

# Transfer Learning with Multi-source Data: High-dimensional Inference for Group Distributionally Robust Models

Zijian Guo\*

## Abstract

The construction of generalizable and transferable models is a fundamental goal of statistical learning. Learning with the multi-source data helps improve model generalizability and is integral to many important statistical problems, including group distributionally robust optimization, minimax group fairness, and maximin projection. This paper considers multiple high-dimensional regression models for the multi-source data. We introduce the covariate shift maximin effect as a group distributionally robust model. This robust model helps transfer the information from the multi-source data to the unlabelled target population. Statistical inference for the covariate shift maximin effect is challenging since its point estimator may have a non-standard limiting distribution. We devise a novel *DenseNet* sampling method to construct valid confidence intervals for the high-dimensional maximin effect. We show that our proposed confidence interval achieves the desired coverage level and attains a parametric length. Our proposed DenseNet sampling method and the related theoretical analysis are of independent interest in addressing other non-regular or non-standard inference problems. We demonstrate the proposed method over a large-scale simulation and genetic data on yeast colony growth under multiple environments.

**KEYWORDS:** Non-standard Inference; Maximin Effect; DenseNet Sampling Methods; Heterogeneous Data; Minimax Group Fairness.

---

\*Z. Guo is an assistant professor at the Department of Statistics, Rutgers University. The research of Z. Guo was supported in part by the NSF DMS 1811857, 2015373 and NIH R01GM140463, R01LM013614.

# 1 Introduction

Integrative analysis from multiple studies is commonly confronted in applications, such as health record data from multiple hospitals (Singh et al., 2021; Rasmy et al., 2018) and genetic data collected from multiple populations (Keys et al., 2020; Sirugo et al., 2019; Kraft et al., 2009) or under different environments (Bloom et al., 2013). The availability of multi-source data helps enhance the generalizability and transferability of the constructed model. For example, synthesizing information from multiple health record data helps build a generalizable prediction model (Cai et al., 2021). Such a model might be safer to transfer to a target population. However, multi-source data often exhibits a high degree of heterogeneity and creates significant statistical analysis challenges. There is a pressing need to develop methods and theories to construct generalizable models with multi-source data.

The objective of the current paper is to leverage  $L$  training data sets  $\{(Y^{(l)}, X^{(l)})\}_{1 \leq l \leq L}$  and build a robust and generalizable model for a target population (denoted as  $\mathbb{Q}$ ). This problem is also referred to as the multi-source transfer learning in the literature (Yao and Doretto, 2010; Ding et al., 2016, e.g.). In transfer learning, due to the labeling cost, there are often no or very limited outcome labels for the target population (Pan and Yang, 2009; Zhuang et al., 2020). Throughout this paper, we focus on the challenging setting that the target population only has covariate observations but no outcome observations. We introduce the covariate shift maximin effect as the robust and generalizable model.

Our proposed maximin effect satisfies the critical group distributional robustness property: it has a robust predictive performance even if the target population is adversarially generated as any mixture of the multiple source populations. Particularly, for the  $l$ -th group with  $1 \leq l \leq L$ , we consider a high-dimensional regression model for the i.i.d. data  $\{Y_i^{(l)}, X_i^{(l)}\}_{1 \leq i \leq n_l}$ ,

$$Y_i^{(l)} = [X_i^{(l)}]^\top b^{(l)} + \epsilon_i^{(l)} \quad \text{with} \quad b^{(l)} \in \mathbb{R}^p \quad \text{and} \quad \mathbf{E}(\epsilon_i^{(l)} \mid X_i^{(l)}) = 0. \quad (1)$$

We allow the distribution of  $X_i^{(l)}$  to change with the group label  $l$  and do not impose any similarity condition on the regression vectors  $\{b^{(l)}\}_{1 \leq l \leq L}$ . For any prediction model  $\beta \in \mathbb{R}^p$ , we define its adversarial reward concerning all observed groups in (1) and the covariate distribution of the target population. This adversarial reward measures a model's worst predictive performance if the conditional outcome distribution of the target population follows any of the  $L$  outcome models in (1). Our proposed covariate shift maximin effect  $\beta^*(\mathbb{Q}) \in \mathbb{R}^p$  is defined to optimize this adversarial reward; see its definition in (7).

When there are no outcome observations for the target population, most transfer learning methods rely on the assumption that the source and target populations share a similar conditional outcome distribution (Zhuang et al., 2020; Pan and Yang, 2009). In contrast, our transfer learning framework does not require such stringent model assumptions but leverages the multi-source data to construct a robust model.

The covariate shift maximin effect generalizes the maximin effect (Meinshausen and Bühlmann, 2015) by allowing for covariate shift between the target and multiple source populations. The allowance for covariate shift is important for constructing group distributionally robust models in multi-source transfer learning (Sagawa et al., 2019; Hu et al., 2018); see Section 2.1 for details. In addition, it accommodates the applications to minimax group fairness (Martinez et al., 2020; Diana et al., 2021) and maximin projection (Shi et al., 2018); see the detailed discussions in Sections A.5 and A.6 in the supplement.

The maximin effect captures the homogeneous components of  $\{b^{(l)}\}_{1 \leq l \leq L}$  (Meinshausen and Bühlmann, 2015). Particularly, for  $1 \leq j \leq p$ , if the regression coefficients  $\{b_j^{(l)}\}_{1 \leq l \leq L}$  are non-zero and share the same sign, then the maximin effect of the  $j$ -th variable is significant with  $\beta_j^*(\mathbb{Q}) \neq 0$ . If the signs of  $\{b_j^{(l)}\}_{1 \leq l \leq L}$  vary or some of  $\{b_j^{(l)}\}_{1 \leq l \leq L}$  are close to zero, the maximin effect  $\beta_j^*(\mathbb{Q})$  is generally not significant. The maximin effect is instrumental in identifying important variables whose effects are of the homogeneous sign under different environments. In Section 7, we demonstrate this important interpretation in the yeast genomic data and provide the evidence for the “KRE33” gene being a vital gene, which has stable effects on the colony growth across different media (Sharma et al., 2015; Cherry et al., 2012). As observed in Figure 6, even though certain SNP has a significant effect in a single medium, the maximin effect of this SNP can be insignificant, which happens if this SNP has an opposite effect or near null effect in other media.

Despite its importance, there is a lack of inference methods for covariate shift maximin effects, including the construction of confidence interval (CI) and hypothesis testing. In this paper, we point out that inference for the maximin effect is a non-standard inference problem and devise a novel sampling inference method effective for non-standard problems.

## 1.1 Our Results and Contribution

There are distinct inference challenges for maximin effects, which appear in low and high dimensions. The covariate shift maximin effect belongs to the convex cone of the regression vectors  $\{b^{(l)}\}_{1 \leq l \leq L}$ . The convex combination weights are simultaneously determined by  $\{b^{(l)}\}_{1 \leq l \leq L}$  and the covariance matrix of the target covariate distribution. As shown in the

following Section 2.3, the limiting distribution of the weight estimator is non-standard in settings with non-regularity and instability. Consequently, the limiting distribution of the maximin effect estimator is not necessarily normal, and we cannot construct CIs for the maximin effects directly based on the asymptotic normality.

To address this, we propose a DenseNet sampling procedure to construct CIs for the linear combination of the coefficients of the high-dimensional maximin effect. The main novelty is to use the sampling method to quantify the uncertainty of the weight estimation. The intuition of our proposed DenseNet sampling is as follows: if we carefully sample a large number of weight vectors near the estimated weight, there exists at least one sampled weight vector approaching the true weight vector at a rate faster than  $1/\sqrt{n}$ . We provide a rigorous statement of this sampling property in Theorem 1. Our proposed CI is shown to achieve the desired coverage level and attain the  $1/\sqrt{n}$  length.

We conduct a large-scale simulation to evaluate the finite-sample performance. In the challenging scenarios with non-regularity and instability, our proposed CI achieves the desired coverage while the CI assuming asymptotic normality is under-coverage; see Section 6.1. In Section B in the supplement, simulation results show that neither subsampling nor bootstrap provides valid inference for the maximin effect in the presence of non-regularity or instability.

To summarize, the contributions of the current paper are three-fold,

1. We introduce the covariate shift maximin effect as a general robust model with applications to group distributional robustness, minimax fairness, and maximin projections.
2. We propose a novel sampling approach to make inferences for the maximin effect in high dimensions. The DenseNet sampling method is useful for addressing other non-regular or non-standard inference problems.
3. We establish the sampling property in Theorem 1 and characterize the dependence of sampling accuracy on the sampling size. The theoretical argument is new and can be of independent interest for studying other sampling methods.

## 1.2 Related Works

The construction of group distributionally robust models has been investigated in Meinshausen and Bühlmann (2015); Bühlmann and Meinshausen (2015); Sagawa et al. (2019); Hu et al. (2018). The focus of these works was on estimation instead of the CI construction. Furthermore, Rothenhäusler et al. (2016) constructed CIs for the maximin effect in low dimensions, relying on asymptotic normality of the maximin effect estimator. In contrast, we point out in Section 2.3 that the maximin effect estimators are not necessarily asymptotically

normal in the presence of non-regularity or instability. The simulation results in Section 6.1 illustrate that CIs based on asymptotic normality are under-coverage.

Inference for the shared component of regression functions was considered under multiple high-dimensional linear models (Liu et al., 2020) and partially linear models (Zhao et al., 2016). As a significant difference, our proposed method does not require  $\{b^{(l)}\}_{1 \leq l \leq L}$  in (1) to share any similarity, and our model is more flexible in modeling the heterogeneity of multi-source data. Peters et al. (2016); Rothenhäusler et al. (2018); Arjovsky et al. (2019) constructed models satisfying certain invariance principles by analyzing the heterogeneous data. Distributional robustness without the pre-specified group structure has been studied in Gao et al. (2017); Sinha et al. (2017) by minimizing a worst-case over a class of distributions.

Inference for the regression coefficients in a single high-dimensional linear model was actively investigated in the recent decade (Zhang and Zhang, 2014; van de Geer et al., 2014; Javanmard and Montanari, 2014; Belloni et al., 2014; Chernozhukov et al., 2015; Farrell, 2015; Chernozhukov et al., 2018). Inference for linear functionals of the high-dimensional regression vector was studied in Cai and Guo (2017); Athey et al. (2018); Zhu and Bradic (2018); Cai et al. (2021). Inference for the maximin effect creates additional challenges as discussed in Section 2.3 and requires novel methods and theories to address these challenges.

Sampling methods have a long history in statistics, such as, bootstrap (Efron, 1979; Efron and Tibshirani, 1994), subsampling (Politis et al., 1999), generalized fiducial inference (Zabell, 1992; Xie and Singh, 2013; Hannig et al., 2016) and repro sampling (Wang and Xie, 2020). In contrast, instead of directly sampling from the original data, we sample the estimator of the regression covariance matrix, which makes our proposed sampling method computationally efficient; see Remark 5.

**Paper Organization.** In Section 2, we introduce the covariate shift maximin effect and the inference challenges. Our proposed method is detailed in Section 3 and the theoretical justification is provided in Section 4. In Section 5, we discuss the stability issue. In Sections 6 and 7, we investigate the numerical performance of our proposed method on simulated and genomic data sets, respectively. We provide conclusion and discussion in Section 8.

**Notations.** Define  $n = \min_{1 \leq l \leq L} \{n_l\}$ . Let  $[p] = \{1, 2, \dots, p\}$ . For a set  $S$ ,  $|S|$  denotes the cardinality of  $S$  and  $S^c$  denotes its complement. For a vector  $x \in \mathbb{R}^p$  and a subset  $S \subset [p]$ ,  $x_S$  is the sub-vector of  $x$  with indices in  $S$  and  $x_{-S}$  is the sub-vector with indices in  $S^c$ . The  $\ell_q$  norm of a vector  $x$  is defined as  $\|x\|_q = (\sum_{l=1}^p |x_l|^q)^{\frac{1}{q}}$  for  $q \geq 0$  with  $\|x\|_0 = |\{1 \leq l \leq p : x_l \neq 0\}|$  and  $\|x\|_\infty = \max_{1 \leq l \leq p} |x_l|$ . For a matrix  $X$ ,  $X_i$  and  $X_{\cdot,j}$  are used to denote its  $i$ -th row and  $j$ -th column, respectively; for index sets  $S_1$  and  $S_2$ ,  $X_{S_1, S_2}$  denotes

the sub-matrix of  $X$  with row and column indices belonging to  $S_1$  and  $S_2$ , respectively;  $X_{S_1}$  denotes the sub-matrix of  $X$  with row indices belonging to  $S_1$ . Let  $X_i \stackrel{\text{i.i.d.}}{\sim} \mathbb{Q}$  for  $1 \leq i \leq n$  denote that the sequence random vectors  $\{X_i\}_{1 \leq i \leq n}$  are i.i.d. following the distribution  $\mathbb{Q}$ . We use  $c$  and  $C$  to denote generic positive constants that may vary from place to place. For two positive sequences  $a_n$  and  $b_n$ ,  $a_n \lesssim b_n$  means that  $\exists C > 0$  such that  $a_n \leq Cb_n$  for all  $n$ ;  $a_n \asymp b_n$  if  $a_n \lesssim b_n$  and  $b_n \lesssim a_n$ , and  $a_n \ll b_n$  if  $\limsup_{n \rightarrow \infty} a_n/b_n = 0$ . For a matrix  $A$ , we use  $\|A\|_F$ ,  $\|A\|_2$  and  $\|A\|_\infty$  to denote its Frobenius norm, spectral norm and element-wise maximum norm, respectively. For a symmetric matrix  $A \in \mathbb{R}^{L \times L}$  with eigen-decomposition  $A = U\Lambda U^\top$ , we use  $\lambda_{\max}(A)$  and  $\lambda_{\min}(A)$  to denote its maximum and minimum eigenvalues, respectively; define  $A_+ = U\Lambda_+U^\top$  and  $A^{1/2} = U\Lambda^{1/2}U^\top$  with  $(\Lambda_+)_{l,l} = \max\{\Lambda_{l,l}, 0\}$  and  $(\Lambda^{1/2})_{l,l} = \sqrt{\Lambda_{l,l}}$  for  $1 \leq l \leq L$ . We use  $A^{-1/2}$  to denote the inverse of  $A^{1/2}$ . For a matrix  $B$ , we use  $\lambda_j(B)$  to denote its  $j$ -th largest singular value. For a symmetric matrix  $D \in \mathbb{R}^{L \times L}$ , we use  $\text{vecl}(D) \in \mathbb{R}^{L(L+1)/2}$  to denote the long vector which stacks the columns of the lower triangle part of  $D$ . Define  $\mathcal{I}_L = \{(l, k) : 1 \leq k \leq l \leq L\}$  as the index set of the lower triangular part of  $D$  and  $[L(L+1)/2]$  as the index set of  $\text{vecl}(D)$ . We define the one-to-one index mapping  $\pi$  from  $\mathcal{I}_L$  to  $[L(L+1)/2]$  as

$$\pi(l, k) = \frac{(2L - k)(k - 1)}{2} + l \quad \text{for} \quad (l, k) \in \mathcal{I}_L := \{(l, k) : 1 \leq k \leq l \leq L\}. \quad (2)$$

For  $(l, k) \in \mathcal{I}_L$ , we have  $[\text{vecl}(D)]_{\pi(l, k)} = D_{l, k}$ .

## 2 Maximin Effects: Definition, Identification and Challenges

In Section 2.1, we motivate the covariate shift maximin effect in the multi-source transfer learning framework. In Section 2.2, we present the identification and interpretation of the maximin effect. In Section 2.3, we discuss the inference challenges for the maximin effect.

### 2.1 Multi-source Transfer Learning: Group Distributional Robustness

We consider  $L$  groups of training data  $\{(Y^{(l)}, X^{(l)})\}_{1 \leq l \leq L}$ , which are generated from multiple source populations. The heterogeneity among these training data sets may occur if they are collected under different environments, e.g., the data sets from different healthcare centers. We aim to build a robust prediction model for a target population (e.g., a new healthcare center) by leveraging the multi-source training data sets.

We first introduce the setup for the multi-source transfer learning problem. For the  $l$ -th group of data  $\{X_i^{(l)}, Y_i^{(l)}\}_{1 \leq i \leq n_l}$  with  $1 \leq l \leq L$ , let  $\mathbb{P}_X^{(l)}$  denote the distribution of  $X_i^{(l)} \in \mathbb{R}^p$

and  $\mathbb{P}_{Y|X}^{(l)}$  denote the conditional distribution of the outcome  $Y_i^{(l)}$  given  $X_i^{(l)}$ . We write

$$X_i^{(l)} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_X^{(l)}, \quad Y_i^{(l)} | X_i^{(l)} \stackrel{\text{i.i.d.}}{\sim} \mathbb{P}_{Y|X}^{(l)} \quad \text{for } 1 \leq i \leq n_l. \quad (3)$$

For the test data  $\{X_i^{\mathbb{Q}}, Y_i^{\mathbb{Q}}\}_{1 \leq i \leq N_{\mathbb{Q}}}$ , we consider the target population,

$$X_i^{\mathbb{Q}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{Q}_X, \quad Y_i^{\mathbb{Q}} | X_i^{\mathbb{Q}} \stackrel{\text{i.i.d.}}{\sim} \mathbb{Q}_{Y|X} \quad \text{for } 1 \leq i \leq N_{\mathbb{Q}}, \quad (4)$$

with  $\mathbb{Q}_X$  and  $\mathbb{Q}_{Y|X}$  denoting the covariate distribution and the conditional outcome model of the target population, respectively. We focus on the setting without any labelled test data, that is, only  $\{X_i^{\mathbb{Q}}\}_{1 \leq i \leq N_{\mathbb{Q}}}$  are observed but  $\{Y_i^{\mathbb{Q}}\}_{1 \leq i \leq N_{\mathbb{Q}}}$  are missing. This is commonly encountered in transfer learning applications, where the labels for the target population are hard to obtain (Zhuang et al., 2020; Pan and Yang, 2009).

In transfer learning, the covariate shift between the source and target populations commonly exists (Zhuang et al., 2020; Pan and Yang, 2009). For example, the new healthcare center is likely to have a different covariate distribution from existing healthcare centers, where the training data sets were collected. To better accommodate this, we allow the target covariate distribution  $\mathbb{Q}_X$  to be different from any of  $\{\mathbb{P}_X^{(l)}\}_{1 \leq l \leq L}$ . In addition, the target conditional outcome model  $\mathbb{Q}_{Y|X}$  might also be different from any of  $\{\mathbb{P}_{Y|X}^{(l)}\}_{1 \leq l \leq L}$ . We assume that the conditional outcome model  $\mathbb{Q}_{Y|X}$  is a mixture of  $\{\mathbb{P}_{Y|X}^{(l)}\}_{1 \leq l \leq L}$ ,

$$\mathbb{Q}_{Y|X} = \sum_{l=1}^L q_l^* \cdot \mathbb{P}_{Y|X}^{(l)}, \quad \text{with } q^* \in \Delta_L := \left\{ q \in \mathbb{R}^L : \sum_{l=1}^L q_l = 1, \min_{1 \leq l \leq L} q_l \geq 0 \right\}. \quad (5)$$

As a special case, if  $q_l^* = 1$  and  $q_j^* = 0$  for  $j \neq l$ , the conditional outcome model  $\mathbb{Q}_{Y|X}$  for the target population is the same as the outcome model  $\mathbb{P}_{Y|X}^{(l)}$  for the  $l$ -th source population.

We motivate the maximin effect concept from the construction of a group distributionally robust model (Sagawa et al., 2019; Hu et al., 2018). Since the weight vector  $q^* \in \mathbb{R}^L$  in (5) is not identifiable with the unlabelled test data, the target distribution can be any mixture of the  $L$  source populations. We shall fit a model such that it has a reasonable predictive performance for any such target population. We define a class of joint distributions on  $(X_i^{\mathbb{Q}}, Y_i^{\mathbb{Q}})$  as,

$$\mathcal{C}(\mathbb{Q}_X) = \left\{ \mathbb{T} = (\mathbb{Q}_X, \mathbb{T}_{Y|X}) : \mathbb{T}_{Y|X} = \sum_{l=1}^L q_l \cdot \mathbb{P}_{Y|X}^{(l)} \quad \text{with } q \in \Delta_L \right\}, \quad (6)$$

where  $\Delta_L$  is defined in (5). The distribution class  $\mathcal{C}(\mathbb{Q}_X)$  contains the true distribution  $\mathbb{Q} = (\mathbb{Q}_X, \mathbb{Q}_{Y|X})$  of the target population, which is not identifiable with the unlabelled test data. In the definition (6), the covariate distribution of the class  $\mathcal{C}(\mathbb{Q}_X)$  is fixed at  $\mathbb{Q}_X$  since  $\mathbb{Q}_X$  is identifiable with the covariate data  $\{X_i^{\mathbb{Q}}\}_{1 \leq i \leq N_{\mathbb{Q}}}$ ; the conditional outcome model of the class  $\mathcal{C}(\mathbb{Q}_X)$  contains any mixture of  $\{\mathbb{P}_{Y|X}^{(l)}\}_{1 \leq l \leq L}$ , including the true conditional outcome model  $\mathbb{Q}_{Y|X}$  for the target population and  $\{\mathbb{P}_{Y|X}^{(l)}\}_{1 \leq l \leq L}$  for  $L$  source populations.

We propose to optimize the worst-case reward defined with respect to the class  $\mathcal{C}(\mathbb{Q}_X)$  and introduce a group distributionally robust optimization problem as

$$\beta^*(\mathbb{Q}) := \arg \max_{\beta \in \mathbb{R}^p} R_{\mathbb{Q}}(\beta) \quad \text{with} \quad R_{\mathbb{Q}}(\beta) = \min_{\mathbb{T} \in \mathcal{C}(\mathbb{Q}_X)} \{ \mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^\top \beta)^2 \}, \quad (7)$$

where  $\mathcal{C}(\mathbb{Q}_X)$  is defined in (6) and  $\mathbf{E}_{\mathbb{T}}$  denotes the expectation taken with respect to the distribution  $\mathbb{T}$  belonging to  $\mathcal{C}(\mathbb{Q}_X)$ . For any model  $\beta \in \mathbb{R}^p$ , the reward function  $\mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^\top \beta)^2$  measures the variance explained by  $\beta$  for the test data  $\{X_i, Y_i\} \sim \mathbb{T}$ . By taking the minimum with respect to the distribution  $\mathbb{T}$ ,  $\min_{\mathbb{T} \in \mathcal{C}(\mathbb{Q}_X)} \{ \mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^\top \beta)^2 \}$  denotes the adversarial reward among the class  $\mathcal{C}(\mathbb{Q}_X)$ . The optimizer  $\beta^*(\mathbb{Q})$  defined in (7) optimizes this adversarial reward. We refer to  $\beta^*(\mathbb{Q})$  as the covariate shift maximin effect.

The covariate shift maximin effect  $\beta^*(\mathbb{Q})$  can be interpreted as the robust linear model guaranteeing an excellent predictive performance even when the target population is adversarially generated from the class  $\mathcal{C}(\mathbb{Q}_X)$ . We further interpret its robustness from the perspective of a two-side game: we select a model  $\beta$ , and the counter agent generates the adversarial target population by taking a mixture of  $L$  conditional outcome models in (1).  $\beta^*(\mathbb{Q})$  guarantees the optimal prediction accuracy for such an adversarial target population. An equivalent definition of  $\beta^*(\mathbb{Q})$  is presented in Section A.2 in the supplement.

