

---

# Sparse Mixed Linear Regression with Guarantees: Taming an Intractable Problem with Invox Relaxation

---

Adarsh Barik<sup>\*1</sup> Jean Honorio<sup>\*1</sup>

## Abstract

In this paper, we study the problem of sparse mixed linear regression on an unlabeled dataset that is generated from linear measurements from two different regression parameter vectors. Since the data is unlabeled, our task is not only to figure out a good approximation of the regression parameter vectors but also to label the dataset correctly. In its original form, this problem is NP-hard. The most popular algorithms to solve this problem (such as Expectation-Maximization) have a tendency to stuck at local minima. We provide a novel invox relaxation for this intractable problem which leads to a solution with provable theoretical guarantees. This relaxation enables exact recovery of data labels. Furthermore, we recover a close approximation of the regression parameter vectors which match the true parameter vectors in support and sign. Our formulation uses a carefully constructed primal dual witnesses framework for the invox problem. Furthermore, we show that the sample complexity of our method is only logarithmic in terms of the dimension of the regression parameter vectors.

## 1. Introduction

In this paper, we study sparse mixed linear regression where the measurements come from one of the two regression models depending upon the unknown label  $z_i^* \in \{0, 1\}$ . The observation model can be described as follows:

$$y_i = z_i^* \langle X_i, \beta_1^* \rangle + (1 - z_i^*) \langle X_i, \beta_2^* \rangle + e_i, \forall i \in \{1, \dots, n\}, \quad (1)$$

where  $X_i \in \mathbb{R}^d$ ,  $y_i \in \mathbb{R}$  and  $e_i \in \mathbb{R}$  is independent additive noise. The regression parameter vectors  $\beta_1^* \in \mathbb{R}^d$ ,  $\beta_2^* \in \mathbb{R}^d$

are sparse vectors with possibly non-overlapping supports.

Mixed linear regression models have been extensively used in a wide range of applications (Grün et al., 2007) which include but are not limited to behavioral health-care (Deb & Holmes, 2000), market segmentation (Wedel & Kamakura, 2000), music perception studies (Viele & Tong, 2002) and vehicle merging (Li et al., 2019). The main task of the problem is to estimate the regression parameter vectors and the unknown labels accurately from linear measurements. However, the problem is NP-hard without any assumptions (Yi et al., 2014). Being such a difficult problem, it also lends itself to be used as a benchmark for many non-convex optimization algorithms (Chaganty & Liang, 2013; Klusowski et al., 2019).

**Related Work.** There have been many approaches to solve the mixed linear regression problem after it was introduced by (Wedel & DeSarbo, 1995). The most popular and natural approach has been to use Expectation-minimization (EM) based alternate minimization algorithms (see Ghosh & Kannan (2020) and references therein). More broadly, the problem can be modeled under the hierarchical mixtures of experts model (Jordan & Jacobs, 1994) and solved using EM based algorithms. All these methods run the risk of getting stuck at local minima (Wu, 1983) without good initialization. (Yi et al., 2014) provides a good initialization for the noiseless case under strict technical conditions, however their method does not provide any guarantees for the noisy case. Based on the recent work of (Anandkumar et al., 2014; Hsu & Kakade, 2013), (Chaganty & Liang, 2013) have proposed an approach which uses a third order moment method based on tensor decomposition. Their approach suffers from high sample complexity (up to  $\mathcal{O}(d^6)$ ) due to tensor decomposition. (Städler et al., 2010) proposed an  $\ell_1$ -regularized approach for the sparse case and showed the existence of a local minimizer with correct support but there are no guarantees that EM achieves this local minima. (Chen et al., 2014) provided a convex relaxation involving nuclear norms for the problem. They do not focus on providing guarantees for exact label recovery and their results only hold for bounded noise and require balanced samples (almost equal number of samples for both labels). Besides, the optimization problems involving nuclear norms are computationally heavy

---

<sup>\*</sup>Equal contribution <sup>1</sup>Department of Computer Science, Purdue University, West Lafayette, Indiana, USA. Correspondence to: Adarsh Barik <abarik@purdue.edu>, Jean Honorio <jhonorio@purdue.edu>.

and slow. The mixed linear regression problem can also be modeled as a subspace clustering problem. But typically these problems require  $\mathcal{O}(d^2)$  measurements to have a unique solution (Vidal et al., 2005; Elhamifar & Vidal, 2013).

**Contribution.** Broadly, we can categorize our contribution in the following points:

- **A Combinatorial Problem:** We view the problem as a combinatorial version of a mixture of sparse linear regressions. The exact label recovery is as important for us as the recovery of regression vectors. This added exact label recovery guarantee comes at no extra cost in terms of the performance.
- **Invox Relaxation:** We solve a non-convex problem which is known to be intractable. We propose a novel relaxation of the combinatorial problem and formally show that this relaxation is invex.
- **Theoretical Guarantees:** Our method solves two sparse linear regressions and a label recovery problem simultaneously with theoretical guarantees. To that end, we recover the true labels and sparse regression parameter vectors which are correct up to the sign of entries with respect to the true parameter vectors. As a side product, we propose a novel primal-dual witness construction for our invex problem and provide theoretical guarantees for recovery. The sample complexity of our method only varies logarithmically with respect to dimension of the regression parameter vector.
- **A Novel Framework:** It should be noted that we are providing a novel framework (not an algorithm) to solve the problem. This opens the door for many algorithms to be used for this problem.

## 2. Problem Setup

In this section, we collect the notations used throughout the paper and define our problem formally. We consider a problem where measurements come from a mixture of two linear regression problems. Let  $y_i \in \mathbb{R}$  be the response variable and  $X_i \in \mathbb{R}^d$  be the observed attributes. Let  $z_i^* \in \{0, 1\}$  denote the unknown label associated with measurement  $i$ . The response  $y_i$  is generated using the observation model (1) where  $e_i \in \mathbb{R}$  is an independent noise term. We collect a total of  $n$  linear measurements with  $n_1$  measurements belonging to label 1 and  $n_2$  measurements belonging to label 0. Clearly,  $n = n_1 + n_2$ . We take  $\|\beta_1^*\|_1 \leq b_1$  and  $\|\beta_2^*\|_1 \leq b_2$ .

Let  $[d]$  denote the set  $\{1, 2, \dots, d\}$ . We assume  $X_i \in \mathbb{R}^d$  to be a zero mean sub-Gaussian random vector (Hsu et al., 2012) with covariance  $\Sigma \in \mathbb{S}_+^d$ , i.e., there exists a  $\rho > 0$ , such that for all  $\tau \in \mathbb{R}^d$  the following holds:

$\mathbb{E}(\exp(\tau^\top X_i)) \leq \exp(\frac{\|\tau\|_2^2 \rho^2}{2})$ . By simply taking  $\tau_j = r$  and  $\tau_k = 0, \forall k \neq j$ , it follows that each entry of  $X_i$  is sub-Gaussian with parameter  $\rho$ . In particular, we will assume that  $\forall j \in [d]$ ,  $\frac{X_{ij}}{\sqrt{\Sigma_{jj}}}$  is a sub-Gaussian random variable with parameter  $\sigma > 0$ . It follows trivially that  $\max_{j \in [d]} \sqrt{\Sigma_{jj}} \sigma \leq \rho$ . We will further assume that  $e_i$  is zero mean independent sub-Gaussian noise with variance  $\sigma_e$ . Our setting works with a variety of random variables as the class of sub-Gaussian random variable includes for instance Gaussian variables, any bounded random variable (e.g., Bernoulli, multinomial, uniform), any random variable with strictly log-concave density, and any finite mixture of sub-Gaussian variables.

The parameter vectors  $\beta_1^* \in \mathbb{R}^d$  and  $\beta_2^* \in \mathbb{R}^d$  are  $s_1$ -sparse and  $s_2$ -sparse respectively, i.e., at most  $s_1$  entries of  $\beta_1^*$  are non-zero whereas at most  $s_2$  entries of  $\beta_2^*$  are non-zero. We receive  $n$  i.i.d. samples of  $X_i \in \mathbb{R}^d$  and  $y_i \in \mathbb{R}$  and collect them in  $X \in \mathbb{R}^{n \times d}$  and  $y \in \mathbb{R}^n$  respectively. Similarly,  $z^* \in \{0, 1\}^n$  collects all the labels. Our goal is to recover  $\beta_1^*, \beta_2^*$  and  $z^*$  using the samples  $(X, y)$ .

We denote a matrix  $A \in \mathbb{R}^{p \times q}$  restricted to the columns and rows in  $P \subseteq [p]$  and  $Q \subseteq [q]$  respectively as  $A_{PQ}$ . Similarly, a vector  $v \in \mathbb{R}^p$  restricted to entries in  $P$  is denoted as  $v_P$ . We use  $\text{eig}_i(A)$  to denote the  $i$ -th eigenvalue (1st being the smallest) of matrix  $A$ . Similarly,  $\text{eig}_{\max}(A)$  denotes the maximum eigenvalue of matrix  $A$ . We use  $\text{diag}(A)$  to denote a vector containing the diagonal element of matrix  $A$ . By overriding the same notation, we use  $\text{diag}(v)$  to denote a diagonal matrix with its diagonal being the entries in vector  $v$ . We denote the inner product between two matrices (or vectors)  $A$  and  $B$  by  $\langle A, B \rangle$ , i.e.,  $\langle A, B \rangle = \text{trace}(A^\top B)$ , where  $\text{trace}$  denotes the trace of a matrix. The notation  $A \geq B$  denotes that  $A - B$  is a positive semidefinite matrix. Similarly,  $A > B$  denotes that  $A - B$  is a positive definite matrix. For vectors,  $\|v\|_p$  denotes the  $\ell_p$ -vector norm of vector  $v \in \mathbb{R}^d$ , i.e.,  $\|v\|_p = (\sum_{i=1}^d |v_i|^p)^{\frac{1}{p}}$ . If  $p = \infty$ , then we define  $\|v\|_\infty = \max_{i=1}^d |v_i|$ . As is the tradition, we used  $\|v\|_0$  to denote number of non-zero entries on vector  $v$ . It should be remembered that  $\ell_0$  is not a proper vector norm. For matrices,  $\|A\|_p$  denotes the induced  $\ell_p$ -matrix norm for matrix  $A \in \mathbb{R}^{p \times q}$ . In particular,  $\|A\|_2$  denotes the spectral norm of  $A$  and  $\|A\|_\infty \triangleq \max_{i \in [p]} \sum_{j=1}^q |A_{ij}|$ . For a matrix  $A \in \mathbb{R}^{p \times q}$ ,  $A(\cdot) \in \mathbb{R}^{pq}$  denotes a vector which collects all entries of the matrix  $A$ . We define an operator  $\text{sign}(A)$  for a matrix (or vector)  $A$ , which returns a matrix (or a vector) with entries being the sign of the entries of  $A$ . A function  $f(x)$  is of order  $\Omega(g(x))$  and denoted by  $f(x) = \Omega(g(x))$ , if there exists a constant  $C > 0$  such that for big enough  $x_0$ ,  $f(x) \geq Cg(x), \forall x \geq x_0$ . Similarly, a function  $f(x)$  is of order  $\mathcal{O}(g(x))$  and denoted by  $f(x) = \mathcal{O}(g(x))$ , if there exists a constant  $C > 0$  such that for big enough  $x_0$ ,  $f(x) \leq Cg(x), \forall x \geq x_0$ . For brevity in our notations, we

treat any quantity independent of  $d, s$  and  $n$  as constant. Detailed proofs for lemmas and theorems are available in the supplementary material.

### 3. A Novel Invox Relaxation

In this section, we introduce a combinatorial formulation for mixed linear regression (MLR) and propose a novel invox relaxation for this problem. Since the measurements come from a true observation model (1), we can write the following optimization problem to estimate  $\beta_1^*, \beta_2^*$  and  $z^*$ .

**Definition 3.1** (Standard MLR).

$$\begin{aligned} \min_{\beta_1 \in \mathbb{R}^d, \beta_2 \in \mathbb{R}^d, z \in \{0,1\}^n} \quad & l(z, \beta_1, \beta_2) \\ \text{such that} \quad & \|\beta_1\|_0 = s_1, \|\beta_2\|_0 = s_2 \end{aligned} \quad (2)$$

where  $l(z, \beta_1, \beta_2) = \frac{1}{n} \sum_{i=1}^n z_i (y_i - X_i^\top \beta_1)^2 + (1 - z_i) (y_i - X_i^\top \beta_2)^2$ .

Even without constraints, optimization problem (2) is a non-convex NP-hard problem (Yi et al., 2014) in its current form. In fact, a continuous relaxation of  $z \in [0, 1]^n$  does not help and it still remains a non-convex problem (See Appendix A). Furthermore, the sparsity constraints make it even difficult to solve. To deal with this intractability, we come up with a novel invox relaxation of the problem.

For ease of notation, we define the following quantities:

$$S_i = \begin{bmatrix} X_i \\ -y_i \end{bmatrix} \begin{bmatrix} X_i^\top & -y_i \end{bmatrix} = \begin{bmatrix} X_i X_i^\top & -X_i y_i \\ -y_i X_i^\top & y_i^2 \end{bmatrix}, \quad (3)$$

We provide the following invox relaxation to the optimization problem (2).

**Definition 3.2** (Invox MLR).

$$\begin{aligned} \min_{t, W, U} \quad & f(t, W, U) + \lambda_1 g(t, W, U) + \lambda_2 h(t, W, U) \\ \text{such that} \quad & W \geq \mathbf{0}, U \geq \mathbf{0} \\ & W_{d+1, d+1} = 1, U_{d+1, d+1} = 1 \\ & \|t\|_\infty \leq 1 \end{aligned} \quad (4)$$

where  $f(t, W, U) = \sum_{i=1}^n \frac{1}{2} \langle S_i, W + U \rangle + \sum_{i=1}^n \frac{1}{2} t_i \langle S_i, W - U \rangle$ ,  $g(t, W, U) = \|W(\cdot)\|_1$  and  $h(t, W, U) = \|U(\cdot)\|_1$  and  $\lambda_1$  and  $\lambda_2$  are positive regularizers.

