

# Convex and Nonconvex Formulations for Mixed Regression with Two Components: Minimax Optimal Rates

Yudong Chen, Xinyang Yi, and Constantine Caramanis, *Member, IEEE*,

**Abstract**—We consider the mixed regression problem with two components, under adversarial and stochastic noise. We give a convex optimization formulation that provably recovers the true solution, as well as a nonconvex formulation that works under more general settings and remains tractable. Upper bounds are provided on the recovery errors for both arbitrary noise and stochastic noise models. We also give matching minimax lower bounds (up to log factors), showing that our algorithm is information-theoretically optimal in a precise sense. Our results represent the first tractable algorithm guaranteeing successful recovery with tight bounds on recovery errors and sample complexity. Moreover, we pinpoint the statistical cost of mixtures: our minimax-optimal results indicate that the mixture poses a fundamentally more difficult problem in the low-SNR regime, where the learning rate changes.

## I. INTRODUCTION

This paper considers the problem of *mixed linear regression*, where each observation of the output variable comes from one of two unknown regression vectors. Formally, we observe  $n$  data points  $(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$ , which satisfies

$$y_i = \begin{cases} \langle \mathbf{x}_i, \beta_1^* \rangle + e_i, & \text{if } z_i = 0, \\ \langle \mathbf{x}_i, \beta_2^* \rangle + e_i, & \text{if } z_i = 1, \end{cases} \quad i = 1, \dots, n,$$

where  $\beta_1^*$  and  $\beta_2^*$  are two unknown regression vectors in  $\mathbb{R}^p$ ,  $e_i$  is the noise, and  $z_i \in \{0, 1\}$  can be thought of as a hidden label determining which regression vector generates the  $i$ -th data point. Our goal is to estimate the pair  $\beta_1^*$  and  $\beta_2^*$ . We consider the setting where the covariates  $\mathbf{x}_i$  and the noise  $e_i$  are independent of the labels, and in particular, no information about the labels can be directly inferred from them. This setting means that predicting labels exactly is impossible, with or without knowing  $\beta_1^*$  and  $\beta_2^*$ .

Manuscript received XXX XX, 20XX; revised XXX XX, 20XX. Y. Chen was supported by NSF grants CCF-1704828, CRII award 1657420, and the School of Operations Research and Information Engineering at Cornell University. X. Yi and C. Caramanis were supported by NSF Grants EEC-1056028, CNS-1302435, CCF-1116955, and the USDOT UTC-D-STOP Center at UT-Austin.

Y. Chen is with the School of Operations Research and Information Engineering, Cornell University, Ithaca, NY 14850 USA (yudong.chen@cornell.edu).

X. Yi and C. Caramanis are with the Department of Electrical and Computer Engineering, the University of Texas at Austin, Austin, TX 78712 USA (yixy@utexas.edu, constantine@utexas.edu).

This work was presented in part at the Conference on Learning Theory, 2014.

Copyright (c) 2017 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to [pubs-permissions@ieee.org](mailto:pubs-permissions@ieee.org).

If the label of each sample is known, the problem decomposes into two standard linear regression problems, and can be easily solved. Without knowing the labels, however, the problem is significantly more difficult. The main challenge of mixture models, and in particular mixed regression falls in the intersection of the *statistical* and *computational* constraints: the problem is difficult when one cares both about an efficient algorithm, and about near-optimal sample complexity and estimation error. Exponential-effort brute force search (over all possibilities of the labels) typically results in statistically near-optimal estimators. On the other hand, recent tensor-based methods give a polynomial-time algorithm, but at the cost of an  $O(p^6)$  sample size (recall  $\beta_1^*, \beta_2^* \in \mathbb{R}^p$ ) instead of the optimal rate  $O(p)$ .<sup>1</sup> The Expectation Maximization (EM) algorithm is computationally very efficient, and widely used in practice for mixture problems. However, there has been only limited understanding of its behavior, and in particular, no general theoretical guarantees on global convergence are known.

*a) Our contributions and the cost of mixtures.:* In this paper, we tackle both statistical and algorithmic objectives at once. The algorithms we give are computationally efficient, specified by solutions of convex optimization problems as well as a tractable nonconvex formulation, which can be solved by polynomial-time, globally convergent procedures. In both the noisy and noiseless settings our results provide better statistical guarantees compared to the best known previous results. In particular, in both the arbitrary noise and stochastic noise regimes, we provide matching estimation error bounds and minimax lower bounds, showing our results are statistically optimal.

An interesting feature of our minimax results is that we pinpoint the statistical cost of dealing with a mixture problem compared to ordinary regression problems. As we detail below in Theorems 3, 4 and 8, we show that in the high SNR regime, there is (up to log factors) no loss, and one can expect to recover the regression parameters at the parametric learning rate (as with ordinary regression). At the low-SNR regime (where, to the best of our knowledge, there are no previous results on mixed regression), the rate changes from  $(1/n)^{1/2}$  to  $(1/n)^{1/4}$ . It is of interest to explore to what we owe this change in rate. On the algorithmic side, our approach is to solve a related low-rank matrix regression problem in the lifted

<sup>1</sup>It may be possible to improve the sample requirement of tensor methods to  $O(p^4)$  for the case of Gaussian design.

space, to estimate (something related to) the tensor product of the two regressors. We show that the regression in this lifted matrix space has an error that decays with  $(1/n)^{1/2}$ . At high SNRs, we prove a perturbation result (Theorem 4) that shows that the top eigenvectors of this matrix inherit the same rate. At low SNRs, however, this comes at a cost, and the rate reduces to  $(1/n)^{1/4}$ . Our matching lower bounds show that this cost of converting an error bound in the lifted space to an error bound in the regressor space is fundamental and thus encapsulates the crux of the challenge. The lower bounds are established by showing that at low SNR, distinguishing one *mixture* of Gaussians from another is more difficult than distinguishing one Gaussian from another, at exactly the rate change indicated above.

Specifically, our contributions are as follows:

- In the arbitrary noise setting where the noise  $e = (e_1, \dots, e_n)^\top$  can be adversarial, we show that under mild technical conditions, as long as the numbers of observations for each regression vector satisfy  $n_1, n_2 \gtrsim p$ , our convex algorithm produces an estimator  $(\hat{\beta}_1, \hat{\beta}_2)$  that satisfies

$$\|\hat{\beta}_b - \beta_b^*\|_2 \lesssim \frac{\|e\|_2}{\sqrt{n}}, \quad \text{for } b = 1, 2.$$

This result immediately implies exact recovery in the noiseless case with  $O(p)$  samples.

- In the stochastic noise setting with sub-Gaussian noise and balanced labels ( $n_1/n_2 \rightarrow 1$ ), if we have  $n_1, n_2 \gtrsim p$  and a Gaussian design matrix, our convex estimator satisfies the following error bound (omitting polylog factors):

$$\|\hat{\beta}_b - \beta_b^*\|_2 \lesssim \begin{cases} \sigma \sqrt{\frac{p}{n}}, & \text{if } \gamma \geq \sigma, \\ \frac{\sigma^2}{\gamma} \sqrt{\frac{p}{n}}, & \text{if } \sigma \left(\frac{p}{n}\right)^{\frac{1}{4}} \leq \gamma \leq \sigma, \\ \sigma \left(\frac{p}{n}\right)^{\frac{1}{4}}, & \text{if } \gamma \leq \sigma \left(\frac{p}{n}\right)^{\frac{1}{4}} \end{cases}$$

for  $b = 1, 2$ , where  $\gamma$  is an lower bound of the signal strength  $\|\beta_1^*\|_2 + \|\beta_2^*\|_2$ , and  $\sigma^2$  is the variance of the noise  $e_i$ .

- In the stochastic noise setting with imbalanced labels, we propose a nonconvex optimization based formulation along with a polynomial-time solver with provable guarantees. Specifically, we show that when  $\min\{n_1/n_2, n_2/n_1\}$  is lower bounded by any constant and  $\gamma/\sigma \gtrsim 1$ , our estimator satisfies the bound

$$\|\hat{\beta}_b - \beta_b^*\|_2 \lesssim \sigma \sqrt{\frac{p}{n}}, \quad \text{for } b = 1, 2.$$

- Finally, in both the arbitrary and stochastic noise settings, we provide minimax lower bounds that match the above upper bounds up to at most polylog factors, thereby showing that the results obtained by the estimates produced by our algorithms are information-theoretically optimal. Particularly in the stochastic setting, the situation is a bit more subtle: the minimax rates in fact depend on the signal-to-noise ratio (SNR)  $\gamma/\sigma$  and exhibit several phases, showing a qualitatively different behavior than in standard regression and many other parametric problems (for which the minimax rate is usually  $\sqrt{1/n}$ ).

## II. RELATED WORK

Mixture models and latent variable modeling are very broadly used in a wide array of contexts far beyond regression. Subspace clustering [19, 26, 32], Gaussian mixture models [3, 21] and  $k$ -means/medians clustering [13] are popular examples of unsupervised learning for mixture models. Arguably the most popular and broadly implemented approach to mixture problems, including mixed regression, is the EM algorithm [18, 22]. In fact, EM has been used for mixed regression for various application domains [20, 31]. Despite its wide use, still little is known about its performance beyond local convergence [4, 34].

One exception is the work in [37], which studies mixed regression in the noiseless setting. They propose an alternating minimization approach initialized by a grid search, and show that their algorithm recovers the regression vector with a sample complexity of  $O(p \log^2 p)$ . Extension to the noisy setting is recently considered by the authors of [4]. Focusing on the stochastic noise setting with sufficiently high SNR (that is, when  $\gamma \gtrsim \sigma$ ; cf. Section I), they show that the EM algorithm with good initialization achieves the error bound  $\|\hat{\beta}_b - \beta_b^*\|_2 \lesssim \sqrt{\gamma^2 + \sigma^2} \sqrt{\frac{p}{n}}$ . In the work [27], EM is adapted to the high-dimensional regression setting, where the regression vectors are known to be sparse and EM is used to solve a penalized (for sparsity) likelihood function. This generalized EM approach achieves support-recovery, though once restricted to that support where the problem becomes a standard mixed regression problem, only convergence to a local optimum can be guaranteed. Very recently for this high-dimensional sparse regression setting, the works in [33] and [36] establish local convergence (that is., assuming that EM is run from a good initialization) of truncated and regularized EM algorithms to a statistically optimal solution.

Mixture problems have been explored using the technology of tensors recently developed in the literature [2, 21]. The authors of [12] consider a tensor-based approach, regressing  $\mathbf{x}^{\otimes 3}$  against  $y^3$  and then using the tensor decomposition techniques to efficiently recover each  $\beta_b^*$ . These methods are not limited to the mixture of only two components, as we are. Yet, even for two components, the tensor approach requires  $O(p^6)$  samples, compared to  $O(p \cdot \text{polylog}(p))$  that our work requires. As noted in their work, the higher sampling requirement seems to be a common difficulty for algorithms based on high order tensors.

In this work we consider the setting with two mixture components. Binary latent factors are common modeling tools for many applications, as they model on/off-type relationships, among others. We refer to the paper [31] for numerous examples of such problems. While the extension to more than two components is of great interest, much is unknown even for two components. In particular, globally convergent algorithms with even near-optimal sample complexity are, to the best of our knowledge, unknown. In fact, in the low-SNR regime (see below for details) we are unaware of any algorithm able to guarantee non-trivial estimation of the parameters. As explained above, our work shows that there may be a fundamental reason for this: the minimax-optimal rate, and hence

the (statistical) difficulty of the mixed regression problem, is different at the high-SNR and low-SNR regimes.

### III. MAIN RESULTS

In this section we present this paper's main results: our algorithms for mixed regression and matching statistical upper and lower bounds. In addition, we describe the precise setup and assumptions, and introduce the basic notation we use.

#### A. Problem Set Up

Suppose that  $\beta_1^*$  and  $\beta_2^*$  are two unknown vectors in  $\mathbb{R}^p$ . We observe  $n$  noisy linear measurements  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ . Let  $\mathcal{I}_1$  and  $\mathcal{I}_2$  denote the (unknown) subsets of the measurements corresponding to  $\beta_1^*$  and  $\beta_2^*$ , respectively, with  $\mathcal{I}_1 \cup \mathcal{I}_2 = [n]$ ,  $\mathcal{I}_1 \cap \mathcal{I}_2 = \emptyset$ ,  $n_1 := |\mathcal{I}_1|$ ,  $n_2 := |\mathcal{I}_2|$  and  $n_1 + n_2 = n$ . The measurements satisfy the following model: for each  $b \in \{1, 2\}$  and  $i \in \mathcal{I}_b$ ,

$$y_i = \langle \mathbf{x}_i, \beta_b^* \rangle + e_i. \quad (1)$$

Given the data  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , the goal is to recover  $\beta_1^*$  and  $\beta_2^*$ .

To measure the estimation error, we use the following symmetric (semi-)metric: for each pair  $\theta = (\beta_1, \beta_2)$  and  $\theta' = (\beta'_1, \beta'_2)$  in  $\mathbb{R}^p \times \mathbb{R}^p$ , define

$$\rho(\theta, \theta') := \min \left\{ \|\beta_1 - \beta'_1\|_2 + \|\beta_2 - \beta'_2\|_2, \|\beta_1 - \beta'_2\|_2 + \|\beta_2 - \beta'_1\|_2 \right\}, \quad (2)$$

which is the total  $\ell_2$  distance up to permutation. For an estimate  $\hat{\theta} = (\hat{\beta}_1, \hat{\beta}_2)$  of the true regression vector pair  $\theta^* := (\beta_1^*, \beta_2^*)$ , we are interested in bounding the symmetric estimation error  $\rho(\hat{\theta}, \theta^*)$ . In the presence of noise, the correct labels (or equivalently, the sets  $\mathcal{I}_1$  and  $\mathcal{I}_2$ ) are in general irrecoverable.

The key insight to our optimization formulation is to utilize a particular lifting to the space of  $p \times p$  matrices, that then allows recovering  $\beta_1$  and  $\beta_2$  with controlled perturbation, without requiring working in the space of 3-tensors. More concretely, defining the true quantities  $\mathbf{K}^*$  and  $\mathbf{g}^*$  as

$$\begin{aligned} \mathbf{K}^* &:= \frac{1}{2} (\beta_1^{* \top} \beta_2^* + \beta_2^{* \top} \beta_1^*) \in \mathbb{R}^{p \times p}, \\ \mathbf{g}^* &:= \frac{1}{2} (\beta_1^* + \beta_2^*) \in \mathbb{R}^p, \end{aligned} \quad (3)$$

we solve a linear regression in the matrix and vector variables  $\mathbf{K}$  and  $\mathbf{g}$  in order to produce an approximation of  $\mathbf{K}^*$  and  $\mathbf{g}^*$ . Given the pair  $(\mathbf{K}^*, \mathbf{g}^*)$ , the true regression vectors  $\beta_1^*$  and  $\beta_2^*$  can be recovered exactly, using the identity

$$\mathbf{J}^* := \mathbf{g}^{* \top} \mathbf{g}^* - \mathbf{K}^* = \frac{1}{4} (\beta_1^* - \beta_2^*) (\beta_1^* - \beta_2^*)^\top.$$

Letting  $\lambda^*$  and  $\mathbf{v}^*$  be the first eigenvalue-eigenvector pair of  $\mathbf{J}^*$ , we have  $\sqrt{\lambda^*} \mathbf{v}^* := \pm \frac{1}{2} (\beta_1^* - \beta_2^*)$  and thus together with  $\mathbf{g}^*$  we can recover  $(\beta_1^*, \beta_2^*)$ . Given an approximation,  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of  $(\mathbf{K}^*, \mathbf{g}^*)$ , this procedure is still well defined, and is described fully in Algorithm 1. In fact, this estimation procedure is stable: if  $\hat{\mathbf{K}}$  and  $\hat{\mathbf{g}}$  are close to true quantities

---

#### Algorithm 1 Estimate $(\beta_1^*, \beta_2^*)$ from $(\mathbf{K}, \mathbf{g})$

---

Input:  $(\hat{\mathbf{K}}, \hat{\mathbf{g}}) \in \mathbb{R}^{p \times p} \times \mathbb{R}^p$ .

Compute the matrix  $\hat{\mathbf{J}} = \hat{\mathbf{g}} \hat{\mathbf{g}}^\top - \hat{\mathbf{K}}$ , and its first eigenvalue-eigenvector pair  $\hat{\lambda}$  and  $\hat{\mathbf{v}}$ .

Compute  $\hat{\beta}_1 = \hat{\mathbf{g}} + \sqrt{\hat{\lambda}} \hat{\mathbf{v}}$  and  $\hat{\beta}_2 = \hat{\mathbf{g}} - \sqrt{\hat{\lambda}} \hat{\mathbf{v}}$ .

Output:  $\hat{\theta} = (\hat{\beta}_1, \hat{\beta}_2)$

---

$\mathbf{K}^*$  and  $\mathbf{g}^*$ , then Algorithm 1 outputs a pair  $(\hat{\beta}_1, \hat{\beta}_2)$  that is close to the true regression vectors  $(\beta_1^*, \beta_2^*)$ .

Below we consider separately the settings with arbitrary noise and stochastic noise in Sections III-B and III-C, and give our algorithms with rigorous sample complexity and estimation error bounds. In Section III-D we further provide matching (up to at most a polylog factor) minimax lower bounds.

a) *Notation*:: We use lower case bold letters to denote vectors, and capital bold-face letters for matrices. For a vector  $\mathbf{u}$ , the notations  $u_i$  and  $u(i)$  both denote its  $i^{\text{th}}$  coordinate. We use standard notation for matrix and vector norms, e.g.,  $\|\cdot\|_*$  to denote the nuclear norm,  $\|\cdot\|_F$  the Frobenius norm, and  $\|\cdot\|$  the operator/spectral norm. We define two quantities that we use repeatedly:

$$\alpha := \frac{\|\beta_1^* - \beta_2^*\|_2^2}{\|\beta_1^*\|_2^2 + \|\beta_2^*\|_2^2} \quad \text{and} \quad \gamma := \|\beta_1^*\|_2 + \|\beta_2^*\|_2. \quad (4)$$

Note that the quantity  $\alpha \in [0, 2]$  measures the angle between  $\beta_1^*$  and  $\beta_2^*$ , and is strictly positive whenever  $\beta_1^* \neq \beta_2^*$ . For a fixed  $\alpha$ , the parameter  $\gamma$  quantifies the signal strength of  $\beta_1^*$  and  $\beta_2^*$ . We say a number  $c$  is a *numerical constant* if  $c$  is independent of the dimension  $p$ , the number of measurements  $n$ , the quantity  $\alpha$  and the magnitude/variance of the noise  $e$ . For ease of parsing, we typically use  $c$  to denote a *large* constant, and  $\frac{1}{c}$  for a *small* constant.

#### B. Arbitrary Noise

We consider first mixed regression with arbitrary noise, with the following specific setting. We take the covariate vectors  $\{\mathbf{x}_i\}$  to have i.i.d. zero-mean and sub-Gaussian entries with sub-Gaussian norm<sup>2</sup> bounded by a numeric constant and satisfy  $\mathbb{E}[(\mathbf{x}_i(l))^2] = 1$ ,  $\mathbb{E}[(\mathbf{x}_i(l))^4] = \mu$  for each  $i \in [n]$  and  $l \in [p]$ .<sup>3</sup> We assume that the forth moment  $\mu$  is a fixed constant and independent of the parameters  $p$ ,  $n$  and  $\alpha$ . If the covariates  $\{\mathbf{x}_i\}$  are standard Gaussian vectors, then these assumptions are satisfied with unit sub-Gaussian norm and  $\mu = 3$ . The only assumption on the noise vector  $\mathbf{e} = (e_1, e_2, \dots, e_n)^\top$  is that it is bounded in  $\ell_2$  norm. The noise  $\mathbf{e}$  is otherwise arbitrary, possibly adversarial, and even potentially depending on  $\{\mathbf{x}_i\}$  and  $(\beta_1^*, \beta_2^*)$ .

<sup>2</sup>The sub-Gaussian norm of a zero-mean random variable  $X$  is defined as  $\|X\|_{\psi_2} := \inf \{b \geq 0 \mid \mathbb{E} \exp(bX) \leq \exp(b^2 t^2/2)\}$ . The variable  $X$  is called sub-Gaussian if  $\|X\|_{\psi_2} < \infty$ .

<sup>3</sup>Recall that, as shown in the paper [37], the general mixed regression problem with deterministic covariates is NP-hard even in the noiseless setting.

Our algorithm is based on the following convex program:

$$\min_{\mathbf{K}, \mathbf{g}} \|\mathbf{K}\|_* \quad (5)$$

$$\text{s.t. } \sum_{i=1}^n \left| -\langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{K} \rangle + 2y_i \langle \mathbf{x}_i, \mathbf{g} \rangle - y_i^2 \right| \leq \eta. \quad (6)$$

The intuition is that in the noiseless setting with  $\mathbf{e} = \mathbf{0}$ , if we substitute the desired solution  $(\mathbf{K}^*, \mathbf{g}^*)$  given by equation (3) into the above program, the LHS of the constraint (6) becomes zero; moreover, the rank of  $\mathbf{K}^*$  is at most 2, and minimizing the nuclear norm term in the objective (5) encourages the optimal solution to have low rank. Our theoretical results give a precise way to set the right hand side,  $\eta$ , of the constraint.

The next two theorems summarize our results for the arbitrary noise setting. Theorem 1 provides guarantees on how close the optimal convex optimization solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  is to the true quantities  $(\mathbf{K}^*, \mathbf{g}^*)$ . Then the companion result, Theorem 2, provides quality bounds on  $(\hat{\beta}_1, \hat{\beta}_2)$ , produced by applying Algorithm 1 to the solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of the convex program.

**Theorem 1** (Arbitrary noise). *There exist numerical positive constants  $\{c_i\}_{i=0}^6$  for which the following holds. Assume that  $\max\left\{\frac{n_1}{n_2}, \frac{n_2}{n_1}\right\} \leq c_0$ . Suppose, moreover, that (i)  $\mu > 1$  and  $\alpha > 0$ , (ii)  $\min\{n_1, n_2\} \geq c_3 \frac{1}{\alpha} p$ , and (iii) the noise satisfies the bound*

$$\|\mathbf{e}\|_2 \leq \frac{\sqrt{\alpha}}{c_5} \sqrt{n} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2).$$

Suppose further that we choose the tuning parameter  $\eta$  to satisfy

$$\eta \geq c_4 \sqrt{n} \|\mathbf{e}\|_2 \|\beta_2^* - \beta_1^*\|_2.$$

Then, with probability at least  $1 - c_1 \exp(-c_2 n)$ , any optimal solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  to the convex program (5)–(6) satisfies the error bounds

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}^*\|_F &\leq c_6 \frac{1}{\sqrt{\alpha n}} \eta, \\ \|\hat{\mathbf{g}} - \mathbf{g}^*\|_2 &\leq c_6 \frac{1}{\sqrt{\alpha n} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2)} \eta. \end{aligned}$$

Given the solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$ , we then estimate  $\theta^* = (\beta_1^*, \beta_2^*)$  by using Algorithm 1, which is stable as shown by the theorem below.

**Theorem 2** (Estimating  $\beta_b^*$ , arbitrary noise). *Suppose that the conditions in Theorem 1 hold, and  $\eta \asymp \sqrt{n} \|\mathbf{e}\|_2 \|\beta_2^* - \beta_1^*\|_2$ . Then with probability at least  $1 - c_1 \exp(-c_2 n)$ , the output  $\hat{\theta} = (\hat{\beta}_1, \hat{\beta}_2)$  of Algorithm 1 satisfies the error bound*

$$\rho(\hat{\theta}, \theta^*) \leq \frac{1}{c_3 \sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}}.$$

Theorem 2 immediately implies exact recovery in the noiseless case.

**Corollary 1** (Exact recovery). *Suppose that  $\mathbf{e} = \mathbf{0}$ , the conditions (i) and (ii) in Theorem 1 hold, and  $\eta = 0$ . Then with probability at least  $1 - c_1 \exp(-c_2 n)$ , Algorithm 1 returns the true regression vectors  $\{\beta_1^*, \beta_2^*\}$ .*

Below we provide several remarks on the above theoretical results.

*a) Discussion of assumptions::* Theorem 1 involves several mild technical assumptions.

- 1) The condition  $\mu > 1$  in Theorem 1 is satisfied, for instance, if  $\{\mathbf{x}_i\}$  is Gaussian (with  $\mu = 3$ ). Moreover, this condition is in general necessary. To see this, suppose that each  $\mathbf{x}_i(l)$  is a Rademacher  $\pm 1$  variable, which has forth moment  $\mu = 1$ , and the true regression vectors  $\beta_1^*$  and  $\beta_2^*$  are in  $\mathbb{R}^2$ . The response variable  $y_i$  must have the form

$$y_i = \pm(\beta_b^*)_1 \pm (\beta_b^*)_2.$$

Consider two possibilities:  $\beta_1^* = -\beta_2^* = (1, 0)^\top$  or  $\beta_1^* = -\beta_2^* = (0, 1)^\top$ . In both cases, the observed data  $(\mathbf{x}_i, y_i)$  takes any one of the values in  $\{\pm 1\}^2 \times \{\pm 1\}$  with equal probability, and hence the problem is unidentifiable as it is impossible to distinguish the two possibilities above.

- 2) The condition  $\alpha > 0$  holds if  $\beta_1^*$  and  $\beta_2^*$  are not equal. Suppose that  $\alpha$  is lower-bounded by a constant. The main assumption on the noise, namely, the condition  $\|\mathbf{e}\|_2 \lesssim \sqrt{n} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2)$  in Theorem 1, cannot be substantially relaxed if we want a bound on  $\|\hat{\mathbf{g}} - \mathbf{g}^*\|_2$ . Indeed, if  $|e_i| \gtrsim \|\beta_b^*\|_2$  for all  $i$ , then an adversary may choose  $e_i$  such that

$$y_i = \mathbf{x}_i^\top \beta_b^* + e_i = 0, \quad \forall i,$$

in which case the convex program (5)–(6) becomes independent of  $\mathbf{g}$ . That said, the case when the noise bound condition is violated can be handled easily. Suppose that  $\|\mathbf{e}\|_2 \geq c_4 \sqrt{\alpha n} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2)$  for any positive constant  $c_4$ . A standard argument for ordinary linear regression shows that the blind estimator  $\hat{\beta} := \arg \min_{\beta} \sum_{i \in \mathcal{I}_1 \cup \mathcal{I}_2} |\mathbf{x}_i^\top \beta - y_i|$  satisfies w.h.p. the bound

$$\max \left\{ \|\hat{\beta} - \beta_1^*\|_2, \|\hat{\beta} - \beta_2^*\|_2 \right\} \lesssim \frac{\|\mathbf{e}\|_2}{\sqrt{n}},$$

and this bound is optimal (see the minimax lower bound in Section III-D). Therefore, the condition (iv) in Theorem 1 is not really restrictive, in the sense that the case when it holds is precisely the interesting setting.

- 3) Finally, note that if  $n_1/n_2 \rightarrow +\infty$ , then a single regression vector  $\beta_1^*$  explains 100% (asymptotically) of the observed data. Moreover, the standard least squares solution provides an accurate estimator of this  $\beta_1^*$ .

*b) Optimality of sample complexity::* The sample complexity bounds of Theorem 2 and Corollary 1 are optimal. The results require the sample size  $n_1$  and  $n_2$  to be  $\Omega(p)$ . Since we are estimating two  $p$  dimensional vectors without any further structure, this result cannot be improved in general.

### C. Stochastic Noise

We now turn to the stochastic noise setting. We assume that the covariates  $\{\mathbf{x}_i\}$  have i.i.d. Gaussian entries with zero mean and unit variance. For the noise, we assume that  $\{e_i\}$  are i.i.d., zero-mean and sub-Gaussian with variance  $\mathbb{E}[e_i^2] = \sigma^2$  and

sub-Gaussian norm  $\|e_i\|_{\psi_2} \leq c\sigma$  for some absolute constant  $c$ , and are independent of  $\{\mathbf{x}_i\}$ .

We discuss two algorithms for consistent estimation of  $(\beta_1^*, \beta_2^*)$ . First, for Gaussian covariates in the *balanced setting* (that is,  $n_1/n_2 \rightarrow 1$ ), we show that by solving a simple convex program, we have asymptotic consistency for any SNR  $\gamma/\sigma$ . The rates we obtain match information-theoretic lower bounds we give in Section III-D, and hence are minimax optimal. An interesting feature we observe is the rate change in the high- and low-SNR regimes mentioned above — a feature that precisely identifies the cost of solving a mixture problem.

Second, in the general *imbalanced setting* (that is,  $n_1/n_2 \not\rightarrow 1$ ), we propose a *nonconvex* yet tractable extension of the above convex program, for which we establish minimax optimal estimation rates under the condition  $\gamma/\sigma \gtrsim 1$ .

1) *Consistent Estimation with Balanced Samples*: For the stochastic noise setting, while one can use the same  $\ell_1$  constraint as we do in arbitrary noise case, it turns out that the analysis is more natural by considering a Lagrangian formulation. In particular, much like in standard regression, the independence assumption on  $\{e_i\}$  makes the least-squares objective analytically convenient. We therefore consider the following formulation, regularizing the squared loss objective with the nuclear norm of  $\mathbf{K}$ :

$$\min_{\mathbf{K}, \mathbf{g}} \sum_{i=1}^n (-\langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{K} \rangle + 2y_i \langle \mathbf{x}_i, \mathbf{g} \rangle - y_i^2 + \sigma^2)^2 + \lambda \|\mathbf{K}\|_*.$$
(7)

We assume that the noise variance  $\sigma^2$  is known and can be estimated.<sup>4</sup>

As with the arbitrary noise case, we first provide a theorem that bounds the distance between the optimal solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of the above program and the true  $(\mathbf{K}^*, \mathbf{g}^*)$ , and then a companion theorem gives error bounds on estimating  $(\beta_1^*, \beta_2^*)$ .

**Theorem 3** (Stochastic noise with nearly balanced samples). *For any constant  $0 < c < 2$ , there exist numerical positive constants  $\{c_i\}_{i=0}^5$ , which might depend on  $c$ , such that the following hold. Assume that  $\max\left\{\frac{n_1}{n_2}, \frac{n_2}{n_1}\right\} \leq c_0$ . Suppose moreover that (i)  $\alpha \geq c$ , (ii)  $\min\{n_1, n_2\} \geq c_4 p$ , and (iii)  $\{\mathbf{x}_i\}$  are Gaussian. For any tuning parameter  $\lambda$  that satisfies*

$$\lambda \geq c_5 \sigma (\|\beta_1^*\|_2 + \|\beta_2^*\|_2 + \sigma) (\sqrt{np} + |n_1 - n_2| \sqrt{p}) \log^3 n,$$

with probability at least  $1 - c_1 n^{-c_2}$ , any optimal solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  to the convex program (7) satisfies the error bounds

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}^*\|_F &\leq c_3 \frac{1}{n} \lambda, \\ \|\hat{\mathbf{g}} - \mathbf{g}^*\|_2 &\leq c_3 \frac{1}{n (\|\beta_1^*\| + \|\beta_2^*\| + \sigma)} \lambda. \end{aligned}$$

The bounds in the above theorem depend on the sample size difference  $|n_1 - n_2|$ . This dependence appears as a result of the objective function in the formulation (7) and is not an

<sup>4</sup>We note that similar assumptions are made in the paper [12]. It is possible to avoid the dependence on  $\sigma$  by using a symmetrized error term in the objective of (7) (see, e.g., [6]).

artifact of our analysis.<sup>5</sup> We later address how to correct for this dependence and handle imbalanced samples. Nevertheless, in the setting where the samples from the two components are approximately balanced in size with  $|n_1 - n_2|$  small, the above result implies consistency with optimal convergence rate. In this case, running Algorithm 1 on the optimal solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of the program (7) to estimate  $\theta^* = (\beta_1^*, \beta_2^*)$ , we have the following guarantees.

