right). \quad (\text{A.19})$$

<sup>6</sup>The notation  $\equiv_d$  means equivalent in distribution.

**Lemma A.11** (Bounding Maximum Norm of Gaussian Vector Sequence). *Let  $\mathbf{x}_1, \dots, \mathbf{x}_p$  be a sequence of (not necessarily independent) standard Gaussian vectors in  $\mathbb{R}^n$ . Then, it holds that*

$$\mathbb{P} \left[ \max_{i \in [p]} \|\mathbf{x}_i\|_2 > \sqrt{2 \log p} + 2\sqrt{n} \right] \leq \exp \left( -\frac{1}{2} n \right). \quad (\text{A.20})$$

*Proof.* Since the function  $\|\cdot\|_2$  is 1-Lipschitz, by Gaussian concentration inequality, for any  $i \in [p]$ , we have

$$\mathbb{P} \left[ \|\mathbf{x}_i\|_2 - \sqrt{\mathbb{E} \|\mathbf{x}_i\|_2^2} > t \right] \leq \mathbb{P} [\|\mathbf{x}_i\|_2 - \mathbb{E} \|\mathbf{x}_i\|_2 > t] \leq \exp \left( -\frac{t^2}{2} \right) \quad (\text{A.21})$$

for all  $t > 0$ . Since  $\mathbb{E} \|\mathbf{x}_i\|_2^2 = n$ , by a simple union bound, we obtain

$$\mathbb{P} \left[ \max_{i \in [p]} \|\mathbf{x}_i\|_2 > \sqrt{n} + t \right] \leq \exp \left( -\frac{t^2}{2} + \log p \right) \quad (\text{A.22})$$

for all  $t > 0$ . Taking  $t = \sqrt{2 \log p} + \sqrt{n}$  and simplifying the terms gives the claimed result.  $\square$

**Lemma A.12** (Covering Number of a Unit Ball). *Let  $\mathcal{B} = \{\mathbf{x} \in \mathbb{R}^n \mid \|\mathbf{x}\|_2 \leq 1\}$  be a unit ball. For any  $\varepsilon \in (0, 1)$ , there exists some  $\varepsilon$  cover of  $\mathcal{B}$  w.r.t. the normal  $\mathbb{R}^n$  metric, denoted as  $\mathcal{N}_\varepsilon$ , such that*

$$|\mathcal{N}_\varepsilon| \leq \left( 1 + \frac{2}{\varepsilon} \right)^n \leq \left( \frac{3}{\varepsilon} \right)^n. \quad (\text{A.23})$$

**Lemma A.13** (Spectrum of Gaussian Matrices, [Ver10]). *Let  $\mathbf{A} \in \mathbb{R}^{p \times n}$  ( $p > n$ ) contain i.i.d. standard normal entries. Then for every  $t \geq 0$ , with probability at least  $1 - 2 \exp(-t^2/2)$ , one has*

$$\sqrt{p} - \sqrt{n} - t \leq \sigma_{\min}(\mathbf{A}) \leq \sigma_{\max}(\mathbf{A}) \leq \sqrt{p} + \sqrt{n} + t. \quad (\text{A.24})$$

**Lemma A.14.** *For any  $\varepsilon \in (0, 1)$ , there exists a constant  $C(\varepsilon) > 1$ , such that provided  $n_1 > C(\varepsilon) n_2$ , the random matrix  $\Phi \in \mathbb{R}^{n_1 \times n_2} \sim_{i.i.d.} \mathcal{N}(0, 1)$  obeys*

$$(1 - \varepsilon) \sqrt{\frac{2}{\pi}} n_1 \|\mathbf{x}\|_2 \leq \|\Phi \mathbf{x}\|_1 \leq (1 + \varepsilon) \sqrt{\frac{2}{\pi}} n_1 \|\mathbf{x}\|_2 \quad \text{for all } \mathbf{x} \in \mathbb{R}^{n_2}, \quad (\text{A.25})$$

with probability at least  $1 - 2 \exp(-c(\varepsilon) n_2)$  for some  $c(\varepsilon) > 0$ .

Geometrically, this lemma roughly corresponds to the well known almost spherical section theorem [FLM77, GG84], see also [GM03]. A slight variant of this version has been proved in [Don06], borrowing ideas from [Pis99].

*Proof.* By homogeneity, it is enough to consider all  $\mathbf{x}$  with unit norms. For a fixed  $\mathbf{x}_0$  with  $\|\mathbf{x}_0\|_2 = 1$ ,  $\Phi \mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . So  $\mathbb{E} \|\Phi \mathbf{x}_0\|_1 = \sqrt{\frac{2}{\pi}} n_1$ . By concentration of measure for Gaussian vectors,

$$\mathbb{P} [|\|\Phi \mathbf{x}_0\|_1 - \mathbb{E} [\|\Phi \mathbf{x}_0\|_1]| > t] \leq 2 \exp \left( -\frac{t^2}{2n_1} \right) \quad (\text{A.26})$$

for any  $t > 0$ . For a fixed  $\delta \in (0, 1)$ ,  $\mathcal{S}^{n_2-1}$  can be covered by a  $\delta$ -net  $N_\delta$  with cardinality  $\#N_\delta \leq (1 + 2/\delta)^{n_2}$ . Now consider the event

$$\mathcal{E} \doteq \left\{ (1 - \delta) \sqrt{\frac{2}{\pi}} n_1 \leq \|\Phi \mathbf{x}\|_1 \leq (1 + \delta) \sqrt{\frac{2}{\pi}} n_1 \quad \forall \mathbf{x} \in N_1 \right\}. \quad (\text{A.27})$$

A simple application of union bound yields

$$\mathbb{P}[\mathcal{E}^c] \leq 2 \exp\left(-\frac{\delta^2 n_1}{\pi} + n_2 \log\left(1 + \frac{2}{\delta}\right)\right). \quad (\text{A.28})$$

Choosing  $\delta$  small enough such that

$$(1 - 3\delta)(1 - \delta)^{-1} \geq 1 - \varepsilon \text{ and } (1 + \delta)(1 - \delta)^{-1} \leq 1 + \varepsilon, \quad (\text{A.29})$$

then conditioned on  $\mathcal{E}$ , we can conclude that

$$(1 - \varepsilon) \sqrt{\frac{2}{\pi}} n_1 \leq \|\Phi \mathbf{x}\|_1 \leq (1 + \varepsilon) \sqrt{\frac{2}{\pi}} n_1 \quad \forall \mathbf{x} \in \mathbb{S}^{n_2-1}. \quad (\text{A.30})$$

Indeed, suppose  $\mathcal{E}$  holds. Then it can easily be seen that any  $\mathbf{z} \in \mathbb{S}^{n_2-1}$  can be written as

$$\mathbf{z} = \sum_{k=0}^{\infty} \lambda_k \mathbf{x}_k, \quad \text{with } |\lambda_k| \leq \delta^k, \mathbf{x}_k \in N_1 \text{ for all } k. \quad (\text{A.31})$$

Hence we have

$$\|\Phi \mathbf{z}\|_1 = \left\| \Phi \sum_{k=0}^{\infty} \lambda_k \mathbf{x}_k \right\|_1 \leq \sum_{k=0}^{\infty} \delta^k \|\Phi \mathbf{x}_k\|_1 \leq (1 + \delta)(1 - \delta)^{-1} \sqrt{\frac{2}{\pi}} n_1. \quad (\text{A.32})$$

Similarly,

$$\|\Phi \mathbf{z}\|_1 = \left\| \Phi \sum_{k=0}^{\infty} \lambda_k \mathbf{x}_k \right\|_1 \geq \left[1 - \delta - \delta(1 + \delta)(1 - \delta)^{-1}\right] \sqrt{\frac{2}{\pi}} n_1 = (1 - 3\delta)(1 - \delta)^{-1} \sqrt{\frac{2}{\pi}} n_1. \quad (\text{A.33})$$

Hence, the choice of  $\delta$  above leads to the claimed result. To make  $\mathbb{P}[\mathcal{E}^c]$  small, it is enough to choose  $C$  such that

$$C\delta^2/\pi > \log\left(1 + \frac{2}{\delta}\right). \quad (\text{A.34})$$

Setting  $C = 2 \log\left(1 + \frac{2}{\delta}\right) \pi / \delta^2$  completes the proof.  $\square$

**Lemma A.15.** Suppose  $n_1 \leq \frac{1}{2} \exp(n_2/2)$ . Fix  $\varepsilon \in (0, 1)$ . Then for any  $\xi$  such that  $\xi^2 > 2 \log(1 + 2\varepsilon)$ . The random matrix  $\Phi \in \mathbb{R}^{n_1 \times n_2} \sim_{i.i.d.} \mathcal{N}(0, 1)$  obeys

$$\|\Phi \mathbf{x}\|_{\infty} \leq \frac{1 + \xi}{1 - \varepsilon} \sqrt{n_2} \|\mathbf{x}\|_2 \quad \text{for all } \mathbf{x} \in \mathbb{R}^{n_2}, \quad (\text{A.35})$$

with probability at least  $1 - \exp(-n_2(\xi^2/2 - \log(1 + 2\varepsilon)))$ .

*Proof.* Again for a fixed  $\mathbf{x}_0 \in \mathbb{S}^{n_2-1}$ ,  $\Phi \mathbf{x}_0 \equiv_d \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ . For any fixed  $\beta > 0$  to be decided later,

$$\mathbb{E}[\beta \|\mathbf{v}\|_{\infty}] = \mathbb{E}\left[\beta \max_{i \in [n_1]} |v_i|\right] = \mathbb{E}\left[\log \max_{i \in [n_1]} \exp(\beta |v_i|)\right] \leq \log \mathbb{E}\left[\max_{i \in [n_1]} \exp(\beta |v_i|)\right] \quad (\text{A.36})$$

$$\leq \log \mathbb{E}\left[\sum_{i=1}^{n_1} \exp(\beta |v_i|)\right] = \log n_1 \mathbb{E}[\exp(\beta |v_1|)] \leq \log 2n_1 \exp(\beta^2/2). \quad (\text{A.37})$$

Hence

$$\mathbb{E}[\|\mathbf{v}\|_{\infty}] \leq \frac{\log 2n_1 \exp(\beta^2/2)}{\beta}. \quad (\text{A.38})$$

Taking  $\beta = \sqrt{2 \log(2n_1)}$ , we obtain

$$\mathbb{E} [\|\mathbf{v}\|_\infty] \leq \sqrt{2 \log(2n_1)}. \quad (\text{A.39})$$

Because the mapping  $\mathbf{v} \mapsto \|\mathbf{v}\|_\infty$  is 1-Lipschitz, by concentration of measure for Gaussian vectors, we obtain

$$\mathbb{P} [\|\Phi \mathbf{x}\|_\infty - \mathbb{E} [\|\Phi \mathbf{x}\|_\infty] > t] \leq \exp\left(-\frac{t^2}{2}\right). \quad (\text{A.40})$$

Taking  $t = \xi \sqrt{n_2}$ , and consider an  $\varepsilon$ -net  $\mathcal{N}_\varepsilon$  that covers  $\mathcal{S}^{n_2-1}$  with cardinality  $|\mathcal{N}_\varepsilon| \leq (1 + 2/\varepsilon)^{n_2}$ , we have the event

$$\mathcal{E} \doteq \{\|\Phi \mathbf{x}\|_\infty \leq (1 + \xi) \sqrt{n_2} \forall \mathbf{x} \in \mathcal{N}_\varepsilon\} \quad (\text{A.41})$$

holds with probability at least  $1 - \exp(-\xi^2 n_2/2 + n_2 \log(1 + 2\varepsilon))$ . Conditioned on  $\mathcal{E}$ , we have

$$\sup_{\|\mathbf{z}\|_2=1} \|\Phi \mathbf{z}\|_\infty \leq \sup_{\mathbf{z}' \in \mathcal{N}_\varepsilon} \|\Phi \mathbf{z}'\|_\infty + \sup_{\|\mathbf{e}\|_2 \leq \varepsilon} \|\Phi \mathbf{e}\|_\infty = \sup_{\mathbf{z}' \in \mathcal{N}_\varepsilon} \|\Phi \mathbf{z}'\|_\infty + \varepsilon \sup_{\|\mathbf{e}\|_2=1} \|\Phi \mathbf{e}\|_\infty. \quad (\text{A.42})$$

Hence we have

$$\sup_{\|\mathbf{z}\|_2=1} \|\Phi \mathbf{z}\|_\infty \leq \frac{1}{1 - \varepsilon} \sup_{\mathbf{z}' \in \mathcal{N}_\varepsilon} \|\Phi \mathbf{z}'\|_\infty = \frac{1 + \xi}{1 - \varepsilon} \sqrt{n_2}, \quad (\text{A.43})$$

completing the proof.  $\square$

## B The Random Basis vs. Its Orthonormalized Version

We consider  $\mathbf{Y}$  obeying the planted sparse model:

$$\mathbf{Y} = [\mathbf{x}_0 \mid \mathbf{G}] \in \mathbb{R}^{p \times n} \quad (\text{B.1})$$

with

$$\mathbf{x}_0 \sim_{i.i.d.} \frac{1}{\sqrt{\theta p}} \text{Ber}(\theta), \mathbf{G} \sim_{i.i.d.} \mathcal{N}\left(0, \frac{1}{p}\right). \quad (\text{B.2})$$

One ‘‘natural/canonical’’ orthonormal basis for the subspace spanned by columns of  $\mathbf{Y}$  is

$$\mathbf{Y}' = \left[ \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \mid \mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G} \left( \mathbf{G}^\top \mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G} \right)^{-1/2} \right]. \quad (\text{B.3})$$

We also write  $\mathbf{G}' \doteq \mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G} \left( \mathbf{G}^\top \mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G} \right)^{-1/2}$  for convenience. In this section, we want to show that the intuition  $\mathbf{Y}'$  well approximating  $\mathbf{Y}^7$  can be made rigorous. These results are needed when we prove Theorem 2.1 for the global optimality of the natural  $\ell^2$  constrained formulation (2.1), as well as when we translate the results for  $\mathbf{Y}$  to quantitative statements about  $\mathbf{Y}'$  in Appendix F.4.

For any realization of  $\mathbf{x}_0$ , let the support (index set of nonzero elements) of  $\mathbf{x}_0$  be  $\mathcal{I}$ . By Hoeffding’s inequality in Lemma A.9, we have the event

$$\mathcal{E}_0 \doteq \left\{ \frac{1}{2} \theta p \leq |\mathcal{I}| \leq 2 \theta p \right\} \quad (\text{B.4})$$

holds with probability at least  $1 - 2 \exp(-p\theta^2/2)$ . Moreover, we show the following:

<sup>7</sup>When  $n$  and  $p$  are large,  $\mathbf{Y}$  has nearly orthonormal columns.



**Lemma B.1.** *The bound*

$$\left|1 - \frac{1}{\|\mathbf{x}_0\|_2}\right| \leq \frac{2\sqrt{2}}{5} \sqrt{\frac{n \log p}{\theta^3 p}} \quad (\text{B.5})$$

holds with probability at least  $1 - 2 \exp(-p\theta^2/2) - 2 \exp(-2n \log p)$ .