To get an intuition behind this formulation, one can think of  $W$  and  $U$  as two rank-1 matrices which are defined as follows:

$$W = \begin{bmatrix} \beta_1 \\ 1 \end{bmatrix} \begin{bmatrix} \beta_1^\top & 1 \end{bmatrix}, U = \begin{bmatrix} \beta_2 \\ 1 \end{bmatrix} \begin{bmatrix} \beta_2^\top & 1 \end{bmatrix} \quad (5)$$

The variable  $t$  is simply a replacement of variable  $z$ , i.e.,  $z_i = \frac{t_i + 1}{2}$ . An analogous transformation exists between  $z_i^*$  and  $t_i^*$ . Then after substituting  $t, W$  and  $U$  in  $f(t, W, U)$ , we get back  $l(z, \beta_1, \beta_2)$ . The  $\ell_1$ -regularization of  $W(\cdot)$  and  $U(\cdot)$  helps us ensure sparsity. Note that for fixed  $t$ , optimization problem (4) is continuous and convex with respect to  $W$  and  $U$ . Specifically, it merges two independent regularized semidefinite programs. Unfortunately, problem (4) is not jointly convex on  $t, W$  and  $U$ , and thus, it might still remain difficult to solve. Next, we will provide arguments that despite being non-convex, optimization problem (4) belongs to a particular class of non-convex functions namely ‘‘invox’’ functions. The ‘‘invoxity’’ of functions can be defined as a generalization of convexity (Hanson, 1981). Invoxity has been recently used by (Barik & Honorio, 2021) to solve fair sparse regression problem with clustering. While, we borrow some definitions from their work to suit our context, we should emphasize that our problem is fundamentally different than their problem. They use two groups in sparse regression which have different means and they try to achieve fairness. While here, we have two groups with the same mean and there is no unfairness in the problem. We also model our parameter vectors with positive semidefinite matrices which is fundamentally different from their approach.

**Definition 3.3** (Invox function (Barik & Honorio, 2021)).

Let  $\phi(t)$  be a function defined on a set  $C$ . Let  $\eta$  be a vector valued function defined in  $C \times C$  such that  $\eta(t_1, t_2)^\top \nabla \phi(t_2)$ , is well defined  $\forall t_1, t_2 \in C$ . Then,  $\phi(t)$  is a  $\eta$ -invox function if  $\phi(t_1) - \phi(t_2) \geq \eta(t_1, t_2)^\top \nabla \phi(t_2)$ ,  $\forall t_1, t_2 \in C$ .

Note that convex functions are  $\eta$ -invox for  $\eta(t_1, t_2) = t_1 - t_2$ . (Hanson, 1981) showed that if the objective function and constraints are both  $\eta$ -invox with respect to same  $\eta$  defined in  $C \times C$ , then Karush-Kuhn-Tucker (KKT) conditions are sufficient for optimality, while it is well-known that KKT conditions are necessary. (Ben-Israel & Mond, 1986) showed a function is invox if and only if each of its stationarity point is a global minimum.

In the next lemma, we show that the relaxed optimization problem (4) is indeed  $\eta$ -invox for a particular  $\eta$  defined in  $C \times C$  and a well defined set  $C$ . Let  $C = \{(t, W, U) \mid t \in [-1, 1]^n, W \geq \mathbf{0}, U \geq \mathbf{0}, W_{d+1, d+1} = 1, U_{d+1, d+1} = 1\}$ .

**Lemma 3.4.** For  $(t, W, U) \in C$ , the functions  $f(t, W, U) = \sum_{i=1}^n \frac{1}{2} \langle S_i, W + U \rangle + \sum_{i=1}^n \frac{1}{2} t_i \langle S_i, W - U \rangle$ ,  $g(t, W, U) = \|W(\cdot)\|_1$  and  $h(t, W, U) = \|U(\cdot)\|_1$  are

$\eta$ -invox for  $\eta(t, \tilde{t}, W, \tilde{W}, U, \tilde{U}) \triangleq \begin{bmatrix} \eta_t \\ \eta_W \\ \eta_U \end{bmatrix}$ , where  $\eta_t = \mathbf{0} \in$

$\mathbb{R}^n$ ,  $\eta_W = -\tilde{W}$  and  $\eta_U = -\tilde{U}$ . We abuse the vector/matrix notation (by ignoring the dimensions) for clarity of presentation, and avoid the vectorization of matrices.

Now that we have established that optimization problem (4) is invex, we are ready to discuss our main results in the next section.

## 4. Main Results

In this section, we present our main results along with the technical assumptions. Our main goal is to show that the solution to optimization problem (4) recovers the labels  $t^*$  exactly and also recovers a good approximation of  $\beta_1^*$  and  $\beta_2^*$ . In that, we will show that the recovered  $\beta_1$  and  $\beta_2$  have the same support and sign as  $\beta_1^*$  and  $\beta_2^*$  respectively and are close to the true vectors in  $\ell_2$ -norm. But before that, we will describe a set of technical assumptions which will help us in our analysis.

### 4.1. Assumptions

Our first assumption ensures that each sample can be assigned only one label. Formally,

**Assumption 4.1** (Identifiability). For  $i \in [n]$ ,  $-\frac{1}{2}(y_i - X_i^\top \beta_1^*)^2 + \frac{1}{2}(y_i - X_i^\top \beta_2^*)^2 \geq \epsilon$  if  $z_i^* = 1$  and  $\frac{1}{2}(y_i - X_i^\top \beta_1^*)^2 - \frac{1}{2}(y_i - X_i^\top \beta_2^*)^2 \geq \epsilon$  if  $z_i^* = 0$  for some  $\epsilon > 0$ .

Clearly, if Assumption 4.1 does not hold for sample  $i$ , then we can reverse the label of sample  $i$  without increasing objective function of optimization problem (2). Another equivalent way of expressing Assumption 4.1 is as following: for  $i \in [n]$ ,  $\langle S_i, W^* \rangle < \langle S_i, U^* \rangle$  if  $z_i^* = 1$  and  $\langle S_i, W^* \rangle > \langle S_i, U^* \rangle$  if  $z_i^* = 0$  where,

$$W^* = \begin{bmatrix} \beta_1^* \\ 1 \end{bmatrix} [\beta_1^{*\top} \quad 1], \quad U^* = \begin{bmatrix} \beta_2^* \\ 1 \end{bmatrix} [\beta_2^{*\top} \quad 1]. \quad (6)$$

Let  $P$  denote the support of  $\beta_1^*$ , i.e.,  $P = \{i \mid \beta_{1i}^* \neq 0, i \in [d]\}$  and let  $Q$  denote the support of  $\beta_2^*$ , i.e.,  $Q = \{i \mid \beta_{2i}^* \neq 0, i \in [d]\}$ . Similarly, we define their complement as  $P^c = \{i \mid \beta_{1i}^* = 0, i \in [d]\}$  and  $Q^c = \{i \mid \beta_{2i}^* = 0, i \in [d]\}$ . We take  $|P| = s_1, |P^c| = d - s_1, |Q| = s_2$  and  $|Q^c| = d - s_2$ . For ease of notation, we define  $H \triangleq \mathbb{E}(X_i X_i^\top) \forall i \in [n]$ . Let  $\mathcal{I}_1 \triangleq \{i \mid z_i^* = 1, i \in [n]\}$  and  $\mathcal{I}_2 \triangleq \{i \mid z_i^* = 0, i \in [n]\}$ . We define  $\hat{H}_1 \triangleq \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_i X_i^\top$  and  $\hat{H}_2 \triangleq \frac{1}{n_2} \sum_{i \in \mathcal{I}_2} X_i X_i^\top$ . As our next assumption, we need the minimum eigenvalue of the population covariance matrix of  $X$  restricted to rows and columns in  $P$  (similarly in  $Q$ ) to be greater than zero.

**Assumption 4.2** (Positive Definiteness of Hessian).  $H_{PP} > \mathbf{0}$  and  $H_{QQ} > \mathbf{0}$  or equivalently  $\min(\text{eig}_{\min}(H_{PP}), \text{eig}_{\min}(H_{QQ})) = C_{\min} > 0$ . We also assume that  $\text{eig}_{\max}(H) = C_{\max} > 0$ . Note that  $\max(\text{eig}_{\max}(H_{PP}), \text{eig}_{\max}(H_{QQ})) \leq C_{\max}$ .

In practice, we only deal with finite samples and not populations. In the next lemma, we will show that with a sufficient

number of samples, a condition similar to Assumption 4.2 holds with high probability in the finite-sample setting.

**Lemma 4.3.** *If Assumption 4.2 holds and  $n_1 = \Omega(\frac{s_1 + \log d}{C_{\min}^2})$  and  $n_2 = \Omega(\frac{s_2 + \log d}{C_{\min}^2})$ , then  $\min(\text{eig}_{\min}(\hat{H}_{1PP}), \text{eig}_{\min}(\hat{H}_{2QQ})) \geq \frac{C_{\min}}{2}$  and  $\max(\text{eig}_{\max}(\hat{H}_{1PP}), \text{eig}_{\max}(\hat{H}_{2QQ})) \leq \frac{3C_{\max}}{2}$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ .*

As the third assumption, we will need to ensure that the variates outside the support of  $\beta_1^*$  and  $\beta_2^*$  do not exert lot of influence on the variates in the support of  $\beta_1^*$  and  $\beta_2^*$  respectively. For this, we use a technical condition commonly known as the mutual incoherence condition. It has been previously used in many problems related to regularized regression such as compressed sensing (Wainwright, 2009b), Markov random fields (Ravikumar et al., 2010), non-parametric regression (Ravikumar et al., 2007), diffusion networks (Daneshmand et al., 2014), among others.

**Assumption 4.4** (Mutual Incoherence).  $\max(\|H_{P^c P} H_{PP}^{-1}\|_\infty, \|H_{Q^c Q} H_{QQ}^{-1}\|_\infty) \leq 1 - \xi$  for some  $\xi \in (0, 1]$ .

Again, we will show that with a sufficient number of samples, a condition similar to Assumption 4.4 holds in the finite-sample setting with high probability.

**Lemma 4.5.** *If Assumption 4.4 holds and  $n_1 = \Omega(\frac{s_1^3(\log s_1 + \log d)}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$  and  $n_2 = \Omega(\frac{s_2^3(\log s_2 + \log d)}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$ , then  $\max(\|\hat{H}_{P^c P} \hat{H}_{PP}^{-1}\|_\infty, \|\hat{H}_{Q^c Q} \hat{H}_{QQ}^{-1}\|_\infty) \leq 1 - \frac{\xi}{2}$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  where  $\tau(C_{\min}, \xi, \sigma, \Sigma)$  is a constant independent of  $n_1, n_2, d, s_1$  and  $s_2$ .*

### 4.2. Main Theorem

Now we are ready to state our main result.

**Theorem 4.6.** *If Assumptions 4.1, 4.2 and 4.4 hold,  $\lambda_1 \geq \frac{64\rho\sigma\epsilon}{\xi} \sqrt{n_1 \log d}$ ,  $\lambda_2 \geq \frac{64\rho\sigma\epsilon}{\xi} \sqrt{n_2 \log d}$  and  $n_1 = \Omega(\frac{s_1^3 \log^2 d}{\tau_0(C_{\min}, \xi, \sigma, \Sigma, \rho)})$  and  $n_2 = \Omega(\frac{s_2^3 \log^2 d}{\tau_0(C_{\min}, \xi, \sigma, \Sigma, \rho)})$ , then with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  the solution to the optimization problem (4) satisfies the following properties:*

1. The labels are recovered exactly, i.e.,

$$t_i = t_i^*, \quad \forall i \in [n] \quad (7)$$

2. The regression parameter vectors are close to the true vectors. Formally,

$$W = \begin{bmatrix} \beta_1 \\ 1 \end{bmatrix} [\beta_1^\top \quad 1], \quad U = \begin{bmatrix} \beta_2 \\ 1 \end{bmatrix} [\beta_2^\top \quad 1] \quad (8)$$

such that  $\beta_1 = [\beta_{1P} \quad \mathbf{0}_{P^c}]^\top$  and  $\beta_2 = [\beta_{2Q} \quad \mathbf{0}_{Q^c}]^\top$



and

$$\begin{aligned} \|\beta_1 - \beta_1^*\|_2 &\leq (2 + b_1) \frac{2\lambda_1\sqrt{s_1}}{C_{\min}n_1} \\ \|\beta_2 - \beta_2^*\|_2 &\leq (2 + b_2) \frac{2\lambda_2\sqrt{s_2}}{C_{\min}n_2}. \end{aligned} \quad (9)$$

In order to prove Theorem 4.6, we will have to show that the labels are recovered exactly. We will also need to show that  $W$  and  $U$  are rank-1 matrices with eigenvectors  $[\beta_1 \ 1]^\top$  and  $[\beta_2 \ 1]^\top$  respectively. Moreover, we will also need to ensure that their supports match supports of the true vectors and they are close to true vectors in  $\ell_2$ -norm.

## 5. Theoretical Analysis

We use primal-dual witness approach to show our results. The primal-dual witness approach was developed by (Wainwright, 2009a) for linear regression problem which has been later used in many convex problems such as Markov random fields (Ravikumar et al., 2010), non-parametric regression (Ravikumar et al., 2007), diffusion networks (Daneshmand et al., 2014) etc. The main idea is to start with a potential solution with certain properties and then later show that these properties are indeed consistent with the final solution. We extend this idea to our invex problem. To that end, we start our proof with a potential solution which has certain ‘‘consistency certificate’’.