**Theorem 4** (Estimating  $\beta_b^*$ , stochastic noise and nearly balanced samples). *Suppose that  $|n_1 - n_2| = O(\sqrt{n \log n})$ , the conditions (i)–(iii) in Theorem 3 hold,  $\lambda \asymp \sigma (\|\beta_1^*\| + \|\beta_2^*\| + \sigma) \sqrt{np} \log^3 n$ , and  $n \geq c_3 p \log^8 n$ . Then with probability at least  $1 - c_1 n^{-c_2}$ , the output  $\hat{\theta} = (\hat{\beta}_1, \hat{\beta}_2)$  of Algorithm 1 satisfies the error bound*

$$\begin{aligned} \rho(\hat{\theta}, \theta^*) &\leq c_4 \sigma \sqrt{\frac{p}{n}} \log^4 n \\ &\quad + c_4 \min \left\{ \frac{\sigma^2}{\|\beta_1^*\|_2 + \|\beta_2^*\|_2} \sqrt{\frac{p}{n}}, \sigma \left( \frac{p}{n} \right)^{\frac{1}{4}} \right\} \log^4 n. \end{aligned}$$

The error bound has three terms, which are proportional to  $\sigma \sqrt{\frac{p}{n}}$ ,  $\frac{\sigma^2}{\|\beta_b^*\|_2} \sqrt{\frac{p}{n}}$  and  $\sigma \left( \frac{p}{n} \right)^{1/4}$ , respectively (ignoring log factors). We show that these three terms match well with the information-theoretic lower bounds given in Section III-D. They represent three phases of the error rate under different SNR; we discuss further in Section III-D.

a) *Discussion of Assumptions*:: The theoretical results above assume Gaussian covariate distribution. This Gaussianity assumption can be relaxed, but using our analysis, it comes at a cost in terms of convergence rate (and hence sample complexity required for bounded error). It can be shown that  $n = \tilde{O}(p\sqrt{p})$  suffices under a general sub-Gaussian assumption on the covariates. We believe that this additional cost is an artifact of our analysis.

2) *Consistent Estimation with Imbalanced Labels*: Now we turn to the general imbalanced setting, where the sample sizes of the two components  $n_1$  and  $n_2$  may be different. To account for the effect of imbalanced samples, we consider the following “corrected” version of the optimization problem (7):

$$\begin{aligned} \min_{\mathbf{K}, \mathbf{g}} \quad & \sum_{i=1}^n (-\langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{K} \rangle + 2y_i \langle \mathbf{x}_i, \mathbf{g} \rangle - y_i^2 + \sigma^2)^2 \\ & - \sum_{i=1}^n 4\sigma^2 (y_i - \langle \mathbf{x}_i, \mathbf{g} \rangle)^2 \\ \text{s.t.} \quad & \|\mathbf{K}\|_* \leq \|\mathbf{K}^*\|_*. \end{aligned} \tag{8}$$

Compared to the previous formulation (7), we replace the nuclear norm penalty in the objective function with its constraint version, which again serves as a surrogate for the low rank structure of  $\mathbf{K}^*$ . More crucially, we add a negative term in the objective function of (8), which corrects for the impact of imbalanced samples on the resulting estimate. This formulation is based on the following intuition: Minimizing the first term results in a solution  $\hat{\mathbf{g}}$  that is biased towards  $\beta_1^*$

<sup>5</sup>Intuitively, if the majority of the observations are generated by one of the  $\beta_b^*$ , then the objective produces a solution biased toward this  $\beta_b^*$  since this solution fits more observations. In Section III-C2, we compensate for such bias by optimizing a modified objective.

when  $n_1 > n_2$ , whereas the second negative term increases when  $\mathbf{g}$  gets closer to  $\beta_1^*$ . Simple calculation shows that in expectation these two effects cancel out. Therefore, jointly minimizing the two terms produces a solution  $\hat{\mathbf{g}}$  that is a consistent estimator of  $(\beta_1^* + \beta_2^*)/2$  even when  $n_1$  and  $n_2$  are imbalanced.

The negative term in the objective function of (8) makes the optimization problem no longer convex. Nevertheless, one can still apply *projected gradient descent* to the nonconvex program (8). To specify the projected gradient descent iteration and ease notation, we use the following shorthand for the objective function:

$$\mathcal{L}_n(\mathbf{K}, \mathbf{g}) := \sum_{i=1}^n \left( -\langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{K} \rangle + 2y_i \langle \mathbf{x}_i, \mathbf{g} \rangle - y_i^2 + \sigma^2 \right)^2 - \sum_{i=1}^n 4\sigma^2 (y_i - \langle \mathbf{x}_i, \mathbf{g} \rangle)^2.$$

Let  $\nabla_{\mathbf{K}} \mathcal{L}_n$  denote the partial gradient of  $\mathcal{L}_n$  over  $\mathbf{K}$ , and  $\nabla_{\mathbf{g}} \mathcal{L}_n$  the partial gradient of  $\mathcal{L}_n$  over  $\mathbf{g}$ . Given any feasible initializer  $(\mathbf{K}^{(0)}, \mathbf{g}^{(0)})$  with  $\|\mathbf{K}^{(0)}\|_* \leq \|\mathbf{K}^*\|_*$ , the projected gradient descent algorithm with step size  $\eta$  is given by the update:

$$\begin{aligned} & (\mathbf{K}^{(t+1)}, \mathbf{g}^{(t+1)}) \\ & \leftarrow \arg \min_{\|\mathbf{K}\|_* \leq \|\mathbf{K}^*\|_*} \left\{ \langle \nabla_{\mathbf{K}} \mathcal{L}_n^{(t)}, \mathbf{K} \rangle + \langle \nabla_{\mathbf{g}} \mathcal{L}_n^{(t)}, \mathbf{g} \rangle \right. \\ & \quad \left. + \frac{\eta}{2} \|\mathbf{K} - \mathbf{K}^{(t)}\|_F^2 + \frac{\eta\gamma^2}{2} \|\mathbf{g} - \mathbf{g}^{(t)}\|_2^2 \right\}, \end{aligned} \quad (9)$$

where  $\mathcal{L}_n^{(t)} := \mathcal{L}_n(\mathbf{K}^{(t)}, \mathbf{g}^{(t)})$ . Recall that the quantity  $\gamma$  is a measure of the signal strength. We need  $\gamma$  in the algorithm because the smoothness constants of  $\mathcal{L}_n(\mathbf{K}, \mathbf{g})$  with respect to  $\mathbf{K}$  and  $\mathbf{g}$  differ by a factor of  $\gamma^2$ . The minimization in (9) can be computed by two simple steps: (1) moving  $\mathbf{K}$  and  $\mathbf{g}$  in the negative gradient direction with step sizes  $1/\eta$  and  $1/(\eta\gamma^2)$ , respectively; (2) projecting the new  $\mathbf{K}$  onto the nuclear norm ball. The second projection step can be done with a singular value decomposition (SVD) followed by shrinking the singular values (see, e.g., [1]).

The optimization formulation (8) is non-convex. Nevertheless, the projected gradient descent algorithm still converges to a statistically accurate solution, as we show below. The intuition is that the gradient descent iterates will stay in the directions along which the objective function is convex-like. In particular, in addition to the statistical error (how close the optimal solution of (8) is to the ground-truth quantity), we can control the optimization error — how close we can get to the optimal solution of (8). Accordingly, Theorem 5 below consists of two parts. Part (a) is an analogue to Theorem 3, and bounds the *statistical error*: the distance between the global optimum  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of (8) and the true pair  $(\mathbf{K}^*, \mathbf{g}^*)$ . Part (b) bounds the *optimization error* of the gradient descent iteration (9), that is, the distance between the iterates  $(\mathbf{K}^{(t)}, \mathbf{g}^{(t)})$  and the global optimum  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$ .

**Theorem 5** (Stochastic noise with imbalanced samples). *Suppose the conditions (i)–(iii) in Theorem 3 hold for some suitable constants. The following results hold for any  $(n_1, n_2)$  that satisfies  $n_1/n_2 = \Theta(1)$ .*

(a) *(Statistical error) There exist positive constants  $\{c_i\}_{i=1}^4$  such that if  $\gamma/\sigma \geq c_1$ , then with probability at least  $1 - c_3 n^{-c_4}$ , the global optimum  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of program (8) satisfies the bounds*

$$\begin{aligned} \|\hat{\mathbf{K}} - \mathbf{K}^*\|_F & \leq c_2 \gamma \sigma \sqrt{\frac{p}{n}} \log^3 n, \\ \|\hat{\mathbf{g}} - \mathbf{g}^*\|_2 & \leq c_2 \sigma \sqrt{\frac{p}{n}} \log^3 n. \end{aligned}$$

(b) *(Optimization error) Let  $\mathbf{H}^{(t)} := \mathbf{K}^{(t)} - \hat{\mathbf{K}}$  and  $\mathbf{h}^{(t)} := \mathbf{g}^{(t)} - \hat{\mathbf{g}}$  be the optimization error terms at step  $t$  of the projected gradient descent algorithm (9) for the program (8). There exist positive constants  $\{c_i\}_{i=1}^4$  such that if  $\gamma/\sigma \geq c_1$  and  $\eta \geq c_2 n (\sqrt{p})^3 \log n$ , then with probability at least  $1 - c_3 n^{-c_4}$ , there holds the bound*

$$\begin{aligned} & \|\mathbf{H}^{(t)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 \\ & \leq \left(1 - \frac{c_3 n}{\eta}\right)^t \left( \|\mathbf{H}^{(0)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(0)}\|_2^2 \right) + c_4 \|\hat{\mathbf{K}} - \mathbf{K}^*\|_F^2 \end{aligned}$$

for every  $t = 1, 2, \dots$

Part (b) establishes geometric convergence of the optimization error. Thus, the iterate  $(\mathbf{K}^{(t)}, \mathbf{g}^{(t)})$  quickly converges to some solution whose distance to  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  is of the same order as the statistical error. Therefore, the distance between  $(\mathbf{K}^{(t)}, \mathbf{g}^{(t)})$  and  $(\mathbf{K}^*, \mathbf{g}^*)$  satisfies the bound in part (a) when  $t$  is sufficiently large by the triangle inequality. If we choose the step size  $\eta \asymp n (\sqrt{p})^3 \log n$ , the contractive factor is roughly  $1 - 1/((\sqrt{p})^3 \log n)$ , hence  $T \asymp (\sqrt{p})^3 \log n$  iterations of projected gradient descent suffice. We can then apply Algorithm 1 to the  $T^{th}$  iterate  $(\mathbf{K}^{(T)}, \mathbf{g}^{(T)})$  to obtain an accurate estimate of the true regression vectors  $\theta^* = (\beta_1^*, \beta_2^*)$ , as shown in the theorem below.

**Theorem 6** (Estimating  $\beta_b^*$ , stochastic noise with imbalanced samples). *Suppose that the projected gradient descent in (9) is initialized with  $\mathbf{K}^{(0)} = \mathbf{0}$  and  $\mathbf{g}^{(0)} = \mathbf{0}$ . Under the setting of Theorem 5, there exist positive constants  $\{c_i\}_{i=1}^7$  such that if  $n \geq c_1 p \log^6 n$ ,  $\eta \geq c_2 n (\sqrt{p})^3 \log n$  and*

$$T \geq c_3 \eta n^{-1} \max \left\{ 1 - c_4 \sigma^2 \gamma^{-2} \frac{p}{n} \log^6 n, 0 \right\},$$

*then with probability at least  $1 - c_5 n^{-c_6}$ , the output  $\hat{\theta} = (\hat{\beta}_1, \hat{\beta}_2)$  of Algorithm 1 applied to the input  $(\mathbf{K}^{(T)}, \mathbf{g}^{(T)})$  satisfies the error bound*

$$\rho(\hat{\theta}, \theta^*) \leq c_7 \sigma \sqrt{\frac{p}{n}} \log^3 n.$$

**Remark 1** (Scalability). *Both the formulations (5)–(6) and (7) can be cast as Semidefinite Programs (SDP). In the arbitrary noise setting, the constraint in the convex program (5)–(6) can be rewritten as a collection of linear constraints through the standard transformation of convex  $\ell_1$  constraints. The Lagrangian formulation (7) in the setting of stochastic noise,*

involves minimizing the sum of a trace norm term and a smooth quadratic term. The computational complexity of solving this regularized quadratic optimization in the matrix space has similar complexity to problems such as matrix completion [9] and PhaseLift [10], and various first order methods can easily be adapted, thus allowing solution of large scale instances of the mixed regression problem.

#### D. Minimax Lower Bounds

We now derive minimax lower bounds on the estimation errors for both the arbitrary and stochastic noise settings, and show that these match our upper bounds. Recall the error (semi)-metric  $\rho(\boldsymbol{\theta}, \boldsymbol{\theta}')$  defined in equation (2) for a pair  $\boldsymbol{\theta} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2)$  and  $\boldsymbol{\theta}' = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)$  in  $\mathbb{R}^p \times \mathbb{R}^p$ . It is straightforward to verify that the metric  $\rho(\cdot, \cdot)$  satisfies the triangle inequality. An estimator

$$\hat{\boldsymbol{\theta}} \equiv \hat{\boldsymbol{\theta}}(\mathbf{X}, \mathbf{y}) = (\hat{\boldsymbol{\beta}}_1(\mathbf{X}, \mathbf{y}), \hat{\boldsymbol{\beta}}_2(\mathbf{X}, \mathbf{y})).$$

of the true regression vectors  $\boldsymbol{\theta}^* = (\boldsymbol{\beta}_1^*, \boldsymbol{\beta}_2^*)$  is any measurable function of the observed data  $(\mathbf{X}, \mathbf{y})$ . For each number  $\underline{\gamma} > 0$ , we consider the following parameter class

$$\Theta(\underline{\gamma}) := \left\{ \boldsymbol{\theta} = (\boldsymbol{\beta}_1, \boldsymbol{\beta}_2) \in \mathbb{R}^p \times \mathbb{R}^p : 2 \|\boldsymbol{\beta}_1 - \boldsymbol{\beta}_2\|_2 \geq \|\boldsymbol{\beta}_1\|_2 + \|\boldsymbol{\beta}_2\|_2 \geq \underline{\gamma} \right\}, \quad (10)$$

that is, pairs of regression vectors whose norms and separation are lower bounded. Recall that  $z_i \in \{0, 1\}$  is the hidden label for the  $i$ -th observation, that is,  $z_i = 0$  if and only if  $i \in \mathcal{I}_1$ , for each  $i = 1, 2, \dots, n$ .

1) *Lower Bounds for Arbitrary Noise*: In the arbitrary noise setting, the noise vector  $\mathbf{e}$  is assumed to lie in the  $\ell_2$ -ball  $\mathbb{B}(\epsilon) := \{\mathbf{u} \in \mathbb{R}^n : \|\mathbf{u}\|_2 \leq \epsilon\}$  and otherwise arbitrary. For this setting we have the following minimax lower bounds.

**Theorem 7** (Lower bound, arbitrary noise). *There exist universal positive constants  $c_0, c_1$  for which the following is true. If  $n \geq c_1 p$ , then for any  $\underline{\gamma} > 0$  and any hidden label vector  $\mathbf{z} \in \{0, 1\}^n$ , there holds the bound*

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \sup_{\mathbf{e} \in \mathbb{B}(\epsilon)} \rho(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \geq c_0 \frac{\epsilon}{\sqrt{n}} \quad (11)$$

with probability at least  $1 - n^{-10}$ , where the probability is with respect to the randomness in  $\mathbf{X}$ .

The lower bound above matches the upper bound given in Theorem 2, thus showing that our convex formulation is minimax optimal and order-wise unimprovable. Informally, Theorems 2 and 7 together establish the following minimax rate of the arbitrary noise setting

$$\rho(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \asymp \frac{\|\mathbf{e}\|_2}{\sqrt{n}},$$

which is valid for  $n \gtrsim p$ .

2) *Lower Bounds for Stochastic Noise*: For the stochastic setting where the noise is i.i.d. Gaussian, we further assume that the two components have equal mixing weights:  $\mathbb{P}(z_i = 0) = \mathbb{P}(z_i = 1) = 1/2$  for each  $i = 1, 2, \dots, n$ . For this setting we have the following minimax lower bounds.

**Theorem 8** (Lower bound, stochastic noise). *Suppose that  $n \geq p \geq 64$ ,  $\mathbf{X} \in \mathbb{R}^{n \times p}$  has i.i.d. standard Gaussian entries,  $\mathbf{e}$  has i.i.d. zero-mean Gaussian entries with variance  $\sigma^2$ , and  $z_i \stackrel{iid}{\sim} \text{Bernoulli}(1/2)$ . The following statements hold for some absolute constants  $0 < c_0, c_1 < 1$ .*

1) *For each  $\underline{\gamma} > \sigma$ , we have*

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\rho(\boldsymbol{\theta}^*, \hat{\boldsymbol{\theta}})] \geq c_0 \sigma \sqrt{\frac{p}{n}}. \quad (12)$$

2) *For each  $c_1 \sigma \left(\frac{p}{n}\right)^{1/4} \leq \underline{\gamma} \leq \sigma$ , we have*

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\rho(\boldsymbol{\theta}^*, \hat{\boldsymbol{\theta}})] \geq c_0 \frac{\sigma^2}{\underline{\gamma}} \sqrt{\frac{p}{n}}. \quad (13)$$

3) *For each  $0 < \underline{\gamma} \leq c_1 \sigma \left(\frac{p}{n}\right)^{1/4}$ , we have*

$$\inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\rho(\boldsymbol{\theta}^*, \hat{\boldsymbol{\theta}})] \geq c_0 \sigma \left(\frac{p}{n}\right)^{1/4}. \quad (14)$$

Here  $\mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\cdot]$  denotes the expectation with respect to the covariate matrix  $\mathbf{X}$ , the hidden label vector  $\mathbf{z}$  and the noise vector  $\mathbf{e}$ .

We see that the three lower bounds in the above theorem match each of the three terms in the upper bound given in Theorem 4 up to at most polylog factors, proving the minimax optimality of the error bounds of our convex formulation. Informally, Theorems 4 and 8 together establish the following minimax error rate (up to a polylog factor) in the stochastic noise setting:

$$\rho(\boldsymbol{\theta}^*, \hat{\boldsymbol{\theta}}) \asymp \begin{cases} \sigma \sqrt{\frac{p}{n}}, & \text{if } \underline{\gamma} \gtrsim \sigma, \\ \frac{\sigma^2}{\underline{\gamma}} \sqrt{\frac{p}{n}}, & \text{if } \sigma \left(\frac{p}{n}\right)^{1/4} \lesssim \underline{\gamma} \lesssim \sigma, \\ \sigma \left(\frac{p}{n}\right)^{1/4}, & \text{if } \underline{\gamma} \lesssim \sigma \left(\frac{p}{n}\right)^{1/4}. \end{cases}$$

Here,  $\underline{\gamma}$  is any lower bound on  $\|\boldsymbol{\beta}_1^*\|_2 + \|\boldsymbol{\beta}_2^*\|_2$  and represents the signal strength (recall the definition of the parameter class in equation (10)). Notice how the scaling of the minimax error rate exhibits three phases depending on the SNR  $\underline{\gamma}/\sigma$ : (i) In the high SNR regime with  $\underline{\gamma} \gtrsim \sigma$ , we see a fast rate — proportional to  $1/\sqrt{n}$  — that is dominated by the error of estimating a single  $\boldsymbol{\beta}_b^*$  and is the same as the rate for standard linear regression. (ii) In the low SNR regime with  $\underline{\gamma} \lesssim \sigma \left(\frac{p}{n}\right)^{1/4}$ , we have a slow rate that is proportional to  $1/n^{1/4}$ , which is associated with the demixing of the two components  $\boldsymbol{\beta}_1^*$  and  $\boldsymbol{\beta}_2^*$ . (iii) In the medium SNR regime, the error rate transitions between the fast and slow phases, and depends in a precise way on the SNR. For related phenomena, see the work in [3, 14].

The lower bounds in Theorem 8 apply to the balanced sample setting. The error upper bound in Theorem 6 for our nonconvex approach in the imbalanced setting, is also near-optimal in the high SNR regime, as this bound matches (up

to logarithmic factors) the minimax lower bound of standard linear regression for estimating a single  $\beta_b^*$  (that is, when the labels are known).

### E. Implications for Phase Retrieval

As an illustration of the power of our results, we discuss an application to the *Phase Retrieval* problem, which has recently received much attention (see, e.g., the work in [6, 7, 11, 16, 23, 15]). Recall that in the real value setting, the phase retrieval problem is essentially a regression problem without sign information. Recent work has mostly focused on the noiseless case. Here, the problem is as follows: we observe  $(\mathbf{x}_i, y_i) \in \mathbb{R}^p \times \mathbb{R}$ ,  $i = 1, 2, \dots, n$ , where

$$y_i = |\mathbf{x}_i^\top \beta^*|.$$

The goal is to recover the unknown vector  $\beta^* \in \mathbb{R}^p$ . The stability of recovery algorithms has also been considered. Most work has focused on the setting where noise is added to the phase-less measurements, that is,

$$y_i = |\mathbf{x}_i^\top \beta^*| + e_i. \quad (15)$$

In many applications, however, it is also natural to consider the setting where the measurement noise is added before the phase is lost. This corresponds to the model

$$y_i = |\mathbf{x}_i^\top \beta^* + e_i|. \quad (16)$$

We may call formulation (16) the *Noisy Phase Model*, as opposed to the *Noisy Magnitude Model* in (15) considered by previous work. This problem can be reduced to a mixed regression problem and solved by our algorithm. The reduction is as follows. We generate  $n$  independent Rademacher random variables  $\epsilon_i, i = 1, \dots, n$ . For each  $i$ , we set  $y'_i = \epsilon_i y_i$ . Let  $s_i := \text{sign}(\mathbf{x}_i^\top \beta^* + e_i)$  and  $e'_i := \epsilon_i s_i e_i$ , where we use the convention that  $\text{sign}(0) = 1$ . Under the noisy phase model (16), we then have

$$y'_i = \epsilon_i \cdot |\mathbf{x}_i^\top \beta^* + e_i| = \epsilon_i s_i (\mathbf{x}_i^\top \beta^* + e_i) = \mathbf{x}_i^\top (\epsilon_i s_i \beta^*) + e'_i.$$

If we let  $\beta_1^* = \beta^*$ ,  $\beta_2^* = -\beta^*$ ,  $\mathcal{I}_1 = \{i : \epsilon_i s_i = 1\}$  and  $\mathcal{I}_2 = \{i : \epsilon_i s_i = -1\}$ , then the model becomes

$$y'_i = \mathbf{x}_i^\top \beta_b^* + e'_i, \quad \forall i \in \mathcal{I}_b,$$

which is precisely the mixed regression model we considered.

Note that with probability at least  $1 - n^{-3}$ ,  $\frac{n}{2} - \sqrt{10n \log n} \leq n_b = |\mathcal{I}_b| \leq \frac{n}{2} + \sqrt{10n \log n}$  for  $b = 1, 2$ , so  $|n_1 - n_2| = O(\sqrt{n \log n})$ . Also note that  $\|e'\|_2 = \|e\|_2$ . Conditioned on  $\{\mathcal{I}_b\}$ , the distribution of  $\{\mathbf{x}_i\}$  is the same as its unconditional distribution. Therefore, applying our arbitrary-noise result from Theorem 2, we immediately get the following guarantees for phase retrieval under the Noisy Phase Model.

**Corollary 2** (Phase retrieval, arbitrary noise). *Consider the Noisy Phase Model in (16). Suppose that the  $\{\mathbf{x}_i\}$  are i.i.d., zero-mean sub-Gaussian with bounded sub-Gaussian norm, unit variance and fourth moment  $\mu > 1$ , and that  $n \gtrsim p$ ,  $n \asymp c_4 \sqrt{n} \|e\|_2 \|\beta^*\|_2$  and the noise is arbitrary but bounded in magnitude:  $\|e\|_2 \lesssim \sqrt{n} \|\beta^*\|_2$ . Then using the reduction*

*described above, the output of the program (5)–(6) followed by Algorithm 1 satisfies the error bound*

$$\min_{b=1,2} \|\hat{\beta}_b - \beta^*\|_2 \lesssim \frac{\|e\|_2}{\sqrt{n}}$$

*with probability at least  $1 - n^{-2}$ .*

In the corollary above we assumed  $\|e\|_2 \lesssim \sqrt{n} \|\beta^*\|_2$ . Similarly as before, the large noise case with  $\|e\|_2 \geq c_4 \sqrt{n} \|\beta^*\|_2$  can be handled easily, using the blind estimator  $\hat{\beta} := \min_{\beta} \sum_{i \in [n]} |\mathbf{x}_i^\top \beta - y_i|$ , which in this case again satisfies the optimal error bound  $\|\hat{\beta} - \beta^*\|_2 \lesssim \frac{\|e\|_2}{\sqrt{n}}$ .

Next, consider the stochastic noise case where  $e_i$  is i.i.d., zero-mean symmetric sub-Gaussian with variance  $\sigma^2$ . Conditioned on  $\{\mathcal{I}_b\}$ , the conditional distributions of  $\{e'_i\}$  and  $\{\mathbf{x}_i\}$  inherit the properties of  $e_i$  and the unconditional  $\mathbf{x}_i$ , and are independent of each other. Applying Theorem 4, we have the following guarantee.

**Corollary 3** (Phase retrieval, stochastic noise). *Consider the Noisy Phase Model in (16). Suppose that the  $\{\mathbf{x}_i\}$  are i.i.d., zero-mean Gaussian with unit variance, and that the noise  $e_i$  is i.i.d., zero-mean symmetric sub-Gaussian with sub-Gaussian norm bounded by  $c_3 \sigma$  and variance equal to  $\sigma^2$ . Suppose further that  $n \gtrsim p$  and  $\lambda \asymp \sigma (\|\beta^*\|_2 + \sigma) \sqrt{np} \log^4 n$ . Then using the reduction described above, the output of the program (7) followed by Algorithm 1 satisfies the error bound*

$$\begin{aligned} \min_{b=1,2} \|\hat{\beta}_b - \beta^*\|_2 \\ \lesssim \sigma \sqrt{\frac{p}{n} \log^4 n} + \min \left\{ \frac{\sigma^2 \sqrt{\frac{p}{n}}}{\|\beta^*\|_2}, \sigma \left( \frac{p}{n} \right)^{\frac{1}{4}} \right\} \log^4 n \end{aligned}$$

*with probability at least  $1 - n^{-2}$ .*

As a passing observation, we note that the error bounds above for phase retrieval are both order-wise optimal. For the deterministic noise setting considered in Corollary 2, we cannot achieve a smaller error even if the phase is not lost. For the stochastic noise setting, note that our corresponding minimax lower bounds for mixed regression in Section III-D2 is in fact derived under the symmetric setting  $\beta_1^* = -\beta_2^*$ . In this case one can reduce a mixed regression problem to a phase retrieval problem by dropping the signs in  $\{y_i\}$ , hence our minimax lower bounds certify the near-optimality of Corollary 3 for phase retrieval.

## IV. PROOFS

We now provide the proofs of the main Theorems. Conceptually, there are three parts to our results. *i) Regression error* (Theorems 1, 3, 5). We prove that the lifted optimizations in the matrix space (for deterministic noise, and balanced and unbalanced stochastic noise) recover good approximations to  $\mathbf{K}^*$  and  $\mathbf{g}^*$ . We note that these results have no dependence on SNR; that is, alone, they do not reveal a change in the rate of convergence. *ii) Decomposition and perturbation error* (Theorems 2, 4, 6). We prove a matrix perturbation result that shows that a good approximation of  $\mathbf{K}^*$  and  $\mathbf{g}^*$  results in a

good approximation of  $\beta_1^*$  and  $\beta_2^*$ . The change in rate in the low-SNR setting, emerges from these decomposition results.

iii) *Minimax lower bounds* (Theorems 7 and 8). We prove minimax (information theoretic) lower bounds matching the upper bounds. For the lower bounds, the change in rate comes from a detection problem that involves distinguishing mixtures of Gaussians with different separation of the component centers.

We prove these results here, and in the interest of readability, defer the proofs of technical lemmas to the appendix.

### A. Notation and Preliminaries

We use  $\beta_{-b}^*$  to denote  $\beta_2^*$  if  $b = 1$  and  $\beta_1^*$  if  $b = 2$ . Let  $\delta_b^* := \beta_b^* - \beta_{-b}^*$ . Without loss of generality, we assume  $\mathcal{I}_1 = \{1, \dots, n_1\}$  and  $\mathcal{I}_2 = \{n_1 + 1, \dots, n\}$ . For  $i = 1, \dots, n_1$ , we define  $\mathbf{x}_{1,i} := \mathbf{x}_i$ ,  $y_{1,i} = y_i$  and  $e_{1,i} = e_i$ ; correspondingly, for  $i = 1, \dots, n_2$ , we define  $\mathbf{x}_{2,i} := \mathbf{x}_{n_1+i}$ ,  $y_{2,i} := y_{n_1+i}$  and  $e_{2,n+i}$ . For each  $b = 1, 2$ , let  $\mathbf{X}_b \in \mathbb{R}^{n_b \times p}$  be the matrix with rows  $\{\mathbf{x}_{b,i}^\top, i = 1, \dots, n_b\}$ . Also let  $\mathbf{e}_b := [e_{b,1} \cdots e_{b,n_b}]^\top \in \mathbb{R}^{n_b}$ .

While the measurements in the original model are given by  $\mathbf{X}$ , in the lifted space, one can regard the measurements as given by rank-one matrices that are quadratic in  $\mathbf{x}_{b,i}$ . Thus it is natural to define the matrices  $\mathbf{A}_{b,i} := \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top$ ,  $i \in [n_b]$  and the mapping  $\mathcal{A}_b : \mathbb{R}^{p \times p} \mapsto \mathbb{R}^{n_b}$  given by

$$(\mathcal{A}_b \mathbf{Z})_i = \frac{1}{n_b} \langle \mathbf{A}_{b,i}, \mathbf{Z} \rangle, \quad \text{for each } i \in [n_b].$$

Because of their quadratic nature, these measurements are not mean-zero. As we detail below, our proofs rely on establishing a restricted isometry-like property of the measurements, but as the measurements are not zero mean, this does not hold. It is convenient, therefore, to define matrices and a mapping that are related to  $\mathbf{A}_{b,i}$  and  $\mathcal{A}_b$ , but with zero mean. To this end, for  $b = 1, 2$  and  $j = 1, \dots, \lfloor n_b/2 \rfloor$ , define the matrix  $\mathbf{B}_{b,j} := \mathbf{x}_{b,2j} \mathbf{x}_{b,2j}^\top - \mathbf{x}_{b,2j-1} \mathbf{x}_{b,2j-1}^\top$ , as well as the vector  $\mathbf{d}_{b,j} = e_{b,2j} \mathbf{x}_{b,2j} - e_{b,2j-1} \mathbf{x}_{b,2j-1}$ . For  $b \in \{1, 2\}$ , define the mapping  $\mathcal{B}_b : \mathbb{R}^{p \times p} \mapsto \mathbb{R}^{\lfloor n_b/2 \rfloor}$  by

$$(\mathcal{B}_b \mathbf{Z})_j = \frac{1}{\lfloor n_b/2 \rfloor} \langle \mathbf{B}_{b,j}, \mathbf{Z} \rangle, \quad \text{for each } j = 1, \dots, \lfloor n_b \rfloor.$$

Since  $y_{b,i} = \mathbf{x}_{b,i}^\top \beta_b^* + e_{b,i}$ ,  $i \in [n_b]$ , we have for any  $\mathbf{Z} \in \mathbb{R}^{p \times p}$ ,  $\mathbf{z} \in \mathbb{R}^p$  and for all  $j = 1, \dots, \lfloor n_b \rfloor$ ,

$$\begin{aligned} & \frac{1}{\lfloor n_b/2 \rfloor} (\langle \mathbf{B}_{b,j}, \mathbf{Z} \rangle - 2\mathbf{d}_{b,j}^\top \mathbf{z}) \\ &= \frac{1}{\lfloor n_b/2 \rfloor} \langle \mathbf{B}_{b,j}, \mathbf{Z} - 2\beta_b^* \mathbf{z}^\top \rangle + (e_{b,2j} \mathbf{x}_{b,2j} - e_{b,2j-1} \mathbf{x}_{b,2j-1})^\top \mathbf{z} \\ &= (\mathcal{B}_b (\mathbf{Z} - 2\beta_b^* \mathbf{z}^\top))_j + (e_{b,2j} \mathbf{x}_{b,2j} - e_{b,2j-1} \mathbf{x}_{b,2j-1})^\top \mathbf{z}. \end{aligned}$$

A key part of the proof is in expressing the error in terms of the operators  $\mathcal{A}_b$  and then  $\mathcal{B}_b$ , and then showing that  $\mathcal{B}_b$  satisfies a restricted isometry property. Also key and common to all the proofs, is to show that the optimization formulations recover a near low-rank matrix  $\hat{\mathbf{K}}$ . For this, we need to control the part of  $\hat{\mathbf{K}}$  that has different column and row space from  $\mathbf{K}^*$ . The following notation and definitions are standard. Let the rank-2 SVD of  $\mathbf{K}^*$  be  $\mathbf{U} \Sigma \mathbf{V}^\top$ . Note that  $\mathbf{U}$  and  $\mathbf{V}$  have the same column space, which equals  $\text{span}(\beta_1^*, \beta_2^*)$ .