*Proof.* Because  $\mathbb{E}[\|\mathbf{x}_0\|_2^2] = 1$ , by Hoeffding's inequality in Lemma A.9, we have

$$\mathbb{P}\left[\left|\|\mathbf{x}_0\|_2^2 - \mathbb{E}[\|\mathbf{x}_0\|_2^2]\right| > t\right] = \mathbb{P}\left[\left|\|\mathbf{x}_0\|_2^2 - 1\right| > t\right] \leq 2 \exp(-2\theta^2 p t^2) \quad (\text{B.6})$$

for all  $t > 0$ , which implies

$$\mathbb{P}\left[|\|\mathbf{x}_0\|_2 - 1| (\|\mathbf{x}_0\|_2 + 1) > t\right] \leq 2 \exp(-2\theta^2 p t^2). \quad (\text{B.7})$$

On the intersection with  $\mathcal{E}_0$ ,  $\|\mathbf{x}_0\|_2 + 1 \leq \sqrt{2} + 1 \leq 5/2$ , and setting  $t = \sqrt{\frac{n \log p}{\theta^3 p}}$ , we obtain

$$\mathbb{P}\left[\left|\|\mathbf{x}_0\|_2 - 1\right| > \frac{2}{5} \sqrt{\frac{n \log p}{\theta^3 p}}\right] \leq 2 \exp(-2n \log p). \quad (\text{B.8})$$

So we obtain that with probability at least  $1 - 2 \exp(-p\theta^2/2) - 2 \exp(-2n \log p)$ ,

$$\left|1 - \frac{1}{\|\mathbf{x}_0\|_2}\right| = \frac{|1 - \|\mathbf{x}_0\|_2|}{\|\mathbf{x}_0\|_2} \leq \frac{2\sqrt{2}}{5} \sqrt{\frac{n \log p}{\theta^3 p}}, \quad (\text{B.9})$$

as desired.  $\square$

Next, let  $\mathbf{M} \doteq \left(\mathbf{G}^\top \mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G}\right)^{-1/2}$ , then  $\mathbf{G}' = \mathbf{G}\mathbf{M} - \frac{\mathbf{x}_0 \mathbf{x}_0^\top}{\|\mathbf{x}_0\|_2^2} \mathbf{G}\mathbf{M}$ , we show the following results hold:

**Lemma B.2.** *Provided  $p \geq Cn \geq 2$  for some large enough constant  $C$ , it holds that*

$$\|\mathbf{M}\| \leq 2, \quad \|\mathbf{M} - \mathbf{I}\| \leq 4\sqrt{\frac{n}{p}} \quad (\text{B.10})$$

with probability at least  $1 - c' \exp(-c'n)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* First observe that

$$\|\mathbf{M}\| = \left(\sigma_{\min}\left(\mathbf{G}^\top \mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G}\right)\right)^{-1/2} = \left(\sigma_{\min}\left(\mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G}\right)\right)^{-1}. \quad (\text{B.11})$$

Now suppose  $\mathbf{B}$  is an orthonormal basis spanning  $\mathbf{x}_0^\perp$ . Then it is not hard to see the spectrum of  $\mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G}$  is the same as that of  $\mathbf{B}^\top \mathbf{G} \in \mathbb{R}^{(p-1) \times (n-1)}$ ; in particular,

$$\sigma_{\min}\left(\mathcal{P}_{\mathbf{x}_0^\perp} \mathbf{G}\right) = \sigma_{\min}\left(\mathbf{B}^\top \mathbf{G}\right). \quad (\text{B.12})$$

Since  $\mathbf{G} \sim_{i.i.d.} \mathcal{N}\left(0, \frac{1}{p}\right)$ , and  $\mathbf{B}^\top$  has orthonormal rows,  $\mathbf{B}^\top \mathbf{G} \sim_{i.i.d.} \mathcal{N}\left(0, \frac{1}{p}\right)$ , we can invoke the spectrum results for Gaussian matrices in Lemma A.13 and obtain that

$$\sqrt{\frac{p-1}{p}} - 2\sqrt{\frac{n-1}{p-1}} \leq \sigma_{\min}\left(\mathbf{B}^\top \mathbf{G}\right) \leq \sigma_{\max}\left(\mathbf{B}^\top \mathbf{G}\right) \leq \sqrt{\frac{p-1}{p}} + 2\sqrt{\frac{n-1}{p-1}} \quad (\text{B.13})$$

with probability at least  $1 - c_1 \exp(-c_2 n)$  for some  $c_1, c_2 > 0$ . Thus, when  $p \geq C_1 n$  for some large constant  $C_1$ , we have

$$\|\mathbf{M}\| = \left( \sqrt{\frac{p-1}{p}} - 2\sqrt{\frac{n-1}{p-1}} \right)^{-1} \leq 2, \quad (\text{B.14})$$

$$\|\mathbf{I} - \mathbf{M}\| = \max(|\sigma_{\max}(\mathbf{M}) - 1|, |\sigma_{\min}(\mathbf{M}) - 1|) \leq 2\sqrt{\frac{n-1}{p-1}} \left( \sqrt{\frac{p-1}{p}} - 2\sqrt{\frac{n-1}{p-1}} \right)^{-1} \leq 4\sqrt{\frac{n}{p}}, \quad (\text{B.15})$$

with probability at least  $1 - c_1 \exp(-c_2 n)$ .  $\square$

**Lemma B.3.** *There exists a constant  $C > 0$ , such that when  $p \geq Cn$ , the following*

$$\|\mathbf{Y}\|_{\ell^2 \rightarrow \ell^1} \leq 3\sqrt{p}, \quad (\text{B.16})$$

$$\|\mathbf{Y}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq 7\sqrt{2\theta p}, \quad (\text{B.17})$$

$$\|\mathbf{Y}_{\mathcal{I}} - \mathbf{Y}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \frac{10}{\theta} \sqrt{n \log p}, \quad (\text{B.18})$$

$$\|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \leq 8\sqrt{n}, \quad (\text{B.19})$$

$$\|\mathbf{Y} - \mathbf{Y}'\|_{\ell^2 \rightarrow \ell^1} \leq \frac{10}{\theta} \sqrt{n \log p} \quad (\text{B.20})$$

hold simultaneously with probability at least  $1 - c' \exp(-c'' n) - 2 \exp(-p\theta^2/2)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* First of all, we have

$$\left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top \mathbf{G} \mathbf{M}}{\|\mathbf{x}_0\|_2^2} \right\|_{\ell^2 \rightarrow \ell^1} \leq \frac{1}{\|\mathbf{x}_0\|_2^2} \|\mathbf{x}_0\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{x}_0^\top \mathbf{G} \mathbf{M}\|_{\ell^2 \rightarrow \ell^2} = \frac{2}{\|\mathbf{x}_0\|_2^2} \|\mathbf{x}_0\|_1 \|\mathbf{x}_0^\top \mathbf{G}\|_2, \quad (\text{B.21})$$

where in the last inequality we have applied the fact  $\|\mathbf{M}\| \leq 2$  from Lemma B.2. Now  $\mathbf{x}_0^\top \mathbf{G}$  is an i.i.d. Gaussian vectors with each entry distributed as  $\mathcal{N}\left(0, \frac{\|\mathbf{x}_0\|_2^2}{p}\right)$ , where  $\|\mathbf{x}_0\|_2^2 = \frac{|\mathcal{I}|}{\theta p}$ . So by measure concentration inequality for Gaussian vectors, we have

$$\|\mathbf{x}_0^\top \mathbf{G}\|_2 \leq 2 \|\mathbf{x}_0\|_2 \sqrt{\frac{n}{p}} \quad (\text{B.22})$$

with probability at least  $1 - \exp(-n/2)$ . On the intersection with  $\mathcal{E}_0$ , this implies

$$\left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top \mathbf{G} \mathbf{M}}{\|\mathbf{x}_0\|_2^2} \right\|_{\ell^2 \rightarrow \ell^1} \leq 4\sqrt{|\mathcal{I}|} \sqrt{\frac{n}{p}} \leq 4\sqrt{2\theta n}, \quad (\text{B.23})$$

with probability at least  $1 - \exp(-n/2) - 2 \exp(-p\theta^2/2)$ . Moreover, when intersected with  $\mathcal{E}_0$ , Lemma A.14 implies that when  $p \geq \Omega(n)$ ,

$$\|\mathbf{G}\|_{\ell^2 \rightarrow \ell^1} \leq \sqrt{p}, \quad \|\mathbf{G}_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \sqrt{2\theta p} \quad (\text{B.24})$$

with probability at least  $1 - c_1 \exp(-c_2 n) - 2 \exp(-p\theta^2/2)$ , for some positive constants  $c_1, c_2$ . So when

$p \geq \Omega(n)$ ,

$$\|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{G} - \mathbf{GM}\|_{\ell^2 \rightarrow \ell^1} + \left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top}{\|\mathbf{x}_0\|_2^2} \mathbf{GM} \right\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{G}\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{I} - \mathbf{M}\| + 4\sqrt{2\theta n} \leq 8\sqrt{n}, \quad (\text{B.25})$$

$$\|\mathbf{Y}\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{x}_0\|_{\ell^2 \rightarrow \ell^1} + \|\mathbf{G}\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{x}_0\|_1 + \sqrt{p} \leq 2\sqrt{\theta p} + \sqrt{p} \leq 3\sqrt{p}, \quad (\text{B.26})$$

$$\|\mathbf{G}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{G}_{\mathcal{I}} \mathbf{M}\|_{\ell^2 \rightarrow \ell^1} + \left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top}{\|\mathbf{x}_0\|_2^2} \mathbf{GM} \right\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{G}_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{M}\| + 4\sqrt{2\theta n} \leq 6\sqrt{2\theta p}, \quad (\text{B.27})$$

$$\|\mathbf{G}_{\mathcal{I}} - \mathbf{G}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{G}_{\mathcal{I}} (\mathbf{I} - \mathbf{M})\|_{\ell^2 \rightarrow \ell^1} + \left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top}{\|\mathbf{x}_0\|_2^2} \mathbf{GM} \right\|_{\ell^2 \rightarrow \ell^1} \leq \|\mathbf{G}_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{I} - \mathbf{M}\| + 4\sqrt{2\theta n} \leq 8\sqrt{2\theta n}, \quad (\text{B.28})$$

$$\|\mathbf{Y}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_{\ell^2 \rightarrow \ell^1} + \|\mathbf{G}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \frac{\|\mathbf{x}_0\|_1}{\|\mathbf{x}_0\|_2} + 6\sqrt{2\theta p} \leq 7\sqrt{2\theta p} \quad (\text{B.29})$$

with probability at least  $1 - c_3 \exp(-c_4 n) - 2 \exp(-p\theta^2/2)$  for some positive constants  $c_3, c_4$ , where we have used the above estimates and the results in Lemma B.2. Finally, by Lemma B.1, we obtain

$$\|\mathbf{Y} - \mathbf{Y}'\|_{\ell^2 \rightarrow \ell^1} \leq \left| 1 - \frac{1}{\|\mathbf{x}_0\|_2} \right| \|\mathbf{x}_0\|_1 + \|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \leq \frac{10}{\theta} \sqrt{n \log p}, \quad (\text{B.30})$$

$$\|\mathbf{Y}_{\mathcal{I}} - \mathbf{Y}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \left| 1 - \frac{1}{\|\mathbf{x}_0\|_2} \right| \|\mathbf{x}_0\|_1 + \|\mathbf{G}_{\mathcal{I}} - \mathbf{G}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \leq \frac{10}{\theta} \sqrt{n \log p}, \quad (\text{B.31})$$

holding with probability at least  $1 - c_5 \exp(-c_6 n) - 2 \exp(-p\theta^2/2)$  for some positive constants  $c_5, c_6$ .  $\square$

**Lemma B.4.** *Provided  $Cn \leq p \leq \exp(n/2)/2$  for some constant  $C > 0$ , the following*

$$\|\mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} \leq 16\sqrt{\frac{n}{\theta p}}, \quad (\text{B.32})$$

$$\|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} \leq \frac{32n}{\sqrt{\theta p}} \quad (\text{B.33})$$

hold simultaneously with probability at least  $1 - c' \exp(-c'' n) - 2 \exp(-p\theta^2/2)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* First of all, we have

$$\left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top}{\|\mathbf{x}_0\|_2^2} \mathbf{GM} \right\|_{\ell^2 \rightarrow \ell^\infty} \leq \frac{1}{\|\mathbf{x}_0\|_2^2} \|\mathbf{x}_0\|_{\ell^2 \rightarrow \ell^\infty} \|\mathbf{x}_0^\top \mathbf{GM}\|_{\ell^2 \rightarrow \ell^2} = \frac{2}{\|\mathbf{x}_0\|_2^2} \|\mathbf{x}_0\|_\infty \|\mathbf{x}_0^\top \mathbf{G}\|_2, \quad (\text{B.34})$$

where at the last inequality we have applied the fact  $\|\mathbf{M}\| \leq 2$  from Lemma B.2. Similar to the proof to Lemma B.3, we have that  $\|\mathbf{x}_0^\top \mathbf{G}\|_2 \leq 2\|\mathbf{x}_0\|_2 \sqrt{n/p}$  with probability at least  $1 - \exp(-n/2)$ . So on the intersection with  $\mathcal{E}_0$ , we obtain that

$$\left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top}{\|\mathbf{x}_0\|_2^2} \mathbf{GM} \right\|_{\ell^2 \rightarrow \ell^\infty} \leq \frac{4\|\mathbf{x}_0\|_\infty}{\|\mathbf{x}_0\|_2} \sqrt{\frac{n}{p}} \leq \frac{4\sqrt{2n}}{\sqrt{\theta p}} \quad (\text{B.35})$$

holds with probability at least  $1 - \exp(-n/2) - 2 \exp(-p\theta^2/2)$ . Now taking  $\xi = 2$  and  $\varepsilon = 1/2$  in Lemma A.15, we have that

$$\|\mathbf{G}\|_{\ell^2 \rightarrow \ell^\infty} \leq 6\sqrt{\frac{n}{p}} \quad (\text{B.36})$$

with probability at least  $1 - \exp(-n(2 - \log 2))$ . Combining with results in Lemma B.2, we obtain

$$\begin{aligned} \|\mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} &\leq \|\mathbf{G}\mathbf{M}\|_{\ell^2 \rightarrow \ell^\infty} + \left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top \mathbf{G}\mathbf{M}}{\|\mathbf{x}_0\|_2^2} \right\|_{\ell^2 \rightarrow \ell^\infty} \\ &\leq \|\mathbf{G}\|_{\ell^2 \rightarrow \ell^\infty} \|\mathbf{M}\| + \frac{4\sqrt{2n}}{\sqrt{\theta p}} \leq 12\sqrt{\frac{n}{p}} + \frac{4\sqrt{2n}}{\sqrt{\theta p}} \leq 16\sqrt{\frac{n}{\theta p}}, \end{aligned} \quad (\text{B.37})$$

$$\|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} \leq \|\mathbf{G}\|_{\ell^2 \rightarrow \ell^\infty} \|\mathbf{I} - \mathbf{M}\| + \left\| \frac{\mathbf{x}_0 \mathbf{x}_0^\top \mathbf{G}\mathbf{M}}{\|\mathbf{x}_0\|_2^2} \right\|_{\ell^2 \rightarrow \ell^\infty} \leq \frac{24n}{p} + \frac{4\sqrt{2n}}{\sqrt{\theta p}} \leq \frac{32n}{\sqrt{\theta p}} \quad (\text{B.38})$$

with probability at least  $1 - c_7 \exp(-c_8 n) - 2 \exp(-p\theta^2/2)$  for some positive constants  $c_7, c_8$ .  $\square$

## C Proof of $\ell^1/\ell^2$ Global Optimality

*Proof.* We will first analyze a canonical version, in which the input basis is  $\mathbf{Y}'$ :

$$\min_{\mathbf{q} \in \mathbb{R}^n} \|\mathbf{Y}'\mathbf{q}\|_1, \quad \text{s.t. } \|\mathbf{q}\|_2 = 1. \quad (\text{C.1})$$

Let  $\mathbf{q} = [q_1; \mathbf{q}_2]$ . For any fixed support  $\mathcal{I}$  of  $\mathbf{x}_0$ , we have

$$\begin{aligned} \|\mathbf{Y}'\mathbf{q}\|_1 &= \|\mathbf{Y}'_{\mathcal{I}}\mathbf{q}\|_1 + \|\mathbf{Y}'_{\mathcal{I}^c}\mathbf{q}\|_1 \\ &\geq |q_1| \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 - \|\mathbf{G}'_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}'_{\mathcal{I}^c}\mathbf{q}_2\|_1 \\ &\geq |q_1| \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 - \|\mathbf{G}_{\mathcal{I}}\mathbf{q}_2\|_1 - \|(\mathbf{G}_{\mathcal{I}} - \mathbf{G}'_{\mathcal{I}})\mathbf{q}_2\|_1 + \|\mathbf{G}_{\mathcal{I}^c}\mathbf{q}_2\|_1 - \|(\mathbf{G}_{\mathcal{I}^c} - \mathbf{G}'_{\mathcal{I}^c})\mathbf{q}_2\|_1 \\ &\geq |q_1| \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 - \|\mathbf{G}_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}_{\mathcal{I}^c}\mathbf{q}_2\|_1 - \|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{q}_2\|_2. \end{aligned} \quad (\text{C.2})$$