### 5.1. Consistency Certificate

We start by taking solutions  $W$  and  $U$  with the following properties which we call consistency certificates:

C1.  $W$  and  $U$  are sparse. In particular, they have the following sparsity structure:

$$\begin{aligned} W &= \begin{bmatrix} W_{PP} & \mathbf{0}_{PP^c} & W_{Pd+1} \\ \mathbf{0}_{P^cP} & \mathbf{0}_{P^cP^c} & \mathbf{0}_{P^cd+1} \\ W_{d+1P} & \mathbf{0}_{d+1P^c} & W_{d+1,d+1} \end{bmatrix} \\ U &= \begin{bmatrix} U_{QQ} & \mathbf{0}_{QQ^c} & U_{Qd+1} \\ \mathbf{0}_{Q^cQ} & \mathbf{0}_{Q^cQ^c} & \mathbf{0}_{Q^cd+1} \\ U_{d+1Q} & \mathbf{0}_{d+1Q^c} & U_{d+1,d+1} \end{bmatrix} \end{aligned} \quad (10)$$

We collect all the non-zero entries of  $W$  and  $U$  in  $\bar{W} \in \mathbb{R}^{s_1+1, s_1+1}$  and  $\bar{U} \in \mathbb{R}^{s_2+1, s_2+1}$ .

It should be noted that the consistency certificate C1 is not another assumption. In that, eventually we will have to show that it holds in the final solution. We can prove that C1 is consistent with the final solution by showing strict dual feasibility for both  $W$  and  $U$  which we do in subsection 5.7.

### 5.2. A Modified Compact Invex Problem

Once we substitute  $W$  and  $U$  from C1 in optimization problem (4), we get a low dimensional optimization problem.

**Definition 5.1** (Compact Invex MLR).

$$\begin{aligned} &\min_{t, \bar{W}, \bar{U}} \quad \bar{f}(t, \bar{W}, \bar{U}) + \lambda_1 \bar{g}(t, \bar{W}, \bar{U}) + \lambda_2 \bar{h}(t, \bar{W}, \bar{U}) \\ &\text{such that} \quad \bar{W} \geq \mathbf{0}, \bar{U} \geq \mathbf{0} \\ &\quad \bar{W}_{s_1+1, s_1+1} = 1, \bar{U}_{s_2+1, s_2+1} = 1 \\ &\quad \|t\|_\infty \leq 1 \end{aligned} \quad (11)$$

where  $\bar{f}(t, \bar{W}, \bar{U}) = \sum_{i=1}^n \frac{1}{2} (\langle \bar{S}_i^P, \bar{W} \rangle + \langle \bar{S}_i^Q, \bar{U} \rangle) + \sum_{i=1}^n \frac{1}{2} t_i (\langle \bar{S}_i^P, \bar{W} \rangle - \langle \bar{S}_i^Q, \bar{U} \rangle)$ ,  $\bar{g}(t, \bar{W}, \bar{U}) = \|\bar{W}(\cdot)\|_1$ ,  $\bar{h}(t, \bar{W}, \bar{U}) = \|\bar{U}(\cdot)\|_1$  and  $\lambda_1$  and  $\lambda_2$  are positive regularizers.

Note that

$$\bar{S}_i^P = \begin{bmatrix} S_{iP,P} & S_{iP,d+1} \\ S_{id+1,P} & S_{id+1,d+1} \end{bmatrix}, \bar{S}_i^Q = \begin{bmatrix} S_{iQ,Q} & S_{iQ,d+1} \\ S_{id+1,Q} & S_{id+1,d+1} \end{bmatrix}. \quad (12)$$

For clarity, we will drop the superscripts from  $\bar{S}_i$  when the context is clear. Next, we list down the necessary and sufficient conditions to solve optimization problem (11).

### 5.3. Necessary and Sufficient KKT Conditions

First, we write the Lagrangian  $L(\Theta)$  for fixed  $\lambda_1 > 0$  and  $\lambda_2 > 0$ , where  $\Theta = (t, \bar{W}, \bar{U}; \Pi, \Lambda, \alpha, \gamma, \nu, \mu)$  is a collection of parameters.

$$\begin{aligned} L(\Theta) &= \bar{f}(t, \bar{W}, \bar{U}) + \lambda_1 \bar{g}(t, \bar{W}, \bar{U}) + \lambda_2 \bar{h}(t, \bar{W}, \bar{U}) \\ &\quad - \langle \Pi, \bar{W} \rangle - \langle \Lambda, \bar{U} \rangle + \alpha (\bar{W}_{s_1+1, s_1+1} - 1) + \\ &\quad \gamma (\bar{U}_{s_2+1, s_2+1} - 1) - \sum_{i=1}^n \nu_i (t_i + 1) + \sum_{i=1}^n \mu_i (t_i - 1) \end{aligned} \quad (13)$$

Here  $\Pi \geq \mathbf{0}, \Lambda \geq \mathbf{0}, \alpha \in \mathbb{R}, \gamma \in \mathbb{R}, \nu_i > 0$  and  $\mu_i > 0$  are the dual variables (of appropriate dimensions) for optimization problem (11). Using this Lagrangian, the KKT conditions at the optimum can be written as:

1. Stationarity conditions:

$$\sum_{i=1}^n \frac{t_i + 1}{2} \bar{S}_i^P + \lambda_1 Z - \Pi + I_\alpha = \mathbf{0} \quad (14)$$

where  $Z$  is an element of the subgradient set of  $\|\bar{W}(\cdot)\|_1$ , i.e.,  $Z \in \frac{\partial \|\bar{W}(\cdot)\|_1}{\partial \bar{W}}$  and  $\|Z(\cdot)\|_\infty \leq 1$  and  $I_\alpha \in \mathbb{R}^{s_1+1, s_1+1}$  has all zero entries except  $(s_1 + 1, s_1 + 1)$  entry which is  $\alpha$ .

$$\sum_{i=1}^n \frac{1 - t_i}{2} \bar{S}_i^Q + \lambda_2 V - \Lambda + I_\gamma = \mathbf{0} \quad (15)$$

where  $V$  is an element of the subgradient set of  $\|\bar{U}(\cdot)\|_1$ , i.e.,  $V \in \frac{\partial \|\bar{U}(\cdot)\|_1}{\partial \bar{U}}$  and  $\|V(\cdot)\|_\infty \leq 1$  and  $I_\gamma \in \mathbb{R}^{s_2+1, s_2+1}$  has all zero entries except  $(s_2+1, s_2+1)$  entry which is  $\gamma$ . For all  $i \in [n]$ , it holds that

$$\frac{1}{2} \langle \bar{S}_i^P, \bar{W} \rangle - \frac{1}{2} \langle \bar{S}_i^Q, \bar{U} \rangle - \nu_i + \mu_i = 0 \quad (16)$$

2. Complementary Slackness conditions:

$$\langle \Pi, \bar{W} \rangle = 0, \quad \langle \Lambda, \bar{U} \rangle = 0 \quad (17)$$

$$\nu_i(t_i + 1) = 0, \quad \mu_i(t_i - 1) = 0 \quad \forall i \in [n] \quad (18)$$

3. Dual Feasibility conditions:

$$\Pi \geq \mathbf{0}, \quad \Lambda \geq \mathbf{0} \quad (19)$$

$$\nu_i \geq 0, \quad \mu_i \geq 0 \quad \forall i \in [n] \quad (20)$$

4. Primal Feasibility conditions:

$$\begin{aligned} \bar{W} \geq \mathbf{0}, \quad \bar{U} \geq \mathbf{0} \\ \bar{W}_{s_1+1, s_1+1} = 1, \quad \bar{U}_{s_2+1, s_2+1} = 1, \quad \|t\|_\infty \leq 1 \end{aligned} \quad (21)$$

Next, we will provide a setting for primal and dual variables which satisfies all the KKT conditions.

#### 5.4. Construction of Primal and Dual Variables

In this subsection, we will provide a construction of primal and dual variables which satisfies the KKT conditions for optimization problem (11). To that end, we provide our first main result.

**Theorem 5.2** (Primal Dual Variables Construction). *If Assumptions 4.1, 4.2 and 4.4 hold,  $\lambda_1 \geq \frac{64\rho\sigma_e}{\xi} \sqrt{n_1 \log d}$ ,  $\lambda_2 \geq \frac{64\rho\sigma_e}{\xi} \sqrt{n_2 \log d}$  and  $n_1 = \Omega(\frac{s_1^3 \log^2 d}{\tau_0(C_{\min}, \xi, \sigma, \Sigma, \rho)})$  and  $n_2 = \Omega(\frac{s_2^3 \log^2 d}{\tau_0(C_{\min}, \xi, \sigma, \Sigma, \rho)})$ , then the following setting of primal and dual variables:*

• *Primal Variables:*

$$t_i = t_i^*, \quad \forall i \in [n]$$

$$\bar{W} = \begin{bmatrix} \tilde{\beta}_1 \\ 1 \end{bmatrix} \begin{bmatrix} \tilde{\beta}_1^\top & 1 \end{bmatrix}, \quad \bar{U} = \begin{bmatrix} \tilde{\beta}_2 \\ 1 \end{bmatrix} \begin{bmatrix} \tilde{\beta}_2^\top & 1 \end{bmatrix}$$

where

$$\tilde{\beta}_1 = \arg \min_{\beta \in \mathbb{R}^{s_1}} \sum_{i=1}^n \frac{t_i^* + 1}{2} (y_i - X_{iP}^\top \beta)^2 +$$

$$\lambda_1 (\|\beta\|_1 + 1)^2$$

and

$$\tilde{\beta}_2 = \arg \min_{\beta \in \mathbb{R}^{s_2}} \sum_{i=1}^n \frac{1 - t_i^*}{2} (y_i - X_{iQ}^\top \beta)^2 +$$

$$\lambda_2 (\|\beta\|_1 + 1)^2$$

• *Dual Variables:*

$$\nu_i = 0, \quad \mu_i = -\frac{1}{2} \langle \bar{S}_i^P, \bar{W} \rangle + \frac{1}{2} \langle \bar{S}_i^Q, \bar{U} \rangle \quad \forall i \in \mathcal{I}_1$$

$$\mu_i = 0, \quad \nu_i = \frac{1}{2} \langle \bar{S}_i^P, \bar{W} \rangle - \frac{1}{2} \langle \bar{S}_i^Q, \bar{U} \rangle \quad \forall i \in \mathcal{I}_2$$

$$\Pi = \sum_{i=1}^n \frac{t_i^* + 1}{2} \bar{S}_i^P + \lambda_1 Z + I_\alpha$$

$$\Lambda = \sum_{i=1}^n \frac{1 - t_i^*}{2} \bar{S}_i^Q + \lambda_2 V + I_\gamma$$

$$\alpha = -\left\langle \sum_{i=1}^n \frac{t_i + 1}{2} \bar{S}_i^P + \lambda_1 Z, \bar{W} \right\rangle$$

$$\gamma = -\left\langle \sum_{i=1}^n \frac{1 - t_i}{2} \bar{S}_i^Q + \lambda_2 V, \bar{U} \right\rangle$$

satisfies all the KKT conditions for optimization problem (11) with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ , where  $\tau_0(C_{\min}, \alpha, \sigma, \Sigma, \rho, \gamma)$  is a constant independent of  $s_1, s_2, d, n_1$  and  $n_2$  and thus, the primal variables are a globally optimal solution for (11). Furthermore, the above solution is also unique.

**Proof Sketch.** The main idea behind our proofs is to verify that the setting of primal and dual variables in Theorem 5.2 satisfies all the KKT conditions described in subsection 5.3. We do this by proving multiple lemmas in subsequent subsections. The outline of the proof is as follows:

- It can be trivially verified that the primal feasibility condition (21) holds. The stationarity conditions (14) and (15) holds by construction of  $\Pi$  and  $\Lambda$  respectively. Similarly, the stationarity condition (16) holds by choice of  $\nu_i$  and  $\mu_i$ . Choice of  $t, \nu_i, \mu_i, \alpha$  and  $\gamma$  ensure that complementary slackness conditions (17) and (18) also hold.
- In subsection 5.5, we use Lemmas 5.3, 5.5 and 5.6 to verify that the dual feasibility conditions (20) and (19) hold. We will also show in subsection subsection 5.5 that our solution is also unique.

#### 5.5. Verifying Dual Feasibility

To verify dual feasibility, first we will show that  $\mu_i \geq 0, \forall i \in \mathcal{I}_1, \nu_i \geq 0, \forall i \in \mathcal{I}_2$ . We define  $\Delta_1 \triangleq \tilde{\beta}_1 - \beta_{1P}^*$  and  $\Delta_2 \triangleq \tilde{\beta}_2 - \beta_{2Q}^*$ . Then, the following lemma holds true.

**Lemma 5.3.** *If Assumptions 4.1, 4.2 and 4.4 hold, and  $\lambda_1 \geq 8\rho\sigma_e \sqrt{n_1 \log d}$ ,  $\lambda_2 \geq 8\rho\sigma_e \sqrt{n_2 \log d}$ ,  $n_1 = \Omega(\frac{s_1^3 \log d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$ , and  $n_2 = \Omega(\frac{s_2^3 \log d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$  then  $\|\Delta_1\|_2 \leq (2 + b_1) \frac{2\lambda_1 \sqrt{s_1}}{C_{\min} n_1}$  and  $\|\Delta_2\|_2 \leq (2 + b_2) \frac{2\lambda_2 \sqrt{s_2}}{C_{\min} n_2}$*

with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  where  $\tau(C_{\min}, \xi, \sigma, \Sigma)$  is a constant independent of  $s_1, s_2, d, n_1$  or  $n_2$ .

Using the result of Lemma 5.3, we are going to prove that the settings for dual variables  $\mu_i$  and  $\nu_i$  works with high probability.