As a special case, we consider the no covariate shift setting with  $\mathbb{P}_X$  denoting the shared covariate distribution, that is  $\mathbb{P}_X^{(l)} = \mathbb{P}_X$  for  $1 \leq l \leq L$  and  $\mathbb{Q}_X = \mathbb{P}_X$ . We simplify  $\beta^*(\mathbb{Q})$  in (7) as  $\beta^* := \arg \max_{\beta \in \mathbb{R}^p} \min_{\mathbb{T} \in \mathcal{C}(\mathbb{P}_X)} \{ \mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^\top \beta)^2 \}$ , where  $\mathcal{C}(\mathbb{P}_X)$  is defined in (6) with  $\mathbb{Q}_X = \mathbb{P}_X$ . The no covariate shift maximin effect  $\beta^*$  is equivalent to the maximin effect introduced in Meinshausen and Bühlmann (2015),

$$\beta^* := \arg \max_{\beta \in \mathbb{R}^p} R(\beta) \quad \text{with} \quad R(\beta) = \min_{1 \leq l \leq L} \{ \mathbf{E}(Y_i^{(l)})^2 - \mathbf{E}(Y_i^{(l)} - [X_i^{(l)}]^\top \beta)^2 \}.$$

See its proof in Section A.1 in the supplement. The maximin effect in (7) can be expressed as an equivalent minimax estimator  $\beta^*(\mathbb{Q}) := \arg \min_{\beta \in \mathbb{R}^p} \max_{\mathbb{T} \in \mathcal{C}(\mathbb{Q}_X)} \{ \mathbf{E}_{\mathbb{T}} \ell(Y_i, X_i^\top \beta) \}$  with



$\ell(Y_i, X_i^\top \beta) = (Y_i - X_i^\top \beta)^2 - Y_i^2$ , which is in the form of the group distributionally robust optimization (Sagawa et al., 2019; Hu et al., 2018).

**Remark 1.** Li et al. (2020); Tian and Feng (2021) focused on multi-source transfer learning with access to outcome labels for the target population. Our current paper focuses on a fundamentally different setting. Since the target population has no outcome labels, we cannot identify its outcome model but focus on the group distributional robust model.

**Remark 2.** The definition of  $\beta^*(\mathbb{Q})$  incorporates two types of heterogeneity: covariate shift between the source and target populations. The conditional outcome distribution for the target population is a mixture of those for the source populations. Under our framework, the target conditional distribution  $\mathbb{Q}_{Y|X}$  in (4) is allowed to be different from any of  $\{\mathbb{P}_{Y|X}^{(l)}\}_{1 \leq l \leq L}$ .

## 2.2 Identification and Interpretation of $\beta^*(\mathbb{Q})$

We now present the identification of  $\beta^*(\mathbb{Q})$  in (7) and focus on linear conditional expectation models for the source populations,

$$\mathbf{E}(Y_i^{(l)} | X_i^{(l)}) = [X_i^{(l)}]^\top b^{(l)}, \quad \text{for } 1 \leq l \leq L, 1 \leq i \leq n_l. \quad (8)$$

The model (8) can be extended to handle non-linear conditional expectation if  $X_i^{(l)}$  contains the basis transformation of the observed covariates. Under (8), we simplify (7) as

$$\beta^*(\mathbb{Q}) = \arg \max_{\beta \in \mathbb{R}^p} R_{\mathbb{Q}}(\beta) \quad \text{with} \quad R_{\mathbb{Q}}(\beta) = \min_{b \in \mathbb{B}} [2b^\top \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta], \quad (9)$$

where  $\mathbb{B} = \{b^{(1)}, \dots, b^{(L)}\}$  denotes the set of  $L$  regression vectors and  $\Sigma^{\mathbb{Q}} = \mathbf{E}X_1^{\mathbb{Q}}(X_1^{\mathbb{Q}})^\top$ . See its proof in Section A.1 in the supplement.

The following proposition shows how to identify the maximin effect  $\beta^*(\mathbb{Q})$ .

**Proposition 1.** *If  $\lambda_{\min}(\Sigma^{\mathbb{Q}}) > 0$ , then  $\beta^*(\mathbb{Q})$  defined in (9) is identified as*

$$\beta^*(\mathbb{Q}) = \sum_{l=1}^L [\gamma^*(\mathbb{Q})]_l b^{(l)} \quad \text{with} \quad \gamma^*(\mathbb{Q}) := \arg \min_{\gamma \in \Delta^L} \gamma^\top \Gamma^{\mathbb{Q}} \gamma \quad (10)$$

where  $\Gamma_{lk}^{\mathbb{Q}} = (b^{(l)})^\top \Sigma^{\mathbb{Q}} b^{(k)}$  for  $1 \leq l, k \leq L$  and  $\Delta^L = \{\gamma \in \mathbb{R}^L : \gamma_j \geq 0, \sum_{j=1}^L \gamma_j = 1\}$  is the simplex over  $\mathbb{R}^L$ . Furthermore,  $\max_{\beta \in \mathbb{R}^p} \min_{b \in \mathbb{B}} [2b^\top \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta] = [\beta^*(\mathbb{Q})]^\top \Sigma^{\mathbb{Q}} \beta^*(\mathbb{Q})$ .

Proposition 1 provides an explicit way of computing  $\beta^*(\mathbb{Q})$  by firstly identifying  $\{b^{(l)}\}_{1 \leq l \leq L}$  and  $\Gamma^{\mathbb{Q}} \in \mathbb{R}^{L \times L}$ . Proposition 1 generalizes Theorem 1 in Meinshausen and Bühlmann (2015),

which was focused on the no covariate shift setting. The covariate shift maximin effect  $\beta^*(\mathbb{Q})$  is a convex combination of  $\{b^{(l)}\}_{1 \leq l \leq L}$  which is closest to the origin. Here, the distance measure depends on the covariance matrix  $\Sigma^{\mathbb{Q}}$  for the target population.

The maximin effect  $\beta^*(\mathbb{Q})$  is not only a group distributionally robust model but is interpreted as the model capturing the shared effects of the heterogeneous vectors  $\{b^{(l)}\}_{1 \leq l \leq L}$ ; see Figure 1 and the related discussion in [Meinshausen and Bühlmann \(2015\)](#). Particularly, for some  $1 \leq j \leq p$ , if the signs of  $\{b_j^{(l)}\}_{1 \leq l \leq L}$  are heterogeneous or some of  $\{b_j^{(l)}\}_{1 \leq l \leq L}$  are scattered around zero, the maximin effect  $\beta_j^*(\mathbb{Q})$  is shrunk towards zero. If  $\{b_j^{(l)}\}_{1 \leq l \leq L}$  are non-zero and of the same sign, the maximin effect is significant with  $\beta_j^*(\mathbb{Q}) \neq 0$ . We further demonstrate this interpretation in our real data analysis in Section 7; see Figure 6.

The maximin significance test  $H_{0,j} : \beta_j^*(\mathbb{Q}) = 0$  for  $1 \leq j \leq p$  is of great importance for both scientific discovery and the construction of robust prediction model. The maximin significance indicates that the effect of the  $j$ -th covariate is homogeneously positive or negative across different environments. This indicates that the  $j$ -th covariate is likely to have a similar effect for a new environment, and the effect of the  $j$ -th variable might have some causal interpretation. In addition, the non-zero maximin effect suggests that it would be helpful to include the  $j$ -th covariate in the prediction model for the target population.

Statistical inference for a linear combination of the maximin effect is well motivated from constructing optimal treatment regime with heterogeneous training data ([Shi et al., 2018](#)). With  $x_{\text{new}}$  denoting a future covariate observation, the sign of  $x_{\text{new}}^T \beta^*(\mathbb{Q})$  is instrumental in designing the optimal treatment regime, which motivates the hypothesis testing problem  $H_0 : x_{\text{new}}^T \beta^*(\mathbb{Q}) < 0$  for any  $x_{\text{new}} \in \mathbb{R}^p$ . See more details in Section A.6 in the supplement.

**Remark 3.** In the covariate shift setting, a collection of works ([Tsuboi et al., 2009](#); [Shimodaira, 2000](#); [Sugiyama et al., 2007](#), e.g.) were focused on the misspecified conditional outcome models. In contrast, we focus on the correctly specified conditional outcome model (8). Consequently, the best linear approximation  $\{b^{(l)}\}_{1 \leq l \leq L}$  does not change with the target population  $\mathbb{Q}_X$ . However, the maximin effect  $\beta^*(\mathbb{Q})$  changes with the target population since the weight  $\gamma^*(\mathbb{Q})$  is determined by the target covariate distribution  $\mathbb{Q}_X$ .

### 2.3 Statistical Inference Challenges: Non-regularity and Instability

The inference challenges arise from the fact that estimators of  $\gamma^*(\mathbb{Q})$  and  $\beta^*(\mathbb{Q})$  may have a non-standard limiting distribution. To demonstrate the challenges, we consider the special

case  $L = 2$  and  $n_1 = n_2 = n$ . The optimal weight is  $\gamma^*(\mathbb{Q}) = (\gamma_1^*, 1 - \gamma_1^*)^\top$  with

$$\gamma_1^* = \min \left\{ \max \left\{ \frac{\Gamma_{22}^{\mathbb{Q}} - \Gamma_{12}^{\mathbb{Q}}}{\Gamma_{11}^{\mathbb{Q}} + \Gamma_{22}^{\mathbb{Q}} - 2\Gamma_{12}^{\mathbb{Q}}}, 0 \right\}, 1 \right\}. \quad (11)$$

With  $\hat{\Gamma}^{\mathbb{Q}}$  defined in the following equation (23), we estimate  $\gamma_1^*$  by  $\hat{\gamma}_1 = \min \{ \max \{ \bar{\gamma}_1, 0 \}, 1 \}$  where  $\bar{\gamma}_1 = \frac{\hat{\Gamma}_{22}^{\mathbb{Q}} - \hat{\Gamma}_{12}^{\mathbb{Q}}}{\hat{\Gamma}_{11}^{\mathbb{Q}} + \hat{\Gamma}_{22}^{\mathbb{Q}} - 2\hat{\Gamma}_{12}^{\mathbb{Q}}}$ . We illustrate two settings where  $\hat{\gamma}_1$  may not have a standard limiting distribution even if  $\hat{\Gamma}^{\mathbb{Q}}$  is unbiased and asymptotically normal.

The first is the non-regularity setting due to the boundary effect. Specifically, the estimation error  $\hat{\gamma}_1 - \gamma_1^*$  is decomposed as a mixture distribution,

$$\sqrt{n}(\bar{\gamma}_1 - \gamma_1^*) \cdot \mathbf{1}\{0 < \bar{\gamma}_1 < 1\} + (-\sqrt{n}\gamma_1^*) \cdot \mathbf{1}\{\bar{\gamma}_1 \leq 0\} + \sqrt{n}(1 - \gamma_1^*) \cdot \mathbf{1}\{\bar{\gamma}_1 \geq 1\},$$

where the last two terms appear due to the boundary constraint  $0 \leq \gamma_1^* \leq 1$ . When  $|(b^{(1)} - b^{(2)})^\top \Sigma^{\mathbb{Q}} b^{(2)}| = c/\sqrt{n}$  for some  $c > 0$ , then  $\gamma_1^* = \gamma_1^*(n, p) \asymp 1/\sqrt{n}$  and  $\sqrt{n}(\hat{\gamma}_1 - \gamma_1^*)$  is likely a mixture of  $\sqrt{n}(\bar{\gamma}_1 - \gamma_1^*)$  and the point mass  $-\sqrt{n}\gamma_1^*$ . It is well known that boundary constraints lead to estimators with non-standard limiting distributions; see [Self and Liang \(1987\)](#); [Andrews \(1999\)](#); [Drton \(2009\)](#) and the references therein. Similarly, the boundary effect for the maximin effect leads to a non-standard distribution for the corresponding maximin effect estimator. Another challenge is the instability, which occurs when some of  $\{b^{(l)}\}_{1 \leq l \leq L}$  are similar to each other. For  $L = 2$ , if  $b^{(1)} \approx b^{(2)}$ , then  $\Gamma_{11}^{\mathbb{Q}} + \Gamma_{22}^{\mathbb{Q}} - 2\Gamma_{12}^{\mathbb{Q}}$  in (11) is close to zero. It is hard to accurately estimate  $\gamma_1^*$  since a small error in estimating  $\Gamma_{11}^{\mathbb{Q}} + \Gamma_{22}^{\mathbb{Q}} - 2\Gamma_{12}^{\mathbb{Q}}$  may lead to a large error of estimating  $\gamma_1^*$ . As illustrated in Figure 2, both non-regularity and instability lead to a non-standard limiting distribution of the maximin estimator, which results in the under-coverage of CIs assuming the asymptotic normality; see the numerical illustration in Section 6.1.

For non-regular settings due to the boundary effect, inference methods based on asymptotic normality or bootstrap fail to work ([Andrews, 2000](#), e.g.). In Section B in the supplement, we illustrate that even for the low-dimensional setting, neither  $m$  out of  $n$  bootstrap nor subsampling provides valid inference for the maximin effect in the presence of non-regularity or instability.

To address the inference challenges, we shall devise a novel sampling approach in Section 3, which is a valid inference procedure even in the presence of non-regularity or instability.

### 3 The DenseNet Sampling Method

In Section 3.1, we discuss the estimation of  $x_{\text{new}}^\top \beta^*(\mathbb{Q})$  for any  $x_{\text{new}} \in \mathbb{R}^p$ . In Section 3.2, we devise a sampling method and construct the CI for  $x_{\text{new}}^\top \beta^*(\mathbb{Q})$ . When there is no confusion, we will omit the dependence on  $\mathbb{Q}$  and write  $\Gamma^\mathbb{Q}, \widehat{\Gamma}^\mathbb{Q}, \beta^*(\mathbb{Q}), \gamma^*(\mathbb{Q})$  as  $\Gamma, \widehat{\Gamma}, \beta^*, \gamma^*$ , respectively.

#### 3.1 Estimation of $x_{\text{new}}^\top \beta^*$

To estimate  $x_{\text{new}}^\top \beta^*$ , we construct debiased estimators of  $\{x_{\text{new}}^\top b^{(l)}\}_{1 \leq l \leq L}$  and estimate the aggregation weight  $\gamma^*$  based on a debiased estimator of  $\Gamma^\mathbb{Q}$ . Specifically, for  $1 \leq l \leq L$ , we follow Cai et al. (2021) and construct the debiased estimator of  $x_{\text{new}}^\top b^{(l)}$  as

$$\widehat{x_{\text{new}}^\top b^{(l)}} = x_{\text{new}}^\top \widehat{b}^{(l)} + [\widehat{v}^{(l)}]^\top \frac{1}{n_l} (X^{(l)})^\top (Y^{(l)} - X^{(l)} \widehat{b}^{(l)}), \quad (12)$$

where  $\widehat{b}^{(l)}$  is the Lasso estimator based on  $(X^{(l)}, Y^{(l)})$  and  $\widehat{v}^{(l)} \in \mathbb{R}^p$  is constructed as

$$\widehat{v}^{(l)} = \arg \min_{v \in \mathbb{R}^p} v^\top \frac{1}{n_l} (X^{(l)})^\top X^{(l)} v \quad \text{s.t.} \quad \max_{w \in \mathcal{F}(x_{\text{new}})} \left| \langle w, \frac{1}{n_l} (X^{(l)})^\top X^{(l)} v - x_{\text{new}} \rangle \right| \leq \eta_l \quad (13)$$

$$\|X^{(l)} v\|_\infty \leq \|x_{\text{new}}\|_2 \tau_l$$

with  $\mathcal{F}(x_{\text{new}}) = \{e_1, \dots, e_p, x_{\text{new}}/\|x_{\text{new}}\|_2\}$ ,  $\eta_l \asymp \|x_{\text{new}}\|_2 \sqrt{\log p/n_l}$ , and  $\tau_l \asymp \sqrt{\log n_l}$ . If  $x_{\text{new}}$  is taken as the  $j$ -th Euclidean basis,  $\widehat{x_{\text{new}}^\top b^{(l)}}$  is reduced to the debiased estimator proposed in Zhang and Zhang (2014); Javanmard and Montanari (2014).  $\widehat{x_{\text{new}}^\top b^{(l)}}$  is shown to be asymptotically unbiased and normal in Corollary 5 in Cai et al. (2021).

In Section 3.3, we propose a bias-corrected estimator  $\widehat{\Gamma}^\mathbb{Q}$  in (23), which is shown to satisfy

$$\text{vecl}(\widehat{\Gamma}^\mathbb{Q} - \Gamma^\mathbb{Q}) \stackrel{d}{\approx} \mathcal{N}(\mathbf{0}, \widehat{\mathbf{V}}), \quad (14)$$

where  $\text{vecl}(\widehat{\Gamma}^\mathbb{Q} - \Gamma^\mathbb{Q})$  is the long vector which stacks the columns of the lower triangular part of  $\widehat{\Gamma}^\mathbb{Q} - \Gamma^\mathbb{Q}$ ,  $\stackrel{d}{\approx}$  stands for approximately equal in distribution and  $\widehat{\mathbf{V}}$  is the estimated covariance matrix defined in the following equation (24). Then we estimate  $\gamma^*$  by replacing  $\Gamma^\mathbb{Q}$  in the definition  $\gamma^* := \arg \min_{\gamma \in \Delta^L} \gamma^\top \Gamma^\mathbb{Q} \gamma$  with  $\widehat{\Gamma}^\mathbb{Q}$ . We construct the point estimator of  $x_{\text{new}}^\top \beta^*$  as

$$\widehat{x_{\text{new}}^\top \beta^*} = \sum_{l=1}^L \widehat{\gamma}_l \cdot \widehat{x_{\text{new}}^\top b^{(l)}} \quad \text{with} \quad \widehat{\gamma} := \arg \min_{\gamma \in \Delta^L} \gamma^\top \widehat{\Gamma}^\mathbb{Q} \gamma, \quad (15)$$

where  $\Delta^L$  is defined in (5).

To construct the CI for  $x_{\text{new}}^\top \beta^*$ , we need to quantify the uncertainty of the estimators  $\{\widehat{x_{\text{new}}^\top b^{(l)}}\}_{1 \leq l \leq L}$  and  $\widehat{\gamma}$ . The main challenge is on the uncertainty quantification of  $\widehat{\gamma}$ , which might not have a standard limiting distribution. We propose the DenseNet sampling to quantify the weight estimation uncertainty in the following subsection.

### 3.2 Inference for $x_{\text{new}}^\top \beta^*$ : The DenseNet Sampling

The key idea of the DenseNet sampling is as follows: we carefully construct many sampled weight vectors near the point estimator  $\widehat{\gamma}$  defined in (15); our construction guarantees that at least one of these sampled weight vectors is nearly the same as the true weight vector  $\gamma^*$ . Then, with this particular sampled weight vector, we only need to quantify the uncertainty due to the debiased Lasso estimators  $\{\widehat{x_{\text{new}}^\top b^{(l)}}\}_{1 \leq l \leq L}$ . The DenseNet sampling is instrumental in quantifying the uncertainty of  $\widehat{\gamma}$ , which might have a non-standard limiting distribution.

We refer to our proposed method as DenseNet by its analogy to fishing with a dense net: if the net is dense enough and cast over a small area known to contain fish, there are high chances to catch at least one fish. We propose an efficient sampling method to build the dense fishing net, and the caught fish is  $\gamma^*$ .

We describe the DenseNet sampling via two steps: sampling and aggregation.

**Step 1: Sampling the weight vectors  $\{\widehat{\gamma}^{[m]}\}_{1 \leq m \leq M}$ .** We sample  $\{\widehat{\Gamma}^{[m]}\}_{1 \leq m \leq M}$  such that  $\widehat{\Gamma}^{[m]} - \widehat{\Gamma}^\mathbb{Q}$  approximately follows the asymptotic distribution  $\mathcal{N}(\mathbf{0}, \widehat{\mathbf{V}})$  in (14) and then construct the sampled weight vectors  $\{\widehat{\gamma}^{[m]}\}_{1 \leq m \leq M}$  by solving the optimization problem,

$$\widehat{\gamma}^{[m]} = \arg \min_{\gamma \in \Delta^L} \gamma^\top \widehat{\Gamma}_+^{[m]} \gamma. \quad (16)$$

In the following, we provide the details on how to sample  $\{\widehat{\Gamma}^{[m]}\}_{1 \leq m \leq M}$ . Conditioning on the observed data, we generate i.i.d. samples  $\{S^{[m]}\}_{1 \leq m \leq M}$  as

$$S^{[m]} \sim \mathcal{N}(\mathbf{0}, \widehat{\mathbf{V}} + d_0/n \cdot \mathbf{I}) \quad \text{with} \quad d_0 = \max \left\{ \tau_0 \cdot \max_{(l,k) \in \mathcal{I}_L} \left\{ n \cdot \widehat{\mathbf{V}}_{\pi(l,k), \pi(l,k)} \right\}, 1 \right\}, \quad (17)$$

where  $M$  is the sampling size (default set as 500),  $\tau_0 > 0$  is a positive tuning parameter (default set as 0.2), and  $\mathbf{I}$  is the identity matrix of the same dimension as  $\widehat{\mathbf{V}}$  defined in (24). We construct the symmetric matrix  $\widehat{\Gamma}^{[m]} \in \mathbb{R}^{L \times L}$  by specifying its lower triangular part as

$$\widehat{\Gamma}_{l,k}^{[m]} = \widehat{\Gamma}_{l,k}^\mathbb{Q} - S_{\pi(l,k)}^{[m]} \quad \text{for} \quad 1 \leq k \leq l \leq L, \quad (18)$$

where the mapping  $\pi$  is defined in (2). The above construction guarantees that  $\text{vecl}(\widehat{\Gamma}^{[m]} - \widehat{\Gamma}^{\mathbb{Q}})$  follows the distribution  $\mathcal{N}(\mathbf{0}, \widehat{\mathbf{V}} + d_0/n \cdot \mathbf{I})$ .

**Step 2: Aggregation.** We construct the CI by aggregating intervals constructed with the sampled weight vectors  $\{\widehat{\gamma}^{[m]}\}_{1 \leq m \leq M}$ . For each  $\widehat{\gamma}^{[m]}$ , we first treat it as being fixed and construct an interval for  $x_{\text{new}}^\top \beta^*$  by only quantifying the uncertainty due to the debiased estimators  $\{\widehat{x_{\text{new}}^\top b^{(l)}}\}_{1 \leq l \leq L}$  in (12). We refer to such an interval as the sampled interval. We take a union of a collection of sampled intervals as our constructed CI.