Define the projection matrix  $\mathbf{P}_U := \mathbf{U} \mathbf{U}^\top = \mathbf{V} \mathbf{V}^\top$  and the subspace  $T := \{\mathbf{P}_U \mathbf{Z} + \mathbf{Y} \mathbf{P}_U : \mathbf{Z}, \mathbf{Y} \in \mathbb{R}^{p \times p}\}$ . Let  $T^\perp$  be the orthogonal subspace of  $T$ . The projections to  $T$  and  $T^\perp$  are given by

$$\mathcal{P}_T \mathbf{Z} := \mathbf{P}_U \mathbf{Z} + \mathbf{Z} \mathbf{P}_U - \mathbf{P}_U \mathbf{Z} \mathbf{P}_U, \quad \mathcal{P}_{T^\perp} \mathbf{Z} := \mathbf{Z} - \mathcal{P}_T \mathbf{Z}.$$

Denote the optimal solution to the optimization problem of interest (either (5) or (7)) as  $(\hat{\mathbf{K}}, \hat{\mathbf{g}}) = (\mathbf{K}^* + \hat{\mathbf{H}}, \mathbf{g}^* + \hat{\mathbf{h}})$ . Let  $\hat{\mathbf{H}}_T := \mathcal{P}_T \hat{\mathbf{H}}$  and  $\hat{\mathbf{H}}_{T^\perp} := \mathcal{P}_{T^\perp} \hat{\mathbf{H}}$ .

The optimization proofs follow a similar spirit and conceptual flow. The first part of the proof asserts that the error (compared to  $(\mathbf{K}^*, \mathbf{g}^*)$ ) must satisfy certain conditions controlled by the operators  $\mathcal{B}_b$ ,  $b = 1, 2$ . For Theorem 1 this is a consequence of the constraints; for Theorem 3 this is a consequence of the objective function. Essentially these results say that the errors must lie in certain directions away from  $(\mathbf{K}^*, \mathbf{g}^*)$ .

The next step comes in using properties of  $\mathcal{B}_b$ . In particular, we show that these operators satisfy a restricted isometry-like (RIP) condition. Together with the characterization of how  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  and  $(\mathbf{K}^*, \mathbf{g}^*)$  can differ, we conclude that along all those directions, the objective function has strong convexity, that is, curvature bounded from below). This allows us to provide bounds on how far  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  can be from  $(\mathbf{K}^*, \mathbf{g}^*)$ .

### B. Proof of Theorem 1

As outlined above, the proof follows from three main steps. First, the  $\ell_1$  error term that in this formulation appears in the LHS of the constraint (6) in the optimization, is naturally related to the operators  $\mathcal{A}_b$ . Using the definitions of  $\mathcal{A}_b$ ,  $\mathcal{B}_b$ , we establish a relation between  $\mathbf{e}$ ,  $\eta$  and the feasibility of the optimal solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}}) = (\mathbf{K}^* + \hat{\mathbf{H}}, \mathbf{g}^* + \hat{\mathbf{h}})$ . Second, we show that the operator  $\mathcal{B}$  is an approximate isometry on low-rank matrices. Finally, this allows us to obtain good upper and lower bounds on the error terms, and hence the accuracy of the solution.

Now for the details. Recall that  $\hat{\mathbf{H}}_T := \mathcal{P}_T \hat{\mathbf{H}}$  and  $\hat{\mathbf{H}}_{T^\perp} := \mathcal{P}_{T^\perp} \hat{\mathbf{H}}$ . Note that  $\hat{\mathbf{H}}_T$  has rank at most 4 and  $\hat{\mathbf{H}}_{T^\perp}$  has rank at most  $p - 4$ . We have

$$\begin{aligned} \|\hat{\mathbf{K}}\|_* - \|\mathbf{K}^*\|_* &\geq \|\mathbf{K}^* + \hat{\mathbf{H}}_{T^\perp}\|_* - \|\hat{\mathbf{H}}_T\|_* - \|\mathbf{K}^*\|_* \\ &= \|\hat{\mathbf{H}}_{T^\perp}\|_* - \|\hat{\mathbf{H}}_T\|_*. \end{aligned} \quad (17)$$

1) *Step (1): Consequence of Feasibility:* For any  $(\mathbf{K}, \mathbf{g}) = (\mathbf{K}^* + \mathbf{H}, \mathbf{g}^* + \mathbf{h})$ , it is easy to check that

$$\begin{aligned} & -\langle \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top, \mathbf{K} \rangle + 2y_{b,i} \langle \mathbf{x}_{b,i}, \mathbf{g} \rangle - y_{b,i}^2 \\ &= -\langle \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top, \mathbf{H} \rangle + 2y_{b,i} \langle \mathbf{x}_{b,i}, \mathbf{h} \rangle - e_{b,i} \mathbf{x}_{b,i}^\top \delta_b^* - e_{b,i}^2. \end{aligned} \quad (18)$$

Therefore, the constraint (6) is equivalent to

$$\begin{aligned} & \sum_{b=1}^2 \sum_{i=1}^{n_b} \left| -\langle \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top, \mathbf{H} \rangle + 2(\mathbf{x}_{b,i}^\top \beta_b^* + e_{b,i}) \langle \mathbf{x}_{b,i}, \mathbf{h} \rangle \right. \\ & \quad \left. - e_{b,i} \mathbf{x}_{b,i}^\top \delta_b^* - e_{b,i}^2 \right| \leq \eta. \end{aligned}$$

Using the notation from Section IV-A, this can be rewritten as

$$\begin{aligned} \sum_b \left\| n_b \mathcal{A}_b \left( -\mathbf{H} + 2\beta_b^* \mathbf{h}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \mathbf{h}) \right\|_1 \\ - \mathbf{e}_b \circ (\mathbf{X}_b \delta_b^*) - \mathbf{e}_b^2 \right\|_1 \leq \eta, \end{aligned} \quad (19)$$

where  $\circ$  denotes the element-wise product and  $\mathbf{e}_b^2 = \mathbf{e}_b \circ \mathbf{e}_b$ .

First, note that  $\mathbf{K}^*$  and  $\mathbf{g}^*$  are feasible. By standard bounds on the spectral norm of random matrices [30], we know that with probability at least  $1 - 2 \exp(-cn_b)$ ,

$$\|\mathbf{X}_b \mathbf{z}\|_2 \lesssim \sqrt{n_b} \|\mathbf{z}\|_2, \forall \mathbf{z} \in \mathbb{R}^p.$$

We thus have

$$\begin{aligned} \left\| -\mathbf{e}_b \circ (\mathbf{X}_b \delta_b^*) - \mathbf{e}_b^2 \right\|_1 &\leq c_1 \left( \sqrt{n_b} \|\mathbf{e}_b\|_2 \|\delta_b^*\|_2 + \|\mathbf{e}\|_2^2 \right) \\ &\stackrel{(a)}{\leq} c_1 \sqrt{n_b} \|\mathbf{e}\|_2 \|\beta_1^* - \beta_2^*\|_2 \\ &\stackrel{(b)}{\leq} \eta, \end{aligned}$$

where we use the assumptions on  $\mathbf{e}$  and  $\eta$  in the steps (a) and (b), respectively. The inequality above implies that (19) holds with  $\mathbf{H} = \mathbf{0}$  and  $\mathbf{h} = \mathbf{0}$ , thus showing the feasibility of  $(\mathbf{K}^*, \mathbf{g}^*)$ .

Since  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  is feasible by assumption, combining the last two displayed equations and (19), we further have

$$\begin{aligned} \sum_b \left\| n_b \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \\ \leq \sum_b \left\| 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_1 + \sum_b \left\| -2\mathbf{e}_b \circ (\mathbf{X}_b \delta_b^*) - \mathbf{e}_b^2 \right\|_1 + \eta \\ \leq c_2 \sum_b \sqrt{n_b} \|\mathbf{e}_b\|_2 \|\hat{\mathbf{h}}\|_2 + 2\eta. \end{aligned} \quad (20)$$

Now from the definition of  $\mathcal{A}_b$  and  $\mathcal{B}_b$ , we have

$$\begin{aligned} \lfloor n_b/2 \rfloor \left\| \mathcal{B}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \\ \leq \sum_{j=1}^{\lfloor n_b/2 \rfloor} \left\| \langle \mathbf{A}_{b,2j}, -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \rangle \right\|_1 \\ + \left\| \langle \mathbf{A}_{b,2j-1}, -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \rangle \right\|_1 \\ \leq n_b \left\| \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1. \end{aligned}$$

It follows from (20) and  $n_1 \asymp n_2 \asymp n$  that

$$\sum_b n \left\| \mathcal{B}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 - c_2 \sum_b \sqrt{n} \|\mathbf{e}_b\|_2 \|\hat{\mathbf{h}}\|_2 \leq 2\eta. \quad (21)$$

This concludes Step (1) of the proof.

2) *Step (2): RIP and Lower Bounds:* The bound in (21) relates the  $\ell_1$ -norm of  $\mathcal{B}$  and  $\eta$ . Since we want a bound on the  $\ell_2$  and Frobenius norms of  $\hat{\mathbf{h}}$  and  $\hat{\mathbf{H}}$  respectively, a major step is the proof of an RIP-like property for  $\mathcal{B}$ :

**Lemma 1.** *The following holds for some numerical constants  $c, \underline{\delta}, \bar{\delta}$ . For  $b = 1, 2$ , if  $\mu > 1$  and  $n_b \geq c \rho p$ , then with probability  $1 - \exp(-n_b)$ , we have*

$$\underline{\delta} \|\mathbf{Z}\|_F \leq \|\mathcal{B}_b \mathbf{Z}\|_1 \leq \bar{\delta} \|\mathbf{Z}\|_F$$

simultaneously for all  $\mathbf{Z} \in \mathbb{R}^{p \times p}$  with  $\text{rank}(\mathbf{Z}) \leq \rho$ .

This lemma follows from the more general Lemma 5 that appears in the proof of Theorem 3, by setting  $\sigma = 0$ .

We now turn to the implications of Lemma 1, in order to get lower bounds on the term  $\left\| \mathcal{B}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1$  from the first term in (21), in terms of  $\|\hat{\mathbf{h}}\|_2$  and  $\|\hat{\mathbf{H}}\|_F$ .

Since we have proved that  $(\mathbf{K}^*, \mathbf{g}^*)$  is feasible, we have  $\|\hat{\mathbf{K}}\|_* \leq \|\mathbf{K}^*\|_*$  by optimality. It follows from (17) that

$$\left\| \hat{\mathbf{H}}_{T^\perp} \right\|_* \leq \left\| \hat{\mathbf{H}}_T \right\|_*. \quad (22)$$

Let  $K = c \frac{1}{\alpha}$  for  $c$  some numeric constant to be chosen later. We can partition  $\hat{\mathbf{H}}_{T^\perp}$  into a sum of  $M := \frac{p-4}{K}$  matrices  $\hat{\mathbf{H}}_1, \dots, \hat{\mathbf{H}}_M$  according to the SVD of  $\hat{\mathbf{H}}_{T^\perp}$ , such that  $\text{rank}(\hat{\mathbf{H}}_i) \leq K$  and the smallest singular value of  $\hat{\mathbf{H}}_i$  is larger than the largest singular value of  $\hat{\mathbf{H}}_{i+1}$  (cf. [24]). By Lemma 1, we get that for each  $b = 1, 2$ ,

$$\begin{aligned} \sum_{i=2}^M \left\| \mathcal{B}_b(\hat{\mathbf{H}}_i) \right\|_1 &\leq \bar{\delta} \sum_{i=2}^M \left\| \hat{\mathbf{H}}_i \right\|_F \leq \bar{\delta} \sum_{i=2}^M \frac{1}{\sqrt{K}} \left\| \hat{\mathbf{H}}_{i-1} \right\|_* \\ &\leq \frac{\bar{\delta}}{\sqrt{K}} \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_* \stackrel{(a)}{\leq} \frac{\bar{\delta}}{\sqrt{K}} \sqrt{4} \left\| \hat{\mathbf{H}}_T \right\|_F, \end{aligned} \quad (23)$$

where (a) follows from (22) and the rank of  $\hat{\mathbf{H}}_T$ . It follows that for  $b = 1, 2$ ,

$$\begin{aligned} &\left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \\ &\stackrel{(a)}{\geq} \left\| \mathcal{B}_b \left( \hat{\mathbf{H}}_T + \hat{\mathbf{H}}_1 - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 - \sum_{i=2}^M \left\| \mathcal{B}_b(\hat{\mathbf{H}}_i) \right\|_1 \\ &\stackrel{(b)}{\geq} \underline{\delta} \left\| \hat{\mathbf{H}}_T + \hat{\mathbf{H}}_1 - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F - 2\bar{\delta} \sqrt{\frac{1}{K}} \left\| \hat{\mathbf{H}}_T \right\|_F \\ &\stackrel{(c)}{\geq} \underline{\delta} \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F + \left\| \hat{\mathbf{H}}_1 \right\|_F - 2\bar{\delta} \sqrt{\frac{1}{K}} \left\| \hat{\mathbf{H}}_T \right\|_F \\ &\geq \underline{\delta} \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F - 2\bar{\delta} \sqrt{\frac{1}{K}} \left\| \hat{\mathbf{H}}_T \right\|_F, \end{aligned}$$

where (a) follows from the triangle inequality, (b) follows from Lemma 1 and (23), and (c) follows from the fact that  $\hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \in T$  and  $\hat{\mathbf{H}}_1 \in T^\perp$ . Summing the above inequality for  $b = 1, 2$ , we obtain

$$\begin{aligned} &\sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \\ &\geq \underline{\delta} \sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F - 4\bar{\delta} \sqrt{\frac{1}{K}} \left\| \hat{\mathbf{H}}_T \right\|_F. \end{aligned} \quad (24)$$

The first term in the RHS of (24) can be bounded using the following lemma, whose proof is deferred to Appendix F-A.

**Lemma 2.** *We have*

$$\begin{aligned} \sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F &\geq \sqrt{\alpha} \left\| \hat{\mathbf{H}}_T \right\|_F, \\ \sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F &\geq \sqrt{\alpha} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2) \|\hat{\mathbf{h}}\|_2. \end{aligned}$$

Combining (24) and the lemma, we obtain

$$\sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \geq \left( \frac{\underline{\delta}\sqrt{\alpha}}{4\bar{\delta}} - 4\bar{\delta}\sqrt{\frac{1}{K}} \right) \left\| \hat{\mathbf{H}}_T \right\|_F$$

and

$$\begin{aligned} \sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \\ \geq \left( \frac{\underline{\delta}}{4\bar{\delta}} - 4\bar{\delta}\sqrt{\frac{1}{\alpha K}} \right) \sum_b \left\| \hat{\mathbf{H}}_T - \beta_b \hat{\mathbf{h}}^\top \right\|_F \\ \geq \left( \frac{\underline{\delta}}{4\bar{\delta}} - 4\bar{\delta}\sqrt{\frac{1}{\alpha K}} \right) \sqrt{\alpha} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2) \|\hat{\mathbf{h}}\|_2. \end{aligned}$$

Recall that  $K = c\frac{1}{\alpha}$ . When  $c$  is sufficiently large, the above inequalities imply that for some numeric constant  $c'$ ,

$$\sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \geq \frac{\sqrt{\alpha}}{c''} \left\| \hat{\mathbf{H}}_T \right\|_F \stackrel{(d)}{\geq} \frac{\sqrt{\alpha}}{c'} \left\| \hat{\mathbf{H}} \right\|_F, \quad (25)$$

$$\sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \geq \frac{\sqrt{\alpha}}{c'} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2) \|\hat{\mathbf{h}}\|_2, \quad (26)$$

where the inequality (d) follows from (22) and  $\text{rank}(\hat{\mathbf{H}}_T) \leq 4$ . This concludes the proof of Step (2).

3) *Step (3): Producing Error Bounds:* We now combine the results from the above steps, in order to obtain bounds on  $\|\hat{\mathbf{h}}\|_2$  and  $\|\hat{\mathbf{H}}\|_F$  in terms of  $\eta$ , and the other parameters of the problem, hence concluding the proof of Theorem 1.

From Step (1), we concluded the bound (21), which we reproduce:

$$\sum_b n \left\| \mathcal{B}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 - c_2 \sum_b \sqrt{n} \|\mathbf{e}_b\|_2 \|\hat{\mathbf{h}}\|_2 \leq 2\eta.$$

Applying (26) to the LHS above, we get

$$\sqrt{n} \sum_b (\sqrt{\alpha} \sqrt{n} \|\beta_b^*\|_2 - \|\mathbf{e}_b\|_2) \|\hat{\mathbf{h}}\|_2 \lesssim 2\eta.$$

Under the assumption  $\|\mathbf{e}\|_2 \leq \frac{1}{c_5} \sqrt{\alpha} \sqrt{n} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2)$  for some  $c_5$  sufficiently large, we obtain the following bound for  $\|\hat{\mathbf{h}}\|_2$ :

$$\|\hat{\mathbf{h}}\|_2 \lesssim \frac{1}{\sqrt{\alpha} n (\|\beta_1^*\|_2 + \|\beta_2^*\|_2)} \eta.$$

To obtain a bound on  $\left\| \hat{\mathbf{H}} \right\|_F$ , we note that

$$\begin{aligned} \sum_b \|\mathbf{e}_b\|_2 \|\hat{\mathbf{h}}\|_2 &\leq \frac{1}{c_5} \sqrt{n} \sum_b \sqrt{\alpha} \|\beta_b^*\|_2 \|\hat{\mathbf{h}}\|_2 \\ &\leq \frac{c'}{c_5} \sqrt{n} \sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1, \end{aligned}$$

where we use the assumption on  $\|\mathbf{e}\|$  and (26) in the two inequalities, respectively. When  $c_5$  is large, we combine the last displayed equation with (21) to obtain

$$n\sqrt{\alpha} \left\| \hat{\mathbf{H}} \right\|_F \lesssim n \sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right) \right\|_1 \lesssim 2\eta,$$

where we use (25) in the last inequality. This implies

$$\left\| \hat{\mathbf{H}} \right\|_F \lesssim \frac{1}{n\sqrt{\alpha}} \eta,$$

completing the proof of Step (3) and thus Theorem 1.

### C. Proof of Theorem 3

We now turn to Theorem 3 for the stochastic noise setting. The main conceptual flow of the proof is quite similar to the deterministic noise case, though some significant additional steps are required. For the deterministic case, the starting point is the constraint, which allows us to bound  $\mathcal{A}_b$  and  $\mathcal{B}_b$  in terms of  $\eta$  using feasibility of  $(\mathbf{K}^*, \mathbf{g}^*)$  and  $(\mathbf{K}^* + \hat{\mathbf{H}}, \mathbf{g}^* + \hat{\mathbf{h}})$ . In the stochastic setup we have a Lagrangian (regularized) formulation, and hence we obtain the analogous result from optimality. Thus, the first step here involves showing that as a consequence of optimality, the solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}}) = (\mathbf{K}^* + \hat{\mathbf{H}}, \mathbf{g}^* + \hat{\mathbf{h}})$  satisfies inequality (29) below, which implies that  $\hat{\mathbf{H}}$  and  $\hat{\mathbf{h}}$  must live in a certain cone. The RIP-like condition for  $\mathcal{B}_b$  in the stochastic case is more demanding. We prove a second RIP-like condition (Lemma 5). We then bound  $\mathcal{A}$  by terms involving  $\mathcal{B}$ , and then invoke the RIP condition and the cone constraint.

Now we turn to each step of the proof. We continue to use the notation given in Section IV-A. We let  $\mathbf{D}_b := (|n_b/2|)^{-1} [\mathbf{d}_{b,1}, \dots, \mathbf{d}_{b,|n_b/2|}]^\top \in \mathbb{R}^{|n_b/2| \times p}$ . Recall that we defined

$$\gamma := \|\beta_1^*\|_2 + \|\beta_2^*\|_2.$$

Since the  $\{\mathbf{x}_i\}$  are assumed to be Gaussian with i.i.d. entries, the statement of the theorem is invariant under rotation of the  $\beta_b^*$ 's. Therefore, it suffices to prove the theorem assuming  $\beta_1^* - \beta_2^*$  is supported on the first coordinate. The follow lemma shows that we can further assume  $\{\mathbf{x}_i\}$  and  $\mathbf{e}$  have bounded entries, since we are interested in results that hold with high probability. This simplifies the subsequent analysis.

**Lemma 3.** *There exists an absolute constant  $c > 0$  such that, if the conclusion of Theorem 3 holds w.h.p. with the additional assumption that*

$$\begin{aligned} \mathbf{x}_i(l) &\leq c\sqrt{\log n}, \forall i \in [n], l \in [p], \\ e_i &\leq c\sigma\sqrt{\log n}, \forall i \in [n], \end{aligned}$$

*then it also holds w.h.p. without this assumption.*

We prove this lemma in Appendix F-B. In the sequel, we therefore assume  $\text{support}(\beta_1^* - \beta_2^*) = \{1\}$ , and the  $\{\mathbf{x}_i\}$  and  $\{\mathbf{e}_i\}$  satisfy the bounds in the above lemma.

1) *Step (1): Consequence of Optimality:* This step uses optimality of the solution  $(\hat{\mathbf{K}}, \hat{\mathbf{g}}) = (\mathbf{K}^* + \hat{\mathbf{H}}, \mathbf{g}^* + \hat{\mathbf{h}})$ , to get a bound on  $\mathcal{A}$ . By optimality, we have

$$\begin{aligned} &\sum_b \sum_{i \in \mathcal{I}_b} \left( -\langle \mathbf{x}_i \mathbf{x}_i^\top, \hat{\mathbf{K}} \rangle + 2y_i \langle \mathbf{x}_i, \hat{\mathbf{g}} \rangle - y_i^2 + \sigma^2 \right)^2 + \lambda \left\| \hat{\mathbf{K}} \right\|_* \\ &\leq \sum_b \sum_{i \in \mathcal{I}_b} \left( -\langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{K}^* \rangle + 2y_i \langle \mathbf{x}_i, \mathbf{g}^* \rangle - y_i^2 + \sigma^2 \right)^2 + \lambda \left\| \mathbf{K}^* \right\|_*. \end{aligned}$$

Using the expression (18), we have

$$\begin{aligned} & \sum_b \sum_{i \in \mathcal{I}_b} \left( -\langle \mathbf{x}_i \mathbf{x}_i^\top, \hat{\mathbf{H}} \rangle + 2(\mathbf{x}_i^\top \boldsymbol{\beta}_b^* + e_i) \langle \mathbf{x}_i, \hat{\mathbf{h}} \rangle \right. \\ & \quad \left. - e_i \mathbf{x}_i^\top \boldsymbol{\delta}_b^* - (e_i^2 - \sigma^2) \right)^2 + \lambda \|\hat{\mathbf{K}}\|_* \\ & \leq \sum_b \sum_{i \in \mathcal{I}_b} (-e_i \mathbf{x}_i^\top \boldsymbol{\delta}_b^* - (e_i^2 - \sigma^2))^2 + \lambda \|\mathbf{K}^*\|_*. \end{aligned}$$

Defining the noise vectors  $\mathbf{w}_{1,b} := -\mathbf{e}_b \circ (\mathbf{X}_b \boldsymbol{\delta}_b^*)$ ,  $\mathbf{w}_{2,b} := -(e_b^2 - \sigma^2 \mathbf{1})$  and  $\mathbf{w}_b = \mathbf{w}_{1,b} - \mathbf{w}_{2,b}$ , we can rewrite the inequality above as

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) + \mathbf{w}_b \right\|_2^2 + \lambda \|\hat{\mathbf{K}}\|_* \\ & \leq \sum_{b=1,2} \|\mathbf{w}_b\|_2^2 + \lambda \|\mathbf{K}^*\|_*. \end{aligned}$$

Expanding the squares and rearranging terms, we obtain

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \leq \sum_b 2 \left\langle \hat{\mathbf{H}} - 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top, n_b \mathcal{A}_b^* \mathbf{w}_b \right\rangle - \sum_b \left\langle \hat{\mathbf{h}}, 4\mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_b \right\rangle \\ & \quad + \lambda \left( \|\mathbf{K}^*\|_* - \|\hat{\mathbf{K}}\|_* \right) \\ & \stackrel{(a)}{\leq} \left( \|\hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_* \right) \cdot P + \|\hat{\mathbf{h}}\|_2 \cdot Q \\ & \quad + \lambda \left( \|\mathbf{K}^*\|_* - \|\hat{\mathbf{K}}\|_* \right) \\ & \stackrel{(b)}{\leq} \left( \|\hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_* \right) \cdot P + \|\hat{\mathbf{h}}\|_2 \cdot Q \\ & \quad + \lambda \left( \|\hat{\mathbf{H}}_T\|_* - \|\hat{\mathbf{H}}_{T^\perp}\|_* \right), \end{aligned} \tag{27}$$

where  $\mathcal{A}_b^*$  is the adjoint operator of  $\mathcal{A}_b$ , in (a) we have defined

$$P := 2 \sum_b \|n_b \mathcal{A}_b^* \mathbf{w}_b\|, \tag{28}$$

$$Q := 4 \sum_b \|\boldsymbol{\beta}_b^*\|_2 \|n_b \mathcal{A}_b^* \mathbf{w}_b\| + \sqrt{p} \left\| \sum_b 4\mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_b \right\|_\infty,$$

and (b) follows from (17). We need the following lemma, which bounds the noise terms  $P$  and  $Q$ . Its proof is a substantial part of the proof to the main result, but quite lengthy. We therefore defer it to Appendix A.

**Lemma 4.** *Under the assumptions of the theorem, we have  $\lambda \geq 2P$  and  $\lambda \geq \frac{1}{\sigma+\gamma}Q$  with high probability.*

Applying the lemma, we get

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \leq \lambda \left( \frac{3}{2} \|\hat{\mathbf{H}}_T\|_* - \frac{1}{2} \|\hat{\mathbf{H}}_{T^\perp}\|_* \right) + \lambda (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2. \end{aligned} \tag{29}$$

Since the left hand side of (29) is non-negative, we obtain the following cone constraint for the optimal solution:

$$\|\hat{\mathbf{H}}_{T^\perp}\|_* \leq \frac{5}{2} \|\hat{\mathbf{H}}_T\|_* + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2. \tag{30}$$

This concludes the proof of Step (1) of the proof.

**2) Step (2): RIP and Lower Bounds:** We can get a lower bound to the expression in the LHS of (29) using  $\mathcal{B}$ , as follows. Similarly as before, let  $K$  be some numeric constant to be chosen later. We partition  $\hat{\mathbf{H}}_{T^\perp}$  into a sum of  $M := \frac{p-4}{K}$  matrices  $\hat{\mathbf{H}}_1, \dots, \hat{\mathbf{H}}_M$  according to the SVD of  $\hat{\mathbf{H}}_{T^\perp}$ , such that  $\text{rank}(\hat{\mathbf{H}}_i) \leq K$  and the smallest singular value of  $\hat{\mathbf{H}}_i$  is larger than the largest singular value of  $\hat{\mathbf{H}}_{i+1}$ . Then we have the following chain of inequalities:

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \stackrel{(a)}{\geq} \sum_b \left\| n_b \mathcal{B}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2n_b \mathbf{D}_b \hat{\mathbf{h}} \right\|_2^2 \\ & \stackrel{(b)}{\geq} \sum_b n_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{D}_b \hat{\mathbf{h}} \right\|_1^2 \\ & \stackrel{(c)}{\gtrsim} n \left( \sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}} + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{D}_b \hat{\mathbf{h}} \right\|_1 \right)^2 \\ & \stackrel{(d)}{\geq} n \left( \sum_b \left\| \mathcal{B}_b \left( \hat{\mathbf{H}}_T + 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top + \hat{\mathbf{H}}_1 \right) + 2\mathbf{D}_b \hat{\mathbf{h}} \right\|_1 \right. \\ & \quad \left. - \sum_b \sum_{i=2}^M \left\| \mathcal{B}_b(\hat{\mathbf{H}}_i) \right\|_1 \right)^2. \end{aligned} \tag{31}$$

Here (a) follows from the definitions of  $\mathcal{A}_b$  and  $\mathcal{B}_b$  and the triangle inequality, (b) follows from  $\|\mathbf{u}\|_2^2 \geq \frac{1}{n_b} \|\mathbf{u}\|_1^2$  for all  $\mathbf{u} \in \mathbb{R}^{n_b}$ , (c) follows from  $n_1 \approx n_2$ , and (d) follows from the triangle inequality.

We see that in order to obtain lower bounds on (31) in terms of  $\|\hat{\mathbf{h}}\|_2$  and  $\|\hat{\mathbf{H}}\|_F$ , we need an extension of the previous RIP-like result from Lemma 1, in order to deal with the first term in (31). The following lemma is proved in Appendix B.

**Lemma 5.** *The following holds for some numerical constants  $c, \underline{c}, \bar{\delta}, \delta$ . For  $b = 1, 2$ , if  $\mu > 1$  and  $n_b \geq cpr$ , then with probability  $1 - \exp(-n_b)$ , we have the following RIP-2:*

$$\begin{aligned} \underline{\delta} (\|\mathbf{Z}\|_F + \sigma \|\mathbf{z}\|_2) & \leq \|\mathcal{B}_b \mathbf{Z} - \mathbf{D}_b \mathbf{z}\|_1 \leq \bar{\delta} (\|\mathbf{Z}\|_F + \sigma \|\mathbf{z}\|_2), \\ \forall \mathbf{z} \in \mathbb{R}^p, \forall \mathbf{Z} \in \mathbb{R}^{p \times p} \text{ with } \text{rank}(\mathbf{Z}) \leq r. \end{aligned}$$

Using this we can now bound the last inequality in (31) above. First, note that for each  $b = 1, 2$ ,

$$\begin{aligned} \sum_{i=2}^M \left\| \mathcal{B}_b(\hat{\mathbf{H}}_i) \right\|_1 & \stackrel{(a)}{\leq} \bar{\delta} \sum_{i=2}^M \|\hat{\mathbf{H}}_i\|_F \leq \bar{\delta} \sum_{i=2}^M \frac{1}{\sqrt{K}} \|\hat{\mathbf{H}}_{i-1}\|_* \\ & \leq \frac{\bar{\delta}}{\sqrt{K}} \|\hat{\mathbf{H}}_{T^\perp}\|_*, \end{aligned} \tag{32}$$

where (a) follows from the upper bound in Lemma 5 with  $\sigma$  set to 0. Then, applying the lower-bound in Lemma 5 to the first term in the parentheses in (31), and (32) to the second

term, we obtain

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \geq n \left( \sum_b \delta \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F + 2\delta\sigma \|\hat{\mathbf{h}}\|_2 - \bar{\delta} \sqrt{\frac{1}{K}} \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_* \right)^2 \\ & \gtrsim n \left( \sum_b \underline{\delta}^2 \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}}^\top \right\|_F^2 + \underline{\delta}^2 \sigma^2 \|\hat{\mathbf{h}}\|_2^2 - \bar{\delta}^2 \frac{1}{K} \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_*^2 \right). \end{aligned}$$

Choosing  $K$  to be sufficiently large, and applying Lemma 2, we obtain

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \gtrsim n \left( \left\| \hat{\mathbf{H}}_T \right\|_F^2 + \gamma^2 \|\hat{\mathbf{h}}\|_2^2 + \sigma^2 \|\hat{\mathbf{h}}\|_2^2 - \frac{1}{100} \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_*^2 \right). \end{aligned} \quad (33)$$

Using (30), we further get

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \gtrsim n \left[ \left\| \hat{\mathbf{H}}_T \right\|_F^2 + \gamma^2 \|\hat{\mathbf{h}}\|_2^2 + \sigma^2 \|\hat{\mathbf{h}}\|_2^2 - \frac{1}{8} \left\| \hat{\mathbf{H}}_T \right\|_*^2 \right. \\ & \quad \left. - \frac{1}{25} (\gamma^2 + \sigma^2) \|\hat{\mathbf{h}}\|_2^2 \right] \\ & \gtrsim \frac{1}{8} n \left( \left\| \hat{\mathbf{H}}_T \right\|_F + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2 \right)^2. \end{aligned} \quad (34)$$

This completes Step (2), and we are ready to combine the results to obtain error bounds, as promised in Step (3) and by the theorem.