By Lemma A.14 and intersecting with  $\mathcal{E}_0$ , we have that as long as  $p \geq \Omega(n)$ , there exists constant  $c_1 > 0$  such that

$$\|\mathbf{G}_{\mathcal{I}}\mathbf{q}_2\|_1 \leq \frac{2\theta p}{\sqrt{p}} \|\mathbf{q}_2\|_2 = 2\theta\sqrt{p} \|\mathbf{q}_2\|_2 \text{ for all } \mathbf{q}_2 \in \mathbb{R}^{n-1}, \quad (\text{C.3})$$

$$\|\mathbf{G}_{\mathcal{I}^c}\mathbf{q}_2\|_1 \geq \frac{1}{2} \frac{p - 2\theta p}{\sqrt{p}} \|\mathbf{q}_2\|_2 = \frac{1}{2} \sqrt{p} (1 - 2\theta) \|\mathbf{q}_2\|_2 \text{ for all } \mathbf{q}_2 \in \mathbb{R}^{n-1}, \quad (\text{C.4})$$

hold with probability at least  $1 - 2 \exp(-c_1 n_2) - 2 \exp(-p\theta^2/2)$ . Moreover, by Lemma B.3,

$$\|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \leq 8\sqrt{n} \quad (\text{C.5})$$

holds with probability at least  $1 - c_2 \exp(-c_3 n) - 2 \exp(-p\theta^2/2)$  when  $p \geq \Omega(n)$ . So we obtain that

$$\|\mathbf{Y}'\mathbf{q}\|_1 \geq |q_1| \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 + \|\mathbf{q}_2\|_2 \left( \frac{1}{2} \sqrt{p} (1 - 2\theta) - 2\theta\sqrt{p} - 8\sqrt{n} \right) \quad (\text{C.6})$$

holds with probability at least  $1 - c_4 \exp(-c_5 n) - 2 \exp(-p\theta^2/2)$  for some positive  $c_4$  and  $c_5$ . Assuming  $\mathcal{E}_0$ , we observe

$$\left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 \leq \sqrt{|\mathcal{I}|} \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_2 \leq \sqrt{2\theta p}. \quad (\text{C.7})$$

So in order to minimize the objective  $\|\mathbf{Y}'\mathbf{q}\|_1$  at  $\mathbf{e}_1$  or  $-\mathbf{e}_1$ , i.e.,  $q_1 = 1$ , subject to the constraint  $q_1^2 + \|\mathbf{q}_2\|_2^2 = 1$ , it suffices to have

$$\sqrt{2\theta p} < \frac{1}{2}\sqrt{p}(1 - 2\theta) - 2\theta\sqrt{p} - 8\sqrt{n}, \quad (\text{C.8})$$

which is satisfied when  $\theta$  is sufficiently small. Thus there exists a universal constant  $\theta_0 > 0$ , such that for all  $1/\sqrt{n} \leq \theta \leq \theta_0$ , when  $p \geq \Omega(n)$ ,  $\pm\mathbf{e}_1$  are the global minimizers of (2.2) with probability at least  $1 - c_4 \exp(-c_5 n) - 2 \exp(-p\theta^2/2)$ , if the input basis is  $\mathbf{Y}'$ . As  $\theta > 1/\sqrt{n}$  by assumption, from (B.4), to make the above probability high, it is enough to make  $p \geq \Omega(n \log n)$ .

Any other input basis can be written as  $\mathbf{Y}'\mathbf{R}$ , for some orthogonal matrix  $\mathbf{R}$ . The program now is written as

$$\min_{\mathbf{q} \in \mathbb{R}^n} \|\mathbf{Y}'\mathbf{R}\mathbf{q}\|_1, \quad \text{s.t. } \|\mathbf{q}\|_2 = 1, \quad (\text{C.9})$$

which is equivalent to

$$\min_{\mathbf{q} \in \mathbb{R}^n} \|\mathbf{Y}'\mathbf{R}\mathbf{q}\|_1, \quad \text{s.t. } \|\mathbf{R}\mathbf{q}\|_2 = 1, \quad (\text{C.10})$$

which is obviously equivalent to the canonical program we analyze above by a simple change of variable, i.e.,  $\bar{\mathbf{q}} \doteq \mathbf{R}\mathbf{q}$ , completing the proof.  $\square$

## D Proof of Main Result

In this appendix, we prove our main result in Theorem 4.1. In particular, we will first show that when the  $\mathbf{Y}'$  defined in (B.3) is the input orthonormal basis, the ‘‘initialization + ADM + LP rounding’’ pipeline recovers  $\mathbf{x}_0$  under the stated technical conditions. Then we will upgrade the recovery result to all orthonormal basis by observing that all three stages are ‘‘invariant’’ to the input orthonormal basis  $\mathbf{Y}$ .

Keep the notation in Section 3, let  $\mathbf{y}^1, \dots, \mathbf{y}^p$  be the transpose of the rows of  $\mathbf{Y}$ , and let  $\mathbf{y}'^1, \dots, \mathbf{y}'^p$  be the transpose of the rows of  $\mathbf{Y}'$ . For  $\mathbf{q} \in \mathbb{S}^{n-1}$ , set

$$\mathbf{Q}(\mathbf{q}) = \frac{1}{p} \sum_{k=1}^p \mathbf{y}^k S_\lambda [\mathbf{q}^\top \mathbf{y}^k], \quad (\text{D.1})$$

$$\mathbf{Q}'(\mathbf{q}) = \frac{1}{p} \sum_{k=1}^p \mathbf{y}'^k S_\lambda [\mathbf{q}^\top \mathbf{y}'^k]. \quad (\text{D.2})$$

Further, we write  $\mathbf{Q}(\mathbf{q}) = [Q_1(\mathbf{q}); \mathbf{Q}_2(\mathbf{q})]$ , where  $Q_1(\mathbf{q})$  is the first coordinate, and define similar notations for  $\mathbf{Q}'(\mathbf{q})$ . In addition, for any  $k = 1, \dots, p$ , set

$$X_k^1(Z_k) = x_{0k} S_\lambda [\mathbf{q}^\top \mathbf{y}^k] = x_{0k} S_\lambda [x_{0k} q_1 + Z_k], \quad (\text{D.3})$$

$$\mathbf{X}_k^2(Z_k) = \mathbf{g}^k S_\lambda [\mathbf{q}^\top \mathbf{y}^k] = \mathbf{g}^k S_\lambda [x_{0k} q_1 + Z_k], \quad (\text{D.4})$$

where  $Z_k = \mathbf{q}_2^\top \mathbf{g}^k \sim \mathcal{N}(0, \sigma^2)$  for  $\sigma = \|\mathbf{q}_2\|_2 / \sqrt{p}$ , and  $x_{0k}$  denotes the  $k$ -th coordinate of  $\mathbf{x}_0$ . Hence we obviously have

$$Q_1 = \frac{1}{p} \sum_{k=1}^p X_k^1, \quad \mathbf{Q}_2 = \frac{1}{p} \sum_{k=1}^p \mathbf{X}_k^2. \quad (\text{D.5})$$

Next we sketch the main technical pieces for establishing the recovery results for  $\mathbf{Y}'$  first. All detailed proofs are deferred to later sections of the appendix. We will assume  $\exp(n/2)/2 \geq p \geq Cn^4 \log n$  for some large constant  $C$  for all the subsequent claims.

1. **Good initialization.** Proposition E.1 in Appendix E shows that with overwhelming probability, at least one of our  $p$  initialization vectors suggested in Section 3, say  $\mathbf{q}_i^{(0)} = \mathbf{y}^i$ , obeys that

$$\left| \left\langle \frac{\mathbf{y}^i}{\|\mathbf{y}^i\|_2}, \mathbf{e}_1 \right\rangle \right| \geq \frac{1}{4\sqrt{\theta n}}. \quad (\text{D.6})$$

2. **Uniform progress away from the equator.** By Proposition F.1 in Appendix F, there exists some constant  $\theta_0 > 0$ , such that for any  $\theta \in \left(\frac{1}{\sqrt{n}}, \theta_0\right)$ ,

$$G'(\mathbf{q}) = \frac{|Q'_1(\mathbf{q})|}{|q_1|} - \frac{\|\mathbf{Q}'_2(\mathbf{q})\|_2}{\|\mathbf{q}\|_2} \geq \frac{1}{4000\theta^2 np} \quad (\text{D.7})$$

holds uniformly for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  in the region  $\frac{1}{4\sqrt{\theta n}} \leq |q_1| \leq 3\sqrt{\theta}$  with overwhelming probability.

3. **No jumps away from the cap.** Proposition G.1 in Appendix G shows that for any  $\theta \in \left(\frac{1}{\sqrt{n}}, \theta_0\right)$ , with overwhelming probability,

$$\frac{Q'_1(\mathbf{q})}{\|\mathbf{Q}'(\mathbf{q})\|_2} \geq 2\sqrt{\theta} \quad (\text{D.8})$$

holds for all  $\mathbf{q}$  with  $|q_1| \geq 3\sqrt{\theta}$ .

4. **Location of the stationary/stopping point.** The first point above ensures that with overwhelming probability at least one starting point  $\mathbf{q}^{(0)}$  will satisfy  $|q_1^{(0)}| \geq \frac{1}{4\sqrt{\theta n}}$ . As shown in Appendix H, the strictly positive gap of the second point ensures that one needs to run at most  $O(n^4 \log n)$  iterations to first encounter an iterate  $\mathbf{q}^{(k)}$  such that  $|q_1^{(k)}| \geq 3\sqrt{\theta}$ . The third point suggests extra iterations will not move away from the cap area, and hence the stationary point  $\bar{\mathbf{q}}$  of the ADM algorithm will satisfy  $|\bar{q}_1| \geq 2\sqrt{\theta}$ . If one enforces a hard stop after  $O(n^4 \log n)$  iterations, the stopping point will similarly stay in the region  $|q_1| \geq 2\sqrt{\theta}$ .

5. **LP Rounding succeeds.** We know that in the LP rounding stage, described in Section 3, will receive a vector  $\mathbf{r} = \bar{\mathbf{q}}$  with its first coordinate  $|r_1| \geq 2\sqrt{\theta}$ . Proposition I.1 in Appendix I proves that with overwhelming probability, the LP rounding (3.6) (operated on  $\mathbf{Y}'$ ) will output a solution  $\mathbf{q}^* = \mathbf{e}_1$ .

In summary, our ADM algorithm in Algorithm 1 using a smart initialization, plus an LP rounding stage (3.6), will output  $\mathbf{q}^* = \pm \mathbf{e}_1$  with overwhelming probability, or  $\mathbf{Y}'\mathbf{q}^*$  as a nontrivial scaled version of  $\mathbf{x}_0$ .

For the general case when the input is an arbitrary orthonormal basis  $\hat{\mathbf{Y}} = \mathbf{Y}'\mathbf{R}$  for a certain orthogonal matrix  $\mathbf{R}$ , the target solution is  $\mathbf{R}^\top \mathbf{e}_1$ . The following technical pieces are perfectly parallel to the above for  $\mathbf{Y}'$ .

- Discussion at the end of Appendix E suggests with overwhelming probability, at least one row of  $\hat{\mathbf{Y}}$  provides an initial point  $\mathbf{q}^{(0)}$  such that  $|\langle \mathbf{q}^{(0)}, \mathbf{R}^\top \mathbf{e}_1 \rangle| \geq \frac{1}{4\sqrt{\theta n}}$ .
- Discussion following Proposition F.1 in Appendix F suggests that for all  $\mathbf{q}$  such that  $\frac{1}{4\sqrt{\theta n}} \leq |\langle \mathbf{q}, \mathbf{R}^\top \mathbf{e}_1 \rangle| \leq 3\sqrt{\theta}$ , there is a strictly positive gap, indicating steady progress towards a point  $\mathbf{q}^{(k)}$  such that  $|\langle \mathbf{q}^{(k)}, \mathbf{R}^\top \mathbf{e}_1 \rangle| \geq 3\sqrt{\theta}$ .
- Discussion at the end of Appendix G indicates once  $\mathbf{q}$  satisfying  $|\langle \mathbf{q}, \mathbf{R}^\top \mathbf{e}_1 \rangle|$ , the next iterate will not move far away from the target:

$$\frac{\langle \mathbf{Q}'(\mathbf{q}; \hat{\mathbf{Y}}), \mathbf{R}^\top \mathbf{e}_1 \rangle}{\|\mathbf{Q}'(\mathbf{q}; \hat{\mathbf{Y}})\|_2} \geq 2\sqrt{\theta}. \quad (\text{D.9})$$

- Repeating the argument in Appendix H for general input  $\widehat{\mathbf{Y}}$  shows it is enough to run the ADM algorithm  $O(n^4 \log n)$  iterations to cross the range  $\frac{1}{4\sqrt{\theta n}} \leq |\langle \mathbf{q}, \mathbf{R}^\top \mathbf{e}_1 \rangle| \leq 3\sqrt{\theta}$ . So the above three points together dictates that with the proposed initialization, with overwhelming probability, we finally obtain a point  $\bar{\mathbf{q}}$  that satisfies  $|\langle \bar{\mathbf{q}}, \mathbf{R}^\top \mathbf{e}_1 \rangle| \geq 2\sqrt{\theta}$ , if we run at least  $O(n^4 \log n)$  iterations.
- Since the ADM returns  $\bar{\mathbf{q}}$  satisfying  $|\langle \bar{\mathbf{q}}, \mathbf{R}^\top \mathbf{e}_1 \rangle| \geq 2\sqrt{\theta}$ , discussion at the end of Appendix I dictates that we will obtain  $\mathbf{q}^* = \mathbf{R}^\top \mathbf{e}_1$  as the optimizer of the rounding program, exactly the target solution.

We complete the proof.

## E Good Initialization

**Proposition E.1.** *Let  $\mathbf{y}^k$  for  $k = 1, \dots, p$  be the transpose of the rows of the orthonormal bases  $\mathbf{Y}'$  defined in (B.3). If  $\theta > 1/\sqrt{n}$  and  $\exp(n/2)/2 \geq p \geq Cn^2$  for some constant  $C > 0$ , it holds that at least one of our  $p$  initialization vectors suggested in Section 3, say  $\mathbf{q}_i^{(0)} = \mathbf{y}^i$ , obeys*

$$\left| \left\langle \frac{\mathbf{y}^i}{\|\mathbf{y}^i\|_2}, \mathbf{e}_1 \right\rangle \right| \geq \frac{1}{4\sqrt{\theta n}}, \quad (\text{E.1})$$

with probability at least  $1 - c' \exp(-c''n)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* Since  $\mathbf{x}_0$  is i.i.d. Bernoulli, with probability at least  $1 - (1 - \theta)^p \geq 1 - \exp(-\theta p)$ , at least one component of  $\mathbf{x}_0$  is nonzero. Without loss of generality (w.l.o.g.), assume the  $k$ -th component of  $\mathbf{x}_0$  is nonzero. Then  $x_{0k} = \frac{1}{\sqrt{\theta p}}$ , and

$$|q_{i1}| = \frac{\frac{1}{\sqrt{\theta p}}}{\|\mathbf{x}_0\|_2 \|\mathbf{y}^i\|_2} \geq \frac{\frac{1}{\sqrt{\theta p}}}{\|\mathbf{x}_0\|_2 (\|\mathbf{x}_0\|_2 \|\mathbf{x}_0\|_2 \|\ell^2 \rightarrow \ell^\infty + \|\mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty})} = \frac{1}{\sqrt{\theta p} \|\mathbf{x}_0\|_\infty + \|\mathbf{x}_0\|_2 \|\mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty}}. \quad (\text{E.2})$$

We know that with probability at least  $1 - \exp(-p\theta^2/2)$ , it holds that

$$\|\mathbf{x}_0\|_2 = \sqrt{|\mathcal{I}| \times \frac{1}{\theta p}} \leq \sqrt{2\theta p \times \frac{1}{\theta p}} = \sqrt{2}. \quad (\text{E.3})$$