In the following, we state the details for the aggregation step. For  $1 \leq m \leq M$ , we compute  $\widehat{x_{\text{new}}^\top \beta}^{[m]} = \sum_{l=1}^L \widehat{\gamma}_l^{[m]} \cdot \widehat{x_{\text{new}}^\top b^{(l)}}$ , and construct the sampled interval as,

$$\text{Int}_\alpha^{[m]}(x_{\text{new}}) = \left( \widehat{x_{\text{new}}^\top \beta}^{[m]} - (1 + \eta_0) z_{\alpha/2} \widehat{\text{se}}^{[m]}(x_{\text{new}}), \widehat{x_{\text{new}}^\top \beta}^{[m]} + (1 + \eta_0) z_{\alpha/2} \widehat{\text{se}}^{[m]}(x_{\text{new}}) \right), \quad (19)$$

where  $\eta_0$  is any positive constant (with default value 0.01),  $z_{\alpha/2}$  is the upper  $\alpha/2$  quantile of the standard normal distribution, and

$$\widehat{\text{se}}^{[m]}(x_{\text{new}}) = \sqrt{\sum_{l=1}^L (\widehat{\sigma}_l^2/n_l^2) \cdot [\widehat{\gamma}_l^{[m]}]^2 [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}} \quad \text{with } \widehat{v}^{(l)} \text{ defined in (13)}.$$

The positive constant  $\eta_0$  in (19) is used to offset the finite-sample bias in high dimensions.

To aggregate the sampled intervals in (19), we introduce an index set  $\mathbb{M}$  to exclude the samples near the boundary of the sampling distribution. Define

$$\mathbb{M} = \left\{ 1 \leq m \leq M : \max_{1 \leq \pi(l,k) \leq L(L+1)/2} \left| \frac{S_{\pi(l,k)}^{[m]}}{\sqrt{\widehat{\mathbf{V}}_{\pi(l,k), \pi(l,k)} + d_0/n}} \right| \leq 1.1 \cdot z_{\alpha_0/[L(L+1)]} \right\}, \quad (20)$$

where  $z_{\alpha_0/[L(L+1)]}$  is the upper  $\alpha_0/[L(L+1)]$  quantile of the standard normal distribution (default value  $\alpha_0 = 0.01$ ). The index set  $\mathbb{M}$  excludes the generated  $S^{[m]}$  if any entry of  $S^{[m]}$  is above certain threshold. We construct the CI for  $x_{\text{new}}^\top \beta^*$  by aggregating the sampled intervals with the index  $m \in \mathbb{M}$ ,

$$\text{CI}_\alpha(x_{\text{new}}^\top \beta^*) = \cup_{m \in \mathbb{M}} \text{Int}_\alpha^{[m]}(x_{\text{new}}), \quad (21)$$

where  $\mathbb{M}$  is defined in (20) and  $\text{Int}_\alpha^{[m]}(x_{\text{new}})$  is defined in (19). Regarding the null hypothesis  $H_0 : x_{\text{new}}^\top \beta^* = 0$ , we propose the test  $\phi_\alpha = \mathbf{1}(0 \notin \text{CI}_\alpha(x_{\text{new}}^\top \beta^*))$ . If  $x_{\text{new}}$  is taken as the  $j$ -th Euclidean basis, we are testing the maximin significance of the  $j$ -th variable.

Our proposed DenseNet sampling is also adequate for the low-dimensional setting; see the extension in Section B.2 in the supplement.

As illustrated in Figure 1, the red interval represents  $\text{CI}_\alpha(x_{\text{new}}^\top \beta^*)$ , which takes a union of sampled intervals with indexes belonging to  $\mathbb{M}$ . Not all of  $\{\text{Int}_\alpha^{[m]}(x_{\text{new}})\}_{m \in \mathbb{M}}$  cover  $x_{\text{new}}^\top \beta^*$  since the uncertainty of  $\hat{\gamma}^{[m]}$  is not quantified in the construction of  $\text{Int}_\alpha^{[m]}(x_{\text{new}})$ .

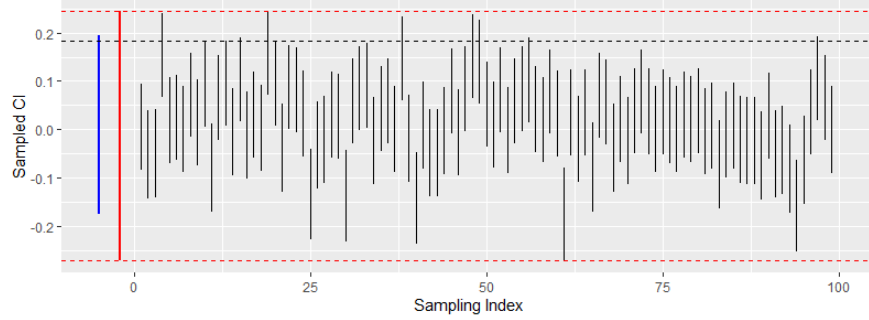


Figure 1: Illustration of  $\text{CI}_\alpha(x_{\text{new}}^\top \beta^*)$  with  $M = 100$  (in red) for setting 2 in Section H.1 in the supplement. The intervals in black denote  $\text{Int}_\alpha^{[m]}(x_{\text{new}})$  for  $m \in \mathbb{M}$ . The interval in blue is the normality CI in (34) with an oracle standard error computed by 500 simulations. The horizontal black dashed line represents the value of  $x_{\text{new}}^\top \beta^*$ .

We provide intuition for the DenseNet sampling method. By (18), we write

$$\text{vecl}(\hat{\Gamma}^{[m]}) \stackrel{d}{\approx} \text{vecl}(\Gamma^{\mathbb{Q}}) + \text{vecl}(\hat{\Gamma}^{\mathbb{Q}} - \Gamma^{\mathbb{Q}}) - S^{[m]} \quad \text{for } 1 \leq m \leq M.$$

For a large  $M$ , there exists  $m^* \in \mathbb{M}$  such that  $S^{[m^*]} \approx \text{vecl}(\hat{\Gamma}^{\mathbb{Q}} - \Gamma^{\mathbb{Q}})$  and then  $\hat{\Gamma}^{[m^*]}$  is almost the same as  $\Gamma^{\mathbb{Q}}$ ; see Theorem 1 for the exact argument. Consequently, the estimation error of  $\hat{\gamma}^{[m^*]}$  is smaller than the parametric rate and the sampled interval  $\text{Int}_\alpha^{[m^*]}(x_{\text{new}})$  in (19) is nearly a  $1 - \alpha$  confidence interval. The aggregation step in (21) is needed since we only know such  $m^*$  exists with a high probability but cannot locate the exact  $m^*$ .

We provide a few important remarks on our proposed sampling method.

**Remark 4** (Construction Details.). In (17), we generate  $S^{[m]}$  with the covariance matrix  $\hat{\mathbf{V}} + d_0/n \cdot \mathbf{I}$ , which is slightly larger than  $\hat{\mathbf{V}}$ . This enlargement ensures that  $\hat{\mathbf{V}} + d_0/n \cdot \mathbf{I}$  is positive definite even for a nearly singular  $\hat{\mathbf{V}}$  and the square-root of the enlarged variance dominates the bias of  $\hat{\Gamma}_{l,k}^{\mathbb{Q}}$ . In addition, we may replace the construction of  $\mathbb{M}$  in (20) by

$$\mathbb{M}' = \left\{ 1 \leq m \leq M : \|(\hat{\mathbf{V}} + d_0/n \cdot \mathbf{I})^{-1/2} S^{[m]}\|_2^2 \leq 1.1 \cdot \chi_{L(L+1)/2, \alpha_0}^2 \right\}, \quad (22)$$

where  $\chi_{L(L+1)/2, \alpha_0}^2$  is the upper  $\alpha_0$  (default value 0.01) quantile of the  $\chi^2$  distribution with  $L(L+1)/2$  degrees of freedom. We then modify (21) by replacing  $\mathbb{M}$  with  $\mathbb{M}'$ . In Section H.4 in the supplement, we compare different ways of constructing the index sets.

**Remark 5** (Computational Efficiency). The proposed method is computationally efficient even though we generate many samples  $\{S^{[m]}\}_{1 \leq m \leq M}$ . For each  $S^{[m]}$ , we solve the  $L$ -dimensional optimization problem in (16), instead of solving a  $p$ -dimensional optimization problem. When the group number  $L$  is much smaller than  $p$ , this significantly reduces the computation cost compared to non-parametric bootstrap, which directly samples the data and requires the implementation of high-dimensional optimization for each sampled data.

**Remark 6** (Non-regular Inference Problems). The DenseNet sampling relies on a different idea from bootstrap and subsampling, which construct the CIs by computing quantiles of the sampled estimates. Our proposed DenseNet sampling takes a union of sampled intervals, where each interval quantifies the uncertainty of  $\{\widehat{x_{\text{new}}^\top b^{(l)}}\}_{1 \leq l \leq L}$  and the union step accounts for the uncertainty of  $\widehat{\gamma}$ . In Section B in the supplement, we show that neither subsampling nor bootstrap provides valid inference in the presence of non-regularity or instability. The proposed DenseNet sampling helps address other non-regular problems.

### 3.3 Estimation of $\Gamma^\mathbb{Q}$ in High Dimensions

Now we present a debiased estimator of the matrix  $\Gamma^\mathbb{Q}$  by generalizing the high-dimensional inference methods for quadratic forms (Verzelen and Gassiat, 2018; Cai and Guo, 2020; Guo et al., 2019) and inner products (Guo et al., 2019). For  $1 \leq l \leq L$ , we randomly split the data  $(X^{(l)}, Y^{(l)})$  into two approximate equal-size subsamples  $(X_{A_l}^{(l)}, Y_{A_l}^{(l)})$  and  $(X_{B_l}^{(l)}, Y_{B_l}^{(l)})$ , where the index sets  $A_l$  and  $B_l$  satisfy  $A_l \cap B_l = \emptyset$ ,  $|A_l| = \lfloor n_l/2 \rfloor$  and  $|B_l| = n_l - \lfloor n_l/2 \rfloor$ . We randomly split the data  $\{X_i^\mathbb{Q}\}_{1 \leq i \leq N_\mathbb{Q}}$  into  $X_A^\mathbb{Q}$  and  $X_B^\mathbb{Q}$ , where the index sets  $A$  and  $B$  satisfy  $A \cap B = \emptyset$ ,  $|A| = \lfloor N_\mathbb{Q}/2 \rfloor$  and  $|B| = N_\mathbb{Q} - \lfloor N_\mathbb{Q}/2 \rfloor$ .

For  $1 \leq l \leq L$ , we construct the Lasso estimator  $\widehat{b}_{init}^{(l)}$  (Tibshirani, 1996) using the subsample  $(Y_{A_l}^{(l)}, X_{A_l}^{(l)})$ ; see its exact definition in (43) in the supplement. We define  $\widehat{\Sigma}^\mathbb{Q} = \frac{1}{|B|} \sum_{i \in B} X_i^\mathbb{Q} (X_i^\mathbb{Q})^\top$  and construct the plug-in estimator  $[\widehat{b}_{init}^{(l)}]^\top \widehat{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)}$ . We estimate  $\Gamma_{l,k}^\mathbb{Q}$  by correcting the bias of this plug-in estimator,

$$\widehat{\Gamma}_{l,k}^\mathbb{Q} = (\widehat{b}_{init}^{(l)})^\top \widehat{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)} + [\widehat{u}^{(l,k)}]^\top \frac{1}{|B_l|} [X_{B_l}^{(l)}]^\top (Y_{B_l}^{(l)} - X_{B_l}^{(l)} \widehat{b}_{init}^{(l)}) + [\widehat{u}^{(k,l)}]^\top \frac{1}{|B_k|} [X_{B_k}^{(k)}]^\top (Y_{B_k}^{(k)} - X_{B_k}^{(k)} \widehat{b}_{init}^{(k)}) \quad (23)$$



where  $\hat{u}^{(l,k)} \in \mathbb{R}^p$  and  $\hat{u}^{(k,l)} \in \mathbb{R}^p$  are the projection directions constructed in equations (44), (45) and (46) in the supplement. More details about constructing  $\hat{\Gamma}_{l,k}^{\mathbb{Q}}$  in (23) are postponed to Section A.7 in the supplement.

For  $1 \leq l_1 \leq k_1 \leq L$  and  $1 \leq l_2 \leq k_2 \leq L$ , we estimate the covariance between  $\hat{\Gamma}_{l_1,k_1}^{\mathbb{Q}} - \Gamma_{l_1,k_1}^{\mathbb{Q}}$  and  $\hat{\Gamma}_{l_2,k_2}^{\mathbb{Q}} - \Gamma_{l_2,k_2}^{\mathbb{Q}}$  by

$$\begin{aligned} \hat{\mathbf{V}}_{\pi(l_1,k_1),\pi(l_2,k_2)} &:= \frac{\hat{\sigma}_{l_1}^2}{|B_{l_1}|} (\hat{u}^{(l_1,k_1)})^\top \hat{\Sigma}^{(l_1)} [\hat{u}^{(l_2,k_2)} \mathbf{1}(l_2 = l_1) + \hat{u}^{(k_2,l_2)} \mathbf{1}(k_2 = l_1)] + \\ &\frac{\hat{\sigma}_{k_1}^2}{|B_{k_1}|} (\hat{u}^{(k_1,l_1)})^\top \hat{\Sigma}^{(k_1)} [\hat{u}^{(l_2,k_2)} \mathbf{1}(l_2 = k_1) + \hat{u}^{(k_2,l_2)} \mathbf{1}(k_2 = k_1)] + \\ &\frac{1}{|B|N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left( (\hat{b}_{init}^{(l_1)})^\top X_i^{\mathbb{Q}} (\hat{b}_{init}^{(k_1)})^\top X_i^{\mathbb{Q}} (\hat{b}_{init}^{(l_2)})^\top X_i^{\mathbb{Q}} (\hat{b}_{init}^{(k_2)})^\top X_i^{\mathbb{Q}} - (\hat{b}_{init}^{(l_1)})^\top \bar{\Sigma}^{\mathbb{Q}} \hat{b}_{init}^{(k_1)} (\hat{b}_{init}^{(l_2)})^\top \bar{\Sigma}^{\mathbb{Q}} \hat{b}_{init}^{(k_2)} \right), \end{aligned} \quad (24)$$

where  $\bar{\Sigma}^{\mathbb{Q}} = \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} (X_i^{\mathbb{Q}})^\top$ ,  $\hat{\Sigma}^{(l)} = \frac{1}{|B_l|} \sum_{i \in B_l} X_i^{(l)} [X_i^{(l)}]^\top$ , and  $\hat{\sigma}_l^2 = \|Y^{(l)} - X^{(l)} \hat{b}^{(l)}\|_2^2 / n_l$  for  $1 \leq l \leq L$ .

**Remark 7** (Sample splitting). Sample splitting is introduced for estimating  $\Gamma^{\mathbb{Q}}$  and the main reason is to create certain independence structure between the random errors  $\{\epsilon_{B_l}^{(l)}, \epsilon_{B_k}^{(k)}, \hat{\Sigma}^{\mathbb{Q}} - \Gamma^{\mathbb{Q}}\}$ , and the projection directions  $\hat{u}^{(k,l)}, \hat{u}^{(l,k)}$ , which are constructed by solving an optimization problem using the data  $X^{(l)}, Y_{A_l}^{(l)}, X^{(k)}, Y_{A_k}^{(k)}$  and  $X_A^{\mathbb{Q}}$ . We believe that the sample splitting is only needed for technical analysis. In simulations, we observe that the procedure without sample splitting performs well and improves the efficiency in comparison to that with sample splitting; see Table S6 in the supplement for details. We shall also remark that sample splitting is not needed for constructing  $\widehat{x_{\text{new}}^\top b^{(l)}}$  in (12). Furthermore, when there is no covariate shift, we construct the estimator of  $\Gamma^{\mathbb{Q}}$  without sample splitting; see more details in Section A.8 in the supplement.

### 3.4 Algorithm: Construction of $\text{CI}_\alpha(x_{\text{new}}^\top \beta^*)$

Algorithm 1 summarizes our proposed CI with the DenseNet sampling method. Regarding the tuning parameter selection, the Lasso estimators  $\{\hat{b}^{(l)}\}_{1 \leq l \leq L}$  are implemented by the R-package `glmnet` (Friedman et al., 2010) with tuning parameters  $\{\lambda_l\}_{1 \leq l \leq L}$  chosen by cross validation; the estimator  $\widehat{x_{\text{new}}^\top b^{(l)}}$  in (12) is implemented using the R-package `SIHR` (Rakshit et al., 2021) with the built-in selection of the tuning parameters  $\eta_l$  and  $\tau_l$ ; the tuning parameter selection for constructing  $\hat{\Gamma}_{l,k}^{\mathbb{Q}}$  in (23) is presented in (49) in the supplement. The code with the built-in tuning parameter selection is submitted together with the current paper.

---

**Algorithm 1** DenseNet Sampling Methods
 

---

**Input:**  $\{X^{(l)}, Y^{(l)}\}_{1 \leq l \leq L}$ ,  $X^{\mathbb{Q}}$ ; loading  $x_{\text{new}} \in \mathbb{R}^p$ ;  $\alpha \in (0, 1/2)$ ; sampling size  $M$ ;  $\tau_0 > 0$ .

**Output:** Confidence interval  $\text{CI}_{\alpha}(x_{\text{new}}^{\top} \beta^*)$

```

1: for  $l \leftarrow 1$  to  $L$  do
2:   Compute  $x_{\text{new}}^{\top} b^{(l)}$  in (12);
3: end for ▷ Estimation of  $\{x_{\text{new}}^{\top} b^{(l)}\}_{1 \leq l \leq L}$ 

4: for  $(l, k) \leftarrow \mathcal{I}_L = \{(l, k) : 1 \leq k \leq l \leq L\}$  do
5:   Compute  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}}$  in (23);
6: end for ▷ Estimation of  $\Gamma^{\mathbb{Q}}$ 

7: for  $(l_1, k_1) \leftarrow \mathcal{I}_L$ , and  $(l_2, k_2) \leftarrow \mathcal{I}_L$  do
8:   Compute  $\widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}$  in (24);
9: end for ▷ Uncertainty quantification of  $\widehat{\Gamma}^{\mathbb{Q}}$ 

10: for  $m \leftarrow 1, 2, \dots, M$  do
11:   Sample  $\widehat{\Gamma}^{[m]}$  in (17) and (18) with  $\widehat{\Gamma}^{\mathbb{Q}}$ ,  $\widehat{\mathbf{V}}$  and  $\tau_0 > 0$ ;
12:   Compute  $\widehat{\gamma}^{[m]}$  in (16);
13:   Construct  $\text{Int}_{\alpha}^{[m]}(x_{\text{new}})$  in (19) with  $\widehat{\gamma}^{[m]}$  and  $\widehat{x_{\text{new}}^{\top} b^{(l)}}$ ;
14: end for ▷ Sampling

15: Construct the index set  $\mathbb{M}$  in (20);
16: Construct  $\text{CI}_{\alpha}(x_{\text{new}}^{\top} \beta^*)$  in (21). ▷ Aggregation

```

---

## 4 Theoretical Justification

Before presenting the main theorems, we introduce the assumptions for the model (1). Define  $s = \max_{1 \leq l \leq L} \|b^{(l)}\|_0$  and  $n = \min_{1 \leq l \leq L} n_l$ .

(A1) For  $1 \leq l \leq L$ ,  $\{X_i^{(l)}, Y_i^{(l)}\}_{1 \leq i \leq n_l}$  are i.i.d. random variables, where  $X_i^{(l)} \in \mathbb{R}^p$  is sub-gaussian with  $\Sigma^{(l)} = \mathbf{E}X_i^{(l)}[X_i^{(l)}]^{\top}$  satisfying  $c_0 \leq \lambda_{\min}(\Sigma^{(l)}) \leq \lambda_{\max}(\Sigma^{(l)}) \leq C_0$  for positive constants  $C_0 > c_0 > 0$ ; the error  $\epsilon_i^{(l)}$  is sub-gaussian with  $\mathbf{E}(\epsilon_i^{(l)} | X_i^{(l)}) = 0$ ,  $\mathbf{E}([\epsilon_i^{(l)}]^2 | X_i^{(l)}) = \sigma_l^2$ , and  $\mathbf{E}([\epsilon_i^{(l)}]^{2+c} | X_i^{(l)}) \leq C$  for some positive constants  $c > 0$  and  $C > 0$ .  $\{X_i^{\mathbb{Q}}\}_{1 \leq i \leq N_{\mathbb{Q}}}$  are i.i.d. sub-gaussian with  $\Sigma^{\mathbb{Q}} = \mathbf{E}X_i^{\mathbb{Q}}[X_i^{\mathbb{Q}}]^{\top}$  satisfying  $c_1 \leq \lambda_{\min}(\Sigma^{\mathbb{Q}}) \leq \lambda_{\max}(\Sigma^{\mathbb{Q}}) \leq C_1$  for positive constants  $C_1 > c_1 > 0$ .

(A2)  $L$  is finite,  $\max_{1 \leq l \leq L} \|b^{(l)}\|_2 \leq C$  for a positive constant  $C > 0$ ,  $n \asymp \max_{1 \leq l \leq L} n_l$ , and the model complexity parameters  $(s, n, p, N_{\mathbb{Q}})$  satisfy  $n \gg (s \log p)^2$  and  $N_{\mathbb{Q}} \gg n^{3/4} [\log \max\{N_{\mathbb{Q}}, p\}]^2$ .

We always consider asymptotic expressions in the limit where both  $n, p \rightarrow \infty$ . Assumption (A1) is commonly assumed for the theoretical analysis of high-dimensional linear models;

c.f. [Bühlmann and van de Geer \(2011\)](#). The positive definite  $\Sigma^{(l)}$  and the sub-gaussianity of  $X_i^{(l)}$  guarantee the restricted eigenvalue condition with a high probability ([Bickel et al., 2009](#); [Zhou, 2009](#)). The sub-gaussian errors are generally required for the theoretical analysis of the Lasso estimator in high dimensions ([Bickel et al., 2009](#); [Bühlmann and van de Geer, 2011](#), e.g.). The moment conditions on  $\epsilon_i^{(l)}$  are needed to establish the asymptotic normality of the debiased estimators of single regression coefficients ([Javanmard and Montanari, 2014](#)). Similarly, they are imposed here to establish the asymptotic normality of  $\widehat{\Gamma}^{\mathbb{Q}}$  and  $\widehat{x_{\text{new}}^{\top} b^{(l)}}$ . The model complexity condition  $n \gg (s \log p)^2$  in (A2) is assumed in the CI construction for high-dimensional linear models ([Zhang and Zhang, 2014](#); [van de Geer et al., 2014](#); [Javanmard and Montanari, 2014](#)) and has been shown in [Cai and Guo \(2017\)](#) as the minimum sample size for constructing adaptive CIs. The condition on  $N_{\mathbb{Q}}$  is mild as there is typically a large amount of unlabelled data for the target population. The boundedness assumptions on  $L$  and  $\|b^{(l)}\|_2$  are mainly imposed to simplify the presentation, so is the assumption  $n \asymp \max_{1 \leq l \leq L} n_l$ .

#### 4.1 The Sampling Property

We justify the sampling step proposed in Section 3.2 and show that there exists at least one sampled weight vector converging to  $\gamma^*$  at a rate faster than  $1/\sqrt{n}$ . To characterize the sample accuracy, we introduce the sampling error ratio  $\text{err}_n(M)$  as a function of the sampling size  $M$ :

$$\text{err}_n(M) = \left[ \frac{4 \log n}{C^*(L, \alpha_0) \cdot M} \right]^{\frac{2}{L(L+1)}}, \quad (25)$$

where  $L$  is the total number of groups,  $\alpha_0 \in (0, 0.01]$  is the pre-specified constant used in the construction of  $\mathbb{M}$  in (20), and  $C^*(L, \alpha_0)$  is a constant defined in (71) in the supplement.