3) *Step (3): Producing Error bounds:* Combining (29) and (34), we get

$$n \left( \left\| \hat{\mathbf{H}}_T \right\|_F + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2 \right)^2 \lesssim \lambda \|\mathbf{H}_T\|_F + \lambda(\gamma + \sigma) \|\hat{\mathbf{h}}\|_2,$$

which implies  $\left\| \hat{\mathbf{H}}_T \right\|_F + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2 \lesssim \frac{\lambda}{n}$ . It follows that  $\|\hat{\mathbf{h}}\|_2 \lesssim \frac{1}{n(\gamma + \sigma)} \lambda$  and

$$\begin{aligned} \left\| \hat{\mathbf{H}} \right\|_F & \leq \left\| \hat{\mathbf{H}}_T \right\|_* + \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_* \\ & \stackrel{(a)}{\leq} \frac{7}{2} \left\| \hat{\mathbf{H}}_T \right\|_* + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2 \\ & \stackrel{(b)}{\leq} \frac{7}{2} \cdot \sqrt{4} \left\| \hat{\mathbf{H}}_T \right\|_F + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2 \\ & \lesssim \frac{1}{n} \lambda, \end{aligned}$$

where we use (30) in (a) and  $\text{rank}(\hat{\mathbf{H}}_T) \leq 4$  in (b). This completes Step (3) and the proof of the theorem.

#### D. Proof of Theorem 5

In this section we prove the error bounds in Theorem 5 for our nonconvex approach.

1) *Statistical Error:* The proof of part (a) of Theorem 5 follows along the same lines as the proof of Theorem 3. In addition, we need to derive RIP and error bounds of the negative term. For the optimal solution  $(\hat{\mathbf{K}}, \hat{\mathbf{h}}) = (\mathbf{K}^* + \hat{\mathbf{H}}, \mathbf{g}^* + \hat{\mathbf{h}})$ , from the constraint we have that

$$\|\mathbf{K}^* + \hat{\mathbf{H}}\|_* \leq \|\mathbf{K}^*\|_*.$$

Using the decomposition  $\hat{\mathbf{H}} = \hat{\mathbf{H}}_T + \hat{\mathbf{H}}_{T^\perp}$ , we have  $\|\mathbf{K}^* + \hat{\mathbf{H}}\|_* = \|\mathbf{K}^* + \hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_* \geq \|\mathbf{K}^*\|_* - \|\hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_*$ . We thus have

$$\|\hat{\mathbf{H}}_{T^\perp}\|_* \leq \|\hat{\mathbf{H}}_T\|_*. \quad (35)$$

Following similar calculations as those in Section IV-C1, one can obtain the consequence of optimality as

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \quad - \underbrace{\sum_b 4\sigma^2 \mathcal{A}_b(\hat{\mathbf{h}} \hat{\mathbf{h}}^\top)}_{S_5} \\ & \leq \sum_b 2 \left\langle \hat{\mathbf{H}} - 2\beta_b^* \hat{\mathbf{h}}^\top, n_b \mathcal{A}_b^* \mathbf{w}_b \right\rangle - \sum_b \left\langle \hat{\mathbf{h}}, 4\mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_b \right\rangle \\ & \quad - \underbrace{\sum_b 4\sigma^2 \mathcal{A}_b(\delta_b^* \hat{\mathbf{h}}^\top)}_{S_6} - \underbrace{\sum_b 8\sigma^2 \langle \mathbf{e}_b, \mathbf{X}_b \hat{\mathbf{h}} \rangle}_{S_7}. \end{aligned} \quad (36)$$

Compared to (27), the above equality does not have the nuclear norm term since we remove the regularization term. The additional terms  $S_5$ ,  $S_6$  and  $S_7$  come from the negative term  $-\sum_{i=1}^n 4\sigma^2(y_i - \langle \mathbf{x}_i, \mathbf{g} \rangle)^2$ .

By standard concentration result, when  $n \gtrsim p/\epsilon^2$ , with probability at least  $1 - \exp(-p)$ ,  $\|\frac{1}{n} \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{I}\| \leq \epsilon$ . Choosing  $\epsilon = 0.1$ , we have w.h.p.

$$S_5 = 4\sigma^2 \sum_{i=1}^n \langle \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{I}, \hat{\mathbf{h}} \hat{\mathbf{h}}^\top \rangle + 4n\sigma^2 \|\hat{\mathbf{h}}\|_2^2 \leq 4.4n\sigma^2 \|\hat{\mathbf{h}}\|_2^2.$$

Using the above result and (34), we obtain

$$\begin{aligned} & \sum_b \left\| n_b \mathcal{A}_b \left( -\hat{\mathbf{H}} + 2\beta_b^* \hat{\mathbf{h}}^\top \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \hat{\mathbf{h}}) \right\|_2^2 \\ & \quad - \underbrace{\sum_b 4\sigma^2 \mathcal{A}_b(\hat{\mathbf{h}} \hat{\mathbf{h}}^\top)}_{S_5} \\ & \gtrsim n \left( \left\| \hat{\mathbf{H}}_T \right\|_F + (\gamma + \sigma) \|\hat{\mathbf{h}}\|_2 \right)^2 - n\sigma^2 \|\hat{\mathbf{h}}\|_2^2 \\ & \gtrsim n \left( \left\| \hat{\mathbf{H}}_T \right\|_F + \gamma \|\hat{\mathbf{h}}\|_2 \right)^2, \end{aligned} \quad (37)$$

where the second inequality follows from the assumption  $\gamma/\sigma \geq c$  for sufficiently large constant  $c$ .

Now we turn to the right hand side of (36). By the Cauchy-Schwarz inequality, it is upper bounded by

$$\left( \|\hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_* \right) \cdot P + \|\hat{\mathbf{h}}\|_2 \cdot Q',$$

where  $P$  is given in (28), and we defined

$$Q' := 4\gamma P + W$$

and

$$\begin{aligned} W := & \sqrt{p} \left\| \sum_b \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_b + \sum_b \sigma^2 \mathbf{X}_b^\top \mathbf{X}_b \boldsymbol{\delta}_b^* \right. \\ & \left. + \sum_b 2\sigma^2 \mathbf{X}_b^\top \mathbf{e}_b \right\|_\infty. \end{aligned} \quad (38)$$

Using Lemma 4, we obtain  $P \lesssim \sigma\gamma\sqrt{np}\log^3 n$ . The following lemma gives an upper bound of the term  $W$ ; see Appendix F-C for the proof.

**Lemma 6.** *Under the assumptions of the Theorem 5, we have that w.h.p.,*

$$W \lesssim \sigma^2\gamma\sqrt{np}\log^2 n.$$

We therefore conclude that the right hand side of (36) is at most on the order of

$$\sigma\gamma\sqrt{np}\log^3 n \left( \|\hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_* \right) + \sigma\gamma^2\sqrt{np}\log^3 n \|\hat{\mathbf{h}}\|_2. \quad (39)$$

Putting (37) and (39) together, we finish our proof by showing

$$\begin{aligned} & n \left( \|\hat{\mathbf{H}}_T\|_F + \gamma\|\hat{\mathbf{h}}\|_2 \right)^2 \\ & \lesssim \sigma\gamma\sqrt{np}\log^3 n \left( \|\hat{\mathbf{H}}_T\|_* + \|\hat{\mathbf{H}}_{T^\perp}\|_* \right) + \sigma\gamma^2\sqrt{np}\log^3 n \|\hat{\mathbf{h}}\|_2 \\ & \lesssim \sigma\gamma\sqrt{np}\log^3 n \|\hat{\mathbf{H}}_T\|_F + \sigma\gamma^2\sqrt{np}\log^3 n \|\hat{\mathbf{h}}\|_2, \end{aligned}$$

where the second inequality follows from the conic constraint (35) and  $\|\hat{\mathbf{H}}_T\|_* \leq 2\|\hat{\mathbf{H}}_T\|_F$ .

2) *Optimization Error:* Recall that  $\mathbf{H}^{(t)} := \mathbf{K}^{(t)} - \hat{\mathbf{K}}$  and  $\mathbf{h}^{(t)} := \mathbf{g}^{(t)} - \hat{\mathbf{g}}$ , which represent the optimization error, and  $\hat{\mathbf{H}} := \hat{\mathbf{K}} - \mathbf{K}^*$  and  $\hat{\mathbf{h}} := \hat{\mathbf{g}} - \mathbf{g}^*$ , which represent the statistical error. We need several auxiliary lemmas for the proof. First, we establish a result that shows that  $\mathbf{H}^{(t)}$  belongs to a cone-like set. See Appendix F-D for the proof.

**Lemma 7.** *From program (8), due to the nuclear norm constraint, we have that for any  $t = 0, 1, \dots$ ,*

$$\|\mathbf{H}_{T^\perp}^{(t)}\|_* \leq 3\|\mathbf{H}_T^{(t)}\|_* + 4\|\hat{\mathbf{H}}_T\|_*.$$

Next, we characterize the curvature of  $\mathcal{L}_n$ . To ease notation, for any  $\mathbf{K}, \mathbf{K}' \in \mathbb{R}^{p \times p}$  and  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^p$ , we define

$$\begin{aligned} Q_n(\mathbf{K}', \mathbf{g}'; \mathbf{K}, \mathbf{g}) &:= \mathcal{L}_n(\mathbf{K}', \mathbf{g}') - \mathcal{L}_n(\mathbf{K}, \mathbf{g}) - \langle \nabla_{\mathbf{K}} \mathcal{L}_n(\mathbf{K}, \mathbf{g}), \mathbf{K}' - \mathbf{K} \rangle \\ &\quad - \langle \nabla_{\mathbf{g}} \mathcal{L}_n(\mathbf{K}, \mathbf{g}), \mathbf{g}' - \mathbf{g} \rangle \\ &= \sum_b \left\| n_b \mathcal{A}_b \left( -(\mathbf{K}' - \mathbf{K}) + 2\boldsymbol{\beta}_b^* (\mathbf{g}' - \mathbf{g})^\top \right) \right. \\ &\quad \left. + 2\mathbf{e}_b \circ (\mathbf{X}_b (\mathbf{g}' - \mathbf{g})) \right\|_2^2 - \sum_{i=1}^n 4\sigma^2 \langle \mathbf{x}_i, \mathbf{g}' - \mathbf{g} \rangle^2. \end{aligned} \quad (40)$$

The next lemma shows the smoothness of  $\mathcal{L}_n$ . Its proof is given in Appendix C.

**Lemma 8.** *Given  $Q_n(\mathbf{K}', \mathbf{g}'; \mathbf{K}, \mathbf{g})$  that is defined in (40), we have that there exist constants  $\{c_i\}_{i=1}^2$  such that with probability at least  $1 - c_2/n$ , the inequality*

$$\begin{aligned} Q_n(\mathbf{K}', \mathbf{g}'; \mathbf{K}, \mathbf{g}) &\leq c_1 n \sqrt{p^3 \log n} \|\mathbf{K}' - \mathbf{K}\|_F + c_1 n \log n \gamma^2 \|\mathbf{g}' - \mathbf{g}\|_2^2 \end{aligned} \quad (41)$$

holds for any  $\mathbf{K}, \mathbf{K}' \in \mathbb{R}^{p \times p}$  and  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^p$ .

Next we turn to the lower bound of  $Q$ . We first let  $\mathbf{K}' = \mathbf{K}^{(t)}$ ,  $\mathbf{g}' = \mathbf{g}^{(t)}$ ,  $\mathbf{K} = \hat{\mathbf{K}}$ ,  $\mathbf{g} = \hat{\mathbf{g}}$ . Then we have

$$\begin{aligned} & Q_n(\hat{\mathbf{K}}, \hat{\mathbf{g}}; \mathbf{K}^{(t)}, \mathbf{g}^{(t)}) \\ &= \sum_b \left\| n_b \mathcal{A}_b \left( -\mathbf{H}^{(t)} + 2\boldsymbol{\beta}_b^* \mathbf{h}^{(t)\top} \right) + 2\mathbf{e}_b \circ (\mathbf{X}_b \mathbf{h}^{(t)}) \right\|_2^2 \\ &\quad - \sum_{i=1}^n 4\sigma^2 \langle \mathbf{x}_i, \mathbf{h}^{(t)} \rangle^2 \\ &\gtrsim n \left( \left\| \mathbf{H}_T^{(t)} \right\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 + \sigma^2 \|\mathbf{h}^{(t)}\|_2^2 - \frac{1}{100} \left\| \mathbf{H}_{T^\perp}^{(t)} \right\|_*^2 \right) \\ &\quad - \sum_b 4\sigma^2 \mathcal{A}_b(\hat{\mathbf{h}} \hat{\mathbf{h}}^\top) \\ &\gtrsim n \left( \left\| \mathbf{H}_T^{(t)} \right\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 + \sigma^2 \|\mathbf{h}^{(t)}\|_2^2 - \frac{1}{100} \left\| \mathbf{H}_{T^\perp}^{(t)} \right\|_*^2 \right) \\ &\quad - n\sigma^2 \|\mathbf{h}^{(t)}\|_2^2, \end{aligned}$$

where the first inequality follows from (33) by replacing the  $(\hat{\mathbf{H}}, \hat{\mathbf{h}})$  there by  $(\mathbf{H}^{(t)}, \mathbf{h}^{(t)})$ ; the second inequality is from the concentration property of operator  $\mathcal{A}_b$  as in the proof of part (a). Using Lemma 7 to bound  $\|\mathbf{H}_{T^\perp}^{(t)}\|_*^2$ , we further have

$$\begin{aligned} & Q_n(\hat{\mathbf{K}}, \hat{\mathbf{g}}; \mathbf{K}^{(t)}, \mathbf{g}^{(t)}) \\ &\gtrsim n \left( 0.38 \left\| \mathbf{H}_T^{(t)} \right\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 + \sigma^2 \|\mathbf{h}^{(t)}\|_2^2 - 0.32 \left\| \hat{\mathbf{H}}_T \right\|_*^2 \right. \\ &\quad \left. - n\sigma^2 \|\mathbf{h}^{(t)}\|_2^2 \right) \\ &\gtrsim n \left\| \mathbf{H}^{(t)} \right\|_F^2 + n\gamma^2 \|\mathbf{h}^{(t)}\|_2^2 - n \left\| \hat{\mathbf{H}} \right\|_F^2. \end{aligned} \quad (42)$$

We need the following the result to connect the established curvature property to the optimization error. Its proof is given in Appendix F-E.

**Lemma 9.** *Suppose function  $Q_n$  satisfies the following two conditions with functions  $\bar{Q}_n$  and  $\underline{Q}_n$  that map  $\mathbb{R}^{p \times p} \times \mathbb{R}^p$  to  $\mathbb{R}$ :*

- For  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  and every  $(\mathbf{K}^{(t)}, \mathbf{g}^{(t)})$ ,  $\underline{Q}_n(\mathbf{K}^{(t)}, \mathbf{g}^{(t)}; \hat{\mathbf{K}}, \hat{\mathbf{g}}) \geq \bar{Q}_n(\mathbf{H}^{(t)}, \mathbf{h}^{(t)})$ ;
- For any  $\mathbf{K}, \mathbf{K}' \in \mathbb{R}^{p \times p}$  and  $\mathbf{g}, \mathbf{g}' \in \mathbb{R}^p$ ,  $\underline{Q}_n(\mathbf{K}', \mathbf{g}'; \mathbf{K}, \mathbf{g}) \leq \bar{Q}_n(\mathbf{K}' - \mathbf{K}, \mathbf{g}' - \mathbf{g})$ .

Then we have

$$\begin{aligned} & \|\mathbf{H}^{(t+1)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t+1)}\|_2^2 \\ &\leq \|\mathbf{H}^{(t)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 - \frac{2}{\eta} \bar{Q}_n(\mathbf{H}^{(t)}, \mathbf{h}^{(t)}) - \|\Delta^{(t)}\|_F^2 \\ &\quad - \gamma^2 \|\boldsymbol{\delta}^{(t)}\|_2^2 + \frac{2}{\eta} \bar{Q}_n(\Delta^{(t)}, \boldsymbol{\delta}^{(t)}), \end{aligned}$$

where we let  $\Delta^{(t)} := \mathbf{K}^{(t+1)} - \mathbf{K}^{(t)}$  and  $\boldsymbol{\delta}^{(t)} := \mathbf{g}^{(t+1)} - \mathbf{g}^{(t)}$ .

Now we are ready to prove Theorem 5 part (b). Plugging the established bounds of the function  $Q_n$ , (41) and (42), into Lemma 9, we obtain that for some constants  $c_1, c_2$ ,

$$\begin{aligned} & \|\mathbf{H}^{(t+1)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t+1)}\|_2^2 \\ &\leq \left( 1 - \frac{c_1 n}{\eta} \right) \left( \|\mathbf{H}^{(t)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 \right) + \frac{c_1 n}{\eta} \|\hat{\mathbf{H}}\|_F^2 \\ &\quad + \left( \frac{c_2 n \log n}{\eta} - 1 \right) \|\boldsymbol{\delta}^{(t)}\|_2^2 + \frac{c_2 n \sqrt{p^3 \log n}}{\eta} \|\Delta^{(t)}\|_F^2 - \|\Delta^{(t)}\|_F^2. \end{aligned}$$

Therefore by setting  $\eta > c_3 n \sqrt{p^3 \log n}$  for sufficiently large constant  $c_3$ , we have

$$\begin{aligned} & \|\mathbf{H}^{(t+1)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t+1)}\|_2^2 \\ & \leq \left(1 - \frac{c_1 n}{\eta}\right) \left(\|\mathbf{H}^{(t)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2\right) + \frac{c_1 n}{\eta} \|\hat{\mathbf{H}}\|_F^2. \end{aligned}$$

Applying this bound recursively to  $t = 0, 1, 2, \dots$  proves our result.

### E. Proofs of Theorems 2, 4 and 6

In this section, we show that an error bound on the input  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  of Algorithm 1 implies an error bound on its output  $(\hat{\boldsymbol{\beta}}_1, \hat{\boldsymbol{\beta}}_2)$ . Recall the quantities  $\hat{\mathbf{J}}$ ,  $\mathbf{J}^*$ ,  $\hat{\lambda}$ ,  $\lambda^*$ ,  $\hat{\mathbf{v}}$  and  $\mathbf{v}^*$  defined in Section III-A and in Algorithm 1.

The main component of the proof is a general perturbation bound. We prove these in the first section below, and then use them to prove Theorems 2, 4 and 6 in the three subsequent sections.

1) *Perturbation Bound:* We require the following perturbation bound.

**Lemma 10.** *If  $\|\hat{\mathbf{J}} - \mathbf{J}^*\|_F \leq \delta$ , then*

$$\left\| \sqrt{\hat{\lambda}} \hat{\mathbf{v}} - \sqrt{\lambda^*} \mathbf{v}^* \right\|_2 \leq 10 \min \left\{ \frac{\delta}{\sqrt{\|\mathbf{J}^*\|}}, \sqrt{\delta} \right\}.$$

*Proof.* By Weyl's inequality, we have

$$|\hat{\lambda} - \lambda^*| \leq \|\hat{\mathbf{J}} - \mathbf{J}^*\| \leq \delta.$$

This implies

$$|\sqrt{\hat{\lambda}} - \sqrt{\lambda^*}| = \left| \frac{\hat{\lambda} - \lambda^*}{\sqrt{\hat{\lambda}} + \sqrt{\lambda^*}} \right| \leq 2 \min \left\{ \frac{\delta}{\sqrt{\lambda^*}}, \sqrt{\delta} \right\}. \quad (43)$$

Using Weyl's inequality and the Davis-Kahan sine theorem, we obtain

$$|\sin \angle(\hat{\mathbf{v}}, \mathbf{v}^*)| \leq \min \left\{ \frac{2\|\hat{\mathbf{K}} - \mathbf{K}^*\|}{\|\mathbf{K}^*\|}, 1 \right\} \leq \min \left\{ \frac{2\delta}{\lambda^*}, 1 \right\}. \quad (44)$$

On the other hand, we have

$$\begin{aligned} & \left\| \hat{\mathbf{v}} \sqrt{\hat{\lambda}} - \mathbf{v}^* \sqrt{\lambda^*} \right\|_2 \\ & \leq \left\| \hat{\mathbf{v}} \sqrt{\hat{\lambda}} - \mathbf{v}^* \sqrt{\hat{\lambda}} \right\|_2 + \left\| \mathbf{v}^* \sqrt{\hat{\lambda}} - \mathbf{v}^* \sqrt{\lambda^*} \right\|_2 \\ & = \sqrt{\hat{\lambda}} \|\hat{\mathbf{v}} - \mathbf{v}^*\|_2 + \|\mathbf{v}^*\|_2 \left| \sqrt{\hat{\lambda}} - \sqrt{\lambda^*} \right| \\ & = \left( \sqrt{\lambda^*} + \sqrt{\hat{\lambda}} - \sqrt{\lambda^*} \right) \|\hat{\mathbf{v}} - \mathbf{v}^*\|_2 + \|\mathbf{v}^*\|_2 \left| \sqrt{\hat{\lambda}} - \sqrt{\lambda^*} \right| \\ & \leq \sqrt{\lambda^*} \|\hat{\mathbf{v}} - \mathbf{v}^*\|_2 + 3 \left| \sqrt{\hat{\lambda}} - \sqrt{\lambda^*} \right|, \end{aligned}$$

where in the last inequality we use the fact that  $\|\mathbf{v}^*\| = \|\hat{\mathbf{v}}\| = 1$ . Elementary calculation shows that

$$\|\hat{\mathbf{v}} - \mathbf{v}^*\|_2 = 2 \left| \sin \frac{1}{2} \angle(\hat{\mathbf{v}}, \mathbf{v}^*) \right| \leq \sqrt{2} |\sin \angle(\hat{\mathbf{v}}, \mathbf{v}^*)|.$$

It follows that

$$\begin{aligned} & \left\| \hat{\mathbf{v}} \sqrt{\hat{\lambda}} - \mathbf{v}^* \sqrt{\lambda^*} \right\|_2 \\ & \leq \sqrt{2} \sqrt{\lambda^*} |\sin \angle(\hat{\mathbf{v}}, \mathbf{v}^*)| + 3 \left| \sqrt{\hat{\lambda}} - \sqrt{\lambda^*} \right| \\ & \leq \sqrt{2} \min \left\{ \frac{2\delta}{\sqrt{\lambda^*}}, \sqrt{\lambda^*} \right\} + 6 \min \left\{ \frac{\delta}{\sqrt{\lambda^*}}, \sqrt{\delta} \right\} \\ & \leq 10 \min \left\{ \frac{\delta}{\sqrt{\lambda^*}}, \sqrt{\delta} \right\}, \end{aligned}$$

where we use (43) and (44) in the second inequality. This concludes the proof.

We can now use this perturbation result to provide guarantees on recovering  $\hat{\boldsymbol{\beta}}_1$  and  $\hat{\boldsymbol{\beta}}_2$  given noisy versions of  $\mathbf{g}^*$  and  $\mathbf{K}^*$ . To this end, suppose we are given  $\hat{\mathbf{K}}$  and  $\hat{\mathbf{g}}$  which satisfy

$$\|\hat{\mathbf{K}} - \mathbf{K}^*\|_F \leq \delta_K, \quad \|\hat{\mathbf{g}} - \mathbf{g}^*\|_2 \leq \delta_g.$$

Then by the triangle inequality we have

$$\|\hat{\mathbf{J}} - \mathbf{J}^*\|_F \leq \delta_K + 2\delta_g \|\mathbf{g}^*\|_2 + \delta_g^2.$$

Therefore, up to relabeling  $b$ , we have

$$\begin{aligned} & \|\hat{\boldsymbol{\beta}}_b - \boldsymbol{\beta}_b^*\|_2 \\ & \leq \|\hat{\mathbf{g}} - \mathbf{g}^*\|_2 + \left\| \sqrt{\hat{\lambda}} \hat{\mathbf{v}} - \sqrt{\lambda^*} \mathbf{v}^* \right\|_2 \\ & \lesssim \delta_g + \min \left\{ \frac{\delta_K + 2\delta_g \|\mathbf{g}^*\|_2 + \delta_g^2}{\|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2}, \sqrt{\delta_K + 2\delta_g \|\mathbf{g}^*\|_2 + \delta_g^2} \right\}, \end{aligned} \quad (45)$$

where the second inequality follows from Lemma 10 and  $\lambda^* = \frac{1}{4} \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2^2$ .

2) *Proof of Theorem 2 (Arbitrary Noise):* In the case of arbitrary noise, as set up above, Theorem 1 guarantees the following:

$$\begin{aligned} \delta_K & \asymp \frac{\sqrt{n} \|\mathbf{e}\|_2 \|\boldsymbol{\beta}_2^* - \boldsymbol{\beta}_1^*\|_2 + \|\mathbf{e}\|_2^2}{\sqrt{\alpha n}} \lesssim \frac{1}{\sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}} \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2, \\ \delta_g & \asymp \frac{\sqrt{n} \|\mathbf{e}\|_2 \|\boldsymbol{\beta}_2^* - \boldsymbol{\beta}_1^*\|_2 + \|\mathbf{e}\|_2^2}{\sqrt{\alpha n} (\|\boldsymbol{\beta}_1^*\|_2 + \|\boldsymbol{\beta}_2^*\|_2)} \lesssim \frac{\|\mathbf{e}\|_2}{\sqrt{n}}. \end{aligned}$$

where we use the assumption  $\|\mathbf{e}\|_2 \leq \frac{\sqrt{\alpha}}{c_4} \sqrt{n} (\|\boldsymbol{\beta}_1^*\|_2 + \|\boldsymbol{\beta}_2^*\|_2) \asymp \frac{1}{c_4} \sqrt{n} \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2$ . Using (45), we get that up to relabeling  $b$ ,

$$\begin{aligned} & \|\hat{\boldsymbol{\beta}}_b - \boldsymbol{\beta}_b^*\|_2 \\ & \lesssim \frac{\|\mathbf{e}\|_2}{\sqrt{n}} + \min \left\{ \frac{1}{\sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}} + \frac{\|\mathbf{e}\|_2^2}{n \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2}, \right. \\ & \quad \left. \sqrt{\frac{1}{\sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}} \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2 + \frac{\|\mathbf{e}\|_2^2}{n}} \right\} \\ & \lesssim \frac{1}{\sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}} + \min \left\{ \frac{\|\mathbf{e}\|_2^2}{n \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2}, \sqrt{\frac{1}{\sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}} \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2} \right\} \\ & \leq \frac{1}{\sqrt{\alpha}} \frac{\|\mathbf{e}\|_2}{\sqrt{n}}. \end{aligned}$$

3) *Proof of Theorem 4 (Stochastic Noise):* Next consider the setting with stochastic noise. Under the assumptions of Theorem 4, Theorem 3 guarantees the following bounds on the errors in recovering  $\mathbf{K}^*$  and  $\mathbf{g}^*$ :

$$\begin{aligned}\delta_K &\asymp \sigma (\|\beta_1^*\|_2 + \|\beta_2^*\|_2 + \sigma) \sqrt{\frac{p}{n} \log^4 n}, \\ \delta_g &\asymp \sigma \sqrt{\frac{p}{n} \log^4 n}.\end{aligned}$$

If we let  $\gamma = \|\beta_1^*\|_2 + \|\beta_2^*\|_2$ , then this means

$$\begin{aligned}\delta_K + 2\delta_g \|\mathbf{g}^*\|_2 + \delta_g^2 \\ \asymp \sigma\gamma \sqrt{\frac{p}{n} \log^4 n} + \sigma^2 \sqrt{\frac{p}{n} \log^4 n} + \sigma^2 \frac{p}{n} \log^8 n \\ \lesssim \sigma\gamma \sqrt{\frac{p}{n} \log^4 n} + \sigma^2 \sqrt{\frac{p}{n} \log^4 n},\end{aligned}$$

where the last inequality follows from the assumption that  $n \geq cp \log^8 n$  for some  $c > 1$ . Combining these with (45), we obtain that up to relabeling of  $b$ ,

$$\begin{aligned}\|\hat{\beta}_b - \beta_b^*\|_2 \\ \lesssim \sigma \sqrt{\frac{p}{n} \log^4 n} \\ + \min \left\{ \frac{\sigma\gamma \sqrt{\frac{p}{n}} + \sigma^2 \sqrt{\frac{p}{n}}}{\sqrt{\alpha\gamma}}, \sqrt{\sigma\gamma \sqrt{\frac{p}{n}} + \sigma^2 \sqrt{\frac{p}{n}}} \right\} \log^4 n \\ \lesssim \sigma \sqrt{\frac{p}{n} \log^4 n} + \min \left\{ \frac{\sigma^2 \sqrt{\frac{p}{n}}}{\gamma}, \sqrt{\sigma\gamma \sqrt{\frac{p}{n}} + \sigma^2 \sqrt{\frac{p}{n}}} \right\} \log^4 n,\end{aligned}$$

where the last inequality follows from  $\alpha$  being lower-bounded by a constant. Observe that the minimization in the last RHS is no larger than  $\sigma\sqrt{\frac{p}{n}}$  if  $\gamma \geq \sigma$ , and equals  $\min \left\{ \frac{\sigma^2 \sqrt{\frac{p}{n}}}{\gamma}, \sigma \left( \frac{p}{n} \right)^{1/4} \right\}$  if  $\gamma < \sigma$ . It follows that

$$\begin{aligned}\|\hat{\beta}_b - \beta_b^*\|_2 \\ \lesssim \sigma \sqrt{\frac{p}{n} \log^4 n} + \min \left\{ \frac{\sigma^2 \sqrt{\frac{p}{n}}}{\gamma}, \sigma \left( \frac{p}{n} \right)^{1/4} \right\} \log^4 n.\end{aligned}$$

4) *Proof of Theorem 6 (Nonconvex Formulation):* Under our assumption, we have  $\|\mathbf{K}^{(T)} - \hat{\mathbf{K}}\|_F \lesssim \|\hat{\mathbf{H}}\|_F, \|\mathbf{g}^{(T)} - \hat{\mathbf{g}}\|_2 \lesssim \|\hat{\mathbf{h}}\|_2$ . From the triangle inequality, we have

$$\begin{aligned}\|\mathbf{K}^{(T)} - \mathbf{K}^*\|_F &\leq \|\mathbf{K}^{(T)} - \hat{\mathbf{K}}\|_F + \|\hat{\mathbf{H}}\|_F \\ &\lesssim \sigma (\|\beta_1^*\|_2 + \|\beta_2^*\|_2) \sqrt{\frac{p}{n} \log^3 n},\end{aligned}$$

and

$$\|\mathbf{g}^{(T)} - \mathbf{g}^*\|_2 \leq \|\mathbf{g}^{(T)} - \hat{\mathbf{g}}\|_2 + \|\hat{\mathbf{h}}\|_2 \lesssim \sigma \sqrt{\frac{p}{n} \log^3 n}.$$

From here the proof follows along the lines of Section IV-E3. We omit the details.