**Lemma 5.4.** *If Assumptions 4.1, 4.2 and 4.4 hold, and  $\lambda_1 \geq 8\rho\sigma_e\sqrt{n_1\log d}$ ,  $\lambda_2 \geq 8\rho\sigma_e\sqrt{n_2\log d}$ ,  $n_1 = \Omega(\frac{s_1^3\log d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$ , and  $n_2 = \Omega(\frac{s_2^3\log d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$  then  $\mu_i \geq 0, \forall i \in \mathcal{I}_1$  and  $\nu_i \geq 0, \forall i \in \mathcal{I}_2$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  where  $\tau(C_{\min}, \xi, \sigma, \Sigma)$  is a constant independent of  $s_1, s_2, d, n_1$  or  $n_2$ .*

Now we will show that  $\Pi \geq \mathbf{0}$  and  $\Lambda \geq \mathbf{0}$ . We will do this in two steps. The first step is to show that both  $\Pi$  and  $\Lambda$  have a zero eigenvalue. In particular,

**Lemma 5.5.** *Both  $\Pi$  and  $\Lambda$  have zero eigenvalues corresponding to eigenvectors  $\begin{bmatrix} \tilde{\beta}_1 \\ 1 \end{bmatrix}$  and  $\begin{bmatrix} \tilde{\beta}_2 \\ 1 \end{bmatrix}$  respectively.*

Next, we show that all the other eigenvalues of both  $\Pi$  and  $\Lambda$  are strictly positive.

**Lemma 5.6.** *If Assumption 4.2 holds and  $n_1 = \Omega(\frac{s_1 + \log d}{C_{\min}^2})$  and  $n_2 = \Omega(\frac{s_2 + \log d}{C_{\min}^2})$ , then the second eigenvalues of  $\Pi$  and  $\Lambda$  are strictly positive with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ , i.e.,  $\text{eig}_2(\Pi) > 0$  and  $\text{eig}_2(\Lambda) > 0$ .*

On the one hand, Lemma 5.6 ensures that  $\Pi \geq 0$  and  $\Lambda \geq 0$ , but on the other it also forces  $\bar{W}$  and  $\bar{U}$  to be rank-1 and unique as both  $\Pi$  and  $\bar{W}$  have to be positive semidefinite and  $\Pi$  has exactly one vector in its nullspace (same with  $\Lambda$  and  $\bar{U}$ ).

## 5.6. Going back to Invox MLR

Now that we have the setting of  $t_i, \bar{W}$  and  $\bar{U}$  for Compact Invox MLR problem (11), we can extend these to the original Invox MLR problem (4). Notice that all the other entries of  $W$  and  $U$  are zeros, thus it readily follows that

$$W = \begin{bmatrix} \beta_1 \\ 1 \end{bmatrix} \begin{bmatrix} \beta_1^\top & 1 \end{bmatrix}, U = \begin{bmatrix} \beta_2 \\ 1 \end{bmatrix} \begin{bmatrix} \beta_2^\top & 1 \end{bmatrix} \quad (22)$$

where  $\beta_1 = \begin{bmatrix} \tilde{\beta}_1 \\ \mathbf{0} \end{bmatrix}$  and  $\beta_2 = \begin{bmatrix} \tilde{\beta}_2 \\ \mathbf{0} \end{bmatrix}$ . Furthermore, result from Lemma 5.3 extends directly and gives us

$$\begin{aligned} \|\beta_1 - \beta_1^*\|_2 &\leq (2 + b_1) \frac{2\lambda_1\sqrt{s_1}}{C_{\min}n_1} \\ \|\beta_2 - \beta_2^*\|_2 &\leq (2 + b_2) \frac{2\lambda_2\sqrt{s_2}}{C_{\min}n_2}. \end{aligned} \quad (23)$$

The last remaining thing is to show that consistency certificate C1 indeed holds which we will do in next subsection.

## 5.7. Validating Consistency Certificate

Observe that once we substitute  $t_i = t_i^*$  in optimization problem (4), it decouples into two independent convex optimization problems involving  $W$  and  $U$  respectively. Furthermore, since we established that  $W$  and  $U$  are rank-1, we can rewrite these independent problems in terms of  $\beta_1$  and  $\beta_2$ . Our task is to show that  $\beta_{1P^c} = \mathbf{0}$  and  $\beta_{2Q^c} = \mathbf{0}$ . It suffices to show it for  $\beta_1$  as arguments for  $\beta_2$  are the same. Below, we consider the simplified optimization problem in terms of  $\beta_1$ :

$$\beta_1 = \arg \min_{\beta \in \mathbb{R}^d} \sum_{i \in \mathcal{I}_1} (X_i^\top \beta - y_i)^2 + \lambda_1 (\|\beta\|_1 + 1)^2 \quad (24)$$

Since we are only dealing with measurements in  $\mathcal{I}_1$ , we can substitute  $y_i = X_i^\top \beta^* + e_i$ . Furthermore,  $\beta_1$  must satisfy stationarity KKT condition which can be written as:

$$\begin{aligned} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_i X_i^\top (\beta_1 - \beta^*) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_i e_i + \\ \frac{1}{n_1} \lambda_1 (\|\beta_1\|_1 + 1) \zeta = \mathbf{0}, \end{aligned} \quad (25)$$

where  $\zeta$  is in the subdifferential set of  $\|\beta_1\|_1$  and  $\|\zeta\|_\infty \leq 1$ . Specifically,  $\zeta_i = \text{sign}(\beta_1)$ ,  $\forall i \in P$  and  $\zeta_i \in [-1, 1]$ ,  $\forall i \in P^c$ . Our task is to show that  $\zeta$  fulfills strict dual feasibility, i.e.,  $\|\zeta_{P^c}\|_\infty < 1$ . We decompose equation (25) into two parts – one corresponding to entries in  $P$  and the other corresponding to entries in  $P^c$ . For entries in  $P$ , we have

$$\begin{aligned} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} X_{i_P}^\top (\beta_{1_P} - \beta_{1_P}^*) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \\ + \frac{1}{n_1} \lambda_1 (\|\beta_1\|_1 + 1) \zeta_P = \mathbf{0} \end{aligned} \quad (26)$$

Similarly, for entries in  $P^c$ , we have

$$\begin{aligned} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} X_{i_{P^c}}^\top (\beta_{1_{P^c}} - \beta_{1_{P^c}}^*) - \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} e_i \\ + \frac{1}{n_1} \lambda_1 (\|\beta_1\|_1 + 1) \zeta_{P^c} = \mathbf{0} \end{aligned}$$

After rearranging the terms and substituting for  $(\beta_{1_P} - \beta_{1_P}^*)$  from equation (26), we get

$$\begin{aligned} \frac{\lambda_1}{n_1} (1 + \|\beta_1\|_1) \zeta_{P^c} = -\hat{H}_{P^c P} \hat{H}_{P P}^{-1} \left( \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i - \right. \\ \left. \frac{1}{n_1} \lambda_1 (\|\beta_1\|_1 + 1) \zeta_P \right) + \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} e_i \end{aligned}$$

Let  $\bar{\lambda}_1 = \frac{\lambda_1}{n_1}$  and note that  $\|\beta_1\|_1 \geq 0$ , using norm inequality

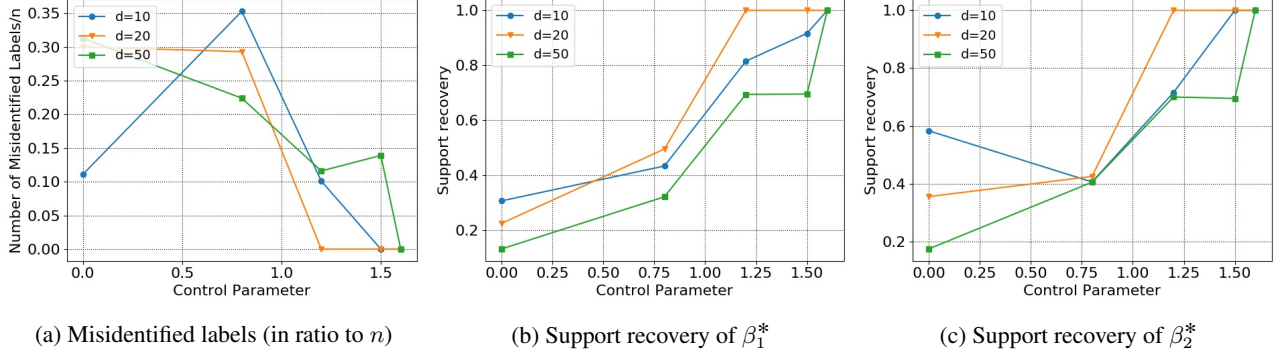


Figure 1. Label and support recovery with control parameter  $C_p$ .

ties we can rewrite the above equation as:

$$\begin{aligned} \|\zeta_{P^c}\|_\infty &\leq \|\widehat{H}_{P^c P} \widehat{H}_{PP}^{-1}\|_\infty \left( \frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right)_\infty + \|\zeta_P\|_\infty \\ &+ \left\| \frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} e_i \right\|_\infty \end{aligned}$$

We know that  $\|\widehat{H}_{P^c P} \widehat{H}_{PP}^{-1}\|_\infty \leq (1 - \frac{\xi}{2})$  for some  $\xi \in (0, 1]$ . The following lemma provides bounds on  $\left\| \frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_\infty$  and  $\left\| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} e_i \right\|_\infty$ .

**Lemma 5.7.** *Let  $\lambda_1 \geq \frac{64\rho\sigma_e}{\xi} \sqrt{n_1 \log d}$ . Then the following holds true:*

$$\begin{aligned} \mathbb{P}\left(\left\| \frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_\infty \geq \frac{\xi}{8 - 4\xi}\right) &\leq \mathcal{O}\left(\frac{1}{d}\right), \\ \mathbb{P}\left(\left\| \frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} e_i \right\|_\infty \geq \frac{\xi}{8}\right) &\leq \mathcal{O}\left(\frac{1}{d}\right) \end{aligned}$$

It follows that  $\|\zeta_{P^c}\|_\infty \leq 1 - \frac{\xi}{4}$  with probability at least  $1 - \mathcal{O}\left(\frac{1}{d}\right)$ . Thus, the consistency certificate C1 indeed holds with high probability.

## 6. Experimental Validation

Note that we are not proposing any new algorithm in our paper. However, to validate our theoretical results we performed experiments on synthetic data. We generated response  $y$  using Gaussian random variables  $X$  and chose regression parameter  $\beta_1^*$  (or  $\beta_2^*$ ) based on the label of the samples. We fixed the sparsity  $s_1 = s_2 = 4$ , however supports were not necessarily the same for both the regression parameter vectors. We varied  $n_1$  and  $n_2$  according to our theorems, i.e., both were varied with  $10^{C_p} \log^2 d$  for  $d = 10, 20$  and  $50$  where  $C_p$  is a control parameters. The regularizers were kept according to our theorem and were varied as  $\mathcal{O}(\sqrt{n_1 \log d})$  and  $\mathcal{O}(\sqrt{n_2 \log d})$ . We measured performance of our algorithm based on the label recovery (in

ratio to supplied  $n$ ) and support recovery for both parameter vectors. The experiments were run three times independently. Note how we make zero mistakes in label recovery as we increase number of samples. Similarly, support recovery (ratio of intersection and union with correct support) for both parameter vectors goes to 1 as we increase sample size. It should be noted that we do not propose any new algorithm but our method is free of any initialization requirement. As for specific algorithm for empirical verification, we use a projected subgradient method (Duchi & Singer, 2009) to check convergence for our problem which is achieved without any requirement on initialization. In fact, any algorithm which converges to a stationary point should work for our framework.

## 7. Concluding Remarks

We provide a novel formulation of invex MLR. We show that invexity of our optimization problem allows for a tractable solution. We provide provable theoretical guarantees for our solution. The sample complexity of our method is polynomial in terms of sparsity and logarithmic in terms of the dimension of the true parameter. Our method helps to identify labels exactly and recovers regression parameter vectors with correct support and correct sign. It would be interesting to think about extending our ideas to mixture of more than two groups of regressions in future.

## Acknowledgements

This material is based upon work supported by the National Science Foundation under Grant No. 2134209-DMS.

## References

Anandkumar, A., Ge, R., Hsu, D., Kakade, S. M., and Telgarsky, M. Tensor decompositions for learning latent variable models. *Journal of machine learning research*, 15:2773–2832, 2014.