The following theorem establishes the rate of convergence for  $\min_{m \in \mathbb{M}} \|\widehat{\gamma}^{[m]} - \gamma^*\|_2$ , which represents the best approximation accuracy among all sampled vectors  $\{\widehat{\gamma}^{[m]}\}_{m \in \mathbb{M}}$ .

**Theorem 1.** *Consider the model (1). Suppose Conditions (A1) and (A2) hold. If  $\text{err}_n(M)$  defined in (25) satisfies  $\text{err}_n(M) \ll \min\{1, c^*(\alpha_0), \lambda_{\min}(\Gamma^{\mathbb{Q}})\}$  where  $c^*(\alpha_0)$  is a constant defined in (71) in the supplement, then*

$$\liminf_{n, p \rightarrow \infty} \mathbf{P} \left( \min_{m \in \mathbb{M}} \|\widehat{\gamma}^{[m]} - \gamma^*\|_2 \leq \frac{\sqrt{2} \text{err}_n(M)}{\lambda_{\min}(\Gamma^{\mathbb{Q}})} \cdot \frac{1}{\sqrt{n}} \right) \geq 1 - \alpha_0,$$

where  $\alpha_0 \in (0, 0.01]$  is the pre-specified constant used in the construction of  $\mathbb{M}$  in (20).

Let  $m^*$  denote one index belonging to  $\mathbb{M}$  such that  $\|\widehat{\gamma}^{[m^*]} - \gamma^*\|_2 = \min_{m \in \mathbb{M}} \|\widehat{\gamma}^{[m]} - \gamma^*\|_2$ . Since  $\text{err}_n(M) \ll \lambda_{\min}(\Gamma^{\mathbb{Q}})$ , Theorem 1 states that  $\|\widehat{\gamma}^{[m^*]} - \gamma^*\|_2$  converges to zero at a faster

rate than  $1/\sqrt{n}$ . Corollary 1 in the supplement shows that  $c^*(\alpha_0)$  and  $C^*(L, \alpha_0)$  are at least of constant orders under regularity conditions. Together with  $\lambda_{\min}(\Gamma^{\mathbb{Q}}) \geq c$  for some positive constant  $c > 0$ , any choice of  $M \gg \log n$  guarantees the condition  $\text{err}_n(M) \ll \min\{1, c^*(\alpha_0), \lambda_{\min}(\Gamma^{\mathbb{Q}})\}$ . In practice, we set  $M = 500$  as the default value and observe reliable inference results.

Theorem 1 covers the important setting with a nearly singular  $\Gamma^{\mathbb{Q}}$ , that is,  $\lambda_{\min}(\Gamma^{\mathbb{Q}}) > 0$  for any given  $p$  but  $\liminf_{p \rightarrow \infty} \lambda_{\min}(\Gamma^{\mathbb{Q}}) = 0$ . This setting will appear if some of  $\{b^{(l)}\}_{1 \leq l \leq L}$  are similar to each other but not exactly the same. Theorem 1 can be applied to this nearly singular setting if we choose a sufficiently large sampling number  $M > 0$  such that  $\text{err}_n(M) \ll \min\{1, c^*(\alpha_0), \lambda_{\min}(\Gamma^{\mathbb{Q}})\}$ . Theorem 1 is not applied to the exactly singular setting  $\lambda_{\min}(\Gamma^{\mathbb{Q}}) = 0$  since the condition  $\lambda_{\min}(\Gamma^{\mathbb{Q}}) \gg \text{err}_n(M)$  is violated. The ridge-type maximin effect introduced in Section 5 is helpful for the exactly singular setting.

## 4.2 Statistical Inference for Maximin Effects

The following decomposition reveals why the DenseNet sampling method works,

$$\widehat{x_{\text{new}}^{\top} \beta}^{[m]} - x_{\text{new}}^{\top} \beta = \sum_{l=1}^L (\widehat{\gamma}_l^{[m]} - \gamma_l^*) \cdot \widehat{x_{\text{new}}^{\top} b^{(l)}} + \sum_{l=1}^L \gamma_l^* \cdot (\widehat{x_{\text{new}}^{\top} b^{(l)}} - x_{\text{new}}^{\top} b^{(l)}). \quad (26)$$

With  $m = m^*$  defined after Theorem 1, the uncertainty of  $\sum_{l=1}^L (\widehat{\gamma}_l^{[m]} - \gamma_l^*) \cdot \widehat{x_{\text{new}}^{\top} b^{(l)}}$  is negligible and we just need to quantify the uncertainty of  $\sum_{l=1}^L \gamma_l^* \cdot (\widehat{x_{\text{new}}^{\top} b^{(l)}} - x_{\text{new}}^{\top} b^{(l)})$ . The following theorem establishes the properties of  $\text{CI}_{\alpha}(x_{\text{new}}^{\top} \beta^*)$  defined in (21).

**Theorem 2.** *Suppose that the conditions of Theorem 1 hold. Then for any positive constant  $\eta_0 > 0$ , the confidence interval  $\text{CI}_{\alpha}(x_{\text{new}}^{\top} \beta^*)$  defined in (21) satisfies*

$$\lim_{n, p \rightarrow \infty} \mathbf{P}(x_{\text{new}}^{\top} \beta^* \in \text{CI}_{\alpha}(x_{\text{new}}^{\top} \beta^*)) \geq 1 - \alpha - \alpha_0, \quad (27)$$

where  $\alpha \in (0, 1/2)$  is the pre-specified significance level and  $\alpha_0 \in (0, 0.01]$  is the pre-specified constant used in the construction of  $\mathbb{M}$  in (20). By further assuming  $N_{\mathbb{Q}} \gtrsim \max\{n, p\}$  and  $\lambda_{\min}(\Gamma^{\mathbb{Q}}) \gtrsim \sqrt{\log p / \min\{n, N_{\mathbb{Q}}\}}$ , then there exists some positive constant  $C > 0$  such that

$$\liminf_{n, p \rightarrow \infty} \mathbf{P}\left(\mathbf{Leng}(\text{CI}_{\alpha}(x_{\text{new}}^{\top} \beta^*)) \leq C \max\left\{1, \frac{z_{\alpha_0/[L(L+1)]}}{\lambda_{\min}(\Gamma^{\mathbb{Q}})}\right\} \cdot \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}\right) = 1, \quad (28)$$

where  $\mathbf{Leng}(\text{CI}(x_{\text{new}}^{\top} \beta^*))$  denotes the interval length and  $z_{\alpha_0/[L(L+1)]}$  is the upper  $\alpha_0/[L(L+1)]$  quantile of the standard normal distribution.

A few remarks are in order for Theorem 2. Firstly, the coverage property uniformly holds for any  $x_{\text{new}} \in \mathbb{R}^p$ , and the validity of the constructed CI does not require the asymptotic normality of  $\widehat{x_{\text{new}}^\top \beta^*}$ . Due to the non-regularity and instability detailed in Section 2.3, the asymptotic distribution of the maximin effect estimator  $\widehat{x_{\text{new}}^\top \beta^*}$  is not necessarily normal. In Sections 6.1 and B.2 in the supplement, we illustrate that the CIs assuming asymptotic normality are under-coverage in both low and high dimensions.

Secondly, in (27), the coverage is only guaranteed to be above  $1 - \alpha - \alpha_0$ , instead of converging to  $1 - \alpha$ . This one-sided coverage guarantee comes with our proposed sampling method. For  $m^*$  defined after Theorem 1 and  $\text{Int}_\alpha^{[m^*]}(x_{\text{new}})$  defined in (19) with  $\text{err}_n(M)/\lambda_{\min}(\Gamma^\mathbb{Q}) \lesssim \eta_0 \ll 1$ , we have a more precise coverage statement,

$$\lim_{n,p \rightarrow \infty} \mathbf{P}(x_{\text{new}}^\top \beta^* \in \text{Int}_\alpha^{[m^*]}(x_{\text{new}})) \in (1 - \alpha - \alpha_0, 1 - \alpha + \alpha_0). \quad (29)$$

See its proof in Section E.3 in the supplement. However, since we cannot locate  $m^*$ , we take a union over  $\mathbb{M}$  to guarantee the coverage property, but this will lead to the one-sided bound for the coverage in (27). We examine the tightness of the coverage inequality (27) throughout the simulation studies; see Table 3 in Section 6 for a summary.

Thirdly, if  $\lambda_{\min}(\Gamma^\mathbb{Q}) \geq c$  for a positive constant  $c > 0$ , then the CI length is of the rate  $\|x_{\text{new}}\|_2/\sqrt{n}$ . In consideration of a single high-dimensional linear model, Cai et al. (2021) showed that, without the knowledge of the sparsity level  $s$  of  $b^{(l)}$ , the optimal length of CIs for  $x_{\text{new}}^\top b^{(l)}$  is  $\|x_{\text{new}}\|_2/\sqrt{n}$  if  $\sqrt{\|x_{\text{new}}\|_0} \lesssim \sqrt{n}/\log p$  and  $s \lesssim \sqrt{n}/\log p$ ; see Corollary 4 in Cai et al. (2021) for the exact details. In Section 6, we evaluate the precision properties of our proposed CI in finite samples; see Table 3 in Section 6 for a summary.

## 5 Stability: Ridge-type Maximin Effect

We generalize the maximin effect  $\beta^*(\mathbb{Q})$  in (10) and propose the following ridge-type maximin effect for  $\delta \geq 0$ ,

$$\beta_\delta^*(\mathbb{Q}) = \sum_{l=1}^L [\gamma_\delta^*(\mathbb{Q})]_l \cdot b^{(l)} \quad \text{with} \quad \gamma_\delta^*(\mathbb{Q}) = \arg \min_{\gamma \in \Delta^L} [\gamma^\top \Gamma^\mathbb{Q} \gamma + \delta \|\gamma\|_2^2], \quad (30)$$

which adds the ridge penalty in constructing the weight vector. When there is no confusion, we omit the dependence on  $\mathbb{Q}$  and write  $\beta_\delta^*(\mathbb{Q})$  and  $\gamma_\delta^*(\mathbb{Q})$  as  $\beta_\delta^*$  and  $\gamma_\delta^*$ , respectively. The ridge-type maximin effect  $\beta_\delta^*(\mathbb{Q})$  generalizes  $\beta^*(\mathbb{Q})$  but  $\beta_\delta^*(\mathbb{Q})$  with a positive  $\delta > 0$  has a different interpretation from  $\beta^*(\mathbb{Q})$ . In Section A.3 in the supplement, we express  $\beta_\delta^*$  as the

solution to a distributionally robust optimization problem, where adding the ridge penalty is equivalent to perturbing the covariate distribution for the target population.

Section 2.3 pointed out that the maximin effect might not be a stable aggregation, especially when some  $\{b^{(l)}\}$  are similar to each other. The ridge-type maximin effect  $\beta_\delta^*$  is proposed to address this instability issue. When  $\beta^*$  is not unique,  $\beta_\delta^*$  is uniquely defined for any positive  $\delta > 0$ , which is implied by Lemma 3 in the supplement. For  $L = 2$ , the solution of (30) is  $\gamma_\delta^* = ([\gamma_\delta^*]_1, 1 - [\gamma_\delta^*]_1)^\top$ , with  $[\gamma_\delta^*]_1 = \min \left\{ \max \left\{ \frac{\Gamma_{22} + \delta - \Gamma_{12}}{\Gamma_{11} + \Gamma_{22} + 2\delta - 2\Gamma_{12}}, 0 \right\}, 1 \right\}$ . Even if  $\Gamma_{11} + \Gamma_{22} - 2\Gamma_{12}$  is near zero, the instability issue is addressed with  $\delta > 0$  since  $1/[\Gamma_{11} + \Gamma_{22} + 2\delta - 2\Gamma_{12}]$  can be accurately estimated.

The estimation and inference methods detailed in Sections 3.1 and 3.2 can be extended to dealing with  $x_{\text{new}}^\top \beta_\delta^*$  for any  $\delta \geq 0$ . Specifically, we generalize the point estimator (15) as

$$\widehat{x_{\text{new}}^\top \beta_\delta^*} = \sum_{l=1}^L [\widehat{\gamma}_\delta]_l \cdot \widehat{x_{\text{new}}^\top b^{(l)}} \quad \text{with} \quad \widehat{\gamma}_\delta := \arg \min_{\gamma \in \Delta^L} \left[ \gamma^\top \widehat{\Gamma}^Q \gamma + \delta \|\gamma\|_2^2 \right]. \quad (31)$$

Regarding the CI construction for  $x_{\text{new}}^\top \beta_\delta^*$ , we replace  $\widehat{\gamma}^{[m]}$  in (16) by

$$\widehat{\gamma}_\delta^{[m]} = \arg \min_{\gamma \in \Delta^L} \gamma^\top (\widehat{\Gamma}^{[m]} + \delta \cdot \mathbf{I})_+ \gamma \quad \text{for} \quad \delta \geq 0. \quad (32)$$

The detailed inference procedure is presented in Section E in the supplement.

For a general  $L \geq 2$ , we propose a stability measure (depending on  $\delta$ ) as

$$\mathbb{I}(\delta) = \frac{\sum_{m=1}^M \|\widehat{\gamma}_\delta^{[m]} - \widehat{\gamma}_\delta\|_2^2}{\sum_{m=1}^M \|\widehat{\Gamma}^{[m]} - \widehat{\Gamma}^Q\|_2^2}, \quad (33)$$

where  $\{\widehat{\Gamma}^{[m]}\}_{1 \leq m \leq M}$  and  $\{\widehat{\gamma}_\delta^{[m]}\}_{1 \leq m \leq M}$  are the resampled estimates defined in (18) and (32), respectively. A large value of  $\mathbb{I}(\delta)$  indicates that a small error in estimating  $\Gamma^Q$  may lead to a much larger error in estimating the weight vector. We claim the weight optimization problem to be stable if  $\mathbb{I}(\delta) < t_0$  (default value  $t_0 = 0.5$ ) and unstable otherwise.

	$\mathbb{I}(\delta)$					
Setting	$\delta = 0$	$\delta = 0.1$	$\delta = 0.5$	$\delta = 1$	$\delta = 2$	$\Gamma_{11}^Q + \Gamma_{22}^Q - 2\Gamma_{12}^Q$
1	5.464	1.966	0.264	0.072	0.019	0.026
4(a)	0.108	0.087	0.044	0.024	0.011	2.007
6(a)	3.305	1.449	0.221	0.065	0.018	0.160

Table 1: The instability measure  $\mathbb{I}(\delta)$  for simulation settings detailed in Section 6.2 and Section H.1 in the supplement. The reported values are averaged over 100 repeated simulations.

In Table 1, we report  $\mathbb{I}(\delta)$  for several simulation settings detailed in Section 6.2 and Section H.1 in the supplement. For  $\delta = 0$ , the weight optimization is unstable for settings 1 and 6(a) but stable for 4(a). A larger  $\delta$  value is effective in stabilizing the optimization problem.

The following proposition controls the worse-off amount  $R_{\mathbb{Q}}[\beta_{\delta}^*] - R_{\mathbb{Q}}[\beta^*]$ .

**Proposition 2.** *Suppose  $\lambda_L(\mathcal{B}) > 0$  with  $\mathcal{B} = (b^{(1)}, \dots, b^{(L)}) \in \mathbb{R}^{p \times L}$ , then the ridge-type minimizer  $\beta_{\delta}^*$  defined in (30) satisfies  $R_{\mathbb{Q}}[\beta_{\delta}^*] \geq R_{\mathbb{Q}}[\beta^*] - 2\delta(\|\gamma_{\delta}^*\|_{\infty} - \|\gamma_{\delta}^*\|_2^2)$ , where  $R_{\mathbb{Q}}[\cdot]$  and  $\beta^*$  are defined in (9) and  $\gamma_{\delta}^*$  is defined in (30).*

A data-dependent way of choosing  $\delta$  is presented in Section A.3 in the supplement, which balances the stability measure  $\mathbb{I}(\delta)$  and the reward ratio  $R_{\mathbb{Q}}[\beta_{\delta}^*]/R_{\mathbb{Q}}[\beta^*]$ .

## 6 Simulation Results

We generate  $\{X^{(l)}, Y^{(l)}\}_{1 \leq l \leq L}$  by (1), where, for the  $l$ -th group,  $\{X_i^{(l)}\}_{1 \leq i \leq n_l} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \Sigma^{(l)})$  and  $\{\epsilon_i^{(l)}\}_{1 \leq i \leq n_l} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \sigma_l^2)$ . For  $1 \leq l \leq L$ , we take  $n_l = n$ ,  $\sigma_l = 1$  and  $\Sigma^{(l)} = \Sigma$ , with  $\Sigma_{j,k} = 0.6^{|j-k|}$  for  $1 \leq j, k \leq p$ . In the covariate shift setting, we have access to  $\{X_i^{\mathbb{Q}}\}_{1 \leq i \leq N_{\mathbb{Q}}} \stackrel{\text{i.i.d.}}{\sim} N(\mathbf{0}, \Sigma^{\mathbb{Q}})$ .  $p$  is set as 500 and  $N_{\mathbb{Q}}$  is set as 2000 by default. We construct  $\hat{\Gamma}^{\mathbb{Q}}$  without the sample splitting in the simulation and real data analysis. A numerical comparison between our proposed methods with and without sample splitting is reported in Table S6 in the supplement. Throughout the simulation, we report the average measures over 500 simulations.

We compare our proposed CI with a normality CI of the form

$$(\widehat{x_{\text{new}}^{\top} \beta_{\delta}^*} - 1.96 \cdot \widehat{\text{SE}}, \widehat{x_{\text{new}}^{\top} \beta_{\delta}^*} + 1.96 \cdot \widehat{\text{SE}}), \quad (34)$$

where  $\widehat{x_{\text{new}}^{\top} \beta_{\delta}^*}$  is defined in (31) and  $\widehat{\text{SE}}$  denotes the sample standard deviation of  $\widehat{x_{\text{new}}^{\top} \beta_{\delta}^*}$  calculated based on 500 simulations. Since  $\widehat{\text{SE}}$  is calculated in an oracle way, this normality CI is not a practical procedure but a favorable implementation of the CI constructed by assuming the asymptotic normality of the point estimator  $\widehat{x_{\text{new}}^{\top} \beta_{\delta}^*}$ . Since  $\widehat{x_{\text{new}}^{\top} \beta_{\delta}^*}$  might have a non-standard limiting distribution, the normality CI in (34) might not be valid.

### 6.1 Simulations for Non-regularity and Instability Settings

We show that the normality CI in (34) suffers from under-coverage in the presence of non-regularity and instability. We generate ten simulation settings, where (I-1) to (I-6) corre-

spond to settings with both non-regularity and instability, (I-7) to (I-9) correspond to the non-regularity settings, and (I-10) corresponds to an easier setting without non-regularity and instability. The detailed settings are reported in Section B.1 in the supplement.

We focus on inference for the maximin effect without the ridge penalty. In Table 2, except for (I-10), the empirical coverages of the normality CI in (34) are between 70% and 85%. Our proposed CI achieves the desired coverage at the expense of a wider interval. The ratio of the average length of our proposed CI to the normality CI is between 1.35 and 2.02.

Setting	$\mathbb{I}(\delta)$	Coverage		Length		Length Ratio
		normality	Proposed	normality	Proposed	
(I-1)	3.368	0.700	0.960	0.352	0.597	1.693
(I-2)	3.707	0.818	0.978	0.320	0.543	1.699
(I-3)	3.182	0.748	0.970	0.352	0.588	1.673
(I-4)	1.732	0.770	0.956	0.520	0.796	1.532
(I-5)	1.857	0.796	0.978	0.445	0.710	1.594
(I-6)	1.987	0.710	0.980	0.480	0.832	1.732
(I-7)	0.029	0.848	0.985	0.250	0.507	2.028
(I-8)	0.031	0.758	0.981	0.262	0.530	2.020
(I-9)	0.010	0.830	0.988	0.690	1.264	1.832
(I-10)	0.030	0.940	0.988	0.232	0.315	1.354

Table 2: Coverage and length of our proposed CI in Algorithm 1 and the normality CI in (34) (with  $\delta = 0$ ). The column indexed with “Coverage” and “Length” represent the empirical coverage and average length for CIs, respectively; the columns indexed with “normality ” and “Proposed” represent the normality CI and our proposed CI, respectively. The column indexed with “Length Ratio” represents the ratio of the average length of our proposed CI to that of the normality CI. The column indexed with “ $\mathbb{I}(\delta)$ ” reports the instability measure defined in (33).

To further investigate the under-coverage of the normality CI, we plot in Figure 2 the histogram of  $\widehat{x_{\text{new}}^T \beta^*}$  and  $\widehat{\gamma}$  in (15) over 500 simulations. The leftmost panel of Figure 2 corresponds to the setting (I-1) with non-regularity and instability. Due to the instability, the histogram of the weight estimates has some concentrations near both 0 and 1, which results in the bias component of  $\widehat{x_{\text{new}}^T \beta^*}$  being comparable to its standard error. Consequently, the empirical coverage of the corresponding normality CI is only 70%. The middle panel of Figure 2 corresponds to the setting (I-8) with non-regularity, where the weight for the first group is left-censored at zero. This censoring at zero leads to the bias of  $\widehat{x_{\text{new}}^T \beta^*}$  being comparable to its standard error and under-coverage of the normality CI. The rightmost panel corresponds to the favorable setting (I-10) without non-regularity and instability. The weight distributions and the maximin effect estimator are nearly normal, and the corresponding normality CI in (34) achieves the 95% coverage level.



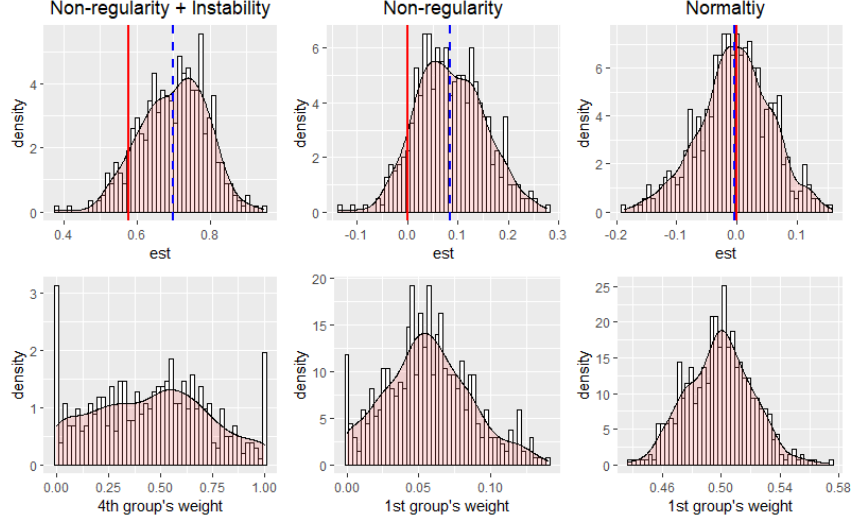


Figure 2: The histogram of the maximin estimator  $\widehat{x_{\text{new}}^\top \beta^*}$  (top) and one coordinate of the weight estimator (bottom) over 500 simulations. The figures from the leftmost to the rightmost correspond to settings (I-1), (I-8) and (I-10). The red solid line denotes the true value  $x_{\text{new}}^\top \beta^*$  while the blue dashed line denotes the sample average over 500 simulations.