#### F. Proof Outlines of Theorems 7 and 8

Next, we provide the main steps of proving the minimax lower bounds in Theorems 7 and 8, and postpone the full proofs to Appendix D and E. The high-level ideas in the proofs of Theorems 7 and 8 are similar: we use a standard argument [5, 38, 35] to convert the estimation problem into a hypothesis testing problem, and then use information-theoretic inequalities to lower bound the error probability in hypothesis testing. In particular, recall the definition of the set  $\Theta(\underline{\gamma})$  of regression vector pairs in (10); we construct a  $\delta$ -packing  $\Theta = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\}$  of  $\Theta(\underline{\gamma})$  in the metric  $\rho$ , and use the following inequality:

$$\inf_{\tilde{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \mathbb{E} \left[ \rho(\tilde{\boldsymbol{\theta}}, \boldsymbol{\theta}^*) \right] \geq \delta \inf_{\tilde{\boldsymbol{\theta}}} \mathbb{P} \left( \tilde{\boldsymbol{\theta}} \neq \boldsymbol{\theta}^* \right), \quad (46)$$

where on the RHS  $\boldsymbol{\theta}^*$  is assumed to be sampled uniformly at random from  $\Theta$ . To lower-bound the minimax expected error by  $\frac{1}{2}\delta$ , it suffices to show that the probability on the last RHS is at least  $\frac{1}{2}$ . By Fano's inequality [17], we have

$$\mathbb{P} \left( \tilde{\boldsymbol{\theta}} \neq \boldsymbol{\theta}^* \right) \geq 1 - \frac{I(\mathbf{y}, \mathbf{X}; \boldsymbol{\theta}^*) + \log 2}{\log M}. \quad (47)$$

It remains to construct a packing set  $\Theta$  with the appropriate separation  $\delta$  and cardinality  $M$ , and to upper-bound the mutual information  $I(\mathbf{y}, \mathbf{X}; \boldsymbol{\theta}^*)$ . We show how to do this for Part 2 of Theorem 8, for which the desired separation is  $\delta = 2c_0 \frac{\sigma^2}{\kappa} \sqrt{\frac{p}{n}}$ , where  $\kappa = \frac{\underline{\gamma}}{2}$ . Let  $\{\xi_1, \dots, \xi_M\}$  be a  $\frac{p-1}{16}$ -packing of  $\{0, 1\}^{p-1}$  in Hamming distance with  $\log M \geq (p-1)/16$ , which exists by the Varshamov-Gilbert bound [29]. We construct  $\Theta$  by setting  $\boldsymbol{\theta}_i := (\beta_i, -\beta_i)$  for  $i = 1, \dots, M$  with

$$\boldsymbol{\beta}_i = \kappa_0 \boldsymbol{\epsilon}_p + \sum_{j=1}^{p-1} (2\xi_i(j) - 1) \tau \boldsymbol{\epsilon}_j,$$

where  $\tau = \frac{4\delta}{\sqrt{p-1}}$ ,  $\kappa_0^2 = \kappa^2 - (p-1)\tau^2$ , and  $\boldsymbol{\epsilon}_j$  is the  $j^{th}$  standard basis in  $\mathbb{R}^p$ . We verify that this  $\Theta$  indeed defines a  $\delta$ -packing of  $\Theta(\underline{\gamma})$ , and moreover satisfies  $\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2 \leq 16\delta^2$  for all  $i \neq i'$ . To bound the mutual information, we observe that by independence between  $\mathbf{X}$  and  $\boldsymbol{\theta}^*$ , we have

$$\begin{aligned}I(\boldsymbol{\theta}^*; \mathbf{X}, \mathbf{y}) &\leq \frac{1}{M^2} \sum_{1 \leq i, i' \leq M} D(\mathbb{P}_i \|\mathbb{P}_{i'}) \\ &= \frac{1}{M} \sum_{1 \leq i, i' \leq M} \sum_{j=1}^n \mathbb{E}_{\mathbf{X}} \left[ D \left( \mathbb{P}_{i, \mathbf{X}}^{(j)} \|\mathbb{P}_{i', \mathbf{X}}^{(j)} \right) \right],\end{aligned}$$

where  $\mathbb{P}_{i, \mathbf{X}}^{(j)}$  denotes the distribution of  $y_j$  conditioned on  $\mathbf{X}$  and  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_i$ . The remaining and crucial step is to obtain sharp upper bounds on the above KL-divergence between two mixtures of one-dimensional Gaussian distributions. This requires some technical calculations, from which we obtain

$$\mathbb{E}_{\mathbf{X}} D \left( \mathbb{P}_{i, \mathbf{X}}^{(j)} \|\mathbb{P}_{i', \mathbf{X}}^{(j)} \right) \leq \frac{c' \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2 \kappa^2}{\sigma^4}.$$

We conclude that  $I(\boldsymbol{\theta}^*; \mathbf{X}, \mathbf{y}) \leq \frac{1}{4} \log M$ . Combining with (46) and (47) proves Part 2 of Theorem 8. Theorem 7 and the remaining part of Theorem 8 are proved in a similar manner.

## V. CONCLUSION

This paper provides a computationally and statistically efficient algorithm for mixed regression with two components. To the best of our knowledge, this is the first efficient algorithm that can provide  $O(p)$  sample complexity guarantees. Under certain conditions, we prove matching lower bounds, thus demonstrating our algorithm achieves the minimax optimal rates. There are several interesting open questions that remain. Most immediate is the issue of understanding the degree to which the assumptions currently required for minimax optimality can be removed or relaxed. The extension to more than two components is important, though how to do this within the current framework is not obvious.

At its core, the approach here is a method of moments, as the convex optimization formulation produces an estimate of the cross moments,  $(\beta_1^* \beta_2^{*\top} + \beta_2^* \beta_1^{*\top})$ . An interesting aspect of these results is the significant improvement in sample complexity guarantees this tailored approach brings, compared to a more generic implementation of the tensor machinery which requires use of third order moments. Given the statistical and also computational challenges related to third order tensors, understanding the connections more carefully seems to be an important future direction.

APPENDIX A  
PROOF OF LEMMA 4

We now move to the proof of Lemma 4, which bounds the noise terms  $P$  and  $Q$ . Note that

$$\begin{aligned} P &= 2 \sum_b \|n_b \mathcal{A}_b^* \mathbf{w}_b\| \\ &\leq 2 \underbrace{\sum_b \|n_b \mathcal{A}_b^* \mathbf{w}_{1,b}\|}_{S_1} + 2 \underbrace{\sum_b \|n_b \mathcal{A}_b^* \mathbf{w}_{2,b}\|}_{S_2}, \end{aligned}$$

and

$$\begin{aligned} Q &= \sum_b 4 \|\beta_b^*\|_2 \|n_b \mathcal{A}_b^* \mathbf{w}_b\| + \sqrt{p} \left\| \sum_b 4 \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_b \right\|_\infty \\ &\leq 4\gamma P + \sqrt{p} \underbrace{\left\| \sum_b 4 \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_{1,b} \right\|_\infty}_{S_3} \\ &\quad + \sqrt{p} \underbrace{\left\| \sum_b 4 \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_{2,b} \right\|_\infty}_{S_4}. \end{aligned}$$

Therefore, the lemma is implied if we can show

$$S_1 + S_2 \leq \frac{\lambda}{2}, \quad S_3 + S_4 \leq \sigma\lambda, \quad \text{w.h.p.}$$

But  $\lambda \gtrsim \sigma(\gamma + \sigma)(\sqrt{np} + |n_1 - n_2| \sqrt{p}) \log^3 n$  by assumption of Theorem 3. Therefore, the lemma follows if each of the following bounds holds w.h.p.

$$\begin{aligned} S_1 &\lesssim \sigma\gamma\sqrt{np}\log^3 n, \\ S_2 &\lesssim \sigma^2\sqrt{np}\log^3 n, \\ S_3 &\lesssim \sigma^2\gamma(\sqrt{np} + |n_1 - n_2| \sqrt{p})\log^2 n, \\ S_4 &\lesssim \sigma^3\sqrt{np}\log^2 n. \end{aligned}$$

We now prove these bounds.

a) *Term  $S_1$ :* Note that  $\gamma \geq \|\beta_1^* - \beta_2^*\|_2$ , so the desired bound on  $S_1$  follows from the lemma below, which is proved in Section F-F.

**Lemma 11.** *Suppose  $\beta_1^* - \beta_2^*$  is supported on the first coordinate. Then w.h.p.*

$$S_1 \lesssim \|\beta_1^* - \beta_2^*\|_2 \sigma\sqrt{np}\log^3 n.$$

b) *Term  $S_2$ :* By definition, we have

$$S_2 = 2 \sum_b \left\| \sum_{i=1}^{n_b} (e_{b,i}^2 - \sigma^2) \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top \right\|.$$

Here each  $e_{b,i}^2 - \sigma^2$  is zero-mean,  $\lesssim \sigma^2 \log n$  almost surely, and has variance  $\lesssim \sigma^4$ . The quantity inside the spectral norm is the sum of independent zero-mean bounded matrices. An application of the Matrix Bernstein inequality [28] gives

$$\left\| \sum_{i=1}^{n_b} (e_{b,i}^2 - \sigma^2) \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top \right\| \lesssim \sigma^2 \sqrt{n_b p} \log^3 n_b,$$

for each  $b = 1, 2$ . The desired bound follows.

c) *Term  $S_3$ :* We have

$$\begin{aligned} S_3/4 &= \sqrt{p} \left\| \sum_b \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) (-\mathbf{e}_b \circ (\mathbf{X}_b \delta_b^*)) \right\|_\infty \\ &= \sqrt{p} \left\| \sum_b \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b^2) \mathbf{X}_b \delta_b^* \right\|_\infty \\ &= \sqrt{p} \max_{l \in [p]} \left| \sum_b (\mathbf{e}_b^2 \circ \mathbf{X}_{b,l})^\top \mathbf{X}_b \delta_b^* \right|, \end{aligned}$$

where  $\mathbf{X}_{b,l}$  is the  $l^{th}$  column of  $\mathbf{X}_b$ . WLOG, we assume  $n_1 \geq n_2$ . Observe that for each  $l \in [p]$ ,

$$\begin{aligned} &\sum_b (\mathbf{e}_b^2 \circ \mathbf{X}_{b,l})^\top \mathbf{X}_b \delta_b^* \\ &= \underbrace{\sum_{i=1}^{n_2} (e_{1,i}^2 \mathbf{x}_{1,i}(l) \mathbf{x}_{1,i}^\top - e_{2,i}^2 \mathbf{x}_{2,i}(l) \mathbf{x}_{2,i}^\top) \delta_1^*}_{S_{3,1,l}} \\ &\quad + \underbrace{\sum_{i=n_2+1}^{n_1} e_{1,i}^2 \mathbf{x}_{1,i}(l) \mathbf{x}_{1,i}^\top \delta_1^*}_{S_{3,2,l}}. \end{aligned}$$

Let  $\epsilon_i$  be the  $i^{th}$  standard basis vector in  $\mathbb{R}^n$ . The term  $S_{3,1,l}$  can be written as

$$S_{3,1,l} = \sum_{i=1}^{n_2} \left( \mathbf{x}_{1,i}^\top (e_{1,i}^2 \epsilon_l \delta_1^{*\top}) \mathbf{x}_{1,i} - \mathbf{x}_{2,i}^\top (e_{2,i}^2 \epsilon_l \delta_1^{*\top}) \mathbf{x}_{2,i} \right) \\ = \mathbf{x}^\top \mathbf{G} \mathbf{x},$$

where we defined

$$\mathbf{x}^\top := [e_{1,1} \mathbf{x}_{1,1}^\top \ e_{1,2} \mathbf{x}_{1,2}^\top \ \cdots \ e_{1,n_2} \mathbf{x}_{1,n_2}^\top \\ e_{2,1} \mathbf{x}_{2,1}^\top \ e_{2,2} \mathbf{x}_{2,2}^\top \ \cdots \ e_{2,n_2} \mathbf{x}_{2,n_2}^\top] \in \mathbb{R}^{2n_2 p}$$

and

$$\mathbf{G} := \text{diag}(\epsilon_l \delta_1^{*\top}, \epsilon_l \delta_1^{*\top}, \dots, \epsilon_l \delta_1^{*\top}, \\ -\epsilon_l \delta_1^{*\top}, -\epsilon_l \delta_1^{*\top}, \dots, -\epsilon_l \delta_1^{*\top}) \in \mathbb{R}^{2n_2 p \times 2n_2 p},$$

in other words,  $\mathbf{G}$  is the block-diagonal matrix with  $\{\pm \epsilon_l \delta_1^{*\top}\}$  on its diagonal. Note that  $\mathbb{E} S_{3,1,l} = 0$ , and the entries of  $\mathbf{x}$  are i.i.d. sub-Gaussian with parameter bounded by  $\sigma \sqrt{\log n}$ . Using the Hanson-Wright inequality (e.g., [25]), we obtain w.h.p.

$$\max_{l \in [p]} |S_{3,1,l}| \lesssim \|\mathbf{G}\|_F \sigma^2 \log^2 n \leq \sigma^2 \sqrt{2n} \gamma \log^2 n.$$

Since  $\delta_1^*$  is supported on the first coordinate, the term  $S_{3,2,l}$  can be bounded w.h.p. by

$$\max_{l \in [p]} |S_{3,2,l}| = \max_{l \in [p]} \left| \sum_{i=n_2+1}^{n_1} e_{1,i}^2 \mathbf{x}_{1,i}(l) \mathbf{x}_{1,i}(1) \delta_1^*(1) \right| \\ \lesssim (n_1 - n_2) \sigma^2 \gamma \log^2 n$$

using Hoeffding's inequality. It follows that w.h.p.,

$$S_3 \lesssim \sqrt{p} \max_{l \in [p]} (|S_{3,1,l}| + |S_{3,2,l}|) \\ \lesssim \sigma^2 \gamma (\sqrt{np} + |n_1 - n_2| \sqrt{p}) \log^2 n.$$

d) Term  $S_4$ :: We have w.h.p.

$$S_4/4 \leq \sqrt{p} \sum_b \|\mathbf{X}^\top (\mathbf{e}_b \circ \mathbf{w}_{2,b})\|_\infty \\ \stackrel{(a)}{\lesssim} \sqrt{p \log n} \sum_b \|\mathbf{e}_b \circ \mathbf{w}_{2,b}\|_2 \\ = \sqrt{p \log n} \sum_b \|\mathbf{e}_b^3 - \sigma^2 \mathbf{e}_b\|_2 \\ \stackrel{(b)}{\lesssim} \sigma^3 \sqrt{np} \log^2 n,$$

where in (a) we use the independence between  $\mathbf{X}$  and  $\mathbf{e}_b \circ \mathbf{w}_{2,b}$  and the standard sub-Gaussian concentration inequality (e.g., [30]), and (b) follows from the boundedness of  $\mathbf{e}$ .

## APPENDIX B PROOF OF LEMMA 5

The proofs for  $b = 1$  and  $2$  are identical, so we omit the subscript  $b$ . WLOG we may assume  $\sigma = 1$ . Our proof generalizes the proof of an RIP-type result in [16].

Fix  $\mathbf{Z}$  and  $\mathbf{z}$ . Let  $\xi_j := \langle \mathbf{B}_j, \mathbf{Z} \rangle$  and  $\nu := \|\mathbf{Z}\|_F$ . We already know that  $\xi_j$  is a sub-exponential random variable with  $\|\xi_j\|_{\psi_1} \leq c_1 \nu$  and  $\|\xi_j - \mathbb{E}[\xi_j]\|_{\psi_1} \leq 2c_1 \nu$ .

Let  $\gamma_j = \langle \mathbf{d}_j, \mathbf{z} \rangle$  and  $\omega := \|\mathbf{z}\|_2$ . It is easy to check that  $\gamma_j$  is sub-Gaussian with  $\|\gamma_j\|_{\psi_2} \leq c_1 \mu$ . It follows that  $\|\xi_j - \gamma_j\|_{\psi_1} \leq c_1 (\nu + \omega)$ .

Note that

$$\|\mathbf{BZ} - \mathbf{Dz}\|_1 = \sum_{j=1}^{n/2} \frac{2}{n} |\xi_j - \gamma_j|.$$

Therefore, applying the Bernstein-type inequality for the sum of sub-exponential variables [30], we obtain

$$\mathbb{P} [|\|\mathbf{BZ} - \mathbf{Dz}\|_1 - \mathbb{E}|\xi_j - \gamma_j|| \geq t] \\ \leq 2 \exp \left[ -c \min \left\{ \frac{t^2}{c_2(\nu + \mu)^2/n}, \frac{t}{c_2(\nu + \mu)/n} \right\} \right].$$

Setting  $t = (\nu + \sigma \omega)/c_3$  for any  $c_3 > 1$ , we get

$$\mathbb{P} \left[ |\|\mathbf{BZ} - \mathbf{Dz}\|_1 - \mathbb{E}|\xi_j - \gamma_j|| \geq \frac{\nu + \omega}{c_3} \right] \leq 2 \exp [-c_4 n]. \quad (48)$$

But sub-exponentiality implies

$$\mathbb{E} [|\xi_j - \gamma_j|] \leq \|\xi_j - \gamma_j\|_{\psi_1} \leq c_2 (\nu + \mu),$$

and

$$\mathbb{P} \left[ \|\mathbf{BZ} - \mathbf{Dz}\|_1 \geq \left( c_2 + \frac{1}{c_3} \right) (\nu + \omega) \right] \leq 2 \exp [-c_4 n].$$

Now, note that

$$\mathbb{E} [|\xi_j - \gamma_j|] \geq \sqrt{\frac{(\mathbb{E} [(\xi_j - \gamma_j)^2])^3}{\mathbb{E} [(\xi_j - \gamma_j)^4]}}.$$

We bound the numerator and denominator. By sub-exponentiality, we have  $\mathbb{E} [(\xi_j - \gamma_j)^4] \leq c_5 (\nu + \omega)^4$ . On the other hand, note that

$$\begin{aligned} & \mathbb{E} (\xi_j - \gamma_j)^2 \\ &= \mathbb{E} (\langle \mathbf{B}_j, \mathbf{Z} \rangle - \langle \mathbf{d}_j, \mathbf{z} \rangle)^2 \\ &= \mathbb{E} \langle \mathbf{B}_j, \mathbf{Z} \rangle^2 + \mathbb{E} \langle \mathbf{d}_j, \mathbf{z} \rangle^2 - 2\mathbb{E} [\langle \mathbf{B}_j, \mathbf{Z} \rangle \langle \mathbf{d}_j, \mathbf{z} \rangle] \\ &= \mathbb{E} \langle \mathbf{B}_j, \mathbf{Z} \rangle^2 + \mathbb{E} \langle \mathbf{d}_j \mathbf{d}_j^\top, \mathbf{z} \mathbf{z}^\top \rangle \\ &\quad - 2\mathbb{E} [\langle \mathbf{B}_j, \mathbf{Z} \rangle \langle \mathbf{e}_{2j} \mathbf{x}_{2j} - \mathbf{e}_{2j-1} \mathbf{x}_{2j-1}, \mathbf{z} \rangle] \\ &= \mathbb{E} \langle \mathbf{B}_j, \mathbf{Z} \rangle^2 + \mathbb{E} \langle \mathbf{d}_j \mathbf{d}_j^\top, \mathbf{z} \mathbf{z}^\top \rangle - 2\mathbb{E} [e_{2j}] \mathbb{E} [\langle \mathbf{B}_j, \mathbf{Z} \rangle \langle \mathbf{e}_{2j}, \mathbf{z} \rangle] \\ &\quad - 2\mathbb{E} [e_{2j-1}] \mathbb{E} [\langle \mathbf{B}_{j-1}, \mathbf{Z} \rangle \langle \mathbf{e}_{2j-1}, \mathbf{z} \rangle] \\ &= \mathbb{E} \langle \mathbf{B}_j, \mathbf{Z} \rangle^2 + \mathbb{E} \langle \mathbf{d}_j \mathbf{d}_j^\top, \mathbf{z} \mathbf{z}^\top \rangle, \end{aligned}$$

where in the last equality we use the fact that  $\{e_i\}$  are independent of  $\{x_i\}$  and  $\mathbb{E}[e_i] = 0$  for all  $i$ . We already know

$$\begin{aligned} \mathbb{E} \langle \mathbf{B}_j, \mathbf{Z} \rangle^2 &= \langle \mathbb{E} [\langle \mathbf{B}_j, \mathbf{Z} \rangle \mathbf{B}_j], \mathbf{Z} \rangle \\ &= 4 \|\mathbf{Z}\|_F^2 + 2(\mu - 3) \|\text{diag}(\mathbf{Z})\|_F^2 \\ &\geq 2(\mu - 1) \|\mathbf{Z}\|_F^2. \end{aligned}$$

Some calculation shows that

$$\begin{aligned} \mathbb{E} \langle \mathbf{d}_j \mathbf{d}_j^\top, \mathbf{z} \mathbf{z}^\top \rangle &= \langle \mathbb{E} [e_{2j}^2 \mathbf{x}_{2j} \mathbf{x}_{2j}^\top + e_{2j}^2 \mathbf{x}_{2j} \mathbf{x}_{2j}^\top], \mathbf{z} \mathbf{z}^\top \rangle \\ &= 2 \langle \mathbf{I}, \mathbf{z} \mathbf{z}^\top \rangle = 2 \|\mathbf{z}\|^2. \end{aligned}$$

It follows that

$$\mathbb{E} (\xi_j - \gamma_j)^2 \geq 2(\mu - 1) \|\mathbf{Z}\|_F^2 + 2 \|\mathbf{z}\|^2 \geq c_6 (\nu^2 + \omega^2),$$

where the inequality holds when  $\mu > 1$ . We therefore obtain

$$\mathbb{E} [|\xi_j - \gamma_j|] \geq c_7 \frac{\sqrt{(\nu^2 + \omega^2)^3}}{(\nu + \omega)^2} \geq c_8(\nu + \omega).$$

Substituting back to (48), we get

$$\mathbb{P} \left[ \|\mathcal{B}\mathbf{Z} - \mathbf{D}\mathbf{z}\|_1 \leq \left( c_8 - \frac{1}{c_3} \right) (\nu + \omega) \right] \leq 2 \exp [-c_4 n].$$

To complete the proof of the lemma, we use an  $\epsilon$ -net argument. Let  $\mathcal{S}_r$  be the set

$$\left\{ (\mathbf{Z}, \mathbf{z}) \in \mathbb{R}^{p \times p} \times \mathbb{R}^p : \text{rank}(\mathbf{Z}) \leq r, \|\mathbf{Z}\|_F^2 + \|\mathbf{z}\|_2^2 = 1 \right\}.$$

We need the following lemma, which is proved in Section F-G.

**Lemma 12.** For each  $\epsilon > 0$  and  $r \geq 1$ , there exists a set  $\mathcal{N}_r(\epsilon)$  with  $|\mathcal{N}_r(\epsilon)| \leq \left(\frac{40}{\epsilon}\right)^{10pr}$  which is an  $\epsilon$ -covering of  $\mathcal{S}_r$ , meaning that for all  $(\mathbf{Z}, \mathbf{z}) \in \mathcal{S}_r$ , there exists  $(\tilde{\mathbf{Z}}, \tilde{\mathbf{z}}) \in \mathcal{N}_r(\epsilon)$  such that

$$\sqrt{\|\tilde{\mathbf{Z}} - \mathbf{Z}\|_F^2 + \|\tilde{\mathbf{z}} - \mathbf{z}\|_2^2} \leq \epsilon.$$

Note that  $\frac{1}{\sqrt{2}} (\|\mathbf{Z}\|_F + \|\mathbf{z}\|_2) \leq \sqrt{\|\mathbf{Z}\|_F^2 + \|\mathbf{z}\|_2^2} \leq \|\mathbf{Z}\|_F + \|\mathbf{z}\|_2$  for all  $\mathbf{Z}$  and  $\mathbf{z}$ . Therefore, up to a change of constant, it suffices to prove Lemma 5 for all  $(\mathbf{Z}, \mathbf{z})$  in  $\mathcal{S}_r$ . By the union bound and Lemma 12, we have

$$\begin{aligned} & \mathbb{P} \left( \max_{(\tilde{\mathbf{Z}}, \tilde{\mathbf{z}}) \in \mathcal{N}_r(\epsilon)} \|\mathcal{B}\tilde{\mathbf{Z}} - \mathbf{D}\tilde{\mathbf{z}}\|_1 \leq 2 \left( c_2 + \frac{1}{c_3} \right) \right) \\ & \geq 1 - |\mathcal{N}_r(\epsilon)| \cdot \exp(-c_4 n) \geq 1 - \exp(-c_4 n/2), \end{aligned}$$

when  $n \geq (2/c_4) \cdot 10pr \log(40/\epsilon)$ . On this event, we have

$$\begin{aligned} \bar{M} &:= \sup_{(\mathbf{Z}, \mathbf{z}) \in \mathcal{S}_r} \|\mathcal{B}\mathbf{Z} - \mathbf{D}\mathbf{z}\|_1 \\ &\leq \max_{(\tilde{\mathbf{Z}}, \tilde{\mathbf{z}}) \in \mathcal{N}_r(\epsilon)} \|\mathcal{B}\tilde{\mathbf{Z}} - \mathbf{D}\tilde{\mathbf{z}}\|_1 \\ &\quad + \sup_{(\mathbf{Z}, \mathbf{z}) \in \mathcal{S}_r} \|\mathcal{B}(\mathbf{Z} - \tilde{\mathbf{Z}}) - \mathbf{D}(\mathbf{z} - \tilde{\mathbf{z}})\|_1 \\ &\leq 2 \left( c_2 + \frac{1}{c_3} \right) + \sup_{\mathbf{Z} \in \mathcal{S}_r} \sqrt{\|\mathbf{Z} - \tilde{\mathbf{Z}}\|_F^2 + \|\mathbf{z} - \tilde{\mathbf{z}}\|_2^2} \\ &\quad \times \sup_{(\mathbf{Z}', \mathbf{z}') \in \mathcal{S}_{2r}} \|\mathcal{B}\mathbf{Z}' - \mathbf{D}\mathbf{z}'\|_1 \\ &\leq 2 \left( c_2 + \frac{1}{c_3} \right) + \epsilon \sup_{(\mathbf{Z}', \mathbf{z}') \in \mathcal{S}_{2r}} \|\mathcal{B}\mathbf{Z}' - \mathbf{D}\mathbf{z}'\|_1. \end{aligned}$$

Note that for  $(\mathbf{Z}', \mathbf{z}') \in \mathcal{S}_{2r}$ , we can write  $\mathbf{Z}' = \mathbf{Z}'_1 + \mathbf{Z}'_2$  such that  $\mathbf{Z}'_1, \mathbf{Z}'_2$  both have rank  $r$  and  $1 = \|\mathbf{Z}'\|_F \geq \max\{\|\mathbf{Z}'_1\|_F, \|\mathbf{Z}'_2\|_F\}$ . So

$$\begin{aligned} & \sup_{(\mathbf{Z}', \mathbf{z}') \in \mathcal{S}_{2r}} \|\mathcal{B}\mathbf{Z}' - \mathbf{D}\mathbf{z}'\|_1 \\ & \leq \sup_{\mathbf{Z}'_1 \in \mathcal{S}_{2r}} \|\mathcal{B}\mathbf{Z}'_1 - \mathbf{D}\mathbf{z}'\|_1 + \sup_{\mathbf{Z}'_2 \in \mathcal{S}_{2r}} \|\mathcal{B}\mathbf{Z}'_2\|_1 \quad (49) \\ & \leq 2\bar{M}. \end{aligned}$$

Combining the last two displayed equations and choosing  $\epsilon = \frac{1}{4}$ , we obtain

$$\bar{M} \leq \bar{\delta} := \frac{2}{1 - 2\epsilon} \left( c_2 + \frac{1}{c_3} \right),$$

with probability at least  $1 - \exp(-c_9 n)$ . Note that  $\bar{\delta}$  is a constant independent of  $p$  and  $r$  (but it might depend on  $\mu := \mathbb{E} [(\mathbf{x}_i)_l^4]$ ).

For a possibly different  $\epsilon'$ , we have

$$\begin{aligned} & \inf_{(\mathbf{Z}, \mathbf{z}) \in \mathcal{S}_r} \|\mathcal{B}\mathbf{Z} - \mathbf{D}\mathbf{z}\|_1 \\ & \geq \min_{(\tilde{\mathbf{Z}}, \tilde{\mathbf{z}}) \in \mathcal{N}_r(\epsilon')} \|\mathcal{B}\tilde{\mathbf{Z}} - \tilde{\mathbf{z}}\|_1 \\ & \quad - \sup_{(\mathbf{Z}, \mathbf{z}) \in \mathcal{S}_r} \|\mathcal{B}(\mathbf{Z} - \tilde{\mathbf{Z}}) - \mathbf{D}(\mathbf{z} - \tilde{\mathbf{z}})\|_1. \end{aligned}$$

By the union bound, we have

$$\begin{aligned} & \mathbb{P} \left( \min_{(\tilde{\mathbf{Z}}, \tilde{\mathbf{z}}) \in \mathcal{N}_r(\epsilon')} \|\mathcal{B}\tilde{\mathbf{Z}} - \tilde{\mathbf{z}}\|_1 \geq \left( c_7 - \frac{1}{c_3} \right) \right) \\ & \geq 1 - \exp(-c_4 n + 10pr \log(40/\epsilon')) \\ & \geq 1 - \exp(-c_4 n/2), \end{aligned}$$

provided  $n \geq (2/c_4) \cdot 10pr \log(40/\epsilon')$ . On this event, we have

$$\begin{aligned} & \inf_{(\mathbf{Z}, \mathbf{z}) \in \mathcal{S}_r} \|\mathcal{B}\mathbf{Z} - \mathbf{D}\mathbf{z}\|_1 \stackrel{(a)}{\geq} \left( c_7 - \frac{1}{c_3} \right) - 2\epsilon' \bar{M} \\ & \stackrel{(b)}{\geq} \left( c_7 - \frac{1}{c_3} \right) - 2\epsilon' \bar{\delta}, \end{aligned}$$

where (a) follows from (49) and (b) follows from the the upper-bound on  $\bar{M}$  we just established. We complete the proof by choosing  $\epsilon'$  to be a sufficiently small constant such that  $\underline{\delta} := \left( c_7 - \frac{1}{c_3} \right) - 2\epsilon' \bar{\delta} > 0$ .

## APPENDIX C PROOF OF LEMMA 8

Let  $\mathbf{H} = \mathbf{K}' - \mathbf{g}'$  and  $\mathbf{h} = \mathbf{g}' - \mathbf{g}$ . We observe that

$$\begin{aligned} & Q_n(\mathbf{K}', \mathbf{g}'; \mathbf{K}, \mathbf{g}) \\ & = \sum_{i=1}^n (-\langle \mathbf{x}_i \mathbf{x}_i, \mathbf{H} \rangle + 2y_i \langle \mathbf{x}_i, \mathbf{h} \rangle)^2 - \sum_{i=1}^n 4\sigma^2 \langle \mathbf{x}_i, \mathbf{h} \rangle^2 \\ & \leq \sum_{i=1}^n 2\langle \mathbf{x}_i \mathbf{x}_i, \mathbf{H} \rangle^2 + 8 \max_{i \in [n]} \{y_i^2\} \left( \sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{h} \rangle^2 \right) \end{aligned}$$

By standard concentration results, we have  $\Pr(\sum_{i=1}^n \langle \mathbf{x}_i, \mathbf{h} \rangle^2 \geq 1.1 \|\mathbf{h}\|_2^2) \leq c_1 \exp(-c_2 p)$  for any  $\mathbf{h} \in \mathbb{R}^p$ . We thus obtain that w.h.p.

$$\begin{aligned} & Q_n(\mathbf{K}', \mathbf{g}'; \mathbf{K}, \mathbf{g}) \\ & \leq \sum_{i=1}^n 2\langle \mathbf{x}_i \mathbf{x}_i, \mathbf{H} \rangle^2 + 9n \cdot \max_{i \in [n]} \{y_i^2\} \cdot \|\mathbf{h}\|_2^2. \end{aligned}$$

We need the following lemma, which is proved in Section F-H, to bound the first term on the right hand side of the above inequality.