Moreover, using Lemma B.4, and Lemma A.15 with  $\varepsilon = 1/16$  and  $\xi = 1/2$ , we know when  $p \geq C_1 n$  for some large  $C_1 > 0$ , it holds that (note that  $\|\mathbf{M}\|$  can be arbitrarily close to 1 for large  $C_1$  in Lemma B.2)

$$\|\mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} \leq \|\mathbf{G}\|_{\ell^2 \rightarrow \ell^\infty} \|\mathbf{M}\| + \frac{4\sqrt{2n}}{\sqrt{\theta p}} \leq \frac{9}{5} \sqrt{\frac{n}{p}} + \frac{4\sqrt{2n}}{\sqrt{\theta p}} \leq 2\sqrt{\frac{n}{p}} \quad (\text{E.4})$$

with probability at least  $1 - c_1 \exp(-c_2 n)$  for some positive constants  $c_1$  and  $c_2$ . Therefore with probability at least  $1 - \exp(-\theta p) - \exp(-\theta^2 p) - c_1 \exp(-c_2 n)$ , it holds that

$$|q_{i1}| \geq \frac{1}{1 + \sqrt{2\theta p} \times 2\sqrt{\frac{n}{p}}} = \frac{1}{1 + 2\sqrt{2}\sqrt{\theta n}}. \quad (\text{E.5})$$

Using the fact that  $\theta \geq 1/\sqrt{n}$ , we obtain  $|q_{i1}| \geq \frac{1}{(1+2\sqrt{2})\sqrt{\theta n}}$ . It is sufficient to set  $p \geq C_2 n^2$  for some large enough  $C_2 > 0$  to make the probability overwhelming, as desired.  $\square$

We will next show that for an arbitrary orthonormal basis  $\widehat{\mathbf{Y}} \doteq \mathbf{Y}'\mathbf{R}$  the initialization still biases towards the target solution. To see this, suppose w.l.o.g.  $(\mathbf{y}^i)^\top$  is a row of  $\mathbf{Y}'$  with nonzero first coordinate. We have

shown above that with high probability  $\left| \left\langle \frac{\mathbf{y}^i}{\|\mathbf{y}^i\|_2}, \mathbf{e}_1 \right\rangle \right| \geq \frac{1}{4\sqrt{\theta n}}$  if  $\mathbf{Y}'$  is the input orthonormal basis. For  $\mathbf{Y}'$ , as  $\mathbf{x}_0 = \mathbf{Y}'\mathbf{e}_1 = \mathbf{Y}'\mathbf{R}\mathbf{R}^\top\mathbf{e}_1$ , we know  $\mathbf{q}^* = \mathbf{R}^\top\mathbf{e}_1$  is the target solution corresponding to  $\widehat{\mathbf{Y}}$ . Observing that

$$\left| \left\langle \mathbf{R}^\top\mathbf{e}_1, \frac{(\mathbf{e}_i^\top\widehat{\mathbf{Y}})^\top}{\|(\mathbf{e}_i^\top\widehat{\mathbf{Y}})^\top\|_2} \right\rangle \right| = \left| \left\langle \mathbf{R}^\top\mathbf{e}_1, \frac{\mathbf{R}^\top(\mathbf{Y}')^\top\mathbf{e}_i}{\|\mathbf{R}^\top(\mathbf{Y}')^\top\mathbf{e}_i\|_2} \right\rangle \right| = \left| \left\langle \mathbf{e}_1, \frac{(\mathbf{Y}')^\top\mathbf{e}_i}{\|(\mathbf{Y}')^\top\mathbf{e}_i\|_2} \right\rangle \right| = \left| \left\langle \mathbf{e}_1, \frac{\mathbf{y}^i}{\|\mathbf{y}^i\|_2} \right\rangle \right| \geq \frac{1}{4\sqrt{n\theta}}, \quad (\text{E.6})$$

corroborating our claim.

## F Lower Bounding Finite Sample Gap $G'(\mathbf{q})$

We will first work with the ‘‘canonical’’ orthonormal basis  $\mathbf{Y}'$ . The task is to lower bound the gap for finite samples

$$G'(\mathbf{q}) = \frac{|Q'_1(\mathbf{q})|}{|q_1|} - \frac{\|\mathbf{Q}'_2(\mathbf{q})\|_2}{\|\mathbf{q}_2\|_2}. \quad (\text{F.1})$$

Since we can deterministically constrain  $|q_1|$  and  $\|\mathbf{q}_2\|_2$  (e.g.,  $\frac{1}{4\sqrt{n\theta}} \leq |q_1| \leq 3\sqrt{\theta}$  and  $\|\mathbf{q}_2\|_2 \geq \frac{1}{4}$ , where the choice of  $\frac{1}{4}$  is arbitrary here, as we can always take a sufficiently small  $\theta$ ), the challenge lies in lower bounding  $|Q'_1(\mathbf{q})|$  and upper bounding  $\|\mathbf{Q}'_2(\mathbf{q})\|_2$ , which depends on the orthonormal basis  $\mathbf{Y}'$ . It turns out to cook up a typical expectation-concentration style argument, the unnormalized basis  $\mathbf{Y}$  is much easier to work with than  $\mathbf{Y}'$ . Hence our proof will follow the observation that

$$|Q'_1(\mathbf{q})| \geq |\mathbb{E}[Q_1(\mathbf{q})]| - |Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| - |Q'_1(\mathbf{q}) - Q_1(\mathbf{q})|, \quad (\text{F.2})$$

$$\|\mathbf{Q}'_2(\mathbf{q})\|_2 \leq \|\mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 + \|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 + \|\mathbf{Q}'_2(\mathbf{q}) - \mathbf{Q}_2(\mathbf{q})\|_2. \quad (\text{F.3})$$

In particular, define the set  $\Gamma = \left\{ \mathbf{q} \in \mathbb{S}^{n-1} : \frac{1}{4\sqrt{n\theta}} \leq |q_1| \leq 3\sqrt{\theta}, \|\mathbf{q}_2\|_2 \geq \frac{1}{4} \right\}$ :

- Appendix F.1 shows that the expected gap is lower bounded for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  with  $|q_1| \leq 3\sqrt{\theta}$ :

$$G(\mathbf{q}) \doteq \frac{|\mathbb{E}[Q_1(\mathbf{q})]|}{|q_1|} - \frac{\|\mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2}{\|\mathbf{q}_2\|_2} \geq \frac{1}{50} \frac{q_1^2}{\theta p}. \quad (\text{F.4})$$

As  $|q_1| \geq \frac{1}{4\sqrt{n\theta}}$ , we have

$$\inf_{\mathbf{q} \in \Gamma} \frac{|\mathbb{E}[Q_1(\mathbf{q})]|}{|q_1|} - \frac{\|\mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2}{\|\mathbf{q}_2\|_2} \geq \frac{1}{800} \frac{1}{\theta^2 np}. \quad (\text{F.5})$$

- Appendix F.2, as summarized in Proposition F.9, shows that whenever  $\exp(n) \geq p \geq \Omega(n^4 \log n)$ , it holds with overwhelmingly probability that

$$\begin{aligned} & \sup_{\mathbf{q} \in \Gamma} \frac{|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]|}{|q_1|} + \frac{\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2}{\|\mathbf{q}_2\|_2} \\ & \leq \frac{4\sqrt{\theta n}}{16000\theta^{5/2}n^{3/2}p} + \frac{4}{16000\theta^2 np} = \frac{1}{2000\theta^2 np}. \end{aligned} \quad (\text{F.6})$$



- Appendix F.4 shows that whenever  $\exp(n/2)/2 \geq p \geq \Omega(n^4 \log n)$ , it holds with overwhelmingly probability that

$$\begin{aligned} & \sup_{\mathbf{q} \in \Gamma} \frac{|Q_1(\mathbf{q}) - Q'_1(\mathbf{q})|}{|q_1|} + \frac{\|\mathbf{Q}_2(\mathbf{q}) - \mathbf{Q}'_2(\mathbf{q})\|_2}{\|\mathbf{q}_2\|_2} \\ & \leq \frac{4\sqrt{\theta n}}{16000\theta^{5/2}n^{3/2}p} + \frac{4}{16000\theta^2np} = \frac{1}{2000\theta^2np}. \end{aligned} \quad (\text{F.7})$$

Observing that

$$\begin{aligned} \inf_{\mathbf{q} \in \Gamma} G'(\mathbf{q}) & \geq \inf_{\mathbf{q} \in \Gamma} \left( \frac{|\mathbb{E}[Q_1(\mathbf{q})]|}{|q_1|} - \frac{\|\mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2}{\|\mathbf{q}_2\|_2} \right) - \sup_{\mathbf{q} \in \Gamma} \left( \frac{|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]|}{|q_1|} + \frac{\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2}{\|\mathbf{q}_2\|_2} \right) \\ & \quad - \sup_{\mathbf{q} \in \Gamma} \left( \frac{|Q_1(\mathbf{q}) - Q'_1(\mathbf{q})|}{|q_1|} + \frac{\|\mathbf{Q}_2(\mathbf{q}) - \mathbf{Q}'_2(\mathbf{q})\|_2}{\|\mathbf{q}_2\|_2} \right), \end{aligned} \quad (\text{F.8})$$

we obtain the following:

**Proposition F.1.** *There exists some constant  $\theta_0 > 0$  such that, for all  $\theta \in (\frac{1}{\sqrt{n}}, \theta_0)$ , when  $\exp(n/2)/2 \geq p \geq Cn^4 \log n$  for some large constant  $C > 0$ , we have*

$$\inf_{\mathbf{q} \in \Gamma} G'(\mathbf{q}) \geq \frac{1}{4000\theta^2np}, \quad (\text{F.9})$$

with probability at least  $1 - c' \exp(-c''n)$  for some positive constants  $c'$  and  $c''$ .

For the general case when the input orthonormal basis is  $\widehat{\mathbf{Y}} = \mathbf{Y}'\mathbf{R}$  with target solution  $\mathbf{q}^* = \mathbf{R}^\top \mathbf{e}_1$ , a straightforward extension of the definition for the gap would be:

$$G'(\mathbf{q}; \widehat{\mathbf{Y}} = \mathbf{Y}'\mathbf{R}) \doteq \frac{|\langle \mathbf{Q}'(\mathbf{q}; \widehat{\mathbf{Y}}), \mathbf{R}^\top \mathbf{e}_1 \rangle|}{|\langle \mathbf{q}, \mathbf{R}^\top \mathbf{e}_1 \rangle|} - \frac{\|(\mathbf{I} - \mathbf{R}^\top \mathbf{e}_1 \mathbf{e}_1^\top \mathbf{R}) \mathbf{Q}'(\mathbf{q}; \widehat{\mathbf{Y}})\|_2}{\|(\mathbf{I} - \mathbf{R}^\top \mathbf{e}_1 \mathbf{e}_1^\top \mathbf{R}) \mathbf{q}\|_2}. \quad (\text{F.10})$$

Since  $\mathbf{Q}'(\mathbf{q}; \widehat{\mathbf{Y}}) = \frac{1}{p} \sum_{k=1}^p \mathbf{R}^\top \mathbf{y}^k S_\lambda(\mathbf{q}^\top \mathbf{R}^\top \mathbf{y}^k)$ , we have

$$\mathbf{R} \mathbf{Q}'(\mathbf{q}; \widehat{\mathbf{Y}}) = \frac{1}{p} \sum_{k=1}^p \mathbf{R} \mathbf{R}^\top \mathbf{y}'^k S_\lambda(\mathbf{q}^\top \mathbf{R}^\top \mathbf{y}'^k) = \frac{1}{p} \sum_{k=1}^p \mathbf{y}'^k S_\lambda[(\mathbf{R} \mathbf{q})^\top \mathbf{y}'^k] = \mathbf{Q}'(\mathbf{R} \mathbf{q}; \mathbf{Y}'). \quad (\text{F.11})$$

Hence we have

$$G'(\mathbf{q}; \widehat{\mathbf{Y}} = \mathbf{Y}'\mathbf{R}) = \frac{|\langle \mathbf{Q}'(\mathbf{R} \mathbf{q}; \mathbf{Y}'), \mathbf{e}_1 \rangle|}{|\langle \mathbf{R} \mathbf{q}, \mathbf{e}_1 \rangle|} - \frac{\|(\mathbf{I} - \mathbf{e}_1 \mathbf{e}_1^\top) \mathbf{Q}'(\mathbf{R} \mathbf{q}; \mathbf{Y}')\|_2}{\|(\mathbf{I} - \mathbf{e}_1 \mathbf{e}_1^\top) \mathbf{R} \mathbf{q}\|_2}. \quad (\text{F.12})$$

Therefore, from Proposition F.1 above, we conclude that under the same technical conditions as therein,

$$\inf_{\mathbf{q} \in \mathbb{S}^{n-1}: \frac{1}{4\sqrt{\theta n}} \leq |\langle \mathbf{R} \mathbf{q}, \mathbf{e}_1 \rangle| \leq 3\sqrt{\theta}} G'(\mathbf{q}; \widehat{\mathbf{Y}}) \geq \frac{1}{4000\theta^2np} \quad (\text{F.13})$$

with overwhelmingly probability.

## F.1 Lower Bounding the Expected Gap $G(\mathbf{q})$

In this section, we provide a nontrivial lower bound for the gap

$$G(\mathbf{q}) = \frac{|\mathbb{E}[Q_1(\mathbf{q})]|}{|q_1|} - \frac{\|\mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2}{\|\mathbf{q}_2\|_2}. \quad (\text{F.14})$$

More specifically, we show that:

**Proposition F.2.** *There exists some numerical constant  $\theta_0 > 0$ , such that for all  $\theta \in (0, \theta_0)$ , it holds that*

$$G(\mathbf{q}) \geq \frac{1}{50} \frac{q_1^2}{\theta p} \quad (\text{F.15})$$

for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  with  $|q_1| \leq 3\sqrt{\theta}$ .

Because the estimation for the gap  $G(\mathbf{q})$  involves dedicated estimation for  $\mathbb{E}[Q_1(\mathbf{q})]$  and  $\mathbb{E}[\mathbf{Q}_2(\mathbf{q})]$ , we sketch the main proof in Appendix F.1.1, and leave those detailed technical calculations in the subsequent subsections.

### F.1.1 Sketch of the Proof

W.l.o.g., we only consider the situation that  $q_1 > 0$ , because the case of  $q_1 < 0$  can be similarly shown by symmetry. By (D.5), (D.3) and (D.4), we have

$$\mathbb{E}[Q_1(\mathbf{q})] = \mathbb{E}[x_0 S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}]], \quad (\text{F.16})$$

$$\mathbb{E}[\mathbf{Q}_2(\mathbf{q})] = \mathbb{E}[\mathbf{g} S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}]], \quad (\text{F.17})$$

where  $\mathbf{q} = [q_1, \mathbf{q}_2^\top]^\top$ ,  $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \frac{1}{p} \mathbf{I})$ , and  $x_0 \sim \frac{1}{\sqrt{\theta p}} \text{Ber}(\theta)$ . Let us decompose

$$\mathbf{g} = \mathbf{g}_\parallel + \mathbf{g}_\perp, \quad (\text{F.18})$$

with  $\mathbf{g}_\parallel = \mathcal{P}_\parallel \mathbf{g} = \frac{\mathbf{q}_2 \mathbf{q}_2^\top}{\|\mathbf{q}_2\|_2^2} \mathbf{g}$ , and  $\mathbf{g}_\perp = (\mathbf{I} - \mathcal{P}_\parallel) \mathbf{g}$ . Therefore, we have

$$\begin{aligned} \mathbb{E}[\mathbf{Q}_2(\mathbf{q})] &= \mathbb{E}[\mathbf{g}_\parallel S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}_\parallel]] + \mathbb{E}[\mathbf{g}_\perp S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}_\parallel]] \\ &= \mathbb{E}[\mathbf{g}_\parallel S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}]] + \mathbb{E}[\mathbf{g}_\perp] \mathbb{E}[S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}]] \\ &= \frac{\mathbf{q}_2}{\|\mathbf{q}_2\|_2^2} \mathbb{E}[\mathbf{q}_2^\top \mathbf{g} S_\lambda [x_0 q_1 + \mathbf{q}_2^\top \mathbf{g}]], \end{aligned} \quad (\text{F.19})$$

where we used the facts that  $\mathbf{q}_2^\top \mathbf{g} = \mathbf{q}_2^\top \mathbf{g}_\parallel$ ,  $\mathbf{g}_\perp$  and  $\mathbf{g}_\parallel$  are uncorrelated Gaussian vectors and therefore independent, and  $\mathbb{E}[\mathbf{g}_\perp] = \mathbf{0}$ . Let  $Z \doteq \mathbf{g}^\top \mathbf{q}_2 \sim \mathcal{N}(0, \sigma^2)$  with  $\sigma^2 = \|\mathbf{q}_2\|_2^2 / p$ , by partial evaluation of the expectations with respect to  $x_0$ , we get

$$\mathbb{E}[Q_1(\mathbf{q})] = \sqrt{\frac{\theta}{p}} \mathbb{E}\left[S_\lambda \left[\frac{q_1}{\sqrt{\theta p}} + Z\right]\right], \quad (\text{F.20})$$

$$\mathbb{E}[\mathbf{Q}_2(\mathbf{q})] = \frac{\theta \mathbf{q}_2}{\|\mathbf{q}_2\|_2^2} \mathbb{E}\left[Z S_\lambda \left[\frac{q_1}{\sqrt{\theta p}} + Z\right]\right] + \frac{(1-\theta) \mathbf{q}_2}{\|\mathbf{q}_2\|_2^2} \mathbb{E}[Z S_\lambda [Z]]. \quad (\text{F.21})$$