## 6.2 Other Simulation Studies

We investigate our proposed method over seven additional settings. Settings 1, 3, and 4 are described in the following, and other settings are described in Section H.1 in the supplement.

*Setting 1* ( $L = 2$  with no covariate shift).  $b_j^{(1)} = j/40$  for  $1 \leq j \leq 10$ ,  $b_j^{(1)} = 1$  for  $j = 22, 23$ ,  $b_j^{(1)} = 0.1$  for  $j = 499, 500$ , and  $b_j^{(1)} = 0$  otherwise;  $b_j^{(2)} = b_j^{(1)}$  for  $1 \leq j \leq 499$  and  $b_{500}^{(2)} = 0.3$ ;  $[x_{\text{new}}]_j = 1$  for  $j = 500$  and  $[x_{\text{new}}]_j = 0$  otherwise.

*Setting 3* ( $L = 2$  with/without covariate shift).  $b^{(1)}$  and  $b^{(2)}$  are the same as setting 1, except for  $b_{498}^{(1)} = 0.5$ ,  $b_j^{(1)} = -0.5$  for  $j = 499, 500$ , and  $b_{500}^{(2)} = 1$ ;  $[x_{\text{new}}]_j = 1$  for  $j = 499, 500$ , and  $[x_{\text{new}}]_j = 0$  otherwise. Setting 3(a) is the covariate shift setting with  $\Sigma_{i,i}^{\mathbb{Q}} = 1.5$  for  $1 \leq i \leq 500$ ,  $\Sigma_{i,j}^{\mathbb{Q}} = 0.6$  for  $1 \leq i \neq j \leq 5$ ,  $\Sigma_{i,j}^{\mathbb{Q}} = -0.9$  for  $499 \leq i \neq j \leq 500$  and  $\Sigma_{i,j}^{\mathbb{Q}} = \Sigma_{i,j}$  otherwise; Setting 3(b) is the no covariate shift setting.

*Setting 4* (varying  $L$ ). Vary  $L$  across  $\{2, 5, 10\}$ , denoted as (4a), (4b) and (4c), respectively.  $b_j^{(1)} = j/40$  for  $1 \leq j \leq 10$ ,  $b_{498}^{(1)} = 0.5$ ,  $b_j^{(1)} = -0.5$  for  $j = 499, 500$ , and  $b_j^{(1)} = 0$  otherwise. For  $2 \leq l \leq L$ ,  $b_{10+l+j}^{(l)} = b_j^{(1)}$  for  $1 \leq j \leq 10$  and  $b_j^{(l)} = b_j^{(1)}/2^{l-1}$  for  $j = 498 \leq j \leq 500$ , and  $b_j^{(l)} = 0$  otherwise.  $[x_{\text{new}}]_j = 1$  for  $498 \leq j \leq 500$ , and  $[x_{\text{new}}]_j = 0$  otherwise;  $\Sigma_{i,i}^{\mathbb{Q}} = 1.5$  for  $1 \leq i \leq 500$ ,  $\Sigma_{i,j}^{\mathbb{Q}} = 0.9$  for  $1 \leq i \neq j \leq 5$  and  $499 \leq i \neq j \leq 500$ , and  $\Sigma_{i,j}^{\mathbb{Q}} = \Sigma_{i,j}$  otherwise.

We compute Coverage Error =  $|\text{Empirical Coverage} - 95\%|$ , with the empirical coverage computed based on 500 simulations. We report the ratio of the average length of our proposed

CI to that of the normality CI in (34). For each simulation setting, we average the coverage error and the length ratio over different combinations of  $\delta \in \{0, 0.1, 0.5, 1, 2\}$  and  $n \in \{200, 300, 500\}$ .<sup>1</sup>

In Table 3, we summarize the average coverage error and length ratio over different settings. Since our proposed CIs generally achieve 95% for  $n \geq 200$ , the coverage errors mainly result from over-coverage instead of under-coverage. For settings 3(a), 4(a), and 5, the empirical coverage of our proposed CI is nearly 95%, and the corresponding length ratios for settings 4(a) and 5 are near 1. For settings 3(b), 6, and 7, our proposed CIs are over-coverage, but the average length ratios are at most 1.864. More detailed results for these settings are presented in the following and Section H in the supplement.

Setting	1	2	3(a)	3(b)	4(a)	4(b)	4(c)	5	6	7
Coverage Error	2.60%	3.45%	0.95%	4.51%	1.64%	2.71%	3.57%	0.75%	4.25%	4.68%
Length Ratio	1.322	1.607	1.516	1.864	1.047	1.356	1.587	1.268	1.554	1.579

Table 3: Average coverage error and length ratio across different settings.

**Dependence on  $n$  and  $\delta$ .** For setting 1, we plot in Figure 3 the empirical coverage and CI length over  $\delta \in \{0, 0.1, 0.5, 1, 2\}$ . Our proposed CIs achieve the desired coverage level for  $n \geq 200$ . The CIs get shorter with increasing  $n$  or  $\delta$ : the lengths of CIs for  $\delta = 2$  are around half of those for  $\delta = 0$ . This shows that a positive  $\delta$  is effective in reducing the length of the CI in setting 1, which is unstable and has a large instability measure as shown in Table 1.

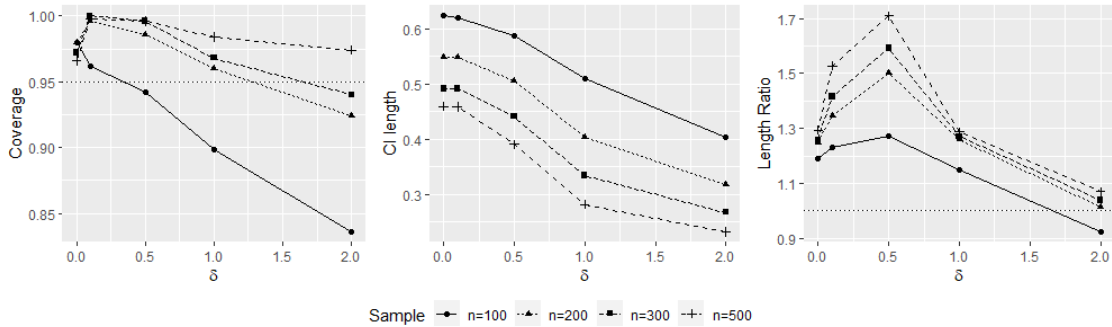
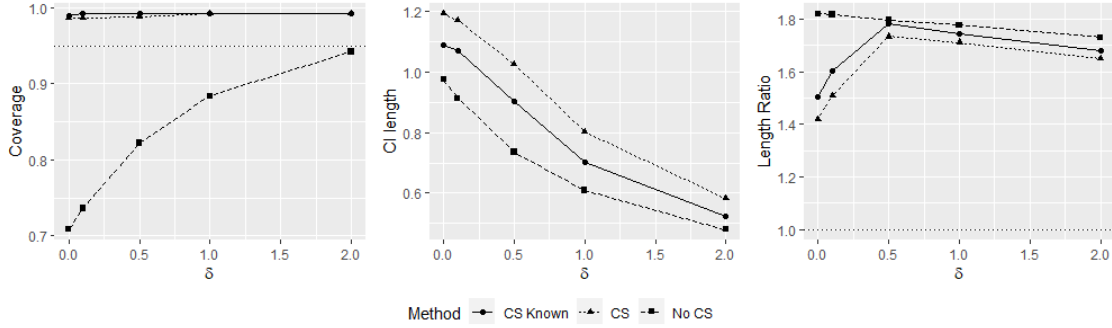


Figure 3: Dependence on  $\delta$  and  $n$  (setting 1). “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).

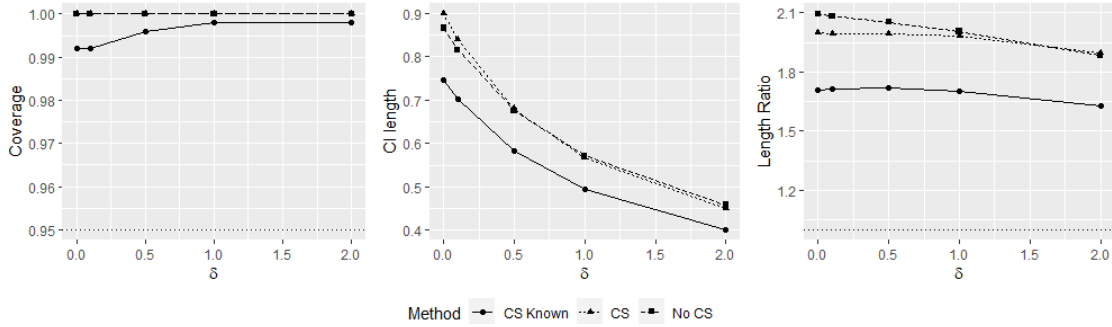
**Covariate shift.** We compare Algorithm 1 with and without covariate shift. For the covariate shift setting, if  $\Sigma^{\mathbb{Q}}$  is known, we present the modification of Algorithm 1 in Section A.7

<sup>1</sup>For setting 5, instead of averaging over  $\delta \in \{0, 0.1, 0.5, 1, 2\}$ , we take an average with respect to the perb parameter; see more details in Section H.1 in the supplement.

in the supplement. The top of Figure 4 corresponds to the simulation settings with covariate shift and  $n = 500$ . The no covariate shift algorithm does not achieve the 95% coverage due to the bias of assuming no covariate shift. In contrast, the covariate shift algorithms (with or without knowing  $\Sigma^Q$ ) achieve the 95% coverage level, and the CI constructed with known  $\Sigma^Q$  is shorter as it does not need to quantify the uncertainty of estimating  $\Sigma^Q$ . The bottom of Figure 4 corresponds to the setting with no covariate shift. All algorithms achieve the desired coverage level. The results for  $n = 200$  is reported in Figure S5 in the supplement.



(a) Setting 3(a) with covariate shift



(b) Setting 3(b) with no covariate shift

Figure 4: Comparison of covariate shift and no covariate shift algorithms ( $n = 500$ ). “CS Known”, “CS” and “No CS” represent Algorithm 1 with known  $\Sigma^Q$ , with covariate shift but unknown  $\Sigma^Q$ , and with no covariate shift, respectively. “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI.

**Varying  $L$  with covariate shift.** For setting 4(b) with  $L = 5$ , we plot in Figure 5 the empirical coverage and CI length over  $\delta \in \{0, 0.1, 0.5, 1, 2\}$ . The results are similar to those in Figure 3. The results for settings 4(a) and 4(b) are presented in Figure S6 in the supplement.

**Additional simulation studies.** In Section H in the supplement, we conduct additional simulation studies, including settings with opposite effects, a larger  $p$ , and different choices of the index set  $\mathbb{M}$ . We also investigate the effect of sample splitting on our proposed CI.

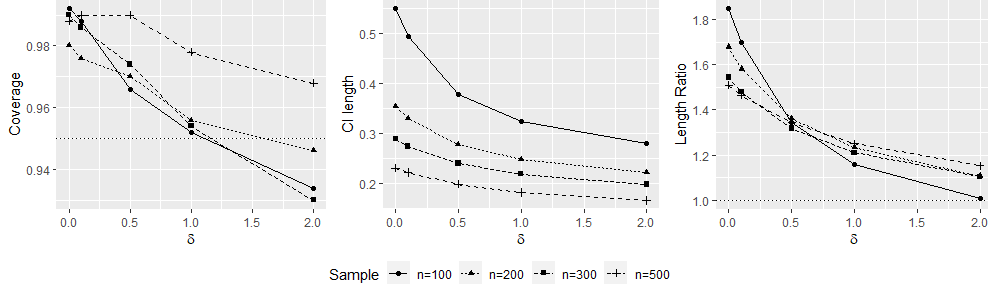


Figure 5: Setting 4(b) with  $L = 5$ . “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).

## 7 Real Data Applications

We apply the DenseNet sampling method to a genome-wide association study (Bloom et al., 2013) on the yeast colony growth under 46 different media. The study is based on  $n = 1008$  *Saccharomyces cerevisiae* segregants crossbred from a laboratory and a wine strain. A set of  $p = 4410$  Single Nucleotide Polymorphisms (SNPs) have been selected out of the total 11623 SNPs (Bloom et al., 2013). The outcome variables of interest are the end-point colony sizes under different growth media. To demonstrate our method, we consider the colony sizes under five growth media: “Ethanol”, “Lactate”, “Lactose”, “Sorbitol” and “Trehalose”. Our model (1) can be applied here with  $L = 5$  and each  $1 \leq l \leq 5$  corresponds to one growth medium (environment). These outcome variables are normalized to have variance 1, with the corresponding variance explained by the genetic markers 0.60, 0.69, 0.68, 0.51, and 0.66. We focus on an index set  $\mathcal{S} \subset [p]$  consisting of 10 pre-selected SNPs, which are relabeled with the indexes  $\{1, 2, \dots, 10\}$ . In Table 4, we report the genes where these SNPs occur. For example, the SNP with index 1 occurs in the KRE33 gene.

Index of SNP	1	2	3	4	5	6	7	8	9	10
Gene Name	KRE33	PUF3	CIR2	WHI2	BUD3	YOR019W	IRA2	HSP78	MKT1	PHM7

Table 4: The corresponding gene where the SNP occurs. The SNP with index 6 occurs at a region near the gene “YOR019W”, but not inside the gene.

On the top panel of Figure 6, we report the CIs for the regression coefficients  $\{b_j^{(l)}\}_{1 \leq l \leq 5, j \in \mathcal{S}}$ . On the middle and bottom panels of Figure 6, we plot our proposed CIs for the maximin effect  $[\beta_\delta^*(\mathbb{Q})]_j$  for  $j \in \mathcal{S}$  and  $\delta \in \{0, 0.1, 0.5, 1\}$ . The middle panel corresponds to the no covariate shift setting and the bottom panel to the covariate shift setting with  $\Sigma^{\mathbb{Q}} = \mathbf{I}$ . In Figure 6, we observe the three patterns, which are summarized in the following.

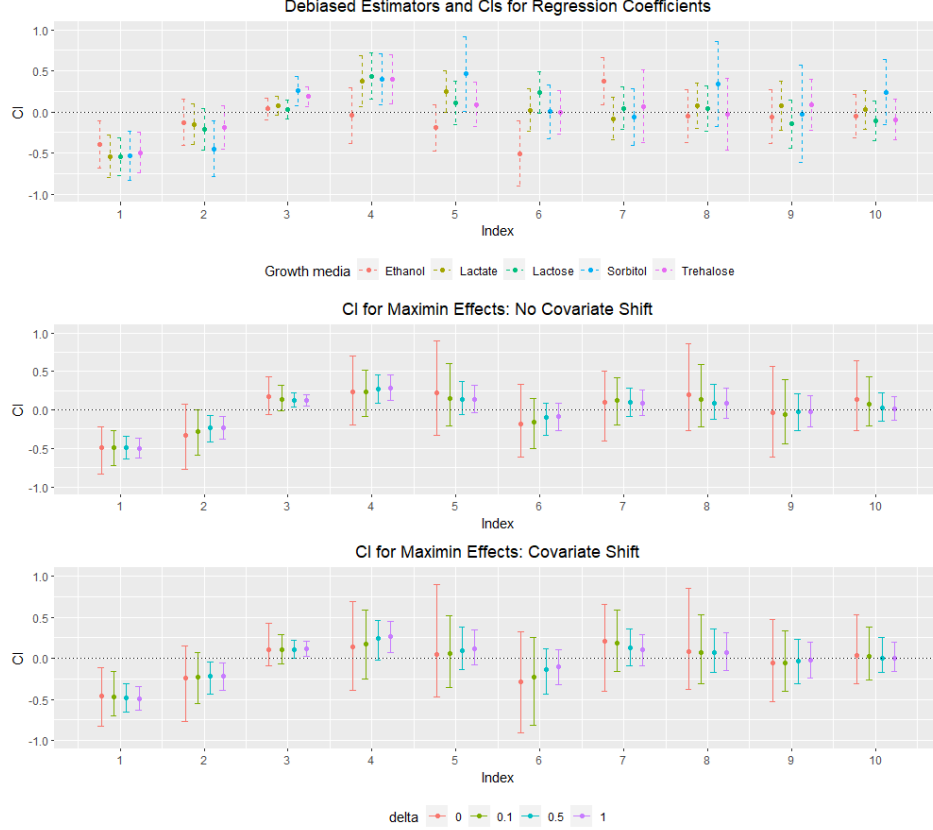


Figure 6: The top panel plots the debiased estimators of  $\{b_j^{(l)}\}_{1 \leq l \leq 5, j \in \mathcal{S}}$  and the corresponding CIs for  $\{b_j^{(l)}\}_{1 \leq l \leq 5, j \in \mathcal{S}}$ , where  $\mathcal{S}$  is the index set of ten pre-selected SNPs in Table 4. The middle and bottom panels plot respectively the CIs for  $\{[\beta_\delta^*(\mathbb{Q})]_j\}_{j \in \mathcal{S}}$  in the no covariate shift setting and the covariance-shift setting, with  $\delta$  varied across  $\{0, 0.1, 0.5, 1\}$ .

1. *Homogeneous effects.* The effects of SNPs with indexes 1 to 4 have (nearly) homogeneous signs across different media. The maximin effects are significant to certain extents, matching with our interpretation that the maximin effect captures the homogeneous effects. The SNP with index 1 is the most significant. The corresponding gene KRE33 is an essential gene for yeast (Cherry et al., 2012), which is a gene absolutely required to maintain life provided that all nutrients are available (Zhang and Lin, 2009). For SNPs with indexes 2, 3, 4,  $[\beta_\delta^*(\mathbb{Q})]_j$  are significant for a larger  $\delta$ .
2. *Heterogeneous effects.* For SNPs with  $j = 5, 6, 7$ , the SNPs have significant effects for one or two growth media but have opposite effects or near null effects in other media. The maximin effects are not significant, matching with the interpretation that the maximin effect tends to shrink the effects with heterogeneous signs towards zero. The corresponding genes for index  $j = 5, 6, 7$  are non-essential genes or genes of unknown function (Cherry et al., 2012).

3. *Near-null Effects.* For  $j = 8, 9, 10$ , the SNPs do not have any significant effect across different growth media. The corresponding maximin effect is not significant.

To conclude, our multi-source data analysis demonstrates that the maximin effect captures the homogeneous effects across different environments. Since the SNPs with indexes 1 to 4 have stable and sign-consistent effects across different growth media, the corresponding genes might have some causal interpretation and warrant further investigation.

## 8 Conclusion and Discussions

Our proposed covariate shift maximin effect transfers the information from multiple source populations to the unlabelled target population and guarantees excellent predictive performance even for the adversarially generated target population. This unique way of constructing robust models with heterogenous data creates additional statistical inference challenges. Our proposed sampling approach effectively addresses these challenges and is helpful for addressing other non-regular inference problems. An interesting direction is to study the maximin effects when the linear models in (1) are possibly misspecified ([Wasserman, 2014](#); [Bühlmann and van de Geer, 2015](#)), which is left for future research.

## Supplement

The supplement contains all proofs and additional methods, theories, and simulation results.

## Acknowledgement

Z. Guo acknowledges Dr. Peter Bühlmann, Dr. Tianxi Cai, Dr. Tony Cai, and Dr. Nicolai Meinshausen for their helpful discussions about the definitions and interpretations of the maximin effect and Dr. Minge Xie for the thought-provoking discussion on the sampling method. Z. Guo is grateful to Mr. Zhenyu Wang for the help with simulation studies, to Mr. Molei Liu, Dr. Hou Jue, and Dr. Rong Ma for their helpful comments about an earlier draft.

## References

Andrews, D. W. (1999). Estimation when a parameter is on a boundary. *Econometrica* 67(6), 1341–1383.

- Andrews, D. W. (2000). Inconsistency of the bootstrap when a parameter is on the boundary of the parameter space. *Econometrica*, 399–405.
- Arjovsky, M., L. Bottou, I. Gulrajani, and D. Lopez-Paz (2019). Invariant risk minimization. *arXiv preprint arXiv:1907.02893*.
- Athey, S., G. W. Imbens, and S. Wager (2018). Approximate residual balancing: debiased inference of average treatment effects in high dimensions. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 80(4), 597–623.
- Belloni, A., V. Chernozhukov, and C. Hansen (2014). High-dimensional methods and inference on structural and treatment effects. *Journal of Economic Perspectives* 28(2), 29–50.
- Belloni, A., V. Chernozhukov, and L. Wang (2011). Square-root lasso: pivotal recovery of sparse signals via conic programming. *Biometrika* 98(4), 791–806.
- Bickel, P. J., Y. Ritov, and A. B. Tsybakov (2009). Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics* 37(4), 1705–1732.
- Bloom, J. S., I. M. Ehrenreich, W. T. Loo, T.-L. V. Lite, and L. Kruglyak (2013). Finding the sources of missing heritability in a yeast cross. *Nature* 494(7436), 234–237.
- Bühlmann, P. and N. Meinshausen (2015). Magging: maximin aggregation for inhomogeneous large-scale data. *Proceedings of the IEEE* 104(1), 126–135.
- Bühlmann, P. and S. van de Geer (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.
- Bühlmann, P. and S. van de Geer (2015). High-dimensional inference in misspecified linear models. *Electronic Journal of Statistics* 9(1), 1449–1473.
- Cai, T., M. Liu, and Y. Xia (2021). Individual data protected integrative regression analysis of high-dimensional heterogeneous data. *Journal of the American Statistical Association*, 1–15.
- Cai, T. and W. Liu (2011). Adaptive thresholding for sparse covariance matrix estimation. *Journal of the American Statistical Association* 106(494), 672–684.
- Cai, T., T. Tony Cai, and Z. Guo (2021). Optimal statistical inference for individualized treatment effects in high-dimensional models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 83(4), 669–719.
- Cai, T. T. and Z. Guo (2017). Confidence intervals for high-dimensional linear regression: Minimax rates and adaptivity. *The Annals of Statistics* 45(2), 615–646.
- Cai, T. T. and Z. Guo (2020). Semisupervised inference for explained variance in high dimensional linear regression and its applications. *Journal of the Royal Statistical Society: Series B* 82(2), 391–419.
- Chernozhukov, V., D. Chetverikov, M. Demirer, E. Duflo, C. Hansen, W. Newey, and J. Robins (2018). Double/debiased machine learning for treatment and structural parameters: Double/debiased machine learning. *The Econometrics Journal* 21(1), C1–C68.