**Lemma 13.** Suppose  $\{\mathbf{x}_i\}_{i=1}^n$  are i.i.d.  $p$ -dimensional centered sub-Gaussian random vectors with norm  $\|\mathbf{x}_i\|_{\psi_2} \leq c$  for some constant  $c$  and identity covariance matrix  $\mathbb{E}(\mathbf{x}_i \mathbf{x}_i^\top) = \mathbf{I}$ . There exist constants  $c_i$  such that for any matrix  $\mathbf{H} \in \mathbb{R}^{p \times p}$ , we have

$$\frac{1}{n} \sum_{i=1}^n \langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{H} \rangle^2 \leq c_1 p \log n \cdot \|\mathbf{H}\|_* \cdot \|\mathbf{H}\|_F,$$

with probability at least  $1 - c_2 n^{-1}$  under condition  $n \geq c_3 p$ .

Note that  $\|\mathbf{H}\|_* \leq \sqrt{p} \|\mathbf{H}\|_F$  for any  $\mathbf{H}$ . We have that  $\sum_{i \in [n]} \langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{H} \rangle^2 \lesssim n \sqrt{p}^3 \log n \|\mathbf{H}\|_F^2$ . Also note that

$$\Pr \left( \max_{i \in [n]} |y_i| > t \right) \leq n \exp \left( 1 - t^2 / (\gamma^2 + \sigma^2) \right).$$

Setting  $t \asymp (\gamma + \sigma) \sqrt{\log n}$  in the above inequality leads to our result.

#### APPENDIX D PROOF OF THEOREM 7

We need some additional notation. Let  $\mathbf{z} := (z_1, z_2, \dots, z_n)^\top \in \{0, 1\}^n$  be the vector of hidden labels with  $z_i = 1$  if and only if  $i \in \mathcal{I}_1$ . We use  $\mathbf{y}(\theta^*, \mathbf{X}, \mathbf{e}, \mathbf{z})$  to denote the value of the response vector  $\mathbf{y}$  given  $\theta^*$ ,  $\mathbf{X}$ ,  $\mathbf{e}$  and  $\mathbf{z}$ , that is,

$$\mathbf{y}(\theta^*, \mathbf{X}, \mathbf{e}, \mathbf{z}) = \mathbf{z} \circ (\mathbf{X} \beta_1^*) + (\mathbf{1} - \mathbf{z}) \circ (\mathbf{X} \beta_2^*) + \mathbf{e},$$

where  $\mathbf{1}$  is the all-ones vector in  $\mathbb{R}^n$  and  $\circ$  denotes element-wise product.

By standard results, we know that with probability at least  $1 - n^{-10}$ ,

$$\|\mathbf{X} \alpha\|_2 \leq 2\sqrt{n} \|\alpha\|_2, \forall \alpha \in \mathbb{R}^p. \quad (50)$$

Hence it suffices to prove (11) in the theorem statement assuming (50) holds.

Let  $\mathbf{v}$  be an arbitrary unit vector in  $\mathbb{R}^p$ . We define  $\delta := c_0 \frac{\epsilon}{\sqrt{n}}$ ,  $\theta_1 := (\frac{1}{2} \underline{\gamma} \mathbf{v}, -\frac{1}{2} \underline{\gamma} \mathbf{v})$  and  $\theta_2 = (\frac{1}{2} \underline{\gamma} \mathbf{v} + \delta \mathbf{v}, -\frac{1}{2} \underline{\gamma} \mathbf{v} - \delta \mathbf{v})$ . Note that  $\theta_1, \theta_2 \in \Theta(\underline{\gamma})$  as long as  $c_0$  is sufficiently small, and  $\rho(\theta_1, \theta_2) = 2\delta$ . We further define  $\mathbf{e}_1 := \mathbf{0}$  and  $\mathbf{e}_2 := -\delta(2\mathbf{z} - \mathbf{1}) \circ (\mathbf{X} \mathbf{v})$ . Note that  $\|\mathbf{e}_2\| \leq 2\sqrt{n}\delta \leq \epsilon$  by (50), so  $\mathbf{e}_1, \mathbf{e}_2 \in \mathbb{B}(\epsilon)$ . If we set  $\mathbf{y}_i = \mathbf{y}(\theta_i, \mathbf{X}, \mathbf{e}_i, \mathbf{z})$  for  $i = 1, 2$ , then we have

$$\begin{aligned} \mathbf{y}_2 &= \mathbf{z} \circ \left( \mathbf{X} \left( \frac{1}{2} \underline{\gamma} \mathbf{v} + \delta \mathbf{v} \right) \right) + (\mathbf{1} - \mathbf{z}) \circ \left( \mathbf{X} \left( -\frac{1}{2} \underline{\gamma} \mathbf{v} - \delta \mathbf{v} \right) \right) + \mathbf{e}_2 \\ &= (2\mathbf{z} - \mathbf{1}) \circ \left( \mathbf{X} \left( \frac{1}{2} \underline{\gamma} \mathbf{v} + \delta \mathbf{v} \right) \right) - \delta(2\mathbf{z} - \mathbf{1}) \circ (\mathbf{X} \mathbf{v}) \\ &= (2\mathbf{z} - \mathbf{1}) \circ \left( \mathbf{X} \left( \frac{1}{2} \underline{\gamma} \mathbf{v} \right) \right) + \mathbf{e}_1 \\ &= \mathbf{y}_1, \end{aligned}$$

which holds for any  $\mathbf{X}$  and  $\mathbf{z}$ . Therefore, for any  $\hat{\theta}$ , we have

$$\begin{aligned} &\sup_{\theta^* \in \Theta(\underline{\gamma})} \sup_{\mathbf{e} \in \mathbb{B}(\epsilon)} \rho \left( \hat{\theta}(\mathbf{X}, \mathbf{y}), \theta^* \right) \\ &\geq \frac{1}{2} \rho \left( \hat{\theta}(\mathbf{X}, \mathbf{y}_1), \theta_1 \right) + \frac{1}{2} \rho \left( \hat{\theta}(\mathbf{X}, \mathbf{y}_2), \theta_2 \right) \\ &= \frac{1}{2} \rho \left( \hat{\theta}(\mathbf{X}, \mathbf{y}_1), \theta_1 \right) + \frac{1}{2} \rho \left( \hat{\theta}(\mathbf{X}, \mathbf{y}_1), \theta_2 \right) \\ &\geq \frac{1}{2} \rho(\theta_1, \theta_2) \\ &= \delta, \end{aligned}$$

where the second inequality holds because  $\rho$  is a metric and satisfies the triangle inequality. Taking the infimum over  $\hat{\theta}$  proves the theorem.

#### APPENDIX E PROOF OF THEOREM 8

Throughout the proof we set  $\kappa := \frac{1}{2} \underline{\gamma}$ .

##### A. Part 1 of the Theorem

We prove the first part of the theorem by establishing a lower-bound for standard linear regression. Set  $\delta_1 := c_0 \sigma \sqrt{\frac{p-1}{n}}$ , and define the (semi)-metric  $\rho_1(\cdot, \cdot)$  by  $\rho_1(\beta, \beta') = \min \{ \|\beta - \beta'\|, \|\beta + \beta'\| \}$ . We begin by constructing a  $\delta_1$ -packing set  $\Phi_1 := \{\beta_1, \dots, \beta_M\}$  of  $\mathbb{G}^p(\kappa) := \{\beta \in \mathbb{R}^p : \|\beta\| \geq \kappa\}$  in the metric  $\rho_1$ . We need a packing set of the hypercube  $\{0, 1\}^{p-1}$  in the Hamming distance.

**Lemma 14.** For  $p \geq 16$ , there exists  $\{\xi_1, \dots, \xi_M\} \subset \{0, 1\}^{p-1}$  such that  $M \geq 2^{(p-1)/16}$  and

$$\min \{ \|\xi_i - \xi_j\|_0, \|\xi_i + \xi_j\|_0 \} \geq \frac{p-1}{16}, \quad \forall 1 \leq i < j \leq M.$$

See Section F-I for the proof. Let  $\tau := 2c_0 \sigma \sqrt{\frac{1}{n}}$  for some absolute constant  $c_0 > 0$  that is sufficiently small, and  $\kappa_0^2 := \kappa^2 - (p-1)\tau^2$ . Note that  $\kappa_0 \geq 0$  since  $\underline{\gamma} \geq \sigma$  by assumption. For  $i = 1, \dots, M$ , we set

$$\beta_i = \kappa_0 \epsilon_p + \sum_{j=1}^{p-1} (2\xi_i(j) - 1) \tau \epsilon_j,$$

where  $\epsilon_j$  is the  $j^{th}$  standard basis in  $\mathbb{R}^p$  and  $\xi_i(j)$  is the  $j^{th}$  coordinate of  $\xi_i$ . Note that  $\|\beta_i\|_2 = \kappa, \forall i \in [M]$ , so  $\Phi_1 = \{\beta_1, \dots, \beta_M\} \subset \mathbb{G}^p(\kappa)$ . We also have that for all  $1 \leq i < j \leq M$ ,

$$\|\beta_i - \beta_j\|_2^2 \leq (p-1)\tau^2 = 4c_0^2 \frac{\sigma^2(p-1)}{n}. \quad (51)$$

Moreover, we have

$$\begin{aligned} \rho^2(\beta_i, \beta_j) &= \min \left\{ \|\beta_i - \beta_j\|_2^2, \|\beta_i + \beta_j\|_2^2 \right\} \\ &\geq 4\tau^2 \min \{ \|\xi_i - \xi_j\|_0, \|\xi_i + \xi_j\|_0 \} \\ &\geq 4 \cdot 4c_0^2 \frac{\sigma^2}{n} \cdot \frac{p-1}{16} = \delta_1^2, \end{aligned} \quad (52)$$

so  $\Phi_1 = \{\beta_1, \dots, \beta_M\}$  is a  $\delta_1$ -packing of  $\mathbb{G}^p(\kappa)$  in the metric  $\rho_1$ .

Suppose  $\beta^*$  is sampled uniformly at random from the set  $\Phi_1$ . For  $i = 1, \dots, M$ , let  $\mathbb{P}_{i, \mathbf{X}}$  denote the distribution of  $\mathbf{y}$  conditioned on  $\beta^* = \beta_i$  and  $\mathbf{X}$ , and  $\mathbb{P}_i$  denote the joint distribution of  $\mathbf{X}$  and  $\mathbf{y}$  conditioned on  $\beta^* = \beta_i$ . Because  $\mathbf{X}$  is independent of  $\mathbf{z}, \mathbf{e}$  and  $\beta^*$ , we have

$$\begin{aligned} D(\mathbb{P}_i \|\mathbb{P}_{i'}) &= \mathbb{E}_{\mathbb{P}_i(\mathbf{X}, \mathbf{y})} \log \frac{p_i(\mathbf{X}, \mathbf{y})}{p_{i'}(\mathbf{X}, \mathbf{y})} \\ &= \mathbb{E}_{\mathbb{P}_i(\mathbf{X}, \mathbf{y})} \log \frac{p_i(\mathbf{y}|\mathbf{X})}{p_{i'}(\mathbf{y}|\mathbf{X})} \\ &= \mathbb{E}_{\mathbb{P}(\mathbf{X})} \left[ \mathbb{E}_{\mathbb{P}_i(\mathbf{y}|\mathbf{X})} \left[ \log \frac{p_i(\mathbf{y}|\mathbf{X})}{p_{i'}(\mathbf{y}|\mathbf{X})} \right] \right] \\ &= \mathbb{E}_{\mathbf{X}} [D(\mathbb{P}_{i, \mathbf{X}} \|\mathbb{P}_{i', \mathbf{X}})]. \end{aligned}$$

Using the above equality and the convexity of the mutual information, we get that

$$\begin{aligned} I(\boldsymbol{\beta}^*; \mathbf{X}, \mathbf{y}) &\leq \frac{1}{M^2} \sum_{1 \leq i, i' \leq M} D(\mathbb{P}_i \parallel \mathbb{P}_{i'}) \\ &= \frac{1}{M^2} \sum_{1 \leq i, i' \leq M} \mathbb{E}_{\mathbf{X}} [D(\mathbb{P}_{i, \mathbf{X}} \parallel \mathbb{P}_{i', \mathbf{X}})] \\ &= \frac{1}{M^2} \sum_{1 \leq i, i' \leq M} \mathbb{E}_{\mathbf{X}} \frac{\|\mathbf{X}\boldsymbol{\beta}_i - \mathbf{X}\boldsymbol{\beta}_{i'}\|^2}{2\sigma^2} \\ &= \frac{1}{M^2} \sum_{1 \leq i, i' \leq M} \frac{n \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2}{2\sigma^2}. \end{aligned}$$

It follows from (51) that

$$I(\boldsymbol{\beta}^*; \mathbf{X}, \mathbf{y}) \leq 8c_0^2 p \leq \frac{1}{2} (\log_2 M) / (\log_2 e) = \frac{1}{4} \log M,$$

provided  $c_0$  is sufficiently small. Following a standard argument [5, 35, 38] to transform the estimation problem into a hypothesis testing problem (cf. Eq. (46) and (47)), we obtain

$$\begin{aligned} \inf_{\hat{\boldsymbol{\beta}}} \sup_{\boldsymbol{\beta}^* \in \mathbb{G}^p(\kappa)} \mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\rho_1(\hat{\boldsymbol{\beta}}, \boldsymbol{\beta}^*)] \\ \geq \delta_1 \left( 1 - \frac{I(\boldsymbol{\beta}^*; \mathbf{X}, \mathbf{y}) + \log 2}{\log M} \right) \\ \geq \frac{1}{2} \delta_1 = \frac{1}{2} c_0 \sigma \sqrt{\frac{p}{n}}. \end{aligned}$$

This establishes a minimax lower bound for standard linear regression. Now observe that given any standard linear regression problem with regression vector  $\boldsymbol{\beta}^* \in \mathbb{G}^p(\kappa)$ , we can reduce it to a mixed regression problem with  $\boldsymbol{\theta}^* = (\boldsymbol{\beta}^*, -\boldsymbol{\beta}^*) \in \Theta(\underline{\gamma})$  by multiplying each  $y_i$  by a Rademacher  $\pm 1$  variable. Part 1 of the theorem hence follows.

### B. Part 2 of the Theorem

Let  $\delta_2 := 2c_0 \frac{\sigma^2}{\kappa} \sqrt{\frac{p-1}{n}}$ . We first construct a  $\delta_2$ -packing set  $\Theta_2 := \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\}$  of  $\Theta(\underline{\gamma})$  in the metric  $\rho(\cdot, \cdot)$ . Set  $\tau := 2c_0 \frac{\sigma^2}{\kappa} \sqrt{\frac{1}{n}}$  and  $\kappa_0^2 := \kappa^2 - (p-1)\tau^2$ . Note that  $\kappa_0 \geq 0$  under the assumption  $\kappa \geq c_1 \sigma \left(\frac{p}{n}\right)^{1/4}$  provided that  $c_0$  is small enough. For  $i = 1, \dots, M$ , we set  $\boldsymbol{\theta}_i := (\boldsymbol{\beta}_i, -\boldsymbol{\beta}_i)$  with

$$\boldsymbol{\beta}_i = \kappa_0 \boldsymbol{\epsilon}_p + \sum_{j=1}^{p-1} (2\xi_i(j) - 1) \tau \boldsymbol{\epsilon}_j,$$

where  $\{\xi_i\}$  are the vectors in Lemma 14. Note that  $\|\boldsymbol{\beta}_i\| = \kappa$  for all  $i$ , so  $\Theta_2 = \{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_M\} \subset \Theta(\underline{\gamma})$ . We also have that for all  $1 \leq i < i' \leq M$ ,

$$\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2 \leq p\tau^2 = 4c_0^2 \frac{\sigma^4 p}{\kappa^2 n}. \quad (53)$$

Moreover, we have

$$\begin{aligned} \rho^2(\boldsymbol{\theta}_i, \boldsymbol{\theta}_{i'}) &= 4 \min \left\{ \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2, \|\boldsymbol{\beta}_i + \boldsymbol{\beta}_{i'}\|^2 \right\} \\ &\geq 16\tau^2 \min \{\|\xi_i - \xi_{i'}\|_0, \|\xi_i + \xi_{i'}\|_0\} \\ &\geq 16 \cdot 4c_0^2 \frac{\sigma^4}{\kappa^2 n} \cdot \frac{p-1}{16} = \delta_2^2, \end{aligned} \quad (54)$$

so  $\Theta_2 = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\}$  forms a  $\delta_2$ -packing of  $\Theta(\underline{\gamma})$  in the metric  $\rho$ .

Suppose  $\boldsymbol{\theta}^*$  is sampled uniformly at random from the set  $\Theta_2$ . For  $i = 1, \dots, M$ , let  $\mathbb{P}_{i, \mathbf{X}}^{(j)}$  denote the distribution of  $\mathbf{y}_j$  conditioned on  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_i$  and  $\mathbf{X}$ ,  $\mathbb{P}_{i, \mathbf{X}}$  denote the distribution of  $\mathbf{y}$  conditioned on  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_i$  and  $\mathbf{X}$ , and  $\mathbb{P}_i$  denote the joint distribution of  $\mathbf{X}$  and  $\mathbf{y}$  conditioned on  $\boldsymbol{\theta}^* = \boldsymbol{\theta}_i$ . We need the following bound on the KL divergence between two mixtures of univariate Gaussians. For any  $a > 0$ , we use  $\mathbb{Q}_a$  to denote the distribution of the equal-weighted mixture of two Gaussian distributions  $\mathcal{N}(a, \sigma^2)$  and  $\mathcal{N}(-a, \sigma^2)$ .

**Lemma 15.** *The following bound holds for any  $u, v \geq 0$ :*

$$\begin{aligned} D(\mathbb{Q}_u \parallel \mathbb{Q}_v) \\ \leq \frac{u^2 - v^2}{2\sigma^4} u^2 + \frac{v^3 \max\{0, v-u\}}{2\sigma^8} (u^4 + 6u^2\sigma^2 + 3\sigma^4). \end{aligned}$$

See Section F-J for the proof. Note that  $\mathbb{P}_{i, \mathbf{X}}^{(j)} = \mathbb{Q}_{|\mathbf{x}_j^\top \boldsymbol{\beta}_i|}$ . Using  $\mathbb{P}_{i, \mathbf{X}} = \otimes_{j=1}^n \mathbb{P}_{i, \mathbf{X}}^{(j)}$  and the above lemma, we have

$$\begin{aligned} &\mathbb{E}_{\mathbf{X}} D(\mathbb{P}_{i, \mathbf{X}} \parallel \mathbb{P}_{i', \mathbf{X}}) \\ &= \sum_{j=1}^n \mathbb{E}_{\mathbf{X}} D\left(\mathbb{P}_{i, \mathbf{X}}^{(j)} \parallel \mathbb{P}_{i', \mathbf{X}}^{(j)}\right) \\ &\leq n \mathbb{E} \frac{|\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 - |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^2}{2\sigma^4} |\mathbf{x}_j^\top \boldsymbol{\beta}_i|^2 \\ &\quad + n \mathbb{E}_{\mathbf{X}} \frac{|\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^3 \max\{0, |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}| - |\mathbf{x}_1^\top \boldsymbol{\beta}_i|\}}{2\sigma^8} \\ &\quad \times \left( |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + 6|\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \sigma^2 + 3\sigma^4 \right). \end{aligned}$$

To bound the expectations in the last RHS, we need a simple technical lemma proved in Section F-K.

**Lemma 16.** *Suppose  $\mathbf{x} \in \mathbb{R}^p$  has i.i.d. standard Gaussian components, and  $\boldsymbol{\alpha}, \boldsymbol{\beta} \in \mathbb{R}^p$  are any fixed vectors with  $\|\boldsymbol{\alpha}\|_2 = \|\boldsymbol{\beta}\|_2$ . There exists an absolute constant  $\bar{c}$  such that for any non-negative integers  $k, l$  with  $k + l \leq 8$ ,*

$$\mathbb{E} |\mathbf{x}^\top \boldsymbol{\alpha}|^k |\mathbf{x}^\top \boldsymbol{\beta}|^l \leq \bar{c} \|\boldsymbol{\alpha}\|^k \|\boldsymbol{\beta}\|^l.$$

Moreover, we have

$$\mathbb{E}_{\mathbf{X}} \left[ (|\mathbf{x}^\top \boldsymbol{\alpha}|^2 - |\mathbf{x}^\top \boldsymbol{\beta}|^2) |\mathbf{x}^\top \boldsymbol{\alpha}|^2 \right] \leq 2 \|\boldsymbol{\alpha}\| \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|^2.$$

$$\mathbb{E} (|\mathbf{x}^\top \boldsymbol{\alpha}|^2 - |\mathbf{x}^\top \boldsymbol{\beta}|^2)^2 \leq \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|^4.$$

Using the above lemma and the fact that  $\|\boldsymbol{\beta}_i\|_2 = \|\boldsymbol{\beta}_{i'}\|_2 = \kappa$  for all  $1 \leq i < i' \leq M$ , we have

$$\mathbb{E}_{\mathbf{X}} \frac{|\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 - |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^2}{2\sigma^4} |\mathbf{x}_j^\top \boldsymbol{\beta}_i|^2 \leq \frac{1}{2\sigma^4} \kappa^2 \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2,$$

and for some universal constant  $c' > 0$ ,

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}} \frac{|\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^3 \max \{0, |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}| - |\mathbf{x}_1^\top \boldsymbol{\beta}_i|\}}{2\sigma^8} \\ & \times \left( |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + 6 |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \sigma^2 + 3\sigma^4 \right) \\ & \leq \frac{1}{2\sigma^8} \mathbb{E}_{\mathbf{X}} \max \left\{ 0, |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^2 - |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \right\} |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^2 \\ & \times \left( |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + 6 |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \sigma^2 + 3\sigma^4 \right) \\ & \stackrel{(a)}{\leq} \frac{1}{2\sigma^4} \sqrt{\mathbb{E}_{\mathbf{X}} \left( |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^2 - |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \right)^2 \cdot \frac{1}{\sigma^8} \mathbb{E}_{\mathbf{X}} |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^4} \\ & \times \sqrt{\left( |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + 6 |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \sigma^2 + 3\sigma^4 \right)^2} \\ & \stackrel{(b)}{\leq} \frac{1}{2\sigma^4} \sqrt{\|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^4 \cdot c'^2 \|\boldsymbol{\beta}_{i'}\|^4} = \frac{c'}{2\sigma^4} \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2 \kappa^2, \end{aligned}$$

where (a) follows from the Cauchy-Schwarz inequality, and (b) follows from the first and third inequalities in Lemma 16 as well as  $\|\boldsymbol{\beta}_i\| = \|\boldsymbol{\beta}_{i'}\| = \kappa \leq \sigma$ . It follows that

$$\mathbb{E}_{\mathbf{X}} D(\mathbb{P}_{i,\mathbf{X}} \|\mathbb{P}_{i',\mathbf{X}}) \leq n \cdot \frac{c' \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|^2 \kappa^2}{\sigma^4} \leq c'' p,$$

where the last inequality follows from (53) and  $c''$  can be made sufficiently small by choosing  $c_0$  small enough. We therefore obtain

$$\begin{aligned} I(\boldsymbol{\theta}^*; \mathbf{X}, \mathbf{y}) & \leq \frac{1}{M^2} \sum_{1 \leq i, i' \leq M} D(\mathbb{P}_i \|\mathbb{P}_{i'}) \\ & = \frac{1}{M} \sum_{1 \leq i, i' \leq M} \mathbb{E}_{\mathbf{X}} [D(\mathbb{P}_{i,\mathbf{X}} \|\mathbb{P}_{i',\mathbf{X}})] \leq c'' p \leq \frac{1}{4} \log M \end{aligned}$$

using  $M \geq 2^{(p-1)/16}$ . Following a standard argument [38, 35, 5] to transform the estimation problem into a hypothesis testing problem (cf. Eq. (46) and (47)), we obtain

$$\begin{aligned} & \inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\rho(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)] \\ & \geq \delta_2 \left( 1 - \frac{I(\boldsymbol{\theta}^*; \mathbf{X}, \mathbf{y}) + \log 2}{\log M} \right) \\ & \geq \frac{1}{2} \delta_2 = c_0 \frac{\sigma^2}{\kappa} \sqrt{\frac{p}{n}}. \end{aligned}$$

### C. Part 3 of the Theorem

The proof follows similar lines as Part 2. Let  $\delta_3 := 2c_0\sigma \left( \frac{p}{n} \right)^{1/4}$ . Again we first construct a  $\delta_3$ -packing set  $\Theta_3 := (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M)$  of  $\Theta(\underline{\gamma})$  in the metric  $\rho(\cdot, \cdot)$ . Set  $\tau := \frac{2c_0\sigma}{\sqrt{p-1}} \left( \frac{p}{n} \right)^{1/4}$ . For  $i = 1, \dots, M$ , we set  $\boldsymbol{\theta}_i = (\boldsymbol{\beta}_i, -\boldsymbol{\beta}_i)$  with

$$\boldsymbol{\beta}_i = \sum_{j=1}^{p-1} (2\xi_i(j) - 1) \tau \boldsymbol{\epsilon}_j,$$

where  $\{\xi_i\}$  are the vectors from Lemma 14. Note that  $\|\boldsymbol{\beta}_i\|_2 = \sqrt{p-1}\tau = 2c_0\sigma \left( \frac{p}{n} \right)^{1/4} \geq c_1\sigma \left( \frac{p}{n} \right)^{1/4} \geq \kappa$  provided  $c_1$  is

sufficiently small, so  $\Theta_3 = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\} \subset \Theta(\underline{\gamma})$ . We also have for all  $1 \leq i < i' \leq M$ ,

$$\begin{aligned} \rho^2(\boldsymbol{\beta}_i, \boldsymbol{\beta}_{i'}) & = 4 \min \left\{ \|\boldsymbol{\beta}_i - \boldsymbol{\beta}_{i'}\|_2^2, \|\boldsymbol{\beta}_i + \boldsymbol{\beta}_{i'}\|_2^2 \right\} \\ & \geq 16\tau^2 \min \{ \|\xi_i - \xi_{i'}\|_0, \|\xi_i + \xi_{i'}\|_0 \} \\ & = 16 \cdot \frac{4c_0^2\sigma^2}{p-1} \sqrt{\frac{p}{n}} \cdot \frac{p-1}{16} \geq \delta_3^2, \end{aligned} \quad (55)$$

so  $\Theta_3 = \{\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_M\}$  is a  $\delta_3$ -packing of  $\Theta(\underline{\gamma})$  in the metric  $\rho$ .

Suppose  $\boldsymbol{\theta}^*$  is sampled uniformly at random from the set  $\Theta_2$ . Define  $\mathbb{P}_{i,\mathbf{X}}, \mathbb{P}_{i',\mathbf{X}}^{(j)}$  and  $\mathbb{P}_i$  as in the proof of Part 2 of the theorem. We have

$$\begin{aligned} & \mathbb{E}_{\mathbf{X}} D(\mathbb{P}_{i,\mathbf{X}} \|\mathbb{P}_{i',\mathbf{X}}) \\ & = \sum_{j=1}^n \mathbb{E}_{\mathbf{X}} D(\mathbb{P}_{i,\mathbf{X}}^{(j)} \|\mathbb{P}_{i',\mathbf{X}}^{(j)}) \\ & \stackrel{(a)}{\leq} n \mathbb{E}_{\mathbf{X}} \frac{|\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 - |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^2}{2\sigma^4} |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \\ & \quad + n \mathbb{E}_{\mathbf{X}} \frac{|\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^3 \max \{0, |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}| - |\mathbf{x}_1^\top \boldsymbol{\beta}_i|\}}{2\sigma^8} \\ & \quad \times \left( |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + 6 |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \sigma^2 + 3\sigma^4 \right) \\ & \leq \frac{n}{2\sigma^4} \mathbb{E}_{\mathbf{X}} |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + \frac{n}{2\sigma^8} \mathbb{E}_{\mathbf{X}} |\mathbf{x}_1^\top \boldsymbol{\beta}_{i'}|^4 \\ & \quad \times \left( |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^4 + 6 |\mathbf{x}_1^\top \boldsymbol{\beta}_i|^2 \sigma^2 + 3\sigma^4 \right) \\ & \stackrel{(b)}{\leq} \frac{n}{2\sigma^4} \bar{c} \|\boldsymbol{\beta}_i\|^4 + \frac{n}{2\sigma^8} \bar{c} \|\boldsymbol{\beta}_{i'}\|^4 \left( \|\boldsymbol{\beta}_i\|^4 + 6\sigma^2 \|\boldsymbol{\beta}_i\|^2 + 9\sigma^4 \right) \\ & \stackrel{(c)}{\leq} c' p, \end{aligned}$$

where (a) follows from Lemma 15, (b) follows from Lemma 16, (c) follows from  $\|\boldsymbol{\beta}_i\| = 2c_0\sigma \left( \frac{p}{n} \right)^{1/4} \leq \sigma, \forall i$ , and  $c'$  is a sufficiently small absolute constant. It follows that

$$I(\boldsymbol{\theta}^*; \mathbf{X}, \mathbf{y}) \leq \frac{1}{M} \sum_{1 \leq i, i' \leq M} \mathbb{E}_{\mathbf{X}} D(\mathbb{P}_i \|\mathbb{P}_{i'}) \leq c' p \leq \frac{1}{4} \log M,$$

since  $M \geq 2^{(p-1)/8}$ . Again transforming the estimation problem into a hypothesis testing problem, we obtain

$$\begin{aligned} & \inf_{\hat{\boldsymbol{\theta}}} \sup_{\boldsymbol{\theta}^* \in \Theta(\underline{\gamma})} \mathbb{E}_{\mathbf{X}, \mathbf{z}, \mathbf{e}} [\rho(\hat{\boldsymbol{\theta}}, \boldsymbol{\theta}^*)] \\ & \geq \delta_3 \left( 1 - \frac{I(\boldsymbol{\theta}^*; \mathbf{X}, \mathbf{y}) + \log 2}{\log M} \right) \\ & \geq \frac{1}{2} \delta_3 = c_0 \sigma \left( \frac{p}{n} \right)^{1/4}. \end{aligned}$$

## APPENDIX F PROOFS OF AUXILIARY RESULTS

### A. Proof of Lemma 2

Simple algebra shows that

$$\begin{aligned} & \sum_b \left\| \hat{\mathbf{H}}_T - 2\boldsymbol{\beta}_b^* \hat{\mathbf{h}}^\top \right\|_F^2 \\ & = 2 \left\| \hat{\mathbf{H}}_T - (\boldsymbol{\beta}_1^* + \boldsymbol{\beta}_2^*) \hat{\mathbf{h}}^\top \right\|_F^2 + 2 \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2^2 \|\hat{\mathbf{h}}\|_2^2 \\ & \geq 2 \|\boldsymbol{\beta}_1^* - \boldsymbol{\beta}_2^*\|_2^2 \|\hat{\mathbf{h}}\|_2^2 \geq \alpha (\|\boldsymbol{\beta}_1^*\|_2 + \|\boldsymbol{\beta}_2^*\|)^2 \|\hat{\mathbf{h}}\|_2^2, \end{aligned}$$

and

$$\begin{aligned}
& \sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}} \right\|_F^2 \\
&= 4 \left( \|\beta_1^*\|_2^2 + \|\beta_2^*\|_2^2 \right) \left\| \hat{\mathbf{h}} - \frac{\hat{\mathbf{H}}_T(\beta_1^* + \beta_2^*)}{2\|\beta_1^*\|_2^2 + 2\|\beta_2^*\|_2^2} \right\|_2^2 \\
&\quad + \frac{2 \left( \|\beta_1^*\|_2^2 + \|\beta_2^*\|_2^2 \right) \left\| \hat{\mathbf{H}}_T \right\|_F^2 - \left\| \hat{\mathbf{H}}_T(\beta_1^* + \beta_2^*) \right\|_2^2}{\|\beta_1^*\|_2^2 + \|\beta_2^*\|_2^2} \\
&\stackrel{(a)}{\geq} \frac{2 \left( \|\beta_1^*\|_2^2 + \|\beta_2^*\|_2^2 \right) \left\| \hat{\mathbf{H}}_T \right\|_F^2 - \left\| \hat{\mathbf{H}}_T \right\|_F^2 \|\beta_1^* + \beta_2^*\|_2^2}{\|\beta_1^*\|_2^2 + \|\beta_2^*\|_2^2} \\
&= \alpha \left\| \hat{\mathbf{H}}_T \right\|_F^2,
\end{aligned}$$

where the inequality (a) follows from  $\left\| \hat{\mathbf{H}}_T \right\| \leq \left\| \hat{\mathbf{H}}_T \right\|_F$ . Combining the last two displayed equations with the simple inequality

$$\sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}} \right\|_F \geq \sqrt{\sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}} \right\|_F^2},$$

we obtain

$$\begin{aligned}
\sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}} \right\|_F &\geq \sqrt{\alpha} (\|\beta_1^*\|_2 + \|\beta_2^*\|_2) \|\hat{\mathbf{h}}\|_2, \\
\sum_b \left\| \hat{\mathbf{H}}_T - 2\beta_b^* \hat{\mathbf{h}} \right\|_F &\geq \sqrt{\alpha} \left\| \hat{\mathbf{H}}_T \right\|_F.
\end{aligned}$$

### B. Proof of Lemma 3

Without loss of generality, we may assume  $\sigma = 1$ . Set  $L := \sqrt{c \log n}$  for some  $c$  sufficiently large. For each  $i \in [n]$ , we define the event  $\mathcal{E}_i = \{|e_i| \leq L\}$  and the truncated random variables

$$\bar{e}_i = e_i \mathbf{1}(\mathcal{E}_i),$$

where  $\mathbf{1}(\cdot)$  is the indicator function and  $c$  is some sufficiently large numeric constant. Let  $m_i := \mathbb{E}[e_i \mathbf{1}(\mathcal{E}_i^c)]$  and  $s_i := \sqrt{\mathbb{E}[e_i^2 \mathbf{1}(\mathcal{E}_i^c)]}$ . WLOG we assume  $m_i \geq 0$ . Note that the following equation holds almost surely:

$$e_i^2 \mathbf{1}(\mathcal{E}_i^c) = |e_i| \cdot |e_i| \mathbf{1}(\mathcal{E}_i^c) \geq L \cdot |e_i| \mathbf{1}(\mathcal{E}_i^c) \geq L \cdot e_i \mathbf{1}(\mathcal{E}_i^c).$$

Taking the expectation of both sides gives  $s_i^2 \geq Lm_i$ . We further define

$$\tilde{e}_i := \bar{e}_i + L\epsilon_i^+ - L\epsilon_i^-,$$

where  $\epsilon_i^+$  and  $\epsilon_i^-$  are independent random variables distributed as  $\text{Ber}(\nu_i^+)$  and  $\text{Ber}(\nu_i^-)$ , respectively, with

$$\nu_i^+ := \frac{1}{2} \left( \frac{m_i}{L} + \frac{s_i^2}{L^2} \right), \quad \nu_i^- := \frac{1}{2} \left( -\frac{m_i}{L} + \frac{s_i^2}{L^2} \right).$$

Note that  $m_i \geq 0$  and  $s_i^2 \geq Lm_i$  implies that  $\nu_i^+, \nu_i^- \geq 0$ . We show below that  $\nu_i^+, \nu_i^- \leq 1$  so the random variables  $\epsilon_i^+$  and  $\epsilon_i^-$  are well-defined.