- Barik, A. and Honorio, J. Fair sparse regression with clustering: An invex relaxation for a combinatorial problem. *Advances in neural information processing systems*, 2021.
- Ben-Israel, A. and Mond, B. What is invexity? *The ANZIAM Journal*, 28(1):1–9, 1986.
- Chaganty, A. T. and Liang, P. Spectral experts for estimating mixtures of linear regressions. In *International Conference on Machine Learning*, pp. 1040–1048. PMLR, 2013.
- Chen, Y., Yi, X., and Caramanis, C. A convex formulation for mixed regression with two components: Minimax optimal rates. In *Conference on Learning Theory*, pp. 560–604. PMLR, 2014.
- Daneshmand, H., Gomez-Rodriguez, M., Song, L., and Schoelkopf, B. Estimating Diffusion Network Structures: Recovery Conditions, Sample Complexity & Soft-Thresholding Algorithm. In *International Conference on Machine Learning*, pp. 793–801, 2014.
- Deb, P. and Holmes, A. M. Estimates of use and costs of behavioural health care: a comparison of standard and finite mixture models. *Health economics*, 9(6):475–489, 2000.
- Duchi, J. and Singer, Y. Efficient online and batch learning using forward backward splitting. *The Journal of Machine Learning Research*, 10:2899–2934, 2009.
- Elhamifar, E. and Vidal, R. Sparse subspace clustering: Algorithm, theory, and applications. *IEEE transactions on pattern analysis and machine intelligence*, 35(11): 2765–2781, 2013.
- Ghosh, A. and Kannan, R. Alternating minimization converges super-linearly for mixed linear regression. In *International Conference on Artificial Intelligence and Statistics*, pp. 1093–1103. PMLR, 2020.
- Grün, B., Leisch, F., et al. Applications of finite mixtures of regression models. URL: <http://cran.r-project.org/web/packages/flexmix/vignettes/regression-examples.pdf>, 2007.
- Hanson, M. A. On sufficiency of the kuhn-tucker conditions. *Journal of Mathematical Analysis and Applications*, 80 (2):545–550, 1981.
- Haynsworth, E. V. Determination of the inertia of a partitioned hermitian matrix. *Linear algebra and its applications*, 1(1):73–81, 1968.
- Horn, R. A. and Johnson, C. R. *Matrix Analysis*. Cambridge university press, 2012.
- Hsu, D. and Kakade, S. M. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pp. 11–20, 2013.
- Hsu, D., Kakade, S., Zhang, T., et al. A tail inequality for quadratic forms of subgaussian random vectors. *Electronic Communications in Probability*, 17, 2012.
- Jordan, M. I. and Jacobs, R. A. Hierarchical mixtures of experts and the em algorithm. *Neural computation*, 6(2): 181–214, 1994.
- Klusowski, J. M., Yang, D., and Brinda, W. Estimating the coefficients of a mixture of two linear regressions by expectation maximization. *IEEE Transactions on Information Theory*, 65(6):3515–3524, 2019.
- Li, G., Pan, Y., Yang, Z., and Ma, J. Modeling vehicle merging position selection behaviors based on a finite mixture of linear regression models. *IEEE Access*, 7: 158445–158458, 2019.
- Ravikumar, P., Liu, H., Lafferty, J., and Wasserman, L. Spam: Sparse Additive Models. In *Proceedings of the 20th International Conference on Neural Information Processing Systems*, pp. 1201–1208. Curran Associates Inc., 2007.
- Ravikumar, P., Wainwright, M. J., Lafferty, J. D., et al. High-dimensional ising model selection using  $l_1$ -regularized logistic regression. *The Annals of Statistics*, 38(3):1287–1319, 2010.
- Ravikumar, P., Wainwright, M. J., Raskutti, G., Yu, B., et al. High-dimensional Covariance Estimation by Minimizing  $L_1$ -Penalized Log-Determinant Divergence. *Electronic Journal of Statistics*, 5:935–980, 2011.
- Städler, N., Bühlmann, P., and Van De Geer, S.  $l_1$ -penalization for mixture regression models. *Test*, 19(2): 209–256, 2010.
- Vershynin, R. How close is the sample covariance matrix to the actual covariance matrix? *Journal of Theoretical Probability*, 25(3):655–686, 2012.
- Vidal, R., Ma, Y., and Sastry, S. Generalized principal component analysis (gpca). *IEEE transactions on pattern analysis and machine intelligence*, 27(12):1945–1959, 2005.
- Viele, K. and Tong, B. Modeling with mixtures of linear regressions. *Statistics and Computing*, 12(4):315–330, 2002.
- Wainwright, M. J. Information-theoretic bounds on sparsity recovery in the high-dimensional and noisy setting. *IEEE Trans. Info. Theory*, 55:5728–5741, December 2009a.

- Wainwright, M. J. Sharp Thresholds for High-Dimensional and Noisy Sparsity Recovery Using L1-Constrained Quadratic Programming (Lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009b.
- Wainwright, M. J. *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge University Press, 2019.
- Wedel, M. and DeSarbo, W. S. A mixture likelihood approach for generalized linear models. *Journal of classification*, 12(1):21–55, 1995.
- Wedel, M. and Kamakura, W. A. *Market segmentation: Conceptual and methodological foundations*. Springer Science & Business Media, 2000.
- Wu, C. J. On the convergence properties of the em algorithm. *The Annals of statistics*, pp. 95–103, 1983.
- Yi, X., Caramanis, C., and Sanghavi, S. Alternating minimization for mixed linear regression. In *International Conference on Machine Learning*, pp. 613–621. PMLR, 2014.

## Supplementary Material: Sparse Mixed Linear Regression with Guarantees: Taming an Intractable Problem with Invox Relaxation

### A. Continuous Relaxation of Standard MLR is non-convex

It suffices to prove that the objective function  $l(z, \beta_1, \beta_2)$  of optimization problem (2) is non-convex when  $z_i$  is allowed to be between 0 and 1. We note that

$$l(z, \beta_1, \beta_2) = \sum_{i=1}^n z_i (y_i - X_i^\top \beta_1)^2 + (1 - z_i) (y_i - X_i^\top \beta_2)^2 \quad (27)$$

where  $\beta_1 \in \mathbb{R}^d, \beta_2 \in \mathbb{R}^d$  and  $z_i \in [0, 1], \forall i \in [n]$ . Let  $\Theta = (z, \beta_1, \beta_2, \bar{z}, \bar{\beta}_1, \bar{\beta}_2)$ . We consider the following quantity:

$$F(\Theta) = f(z, \beta_1, \beta_2) - f(\bar{z}, \bar{\beta}_1, \bar{\beta}_2) - \sum_{i=1}^n \frac{\partial f}{\partial \bar{z}_i} (z_i - \bar{z}_i) - \frac{\partial f}{\partial \beta_1}^\top (\beta_1 - \bar{\beta}_1) - \frac{\partial f}{\partial \beta_2}^\top (\beta_2 - \bar{\beta}_2) \quad (28)$$

where

$$\begin{aligned} \frac{\partial f}{\partial \bar{z}_i} &= (y_i - X_i^\top \bar{\beta}_1)^2 - (y_i - X_i^\top \bar{\beta}_2)^2, \quad \forall i \in [n] \\ \frac{\partial f}{\partial \beta_1} &= \sum_{i=1}^n -2\bar{z}_i X_i (y_i - X_i^\top \bar{\beta}_1) \\ \frac{\partial f}{\partial \beta_2} &= \sum_{i=1}^n -2(1 - \bar{z}_i) X_i (y_i - X_i^\top \bar{\beta}_2) \end{aligned} \quad (29)$$

It suffices to show that  $F(\Theta)$  changes sign for different feasible values of  $\Theta$ . We choose the following variables:

$$\begin{aligned} z_i &= 0, \bar{z}_i = \frac{1}{2} \quad \forall i \in [n] \\ \beta_{1_k} &= u_1, \beta_{1_j} = 0, \forall j \neq k \\ \beta_{2_l} &= u_2, \beta_{2_j} = 0, \forall j \neq l \\ \bar{\beta}_{1_k} &= w_1, \bar{\beta}_{1_j} = 0, \forall j \neq k \\ \bar{\beta}_{2_l} &= w_2, \bar{\beta}_{2_j} = 0, \forall j \neq l \\ u_1 &= w_1 - (u_2 - w_2) \end{aligned} \quad (30)$$

Note that choice of  $w_1, u_2$  and  $w_2$  can be arbitrary. This simplifies  $F(\Theta)$ :

$$F(\Theta) = \sum_{i=1}^n (y_i - u_2 X_{il})^2 - \sum_{i=1}^n (y_i - w_2 X_{il})^2 \quad (31)$$

Consider the case when  $X_{il} > 0, \forall i \in [n]$ . Then choosing  $u_2 < w_2$  makes  $F(\Theta) > 0$  while choosing  $u_2 > w_2$  makes  $F(\Theta) < 0$ . This proves our claim.

### B. Proof of Lemma 3.4

**Lemma 3.4** For  $(t, W, U) \in C$ , the functions  $f(t, W, U) = \sum_{i=1}^n \frac{1}{2} \langle S_i, W + U \rangle + \sum_{i=1}^n \frac{1}{2} t_i \langle S_i, W - U \rangle, g(t, W, U) = \|W(\cdot)\|_1$  and  $h(t, W, U) = \|U(\cdot)\|_1$  are  $\eta$ -invex for  $\eta(t, \tilde{t}, W, \tilde{W}, U, \tilde{U}) \triangleq \begin{bmatrix} \eta_t \\ \eta_W \\ \eta_U \end{bmatrix}$ , where  $\eta_t = \mathbf{0} \in \mathbb{R}^n, \eta_W = -\tilde{W}$  and

$\eta_U = -\tilde{U}$ . We abuse the vector/matrix notation (by ignoring the dimensions) for clarity of presentation, and avoid the vectorization of matrices.

*Proof.* We know  $f(t, W, U) = \sum_{i=1}^n \frac{t_i+1}{2} \langle S_i, W \rangle + \frac{1-t_i}{2} \langle S_i, U \rangle$ . Then,

$$\begin{aligned} \frac{\partial f}{\partial t_i} &= \frac{1}{2} \langle S_i, W - U \rangle \\ \frac{\partial f}{\partial W} &= \sum_{i=1}^n \frac{t_i+1}{2} S_i \\ \frac{\partial f}{\partial U} &= \sum_{i=1}^n \frac{1-t_i}{2} S_i \end{aligned} \quad (32)$$

To prove that  $f(t, W, U)$  is invex, we need to show that

$$f(t, W, U) - f(\tilde{t}, \tilde{W}, \tilde{U}) - \sum_{i=1}^n \eta_{t_i} \frac{\partial f}{\partial t_i} - \langle \eta_W, \frac{\partial f}{\partial W} \rangle - \langle \eta_U, \frac{\partial f}{\partial U} \rangle \geq 0 \quad (33)$$

We take  $\eta_t = \mathbf{0} \in \mathbb{R}^n$ ,  $\eta_W = -\tilde{W}$  and  $\eta_U = -\tilde{U}$  and expand LHS of equation (33) as follows:

$$\begin{aligned} & \sum_{i=1}^n \frac{t_i+1}{2} \langle S_i, W \rangle + \frac{1-t_i}{2} \langle S_i, U \rangle - \sum_{i=1}^n \frac{\tilde{t}_i+1}{2} \langle S_i, \tilde{W} \rangle - \frac{1-\tilde{t}_i}{2} \langle S_i, \tilde{U} \rangle + \langle \tilde{W}, \sum_{i=1}^n \frac{\tilde{t}_i+1}{2} S_i \rangle + \langle \tilde{U}, \sum_{i=1}^n \frac{1-\tilde{t}_i}{2} S_i \rangle \\ &= \sum_{i=1}^n \frac{t_i+1}{2} \langle S_i, W \rangle + \frac{1-t_i}{2} \langle S_i, U \rangle \\ &\geq 0 \end{aligned} \quad (34)$$

The last inequality holds because  $S_i$ ,  $W$  and  $U$  are all positive semidefinite and  $t_i \in [-1, 1]$ .

Similarly,

$$\begin{aligned} & g(t, W, U) - g(\tilde{t}, \tilde{W}, \tilde{U}) - \sum_{i=1}^n \eta_{t_i} \frac{\partial g}{\partial t_i} - \langle \eta_W, \frac{\partial g}{\partial W} \rangle - \langle \eta_U, \frac{\partial g}{\partial U} \rangle \\ &= \|W(\cdot)\|_1 - \|\tilde{W}(\cdot)\|_1 + \|\tilde{W}(\cdot)\|_1 \geq 0 \end{aligned} \quad (35)$$

and

$$\begin{aligned} & h(t, W, U) - h(\tilde{t}, \tilde{W}, \tilde{U}) - \sum_{i=1}^n \eta_{t_i} \frac{\partial h}{\partial t_i} - \langle \eta_W, \frac{\partial h}{\partial W} \rangle - \langle \eta_U, \frac{\partial h}{\partial U} \rangle \\ &= \|U(\cdot)\|_1 - \|\tilde{U}(\cdot)\|_1 + \|\tilde{U}(\cdot)\|_1 \geq 0 \end{aligned} \quad (36)$$

□

### C. Proof of Lemma 4.3

**Lemma 4.3** *If Assumption 4.2 holds and  $n_1 = \Omega(\frac{s_1 + \log d}{C_{\min}^2})$  and  $n_2 = \Omega(\frac{s_2 + \log d}{C_{\min}^2})$ , then  $\min(\text{eig}_{\min}(\hat{H}_{1PP}), \text{eig}_{\min}(\hat{H}_{2QQ})) \geq \frac{C_{\min}}{2}$  and  $\max(\text{eig}_{\max}(\hat{H}_{1PP}), \text{eig}_{\max}(\hat{H}_{2QQ})) \leq \frac{3C_{\max}}{2}$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ .*

*Proof.* We prove the Lemma for a general support  $S$  and samples  $n$ . The results follow when we substitute  $S$  by  $P$  and  $Q$  and  $n$  by  $n_1$  or  $n_2$  based on the context. By the Courant-Fischer variational representation (Horn & Johnson, 2012):

$$\begin{aligned} \text{eig}_{\min}(\mathbb{E}(X_i X_i^\top)_{SS}) &= \min_{\|y\|_2=1} y^\top \mathbb{E}(X_i X_i^\top)_{SS} y = \min_{\|y\|_2=1} y^\top (\mathbb{E}(X_i X_i^\top)_{SS} - \frac{1}{n} X_S^\top X_S + \frac{1}{n} X_S^\top X_S) y \\ &\leq y^\top (\mathbb{E}(X_i X_i^\top)_{SS} - \frac{1}{n} X_S^\top X_S + \frac{1}{n} X_S^\top X_S) y \\ &= y^\top (\mathbb{E}(X_i X_i^\top)_{SS} - \frac{1}{n} X_S^\top X_S) y + y^\top \frac{1}{n} X_S^\top X_S y \end{aligned} \quad (37)$$



It follows that

$$\text{eig}_{\min}\left(\frac{1}{n}X_S^\top X_S\right) \geq C_{\min} - \|\mathbb{E}(X_i X_i^\top)_{SS} - \frac{1}{n}X_S^\top X_S\|_2 \quad (38)$$

The term  $\|\mathbb{E}(X_i X_i^\top)_{SS} - \frac{1}{n}X_S^\top X_S\|_2$  can be bounded using Proposition 2.1 in (Vershynin, 2012) for sub-Gaussian random variables. In particular,

$$\mathbb{P}\left(\|\mathbb{E}(X_i X_i^\top)_{SS} - \frac{1}{n}X_S^\top X_S\|_2 \geq \epsilon\right) \leq 2 \exp(-c\epsilon^2 n + s) \quad (39)$$

for some constant  $c > 0$ . Taking  $\epsilon = \frac{C_{\min}}{2}$ , we show that  $\text{eig}_{\min}\left(\frac{1}{n}X_S^\top X_S\right) \geq \frac{C_{\min}}{2}$  with probability at least  $1 - 2 \exp(-\frac{cC_{\min}^2 n}{4} + |S|)$ . The specific results for  $n_1$  and  $n_2$  follow directly.