Straightforward integration based on Lemma A.1 gives a explicit form of the expectations as follows

$$\mathbb{E}[Q_1(\mathbf{q})] = \sqrt{\frac{\theta}{p}} \left\{ \left[ \alpha \Psi\left(-\frac{\alpha}{\sigma}\right) + \beta \Psi\left(\frac{\beta}{\sigma}\right) \right] + \sigma \left[ \psi\left(-\frac{\beta}{\sigma}\right) - \psi\left(-\frac{\alpha}{\sigma}\right) \right] \right\}, \quad (\text{F.22})$$

$$\mathbb{E}[\mathbf{Q}_2(\mathbf{q})] = \left\{ \frac{2(1-\theta)}{p} \Psi\left(-\frac{\lambda}{\sigma}\right) + \frac{\theta}{p} \left[ \Psi\left(-\frac{\alpha}{\sigma}\right) + \Psi\left(\frac{\beta}{\sigma}\right) \right] \right\} \mathbf{q}_2, \quad (\text{F.23})$$

where the scalars  $\alpha$  and  $\beta$  are defined as

$$\alpha = \frac{q_1}{\sqrt{\theta p}} + \lambda, \quad \beta = \frac{q_1}{\sqrt{\theta p}} - \lambda, \quad (\text{F.24})$$

and  $\psi(t)$  and  $\Psi(t)$  are *pdf* and *cdf* for standard normal distribution, respectively, as defined in Lemma A.1. Plugging (F.22) and (F.23) into (F.14), by some simplifications, we obtain

$$\begin{aligned} G(\mathbf{q}) &= \frac{1}{q_1} \sqrt{\frac{\theta}{p}} \left[ \alpha \Psi\left(-\frac{\alpha}{\sigma}\right) + \beta \Psi\left(\frac{\beta}{\sigma}\right) - \frac{2q_1}{\sqrt{\theta p}} \Psi\left(-\frac{\lambda}{\sigma}\right) \right] - \frac{\theta}{p} \left[ \Psi\left(-\frac{\alpha}{\sigma}\right) + \Psi\left(\frac{\beta}{\sigma}\right) - 2\Psi\left(-\frac{\lambda}{\sigma}\right) \right] \\ &\quad + \frac{\sigma}{q_1} \sqrt{\frac{\theta}{p}} \left[ \psi\left(\frac{\beta}{\sigma}\right) - \psi\left(-\frac{\alpha}{\sigma}\right) \right]. \end{aligned} \quad (\text{F.25})$$

With  $\lambda = 1/\sqrt{p}$  and  $\sigma^2 = \|\mathbf{q}_2\|_2^2/p = (1 - q_1^2)/p$ , we have

$$-\frac{\alpha}{\sigma} = -\frac{\delta + 1}{\sqrt{1 - q_1^2}}, \quad \frac{\beta}{\sigma} = \frac{\delta - 1}{\sqrt{1 - q_1^2}}, \quad \frac{\lambda}{\sigma} = \frac{1}{\sqrt{1 - q_1^2}}, \quad (\text{F.26})$$

where  $\delta = q_1/\sqrt{\theta}$  for  $q_1 \leq 3\sqrt{\theta}$ . To proceed, it is natural to consider estimating the gap  $G(\mathbf{q})$  by Taylor's expansion. More specifically, we approximate  $\Psi(-\frac{\alpha}{\sigma})$  and  $\psi(-\frac{\alpha}{\sigma})$  around  $-1 - \delta$ , and approximate  $\Psi(\frac{\beta}{\sigma})$  and  $\psi(\frac{\beta}{\sigma})$  around  $-1 + \delta$ . Applying the estimates for the relevant quantities established in Lemma F.3, we obtain

$$\begin{aligned} G(\mathbf{q}) &\geq \frac{1 - \theta}{p} \Phi_1(\delta) - \frac{1}{\delta p} \Phi_2(\delta) + \frac{1 - \theta}{p} \psi(-1) q_1^2 + \frac{1}{p} \left( \sigma\sqrt{p} + \frac{\theta}{2} - 1 \right) \eta_2(\delta) q_1^2 \\ &\quad + \frac{1}{2\delta p} [1 + \delta^2 - \theta\delta^2 - \sigma(1 + \delta^2)\sqrt{p}] q_1^2 \eta_1(\delta) + \frac{\sigma}{\delta\sqrt{p}} \eta_1(\delta) - \frac{5C_T\sqrt{\theta}q_1^3}{p} (\delta + 1)^3, \end{aligned} \quad (\text{F.27})$$

where we define

$$\Phi_1(\delta) = \Psi(-1 - \delta) + \Psi(-1 + \delta) - 2\Psi(-1), \quad \Phi_2(\delta) = \Psi(-1 + \delta) - \Psi(-1 - \delta), \quad (\text{F.28})$$

$$\eta_1(\delta) = \psi(-1 + \delta) - \psi(-1 - \delta), \quad \eta_2(\delta) = \psi(-1 + \delta) + \psi(-1 - \delta), \quad (\text{F.29})$$

and  $C_T$  is as defined in Lemma F.3. Since  $1 - \sigma\sqrt{p} \geq 0$ , dropping those small positive terms  $\frac{q_1^2}{p}(1 - \theta)\psi(-1)$ ,  $\frac{\theta q_1^2}{2p}\eta_2(\delta)$ , and  $(1 + \delta^2)(1 - \sigma\sqrt{p})q_1^2\eta_1(\delta)/(2\delta p)$ , and using the fact that  $\delta = q_1/\sqrt{\theta}$ , we obtain

$$\begin{aligned} G(\mathbf{q}) &\geq \frac{1 - \theta}{p} \Phi_1(\delta) - \frac{1}{\delta p} [\Phi_2(\delta) - \sigma\sqrt{p}\eta_1(\delta)] - \frac{q_1^2}{p} (1 - \sigma\sqrt{p}) \eta_2(\delta) - \frac{\sqrt{\theta}}{2p} q_1^3 \eta_1(\delta) - \frac{C_1\sqrt{\theta}q_1^3}{p} \max\left(\frac{q_1^3}{\theta^{3/2}}, 1\right) \\ &\geq \frac{1 - \theta}{p} \Phi_1(\delta) - \frac{1}{\delta p} [\Phi_2(\delta) - \eta_1(\delta)] - \frac{q_1^2}{p} \frac{\eta_1(\delta)}{\delta} - \frac{q_1^2}{p\theta} \frac{2}{\sqrt{2\pi}} \theta - \frac{q_1^2}{\theta p} \frac{3\theta^2}{2\sqrt{2\pi}} - \frac{q_1^2}{\theta p} (C_1\theta^2), \end{aligned} \quad (\text{F.30})$$

for some constant  $C_1 > 0$ , where we have used  $q_1 \leq 3\sqrt{\theta}$  to simplify the bounds and the fact  $\sigma\sqrt{p} = \sqrt{1 - q_1^2} \geq 1 - q_1^2$  to simplify the expression. Substituting the estimates in Lemma F.5 and use the fact  $\delta \mapsto \eta_1(\delta)/\delta$  is bounded, we obtain

$$G(p) \geq \frac{1}{p} \left( \frac{1}{40} - \frac{1}{2\sqrt{2\pi}} \theta \right) \delta^2 - \frac{q_1^2}{\theta p} (c_1\theta + c_2\theta^2) \quad (\text{F.31})$$

$$\geq \frac{q_1^2}{\theta p} \left( \frac{1}{40} - \frac{1}{\sqrt{2\pi}} \theta - c_1\theta - c_2\theta^2 \right) \quad (\text{F.32})$$

for some positive constants  $c_1$  and  $c_2$ . We obtain the claimed result once  $\theta_0$  is made sufficiently small.

### F.1.2 Auxiliary Results Used in the Proof

**Lemma F.3.** Let  $\delta \doteq q_1/\sqrt{\theta}$ . There exists some universal constant  $C_T > 0$  such that we have the follow polynomial approximations hold for all  $q_1 \in (0, \frac{1}{2})$ :

$$\left| \psi\left(-\frac{\alpha}{\sigma}\right) - \left[1 - \frac{1}{2}(1+\delta)^2 q_1^2\right] \psi(-1-\delta) \right| \leq C_T (1+\delta)^2 q_1^4, \quad (\text{F.33})$$

$$\left| \psi\left(\frac{\beta}{\sigma}\right) - \left[1 - \frac{1}{2}(\delta-1)^2 q_1^2\right] \psi(\delta-1) \right| \leq C_T (\delta-1)^2 q_1^4, \quad (\text{F.34})$$

$$\left| \Psi\left(-\frac{\alpha}{\sigma}\right) - \left[\Psi(-1-\delta) - \frac{1}{2}\psi(-1-\delta)(1+\delta)q_1^2\right] \right| \leq C_T (1+\delta)^2 q_1^4, \quad (\text{F.35})$$

$$\left| \Psi\left(\frac{\beta}{\sigma}\right) - \left[\Psi(\delta-1) + \frac{1}{2}\psi(\delta-1)(\delta-1)q_1^2\right] \right| \leq C_T (\delta-1)^2 q_1^4, \quad (\text{F.36})$$

$$\left| \Psi\left(-\frac{\lambda}{\sigma}\right) - \left[\Psi(-1) - \frac{1}{2}\psi(-1)q_1^2\right] \right| \leq C_T q_1^4. \quad (\text{F.37})$$

*Proof.* First observe that for any  $q_1 \in (0, \frac{1}{2})$  it holds that

$$0 \leq \frac{1}{\sqrt{1-q_1^2}} - \left(1 + \frac{q_1^2}{2}\right) \leq q_1^4. \quad (\text{F.38})$$

Hence we have

$$-(1+\delta) \left(1 + \frac{1}{2}q_1^2 + q_1^4\right) \leq -\frac{\alpha}{\sigma} \leq -(1+\delta) \left(1 + \frac{1}{2}q_1^2\right), \quad (\text{F.39})$$

$$(\delta-1) \left(1 + \frac{1}{2}q_1^2\right) \leq \frac{\beta}{\sigma} \leq (\delta-1) \left(1 + \frac{1}{2}q_1^2 + q_1^4\right), \quad \text{when } \delta \geq 1$$

$$(\delta-1) \left(1 + \frac{1}{2}q_1^2 + q_1^4\right) \leq \frac{\beta}{\sigma} \leq (\delta-1) \left(1 + \frac{1}{2}q_1^2\right), \quad \text{when } \delta \leq 1. \quad (\text{F.40})$$

So we have

$$\psi\left(-\frac{\alpha}{\sigma}\right) \leq \psi\left(-\frac{\alpha}{\sigma}\right) \leq \psi\left(-\frac{\alpha}{\sigma}\right). \quad (\text{F.41})$$

By Taylor expansion of the left and right sides of the above two-side inequality around  $-1-\delta$  using Lemma A.2, we obtain

$$\left| \psi\left(-\frac{\alpha}{\sigma}\right) - \psi(-1-\delta) - \frac{1}{2}(1+\delta)^2 q_1^2 \psi(-1-\delta) \right| \leq C_T (1+\delta)^2 q_1^4, \quad (\text{F.42})$$

for some numerical constant  $C_T > 0$  sufficiently large. In the same way, we can obtain other claimed results.  $\square$

**Lemma F.4.** For any  $\delta \in [0, 3]$ , it holds that

$$\Phi_2(\delta) - \eta_1(\delta) \geq \frac{\eta_1(3)}{9} \delta^3 \geq \frac{1}{20} \delta^3. \quad (\text{F.43})$$

*Proof.* Let us define

$$h(\delta) = \Phi_2(\delta) - \eta_1(\delta) - C\delta^3 \quad (\text{F.44})$$

for some  $C > 0$  to be determined later. Then it is obvious that  $h(0) = 0$ . Direct calculation shows that

$$\frac{d}{d\delta}\Phi_1(\delta) = \eta_1(\delta), \quad \frac{d}{d\delta}\Phi_2(\delta) = \eta_2(\delta), \quad \frac{d}{d\delta}\eta_1(\delta) = \eta_2(\delta) - \delta\eta_1(\delta). \quad (\text{F.45})$$

Thus, to show (F.43), it is sufficient to show that  $h'(\delta) \geq 0$  for all  $\delta \in [0, 3]$ . By differentiating  $h(\delta)$  with respect to  $\delta$  and use the results in (F.45), it is sufficient to have

$$h'(\delta) = \delta\eta_1(\delta) - 3C\delta^2 \geq 0 \iff \eta_1(\delta) \geq 3C\delta \quad (\text{F.46})$$

for all  $\delta \in [0, 3]$ . We obtain the claimed result by observing that  $\delta \mapsto \eta_1(\delta)/3\delta$  is monotonically decreasing over  $\delta \in [0, 3]$  as justified below.

Consider the function

$$p(\delta) \doteq \frac{\eta_1(\delta)}{3\delta} = \frac{1}{3\sqrt{2\pi}} \exp\left(-\frac{\delta^2 + 1}{2}\right) \frac{e^\delta - e^{-\delta}}{\delta}. \quad (\text{F.47})$$

To show it is monotonically decreasing, it is enough to show  $p'(\delta)$  is always nonpositive for  $\delta \in (0, 3)$ , or equivalently

$$g(\delta) \doteq (e^\delta + e^{-\delta})\delta - (\delta^2 + 1)(e^\delta - e^{-\delta}) \leq 0 \quad (\text{F.48})$$

for all  $\delta \in (0, 3)$ , which can be easily verified by noticing that  $g(0) = 0$  and  $g'(\delta) \leq 0$  for all  $\delta \geq 0$ .  $\square$

**Lemma F.5.** For any  $\delta \in [0, 3]$ , we have

$$(1 - \theta)\Phi_1(\delta) - \frac{1}{\delta}[\Phi_2(\delta) - \eta_1(\delta)] \geq \left(\frac{1}{40} - \frac{1}{\sqrt{2\pi}}\theta\right)\delta^2. \quad (\text{F.49})$$

*Proof.* Let us define

$$g(\delta) = (1 - \theta)\Phi_1(\delta) - \frac{1}{\delta}[\Phi_2(\delta) - \eta_1(\delta)] - c_0(\theta)\delta^2, \quad (\text{F.50})$$

where  $c_0(\theta) > 0$  is a function of  $\theta$ . Thus, by the results in (F.45) and L'Hospital's rule, we have

$$\lim_{\delta \rightarrow 0} \frac{\Phi_2(\delta)}{\delta} = \lim_{\delta \rightarrow 0} \eta_2(\delta) = 2\psi(-1), \quad \lim_{\delta \rightarrow 0} \frac{\eta_1(\delta)}{\delta} = \lim_{\delta \rightarrow 0} [\eta_2(\delta) - \delta\eta_1(\delta)] = 2\psi(-1). \quad (\text{F.51})$$

Combined that with the fact that  $\Phi_1(0) = 0$ , we conclude  $g(0) = 0$ . Hence, to show (F.49), it is sufficient to show that  $g'(\delta) \geq 0$  for all  $\delta \in [0, 3]$ . Direct calculation using the results in (F.45) shows that

$$g'(\delta) = \frac{1}{\delta^2}[\Phi_2(\delta) - \eta_1(\delta)] - \theta\eta_1(\delta) - 2c_0(\theta)\delta. \quad (\text{F.52})$$

Since  $\eta_1(\delta)/\delta$  is monotonically decreasing as shown in Lemma F.4, we have that for all  $\delta \in (0, 3)$

$$\eta_1(\delta) \leq \delta \lim_{\delta \rightarrow 0} \frac{\eta_1(\delta)}{\delta} \leq \frac{2}{\sqrt{2\pi}}\delta. \quad (\text{F.53})$$

Using the above bound and the main result from Lemma F.4 again, we obtain

$$g'(\delta) \geq \frac{1}{20}\delta - \frac{2}{\sqrt{2\pi}}\theta\delta - 2c_0\delta. \quad (\text{F.54})$$

Choosing  $c_0(\theta) = \frac{1}{40} - \frac{1}{\sqrt{2\pi}}\theta$  completes the proof.  $\square$

## F.2 Finite Sample Concentration

In the following two subsections, we estimate the deviations around the expectations  $\mathbb{E}Q_1$  and  $\mathbb{E}\mathbf{Q}_2$ , i.e.,  $|Q_1 - \mathbb{E}Q_1|$  and  $\|\mathbf{Q}_2 - \mathbb{E}\mathbf{Q}_2\|_2$ , and show that the total deviations fit into the gap  $G(\mathbf{q})$  we derived in Appendix F.1. Our analysis is based on the scalar and vector Bernstein's inequalities with moment conditions. Finally, in Appendix F.3, we uniform the bound by applying the classical discretization argument.