- Chernozhukov, V., C. Hansen, and M. Spindler (2015). Valid post-selection and post-regularization inference: An elementary, general approach. *Annu. Rev. Econ.* 7(1), 649–688.
- Cherry, J. M., E. L. Hong, C. Amundsen, R. Balakrishnan, G. Binkley, E. T. Chan, K. R. Christie, M. C. Costanzo, S. S. Dwight, S. R. Engel, et al. (2012). Saccharomyces genome database: the genomics resource of budding yeast. *Nucleic acids research* 40(D1), D700–D705.
- Diana, E., W. Gill, M. Kearns, K. Kenthapadi, and A. Roth (2021). Minimax group fairness: Algorithms and experiments. In *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, pp. 66–76.
- Ding, Z., M. Shao, and Y. Fu (2016). Incomplete multisource transfer learning. *IEEE transactions on neural networks and learning systems* 29(2), 310–323.
- Drton, M. (2009). Likelihood ratio tests and singularities. *The Annals of Statistics* 37(2), 979–1012.
- Dwork, C., M. Hardt, T. Pitassi, O. Reingold, and R. Zemel (2012). Fairness through awareness. In *Proceedings of the 3rd innovations in theoretical computer science conference*, pp. 214–226.
- Efron, B. (1979). Bootstrap methods: Another look at the jackknife. *The Annals of Statistics* 7(1), 1–26.
- Efron, B. and R. J. Tibshirani (1994). *An introduction to the bootstrap*. CRC press.
- Farrell, M. H. (2015). Robust inference on average treatment effects with possibly more covariates than observations. *Journal of Econometrics* 189(1), 1–23.
- Friedman, J., T. Hastie, and R. Tibshirani (2010). Regularization paths for generalized linear models via coordinate descent. *Journal of statistical software* 33(1), 1.
- Gao, R., X. Chen, and A. J. Kleywegt (2017). Wasserstein distributional robustness and regularization in statistical learning. *arXiv preprint arXiv:1712.06050*.
- Guo, Z., C. Renaux, P. Bühlmann, and T. T. Cai (2019). Group inference in high dimensions with applications to hierarchical testing. *arXiv preprint arXiv:1909.01503*.
- Guo, Z., W. Wang, T. T. Cai, and H. Li (2019). Optimal estimation of genetic relatedness in high-dimensional linear models. *Journal of the American Statistical Association* 114(525), 358–369.
- Hannig, J., H. Iyer, R. C. S. Lai, and T. C. M. Lee (2016). Generalized fiducial inference: A review and new results. *Journal of American Statistical Association*. To appear. Accepted in March 2016. doi:10.1080/01621459.2016.1165102.
- Hu, W., G. Niu, I. Sato, and M. Sugiyama (2018). Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning*, pp. 2029–2037. PMLR.
- Javanmard, A. and A. Montanari (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research* 15(1), 2869–2909.



- Keys, K. L., A. C. Mak, M. J. White, W. L. Eckalbar, A. W. Dahl, J. Mefford, A. V. Mikhaylova, M. G. Contreras, J. R. Elhawary, C. Eng, et al. (2020). On the cross-population generalizability of gene expression prediction models. *PLoS genetics* 16(8), e1008927.
- Kraft, P., E. Zeggini, and J. P. Ioannidis (2009). Replication in genome-wide association studies. *Statistical science: a review journal of the Institute of Mathematical Statistics* 24(4), 561.
- Li, S., T. T. Cai, and H. Li (2020). Transfer learning for high-dimensional linear regression: Prediction, estimation, and minimax optimality. *arXiv preprint arXiv:2006.10593*.
- Liu, M., Y. Xia, T. Cai, and K. Cho (2020). Integrative high dimensional multiple testing with heterogeneity under data sharing constraints. *arXiv preprint arXiv:2004.00816*.
- Martinez, N., M. Bertran, and G. Sapiro (2020). Minimax pareto fairness: A multi objective perspective. In *International Conference on Machine Learning*, pp. 6755–6764. PMLR.
- Meinshausen, N. and P. Bühlmann (2015). Maximin effects in inhomogeneous large-scale data. *The Annals of Statistics* 43(4), 1801–1830.
- Pan, S. J. and Q. Yang (2009). A survey on transfer learning. *IEEE Transactions on knowledge and data engineering* 22(10), 1345–1359.
- Peters, J., P. Bühlmann, and N. Meinshausen (2016). Causal inference by using invariant prediction: identification and confidence intervals. *Journal of the Royal Statistical Society. Series B (Statistical Methodology)*, 947–1012.
- Politis, D. N., J. P. Romano, and M. Wolf (1999). *Subsampling*. Springer Science & Business Media.
- Rakshit, P., T. T. Cai, and Z. Guo (2021). SihR: An R package for statistical inference in high-dimensional linear and logistic regression models. *arXiv preprint arXiv:2109.03365*.
- Rasmy, L., Y. Wu, N. Wang, X. Geng, W. J. Zheng, F. Wang, H. Wu, H. Xu, and D. Zhi (2018). A study of generalizability of recurrent neural network-based predictive models for heart failure onset risk using a large and heterogeneous ehr data set. *Journal of biomedical informatics* 84, 11–16.
- Rawls, J. (2001). *Justice as fairness: A restatement*. Harvard University Press.
- Rothenhäusler, D., N. Meinshausen, and P. Bühlmann (2016). Confidence intervals for maximin effects in inhomogeneous large-scale data. In *Statistical Analysis for High-Dimensional Data*, pp. 255–277. Springer.
- Rothenhäusler, D., N. Meinshausen, P. Bühlmann, and J. Peters (2018). Anchor regression: heterogeneous data meets causality. *arXiv preprint arXiv:1801.06229*.
- Sagawa, S., P. W. Koh, T. B. Hashimoto, and P. Liang (2019). Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*.

- Self, S. G. and K.-Y. Liang (1987). Asymptotic properties of maximum likelihood estimators and likelihood ratio tests under nonstandard conditions. *Journal of the American Statistical Association* 82(398), 605–610.
- Sharma, S., J.-L. Langhendries, P. Watzinger, P. Kötter, K.-D. Entian, and D. L. Lafontaine (2015). Yeast kre33 and human nat10 are conserved 18s rna cytosine acetyltransferases that modify trnas assisted by the adaptor tan1/thumpd1. *Nucleic acids research* 43(4), 2242–2258.
- Shi, C., R. Song, W. Lu, and B. Fu (2018). Maximin projection learning for optimal treatment decision with heterogeneous individualized treatment effects. *Journal of the Royal Statistical Society. Series B, Statistical methodology* 80(4), 681.
- Shimodaira, H. (2000). Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of statistical planning and inference* 90(2), 227–244.
- Singh, H., V. Mhasawade, and R. Chunara (2021). Generalizability challenges of mortality risk prediction models: A retrospective analysis on a multi-center database. *medRxiv*.
- Sinha, A., H. Namkoong, R. Volpi, and J. Duchi (2017). Certifying some distributional robustness with principled adversarial training. *arXiv preprint arXiv:1710.10571*.
- Sirugo, G., S. M. Williams, and S. A. Tishkoff (2019). The missing diversity in human genetic studies. *Cell* 177(1), 26–31.
- Sugiyama, M., M. Krauledat, and K.-R. Mäzller (2007). Covariate shift adaptation by importance weighted cross validation. *Journal of Machine Learning Research* 8(May), 985–1005.
- Sun, T. and C.-H. Zhang (2012). Scaled sparse linear regression. *Biometrika* 101(2), 269–284.
- Tian, Y. and Y. Feng (2021). Transfer learning under high-dimensional generalized linear models. *arXiv preprint arXiv:2105.14328*.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* 58(1), 267–288.
- Tsuboi, Y., H. Kashima, S. Hido, S. Bickel, and M. Sugiyama (2009). Direct density ratio estimation for large-scale covariate shift adaptation. *Journal of Information Processing* 17, 138–155.
- van de Geer, S., P. Bühlmann, Y. Ritov, and R. Dezeure (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics* 42(3), 1166–1202.
- Vershynin, R. (2012). Introduction to the non-asymptotic analysis of random matrices. In Y. Eldar and G. Kutyniok (Eds.), *Compressed Sensing: Theory and Applications*, pp. 210–268. Cambridge University Press.
- Verzelen, N. and E. Gassiat (2018). Adaptive estimation of high-dimensional signal-to-noise ratios. *Bernoulli* 24(4B), 3683–3710.
- Wang, P. and M. Xie (2020). Repro sampling method for statistical inference of high dimensional linear models. *Research Manuscript*.

- Wasserman, L. (2014). Discussion:” a significance test for the lasso”. *The Annals of Statistics* 42(2), 501–508.
- Xie, M. and K. Singh (2013). Confidence distribution, the frequentist distribution estimator of a parameter: a review. *International Statistical Review* 81, 3–39.
- Yao, Y. and G. Doretto (2010). Boosting for transfer learning with multiple sources. In *2010 IEEE computer society conference on computer vision and pattern recognition*, pp. 1855–1862. IEEE.
- Ye, F. and C.-H. Zhang (2010). Rate minimaxity of the lasso and dantzig selector for the  $l_q$  loss in  $l_r$  balls. *The Journal of Machine Learning Research* 11, 3519–3540.
- Zabell, S. L. (1992). Ra fisher and fiducial argument. *Statistical Science* 7(3), 369–387.
- Zhang, C.-H. and S. S. Zhang (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76(1), 217–242.
- Zhang, R. and Y. Lin (2009). Deg 5.0, a database of essential genes in both prokaryotes and eukaryotes. *Nucleic acids research* 37(suppl\_1), D455–D458.
- Zhao, T., G. Cheng, and H. Liu (2016). A partially linear framework for massive heterogeneous data. *Annals of statistics* 44(4), 1400.
- Zhou, S. (2009). Restricted eigenvalue conditions on subgaussian random matrices. *arXiv preprint arXiv:0912.4045*.
- Zhu, Y. and J. Bradic (2018). Linear hypothesis testing in dense high-dimensional linear models. *Journal of the American Statistical Association*, 1–18.
- Zhuang, F., Z. Qi, K. Duan, D. Xi, Y. Zhu, H. Zhu, H. Xiong, and Q. He (2020). A comprehensive survey on transfer learning. *Proceedings of the IEEE* 109(1), 43–76.

# Supplement to “Transfer Learning with Multi-source Data: High-dimensional Inference for Group Distributionally Robust Models”

The supplementary materials are organized as follows,

1. In Section A, we provide additional discussions, methods, and theories.
2. In Section B, we further discuss the non-regularity and instability challenges for confidence interval construction with bootstrap and subsampling methods.
3. In Section C, we present the proof of Theorem 1.
4. In Section D, we collect all proofs about the theoretical properties of  $\hat{\Gamma}^{\mathbb{Q}}$ .
5. In Section E, we present the proof of Theorem 2.
6. In Section F, we present the proofs of Propositions 1, 2, and 3.
7. In Section G, we provide the proofs of extra lemmas.
8. In Section H, we present additional simulations.

## A Additional Discussions

- In Section A.1, we provide the proof of (9) in the main paper.
- In Section A.2, we introduce an equivalent definition of the covariate-shift maximin effect from the hypothetical outcome perspective.
- In Section A.3, we show that the ridge-type maximin effect is the solution to a distributional robustness optimization problem.
- In Section H.2, we present a data-dependent way of choosing  $\delta$ .
- In Section A.5, we discuss the connection to minimax group fairness.
- In Section A.6, we formulate the maximin projection as a form of the covariate-shift maximin effect.
- In Section A.7, we provide the details on inference for  $\Gamma^{\mathbb{Q}}$  in high dimensions.
- In Section A.8, we consider the setting with no covariate shift and present the inference method for  $\Gamma^{\mathbb{Q}}$  in high dimensions.

### A.1 Proof of (9) in the main paper

We start with the proof of (9) in the main paper. For any  $\mathbb{T} \in \mathcal{C}(\mathbb{Q}_X)$ , we express its conditional outcome model as  $\mathbb{T}_{Y|X} = \sum_{l=1}^L q_l \cdot \mathbb{P}_{Y|X}^{(l)}$  for some weight vector  $q \in \Delta^L$ . Then we have

$$\begin{aligned} \mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^\top \beta)^2 &= \mathbf{E}_{\mathbb{T}} [2Y_i X_i^\top \beta - \beta^\top X_i X_i^\top \beta] \\ &= \mathbf{E}_{X_i \sim \mathbb{Q}_X} \left[ \sum_{l=1}^L q_l \cdot \mathbf{E}_{Y_i|X_i \sim \mathbb{P}_{Y|X}^{(l)}} [2Y_i X_i^\top \beta - \beta^\top X_i X_i^\top \beta] \right] \\ &= \sum_{l=1}^L q_l \cdot \mathbf{E}_{X_i \sim \mathbb{Q}_X} \mathbf{E}_{Y_i|X_i \sim \mathbb{P}_{Y|X}^{(l)}} [2Y_i X_i^\top \beta - \beta^\top X_i X_i^\top \beta]. \end{aligned} \quad (35)$$

By the outcome model (1), we have

$$\mathbf{E}_{X_i \sim \mathbb{Q}_X} \mathbf{E}_{Y_i|X_i \sim \mathbb{P}_{Y|X}^{(l)}} [2Y_i X_i^\top \beta - \beta^\top X_i X_i^\top \beta] = 2b^{(l)} \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta,$$

where  $\Sigma^{\mathbb{Q}} = \mathbf{E} X_1^{\mathbb{Q}} (X_1^{\mathbb{Q}})^\top$ . Together with (35) and the definition  $\mathbb{B} = \{b^{(1)}, \dots, b^{(L)}\}$ , we have

$$R_{\mathbb{Q}}(\beta) = \min_{q \in \Delta^L} \sum_{l=1}^L q_l \cdot [2b^{(l)} \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta] = \min_{b \in \mathbb{B}} [2b^\top \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta], \quad (36)$$

where the last equality follows from the fact that the extreme value of the linear function in  $q$  is achieved at the boundary points of  $\Delta^L$ . Hence, we show that the covariate-shift maximin effect definitions in (7) and (9) in the main paper are equivalent.

For the no-covariate shift setting, we write  $\Sigma^{\mathbb{Q}} = \Sigma$  and  $\Sigma^{(l)} = \Sigma$  for  $1 \leq l \leq L$ . Then we have

$$\min_{1 \leq l \leq L} \left\{ \mathbf{E}(Y_i^{(l)})^2 - \mathbf{E}(Y_i^{(l)} - [X_i^{(l)}]^\top \beta)^2 \right\} = \min_{1 \leq l \leq L} [2[b^{(l)}]^\top \Sigma \beta - \beta^\top \Sigma \beta].$$

Together with (36), in the no covariate-shift setting, we establish the equivalence between the maximin effect in Meinshausen and Bühlmann (2015),

$$\beta^* := \arg \max_{\beta \in \mathbb{R}^p} R(\beta) \quad \text{with} \quad R(\beta) = \min_{1 \leq l \leq L} \left\{ \mathbf{E}(Y_i^{(l)})^2 - \mathbf{E}(Y_i^{(l)} - [X_i^{(l)}]^\top \beta)^2 \right\},$$

and the group distributionally robust definition of the maximin effect,

$$\beta^* := \arg \max_{\beta \in \mathbb{R}^p} \min_{\mathbb{T} \in \mathcal{C}(\mathbb{P}_X)} \left\{ \mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^\top \beta)^2 \right\}.$$

## A.2 Covariate-shift Maximin Effect: Definition by Hypothetical Outcome

We now introduce another equivalent definition of the covariate-shift maximin effect defined in (7). Consider the hypothetical data generation mechanism,

$$Y_1^{*,(l)} = [X_1^{\mathbb{Q}}]^{\top} b^{(l)} + \epsilon_1^{(l)} \quad \text{for } 1 \leq l \leq L, \quad (37)$$

where  $b^{(l)}$  and  $\epsilon_1^{(l)}$  are the same as those in (1) and  $X_1^{\mathbb{Q}} \sim \mathbb{Q}_X$ . We use the super-index  $*$  in  $Y_1^{*,(l)}$  to denote that the outcome variable is hypothetical instead of being observed.  $Y_1^{*,(l)}$  stands for the hypothetical outcome that we will observe for a subject with  $X_1^{\mathbb{Q}}$  being generated from the  $l$ -th source population. Under (37), we define the maximin effect with respect to the covariate distribution  $\mathbb{Q}$  as

$$\beta^*(\mathbb{Q}) = \arg \max_{\beta \in \mathbb{R}^p} R_{\mathbb{Q}}(\beta) \quad \text{with} \quad R_{\mathbb{Q}}(\beta) = \min_{1 \leq l \leq L} \left[ \mathbf{E}^*(Y_1^{*,(l)})^2 - \mathbf{E}^*(Y_1^{*,(l)} - (X_1^{\mathbb{Q}})^{\top} \beta)^2 \right], \quad (38)$$

where  $\mathbf{E}^*$  is the expectation taken over the joint distribution of  $(Y_1^{*,(l)}, X_1^{\mathbb{Q}})$  defined in (37). Note that

$$\min_{1 \leq l \leq L} \left[ \mathbf{E}^*(Y_1^{*,(l)})^2 - \mathbf{E}^*(Y_1^{*,(l)} - (X_1^{\mathbb{Q}})^{\top} \beta)^2 \right] = \min_{b \in \mathbb{B}} \left[ 2b^{\top} \Sigma^{\mathbb{Q}} \beta - \beta^{\top} \Sigma^{\mathbb{Q}} \beta \right].$$

Hence, it follows from (9) in the main paper that the definition in (38) is equivalent to that in (7) in the main paper.

## A.3 Additional Discussion on Ridge-type Maximin Effect

Consider the setting  $p \geq L$  and  $\lambda_L(\mathcal{B}) > 0$  with  $\mathcal{B} = (b^{(1)}, \dots, b^{(L)}) \in \mathbb{R}^{p \times L}$ . We conduct the singular value decomposition  $\mathcal{B} = U_{p \times L} \Lambda_{L \times L} V_{L \times L}^{\top}$ . For  $1 \leq i \leq N_{\mathbb{Q}}$ , the noise vector  $W_i \in \mathbb{R}^p$  is generated as  $W_i = \sqrt{\delta} \cdot U W_i^0$  with  $W_i^0 \sim \mathcal{N}(\mathbf{0}, \Lambda^{-2})$  and  $W_i^0 \in \mathbb{R}^L$  being independent of  $X_i^{\mathbb{Q}}$ . Note that the distribution of  $W_i$  depends on the penalty level  $\delta$  and the coefficient matrix  $\mathcal{B}$ . Then we show that the ridge-type maximin effect in (30) in the main paper is the solution to a distributional robust optimization problem,

$$\beta_{\delta}^*(\mathbb{Q}) := \arg \max_{\beta \in \mathbb{R}^p} \min_{\mathbb{T} \in \mathcal{C}(\mathbb{Q}_X^{\delta})} \left\{ \mathbf{E}_{\mathbb{T}} Y_i^2 - \mathbf{E}_{\mathbb{T}} (Y_i - X_i^{\top} \beta)^2 \right\}, \quad (39)$$

where  $\mathbb{Q}_X^{\delta}$  denotes the distribution of  $X_i^{\mathbb{Q}} + W_i$  and  $\mathcal{C}(\mathbb{Q}_X^{\delta})$  is defined in (6) with  $\mathbb{Q}_X$  replaced by  $\mathbb{Q}_X^{\delta}$ . We now show that the distributional robustness definition of  $\beta_{\delta}^*(\mathbb{Q})$  in (39) is the

same as the definition in (30). By applying Proposition 1, we have

$$\beta_\delta^*(\mathbb{Q}) = \sum_{l=1}^L [\gamma_\delta^*(\mathbb{Q})]_l b^{(l)} \quad \text{with} \quad \gamma_\delta^*(\mathbb{Q}) = \arg \min_{\gamma \in \Delta^L} \gamma^\top \Gamma^{\mathbb{Q}^\delta} \gamma, \quad (40)$$

where  $\Gamma^{\mathbb{Q}^\delta} = [\mathcal{B}]^\top \mathbf{E}(X_i + W_i)(X_i + W_i)^\top \mathcal{B}$ . Since  $W_i \in \mathbb{R}^p$  is generated as  $W_i = \sqrt{\delta} \cdot UW_i^0$  with  $W_i^0 \sim N(\mathbf{0}, \Lambda^{-2})$  and  $W_i^0 \in \mathbb{R}^L$  being independent of  $X_i$ , we further have

$$\begin{aligned} \Gamma^{\mathbb{Q}^\delta} &= [\mathcal{B}]^\top \mathbf{E}(X_i + W_i)(X_i + W_i)^\top \mathcal{B} \\ &= [\mathcal{B}]^\top \Sigma^{\mathbb{Q}} \mathcal{B} + \delta \cdot [\mathcal{B}]^\top U \Lambda^{-2} U^\top \mathcal{B} \\ &= [\mathcal{B}]^\top \Sigma^{\mathbb{Q}} \mathcal{B} + \delta \cdot \mathbf{I}. \end{aligned}$$

Combined with (40), we show that the definitions (39) and (30) in the main paper are equal.

#### A.4 Choice of $\delta$ in the Ridge-type Maximin Effect

Recall the stability measure (depending on  $\delta$ ) introduced in Section 5 in the main paper,

$$\mathbb{I}(\delta) = \frac{\sum_{m=1}^M \|\hat{\gamma}_\delta^{[m]} - \hat{\gamma}_\delta\|_2^2}{\sum_{m=1}^M \|\hat{\Gamma}^{[m]} - \hat{\Gamma}^{\mathbb{Q}}\|_2^2},$$

where  $\{\hat{\Gamma}^{[m]}\}_{1 \leq m \leq M}$  and  $\{\hat{\gamma}_\delta^{[m]}\}_{1 \leq m \leq M}$  are the resampled estimates defined in (18) and (31) in the main paper, respectively. A large value of  $\mathbb{I}(\delta)$  indicates that a small error in estimating  $\Gamma^{\mathbb{Q}}$  may lead to a much larger error in estimating the weight vector.

If  $\mathbb{I}(0) < 0.5$ , the maximin effect without the ridge penalty is a stable aggregation and we simply apply the maximin effect without the ridge penalty. If  $\mathbb{I}(0) > 0.5$ , our proposed ridge-type effect with  $\delta > 0$  leads to a more stable aggregation. We shall choose  $\delta > 0$  to balance the stability and the prediction accuracy, which is measured by the reward  $R_{\mathbb{Q}}(\beta)$  defined in (9). For the aggregation weight  $\hat{\gamma}_\delta$  defined in (15), we estimate its reward by  $\hat{R}(\hat{\gamma}_\delta) = \min_{\gamma \in \Delta^L} [2\gamma^\top \hat{\Gamma}^{\mathbb{Q}} \hat{\gamma}_\delta - \hat{\gamma}_\delta^\top \hat{\Gamma}^{\mathbb{Q}} \hat{\gamma}_\delta]$ , and recommend the largest  $\delta \in [0, 2]$  such that  $\hat{R}(\hat{\gamma}_\delta)/\hat{R}(\hat{\gamma}_{\delta=0}) \geq 0.95$ . See the detailed simulation results in Section H.2 in the supplement.

#### A.5 Minimax Group Fairness and Rawlsian Max-min Principle

Fairness is an important consideration for designing the machine learning algorithm. In particular, the algorithm trained to maximize average performance on the training data set might under-serve or even cause harm to a sub-population of individuals (Dwork et al.,

2012). The goal of the *group fairness* is to build a model satisfying a certain fairness notation (e.g. statistical parity) across predefined sub-populations. However, such fairness notation can be typically achieved by downgrading the performance on the benefitted groups without improving the disadvantaged ones (Martinez et al., 2020; Diana et al., 2021). To address this, Martinez et al. (2020); Diana et al. (2021) proposed the minimax group fairness algorithm which ensures certain fairness principle and also maximizes the utility for each sub-population.