Remark: Similarly, it can be shown that  $\text{eig}_{\max}\left(\frac{1}{n}X_S^\top X_S\right) \leq \frac{3C_{\max}}{2}$  with probability at least  $1 - 2 \exp(-\frac{cC_{\max}^2 n}{4} + |S|)$ .  $\square$

## D. Proof of Lemma 4.5

**Lemma 4.5** *If Assumption 4.4 holds and  $n_1 = \Omega\left(\frac{s_1^3(\log s_1 + \log d)}{\tau(C_{\min}, \xi, \sigma, \Sigma)}\right)$  and  $n_2 = \Omega\left(\frac{s_2^3(\log s_2 + \log d)}{\tau(C_{\min}, \xi, \sigma, \Sigma)}\right)$ , then  $\max(\|\hat{H}_{P^c P} \hat{H}_{P^c P}^{-1}\|_\infty, \|\hat{H}_{Q^c Q} \hat{H}_{Q^c Q}^{-1}\|_\infty) \leq 1 - \frac{\xi}{2}$  with probability at least  $1 - \mathcal{O}\left(\frac{1}{d}\right)$  where  $\tau(C_{\min}, \xi, \sigma, \Sigma)$  is a constant independent of  $n_1, n_2, d, s_1$  and  $s_2$ .*

*Proof.* We prove the Lemma for a general support  $S$  (and corresponding non-support  $S^c$ ) and samples  $n$ . The results follow when we substitute  $S$  by  $P$  and  $Q$  and  $n$  by  $n_1$  or  $n_2$  based on the context. Let  $|S| = s$  and  $|S^c| = d - s$ . Before we prove the result of Lemma 4.5, we will prove a helper lemma.

**Lemma D.1.** *If Assumption 4.4 holds then for some  $\delta > 0$ , the following inequalities hold:*

$$\begin{aligned} \mathbb{P}(\|\hat{H}_{S^c S} - H_{S^c S}\|_\infty \geq \delta) &\leq 4(d-s)s \exp\left(-\frac{n\delta^2}{128s^2(1+4\sigma^2)\max_l \Sigma_{ll}^2}\right) \\ \mathbb{P}(\|\hat{H}_{SS} - H_{SS}\|_\infty \geq \delta) &\leq 4s^2 \exp\left(-\frac{n\delta^2}{128s^2(1+4\sigma^2)\max_l \Sigma_{ll}^2}\right) \\ \mathbb{P}(\|(\hat{H}_{SS})^{-1} - (H_{SS})^{-1}\|_\infty \geq \delta) &\leq 2 \exp\left(-\frac{c\delta^2 C_{\min}^4 n}{4s} + s\right) + 2 \exp\left(-\frac{cC_{\min}^2 n}{4} + s\right) \end{aligned} \quad (40)$$

*Proof.* Let  $A_{ij}$  be  $(i, j)$ -th entry of  $\hat{H}_{S^c S} - H_{S^c S}$ . Clearly,  $\mathbb{E}(A_{ij}) = 0$ . By using the definition of the  $\|\cdot\|_\infty$  norm, we can write:

$$\begin{aligned} \mathbb{P}(\|\hat{H}_{S^c S} - H_{S^c S}\|_\infty \geq \delta) &= \mathbb{P}\left(\max_{i \in S^c} \sum_{j \in S} |A_{ij}| \geq \delta\right) \\ &\leq (d-s) \mathbb{P}\left(\sum_{j \in S} |A_{ij}| \geq \delta\right) \\ &\leq (d-s)s \mathbb{P}\left(|A_{ij}| \geq \frac{\delta}{s}\right) \end{aligned} \quad (41)$$

where the second last inequality comes as a result of the union bound across entries in  $S^c$  and the last inequality is due to the union bound across entries in  $S$ . Recall that  $X_i, i \in [d]$  are zero mean random variables with covariance  $\Sigma$  and each  $\frac{X_i}{\sqrt{\Sigma_{ii}}}$  is a sub-Gaussian random variable with parameter  $\sigma$ . Using the results from Lemma 1 of (Ravikumar et al., 2011), for some  $\delta \in (0, s \max_l \Sigma_{ll} 8(1+4\sigma^2))$ , we can write:

$$\mathbb{P}\left(|A_{ij}| \geq \frac{\delta}{s}\right) \leq 4 \exp\left(-\frac{n\delta^2}{128s^2(1+4\sigma^2)\max_l \Sigma_{ll}^2}\right) \quad (42)$$

Therefore,

$$\mathbb{P}(\|\hat{H}_{S^c S} - H_{S^c S}\|_\infty \geq \delta) \leq 4(d-s)s \exp\left(-\frac{n\delta^2}{128s^2(1+4\sigma^2)\max_l \Sigma_{ll}^2}\right) \quad (43)$$

Similarly, we can show that

$$\mathbb{P}(\|\hat{H}_{SS} - H_{SS}\|_\infty \geq \delta) \leq 4s^2 \exp\left(-\frac{n\delta^2}{128s^2(1+4\sigma^2)\max_l \Sigma_{ll}^2}\right) \quad (44)$$

Next, we will show that the third inequality in (40) holds. Note that

$$\begin{aligned} \|(\hat{H}_{S^cS})^{-1} - (H_{S^cS})^{-1}\|_\infty &= \|(H_{SS})^{-1}(H_{SS} - \hat{H}_{SS})(\hat{H}_{SS})^{-1}\|_\infty \\ &\leq \sqrt{s}\|(H_{SS})^{-1}(H_{SS} - \hat{H}_{SS})(\hat{H}_{SS})^{-1}\|_2 \\ &\leq \sqrt{s}\|(H_{SS})^{-1}\|_2\|(H_{SS} - \hat{H}_{SS})\|_2\|(\hat{H}_{SS})^{-1}\|_2 \end{aligned} \quad (45)$$

Note that  $\|H_{SS}\|_2 \geq C_{\min}$ , thus  $\|(H_{SS})^{-1}\|_2 \leq \frac{1}{C_{\min}}$ . Similarly,  $\|H_{SS}\|_2 \geq \frac{C_{\min}}{2}$  with probability at least  $1 - 2\exp(-\frac{cC_{\min}^2 n}{4} + s)$ . We also have  $\|(H_{SS} - \hat{H}_{SS})\|_2 \leq \epsilon$  with probability at least  $1 - 2\exp(-c\epsilon^2 n + s)$ . Taking  $\epsilon = \delta \frac{C_{\min}^2}{2\sqrt{s}}$ , we get

$$\mathbb{P}(\|(H_{SS} - \hat{H}_{SS})\|_2 \geq \delta \frac{C_{\min}^2}{2\sqrt{s}}) \leq 2\exp\left(-\frac{c\delta^2 C_{\min}^4 n}{4s} + s\right) \quad (46)$$

It follows that  $\|(\hat{H}_{SS})^{-1} - (H_{SS})^{-1}\|_\infty \leq \delta$  with probability at least  $1 - 2\exp(-\frac{c\delta^2 C_{\min}^4 n}{4s} + s) - 2\exp(-\frac{cC_{\min}^2 n}{4} + s)$ .  $\square$

Now we are ready to show that the statement of Lemma 4.5 holds using the results from Lemma D.1. We will rewrite  $\hat{H}_{S^cS}(\hat{H}_{SS})^{-1}$  as the sum of four different terms:

$$\hat{H}_{S^cS}(\hat{H}_{SS})^{-1} = T_1 + T_2 + T_3 + T_4, \quad (47)$$

where

$$\begin{aligned} T_1 &\triangleq \hat{H}_{S^cS}((\hat{H}_{SS})^{-1} - (H_{SS})^{-1}) \\ T_2 &\triangleq (\hat{H}_{S^cS} - H_{S^cS})(H_{SS})^{-1} \\ T_3 &\triangleq (\hat{H}_{S^cS} - H_{S^cS})((\hat{H}_{SS})^{-1} - (H_{SS})^{-1}) \\ T_4 &\triangleq H_{S^cS}(H_{SS})^{-1}. \end{aligned} \quad (48)$$

Then it follows that  $\|\hat{H}_{S^cS}(\hat{H}_{SS})^{-1}\|_\infty \leq \|T_1\|_\infty + \|T_2\|_\infty + \|T_3\|_\infty + \|T_4\|_\infty$ . Now, we will bound each term separately. First, recall that Assumption 4.4 ensures that  $\|T_4\|_\infty \leq 1 - \xi$ .

**Controlling  $T_1$ .** We can rewrite  $T_1$  as,

$$T_1 = -H_{S^cS}(H_{SS})^{-1}(\hat{H}_{SS} - H_{SS})(\hat{H}_{SS})^{-1} \quad (49)$$

then,

$$\begin{aligned} \|T_1\|_\infty &= \|H_{S^cS}(H_{SS})^{-1}(\hat{H}_{SS} - H_{SS})(\hat{H}_{SS})^{-1}\|_\infty \\ &\leq \|H_{S^cS}(H_{SS})^{-1}\|_\infty\|\hat{H}_{SS} - H_{SS}\|_\infty\|(\hat{H}_{SS})^{-1}\|_\infty \\ &\leq (1 - \xi)\|\hat{H}_{SS} - H_{SS}\|_\infty\sqrt{s}\|(\hat{H}_{SS})^{-1}\|_2 \\ &\leq (1 - \xi)\|\hat{H}_{SS} - H_{SS}\|_\infty\frac{2\sqrt{s}}{C_{\min}} \\ &\leq \frac{\xi}{6} \end{aligned} \quad (50)$$

The last inequality holds with probability at least  $1 - 2\exp(-\frac{cC_{\min}^2 n}{4} + s) - 4s^2 \exp(-\frac{nC_{\min}^2 \xi^2}{18432(1-\xi)^2 s^3(1+4\sigma^2)\max_l \Sigma_{ll}^2})$  by taking  $\delta = \frac{C_{\min} \xi}{12(1-\xi)\sqrt{s}}$ .

**Controlling  $T_2$ .** Recall that  $T_2 = (\hat{H}_{S^cS} - H_{S^cS})(H_{SS})^{-1}$ . Thus,

$$\begin{aligned} \|T_2\|_\infty &\leq \sqrt{s} \|(H_{SS})^{-1}\|_2 \|(\hat{H}_{S^cS} - H_{S^cS})\|_\infty \\ &\leq \frac{\sqrt{s}}{C_{\min}} \|(\hat{H}_{S^cS} - H_{S^cS})\|_\infty \\ &\leq \frac{\xi}{6} \end{aligned} \quad (51)$$

The last inequality holds with probability at least  $1 - 4(d-s)s \exp(-\frac{nC_{\min}^2\xi^2}{4608s^3(1+4\sigma^2)\max_l \Sigma_l^2})$  by choosing  $\delta = \frac{C_{\min}\xi}{6\sqrt{s}}$ .

**Controlling  $T_3$ .** Note that,

$$\begin{aligned} \|T_3\|_\infty &\leq \|(\hat{H}_{S^cS} - H_{S^cS})\|_\infty \|((\hat{H}_{SS})^{-1} - (H_{SS})^{-1})\|_\infty \\ &\leq \frac{\xi}{6} \end{aligned} \quad (52)$$

The last inequality holds with probability at least  $1 - 4(d-s)s \exp(-\frac{n\xi}{768s^2(1+4\sigma^2)\max_l \Sigma_l^2}) - 2 \exp(-\frac{c\xi C_{\min}^4 n}{24s} + s) - 2 \exp(-\frac{cC_{\min}^2 n}{4} + s)$  by choosing  $\delta = \sqrt{\frac{\xi}{6}}$  in the first and third inequality of equation (40). By combining all the above results, we prove Lemma 4.5. The specific results for  $n_1$  and  $n_2$  follow directly.  $\square$

## E. Proof of Lemma 5.3

**Lemma 5.3.** *If Assumptions 4.1, 4.2 and 4.4 hold, and  $\lambda_1 \geq 8\rho\sigma_e\sqrt{n_1 \log d}$ ,  $\lambda_2 \geq 8\rho\sigma_e\sqrt{n_2 \log d}$ ,  $n_1 = \Omega(\frac{s_1^3 \log d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$ , and  $n_2 = \Omega(\frac{s_2^3 \log d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$  then  $\|\Delta_1\|_2 \leq (2 + b_1) \frac{2\lambda_1\sqrt{s_1}}{C_{\min}n_1}$  and  $\|\Delta_2\|_2 \leq (2 + b_2) \frac{2\lambda_2\sqrt{s_2}}{C_{\min}n_2}$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  where  $\tau(C_{\min}, \xi, \sigma, \Sigma)$  is a constant independent of  $s_1, s_2, d, n_1$  or  $n_2$ .*

*Proof.* It suffices to prove the result for  $\Delta_1$  as the result for  $\Delta_2$  follows in the same way. Note,

$$\begin{aligned} \tilde{\beta}_1 &= \arg \min_{\beta \in \mathbb{R}^{s_1}} \sum_{i=1}^{t_i^*+1} \frac{t_i^*+1}{2} (y_i - X_{i_P}^\top \beta)^2 + \lambda_1 (\|\beta\|_1 + 1)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^{s_1}} \sum_{i \in \mathcal{I}_1} (y_i - X_{i_P}^\top \beta)^2 + \lambda_1 (\|\beta\|_1 + 1)^2 \end{aligned}$$