### F.2.1 Concentration for $Q_1$

**Lemma F.6** (Bounding  $|Q_1 - \mathbb{E}[Q_1(\mathbf{q})]|$ ). *For each  $\mathbf{q} \in \mathbb{S}^{n-1}$ , it holds for all  $t > 0$  that*

$$\mathbb{P}[|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| \geq t] \leq 2 \exp\left(-\frac{\theta p^3 t^2}{8 + 4pt}\right). \quad (\text{F.55})$$

*Proof.* By (D.3) and (D.5), we know that

$$Q_1(\mathbf{q}) = \frac{1}{p} \sum_{k=1}^p X_k^1, \quad X_k^1 = x_{0k} \mathcal{S}_\lambda[x_{0k} q_1 + Z_k] \quad (\text{F.56})$$

where  $Z_k \sim \mathcal{N}\left(0, \frac{\|\mathbf{q}_2\|_2^2}{p}\right)$ . Thus, for any  $m \geq 2$ , by Lemma A.4, we have

$$\begin{aligned} \mathbb{E}\left[|X_k^1|^m\right] &\leq \theta \left(\frac{1}{\sqrt{\theta p}}\right)^m \mathbb{E}\left[\left|\frac{q_1}{\sqrt{\theta p}} + Z_k\right|^m\right] \\ &= \theta \left(\frac{1}{\sqrt{\theta p}}\right)^m \sum_{l=0}^m \binom{m}{l} \left(\frac{q_1}{\sqrt{\theta p}}\right)^l \mathbb{E}\left[|Z_k|^{m-l}\right] \\ &= \theta \left(\frac{1}{\sqrt{\theta p}}\right)^m \sum_{l=0}^m \binom{m}{l} \left(\frac{q_1}{\sqrt{\theta p}}\right)^l (m-l-1)!! \left(\frac{\|\mathbf{q}_2\|_2}{\sqrt{p}}\right)^{m-l} \\ &\leq \frac{m!}{2} \theta \left(\frac{1}{\sqrt{\theta p}}\right)^m \left(\frac{q_1}{\sqrt{\theta p}} + \frac{\|\mathbf{q}_2\|_2}{\sqrt{p}}\right)^m \\ &\leq \frac{m!}{2} \theta \left(\frac{2}{\theta p}\right)^m = \frac{m!}{2} \frac{4}{\theta p^2} \left(\frac{2}{\theta p}\right)^{m-2} \end{aligned} \quad (\text{F.57})$$

let  $\sigma_X^2 = 4/(\theta p^2)$  and  $R = 2/(\theta p)$ , apply Lemma A.7, we get

$$\mathbb{P}[|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| \geq t] \leq 2 \exp\left(-\frac{\theta p^3 t^2}{8 + 4pt}\right). \quad (\text{F.58})$$

as desired.  $\square$

### F.2.2 Concentration for $\mathbf{Q}_2$

**Lemma F.7** (Bounding  $\|\mathbf{Q}_2 - \mathbb{E}[\mathbf{Q}_2]\|_2$ ). *For each  $\mathbf{q} \in \mathbb{S}^{n-1}$ , it holds for all  $t > 0$  that*

$$\mathbb{P}[\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 > t] \leq 2(n+1) \exp\left(-\frac{\theta p^3 t^2}{128n + 16\sqrt{\theta n p t}}\right). \quad (\text{F.59})$$

Before proving Lemma F.7, we record the following useful results.

**Lemma F.8.** *For any positive integer  $s, l > 0$ , we have*

$$\mathbb{E}\left[\|\mathbf{g}^k\|_2^s |\mathbf{q}_2^\top \mathbf{g}^k|^l\right] \leq \frac{(l+s)!}{2} \|\mathbf{q}_2\|_2^l \frac{(2\sqrt{n})^s}{(\sqrt{p})^{s+l}} \quad (\text{F.60})$$

In particular, when  $s = l$ , we have

$$\mathbb{E} \left[ \|\mathbf{g}^k\|_2^l |\mathbf{q}_2^\top \mathbf{g}^k|^l \right] \leq \frac{l!}{2} \|\mathbf{q}_2\|_2^l \left( \frac{4\sqrt{n}}{p} \right)^l \quad (\text{F.61})$$

*Proof.* Let  $\mathcal{P}_{\mathbf{q}_2} = \frac{\mathbf{q}_2 \mathbf{q}_2^\top}{\|\mathbf{q}_2\|_2^2}$  and  $\mathcal{P}_{\mathbf{q}_2^\perp} = \left( \mathbf{I} - \frac{1}{\|\mathbf{q}_2\|_2^2} \mathbf{q}_2 \mathbf{q}_2^\top \right)$  denote the projection operators onto  $\mathbf{q}_2$  and its orthogonal complement, respectively. By Lemma A.4, we have

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{g}^k\|_2^s |\mathbf{q}_2^\top \mathbf{g}^k|^l \right] &\leq \mathbb{E} \left[ \left( \|\mathcal{P}_{\mathbf{q}_2} \mathbf{g}^k\|_2 + \|\mathcal{P}_{\mathbf{q}_2^\perp} \mathbf{g}^k\|_2 \right)^s |\mathbf{q}_2^\top \mathbf{g}^k|^l \right] \\ &= \sum_{i=0}^s \binom{s}{i} \mathbb{E} \left[ \|\mathcal{P}_{\mathbf{q}_2^\perp} \mathbf{g}^k\|_2^i \right] \mathbb{E} \left[ |\mathbf{q}_2^\top \mathbf{g}^k|^l \|\mathcal{P}_{\mathbf{q}_2} \mathbf{g}^k\|_2^{s-i} \right] \\ &= \sum_{i=0}^s \binom{s}{i} \mathbb{E} \left[ \|\mathcal{P}_{\mathbf{q}_2^\perp} \mathbf{g}^k\|_2^i \right] \mathbb{E} \left[ |\mathbf{q}_2^\top \mathbf{g}^k|^{l+s-i} \right] \frac{1}{\|\mathbf{q}_2\|_2^{s-i}} \\ &\leq \|\mathbf{q}_2\|_2^l \sum_{i=0}^s \binom{s}{i} \mathbb{E} \left[ \|\mathcal{P}_{\mathbf{q}_2^\perp} \mathbf{g}^k\|_2^i \right] \left( \frac{1}{\sqrt{p}} \right)^{l+s-i} (l+s-i-1)!! \end{aligned} \quad (\text{F.62})$$

Using Lemma A.5 and the fact that  $\|\mathcal{P}_{\mathbf{q}_2^\perp} \mathbf{g}^k\|_2^2 \leq \|\mathbf{g}^k\|_2^2$ , we obtain

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{g}^k\|_2^s |\mathbf{q}_2^\top \mathbf{g}^k|^l \right] &\leq \|\mathbf{q}_2\|_2^l \sum_{i=0}^s \binom{s}{i} \left( \frac{\sqrt{n}}{\sqrt{p}} \right)^i i! \left( \frac{1}{\sqrt{p}} \right)^{l+s-i} (l+s-i-1)!! \\ &\leq \|\mathbf{q}_2\|_2^l \left( \frac{1}{\sqrt{p}} \right)^l \frac{(l+s)!}{2} \left( \frac{\sqrt{n}}{\sqrt{p}} + \frac{1}{\sqrt{p}} \right)^s \\ &\leq \frac{(l+s)!}{2} \|\mathbf{q}_2\|_2^l \frac{(2\sqrt{n})^s}{(\sqrt{p})^{s+l}}. \end{aligned} \quad (\text{F.63})$$

□

Now, we are ready to prove Lemma F.7,

*Proof.* By (D.5) and (D.4), note that

$$\mathbf{Q}_2 = \frac{1}{p} \sum_{k=1}^p \mathbf{X}_k^2, \quad \mathbf{X}_k^2 = \mathbf{g}^k \mathcal{S}_\lambda [x_{0k} q_1 + Z_k] \quad (\text{F.64})$$

where  $Z_k = \mathbf{q}_2^\top \mathbf{g}^k$ . Thus, for any  $m \geq 2$ , by Lemma F.8, we have

$$\begin{aligned} \mathbb{E} \left[ \|\mathbf{X}_k^2\|_2^m \right] &\leq \theta \mathbb{E} \left[ \|\mathbf{g}^k\|_2^m \left| \frac{q_1}{\sqrt{\theta p}} + \mathbf{q}_2^\top \mathbf{g}^k \right|^m \right] + (1-\theta) \mathbb{E} \left[ \|\mathbf{g}^k\|_2^m |\mathbf{q}_2^\top \mathbf{g}^k|^m \right] \\ &\leq \theta \sum_{l=0}^m \binom{m}{l} \mathbb{E} \left[ |\mathbf{q}_2^\top \mathbf{g}^k|^l \|\mathbf{g}^k\|_2^m \right] \left| \frac{q_1}{\sqrt{\theta p}} \right|^{m-l} + (1-\theta) \mathbb{E} \left[ \|\mathbf{g}^k\|_2^m |\mathbf{q}_2^\top \mathbf{g}^k|^m \right] \\ &\leq \theta \left( \frac{2\sqrt{n}}{\sqrt{p}} \right)^m \sum_{l=0}^m \binom{m}{l} \frac{(m+l)!}{2} \left( \frac{\|\mathbf{q}_2\|_2}{\sqrt{p}} \right)^l \left| \frac{q_1}{\sqrt{\theta p}} \right|^{m-l} + (1-\theta) \frac{m!}{2} \|\mathbf{q}_2\|_2^m \left( \frac{4\sqrt{n}}{p} \right)^m \\ &\leq \theta \frac{m!}{2} \left( \frac{4\sqrt{n}}{\sqrt{p}} \right)^m \left( \frac{\|\mathbf{q}_2\|_2}{\sqrt{p}} + \frac{q_1}{\sqrt{\theta p}} \right)^m + (1-\theta) \frac{m!}{2} \|\mathbf{q}_2\|_2^m \left( \frac{4\sqrt{n}}{p} \right)^m \end{aligned} \quad (\text{F.65})$$

$$\leq \frac{m!}{2} \left( \frac{8\sqrt{n}}{\sqrt{\theta p}} \right)^m. \quad (\text{F.66})$$

Taking  $\sigma_X^2 = 64n/(\theta p^2)$  and  $R = 8\sqrt{n}/(\sqrt{\theta}p)$  and using vector Bernstein's inequality in Lemma A.8, we obtain

$$\mathbb{P}[\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 \geq t] \leq 2(n+1) \exp\left(-\frac{\theta p^3 t^2}{128n + 16\sqrt{\theta} n p t}\right), \quad (\text{F.67})$$

as desired.  $\square$

### F.3 Union Bound

**Proposition F.9** (Uniformizing the Bounds). *Suppose that  $\theta > \frac{1}{\sqrt{n}}$ . Given any  $\xi > 0$ , there exists some constant  $C(\xi)$ , such that whenever  $\exp(n) \geq p \geq C(\xi) n^4 \log n$ , we have*

$$|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| \leq \frac{2\xi}{\theta^{5/2} n^{3/2} p}, \quad (\text{F.68})$$

$$\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 \leq \frac{2\xi}{\theta^2 n p} \quad (\text{F.69})$$

hold uniformly for all  $\mathbf{q} \in \mathbb{S}^{n-1}$ , with probability at least  $1 - c' \exp(-c''n)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* We apply the standard covering argument. For any  $\varepsilon \in (0, 1)$ , by Lemma A.12, the unit hemisphere of interest can be covered by an  $\varepsilon$ -net  $\mathcal{N}_\varepsilon$  of cardinality at most  $(3/\varepsilon)^n$ . For any  $\mathbf{q} \in \mathbb{S}^{n-1}$ , it can be written as

$$\mathbf{q} = \mathbf{q}' + \mathbf{e} \quad (\text{F.70})$$

where  $\mathbf{q}' \in \mathcal{N}_\varepsilon$  and  $\|\mathbf{e}\|_2 \leq \varepsilon$ . Let  $\mathbf{y}^k = [x_{0k}, \mathbf{g}^k]^\top$  be a row of  $\mathbf{Y}$ , by (D.3) and (D.5), we have

$$\begin{aligned} & |Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| \\ &= \left| \frac{1}{p} \sum_{k=1}^p \{x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' + \mathbf{e} \rangle] - \mathbb{E}[x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' + \mathbf{e} \rangle]]\} \right| \\ &\leq \left| \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' + \mathbf{e} \rangle] - \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' \rangle] \right| + \left| \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' \rangle] - \mathbb{E}[x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' \rangle]] \right| \\ &\quad + |\mathbb{E}[x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' \rangle]] - \mathbb{E}[x_{0k} \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' + \mathbf{e} \rangle]]|. \end{aligned} \quad (\text{F.71})$$

Using Cauchy-Schwarz inequality and the fact that  $\mathcal{S}_\lambda[\cdot]$  is a nonexpansive operator, we have

$$|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| \leq |Q_1(\mathbf{q}') - \mathbb{E}[Q_1(\mathbf{q}')]| + \left( \frac{1}{p} \sum_{k=1}^p |x_{0k}| \|\mathbf{y}^k\|_2 + \mathbb{E}[|x_{0k}| \|\mathbf{y}^k\|_2] \right) \|\mathbf{e}\|_2 \quad (\text{F.72})$$

$$\leq |Q_1(\mathbf{q}') - \mathbb{E}[Q_1(\mathbf{q}')]| + 2\varepsilon \frac{1}{\sqrt{\theta} p} \left( \frac{1}{\sqrt{\theta} p} + \max_{k \in [p]} \|\mathbf{g}^k\|_2 \right). \quad (\text{F.73})$$

By Lemma A.11 and the assumption that  $p \leq \exp(n)$ , we have that  $\max_{k \in [p]} \|\mathbf{g}^k\|_2 \leq 4\sqrt{n/p}$  with probability at least  $1 - \exp(-n/2)$ . Taking  $t = \xi \theta^{-5/2} n^{-3/2} p^{-1}$  in Lemma F.6 and applying a union bound, setting  $\varepsilon = \xi \theta^{-2} n^{-2}/10$  and combining with the above estimate, we obtain that

$$|Q_1(\mathbf{q}) - \mathbb{E}[Q_1(\mathbf{q})]| \leq \frac{\xi}{\theta^{5/2} n^{3/2} p} + \frac{\xi}{5} \frac{1}{\theta^2 n^2} \frac{5\sqrt{n}}{\sqrt{\theta} p} \leq \frac{2\xi}{\theta^{5/2} n^{3/2} p} \quad (\text{F.74})$$

holds for all  $\mathbf{q} \in \mathbb{S}^{n-1}$ , with probability at least  $1 - \exp(-n/2) - \exp\left(-\frac{c_1(\xi)p}{\theta^4 n^3} + c_2(\xi) n \log n\right)$  for some numerical constants  $c_1(\xi)$  and  $c_2(\xi)$ .