With this setup, we now characterize the distribution of  $\tilde{e}_i$ . Note that

$$\begin{aligned}
\mathbb{E}[L\epsilon_i^+ - L\epsilon_i^-] &= m_i, \\
\mathbb{E}[(L\epsilon_i^+)^2 + (L\epsilon_i^-)^2] &= s_i^2,
\end{aligned}$$

which means

$$\begin{aligned}
\mathbb{E}[\tilde{e}_i] &= \mathbb{E}[\bar{e}_i] + \mathbb{E}[e_i \mathbf{1}(\mathcal{E}_i^c)] = \mathbb{E}[e_i] = 0, \\
\text{Var}[\tilde{e}_i^2] &= \mathbb{E}[\bar{e}_i^2] + \mathbb{E}[e_i^2 \mathbf{1}(\mathcal{E}_i^c)] = \mathbb{E}[e_i^2] = 1.
\end{aligned}$$

Moreover,  $\tilde{e}_i$  is bounded by  $3L$  almost surely, which means it is sub-Gaussian with sub-Gaussian norm at most  $3L$ . Also note that

$$\begin{aligned}
m_i &\leq \mathbb{E}[|e_i \mathbf{1}(\mathcal{E}_i^c)|] \\
&= \int_0^\infty \mathbb{P}(|e_i \mathbf{1}(\mathcal{E}_i^c)| \geq t) dt \\
&= L \cdot \mathbb{P}(|e_i| \geq L) + \int_L^\infty \mathbb{P}(|e_i| \geq t) dt \\
&\leq \sqrt{c \log n} \frac{1}{n^{c_1}} + \int_L^\infty e^{1-t^2} dt \leq \frac{4}{n^{c_2}}
\end{aligned}$$

for some large constant  $c_1$  and  $c_2$  by sub-Gaussianity of  $e_i$ . A similar calculation gives

$$s_i^2 = \mathbb{E}[e_i^2 \mathbf{1}(\mathcal{E}_i^c)] \lesssim \frac{1}{n^{c_2}}.$$

This implies  $\nu_i^+, \nu_i^- \lesssim \frac{1}{n^{c_2}}$ , or equivalently  $L\epsilon_i^+ - L\epsilon_i^- = 0$  w.h.p. We also have  $\bar{e}_i = e_i$  w.h.p. by sub-Gaussianity of  $e_i$ . It follows that  $\tilde{e}_i = \bar{e}_i + L\epsilon_i^+ - L\epsilon_i^- = e_i$  w.h.p. Moreover,  $\tilde{e}_i$  and  $e_i$  have the same mean and variance.

We define the variables  $\{(\tilde{\mathbf{x}}_i)_l, i \in [n], l \in [p]\}$  in a similar manner. Each  $(\tilde{\mathbf{x}}_i)_l$  is sub-Gaussian, bounded by  $L$  a.s., has mean 0 and variance 1, and equals  $(\mathbf{x}_i)_l$  w.h.p.

Now suppose the conclusion of Theorem 3 holds w.h.p. for the program (7) with  $\{(\tilde{\mathbf{x}}_i, \tilde{y}_i)\}$  as the input, where  $\tilde{y}_i = \tilde{\mathbf{x}}_i^\top \beta_b^* + \tilde{e}_i$  for all  $i \in \mathcal{I}_b$  and  $b = 1, 2$ . We know that  $\mathbf{e} = \tilde{\mathbf{e}}$  and  $\mathbf{x}_i = \tilde{\mathbf{x}}_i, \forall i$  with high probability. On this event, the program above is identical to the original program with  $\{(\mathbf{x}_i, y_i)\}$  as the input. Therefore, the conclusion of the theorem also holds w.h.p. for the original program.

### C. Proof of Lemma 6

We need to bound

$$\begin{aligned}
W &= \sqrt{p} \left\| \sum_b \mathbf{X}_b^\top \text{diag}(\mathbf{e}_b) \mathbf{w}_b + \sum_b \sigma^2 \mathbf{X}_b^\top \mathbf{X}_b \delta_b^* \right. \\
&\quad \left. + \sum_b 2\sigma^2 \mathbf{X}_b^\top \mathbf{e}_b \right\|_\infty \\
&\leq \sqrt{p} \sum_b \left\| \sum_{i \in \mathcal{I}_b} (\sigma^2 - e_i^2) \langle \mathbf{x}_i, \delta_b^* \rangle \mathbf{x}_i + \sum_{i \in \mathcal{I}_b} (e_i^2 + \sigma^2) e_i \mathbf{x}_i \right\|_\infty \\
&\leq \sqrt{p} \sum_b \underbrace{\left\| \sum_{i \in \mathcal{I}_b} (\sigma^2 - e_i^2) \langle \mathbf{x}_i, \delta_b^* \rangle \mathbf{x}_i \right\|_\infty}_{S_9} \\
&\quad + \underbrace{\sqrt{p} \sum_b \left\| \sum_{i \in \mathcal{I}_b} (e_i^2 + \sigma^2) e_i \mathbf{x}_i \right\|_\infty}_{S_{10}}.
\end{aligned}$$

Lemma 3 makes it possible to assume the boundedness of  $e_i$  with the loss of only small probability, we thus have

that  $(e_i^2 + \sigma^2)e_i \mathbf{x}_i$  is sub-Gaussian random vector with Orlicz norm  $O(\sigma^3 \sqrt{\log n})$ . Therefore the standard sub-Gaussian concentration result leads to that the following inequality holds with high probability

$$S_{10} \lesssim \sigma^3 \sqrt{n_1} \log^2 n_1 + \sigma^3 \sqrt{n_2} \log^2 n_2 \lesssim \sigma^3 \sqrt{n} \log^2 n,$$

where the last inequality follows from  $\min\{n_1/n_2, n_2/n_1\} = \Omega(1)$ .

Since  $\langle \mathbf{x}_i, \delta_b^* \rangle$  is sub-Gaussian random variable with Orlicz norm  $O(\gamma)$  and  $(\sigma^2 - e_i^2)\mathbf{x}_i$  is sub-Gaussian random vector with Orlicz norm  $O(\sigma^2 \log n)$ , their product  $(\sigma^2 - e_i^2)\langle \mathbf{x}_i, \delta_b^* \rangle \mathbf{x}_i$  is centered sub-exponential random vector with Orlicz norm  $O(\gamma \sigma^2 \log n)$ . By standard concentration inequality of sub-exponential random variables (see, e.g., Corollary 5.17 in [30]), we obtain that for every  $\epsilon > 0$  and  $K = \gamma \sigma^2 \log n$ ,

$$\begin{aligned} \Pr \left( \left\| \sum_{i \in \mathcal{I}_b} (\sigma^2 - e_i^2) \langle \mathbf{x}_i, \delta_b^* \rangle \mathbf{x}_i + \sum_{i \in \mathcal{I}_b} (e_i^2 + \sigma^2) e_i \mathbf{x}_i \right\|_\infty > \epsilon n_b \right) \\ \leq 2p \exp \left( -c \min \left\{ \frac{\epsilon^2}{K^2}, \frac{\epsilon}{K} \right\} n_b \right). \end{aligned}$$

For each  $b = 1, 2$ , by setting  $\epsilon = c_1 K \sqrt{\log p/n_b}$ , we have that when  $n_b \gtrsim \log p$ , with probability at least  $1 - 1/p$ ,

$$S_9 \lesssim \gamma \sigma^2 (\sqrt{n_1} + \sqrt{n_2}) \sqrt{\log p} \log n \lesssim \gamma \sigma^2 \sqrt{n} \log^2 n.$$

Putting all ingredients together, we have that  $W \lesssim \gamma \sigma^2 \sqrt{pn} \log^2 n$  with high probability.

#### D. Proof of Lemma 7

Using the inequalities  $\|\mathbf{K}^{(t)}\|_* \leq \|\mathbf{K}^*\|_*$  and  $\|\hat{\mathbf{K}}\|_* \leq \|\mathbf{K}^*\|_*$ , we have

$$\begin{aligned} \left\| \left( \mathbf{K}^{(t)} - \mathbf{K}^* \right)_{T^\perp} \right\|_* &\leq \left\| \left( \mathbf{K}^{(t)} - \mathbf{K}^* \right)_T \right\|_*, \quad \text{and} \\ \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_* &\leq \left\| \hat{\mathbf{H}}_T \right\|_*. \end{aligned} \tag{56}$$

For the term  $\mathbf{H}^{(t)}$  we have

$$\begin{aligned} \left\| \mathbf{H}^{(t)} \right\|_* &\leq \left\| \mathbf{K}^{(t)} - \mathbf{K}^* \right\|_* + \left\| \hat{\mathbf{H}} \right\|_* \\ &\stackrel{(a)}{=} \left\| \left( \mathbf{K}^{(t)} - \mathbf{K}^* \right)_{T^\perp} \right\|_* + \left\| \left( \mathbf{K}^{(t)} - \mathbf{K}^* \right)_T \right\|_* \\ &\quad + \left\| \hat{\mathbf{H}}_{T^\perp} \right\|_* + \left\| \hat{\mathbf{H}}_T \right\|_* \\ &\stackrel{(b)}{\leq} 2 \left\| \left( \mathbf{K}^{(t)} - \mathbf{K}^* \right)_T \right\|_* + 2 \left\| \hat{\mathbf{H}}_T \right\|_* \\ &\stackrel{(c)}{\leq} 2 \left\| \mathbf{H}_T^{(t)} \right\|_* + 4 \left\| \hat{\mathbf{H}}_T \right\|_*, \end{aligned}$$

where the inequality (b) follows from (56); (a) and (c) are from the triangle inequality.

Also note that  $\|\mathbf{H}^{(t)}\|_* \geq \left\| \mathbf{H}_{T^\perp}^{(t)} \right\|_* - \left\| \mathbf{H}_T^{(t)} \right\|_*$ . Putting the lower and upper bounds of  $\|\mathbf{H}^{(t)}\|_*$  together, we complete our proof.

#### E. Proof of Lemma 9

Let

$$\begin{aligned} \Psi_t(\mathbf{K}, \mathbf{g}) \\ := \left\langle \nabla_{\mathbf{K}} \mathcal{L}_n^{(t)}, \mathbf{K} \right\rangle + \left\langle \nabla_{\mathbf{g}} \mathcal{L}_n^{(t)}, \mathbf{g} \right\rangle + \frac{\eta}{2} \left\| \mathbf{K} - \mathbf{K}^{(t)} \right\|_F^2 \\ + \frac{\eta \gamma^2}{2} \left\| \mathbf{g} - \mathbf{g}^{(t)} \right\|_2^2. \end{aligned}$$

From the optimality of  $(\hat{\mathbf{K}}, \hat{\mathbf{g}})$ , we have

$$\begin{aligned} &\left\langle \nabla_{\mathbf{K}} \Psi_t(\mathbf{K}^{(t+1)}, \mathbf{g}^{(t+1)}), \hat{\mathbf{K}} - \mathbf{K}^{(t+1)} \right\rangle \\ &+ \left\langle \nabla_{\mathbf{g}} \Psi_t(\mathbf{K}^{(t+1)}, \mathbf{g}^{(t+1)}), \hat{\mathbf{g}} - \mathbf{g}^{(t+1)} \right\rangle \geq 0. \end{aligned}$$

We thus have

$$\begin{aligned} &\left\langle \nabla_{\mathbf{K}} \mathcal{L}_n^{(t)}, \hat{\mathbf{K}} - \mathbf{K}^{(t+1)} \right\rangle + \left\langle \nabla_{\mathbf{g}} \mathcal{L}_n^{(t)}, \hat{\mathbf{g}} - \mathbf{g}^{(t+1)} \right\rangle \\ &\geq \eta \left\langle \mathbf{K}^{(t)} - \mathbf{K}^{(t+1)}, \hat{\mathbf{K}} - \mathbf{K}^{(t+1)} \right\rangle \\ &\quad + \eta \gamma^2 \left\langle \mathbf{g}^{(t)} - \mathbf{g}^{(t+1)}, \hat{\mathbf{g}} - \mathbf{g}^{(t+1)} \right\rangle. \end{aligned} \tag{57}$$

Using the first condition (lower bound), we have

$$\begin{aligned} &\mathcal{L}_n(\hat{\mathbf{K}}, \hat{\mathbf{g}}) \\ &\geq \mathcal{L}_n^{(t)} + \left\langle \nabla_{\mathbf{K}} \mathcal{L}_n^{(t)}, \hat{\mathbf{K}} - \mathbf{K}^{(t)} \right\rangle + \left\langle \nabla_{\mathbf{g}} \mathcal{L}_n^{(t)}, \hat{\mathbf{g}} - \mathbf{g}^{(t)} \right\rangle \\ &\quad + \underline{Q}_n(\mathbf{H}^{(t)}, \mathbf{h}^{(t)}) \\ &\stackrel{(a)}{\geq} \mathcal{L}_n^{(t)} + \left\langle \nabla_{\mathbf{K}} \mathcal{L}_n^{(t)}, \mathbf{K}^{(t+1)} - \mathbf{K}^{(t)} \right\rangle \\ &\quad + \left\langle \nabla_{\mathbf{g}} \mathcal{L}_n^{(t)}, \mathbf{g}^{(t+1)} - \mathbf{g}^{(t)} \right\rangle + \underline{Q}_n(\mathbf{H}^{(t)}, \mathbf{h}^{(t)}) \\ &\quad + \eta \left\langle \mathbf{K}^{(t)} - \mathbf{K}^{(t+1)}, \hat{\mathbf{K}} - \mathbf{K}^{(t+1)} \right\rangle \\ &\quad + \eta \gamma^2 \left\langle \mathbf{g}^{(t)} - \mathbf{g}^{(t+1)}, \hat{\mathbf{g}} - \mathbf{g}^{(t+1)} \right\rangle. \end{aligned}$$

Now applying the second condition (upper bound), we have

$$\begin{aligned} &\mathcal{L}_n(\hat{\mathbf{K}}, \hat{\mathbf{g}}) \\ &\geq \mathcal{L}_n^{(t+1)} - \overline{Q}_n(\Delta^{(t)}, \delta^{(t)}) + \underline{Q}_n(\mathbf{H}^{(t)}, \mathbf{h}^{(t)}) \\ &\quad + \eta \left\langle \mathbf{K}^{(t)} - \mathbf{K}^{(t+1)}, \hat{\mathbf{K}} - \mathbf{K}^{(t+1)} \right\rangle \\ &\quad + \eta \gamma^2 \left\langle \mathbf{g}^{(t)} - \mathbf{g}^{(t+1)}, \hat{\mathbf{g}} - \mathbf{g}^{(t+1)} \right\rangle. \end{aligned}$$

Applying  $\mathcal{L}_n^{(t+1)} \geq \mathcal{L}_n(\hat{\mathbf{K}}, \hat{\mathbf{g}})$  and rearranging terms yield

$$\begin{aligned} &\eta \left\langle \mathbf{K}^{(t)} - \mathbf{K}^{(t+1)}, \hat{\mathbf{K}} - \mathbf{K}^{(t)} \right\rangle + \eta \gamma^2 \left\langle \mathbf{g}^{(t)} - \mathbf{g}^{(t+1)}, \hat{\mathbf{g}} - \mathbf{g}^{(t)} \right\rangle \\ &\leq \overline{Q}_n(\Delta^{(t)}, \delta^{(t)}) - \underline{Q}_n(\mathbf{H}^{(t)}, \mathbf{h}^{(t)}) - \eta \|\Delta^{(t)}\|_F^2 \\ &\quad - \eta \gamma^2 \|\delta^{(t)}\|_2^2. \end{aligned} \tag{58}$$

Finally, by expanding  $\|\mathbf{H}^{(t+1)}\|_F^2$  and  $\|\mathbf{h}^{(t+1)}\|_2^2$  through

$$\begin{aligned} &\|\mathbf{H}^{(t+1)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t+1)}\|_2^2 \\ &\leq \|\mathbf{H}^{(t)}\|_F^2 + \gamma^2 \|\mathbf{h}^{(t)}\|_2^2 + 2 \left\langle \mathbf{K}^{(t)} - \mathbf{K}^{(t+1)}, \hat{\mathbf{K}} - \mathbf{K}^{(t)} \right\rangle \\ &\quad + 2 \gamma^2 \left\langle \mathbf{g}^{(t)} - \mathbf{g}^{(t+1)}, \hat{\mathbf{g}} - \mathbf{g}^{(t)} \right\rangle + \|\Delta^{(t)}\|_F^2 + \gamma^2 \|\delta^{(t)}\|_2^2 \end{aligned}$$

and applying (58), we complete the proof.

### F. Proof of Lemma 11

We need to bound

$$S_1 = 2 \sum_b \left\| \sum_{i \in \mathcal{I}_b} e_{b,i} \mathbf{x}_{b,i} \mathbf{x}_{b,i}^\top \cdot \mathbf{x}_{b,i}^\top (\boldsymbol{\beta}_b^* - \boldsymbol{\beta}_{-b}^*) \right\|,$$

where  $\boldsymbol{\beta}_b^* - \boldsymbol{\beta}_{-b}^*$  is supported on the first coordinate. Because  $n_1 \asymp n_2 \asymp n$  and  $\{(e_{b,i}, \mathbf{x}_{b,i})\}$  are identically distributed, it suffices to prove w.h.p.

$$\|\mathbf{E}\| := \left\| \sum_{i=1}^n e_i \mathbf{x}_i \mathbf{x}_i^\top \cdot \mathbf{x}_i^\top \boldsymbol{\delta}_1^* \right\| \lesssim \sigma \|\boldsymbol{\delta}_1^*\|_2 \sqrt{np} \log^3 n. \quad (59)$$

Let  $\bar{\mathbf{x}}_i \in \mathbb{R}^1$  and  $\underline{\mathbf{x}}_i \in \mathbb{R}^{p-1}$  be the subvectors of  $\mathbf{x}_i$  corresponding to the first and the last  $p-1$  coordinates, respectively. We define  $\bar{\boldsymbol{\delta}}_1^*$  similarly; note that  $\|\bar{\boldsymbol{\delta}}_1^*\| = \|\boldsymbol{\delta}_1^*\|$ .

Note that  $\mathbf{E} := \sum_i e_i \mathbf{x}_i \mathbf{x}_i^\top \cdot \bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*$  due to the support of  $\boldsymbol{\delta}_1^*$ . We partition  $\mathbf{E} \in \mathbb{R}^{p \times p}$  as

$$\mathbf{E} = \begin{bmatrix} \mathbf{E}_1 & \mathbf{E}_{12} \\ \mathbf{E}_{12}^\top & \mathbf{E}_2 \end{bmatrix},$$

where  $\mathbf{E}_1 \in \mathbb{R}^{1 \times 1}$ ,  $\mathbf{E}_2 \in \mathbb{R}^{(p-1) \times (p-1)}$  and  $\mathbf{E}_{12} \in \mathbb{R}^{1 \times p}$ . We have

$$\|\mathbf{E}\| \leq \|\mathbf{E}_1\| + \|\mathbf{E}_2\| + 2\|\mathbf{E}_{12}\|.$$

We bound each term separately.

Consider  $\mathbf{E}_1 = \sum_i e_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \cdot \bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*$ . We condition on  $\{\bar{\mathbf{x}}_i\}$ . Note that  $\|\bar{\mathbf{x}}_i\|_2 \lesssim \sqrt{\log n}$  and  $|\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*| \lesssim \|\bar{\boldsymbol{\delta}}_1^*\| \sqrt{\log n}$  a.s. by boundedness of  $\mathbf{x}_i$ . Since  $\{e_i\}$  are independent of  $\{\bar{\mathbf{x}}_i\}$ , we have

$$\mathbb{P} [\|\mathbf{E}_1\| \lesssim \sigma \|\boldsymbol{\delta}_1^*\| \sqrt{n} \log^2 n | \{\bar{\mathbf{x}}_i\}] \geq 1 - n^{-10},$$

w.h.p. using Hoeffding's inequality. Integrating over  $\{\bar{\mathbf{x}}_i\}$  proves  $\|\mathbf{E}_1\| \lesssim \sigma \|\boldsymbol{\delta}_1^*\| \sqrt{n} \log^2 n$ , w.h.p.

Consider  $\mathbf{E}_2 = \sum_i e_i \underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top \cdot \bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*$ . We condition on the event  $\mathcal{F} := \{\forall i : |\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*| \lesssim \|\bar{\boldsymbol{\delta}}_1^*\| \sqrt{\log n}\}$ , which occurs with high probability and is independent of  $e_i$  and  $\underline{\mathbf{x}}_i$ . We shall apply the matrix Bernstein inequality [28]; to this end, we compute:

$$\|e_i \underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top \cdot \bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*\| \lesssim \sigma p \|\boldsymbol{\delta}_1^*\| \log^2 n, \quad \text{a.s.}$$

by boundedness, and

$$\begin{aligned} & \left\| \sum_i \mathbb{E} e_i^2 (\underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top)^2 \cdot (\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*)^2 \right\| \\ & \leq n \sigma^2 \max_i |\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*|^2 \left\| \mathbb{E} (\underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top)^2 \right\| \\ & \leq np \sigma^2 \|\boldsymbol{\delta}_1^*\|^2 \log n. \end{aligned}$$

Applying the Matrix Bernstein inequality then gives

$$\|\mathbf{E}_2\| \lesssim \sigma \|\boldsymbol{\delta}_1^*\| (p + \sqrt{np}) \log^2 n \leq \sigma \|\boldsymbol{\delta}_1^*\| \sqrt{np} \log^3 n,$$

w.h.p., where we use  $n \gtrsim p$  in the last inequality.

Consider  $\mathbf{E}_{12} = \sum_i e_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \cdot \bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*$ . We again condition on the event  $\mathcal{F}$  and use the matrix Bernstein inequality. Observe that

$$\|e_i \bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top \cdot \bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*\| \lesssim \sigma \sqrt{p} \|\boldsymbol{\delta}_1^*\| \log^2 n, \quad \text{a.s.}$$

by boundedness. Moreover, we have

$$\begin{aligned} & \left\| \sum_i \mathbb{E} e_i^2 (\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*)^2 (\underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top) (\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top) \right\| \\ & \leq n \sigma^2 \max_i |\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*|^2 \|\bar{\mathbf{x}}_i\|^2 \|\mathbb{E} \underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top\| \\ & \lesssim n \sigma^2 \|\boldsymbol{\delta}_1^*\|^2 \log^2 n \end{aligned}$$

and

$$\begin{aligned} & \left\| \sum_i \mathbb{E} e_i^2 (\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*)^2 (\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top) (\underline{\mathbf{x}}_i \underline{\mathbf{x}}_i^\top) \right\| \\ & \leq n \sigma^2 \max_i |\bar{\mathbf{x}}_i^\top \bar{\boldsymbol{\delta}}_1^*|^2 \|\bar{\mathbf{x}}_i \bar{\mathbf{x}}_i^\top\| \mathbb{E} [\underline{\mathbf{x}}_i^\top \underline{\mathbf{x}}_i] \\ & \lesssim np \sigma^2 \|\boldsymbol{\delta}_1^*\|^2 \log^2 n. \end{aligned}$$

Applying the Matrix Bernstein inequality then gives

$$\|\mathbf{E}_{12}\| \lesssim \sigma \|\boldsymbol{\delta}_1^*\| \sqrt{np} \log^3 n.$$

Combining these bounds on  $\|\mathbf{E}_i\|$ ,  $i = 1, 2, 3$ , we conclude that (59) holds w.h.p., which completes the proves of the lemma.

### G. Proof of Lemma 12

Define the sphere

$$\mathcal{T}_r(b) := \{\mathbf{Z} \in \mathbb{R}^{p \times p} : \text{rank}(\mathbf{Z}) \leq r, \|\mathbf{Z}\|_F = b\}.$$

Let  $\mathcal{M}_r(\epsilon/2, 1)$  be the smallest  $\epsilon/2$ -net of  $\mathcal{T}_r'(1)$ . We know  $|\mathcal{M}_r(\epsilon/2, 1)| \leq \left(\frac{20}{\epsilon}\right)^{6pr}$  by [8]. For any  $0 \leq b \leq 1$ , we know  $\mathcal{M}_r(\epsilon/2, b) := \{b\mathbf{Z} : \mathbf{Z} \in \mathcal{M}(\epsilon/2, 1)\}$  is an  $\epsilon/2$ -net of  $\mathcal{T}_r'(b)$ , with  $|\mathcal{M}_r(\epsilon/2, b)| = |\mathcal{M}_r(\epsilon/2, 1)| \leq \left(\frac{20}{\epsilon}\right)^{6pr}$ . Let  $k := \lfloor 2/\epsilon \rfloor \leq 2/\epsilon$ . Consider the set  $\bar{\mathcal{M}}_r(\epsilon) = \{\mathbf{0}\} \cup \bigcup_{i=1}^k \mathcal{M}_r(\epsilon/2, i\epsilon/2)$ . We claim that  $\bar{\mathcal{M}}_r(\epsilon)$  is an  $\epsilon$ -net of the ball  $\bar{\mathcal{T}}_r := \{\mathbf{Z} \in \mathbb{R}^{p \times p} : \text{rank}(\mathbf{Z}) \leq r, \|\mathbf{Z}\|_F \leq 1\}$ , with the additional property that every  $\mathbf{Z}$ 's nearest neighbor  $\tilde{\mathbf{Z}}$  in  $\bar{\mathcal{M}}_r(\epsilon)$  satisfies  $\|\tilde{\mathbf{Z}}\|_F \leq \|\mathbf{Z}\|_F$ . To see this, note that for any  $\mathbf{Z} \in \bar{\mathcal{T}}_r$ , there must be some  $0 \leq i \leq k$  such that  $i\epsilon/2 \leq \|\mathbf{Z}\|_F \leq (i+1)\epsilon/2$ . Define  $\mathbf{Z}' := i\epsilon\mathbf{Z}/(2\|\mathbf{Z}\|_F)$ , which is in  $\mathcal{T}_r(i\epsilon/2)$ . We choose  $\tilde{\mathbf{Z}}$  to be the point in  $\mathcal{M}_r(\epsilon/2, i\epsilon/2)$  that is closest to  $\mathbf{Z}'$ . We have

$$\begin{aligned} \|\tilde{\mathbf{Z}} - \mathbf{Z}\|_F & \leq \|\tilde{\mathbf{Z}} - \mathbf{Z}'\|_F + \|\mathbf{Z}' - \mathbf{Z}\|_F \\ & \leq \epsilon/2 + (\|\mathbf{Z}\|_F - i\epsilon/2) \leq \epsilon, \end{aligned}$$

and  $\|\tilde{\mathbf{Z}}\|_F = i\epsilon/2 \leq \|\mathbf{Z}\|_F$ . The cardinality of  $\bar{\mathcal{M}}_r(\epsilon)$  satisfies

$$\begin{aligned} |\bar{\mathcal{M}}_r(\epsilon)| & \leq 1 + \sum_{i=1}^k |\mathcal{M}_r(\epsilon/2, k\epsilon/2)| \\ & \leq 1 + \frac{1}{\epsilon} \left(\frac{20}{\epsilon}\right)^{6pr} \leq \left(\frac{20}{\epsilon}\right)^{7pr}. \end{aligned}$$

We know that the smallest  $\epsilon/2$ -net  $\mathcal{M}'(\epsilon/2, 1)$  of the sphere  $\mathcal{T}'(1) := \{\mathbf{z} \in \mathbb{R}^p : \|\mathbf{z}\| = 1\}$  satisfies  $|\mathcal{M}'(\epsilon/2, 1)| \leq \left(\frac{20}{\epsilon}\right)^p$ . It follows from an argument similar to above that there is an  $\epsilon$ -covering  $\bar{\mathcal{M}}'(\epsilon)$  of the ball  $\bar{\mathcal{T}}' := \{\mathbf{z} \in \mathbb{R}^p : \|\mathbf{z}\| \leq 1\}$  with cardinality  $|\bar{\mathcal{M}}'(\epsilon)| \leq \left(\frac{20}{\epsilon}\right)^{2p}$  and the property that every  $\mathbf{z}$ 's nearest neighbor  $\tilde{\mathbf{z}}$  in  $\bar{\mathcal{M}}'(\epsilon)$  satisfies  $\|\tilde{\mathbf{z}}\|_2 \leq \|\mathbf{z}\|_2$ .