We assume that we have access to the i.i.d data  $\{Y_i, X_i, A_i\}_{1 \leq i \leq n}$ , where for the  $i$ -th observation,  $Y_i$  and  $X_i \in \mathbb{R}^p$  denote the outcome and the covariates, respectively, and  $A_i$  denotes the sensitive variable (e.g. age or sex). The training data can be separated into different sub-groups depending on the value of the sensitive variable  $A_i$ . For a discrete  $A_i$ , we use  $\mathcal{A}$  to denote the set of all possible values that  $A_i$  can take. Then the minimax group fairness can be defined as

$$\beta^{\text{mm-fair}} := \arg \min_{\beta} \max_{a \in \mathcal{A}} \mathbf{E}_{Y_i, X_i | A_i=a} \ell(Y_i, X_i^T \beta) = \arg \max_{\beta} \min_{a \in \mathcal{A}} \mathbf{E}_{Y_i, X_i | A_i=a} [-\ell(Y_i, X_i^T \beta)]$$

where  $\ell(Y_i, X_i^T \beta)$  denotes a loss function and  $-\ell(Y_i, X_i^T \beta)$  can be viewed as a reward/utility function. In terms of utility maximization, the idea of minimax group fairness estimator dates at least back to Rawlsian max-min fairness (Rawls, 2001).

When the distribution of  $X_i$  does not change with the value of  $A_i$ , then minimax group fairness estimator  $\beta^{\text{mm-fair}}$  is equivalent to the minimax or the maximin estimator, where the group label is determined by the value of  $A_i$ .

## A.6 Individualized Treatment Effect: Maximin Projection

Shi et al. (2018) proposed the maximin projection algorithm to construct the optimal treatment regime for new patients by leveraging training data from different groups with heterogeneity in optimal treatment decision. As explained in Shi et al. (2018), the heterogeneity in optimal treatment decision might come from patients' different enrollment periods and/or the treatment quality from different healthcare centers. Particularly, Shi et al. (2018) considered that the data is collected from  $L$  heterogeneous groups. For the  $l$ -th group with  $1 \leq l \leq L$ , let  $Y_i^{(l)} \in \mathbb{R}$ ,  $A_i^{(l)}$  and  $X_i^{(l)} \in \mathbb{R}^p$  denote the outcome, the treatment and the baseline covariates, respectively. Shi et al. (2018) considered the following model for the



data in group  $l$ ,

$$Y_i^{(l)} = h_l(X_i^{(l)}) + A_i^{(l)} \cdot [(b^{(l)})^\top X_i^{(l)} + c] + e_i^{(l)} \quad \text{with} \quad \mathbf{E}(e_i^{(l)} \mid X_i^{(l)}, A_i^{(l)}) = 0,$$

where  $h_l : \mathbb{R}^p \rightarrow \mathbb{R}$  denotes the unknown baseline function for the group  $l$  and the vector  $b^{(l)} \in \mathbb{R}^p$  describes the individualized treatment effect. To address the heterogeneity in optimal treatment regimes, Shi et al. (2018) has proposed the maximum projection

$$\beta^{*,\text{MP}} = \arg \max_{\|\beta\|_2 \leq 1} \min_{1 \leq l \leq L} \beta^\top b^{(l)}. \quad (41)$$

After identifying  $\beta^{*,\text{MP}}$ , we may construct the treatment regime for a new patient with covariates  $x_{\text{new}}$  by testing

$$H_0 : x_{\text{new}}^\top \beta^{*,\text{MP}} + c < 0. \quad (42)$$

The following Proposition 3 identifies the maximin projection  $\beta^{*,\text{MP}}$ . Through comparing it with Proposition 1, we note that the maximin projection is proportional to the general maximin effect defined in (9) with  $\Sigma^\mathbb{Q} = \mathbf{I}$  and hence the identification of  $\beta^*(\mathbb{I})$  is instrumental in identifying  $\beta^{*,\text{MP}}$ , which provides the strong motivation for statistical inference for  $x_{\text{new}}^\top \beta^*(\mathbb{I})$ .

**Proposition 3.** *The maximum projection  $\beta^{*,\text{MP}}$  in (41) satisfies*

$$\beta^{*,\text{MP}} = \frac{1}{\|\beta^*(\mathbb{I})\|_2} \beta^*(\mathbb{I}) \quad \text{with} \quad \beta^*(\mathbb{I}) = \sum_{l=1}^L [\gamma^*(\mathbb{I})]_l b^{(l)}$$

where  $\gamma^*(\mathbb{I}) = \arg \min_{\gamma \in \Delta^L} \gamma^\top \Gamma^\mathbb{I} \gamma$  and  $\Gamma_{lk}^\mathbb{I} = (b^{(l)})^\top b^{(k)}$  for  $1 \leq l, k \leq L$ .

We refer to Shi et al. (2018) for more details on the maximin projection in the low-dimensional setting. Our proposed DenseNet sampling method is useful in devising statistical inference methods for  $\beta^{*,\text{MP}}$  in high dimensions.

## A.7 Debiased Estimators of $\Gamma^\mathbb{Q}$ : Covariate Shift setting

We present the details about constructing the debiased estimator  $\hat{\Gamma}_{l,k}^\mathbb{Q}$  in (23). We use the same notations as in Section 3.3.

We estimate  $\{b^{(l)}\}_{1 \leq l \leq L}$  by applying Lasso ([Tibshirani, 1996](#)) to the sub-sample with the index set  $A_l$ :

$$\hat{b}_{init}^{(l)} = \arg \min_{b \in \mathbb{R}^p} \frac{\|Y_{A_l}^{(l)} - X_{A_l}^{(l)} b\|_2^2}{2|A_l|} + \lambda_l \sum_{j=1}^p \frac{\|X_{A_l,j}^{(l)}\|_2}{\sqrt{|A_l|}} |b_j|, \text{ with } \lambda_l = \sqrt{\frac{(2+c) \log p}{|A_l|}} \sigma_l \quad (43)$$

for some constant  $c > 0$ . The Lasso estimators  $\{\hat{b}_{init}^{(l)}\}_{1 \leq l \leq L}$  are implemented by the R-package `glmnet` ([Friedman et al., 2010](#)) with tuning parameters  $\{\lambda_l\}_{1 \leq l \leq L}$  chosen by cross validation. We may also construct the initial estimator  $\hat{b}_{init}^{(l)}$  by tuning-free penalized estimators ([Sun and Zhang, 2012](#); [Belloni et al., 2011](#)). Define

$$\hat{\Sigma}^Q = \frac{1}{|B|} \sum_{i \in B} X_i^Q (X_i^Q)^\top, \quad \hat{\Sigma}^{(l)} = \frac{1}{|B_l|} \sum_{i \in B_l} X_i^{(l)} [X_i^{(l)}]^\top, \quad \tilde{\Sigma}^Q = \frac{1}{|A|} \sum_{i \in A} X_i^Q (X_i^Q)^\top.$$

For  $1 \leq l, k \leq L$ , the plug-in estimator  $[\hat{b}_{init}^{(l)}]^\top \hat{\Sigma}^Q \hat{b}_{init}^{(k)}$  has the error decomposition:

$$\begin{aligned} (\hat{b}_{init}^{(l)})^\top \hat{\Sigma}^Q \hat{b}_{init}^{(k)} - (b^{(l)})^\top \Sigma^Q b^{(k)} &= (\hat{b}_{init}^{(k)})^\top \hat{\Sigma}^Q (\hat{b}_{init}^{(l)} - b^{(l)}) + (\hat{b}_{init}^{(l)})^\top \hat{\Sigma}^Q (\hat{b}_{init}^{(k)} - b^{(k)}) \\ &\quad - (\hat{b}_{init}^{(l)} - b^{(l)})^\top \hat{\Sigma}^Q (\hat{b}_{init}^{(k)} - b^{(k)}) + (b^{(l)})^\top (\hat{\Sigma}^Q - \Sigma^Q) b^{(k)}. \end{aligned}$$

Our proposed estimator in (23) is to correct the plug-in estimator by accurately estimating  $(\hat{b}_{init}^{(k)})^\top \hat{\Sigma}^Q (\hat{b}_{init}^{(l)} - b^{(l)})$  with  $-\frac{1}{|B_l|} [\hat{u}^{(l,k)}]^\top [X_{B_l}^{(l)}]^\top (Y_{B_l}^{(l)} - X_{B_l}^{(l)} \hat{b}_{init}^{(l)})$ , where the projection direction  $\hat{u}^{(l,k)}$  is constructed as follows,

$$\hat{u}^{(l,k)} = \arg \min_{u \in \mathbb{R}^p} u^\top \hat{\Sigma}^{(l)} u \quad \text{subject to } \|\hat{\Sigma}^{(l)} u - \omega^{(k)}\|_\infty \leq \|\omega^{(k)}\|_2 \mu_l \quad (44)$$

$$\left| [\omega^{(k)}]^\top \hat{\Sigma}^{(l)} u - \|\omega^{(k)}\|_2^2 \right| \leq \|\omega^{(k)}\|_2^2 \mu_l \quad (45)$$

$$\|X_{B_l} u\|_\infty \leq \|\omega^{(k)}\|_2 \tau_l \quad (46)$$

with  $\mu_l \asymp \sqrt{\log p / |B_l|}$ ,  $\tau_l \asymp \sqrt{\log n_l}$ , and

$$\omega^{(k)} = \tilde{\Sigma}^Q \hat{b}_{init}^{(k)} \in \mathbb{R}^p. \quad (47)$$

Similarly, we approximate the bias  $(\hat{b}_{init}^{(l)})^\top \hat{\Sigma}^Q (\hat{b}_{init}^{(k)} - b^{(k)})$  by  $-\frac{1}{|B_k|} [\hat{u}^{(k,l)}]^\top [X_{B_k}^{(k)}]^\top (Y_{B_k}^{(k)} - X_{B_k}^{(k)} \hat{b}_{init}^{(k)})$ .

In the following, we detail why the bias component  $(\hat{b}_{init}^{(k)})^\top \hat{\Sigma}^Q (\hat{b}_{init}^{(l)} - b^{(l)})$  can be accurately approximated by  $-\frac{1}{|B_l|} [\hat{u}^{(l,k)}]^\top [X_{B_l}^{(l)}]^\top (Y_{B_l}^{(l)} - X_{B_l}^{(l)} \hat{b}_{init}^{(l)})$ . The corresponding approximation error

is

$$-\frac{1}{|B_l|}[\widehat{u}^{(l,k)}]^\top [X_{B_l}^{(l)}]^\top \epsilon_{B_l}^{(l)} + [\widehat{\Sigma}^{(l)}\widehat{u}^{(l,k)} - \widehat{\Sigma}^{\mathbb{Q}}\widehat{b}_{init}^{(k)}]^\top (\widehat{b}_{init}^{(l)} - b^{(l)}). \quad (48)$$

In the following, we provide intuitions on why the projection direction  $\widehat{u}^{(l,k)}$  proposed in (44), (45) and (46) ensures a small approximation error in (48). The objective  $u^\top \widehat{\Sigma}^{(l)} u$  in (44) is proportional to the variance of the first term in (48). The constraint set in (44) implies  $\widehat{\Sigma}^{(l)}\widehat{u}^{(l,k)} - \widehat{\Sigma}^{\mathbb{Q}}\widehat{b}_{init}^{(k)} \approx \widehat{\Sigma}^{(l)}\widehat{u}^{(l,k)} - \omega^{(k)} \approx \mathbf{0}$ , which guarantees the second term of (48) to be small. The additional constraint (45) is seemingly useless to control the approximation error in (48). However, this additional constraint ensures that the first term in (48) dominates the second term in (48), which is critical in constructing an asymptotically normal estimator of  $\Gamma_{l,k}^{\mathbb{Q}}$ . The additional constraint (45) is particularly useful in the covariate shift setting, that is,  $\Sigma^{(l)} \neq \Sigma^{\mathbb{Q}}$  for some  $1 \leq l \leq L$ . The last constraint (46) is useful in establishing the asymptotic normality of the debiased estimator for non-Gaussian errors. We believe that the constraint (46) is only needed for technical reasons.

To construct  $\widehat{u}^{(l,k)}$  in (23), we solve the dual problem of (44) and (45),

$$\widehat{h} = \arg \min_{h \in \mathbb{R}^{p+1}} h^\top H^\top \widehat{\Sigma}^{(l)} H h / 4 + (\omega^{(k)})^\top H h / \|\omega^{(k)}\|_2 + \lambda \|h\|_1 \text{ with } H = \left[ \omega^{(k)} / \|\omega^{(k)}\|_2, \mathbf{I}_{p \times p} \right], \quad (49)$$

where we adopt the notation  $0/0 = 0$ . The objective value of this dual problem is unbounded from below when  $H^\top \widehat{\Sigma}^{(l)} H$  is singular and  $\lambda$  is near zero. We choose the smallest  $\lambda > 0$  such that the dual problem is bounded from below and construct  $\widehat{u}^{(l,k)} = -\frac{1}{2}(\widehat{h}_{-1} + \widehat{h}_1 \omega^{(k)} / \|\omega^{(k)}\|_2)$ .

**Remark 8** (Known  $\Sigma^{\mathbb{Q}}$ ). If  $\Sigma^{\mathbb{Q}}$  is known, we modify  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}}$  in (23) by replacing  $\widehat{\Sigma}^{\mathbb{Q}}$  by  $\Sigma^{\mathbb{Q}}$  and  $\omega^{(k)}$  in (47) by  $\omega^{(k)} = \Sigma^{\mathbb{Q}}\widehat{b}_{init}^{(k)}$ . This estimator (with known  $\Sigma^{\mathbb{Q}}$ ) is of a smaller variance as there is no uncertainty of estimating  $\Sigma^{\mathbb{Q}}$ ; see Figure 4 in the main paper for numerical comparisons.

In the following, we provide the theoretical guarantee of our proposed estimator  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}}$ . Define

$$\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)} = \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} + \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)}, \quad (50)$$

with

$$\begin{aligned} \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} &= \frac{\sigma_{l_1}^2}{|B_{l_1}|} (\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} [\widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = l_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = l_1)] \\ &\quad + \frac{\sigma_{k_1}^2}{|B_{k_1}|} (\widehat{u}^{(k_1, l_1)})^\top \widehat{\Sigma}^{(k_1)} [\widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = k_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = k_1)], \end{aligned} \quad (51)$$

and

$$\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} = \frac{1}{|B|} (\mathbf{E}[b^{(l_1)}]^\top X_i^\mathbb{Q} [b^{(k_1)}]^\top X_i^\mathbb{Q} [b^{(l_2)}]^\top X_i^\mathbb{Q} [b^{(k_2)}]^\top X_i^\mathbb{Q} - (b^{(l_1)})^\top \Sigma^\mathbb{Q} b^{(k_1)} (b^{(l_2)})^\top \Sigma^\mathbb{Q} b^{(k_2)}). \quad (52)$$

We introduce some extra notations. For random objects  $X_1$  and  $X_2$ , we use  $X_1 \stackrel{d}{=} X_2$  to denote that they are equal in distribution. For a sequence of random variables  $X_n$  indexed by  $n$ , we use  $X_n \xrightarrow{p} X$  and  $X_n \xrightarrow{d} X$  to represent that  $X_n$  converges to  $X$  in probability and in distribution, respectively.

The following proposition shows that the entry-wise estimation error  $\hat{\Gamma}_{l,k}^\mathbb{Q} - \Gamma_{l,k}^\mathbb{Q}$  can be approximated by a normal random variable. The proof of the following proposition can be found in Section D.5.

**Proposition 4.** *Consider the model (1). Suppose Condition (A1) holds,  $\frac{s \log p}{\min\{n, N_\mathbb{Q}\}} \rightarrow 0$  with  $n = \min_{1 \leq l \leq L} n_l$  and  $s = \max_{1 \leq l \leq L} \|b^{(l)}\|_0$ . Then the proposed estimator  $\hat{\Gamma}_{l,k}^\mathbb{Q} \in \mathbb{R}^{L \times L}$  in (23) in the main paper satisfies  $\hat{\Gamma}_{l,k}^\mathbb{Q} - \Gamma_{l,k}^\mathbb{Q} = D_{l,k} + \text{Rem}_{l,k}$ , where*

$$\frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}}} \xrightarrow{d} \mathcal{N}(0, 1), \quad (53)$$

with  $\mathbf{V}$  defined in (50); for  $1 \leq l, k \leq L$ , with probability larger than  $1 - \min\{n, p\}^{-c}$  for a constant  $c > 0$ , the reminder term  $\text{Rem}_{l,k}$  satisfies

$$|\text{Rem}_{l,k}| \lesssim (1 + \|\omega^{(k)}\|_2 + \|\omega^{(l)}\|_2) \frac{s \log p}{n} + (\|b^{(k)}\|_2 + \|b^{(l)}\|_2) \sqrt{\frac{s(\log p)^2}{n N_\mathbb{Q}}}, \quad (54)$$

where  $c > 0$  is a positive constant and  $\omega^{(k)}$  and  $\omega^{(l)}$  are defined in (47).

We control the diagonal of the covariance matrix  $\mathbf{V}$  in the following proposition, whose proof can be found in Section D.4.

**Proposition 5.** *Suppose that the assumptions of Proposition 4 hold. Then with probability larger than  $1 - \min\{n, p\}^{-c}$ , the diagonal element  $\mathbf{V}_{\pi(l,k), \pi(l,k)}$  in (50) for  $(l, k) \in \mathcal{I}_L$  satisfies,*

$$\frac{\|\omega^{(l)}\|_2^2}{n_k} + \frac{\|\omega^{(k)}\|_2^2}{n_l} \lesssim \mathbf{V}_{\pi(l,k), \pi(l,k)}^{(a)} \lesssim \frac{\|\omega^{(l)}\|_2^2}{n_k} + \frac{\|\omega^{(k)}\|_2^2}{n_l}, \quad \mathbf{V}_{\pi(l,k), \pi(l,k)}^{(b)} \lesssim \frac{\|b^{(l)}\|_2^2 \|b^{(k)}\|_2^2}{N_\mathbb{Q}}, \quad (55)$$

where  $c > 0$  is a positive constant and  $\omega^{(l)}$  and  $\omega^{(k)}$  are defined in (47).

- If  $\Sigma^{\mathbb{Q}}$  is known, then with probability larger than  $1 - \min\{n, p\}^{-c}$ ,

$$n \cdot \mathbf{V}_{\pi(l,k), \pi(l,k)} \lesssim \|b^{(k)}\|_2^2 + \|b^{(l)}\|_2^2 + s \log p/n.$$

- If  $\Sigma^{\mathbb{Q}}$  is unknown, then with probability larger than  $1 - \min\{n, p\}^{-c}$ ,

$$n \cdot \mathbf{V}_{\pi(l,k), \pi(l,k)} \lesssim \left(1 + \frac{p}{N_{\mathbb{Q}}}\right)^2 \left(\|b^{(k)}\|_2^2 + \|b^{(l)}\|_2^2 + s \frac{\log p}{n}\right) + \frac{n}{N_{\mathbb{Q}}} \|b^{(l)}\|_2^2 \|b^{(k)}\|_2^2. \quad (56)$$

The above proposition controls the variance  $\mathbf{V}_{\pi(l,k), \pi(l,k)}$  of  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}}$  for the covariate shift setting. For known  $\Sigma^{\mathbb{Q}}$ , we show that the diagonal elements of  $\mathbf{V}$  is of order  $1/n$  if  $\max_{1 \leq l \leq L} \|b^{(l)}\|_2$  is bounded and  $s \lesssim n/\log p$ . If the matrix  $\Sigma^{\mathbb{Q}}$  is estimated from the data, then (56) shows that the diagonal elements of  $\mathbf{V}$  is of order  $1/n$  under the additional assumption  $N_{\mathbb{Q}} \gtrsim \max\{p, n\}$ . This requires a relatively large sample size  $N_{\mathbb{Q}}$  of the unlabelled covariate data for the target population while  $\{n_l\}_{1 \leq l \leq L}$  are allowed to be much smaller than  $p$ .

**Remark 9** (Related literature). We shall discuss briefly a few related works on quadratic functional inference in high dimensions. [Verzelen and Gassiat \(2018\)](#); [Cai and Guo \(2020\)](#); [Guo et al. \(2019\)](#) considered inference for quadratic functionals in a single high-dimensional linear model while [Guo et al. \(2019\)](#) proposed debiased estimators of  $[b^{(j)}]^{\top} b^{(k)}$  for  $1 \leq j, k \leq L$ . The main challenge in the covariate shift setting is that  $\Sigma^{(l)} \neq \Sigma^{\mathbb{Q}}$  for some  $1 \leq l \leq L$  and  $\omega^{(k)} = \widetilde{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k)} \in \mathbb{R}^p$  can be an arbitrarily dense vector. To address this, the additional constraint (45) is proposed to construct a valid projection direction in (44).

## A.8 Debiased Estimators of $\Gamma^{\mathbb{Q}}$ : No Covariate Shift setting

A few simplifications can be made when there is no covariate shift. For  $1 \leq l \leq L$ , we estimate  $b^{(l)}$  by applying Lasso to the whole data set  $(X^{(l)}, Y^{(l)})$ :

$$\widehat{b}^{(l)} = \arg \min_{b \in \mathbb{R}^p} \|Y^{(l)} - X^{(l)}b\|_2^2 / (2n_l) + \lambda_l \sum_{j=1}^p \|X_{:,j}^{(l)}\|_2 / \sqrt{n_l} \cdot |b_j| \quad (57)$$

with  $\lambda = \sqrt{(2+c) \log p / n_l \sigma_l}$  for some constant  $c > 0$ . Since  $\Sigma^{(l)} = \Sigma^{\mathbb{Q}}$  for  $1 \leq l \leq L$ , we define

$$\widehat{\Sigma} = \frac{1}{\sum_{l=1}^L n_l + N_{\mathbb{Q}}} \left( \sum_{l=1}^L \sum_{i=1}^{n_l} X_i^{(l)} [X_i^{(l)}]^{\top} + \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{(l)} [X_i^{(l)}]^{\top} \right)$$

and estimate  $\Gamma_{l,k}$  by

$$\widehat{\Gamma}_{l,k}^{\mathbb{Q}} = (\widehat{b}^{(l)})^{\top} \widehat{\Sigma} \widehat{b}^{(k)} + (\widehat{b}^{(l)})^{\top} \frac{1}{n_k} [X^{(k)}]^{\top} (Y^{(k)} - X^{(k)} \widehat{b}^{(k)}) + (\widehat{b}^{(k)})^{\top} \frac{1}{n_l} [X^{(l)}]^{\top} (Y^{(l)} - X^{(l)} \widehat{b}^{(l)}). \quad (58)$$

This estimator can be viewed as a special case of (23) by taking  $\widehat{u}^{(l,k)}$  and  $\widehat{u}^{(k,l)}$  as  $\widehat{b}^{(k)}$  and  $\widehat{b}^{(l)}$ , respectively. Neither the optimization in (44) and (45) nor the sample splitting is needed for constructing the debiased estimator in the no covariate shift setting.