Similarly, by (D.3) and (D.5), we have

$$\begin{aligned}
\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 &= \left\| \frac{1}{p} \sum_{k=1}^p \{\mathbf{g}^k \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' + \mathbf{e} \rangle] - \mathbb{E}[\mathbf{g}^k \mathcal{S}_\lambda[\langle \mathbf{y}^k, \mathbf{q}' + \mathbf{e} \rangle]]\} \right\|_2 \\
&\leq \|\mathbf{Q}_2(\mathbf{q}') - \mathbb{E}[\mathbf{Q}_2(\mathbf{q}')]\|_2 + \left( \frac{1}{p} \sum_{k=1}^p \|\mathbf{g}^k\|_2 \|\mathbf{y}^k\|_2 + \mathbb{E}[\|\mathbf{g}^k\|_2 \|\mathbf{y}^k\|_2] \right) \|\mathbf{e}\|_2 \\
&\leq \|\mathbf{Q}_2(\mathbf{q}') - \mathbb{E}[\mathbf{Q}_2(\mathbf{q}')]\|_2 + 2\varepsilon \max_{k \in [p]} \|\mathbf{g}^k\|_2 \left( \frac{1}{\sqrt{\theta p}} + \max_{k \in [p]} \|\mathbf{g}^k\|_2 \right). \tag{F.75}
\end{aligned}$$

Applying the above estimates for  $\max_{k \in [p]} \|\mathbf{g}^k\|_2$ , and taking  $t = \xi \theta^{-2} n^{-1} p^{-1}$  in Lemma F.7 and applying a union bound, then setting  $\varepsilon = \xi \theta^{-2} n^{-2} / 40$ , we obtain that

$$\|\mathbf{Q}_2(\mathbf{q}) - \mathbb{E}[\mathbf{Q}_2(\mathbf{q})]\|_2 \leq \frac{\xi}{\theta^2 n p} + \frac{\xi}{20 \theta^2 n^2} 4 \sqrt{\frac{n}{p}} \left( \frac{1}{\sqrt{\theta p}} + 4 \sqrt{\frac{n}{p}} \right) \leq \frac{2\xi}{\theta^2 n p} \tag{F.76}$$

holds for all  $\mathbf{q} \in \mathbb{S}^{n-1}$ , with probability at least  $1 - \exp(-n/2) - \exp\left(-\frac{c_3(\xi)p}{\theta^3 n^3} + c_4(\xi) n \log n\right)$ .

Overall, it is enough to take  $p \geq C n^4 \log n$  for some large  $C$  to make the above events to hold with overwhelming probability, as desired.  $\square$

#### F.4 $\mathbf{Q}'(\mathbf{q})$ approximates $\mathbf{Q}(\mathbf{q})$

**Proposition F.10.** *Suppose  $\theta > \frac{1}{\sqrt{n}}$ . For any  $\xi > 0$ , there exists some constant  $C(\xi)$ , such that whenever  $\exp(n/2)/2 \geq p \geq C(\xi) n^4 \log n$ , the following bounds*

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q'_1(\mathbf{q}) - Q_1(\mathbf{q})| \leq \frac{\xi}{\theta^{5/2} n^{3/2} p} \tag{F.77}$$

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\mathbf{Q}'_2(\mathbf{q}) - \mathbf{Q}_2(\mathbf{q})\|_2 \leq \frac{\xi}{\theta^2 n p}, \tag{F.78}$$

hold for all  $\mathbf{q} \in \mathbb{S}^{n-1}$ , with probability at least  $1 - c' \exp(-c'' n)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* First, for any  $\mathbf{q} \in \mathbb{S}^{n-1}$ , from D.5, we know that

$$\begin{aligned}
&|Q_1(\mathbf{q}) - Q'_1(\mathbf{q})| \\
&= \left| \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] - \frac{1}{p} \sum_{k=1}^p \frac{x_{0k}}{\|\mathbf{x}_0\|_2} \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k] \right| \\
&\leq \left| \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] - \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k] \right| + \left| \frac{1}{p} \sum_{k=1}^p x_{0k} \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k] - \frac{1}{p} \sum_{k=1}^p \frac{x_{0k}}{\|\mathbf{x}_0\|_2} \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k] \right| \\
&\leq \frac{1}{p} \sum_{k=1}^p |x_{0k}| |\mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] - \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k]| + \frac{1}{p} \sum_{k=1}^p |x_{0k}| \left| 1 - \frac{1}{\|\mathbf{x}_0\|_2} \right| |\mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k]|. \tag{F.79}
\end{aligned}$$

Let  $\mathcal{I} = \text{supp}(\mathbf{x}_0)$ . Conditioned on the support, using the facts that  $\mathcal{S}_\lambda[\cdot]$  is a nonexpansive operator, we obtain

$$\begin{aligned}
\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q_1(\mathbf{q}) - Q'_1(\mathbf{q})| &\leq \frac{1}{p} \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \sum_{k \in \mathcal{I}} |x_{0k}| \|\mathbf{q}^\top (\mathbf{y}^k - \mathbf{y}'^k)\|_2 + \left| 1 - \frac{1}{\|\mathbf{x}_0\|_2} \right| \frac{1}{p} \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \sum_{k \in \mathcal{I}} |x_{0k}| |\mathbf{q}^\top \mathbf{y}'^k| \\
&= \frac{1}{\sqrt{\theta} p^{3/2}} \left( \|\mathbf{Y}_{\mathcal{I}} - \mathbf{Y}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} + \left| 1 - \frac{1}{\|\mathbf{x}_0\|_2} \right| \|\mathbf{Y}'_{\mathcal{I}}\|_{\ell^2 \rightarrow \ell^1} \right). \tag{F.80}
\end{aligned}$$

By Lemma B.1 and Lemma B.3 in Appendix B, we have the following holds

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q_1(\mathbf{q}) - Q'_1(\mathbf{q})| \leq \frac{1}{\sqrt{\theta} p^{3/2}} \left( \frac{10}{\theta} \sqrt{n \log p} + \frac{2\sqrt{2}}{5} \sqrt{\frac{n \log p}{\theta^3 p}} \times 7\sqrt{2\theta p} \right) \leq \frac{16}{\theta^{3/2} p^{3/2}} \sqrt{n \log p}, \quad (\text{F.81})$$

with probability at least  $1 - c_1 \exp(-c_2 n)$  for some positive constants  $c_1$  and  $c_2$ . Since the above holds uniformly for any support pattern  $\mathcal{I}$ , we conclude that

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q_1(\mathbf{q}) - Q'_1(\mathbf{q})| \leq \frac{16}{\theta^{3/2} p^{3/2}} \sqrt{n \log p} \quad (\text{F.82})$$

with probability at least  $1 - c_1 \exp(-c_2 n)$ . Now it is sufficient to let  $p \geq C(\xi) n^4 \log n$  for some  $C(\xi) > 0$  to obtain the claimed result.

Similarly, by Lemma B.3 and Lemma B.4 in Appendix B, we have

$$\begin{aligned} & \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\mathbf{Q}_2(\mathbf{q}) - \mathbf{Q}'_2(\mathbf{q})\|_2 \\ &= \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left\| \frac{1}{p} \sum_{k=1}^p \mathbf{g}^k \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] - \frac{1}{p} \sum_{k=1}^p \mathbf{g}'^k \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k] \right\|_2 \\ &\leq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \left\| \frac{1}{p} \sum_{k=1}^p \mathbf{g}^k \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] - \frac{1}{p} \sum_{k=1}^p \mathbf{g}'^k \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] \right\|_2 + \left\| \frac{1}{p} \sum_{k=1}^p \mathbf{g}'^k \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}^k] - \frac{1}{p} \sum_{k=1}^p \mathbf{g}'^k \mathcal{S}_\lambda[\mathbf{q}^\top \mathbf{y}'^k] \right\|_2 \\ &\leq \frac{1}{p} \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \sum_{k=1}^p \|\mathbf{g}^k - \mathbf{g}'^k\|_2 |\mathbf{q}^\top \mathbf{y}^k| + \frac{1}{p} \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \sum_{k=1}^p \|\mathbf{g}'^k\|_2 |\mathbf{q}^\top (\mathbf{y}^k - \mathbf{y}'^k)| \\ &\leq \frac{1}{p} (\|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} \|\mathbf{Y}\|_{\ell^2 \rightarrow \ell^1} + \|\mathbf{G}'\|_{\ell^2 \rightarrow \ell^\infty} \|\mathbf{Y} - \mathbf{Y}'\|_{\ell^2 \rightarrow \ell^1}) \\ &\leq \frac{1}{p} \left( \frac{32n}{\sqrt{\theta} p} \times 3\sqrt{p} + 16\sqrt{\frac{n}{\theta p}} \times \frac{10}{\theta} \sqrt{n \log p} \right) \leq \frac{192n\sqrt{\log p}}{\theta^{3/2} p^{3/2}} \end{aligned} \quad (\text{F.83})$$

holds conditioned on any support pattern  $\mathcal{I}$ , with probability at least  $1 - c_3 \exp(-c_4 n)$  for some positive constants  $c_3$  and  $c_4$ , which similarly implies the bound holds uniformly, regardless of the support, with the same probability. It is sufficiently to have  $\exp(n/2)/2 \geq p \geq C_2(\xi) n^4 \log n$  to obtain the claimed result.  $\square$

## G Large $|q_1|$ Iterates Staying in Safe Region for Rounding

**Proposition G.1.** *There exists a constant  $\theta_0 > 0$ , such that for any  $\theta \in (\frac{1}{\sqrt{n}}, \theta_0)$ , whenever  $\exp(n) \geq p \geq C n^4 \log n$  for some large constant  $C > 0$ , we have*

$$\frac{|Q'_1(\mathbf{q})|}{\|\mathbf{Q}'(\mathbf{q})\|_2} \geq 2\sqrt{\theta}, \quad (\text{G.1})$$

for all  $\mathbf{q} \in \mathbb{S}^{n-1}$  satisfying  $|q_1| > 3\sqrt{\theta}$ , with probability at least  $1 - c' \exp(-c'' n)$  for some positive constants  $c'$  and  $c''$ .

*Proof.* For notational simplicity, w.l.o.g. we will proceed to prove assuming  $q_1 > 0$ . The proof for  $q_1 < 0$  is similar by symmetry. It is equivalent to show that

$$\frac{\|\mathbf{Q}'_2(\mathbf{q})\|_2}{|Q'_1(\mathbf{q})|} < \sqrt{\frac{1}{4\theta} - 1}, \quad (\text{G.2})$$

which is implied by

$$\mathcal{L}(\mathbf{q}) \doteq \frac{\|\mathbb{E}\mathbf{Q}_2(\mathbf{q})\|_2 + \|\mathbf{Q}'_2(\mathbf{q}) - \mathbb{E}\mathbf{Q}_2(\mathbf{q})\|_2}{\mathbb{E}Q_1(\mathbf{q}) - |Q'_1(\mathbf{q}) - \mathbb{E}Q_1(\mathbf{q})|} < \sqrt{\frac{1}{4\theta} - 1} \quad (\text{G.3})$$

for any  $\mathbf{q} \in \mathbb{S}^{n-1}$  satisfying  $q_1 > 3\sqrt{\theta}$ . Recall from (F.22) that

$$\mathbb{E}Q_1(\mathbf{q}) = \sqrt{\frac{\theta}{p}} \left\{ \left[ \alpha \Psi\left(-\frac{\alpha}{\sigma}\right) + \beta \Psi\left(\frac{\beta}{\sigma}\right) \right] + \sigma \left[ \psi\left(\frac{\beta}{\sigma}\right) - \psi\left(-\frac{\alpha}{\sigma}\right) \right] \right\}, \quad (\text{G.4})$$

where

$$\alpha = \frac{1}{\sqrt{p}} \left( \frac{q_1}{\sqrt{\theta}} + 1 \right), \quad \beta = \frac{1}{\sqrt{p}} \left( \frac{q_1}{\sqrt{\theta}} - 1 \right), \quad \sigma = \|\mathbf{q}_2\|_2 / \sqrt{p}. \quad (\text{G.5})$$

Noticing the fact that

$$\psi\left(\frac{\beta}{\sigma}\right) - \psi\left(-\frac{\alpha}{\sigma}\right) \geq 0, \quad (\text{G.6})$$

$$\Psi\left(\frac{\beta}{\sigma}\right) = \Psi\left(\frac{1}{\sqrt{1-q_1^2}} \left( \frac{q_1}{\sqrt{\theta}} - 1 \right)\right) \geq \Psi(2) \geq \frac{19}{20} \quad \text{for } q_1 > 3\sqrt{\theta}, \quad (\text{G.7})$$

we have

$$\mathbb{E}Q_1(\mathbf{q}) \geq \frac{\sqrt{\theta}}{p} \left\{ \frac{q_1}{\sqrt{\theta}} \left[ \Psi\left(-\frac{\alpha}{\sigma}\right) + \Psi\left(\frac{\beta}{\sigma}\right) \right] + \Psi\left(-\frac{\alpha}{\sigma}\right) - \Psi\left(\frac{\beta}{\sigma}\right) \right\} \geq \frac{2\sqrt{\theta}}{p} \Psi\left(\frac{\beta}{\sigma}\right) \geq \frac{19\sqrt{\theta}}{10p}. \quad (\text{G.8})$$

Moreover, from (F.23), we have

$$\|\mathbb{E}\mathbf{Q}_2(\mathbf{q})\|_2 = \|\mathbf{q}_2\|_2 \left\{ \frac{2(1-\theta)}{p} \Psi\left(-\frac{\lambda}{\sigma}\right) + \frac{\theta}{p} \left[ \Psi\left(-\frac{\alpha}{\sigma}\right) + \Psi\left(\frac{\beta}{\sigma}\right) \right] \right\} \quad (\text{G.9})$$

$$\leq \frac{2(1-\theta)}{p} \Psi(-1) + \frac{\theta}{p} [\Psi(-1) + 1] \leq \frac{2}{p} \Psi(-1) + \frac{\theta}{p} \leq \frac{2}{5p} + \frac{\theta}{p}, \quad (\text{G.10})$$

where we have used the fact that  $-\lambda/\sigma \leq -1$  and  $-\alpha/\sigma \leq -1$ . Moreover, from results in Proposition F.9 and Proposition F.10 in Appendix F, we know that

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q'_1(\mathbf{q}) - \mathbb{E}Q_1(\mathbf{q})| \leq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q'_1(\mathbf{q}) - Q_1(\mathbf{q})| + \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} |Q_1(\mathbf{q}) - \mathbb{E}Q_1(\mathbf{q})| \leq \frac{1}{8000\theta^{5/2}n^{3/2}p}, \quad (\text{G.11})$$

$$\sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\mathbf{Q}'(\mathbf{q}) - \mathbb{E}\mathbf{Q}(\mathbf{q})\|_2 \leq \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\mathbf{Q}'(\mathbf{q}) - \mathbf{Q}(\mathbf{q})\|_2 + \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \|\mathbf{Q}(\mathbf{q}) - \mathbb{E}\mathbf{Q}(\mathbf{q})\|_2 \leq \frac{1}{8000\theta^2np} \quad (\text{G.12})$$

hold with probability at least  $1 - c' \exp(-c''n)$  for some positive constants  $c'$  and  $c''$  when  $p \geq \Omega(n^4 \log n)$ . Hence, with overwhelming probability, we have

$$\mathcal{L}(\mathbf{q}) \leq \frac{\frac{2}{5p} + \frac{\theta}{p} + \frac{1}{8000\theta^2np}}{\frac{19\sqrt{\theta}}{10p} - \frac{1}{8000\theta^{5/2}n^{3/2}p}} \leq \frac{\frac{3}{5}}{\frac{18\sqrt{\theta}}{10}} \leq \frac{1}{3\sqrt{\theta}} < \sqrt{\frac{1}{4\theta} - 1}, \quad (\text{G.13})$$

whenever  $\theta$  is sufficiently small. This completes the proof.  $\square$

Now, keep the notation in Appendix F for general orthonormal basis. For any current iterate  $\mathbf{q} \in \mathbb{S}^{n-1}$  that is close enough to the target solution, i.e.,  $|\langle \mathbf{q}, \mathbf{R}^\top \mathbf{e}_1 \rangle| = |\langle \mathbf{R}\mathbf{q}, \mathbf{e}_1 \rangle| \geq 3\sqrt{\theta}$ , we have

$$\frac{|\langle \mathbf{Q}'(\mathbf{q}; \hat{\mathbf{Y}}), \mathbf{R}^\top \mathbf{e}_1 \rangle|}{\|\mathbf{Q}'(\mathbf{q}; \hat{\mathbf{Y}})\|_2} = \frac{|\langle \mathbf{R}\mathbf{Q}'(\mathbf{q}; \hat{\mathbf{Y}}), \mathbf{e}_1 \rangle|}{\|\mathbf{R}\mathbf{Q}'(\mathbf{q}; \hat{\mathbf{Y}})\|_2} = \frac{|\langle \mathbf{Q}'(\mathbf{R}\mathbf{q}; \mathbf{Y}'), \mathbf{e}_1 \rangle|}{\|\mathbf{Q}'(\mathbf{R}\mathbf{q}; \mathbf{Y}')\|_2}, \quad (\text{G.14})$$

where we have applied the identity proved in (F.11). Taking  $\mathbf{Rq} \in \mathbb{S}^{n-1}$  as the object of interest, by Proposition G.1, we conclude that

$$\frac{|\langle \mathbf{Q}'(\mathbf{Rq}; \mathbf{Y}'), \mathbf{e}_1 \rangle|}{\|\mathbf{Q}'(\mathbf{Rq}; \mathbf{Y}')\|_2} \geq 2\sqrt{\theta} \quad (\text{G.15})$$

with overwhelming probability.