Let  $\bar{\mathcal{S}}_r$  be the set

$$\{(\mathbf{Z}, \mathbf{z}) \in \mathbb{R}^{p \times p} \times \mathbb{R}^p : \text{rank}(\mathbf{Z}) \leq r, \|\mathbf{Z}\|_F^2 + \|\mathbf{z}\|_2^2 \leq 1\}.$$

We claim that  $\bar{\mathcal{N}}_r(\sqrt{2}\epsilon) := (\bar{\mathcal{M}}_r(\epsilon) \times \bar{\mathcal{M}}'(\epsilon)) \cap \bar{\mathcal{S}}_r$  is an  $\sqrt{2}\epsilon$ -net of  $\bar{\mathcal{S}}_r$ . To see this, for any  $(\mathbf{Z}, \mathbf{z}) \in \bar{\mathcal{S}}_r \subset \bar{\mathcal{T}}(r) \times \bar{\mathcal{T}}'$ , we let  $\tilde{\mathbf{Z}}(\tilde{\mathbf{z}}, \text{resp.})$  be the point in  $\bar{\mathcal{M}}_r(\epsilon)$  ( $\bar{\mathcal{M}}'(\epsilon)$ , resp.) closest to  $\mathbf{Z}$  ( $\mathbf{z}$ , resp.) We have

$$\sqrt{\|\tilde{\mathbf{Z}} - \mathbf{Z}\|_F^2 + \|\tilde{\mathbf{z}} - \mathbf{z}\|_2^2} \leq \sqrt{\epsilon^2 + \epsilon^2} = \sqrt{2}\epsilon,$$

and  $\|\tilde{\mathbf{Z}}\|_F^2 + \|\tilde{\mathbf{z}}\|_2^2 \leq \|\mathbf{Z}\|_F^2 + \|\mathbf{z}\|_2^2 \leq 1$ .

Let  $\mathcal{N}_r(\sqrt{2}\epsilon)$  be the projection of the set  $\bar{\mathcal{N}}_r(\sqrt{2}\epsilon)$  onto the sphere  $\mathcal{S}_r$ . Since projection does not increase distance, we are guaranteed that  $\mathcal{N}_r(\sqrt{2}\epsilon)$  is an  $\sqrt{2}\epsilon$ -net of  $\mathcal{S}_r$ . Moreover,

$$\begin{aligned} & |\mathcal{N}_r(\sqrt{2}\epsilon)| \\ & \leq |\bar{\mathcal{N}}_r(\sqrt{2}\epsilon)| \leq |\bar{\mathcal{M}}_r(\epsilon)| \times |\bar{\mathcal{M}}'(\epsilon)| \leq \left(\frac{20}{\epsilon}\right)^{10pr}. \end{aligned}$$

#### H. Proof of Lemma 13

For simplicity, let's assume  $\mathbf{H}$  is symmetric. Later on one can check our proof works for any general matrix  $\mathbf{H}$  as well. Suppose  $\mathbf{H} = \sum_{j \in [p]} \sigma_j \mathbf{u}_j \mathbf{u}_j^\top$ , where  $\sigma_j$  is  $j^{\text{th}}$  eigenvalue and  $\mathbf{u}_j$  is the corresponding eigenvector with unit Euclidean norm. We have

$$\begin{aligned} & \frac{1}{n} \sum_{i \in [n]} \langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{H} \rangle^2 \\ & = \frac{1}{n} \sum_{i \in [n]} \left( \sum_{j \in [p]} \sigma_j \langle \mathbf{x}_i, \mathbf{u}_j \rangle^2 \right)^2 \\ & = \frac{1}{n} \sum_{i \in [n]} \left( \sum_{j, k \in [p]} \sigma_j \sigma_k \langle \mathbf{x}_i, \mathbf{u}_j \rangle^2 \langle \mathbf{x}_i, \mathbf{u}_k \rangle^2 \right) \\ & \leq \frac{1}{n} \sum_{i \in [n]} \left( \sum_{j \in [p]} |\sigma_j| \langle \mathbf{x}_i, \mathbf{u}_j \rangle^2 \cdot \left\langle \mathbf{x}_i \mathbf{x}_i^\top, \sum_{k \in [p]} |\sigma_k| \mathbf{u}_k \mathbf{u}_k^\top \right\rangle \right) \\ & \leq \frac{1}{n} \max_{i \in [n]} \left\{ \|\mathbf{x}_i\|_2^2 \right\} \cdot \|\mathbf{H}\|_F \cdot \sum_{i \in [n]} \left( \sum_{j \in [p]} |\sigma_j| \langle \mathbf{x}_i, \mathbf{u}_j \rangle^2 \right). \end{aligned} \tag{60}$$

By using standard concentration result, for some constants  $c, c', c''$ ,  $\|n^{-1} \sum_{i \in [n]} \mathbf{x}_i \mathbf{x}_i^\top - \mathbf{I}\| \leq 0.1$  with probability at least  $1 - c \exp(-c'p)$  under assumption  $n \geq c''p$ . We thus have that w.h.p.  $\frac{1}{n} \sum_{i \in [n]} \langle \mathbf{x}_i, \mathbf{u}_j \rangle^2 \leq 1.1$  for all  $j \in [p]$ . Using this result, we continue (60) with

$$\begin{aligned} & \frac{1}{n} \sum_{i \in [n]} \langle \mathbf{x}_i \mathbf{x}_i^\top, \mathbf{H} \rangle^2 \\ & \leq 1.1 \max_{i \in [n]} \left\{ \|\mathbf{x}_i\|_2^2 \right\} \cdot \|\mathbf{H}\|_F \left( \sum_{j \in [p]} |\sigma_j| \right) \\ & = 1.1 \max_{i \in [n]} \left\{ \|\mathbf{x}_i\|_2^2 \right\} \cdot \|\mathbf{H}\|_F \cdot \|\mathbf{H}\|_*. \end{aligned}$$

For each  $j \in [p]$ , let  $x_{ij}$  denote the  $j$ th coordinate of  $\mathbf{x}_i$ . By union bound, for some constant  $c_1, c_2$ ,

$$\Pr \left( \max_{i \in [n], j \in [p]} |x_{ij}| > t \right) \leq np \cdot \exp(1 - c_2 t^2).$$

Choosing  $t \asymp \sqrt{\log(np)}$ , we thus have w.h.p.

$$\max_{i \in [n]} \left\{ \|\mathbf{x}_i\|_2^2 \right\} \lesssim p \cdot \log(np),$$

which completes our proof.

#### I. Proof of Lemma 14

We need a standard result on packing the unit hypercube.

**Lemma 17** (Varshamov-Gilbert Bound, [29]). *For  $p \geq 15$ , there exists a set  $\Omega_0 = \{\xi_1, \dots, \xi_{M_0}\} \subset \{0, 1\}^{p-1}$  such that  $M \geq 2^{(p-1)/8}$  and  $\|\xi_i - \xi_j\|_0 \geq \frac{p-1}{8}$ ,  $\forall 1 \leq i < j \leq M_0$ .*

We claim that for  $i \in [M_0]$ , there is at most one  $\bar{i} \in [M_0]$  with  $\bar{i} \neq i$  such that

$$\|\xi_i - (-\xi_{\bar{i}})\|_0 < \frac{p-1}{16}, \tag{61}$$

otherwise if there are two distinct  $i_1, i_2$  that satisfy the above inequality, then they also satisfy

$$\|\xi_{i_1} - \xi_{i_2}\|_0 \leq \|\xi_{i_1} - (-\xi_{i_1})\|_0 + \|\xi_{i_2} - (-\xi_{i_1})\|_0 < \frac{p-1}{8},$$

which contradicts Lemma 17. Consequently, for each  $i \in [M_0]$ , we use  $\bar{i}$  to denote the unique index in  $[M_0]$  that satisfies (61) if such an index exists.

We construct a new set  $\Omega \subseteq \Omega_0$  by deleting elements from  $\Omega_0$ : Sequentially for  $i = 1, 2, \dots, M$ , we delete  $\xi_{\bar{i}}$  from  $\Omega_0$  if  $\bar{i}$  exists and both  $\xi_i$  and  $\xi_{\bar{i}}$  have not been deleted. Note that at most half of the elements in  $\Omega$  are deleted in this procedure. The resulting  $\Omega = \{\xi_1, \xi_2, \dots, \xi_M\}$  thus satisfies  $M \geq 2^{(p-1)/16}$  and

$$\min \{ \|\xi_i - \xi_j\|_0, \|\xi_i + \xi_j\|_0 \} \geq \frac{p-1}{16}, \forall 1 \leq i < j \leq M.$$

#### J. Proof of Lemma 15

By rescaling, it suffices to prove the lemma for  $\sigma = 1$ . Let  $\psi(x) := \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$  be the density function of the standard Normal distribution. The density function of  $\mathbb{Q}_u$  is

$$f_u(x) = \frac{1}{2} \psi(x-u) + \frac{1}{2} \psi(x+u),$$

and the density of  $\mathbb{Q}_v$  is given similarly. We compute

$$\begin{aligned}
D(\mathbb{Q}_u \parallel \mathbb{Q}_v) &= \int_{-\infty}^{\infty} f_u(x) \log \frac{f_u(x)}{f_v(x)} dx \\
&= \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x-u)] \\
&\quad \times \log \left[ \frac{\exp\left(-\frac{(x-u)^2}{2}\right) + \exp\left(-\frac{(x+u)^2}{2}\right)}{\exp\left(-\frac{(x-v)^2}{2}\right) + \exp\left(-\frac{(x+v)^2}{2}\right)} \right] dx \\
&= \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x-u)] \\
&\quad \times \log \left[ \frac{\exp\left(xu - \frac{u^2}{2}\right) + \exp\left(-xu - \frac{u^2}{2}\right)}{\exp\left(xv - \frac{v^2}{2}\right) + \exp\left(-xv - \frac{v^2}{2}\right)} \right] dx \\
&= \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x-u)] \\
&\quad \times \log \left[ \exp\left(-\frac{u^2 - v^2}{2}\right) \frac{\exp(xu) + \exp(-xu)}{\exp(xv) + \exp(-xv)} \right] dx \\
&= \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x-u)] \\
&\quad \times \left[ -\frac{u^2 - v^2}{2} + \log \frac{\cosh(xu)}{\cosh(xv)} \right] dx \\
&= -\frac{u^2 - v^2}{2} \\
&\quad + \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x-u)] \log \frac{\cosh(xu)}{\cosh(xv)} dx. \tag{62}
\end{aligned}$$

By Taylor's Theorem, the expansion of  $\log \cosh(y)$  at the point  $a$  satisfies

$$\begin{aligned}
\log \cosh(y) &= \log \cosh(a) + (y-a) \tanh(a) + \frac{1}{2}(y-a)^2 \operatorname{sech}^2(u) \\
&\quad - \frac{1}{3}(y-a)^3 \tanh(\xi) \operatorname{sech}^2(\xi)
\end{aligned}$$

for some number  $\xi$  between  $a$  and  $y$ . Let  $w := \frac{u+v}{2}$ . We expand  $\log \cosh(xu)$  and  $\log \cosh(xv)$  separately using the above equation, which gives that for some  $\xi_1$  between  $u$  and  $w$ , and some  $\xi_2$  between  $v$  and  $w$ ,

$$\begin{aligned}
&\log \cosh(xu) - \log \cosh(xv) \\
&= x(u-v) \tanh(xw) + \frac{x^2[(u-w)^2 - (v-w)^2]}{2} \operatorname{sech}^2(xw) \\
&\quad - \frac{x^3(u-w)^3}{3} \tanh(x\xi_1) \operatorname{sech}^2(x\xi_1) \\
&\quad + \frac{x^3(v-w)^3}{3} \tanh(x\xi_2) \operatorname{sech}^2(x\xi_2) \\
&= x(u-v) \tanh\left(\frac{x(u+v)}{2}\right) + \frac{-x^3}{3} \left(\frac{u-v}{2}\right)^3 \\
&\quad \times [\tanh(x\xi_1) \operatorname{sech}^2(x\xi_1) + \tanh(x\xi_2) \operatorname{sech}^2(x\xi_2)], \tag{63}
\end{aligned}$$

where the last equality follows from  $u-w = w-v = \frac{u-v}{2}$ . We bound the RHS of (63) by distinguishing two cases.

*Case 1:  $u \geq v \geq 0$ .* Because  $\tanh(x\xi_1)$  and  $\tanh(x\xi_2)$  have the same sign as  $x^3$ , the second term in (63) is negative. Moreover, we have  $x \tanh\left(\frac{x(u+v)}{2}\right) \leq x \cdot \frac{x(u+v)}{2}$  since  $\frac{u+v}{2} \geq 0$ . It follows that

$$\log \cosh(xu) - \log \cosh(xv) \leq \frac{x^2(u-v)(u+v)}{2},$$

Substituting back to (62), we obtain

$$\begin{aligned}
D(\mathbb{Q}_u \parallel \mathbb{Q}_v) &\leq -\frac{u^2 - v^2}{2} + \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x+u)] \cdot \frac{x^2(u^2 - v^2)}{2} dx \\
&= -\frac{u^2 - v^2}{2} + \frac{u^2 - v^2}{2}(u^2 + 1) = \frac{u^2 - v^2}{2}u^2.
\end{aligned}$$

*Case 2:  $v \geq u \geq 0$ .* Let  $h(y) := \tanh(y) - y + \frac{y^3}{3}$ . Taking the first order Taylor's expansion at the origin, we know that for any  $y \geq 0$  and some  $0 \leq \xi \leq y$ ,  $h(y) = -2(\tanh(\xi) \operatorname{sech}^2(\xi) - \xi) y^2 \geq 0$  since  $\tanh(\xi) \operatorname{sech}^2(\xi) \leq \xi \cdot 1^2$  for all  $\xi \geq 0$ . This means  $\tanh(y) \geq y - \frac{y^3}{3}, \forall y \geq 0$ . Since  $u-v \leq 0$  and  $\tanh(\cdot)$  is an odd function, we have

$$\begin{aligned}
&x(u-v) \tanh(x(u+v)) \\
&\leq x(u-v) \left[ x(u+v) - \frac{1}{3}(xx(u+v))^3 \right].
\end{aligned}$$

On the other hand, we have

$$\begin{aligned}
&x [\tanh(x\xi_1) \operatorname{sech}^2(x\xi_1) + \tanh(x\xi_2) \operatorname{sech}^2(x\xi_2)] \\
&\stackrel{(a)}{\leq} x(x\xi_1 + x\xi_2) \stackrel{(b)}{\leq} x \cdot 2vx,
\end{aligned}$$

where (a) follows from  $\operatorname{sech}^2(y) \leq 1$  and  $0 \leq y \tanh(y) \leq y^2$  for all  $y$ , and (b) follows from  $\xi_1, \xi_2 \leq v$  since  $v \geq w \geq u \geq 0$ . Combining the last two displayed equations with (63), we obtain

$$\begin{aligned}
&\log \cosh(xu) - \log \cosh(xv) \\
&\leq x(u-v) \left[ \frac{x(u+v)}{2} - \frac{1}{3} \left( \frac{x(u+v)}{2} \right)^3 \right] + \frac{x^3}{3} \left( \frac{v-u}{2} \right)^3 (2vx).
\end{aligned}$$

When  $u \leq v$ , we get

$$\begin{aligned}
D(\mathbb{Q}_u \parallel \mathbb{Q}_v) &\leq -\frac{u^2 - v^2}{2} + \frac{1}{2} \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x+v)] \\
&\quad \times \left[ \frac{u^2 - v^2}{2} x^2 + \frac{v-u}{3} \left( \frac{u+v}{2} \right)^3 x^4 + \frac{2v}{3} \left( \frac{v-u}{2} \right)^3 x^4 \right] dx \\
&= -\frac{u^2 - v^2}{2} + \frac{u^2 - v^2}{2}(u^2 + 1) + \left[ \frac{(v-u)(u+v)^3}{48} + \frac{v(v-u)^3}{24} \right] \\
&\quad \times \int_{-\infty}^{\infty} [\psi(x-u) + \psi(x+u)] x^4 dx \\
&= \frac{u^2 - v^2}{2} u^2 + \left[ \frac{(v-u)(u+v)^3}{24} + \frac{2v(v-u)^3}{24} \right] (u^4 + 6u^2 + 3) \\
&\leq \frac{u^2 - v^2}{2} u^2 + (v-u) \left[ \frac{(2v)^3}{24} + \frac{2v(v)^2}{24} \right] (u^4 + 6u^2 + 3) \\
&\leq \frac{u^2 - v^2}{2} u^2 + (v-u) \frac{v^3}{2} (u^4 + 6u^2 + 3).
\end{aligned}$$

Combining the two cases, we conclude that

$$D(\mathbb{Q}_u \parallel \mathbb{Q}_v) \leq \frac{u^2 - v^2}{2} u^2 + \frac{v^3 \max\{0, v - u\}}{2} (u^4 + 6u^2 + 3).$$

### K. Proof of Lemma 16

We recall that for any standard Gaussian variable  $z \sim \mathcal{N}(0, 1)$ , there exists a universal constant  $\bar{c}$  such that  $\mathbb{E}[|z|^k] \leq \bar{c}$  for all  $k \leq 16$ . Now observe that  $\mu := \mathbf{x}^\top \boldsymbol{\alpha} \sim \mathcal{N}(0, \|\boldsymbol{\alpha}\|^2)$  and  $\nu := \mathbf{x}^\top \boldsymbol{\beta} \sim \mathcal{N}(0, \|\boldsymbol{\beta}\|^2)$ . Because  $\mathbf{x}^\top \boldsymbol{\alpha} / \|\boldsymbol{\alpha}\| \sim \mathcal{N}(0, 1)$  and  $\mathbf{x}^\top \boldsymbol{\beta} / \|\boldsymbol{\beta}\| \sim \mathcal{N}(0, 1)$ , it follows from the Cauchy-Schwarz inequality,

$$\begin{aligned} & \mathbb{E}[|\mathbf{x}^\top \boldsymbol{\alpha}|^k |\mathbf{x}^\top \boldsymbol{\beta}|^l] \\ & \leq \|\boldsymbol{\alpha}\|^k \|\boldsymbol{\beta}\|^l \sqrt{\mathbb{E}\left[\frac{\mathbf{x}^\top \boldsymbol{\alpha}}{\|\boldsymbol{\alpha}\|}\right]^{2k} \mathbb{E}\left[\frac{\mathbf{x}^\top \boldsymbol{\beta}}{\|\boldsymbol{\beta}\|}\right]^{2l}} \\ & \leq \bar{c} \|\boldsymbol{\alpha}\|^k \|\boldsymbol{\beta}\|^l. \end{aligned}$$

This proves the first inequality in the lemma.

For the second inequality in the lemma, note that

$$\begin{aligned} & \mathbb{E}\left[\left(|\mathbf{x}^\top \boldsymbol{\alpha}|^2 - |\mathbf{x}^\top \boldsymbol{\beta}|^2\right) |\mathbf{x}^\top \boldsymbol{\alpha}|^2\right] \\ & = \mathbb{E}|\mathbf{x}^\top \boldsymbol{\alpha}|^4 - \mathbb{E}|\mathbf{x}^\top \boldsymbol{\alpha}|^2 |\mathbf{x}^\top \boldsymbol{\beta}|^2 \\ & = 3\|\boldsymbol{\alpha}\|^4 - \mathbb{E}|\mathbf{x}^\top \boldsymbol{\alpha}|^2 |\mathbf{x}^\top \boldsymbol{\beta}|^2. \end{aligned}$$

But

$$\begin{aligned} & \mathbb{E}|\mathbf{x}^\top \boldsymbol{\alpha}|^2 |\mathbf{x}^\top \boldsymbol{\beta}|^2 \\ & = \mathbb{E}(\alpha_1 x_1 + \cdots + \alpha_p x_p)^2 (x_1 \beta_1 + \cdots + x_p \beta_p)^2 \\ & = \mathbb{E} \sum_{i=1}^p x_i^4 \alpha_i^2 \beta_i^2 + \mathbb{E} \sum_{i \neq j} x_i^2 x_j^2 \alpha_i^2 \beta_j^2 + 2\mathbb{E} \sum_{i \neq j} x_i^2 x_j^2 \alpha_i \alpha_j \beta_i \beta_j \\ & = 3 \sum_{i=1}^p \alpha_i^2 \beta_i^2 + \sum_{i \neq j} \alpha_i^2 \beta_j^2 + 2 \sum_{i \neq j} \alpha_i \alpha_j \beta_i \beta_j \\ & = 2 \sum_{i=1}^p \alpha_i^2 \beta_i^2 + \sum_{i,j} \alpha_i^2 \beta_j^2 + 2 \sum_{i \neq j} \alpha_i \alpha_j \beta_i \beta_j \\ & = \|\boldsymbol{\alpha}\|^2 \|\boldsymbol{\beta}\|^2 + 2 \sum_{i,j} \alpha_i \alpha_j \beta_i \beta_j \\ & = \|\boldsymbol{\alpha}\|^2 \|\boldsymbol{\beta}\|^2 + 2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2. \end{aligned} \tag{64}$$

It follows that

$$\begin{aligned} & \mathbb{E}\left[\left(|\mathbf{x}^\top \boldsymbol{\alpha}|^2 - |\mathbf{x}^\top \boldsymbol{\beta}|^2\right) |\mathbf{x}^\top \boldsymbol{\alpha}|^2\right] \\ & = 3\|\boldsymbol{\alpha}\|^4 - \|\boldsymbol{\alpha}\|^2 \|\boldsymbol{\beta}\|^2 - 2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 \\ & = 2\|\boldsymbol{\alpha}\|^4 - 2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 \\ & \leq 2\|\boldsymbol{\alpha}\|^4 + 2 \left(\|\boldsymbol{\alpha}\|^2 - \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle\right)^2 - 2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 \\ & = 4\|\boldsymbol{\alpha}\|^4 - 4\|\boldsymbol{\alpha}\|^2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle \\ & = 2\|\boldsymbol{\alpha}\|^2 \left(\|\boldsymbol{\alpha}\|^2 - 2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle + \|\boldsymbol{\beta}\|^2\right) \\ & \leq 2\|\boldsymbol{\alpha}\|^2 \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|^2. \end{aligned}$$

For the third inequality in the lemma, we use the equality (64) to obtain

$$\begin{aligned} & \mathbb{E}\left(|\mathbf{x}^\top \boldsymbol{\alpha}|^2 - |\mathbf{x}^\top \boldsymbol{\beta}|^2\right)^2 \\ & = \mathbb{E}|\mathbf{x}^\top \boldsymbol{\alpha}|^4 - 2\mathbb{E}|\mathbf{x}^\top \boldsymbol{\alpha}|^2 |\mathbf{x}^\top \boldsymbol{\beta}|^2 + \mathbb{E}|\mathbf{x}^\top \boldsymbol{\beta}|^4 \\ & = 6\|\boldsymbol{\alpha}\|^4 - 2\|\boldsymbol{\alpha}\|^2 \|\boldsymbol{\beta}\|^2 - 4\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 = 4\|\boldsymbol{\alpha}\|^4 - 4\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 \\ & \leq 4\|\boldsymbol{\alpha}\|^4 - 4\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 + 2\left(\|\boldsymbol{\alpha}\|^2 - 2\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle\right)^2 \\ & = 5\|\boldsymbol{\alpha}\|^4 + 4\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 - 8\|\boldsymbol{\alpha}\|^2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle \\ & \leq 4\left[\|\boldsymbol{\alpha}\|^4 + \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle^2 - 2\|\boldsymbol{\alpha}\|^2 \langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle\right] \\ & = \left(2\|\boldsymbol{\alpha}\|^2 - 2\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle\right)^2 \\ & = \left(\|\boldsymbol{\alpha}\|^2 - 2\langle \boldsymbol{\alpha}, \boldsymbol{\beta} \rangle + \|\boldsymbol{\beta}\|^2\right)^2 = \|\boldsymbol{\alpha} - \boldsymbol{\beta}\|^4. \end{aligned}$$

## REFERENCES

- [1] Alekh Agarwal, Sahand Negahban, and Martin J. Wainwright. Fast global convergence of gradient methods for high-dimensional statistical recovery. *Ann. Statist.*, 40(5):2452–2482, 10 2012.
- [2] Animashree Anandkumar, Rong Ge, Daniel Hsu, Sham M. Kakade, and Matus Telgarsky. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- [3] Martin Azizyan, Aarti Singh, and Larry Wasserman. Minimax theory for high-dimensional gaussian mixtures with sparse mean separation. *arXiv preprint arXiv:1306.2035*, 2013.
- [4] Sivaraman Balakrishnan, Martin J. Wainwright, and Bin Yu. Statistical guarantees for the EM algorithm: From population to sample-based analysis. *arXiv preprint arXiv:1408.2156*, 2014.
- [5] Lucien Birgé. Approximation dans les espaces métriques et théorie de l'estimation. *Z. Wahrscheinlichkeitstheorie verw. Gebiete*, 65(2):181–237, 1983.
- [6] T. Tony Cai and Anru Zhang. Rop: Matrix recovery via rank-one projections. *Ann. Statist.*, 43(1):102–138, 02 2015.
- [7] Emmanuel J. Candès and Xiaodong Li. Solving quadratic equations via phaselift when there are about as many equations as unknowns. *Foundations of Computational Mathematics*, 14(5):1017–1026, 2014.
- [8] Emmanuel J. Candès and Yaniv Plan. Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Transactions on Information Theory*, 57(4):2342–2359, 2011.
- [9] Emmanuel J. Candès and Benjamin Recht. Exact matrix completion via convex optimization. *Foundations of Computational Mathematics*, 9(6):717–772, 2009.
- [10] Emmanuel J. Candès, Thomas Strohmer, and Vladislav Voroninski. PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming. *arXiv preprint arXiv:1109.4499*, 2011.
- [11] Emmanuel J. Candès, Thomas Strohmer, and Vladislav Voroninski. PhaseLift: Exact and stable signal recovery from magnitude measurements via convex programming.

*Communications on Pure and Applied Mathematics*, 66(8):1241–1274, 2013.

[12] Arun Chaganty and Percy Liang. Spectral experts for estimating mixtures of linear regressions. In *International Conference on Machine Learning (ICML)*, 2013.

[13] Moses Charikar, Sudipto Guha, Éva Tardos, and David B. Shmoys. A constant-factor approximation algorithm for the k-median problem (extended abstract). In *Proceedings of the Thirty-first Annual ACM Symposium on Theory of Computing*, STOC ’99, pages 1–10, New York, NY, USA, 1999. ACM.

[14] Jiahua Chen. Optimal rate of convergence for finite mixture models. *The Annals of Statistics*, pages 221–233, 1995.

[15] Yuxin Chen and Emmanuel J. Candès. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *arXiv preprint arXiv:1505.05114*, 2015.

[16] Yuxin Chen, Yuejie Chi, and Andrea Goldsmith. Exact and stable covariance estimation from quadratic sampling via convex programming. *Information Theory, IEEE Transactions on*, 61(7):4034–4059, 2015.

[17] Thomas M. Cover and Joy A. Thomas. *Elements of information theory*. Wiley, 2012.

[18] Arthur P. Dempster, Nan M. Laird, and Donald B. Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–38, 1977.

[19] Ehsan Elhamifar and René Vidal. Sparse subspace clustering: Algorithm, theory, and applications. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 35(11):2765–2781, 2013.

[20] Bettina Grün and Friedrich Leisch. Applications of finite mixtures of regression models. URL: <http://cran.r-project.org/web/packages/flexmix/vignettes/regression-examples.pdf>, 2007.

[21] Daniel Hsu and Sham M. Kakade. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 11–20. ACM, 2013.

[22] Geoffrey McLachlan and David Peel. *Finite Mixture Models*. Wiley series in probability and statistics: Applied probability and statistics. Wiley, 2004.

[23] Praneeth Netrapalli, Prateek Jain, and Sujay Sanghavi. Phase retrieval using alternating minimization. *Signal Processing, IEEE Transactions on*, 63(18):4814–4826, 2015.

[24] Benjamin Recht, Maryam Fazel, and Pablo A. Parrilo. Guaranteed Minimum-Rank Solutions of Linear Matrix Equations via Nuclear Norm Minimization. *SIAM Review*, 52(471), 2010.

[25] Mark Rudelson and Roman Vershynin. Hanson-wright inequality and sub-gaussian concentration. *Electron. Commun. Probab.*, 18:no. 82, 1–9, 2013.

[26] Mahdi Soltanolkotabi, Ehsan Elhamifar, and Emmanuel J. Candès. Robust subspace clustering. *Ann. Statist.*, 42(2):669–699, 04 2014.

[27] Nicolas Stadler, Peter Bühlmann, and Sara Geer. L1-penalization for mixture regression models. *TEST*, 19(2):209–256, 2010.

[28] Joel A. Tropp. User-friendly tail bounds for sums of random matrices. *Foundations of Computational Mathematics*, 12(4):389–434, 2012.

[29] Alexandre B. Tsybakov. *Introduction to Nonparametric Estimation*. Springer Series in Statistics. Springer, 2009.

[30] Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. *Arxiv preprint arXiv:1011.3027*, 2010.

[31] Kert Viele and Barbara Tong. Modeling with mixtures of linear regressions. *Statistics and Computing*, 12(4), 2002.

[32] Yu-Xiang Wang and Huan Xu. Noisy sparse subspace clustering. In *Proceedings of The 30th International Conference on Machine Learning*, pages 89–97, 2013.

[33] Zhaoran Wang, Quanquan Gu, Yang Ning, and Han Liu. High dimensional expectation-maximization algorithm: Statistical optimization and asymptotic normality. *arXiv preprint arXiv:1412.8729*, 2014.

[34] C. F. Jeff Wu. On the convergence properties of the EM algorithm. *The Annals of Statistics*, 11(1):95–103, 1983.

[35] Yuhong Yang and Andrew Barron. Information-theoretic determination of minimax rates of convergence. *The Annals of Statistics*, 27(5):1564–1599, 1999.

[36] Xinyang Yi and Constantine Caramanis. Regularized em algorithms: A unified framework and statistical guarantees. In *Advances in Neural Information Processing Systems*, pages 1567–1575, 2015.

[37] Xinyang Yi, Constantine Caramanis, and Sujay Sanghavi. Alternating minimization for mixed linear regression. In *Proceedings of The 31st International Conference on Machine Learning*, pages 613–621, 2014, Arxiv preprint arXiv:1310.3745.

[38] Bin Yu. Assouad, Fano, and Le Cam. In *Festschrift for Lucien Le Cam*, pages 423–435. Springer, 1997.

**Yudong Chen** is an Assistant Professor at the School of Operations Research and Information Engineering at Cornell University. From 2013 to 2015 he was a postdoctoral scholar at the Department of Electrical Engineering and Computer Sciences at University of California, Berkeley. He received his Ph.D. in electrical and computer engineering from the University of Texas at Austin in 2013. He obtained his BS and MS degrees from Tsinghua University, Beijing, China. His research interests include machine learning, high-dimensional and robust statistics, convex and non-convex optimization, and applications in networks and financial systems.

**Xinyang Yi** graduated with a Ph.D. in Electrical and Computer Engineering from The University of Texas at Austin. He obtained his B.E. degree from Tsinghua University in 2012. His research interests include machine learning and non-convex optimization.

**Constantine Caramanis** (M06) received his Ph.D. in electrical engineering and computer science from the Massachusetts Institute of Technology, and his A.B. in Mathematics from Harvard. Since 2006, he has been on the faculty in the Department of Electrical and Computer Engineering at The University of Texas at Austin. He received the NSF CAREER award in 2011. His current research interests include robust and large scale optimization and control, machine learning and high-dimensional statistics, with applications to large scale networks.