The optimal  $\tilde{\beta}_1$  must satisfy stationarity KKT condition at the optimum, i.e.,

$$\sum_{i \in \mathcal{I}_1} X_{i_P} (-y_i + X_{i_P}^\top \tilde{\beta}_1) + z \lambda_1 (z^\top \tilde{\beta}_1 + 1) = \mathbf{0}$$

where  $\|\tilde{\beta}_1\|_1 = z^\top \tilde{\beta}_1$  and  $z$  is in the subdifferential set of  $\|\tilde{\beta}_1\|_1$  and  $\|z\|_\infty \leq 1$ . Since  $i \in \mathcal{I}_1$ , we can substitute  $y_i = X_{i_P}^\top \beta_{1_P}^* + e_i$ .

$$\left( \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} X_{i_P}^\top + \frac{1}{n_1} \lambda_1 z z^\top \right) (\beta_{1_P}^* - \tilde{\beta}_1) + \frac{1}{n_1} \left( \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right) + \frac{1}{n_1} \lambda_1 z z^\top \beta_{1_P}^* + \frac{1}{n_1} \lambda_1 z = \mathbf{0}$$

Note that  $\hat{H}_{1_{PP}} = \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} X_{i_P}^\top$ . Using norm-inequalities:

$$\|\Delta_1\|_2 \leq \|(\hat{H}_{1_{PP}} + \frac{\lambda_1}{n_1} z z^\top)^{-1}\|_2 \left( \frac{1}{n_1} \left\| \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2 + \left\| \frac{1}{n_1} \lambda_1 z z^\top \beta_{1_P}^* \right\|_2 + \left\| \frac{1}{n_1} \lambda_1 z \right\|_2 \right) \quad (53)$$

Using Lemma 4.3,  $\text{eig}_{\min}(\hat{H}_{1_{PP}}) \geq \frac{C_{\min}}{2}$ , and using Weyl's inequality  $\text{eig}_{\min}(\hat{H}_{1_{PP}} + \frac{\lambda_1}{n_1} z z^\top) \geq \frac{C_{\min}}{2}$ . It follows that  $\|(\hat{H}_{1_{PP}} + \frac{\lambda_1}{n_1} z z^\top)^{-1}\|_2 \leq \frac{2}{C_{\min}}$ .

$$\|\Delta_1\|_2 \leq \frac{2}{C_{\min}} \left( \frac{1}{n_1} \left\| \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2 + \frac{\lambda_1 \sqrt{s_1}}{n_1} \|\beta_{1_P}^*\|_1 + \frac{\lambda_1 \sqrt{s_1}}{n_1} \right) \quad (54)$$

We know that  $\|\beta_1^*\|_1 \leq b_1$ . Thus,

$$\|\Delta_1\|_2 \leq \frac{2}{C_{\min}} \left( \frac{1}{n_1} \left\| \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2 + \frac{\lambda_1 \sqrt{s_1}}{n_1} b_1 + \frac{\lambda_1 \sqrt{s_1}}{n_1} \right) \quad (55)$$

It only remains to bound  $\left\| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2$  which we do in the following lemma.

**Lemma E.1.** *If  $\lambda_1 \geq 8\rho\sigma_e\sqrt{n_1 \log d}$ , then  $\left\| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2 \leq \sqrt{s_1} \frac{\lambda_1}{n_1}$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$*

Thus, it follows that

$$\|\Delta_1\|_\infty \leq \|\Delta_1\|_2 \leq (2 + b_1) \frac{2\lambda_1 \sqrt{s_1}}{C_{\min} n_1} \quad (56)$$

□

## F. Proof of Lemma E.1

**Lemma E.1** *If  $\lambda_1 \geq 8\rho\sigma_e\sqrt{n_1 \log d}$ , then  $\left\| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2 \leq \sqrt{s_1} \frac{\lambda_1}{n_1}$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ .*

*Proof.* We will start with  $\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i$ . We take the  $i$ -th entry of  $\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i$  for some  $i \in P$ , i.e.,  $|\frac{1}{n} \sum_{j \in \mathcal{I}_1} X_{ji} e_j|$ . Recall that  $X_{ji}$  is a sub-Gaussian random variable with parameter  $\rho^2$  and  $e_j$  is a sub-Gaussian random variable with parameter  $\sigma_e^2$ . Then,  $\frac{X_{ji} e_j}{\rho \sigma_e}$  is a sub-exponential random variable with parameters  $(4\sqrt{2}, 2)$ . Using the concentration bounds for the sum of independent sub-exponential random variables (Wainwright, 2019), we can write:

$$\mathbb{P}\left( \left| \frac{1}{n_1} \sum_{j \in \mathcal{I}_1} \frac{X_{ji} e_j}{\rho \sigma_e} \right| \geq t \right) \leq 2 \exp\left(-\frac{n_1 t^2}{64}\right), \quad 0 \leq t \leq 8 \quad (57)$$

Taking a union bound across  $i \in P$ :

$$\begin{aligned} \mathbb{P}(\exists i \in P \mid \left| \frac{1}{n_1} \sum_{j \in \mathcal{I}_1} \frac{X_{ji} e_j}{\rho \sigma_e} \right| \geq t) &\leq 2s_1 \exp\left(-\frac{n_1 t^2}{64}\right) \\ 0 \leq t &\leq 8 \end{aligned} \quad (58)$$

It follows that  $\left\| \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i \right\|_2 \leq \sqrt{st}$  with probability at least  $1 - 2s \exp(-\frac{nt^2}{64\rho^2\sigma_e^2})$  for some  $0 \leq t \leq 8\rho\sigma_e$ . Taking  $t = \frac{\lambda_1}{n}$ , we get the desired result. □

## G. Proof of Lemma 5.4

**Lemma 5.4** *If Assumptions 4.1, 4.2 and 4.4 hold, and  $\lambda_1 \geq 8\rho\sigma_e\sqrt{n_1 \log d}$ ,  $\lambda_2 \geq 8\rho\sigma_e\sqrt{n_2 \log d}$ ,  $n_1 = \Omega(\frac{s_1^3 \log^2 d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$ , and  $n_2 = \Omega(\frac{s_2^3 \log^2 d}{\tau(C_{\min}, \xi, \sigma, \Sigma)})$  then  $\mu_i \geq 0, \forall i \in \mathcal{I}_1$  and  $\nu_i \geq 0, \forall i \in \mathcal{I}_2$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  where  $\tau(C_{\min}, \xi, \sigma, \Sigma)$  is a constant independent of  $s_1, s_2, d, n_1$  or  $n_2$ .*

*Proof.* We start with the setting of  $\mu_i$  when  $i$  is in  $\mathcal{I}_1$ .

$$\begin{aligned} \mu_i &= -\frac{1}{2} \langle \bar{S}_i^P, \bar{W} \rangle + \frac{1}{2} \langle \bar{S}_i^Q, \bar{U} \rangle \\ &= -\frac{1}{2} (y_i - X_{i_P}^\top \tilde{\beta}_1)^2 + \frac{1}{2} (y_i - X_{i_Q}^\top \tilde{\beta}_2)^2 \\ &= -\frac{1}{2} (y_i - X_{i_P}^\top \beta_{1_P}^* + X_{i_P}^\top (\tilde{\beta}_1 - \beta_{1_P}^*))^2 + \frac{1}{2} (y_i - X_{i_Q}^\top \beta_{2_Q}^* + X_{i_Q}^\top (\tilde{\beta}_2 - \beta_{2_Q}^*))^2 \\ &= -\frac{1}{2} ((y_i - X_{i_P}^\top \beta_{1_P}^*)^2 + (\tilde{\beta}_1 - \beta_{1_P}^*)^\top X_{i_P} X_{i_P}^\top (\tilde{\beta}_1 - \beta_{1_P}^*) + 2(y_i - X_{i_P}^\top \beta_{1_P}^*) X_{i_P}^\top (\tilde{\beta}_1 - \beta_{1_P}^*)) \\ &\quad + \frac{1}{2} ((y_i - X_{i_Q}^\top \beta_{2_Q}^*)^2 + (\tilde{\beta}_2 - \beta_{2_Q}^*)^\top X_{i_Q} X_{i_Q}^\top (\tilde{\beta}_2 - \beta_{2_Q}^*) + 2(y_i - X_{i_Q}^\top \beta_{2_Q}^*) X_{i_Q}^\top (\tilde{\beta}_2 - \beta_{2_Q}^*)) \end{aligned} \quad (59)$$



Since  $i \in \mathcal{I}_1$ , we can substitute  $y_i = X_{i_P}^\top \beta_{1_P}^* + e_i$ .

$$\begin{aligned} \mu_i &= -\frac{1}{2}(y_i - X_{i_P}^\top \beta_{1_P}^*)^2 + \frac{1}{2}(y_i - X_{i_Q}^\top \beta_{2_Q}^*)^2 - \frac{1}{2}\Delta_1^\top X_{i_P} X_{i_P}^\top \Delta_1 + \frac{1}{2}\Delta_2^\top X_{i_Q} X_{i_Q}^\top \Delta_2 - e_i X_{i_P}^\top \Delta_1 + \\ &\quad (X_{i_P}^\top \beta_{1_P}^* + e_i - X_{i_Q}^\top \beta_{2_Q}^*) X_{i_Q}^\top \Delta_2 \\ &= -\frac{1}{2}(y_i - X_{i_P}^\top \beta_{1_P}^*)^2 + \frac{1}{2}(y_i - X_{i_Q}^\top \beta_{2_Q}^*)^2 - \frac{1}{2}\Delta_1^\top X_{i_P} X_{i_P}^\top \Delta_1 + \frac{1}{2}\Delta_2^\top X_{i_Q} X_{i_Q}^\top \Delta_2 - e_i X_{i_P}^\top \Delta_1 + \\ &\quad e_i X_{i_Q}^\top \Delta_2 + (\beta_{1_P}^* - \beta_{2_Q}^*)^\top X_{i_Q} X_{i_Q}^\top (\beta_{2_Q}^* - \beta_{1_P}^*) \end{aligned} \quad (60)$$

Using bounds on the eigenvalue of data matrix, Assumption 4.1 and bounds on  $\|\Delta_1\|_2$  and  $\|\Delta_2\|_2$ , we can place a bound on  $\mu_i$ .

$$\mu_i \geq \epsilon - \frac{n_1 3C_{\max}}{2} \|\Delta_1\|_2^2 + \frac{n_2 C_{\min}}{2} \|\Delta_2\|_2^2 - |e_i X_{i_P}^\top \Delta_1| - |e_i X_{i_Q}^\top \Delta_2| - \frac{n 3C_{\max}}{2} \|(\beta_{1_P}^* - \beta_{2_Q}^*)\|_2 \|\Delta_2\|_2 \quad (61)$$

We still need bound to bound  $|e_i X_{i_P}^\top \Delta_1|$  and  $|e_i X_{i_Q}^\top \Delta_2|$  which we do in the following lemma.

**Lemma G.1.** *The following holds:*

1. For fixed  $\|\Delta_1\|_2$ ,  $\mathbb{P}(|e_i X_{i_P}^\top \Delta_1| \leq \frac{\epsilon}{4})$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ .
2. For fixed  $\|\Delta_2\|_2$ ,  $\mathbb{P}(|e_i X_{i_Q}^\top \Delta_2| \leq \frac{\epsilon}{4})$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ .

*Proof.* Recall that  $X_{i_P}^\top \Delta_1$  is a sub-Gaussian random variable with parameter  $\rho^2 \|\Delta_1\|_2^2$  and  $e_i$  is a sub-Gaussian random variable with parameter  $\sigma_e^2$ . Then,  $\frac{X_{i_P}^\top \Delta_1}{\rho \|\Delta_1\|_2} \frac{e_i}{\sigma_e}$  is a sub-exponential random variable with parameters  $(4\sqrt{2}, 2)$ . Using the concentration bounds for the sum of independent sub-exponential random variables (Wainwright, 2019), we can write:

$$\mathbb{P}\left(\left|\frac{X_{i_P}^\top \Delta_1}{\rho \|\Delta_1\|_2} \frac{e_i}{\sigma_e}\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{64}\right), \quad 0 \leq t \leq 8 \quad (62)$$

Taking  $t = \frac{t}{\rho \|\Delta_1\|_2 \sigma_e}$ , we get

$$\mathbb{P}(|e_i X_{i_P}^\top \Delta_1| \geq t) \leq 2 \exp\left(-\frac{t^2}{64\rho^2 \|\Delta_1\|_2^2 \sigma_e^2}\right), \quad 0 \leq t \leq 8\rho \|\Delta_1\|_2 \sigma_e \quad (63)$$

We take  $t = \frac{\epsilon}{4}$ , then

$$\mathbb{P}(|e_i X_{i_P}^\top \Delta_1| \geq \frac{\epsilon}{4}) \leq 2 \exp\left(-\frac{\epsilon^2}{16 \times 64\rho^2 \|\Delta_1\|_2^2 \sigma_e^2}\right), \quad 0 \leq \epsilon \leq 32\rho \|\Delta_1\|_2 \sigma_e \quad (64)$$

Since  $\|\Delta_1\|_2$  is upper bounded with  $\mathcal{O}(\frac{\lambda_1}{n_1} \sqrt{s_1})$  and  $n_1$  is of order  $\mathcal{O}(s_1^3 \log^2 d)$ , thus  $\mathbb{P}(|e_i X_{i_P}^\top \Delta_1| \leq \frac{\epsilon}{4})$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ . Similarly,  $\mathbb{P}(|e_i X_{i_Q}^\top \Delta_2| \leq \frac{\epsilon}{4})$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ .  $\square$

Till now, we have considered  $\|\Delta_1\|_2$  and  $\|\Delta_2\|_2$  to be fixed quantity, however they are also upper bounded by  $\mathcal{O}(\frac{\lambda_1}{n_1} \sqrt{s_1})$  with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ , thus the overall probability that  $u_i \geq 0, i \in \mathcal{I}_1$  is at least  $1 - \mathcal{O}(\frac{1}{d})$  as long as

$$\epsilon \geq 3n_1 C_{\max} \|\Delta_1\|_2^2 - n_2 C_{\min} \|\Delta_2\|_2^2 + 3n C_{\max} \|\beta_{1_P}^* - \beta_{2_Q}^*\|_2 \|\Delta_2\|_2 \quad (65)$$

We need to take a union bound across entries in  $\mathcal{I}_1$  which changes the probability to at least  $1 - \mathcal{O}(\exp(-\log d + \log n_1))$  which is still dominated by  $1 - \mathcal{O}(\frac{1}{d})$ .  $\square$

## H. Proof of Lemma 5.5

**Lemma 5.5** Both  $\Pi$  and  $\Lambda$  have zero eigenvalues corresponding to eigenvectors  $\begin{bmatrix} \tilde{\beta}_1 \\ 1 \end{bmatrix}$  and  $\begin{bmatrix} \tilde{\beta}_2 \\ 1 \end{bmatrix}$  respectively.