## H Bounding Iteration Complexity

**Proposition H.1.** *There is a constant  $\theta_0 > 0$ , such that for any  $\theta \in (\frac{1}{\sqrt{n}}, \theta_0)$ , with probability at least  $1 - c' \exp(-c''n)$  ( $c'$  and  $c''$  are positive constants), the ADM algorithm in Algorithm 1, with any initialization  $\mathbf{q}^{(0)} \in \mathbb{S}^{n-1}$  satisfying  $|q_1^{(0)}| \geq \frac{1}{4\sqrt{\theta n}}$ , will produce some iterate  $\bar{\mathbf{q}}$  with  $|\bar{q}_1| > 3\sqrt{\theta}$  at least once in at most  $O(n^4 \log n)$  iterations, provided  $\exp(n) \geq p \geq Cn^4 \log n$  for some large constant  $C$ .*

*Proof.* Recall from Proposition F.1 in Appendix F, the gap

$$G'(\mathbf{q}) = \frac{|Q'_1(\mathbf{q})|}{|q_1|} - \frac{\|\mathbf{Q}'_2(\mathbf{q})\|_2}{\|\mathbf{q}\|_2} \geq \frac{1}{4000\theta^2 np} \quad (\text{H.1})$$

holds uniformly over  $\mathbf{q} \in \mathbb{S}^{n-1}$  satisfying  $\frac{1}{4\sqrt{\theta p}} \leq |q_1| \leq 3\sqrt{\theta}$  with probability at least  $1 - c_1 \exp(-c_2 n)$  for positive constants  $c_1$  and  $c_2$ , provided  $p \geq \Omega(n^4 \log n)$ . The gap  $G'(\mathbf{q})$  implies that

$$|\tilde{Q}'_1(\mathbf{q})| \doteq \frac{|Q'_1(\mathbf{q})|}{\|\mathbf{Q}'(\mathbf{q}_2)\|_2} \geq \frac{|q_1| \|\mathbf{Q}'_2(\mathbf{q})\|_2}{\|\mathbf{q}\|_2 \|\mathbf{Q}'(\mathbf{q})\|_2} + \frac{|q_1|}{4000\theta^2 np \|\mathbf{Q}'(\mathbf{q})\|_2} \quad (\text{H.2})$$

$$\Leftrightarrow |\tilde{Q}'_1(\mathbf{q})| \geq \frac{|q_1|}{\|\mathbf{q}_2\|_2} \sqrt{1 - |\tilde{Q}'_1(\mathbf{q})|^2} + \frac{|q_1|}{4000\theta^2 np \|\mathbf{Q}'(\mathbf{q})\|_2} \quad (\text{H.3})$$

$$\Rightarrow |\tilde{Q}'_1(\mathbf{q})|^2 \geq |q_1|^2 \left( 1 + \frac{\|\mathbf{q}_2\|_2^2}{4000^2 \theta^4 n^2 p^2 \|\mathbf{Q}'(\mathbf{q})\|_2^2} \right). \quad (\text{H.4})$$

Now we know that

$$\sup_{\mathbf{q} \in \Gamma} \|\mathbf{Q}'(\mathbf{q})\|_2 \leq \sup_{\mathbf{q} \in \Gamma} |Q'_1(\mathbf{q})| + \sup_{\mathbf{q} \in \Gamma} \|\mathbf{Q}'_2(\mathbf{q})\|_2 \quad (\text{H.5})$$

$$= \sup_{\mathbf{q} \in \Gamma} \left\| \frac{1}{p} \sum_{k=1}^p x_{0k} S_\lambda [x_{0k} q_1 + \mathbf{q}_2^\top \mathbf{g}^k] \right\| + \sup_{\mathbf{q} \in \Gamma} \left\| \frac{1}{p} \sum_{k=1}^p \mathbf{g}^k S_\lambda [x_{0k} q_1 + \mathbf{q}_2^\top \mathbf{g}^k] \right\|_2 \quad (\text{H.6})$$

$$\leq \frac{1}{p} \left( \sup_{\mathbf{q} \in \Gamma} \sum_{k=1}^p |x_{0k}| |x_{0k} q_1 + \mathbf{q}_2^\top \mathbf{g}^k| + \sup_{\mathbf{q} \in \Gamma} \sum_{k=1}^p \|\mathbf{g}^k\|_2 |x_{0k} q_1 + \mathbf{q}_2^\top \mathbf{g}^k| \right) \quad (\text{H.7})$$

$$\leq 2 \left( \frac{1}{\sqrt{\theta p}} + \sup_{k \in [p]} \|\mathbf{g}^k\|_2 \right) \sup_{\mathbf{q} \in \Gamma} \left( \frac{|q_1|}{\sqrt{\theta p}} + \|\mathbf{q}_2\|_2 \sup_{k \in [p]} \|\mathbf{g}^k\|_2 \right) \quad (\text{H.8})$$

$$\leq 2 \left( \frac{1}{\sqrt{\theta p}} + \sup_{k \in [p]} \|\mathbf{g}^k\|_2 \right)^2. \quad (\text{H.9})$$

From Lemma A.11, we know that  $\sup_{k \in [p]} \|\mathbf{g}^k\|_2 \leq \sqrt{2 \log p} / \sqrt{p} + 2\sqrt{n} / \sqrt{p}$  with probability at least  $1 - \exp(-n/2)$ . Then provided  $p \leq \exp(n)$  and  $\theta \geq 1/\sqrt{n}$ , we obtain

$$\sup_{\mathbf{q} \in \Gamma} \|\mathbf{Q}'(\mathbf{q})\|_2 \leq \frac{50n}{p}. \quad (\text{H.10})$$

So we conclude that

$$\frac{|\tilde{Q}'_1(\mathbf{q})|}{|q_1|} \geq \sqrt{1 + \frac{1-9\theta}{4000^2 \times 50^2 \times \theta^4 n^4}}. \quad (\text{H.11})$$

Therefore, starting with any  $\mathbf{q} \in \mathbb{S}^{n-1}$  such that  $|q_1| \geq \frac{1}{4\sqrt{\theta n}}$ , we will need at most

$$T = \frac{2 \log \left( 3\sqrt{\theta} / \frac{1}{4\sqrt{\theta n}} \right)}{\log \left( 1 + \frac{1-9\theta}{4000^2 \times 50^2 \times \theta^4 n^4} \right)} = \frac{2 \log (12\theta\sqrt{n})}{\log \left( 1 + \frac{1-9\theta}{4000^2 \times 50^2 \times \theta^4 n^4} \right)} \leq \frac{2 \log (12\theta\sqrt{n})}{(\log 2) \frac{1-9\theta}{4000^2 \times 50^2 \times \theta^4 n^4}} \leq C n^4 \log n \quad (\text{H.12})$$

steps to arrive at a  $\bar{\mathbf{q}} \in \mathbb{S}^{n-1}$  with  $|\bar{q}_1| \geq 3\sqrt{\theta}$  for the first time, where  $C > 0$  is a numerical constant, and we assume  $\theta_0 < 1/9$  and used the fact that  $\log(1+x) \geq (\log 2)x$  for  $x \in [0, 1]$  to simplify the final result.  $\square$

## I Rounding to the Desired Solution

For convenience, we will assume the notations we used in Appendix B. Then the rounding scheme can be written as

$$\min_{\mathbf{q}} \|\mathbf{Y}'\mathbf{R}\mathbf{q}\|_1, \quad \text{s.t. } \langle \bar{\mathbf{q}}, \mathbf{q} \rangle = 1, \quad (\text{I.1})$$

for some orthogonal matrix  $\mathbf{R}$ . We will show the rounding procedure get us to the desired solution with overwhelming probability, regardless of the particular orthonormal basis used.

**Proposition I.1.** *Suppose the input basis is  $\mathbf{Y}'$  defined in (B.3) and the ADM algorithm produces  $\bar{\mathbf{q}} \in \mathbb{S}^{n-1}$  with  $\bar{q}_1 > 2\sqrt{\theta}$ . Then there exists some constants  $C, \theta_0 > 0$ , such that when  $p \geq Cn^2$  and  $\theta \in \left(\frac{1}{\sqrt{n}}, \theta_0\right)$ , the rounding procedure with  $\mathbf{r} = \bar{\mathbf{q}}$  returns the desired solution  $\mathbf{e}_1$  with probability at least  $1 - c \exp(-c'n)$  for some numerical constants  $c, c' > 0$ .*

*Proof.* The rounding program (I.1) can be written as

$$\inf_{\mathbf{q}} \|\mathbf{Y}'\mathbf{q}\|_1, \quad \text{s.t. } \bar{q}_1 q_1 + \langle \bar{\mathbf{q}}_2, \mathbf{q}_2 \rangle = 1. \quad (\text{I.2})$$

Consider its relaxation

$$\inf_{\mathbf{q}} \|\mathbf{Y}'\mathbf{q}\|_1, \quad \text{s.t. } \bar{q}_1 q_1 + \|\bar{\mathbf{q}}_2\|_2 \|\mathbf{q}_2\|_2 \geq 1. \quad (\text{I.3})$$

It is obvious that the feasible set of (I.3) contains that of (I.2). So if  $\mathbf{e}_1$  is the unique optimal solution (UOS) of (I.3), it is also the UOS of (I.2). Let  $\mathcal{I} = \text{supp}(\mathbf{x}_0)$ , and consider a modified problem

$$\inf_{\mathbf{q}} \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 |q_1| - \|\mathbf{G}'_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}'_{\mathcal{I}^c}\mathbf{q}_2\|_1, \quad \text{s.t. } \bar{q}_1 q_1 + \|\bar{\mathbf{q}}_2\|_2 \|\mathbf{q}_2\|_2 \geq 1. \quad (\text{I.4})$$

The objective value of (I.4) lower bounds the objective value of (I.3), and are equal when  $\mathbf{q} = \mathbf{e}_1$ . So if  $\mathbf{q} = \mathbf{e}_1$  is the UOS to (I.4), it is also UOS to (I.3), and hence UOS to (I.2) by the argument above. Now

$$- \|\mathbf{G}'_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}'_{\mathcal{I}^c}\mathbf{q}_2\|_1 \geq - \|\mathbf{G}_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}_{\mathcal{I}^c}\mathbf{q}_2\|_1 - \|(\mathbf{G} - \mathbf{G}')\mathbf{q}_2\|_1 \quad (\text{I.5})$$

$$\geq - \|\mathbf{G}_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}_{\mathcal{I}^c}\mathbf{q}_2\|_1 - \|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{q}_2\|_2. \quad (\text{I.6})$$

When  $p \geq \Omega(n^2)$ , by Lemma A.14 and Lemma B.3, we know that

$$\begin{aligned} & - \|\mathbf{G}_{\mathcal{I}}\mathbf{q}_2\|_1 + \|\mathbf{G}_{\mathcal{I}^c}\mathbf{q}_2\|_1 - \|\mathbf{G} - \mathbf{G}'\|_{\ell^2 \rightarrow \ell^1} \|\mathbf{q}_2\|_2 \\ & \geq -\frac{6}{5} \sqrt{\frac{2}{\pi}} 2\theta\sqrt{p} \|\mathbf{q}_2\|_2 + \frac{24}{25} \sqrt{\frac{2}{\pi}} (1-2\theta) \sqrt{p} \|\mathbf{q}_2\|_2 - 8\sqrt{n} \|\mathbf{q}_2\|_2 \doteq \zeta \|\mathbf{q}_2\|_2 \end{aligned} \quad (\text{I.7})$$

holds with probability at least  $1 - c_1 \exp(-c_2 n)$  for some positive constants  $c_1$  and  $c_2$ . Thus, we make a further relaxation of problem (I.2) by

$$\inf_{\mathbf{q}} \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 |q_1| + \zeta \|\mathbf{q}_2\|_2, \quad \text{s.t. } \bar{q}_1 q_1 + \|\bar{\mathbf{q}}_2\|_2 \|\mathbf{q}_2\|_2 \geq 1, \quad (\text{I.8})$$

whose objective value lower bounds that of (I.4). By similar arguments, if  $\mathbf{e}_1$  is UOS to (I.8), it is UOS to (I.2). At the optimal solution to (I.8), notice that it is necessary to have  $\text{sign}(q_1) = \text{sign}(\bar{q}_1)$  and  $\bar{q}_1 q_1 + \|\bar{\mathbf{q}}_2\|_2 \|\mathbf{q}_2\|_2 = 1$ . So (I.8) is equivalent to

$$\inf_{\mathbf{q}} \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 |q_1| + \zeta \|\mathbf{q}_2\|_2, \quad \text{s.t. } \bar{q}_1 q_1 + \|\bar{\mathbf{q}}_2\|_2 \|\mathbf{q}_2\|_2 = 1. \quad (\text{I.9})$$

which is further equivalent to

$$\inf_{q_1} \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 |q_1| + \zeta \frac{1 - |\bar{q}_1| |q_1|}{\|\bar{\mathbf{q}}_2\|_2}, \quad \text{s.t. } |q_1| \leq \frac{1}{|\bar{q}_1|}. \quad (\text{I.10})$$

Notice that the problem in (I.10) is linear in  $|q_1|$  with a compact feasible set, which indicates that the optimal solution only occur at the boundary points  $|q_1| = 0$  and  $|q_1| = 1/|\bar{q}_1|$ . Therefore,  $\mathbf{q} = \mathbf{e}_1$  is the UOS of (I.10) if and only if

$$\frac{1}{|\bar{q}_1|} \left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 < \frac{\zeta}{\|\bar{\mathbf{q}}_2\|_2}. \quad (\text{I.11})$$

Since  $\left\| \frac{\mathbf{x}_0}{\|\mathbf{x}_0\|_2} \right\|_1 \leq \sqrt{2\theta p}$  conditioned on  $\mathcal{E}_0$ , it is sufficient to have

$$\frac{\sqrt{2\theta p}}{2\sqrt{\theta}} \leq \zeta = \frac{24}{25} \sqrt{\frac{2}{\pi}} \sqrt{p} \left( 1 - \frac{9}{2}\theta - \frac{25}{3} \sqrt{\frac{n}{p}} \right). \quad (\text{I.12})$$

Therefore there exists a constant  $\theta_0 > 0$ , such that whenever  $\theta \leq \theta_0$ , the rounding returns  $\mathbf{e}_1$ , completing the proof.  $\square$

When the input basis is  $\mathbf{Y}'\mathbf{R}$  for some  $\mathbf{R} \neq \mathbf{I}$ , if the ADM algorithm produces some  $\bar{\mathbf{q}} = \mathbf{R}^\top \mathbf{q}'$ , such that  $q'_1 > 2\sqrt{\theta}$ . It is not hard to see that now the rounding (I.1) is equivalent to

$$\min_{\mathbf{q}} \|\mathbf{Y}'\mathbf{R}\mathbf{q}\|_1, \quad \text{s.t. } \langle \mathbf{q}', \mathbf{R}\mathbf{q} \rangle = 1. \quad (\text{I.13})$$

Renaming  $\mathbf{R}\mathbf{q}$ , it follows from the above argument that at optimum  $\mathbf{q}^*$  it holds that  $\mathbf{R}\mathbf{q}^* = \mathbf{e}_1$  with overwhelming probability.