*Proof.* It suffices to prove the result for  $\Pi$  as the result for  $\Lambda$  follows in the same way. Note,

$$\begin{aligned} \tilde{\beta}_1 &= \arg \min_{\beta \in \mathbb{R}^{s_1}} \sum_{i=1}^{t_i^*+1} \frac{t_i^*+1}{2} (y_i - X_{i_P}^\top \beta)^2 + \lambda_1 (\|\beta\|_1 + 1)^2 \\ &= \arg \min_{\beta \in \mathbb{R}^{s_1}} \sum_{i \in \mathcal{I}_1} (y_i - X_{i_P}^\top \beta)^2 + \lambda_1 (\|\beta\|_1 + 1)^2 \end{aligned}$$

The optimal  $\tilde{\beta}_1$  must satisfy stationarity KKT condition at the optimum, i.e.,

$$\sum_{i \in \mathcal{I}_1} X_{i_P} (-y_i + X_{i_P}^\top \tilde{\beta}_1) + z \lambda_1 (z^\top \tilde{\beta}_1 + 1) = \mathbf{0}$$

By little algebraic manipulation, we can rewrite the above as following:

$$\left( \sum_{i=1}^n \frac{t_i^*+1}{2} \bar{S}_i^P + \lambda_1 Z + I_\alpha \right) \begin{bmatrix} \tilde{\beta}_1 \\ 1 \end{bmatrix} = \mathbf{0}$$

where  $Z = \begin{bmatrix} z \\ 1 \end{bmatrix} [z^\top \quad 1] = \text{sign}(\bar{W})$ . Clearly,

$$\Pi \begin{bmatrix} \tilde{\beta}_1 \\ 1 \end{bmatrix} = \mathbf{0}$$

Similarly, we can show

$$\Lambda \begin{bmatrix} \tilde{\beta}_2 \\ 1 \end{bmatrix} = \mathbf{0}$$

□

## I. Proof of Lemma 5.6

**Lemma 5.6** If Assumption 4.2 holds and  $n_1 = \Omega(\frac{s_1 + \log d}{C_{\min}^2})$  and  $n_2 = \Omega(\frac{s_2 + \log d}{C_{\min}^2})$ , then the second eigenvalues of  $\Pi$  and  $\Lambda$  are strictly positive with probability at least  $1 - \mathcal{O}(\frac{1}{d})$ , i.e.,  $\text{eig}_2(\Pi) > 0$  and  $\text{eig}_2(\Lambda) > 0$ .

*Proof.* It suffices to prove the result for  $\Pi$  as similar arguments can be used to prove the result for  $\Lambda$ . We know

$$\begin{aligned} \Pi &= \sum_{i \in \mathcal{I}_1} S_i^P + \lambda_1 Z + I_\alpha \\ &= \sum_{i \in \mathcal{I}_1} \begin{bmatrix} X_i X_i^\top & -X_i y_i \\ -y_i X_i^\top & y_i^2 \end{bmatrix} + \lambda_1 \begin{bmatrix} z z^\top & z \\ z^\top & 1 \end{bmatrix} + I_\alpha \\ &= \begin{bmatrix} \sum_{i \in \mathcal{I}_1} X_i X_i^\top + \lambda_1 z z^\top & \sum_{i \in \mathcal{I}_1} -X_i y_i + \lambda_1 z \\ \sum_{i \in \mathcal{I}_1} -y_i X_i^\top + \lambda_1 z^\top & \sum_{i \in \mathcal{I}_1} y_i^2 + \lambda_1 + \alpha \end{bmatrix} \end{aligned} \tag{66}$$

Also note that  $\alpha = -\langle \sum_{i \in \mathcal{I}_1} S_i^P + \lambda_1 Z, \bar{W} \rangle = -\sum_{i \in \mathcal{I}_1} (y_i - X_{i_P}^\top \tilde{\beta}_1)^2 + \lambda_1 (\|\tilde{\beta}_1\|_1 + 1)^2$ . We also know that  $\tilde{\beta}_1$  satisfies the stationarity KKT condition, i.e.,

$$\begin{aligned} \sum_{i \in \mathcal{I}_1} X_{i_P} (-y_i + X_{i_P}^\top \tilde{\beta}_1) + z \lambda_1 (z^\top \tilde{\beta}_1 + 1) &= \mathbf{0} \\ \tilde{\beta}_1 &= -\left( \sum_{i \in \mathcal{I}_1} X_i X_i^\top + \lambda_1 z z^\top \right)^{-1} \left( \sum_{i \in \mathcal{I}_1} -X_i y_i + \lambda_1 z \right) \end{aligned}$$

Using the stationarity KKT condition, we can simplify objective function value of optimization problem (11) at  $\tilde{\beta}_1$  to  $\sum_{i \in \mathcal{I}_1} y_i^2 + (\sum_{i \in \mathcal{I}_1} -y_i X_i^\top + \lambda_1 z^\top) \tilde{\beta}_1 + \lambda_1$ . Now, we invoke Haynesworth's inertia additivity formula (Haynesworth, 1968) to prove our claim. Let  $R$  be a block matrix of the form  $R = \begin{bmatrix} A & B \\ B^\top & C \end{bmatrix}$ , then inertia of matrix  $R$ , denoted by  $\text{In}(R)$ , is defined as the tuple  $(|\text{eig}_+(R)|, |\text{eig}_-(R)|, |\text{eig}_0(R)|)$  where  $|\text{eig}_+(R)|$  is the number of positive eigenvalues,  $|\text{eig}_-(R)|$  is the number of negative eigenvalues and  $|\text{eig}_0(R)|$  is the number of zero eigenvalues of matrix  $R$ . Haynesworth's inertia additivity formula is given as:

$$\text{In}(R) = \text{In}(A) + \text{In}(C - B^\top A^{-1} B) \quad (67)$$

We take  $A = \sum_{i \in \mathcal{I}_1} X_i X_i^\top + \lambda_1 z z^\top$ ,  $B = \sum_{i \in \mathcal{I}_1} -X_i y_i + \lambda_1 z$  and  $C = \sum_{i \in \mathcal{I}_1} y_i^2 + \lambda_1 + \alpha$ . It should be noted that  $C - B^\top A^{-1} B$  evaluates to zero. Thus,

$$\text{In}(\Pi) = \text{In}\left(\sum_{i \in \mathcal{I}_1} X_i X_i^\top + \lambda_1 z z^\top\right) + \text{In}(0) \quad (68)$$

We note that 0 has precisely one zero eigenvalue and no other eigenvalues. Moreover, from Lemma 4.3 and Weyl's inequality:

$$\text{eig}_{\min}\left(\sum_{i \in \mathcal{I}_1} X_i X_i^\top + \lambda_1 z z^\top\right) \geq \frac{C_{\min}}{2} > 0 \quad (69)$$

with probability at least  $1 - \mathcal{O}(\frac{1}{d})$  as long as  $n_1 = \Omega(\frac{s_1 + \log d}{C_{\min}^2})$ . It follows that the second eigenvalue of  $\Pi$  is strictly positive. Similar, arguments can be made for  $\Lambda$ .  $\square$

## J. Proof of Lemma 5.7

**Lemma 5.7** *Let  $\lambda_1 \geq \frac{64\rho\sigma_e}{\xi} \sqrt{n_1 \log d}$ . Then the following holds true:*

$$\begin{aligned} \mathbb{P}\left(\left\|\frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i\right\|_\infty \geq \frac{\xi}{8 - 4\xi}\right) &\leq \mathcal{O}\left(\frac{1}{d}\right), \\ \mathbb{P}\left(\left\|\frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{P^c}} e_i\right\|_\infty \geq \frac{\xi}{8}\right) &\leq \mathcal{O}\left(\frac{1}{d}\right) \end{aligned}$$

*Proof.* We will start with  $\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i$ . We take the  $i$ -th entry of  $\frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i$  for some  $i \in P$ , i.e.,  $|\frac{1}{n} \sum_{j \in \mathcal{I}_1} X_{ji} e_j|$ . Recall that  $X_{ji}$  is a sub-Gaussian random variable with parameter  $\rho^2$  and  $e_j$  is a sub-Gaussian random variable with parameter  $\sigma_e^2$ . Then,  $\frac{X_{ji} e_j}{\rho \sigma_e}$  is a sub-exponential random variable with parameters  $(4\sqrt{2}, 2)$ . Using the concentration bounds for the sum of independent sub-exponential random variables (Wainwright, 2019), we can write:

$$\mathbb{P}\left(\left|\frac{1}{n_1} \sum_{j \in \mathcal{I}_1} \frac{X_{ji} e_j}{\rho \sigma_e}\right| \geq t\right) \leq 2 \exp\left(-\frac{n_1 t^2}{64}\right), \quad 0 \leq t \leq 8 \quad (70)$$

Taking a union bound across  $i \in P$ :

$$\begin{aligned} \mathbb{P}(\exists i \in P \mid \left|\frac{1}{n_1} \sum_{j \in \mathcal{I}_1} \frac{X_{ji} e_j}{\rho \sigma_e}\right| \geq t) &\leq 2s_1 \exp\left(-\frac{n_1 t^2}{64}\right) \\ 0 \leq t &\leq 8 \end{aligned} \quad (71)$$

Taking  $t = \frac{\bar{\lambda}_1 t}{\rho \sigma_e}$ , we get:

$$\begin{aligned} \mathbb{P}(\exists i \in P \mid \left|\frac{1}{\lambda_1} \frac{1}{n_1} \sum_{j=1}^n X_{ji} e_j\right| \geq t) &\leq 2s_1 \exp\left(-\frac{n_1 \bar{\lambda}_1^2 t^2}{64\rho^2 \sigma_e^2}\right) \\ 0 \leq t &\leq 8 \frac{\rho \sigma_e}{\lambda_1} \end{aligned} \quad (72)$$

It follows that  $\|\frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_P} e_i\|_\infty \leq t$  with probability at least  $1 - 2s_1 \exp(-\frac{n\bar{\lambda}_1^2 t^2}{64\rho^2\sigma_e^2})$ .

Using a similar argument, we can show that  $\|\frac{1}{\lambda_1} \frac{1}{n_1} \sum_{i \in \mathcal{I}_1} X_{i_{Pc}} e_i\|_\infty \leq t$  with probability at least  $1 - 2(d - s_1) \exp(-\frac{n_1\bar{\lambda}_1^2 t^2}{64\rho^2\sigma_e^2})$ .

Taking  $t = \frac{\xi}{8-4\xi}$  and  $\frac{\xi}{8}$  in the first and second inequality of Lemma 5.7 and choosing the provided setting of  $\lambda_1$  and  $n_1$  completes our proof.  $\square$

## K. Additional Experiments

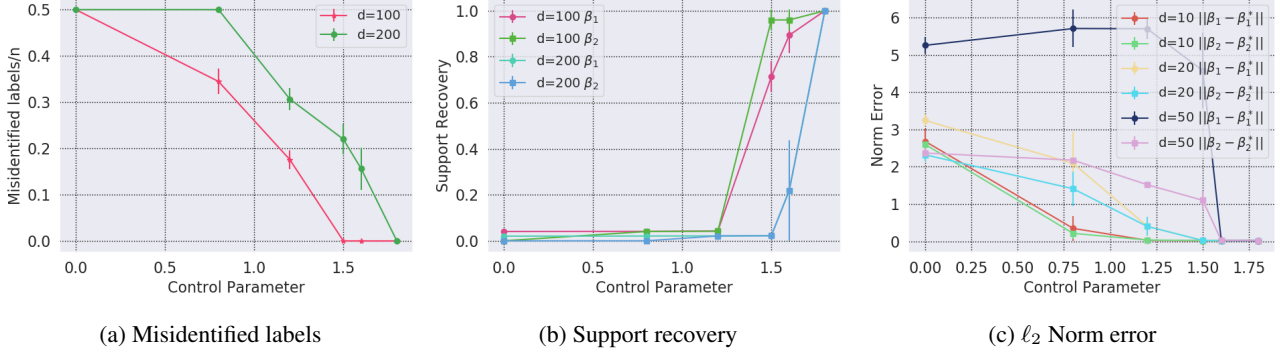


Figure 2. Label and support recovery with control parameter  $C_p$  for high dimensional case  $d = 100, 200$ ). We also show that the norm error indeed goes towards 0.

Following the setting mentioned in Section 6, we conduct further high dimensional ( $d = 100, 200$ ) experiments to validate our theoretical results. We observe a similar trend, i.e., as we increase number of samples, we make zero mistakes in label recovery and achieve 100% correct support recovery for both parameter vectors. Additionally, we also show that the norm error, i.e.,  $\|\beta_j - \beta_j^*\|_2, j \in \{1, 2\}$  goes towards zero in our experiments.