Define  $N = \sum_{l=1}^L n_l + N_{\mathbb{Q}}$  and the covariance matrices  $\mathbf{V} = (\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)})_{(l_1, k_1) \in \mathcal{I}_L, (l_2, k_2) \in \mathcal{I}_L} \in \mathbb{R}^{L(L+1)/2 \times L(L+1)/2}$  for  $j = 1, 2$  as

$$\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)} = \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(1)} + \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(2)}, \quad (59)$$

with

$$\begin{aligned} \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(1)} &= \frac{\sigma_{l_1}^2}{n_{l_1}^2} [b^{(k_1)}]^{\top} \mathbf{E} \left( [X^{(l_1)}]^{\top} [X^{(l_2)} b^{(k_2)} \mathbf{1}(l_2 = l_1) + X^{(k_2)} b^{(l_2)} \mathbf{1}(k_2 = l_1)] \right) \\ &\quad + \frac{\sigma_{k_1}^2}{n_{l_1}^2} [b^{(l_1)}]^{\top} \mathbf{E} \left( [X^{(k_1)}]^{\top} [X^{(l_2)} b^{(k_2)} \mathbf{1}(l_2 = k_1) + X^{(k_2)} b^{(l_2)} \mathbf{1}(k_2 = k_1)] \right), \end{aligned} \quad (60)$$

$$\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(2)} = \frac{\mathbf{E}[b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} [b^{(k_1)}]^{\top} X_i^{\mathbb{Q}} [b^{(l_2)}]^{\top} X_i^{\mathbb{Q}} [b^{(k_2)}]^{\top} X_i^{\mathbb{Q}} - (b^{(l_1)})^{\top} \Sigma^{\mathbb{Q}} b^{(k_1)} (b^{(l_2)})^{\top} \Sigma^{\mathbb{Q}} b^{(k_2)}}{N}. \quad (61)$$

In the following proposition, we establish the properties of the estimator  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}}$  in (58) for no covariate shift setting, whose proof is presented in Section D.6.

**Proposition 6.** *Consider the model (1). Suppose Condition (A1) holds and  $s \log p/n \rightarrow 0$  with  $n = \min_{1 \leq l \leq L} n_l$  and  $s = \max_{1 \leq l \leq L} \|b^{(l)}\|_0$ . If  $\{X_i^{(l)}\}_{1 \leq i \leq n_l} \stackrel{i.i.d.}{\sim} \mathbb{Q}_X$  for  $1 \leq l \leq L$ , then the estimator  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}}$  defined in (58) satisfies*

$$\widehat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}} = D_{l,k} + \text{Rem}_{l,k},$$

where  $\frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}}} \xrightarrow{d} \mathcal{N}(0, 1)$ ; for  $1 \leq l, k \leq L$ , with probability larger than  $1 - \min\{n, p\}^{-c}$  for a constant  $c > 0$ , the reminder term  $\text{Rem}_{l,k}$  satisfies

$$|\text{Rem}_{l,k}| \lesssim (1 + \|b^{(k)}\|_2 + \|b^{(l)}\|_2) \cdot \frac{s \log p}{n}. \quad (62)$$

With probability larger than  $1 - p^{-c}$  for a positive constant  $c > 0$ , for any  $(l, k) \in \mathcal{I}_L$ ,

$$\mathbf{V}_{\pi(l,k),\pi(l,k)}^{(1)} \asymp \frac{\|b^{(k)}\|_2^2 + \|b^{(k)}\|_2^2}{n} \quad \text{and} \quad \mathbf{V}_{\pi(l,k),\pi(l,k)}^{(2)} \lesssim \frac{\|b^{(l)}\|_2^2 \|b^{(k)}\|_2^2}{\sum_{l=1}^L n_l + N_{\mathbb{Q}}}. \quad (63)$$

We estimate the covariance between  $\hat{\Gamma}_{l_1,k_1}^{\mathbb{Q}} - \Gamma_{l_1,k_1}^{\mathbb{Q}}$  and  $\hat{\Gamma}_{l_2,k_2}^{\mathbb{Q}} - \Gamma_{l_2,k_2}^{\mathbb{Q}}$  by

$$\hat{\mathbf{V}}_{\pi(l_1,k_1),\pi(l_2,k_2)} = \hat{\mathbf{V}}_{\pi(l_1,k_1),\pi(l_2,k_2)}^{(1)} + \hat{\mathbf{V}}_{\pi(l_1,k_1),\pi(l_2,k_2)}^{(2)} \quad (64)$$

where

$$\begin{aligned} \hat{\mathbf{V}}_{\pi(l_1,k_1),\pi(l_2,k_2)}^{(1)} &= \frac{\hat{\sigma}_{l_1}^2}{n_{l_1}^2} [b^{(l_1)}]^\top [X^{(l_1)}]^\top X^{(l_1)} [b^{(l_2)} \mathbf{1}(l_2 = l_1) + b^{(k_2)} \mathbf{1}(k_2 = l_1)] \\ &\quad + \frac{\hat{\sigma}_{k_1}^2}{n_{k_1}^2} [b^{(k_1)}]^\top [X^{(k_1)}]^\top X^{(k_1)} [b^{(l_2)} \mathbf{1}(l_2 = k_1) + b^{(k_2)} \mathbf{1}(k_2 = k_1)] \\ \hat{\mathbf{V}}_{\pi(l_1,k_1),\pi(l_2,k_2)}^{(2)} &= \frac{\sum_{i=1}^{N_{\mathbb{Q}}} \left( (\hat{b}^{(l_1)})^\top X_i^{\mathbb{Q}} (\hat{b}^{(k_1)})^\top X_i^{\mathbb{Q}} (\hat{b}^{(l_2)})^\top X_i^{\mathbb{Q}} (\hat{b}^{(k_2)})^\top X_i^{\mathbb{Q}} - (\hat{b}^{(l_1)})^\top \hat{\Sigma} \hat{b}^{(k_1)} (\hat{b}^{(l_2)})^\top \hat{\Sigma} \hat{b}^{(k_2)} \right)}{(\sum_{l=1}^L n_l + N_{\mathbb{Q}})^2} \\ &\quad + \frac{\sum_{l=1}^L \sum_{i=1}^{n_l} \left( (\hat{b}^{(l_1)})^\top X_i^{(l)} (\hat{b}^{(k_1)})^\top X_i^{(l)} (\hat{b}^{(l_2)})^\top X_i^{(l)} (\hat{b}^{(k_2)})^\top X_i^{(l)} - (\hat{b}^{(l_1)})^\top \hat{\Sigma} \hat{b}^{(k_1)} (\hat{b}^{(l_2)})^\top \hat{\Sigma} \hat{b}^{(k_2)} \right)}{(\sum_{l=1}^L n_l + N_{\mathbb{Q}})^2} \end{aligned}$$

with

$$\hat{\Sigma} = \frac{1}{\sum_{l=1}^L n_l + N_{\mathbb{Q}}} \left( \sum_{l=1}^L \sum_{i=1}^{n_l} X_i^{(l)} [X_i^{(l)}]^\top + \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top \right).$$

## B Inference Challenges with Bootstrap and Subsampling

We demonstrate the challenges of confidence interval construction for the maximin effects with bootstrap and subsampling methods.

### B.1 Simulation Settings (I-1) to (I-10)

We focus on the no covariate shift setting with  $\Sigma^{(l)} = \mathbf{I}_p$  for  $1 \leq l \leq L$  and  $\Sigma^{\mathbb{Q}} = \mathbf{I}_p$ . In the following, we describe how to generate the settings (I-1) to (I-6) with non-regularity and instability. We set  $L = 4$ . For  $1 \leq l \leq L$ , we generate  $b^{(l)}$  as  $b_j^{(l)} = j/20 + \kappa_j^{(l)}$  for  $1 \leq j \leq 5$  with  $\{\kappa_j^{(l)}\}_{1 \leq j \leq 5} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, \sigma_{\text{irr}}^2)$ ,  $b_j^{(l)} = j/20$  for  $6 \leq j \leq 10$ , and  $b_j^{(l)} = 0$  for  $11 \leq j \leq p$ . Set

$[x_{\text{new}}]_j = 1$  for  $1 \leq j \leq 5$  and zero otherwise. We choose the following six combinations of  $\sigma_{\text{irr}}$  and the random seed for generating  $\kappa_j^{(l)}$ ,

(I-1)  $\sigma_{\text{irr}} = 0.05$ , seed = 42; (I-2)  $\sigma_{\text{irr}} = 0.05$ , seed = 20; (I-3)  $\sigma_{\text{irr}} = 0.10$ , seed = 36;

(I-4)  $\sigma_{\text{irr}} = 0.15$ , seed = 17; (I-5)  $\sigma_{\text{irr}} = 0.20$ , seed = 12; (I-6)  $\sigma_{\text{irr}} = 0.25$ , seed = 31.

In addition, we generate the following non-regular settings:

(I-7)  $L = 2$ ;  $b_1^{(1)} = 2$ ,  $b_j^{(1)} = j/40$  for  $2 \leq j \leq 10$  and  $b_j^{(1)} = 0$  otherwise;  $b_1^{(2)} = -0.03$ ,  $b_j^{(2)} = j/40$  for  $2 \leq j \leq 10$  and  $b_j^{(2)} = 0$  otherwise;  $x_{\text{new}} = e_1$ .

(I-8) Same as (I-7) except for  $b_j^{(l)} = (10 - j)/40$  for  $11 \leq j \leq 20$  and  $l = 1, 2$ ;

(I-9) Same as (I-7) except for  $b_j^{(l)} = 1$  for  $2 \leq j \leq 30$  and  $l = 1, 2$ .

Finally, we generate (I-10) as a favorable setting without non-regularity or instability.

(I-10)  $L = 2$ ;  $b_j^{(1)} = j/20$  for  $1 \leq j \leq 10$ ,  $b_j^{(2)} = -j/20$  for  $1 \leq j \leq 10$ ;  $[x_{\text{new}}]_j = j/5$  for  $1 \leq j \leq 5$  and  $[x_{\text{new}}]_j = 0$  otherwise.

## B.2 Challenges for Bootstrap and Subsampling: Numerical Evidence

In Section 6.1, we have reported the under-coverage of the normality CIs in high dimensions. In the following, we explore a low dimensional setting with  $p = 30$  and  $n_1 = \dots = n_L = n = 1000$ . We shall compare our proposed CI, the CI assuming asymptotic normality (Rothenhäusler et al., 2016), and CIs by the subsampling or bootstrap methods. We describe these methods in the following.

**DenseNet sampling for low dimensions.** Our proposed methods can be extended to the low dimensional setting by replacing the projection defined in  $\hat{u}^{(l,k)}$  defined in (44) to (46) by  $\hat{u}^{(l,k)} = [\hat{\Sigma}^{(l)}]^{-1}\omega^{(k)}$ . In low dimensions, we can also simplify the point estimator in (15) by

$$\widehat{x_{\text{new}}^\top \beta^*} = \sum_{l=1}^L \hat{\gamma}_l \cdot x_{\text{new}}^\top \hat{b}_{\text{OLS}}^{(l)}, \quad (65)$$

where  $\hat{\gamma}$  is defined in (15) and  $\hat{b}_{\text{OLS}}^{(l)}$  is the ordinary least square estimator computed by  $(X^{(l)}, Y^{(l)})$ .



**Magging estimator and CI assuming asymptotic normality.** In low-dimensional setting, the Magging estimator has been proposed in [Bühlmann and Meinshausen \(2015\)](#) to estimate the maximin effect. In low dimensions, the regression vector  $b^{(l)}$  is estimated by the ordinary least square estimator  $\hat{b}_{\text{OLS}}^{(l)}$  for  $1 \leq l \leq L$  and the covariance matrix  $\Sigma$  is estimated by the sample covariance matrix

$$\hat{\Sigma} = \frac{1}{\sum_{l=1}^L n_l} \sum_{l=1}^L \sum_{i=1}^n X_i^{(l)} [X_i^{(l)}]^\top.$$

Then the Magging estimator in low dimension is of the form,

$$\hat{\beta}^{\text{magging}} = \sum_{l=1}^L \hat{\gamma}_l \hat{b}_{\text{OLS}}^{(l)} \quad \text{with} \quad \hat{\gamma} := \arg \min_{\gamma \in \Delta^L} \gamma^\top \hat{\Gamma} \gamma \quad (66)$$

where  $\hat{\Gamma}_{lk} = (\hat{b}_{\text{OLS}}^{(l)})^\top \hat{\Sigma} \hat{b}_{\text{OLS}}^{(k)}$  for  $1 \leq l, k \leq L$  and  $\Delta^L = \{\gamma \in \mathbb{R}^L : \gamma_j \geq 0, \sum_{j=1}^L \gamma_j = 1\}$  is the simplex over  $\mathbb{R}^L$ . Implied by Proposition 1, if  $\lambda_{\min}(\hat{\Sigma}) > 0$ , then the Magging estimator can be expressed in the equivalent form,

$$\hat{\beta}^{\text{magging}} = \arg \max_{\beta \in \mathbb{R}^p} \hat{R}(\beta) \quad \text{with} \quad \hat{R}(\beta) = \min_{b \in \hat{\mathbb{B}}} \left[ 2b^\top \hat{\Sigma} \beta - \beta^\top \hat{\Sigma} \beta \right], \quad (67)$$

where  $\hat{\mathbb{B}} = \{\hat{b}_{\text{OLS}}^{(1)}, \dots, \hat{b}_{\text{OLS}}^{(L)}\}$ . In the expression (67), the Magging estimator can be viewed as optimizing a plug-in estimator of the reward  $R(\beta)$ , where  $\Sigma$  is estimated by  $\hat{\Sigma}$  and  $b^{(l)}$  is estimated by  $\hat{b}_{\text{OLS}}^{(l)}$  for  $1 \leq l \leq L$ .

[Rothenhäusler et al. \(2016\)](#) have established the asymptotic normality of the magging estimator under certain conditions, which essentially ruled out the non-regularity and instability settings. In the following, we shall show that the CI assuming asymptotic normality fails to provide valid inference for the low-dimensional maximin effects in the presence of non-regularity or instability. In particular, we construct a normality CI of the form

$$(x_{\text{new}}^\top \hat{\beta}^{\text{magging}} - 1.96 \cdot \widehat{\text{SE}}, x_{\text{new}}^\top \hat{\beta}^{\text{magging}} + 1.96 \cdot \widehat{\text{SE}}), \quad (68)$$

where  $\widehat{\text{SE}}$  denotes the sample standard deviation of  $x_{\text{new}}^\top \hat{\beta}^{\text{magging}}$  calculated based on 500 simulations. Since  $\widehat{\text{SE}}$  is calculated in an oracle way, this normality CI is a favorable implementation of the CI construction in [Rothenhäusler et al. \(2016\)](#).

**Bootstrap and subsampling.** We briefly describe the implementation of the bootstrap and subsampling methods. We compute the point estimator for the original data as in (66), denoted as  $\hat{\theta}$ . For  $1 \leq l \leq L$ , we randomly sample  $m$  observations (with/without replacement) from  $(X^{(l)}, Y^{(l)})$  and use this generated sample to compute the point estimator  $\hat{\theta}^{m,j}$  as in (66). We conduct the random sampling 500 times to obtain  $\{\hat{\theta}^{m,j}\}_{1 \leq j \leq 500}$  and define the empirical CDF as,

$$L_n(t) = \frac{1}{500} \sum_{j=1}^{500} \mathbf{1} \left( \sqrt{m}(\hat{\theta}^{m,j} - \hat{\theta}) \leq t \right),$$

where  $\mathbf{1}$  denotes the indicator function. Define  $\hat{t}_{\alpha/2}$  as the minimum  $t$  value such that  $L_n(t) \geq \alpha/2$  and  $\hat{t}_{1-\alpha/2}$  as the minimum  $t$  value such that  $L_n(t) \geq 1 - \alpha/2$ . We construct the bootstrap/subsampling confidence interval as

$$\left[ \hat{\theta} - \frac{\hat{t}_{1-\alpha/2}}{\sqrt{n}}, \hat{\theta} - \frac{\hat{t}_{\alpha/2}}{\sqrt{n}} \right].$$

In Table S1, we report the empirical coverage of the normality CI in (68), the CI by subsampling, the CI by  $m$  out of  $n$  bootstrap, and our proposed CI. The normality CI in (68) and the CIs by subsampling and bootstrap methods are in general under-coverage for settings (I-1) to (I-9). Our proposed CI achieves the desired coverage level at the expense of a wider interval. For the favorable setting (I-10), bootstrap methods achieve the desired coverage level while subsampling methods only work for a small  $m$ , which is an important requirement for the validity of subsampling methods (Politis et al., 1999).

Setting $p = 30$	normality	m-out-of-n subsampling				m-out-of-n bootstrap					Proposed	
		$m = 200$	$m = 300$	$m = 400$	$m = 500$	$m = 200$	$m = 300$	$m = 400$	$m = 500$	$m = 1000$	Cov	L-ratio
(I-1)	0.686	0.380	0.390	0.380	0.408	0.432	0.460	0.502	0.490	0.536	0.956	1.678
(I-2)	0.808	0.418	0.456	0.474	0.446	0.456	0.510	0.532	0.546	0.602	0.990	1.723
(I-3)	0.770	0.392	0.494	0.464	0.440	0.482	0.480	0.516	0.544	0.634	0.984	1.766
(I-4)	0.816	0.620	0.668	0.670	0.668	0.672	0.694	0.710	0.700	0.794	0.990	1.686
(I-5)	0.790	0.612	0.626	0.626	0.594	0.628	0.664	0.702	0.732	0.732	1.000	1.843
(I-6)	0.806	0.590	0.636	0.654	0.632	0.626	0.682	0.698	0.712	0.760	0.994	1.833
(I-7)	0.824	0.914	0.932	0.908	0.888	0.890	0.950	0.934	0.950	0.952	0.996	5.078
(I-8)	0.888	0.912	0.934	0.884	0.856	0.914	0.916	0.932	0.932	0.958	0.996	4.017
(I-9)	0.900	0.834	0.822	0.788	0.778	0.914	0.914	0.916	0.862	0.914	0.996	2.061
(I-10)	0.954	0.956	0.904	0.866	0.836	0.956	0.948	0.962	0.964	0.942	1.000	1.725

Table S1: Empirical coverage of the normality CI in (68), the CI by subsampling, the CI by  $m$  out of  $n$  bootstrap and our proposed CI, where the column indexed with L-ratio denoting the ratio of the average length of our proposed CI to that of the normality CI.

In Section B.3, we discuss why subsampling methods fail to provide valid inference for the maximin effect in the presence of non-regularity.

### B.3 Challenges for Bootstrap and Subsampling Methods: A Theoretical View

We illustrate the challenge of bootstrap and subsampling methods for the maximin effects in non-regular settings. As a remark, the following argument is not a rigorous proof but of a similar style to the discussion of Andrews (2000), explaining why the subsampling methods do not completely solve the non-regular inference problems. The main difficulty appears in a near boundary setting with

$$\gamma_1 = \frac{\mu_1}{\sqrt{n}} \quad \text{for a positive constant } \mu_1 > 0, \quad (69)$$

where  $\gamma_1$  denotes the weight of the first group.

To illustrate the problem, we consider the special setting  $L = 2$ ,  $n_1 = n_2 = n$  and  $b^{(1)}$  and  $b^{(2)}$  are known. The first coefficient of the maximin effect  $\beta^*$  can be expressed as

$$\beta_1^* = b_1^{(1)} \cdot \gamma_1 + b_1^{(2)} \cdot (1 - \gamma_1) = (b_1^{(1)} - b_1^{(2)}) \cdot \gamma_1 + b_1^{(2)}.$$

For this special scenario, the only uncertainty is from estimating  $\gamma_1$  since  $\Sigma^{\mathbb{Q}}$  is unknown. We estimate  $\gamma_1$  by

$$\hat{\gamma}_1 = \max\{\bar{\gamma}_1, 0\} \quad \text{with} \quad \bar{\gamma}_1 = \frac{\hat{\Gamma}_{22} - \hat{\Gamma}_{12}}{\hat{\Gamma}_{11} + \hat{\Gamma}_{22} - 2\hat{\Gamma}_{12}},$$

where  $\hat{\Gamma}_{12} = [b^{(1)}]^\top \hat{\Gamma}^{\mathbb{Q}} b^{(2)}$ ,  $\hat{\Gamma}_{11} = [b^{(1)}]^\top \hat{\Gamma}^{\mathbb{Q}} b^{(1)}$ , and  $\hat{\Gamma}_{22} = [b^{(2)}]^\top \hat{\Gamma}^{\mathbb{Q}} b^{(2)}$ . In the definition of  $\hat{\gamma}_1$ , we do not restrict it to be smaller than 1 as this happens with a high probability under our current setting (69). We then estimate  $\beta_1^*$  by

$$\hat{\beta}_1 = (b_1^{(1)} - b_1^{(2)}) \cdot \hat{\gamma}_1 + b_1^{(2)}.$$

We separately subsample  $\{X_i^{(1)}, Y_i^{(1)}\}_{1 \leq i \leq n}$  and  $\{X_i^{(2)}, Y_i^{(2)}\}_{1 \leq i \leq n}$  and use  $m$  to denote the subsample size. For  $1 \leq t \leq T$  with a positive integer  $T > 0$ , denote the  $t$ -th subsampled data as  $\{X_i^{(*,t,1)}, Y_i^{(*,t,1)}\}_{1 \leq i \leq m}$  and  $\{X_i^{(*,t,2)}, Y_i^{(*,t,2)}\}_{1 \leq i \leq m}$ . We apply these subsampled data

sets to compute the sample covariance matrix  $\widehat{\Sigma}^{(*,t)}$ . Then we compute  $\bar{\gamma}_1^{(*,t)}$  as

$$\hat{\gamma}_1^{(*,t)} = \max\{\bar{\gamma}_1^{(*,t)}, 0\} \quad \text{with} \quad \bar{\gamma}_1^{(*,t)} = \frac{\widehat{\Gamma}_{22}^{(*,t)} - \widehat{\Gamma}_{12}^{(*,t)}}{\widehat{\Gamma}_{11}^{(*,t)} + \widehat{\Gamma}_{22}^{(*,t)} - 2\widehat{\Gamma}_{12}^{(*,t)}},$$

with  $\widehat{\Gamma}_{12}^{(*,t)} = [b^{(1)}]^\top \widehat{\Sigma}^{(*,t)} b^{(2)}$ ,  $\widehat{\Gamma}_{11}^{(*,t)} = [b^{(1)}]^\top \widehat{\Sigma}^{(*,t)} b^{(1)}$ , and  $\widehat{\Gamma}_{22}^{(*,t)} = [b^{(2)}]^\top \widehat{\Sigma}^{(*,t)} b^{(2)}$ . Then we construct the subsampling estimator

$$\widehat{\beta}_1^{(*,t)} = (b_1^{(1)} - b_1^{(2)}) \cdot \hat{\gamma}_1^{(*,t)} + b_1^{(2)}.$$

We assume  $\sqrt{n}(\bar{\gamma}_1 - \gamma_1)$  and  $\sqrt{n}(\bar{\gamma}_1^{(*,t)} - \bar{\gamma}_1)$  share the same limiting normal distribution. Specifically, we assume  $\sqrt{n}(\bar{\gamma}_1 - \gamma_1) \xrightarrow{d} Z$  with  $Z \sim N(0, V_\gamma)$  and conditioning on the observed data,  $\sqrt{n}(\bar{\gamma}_1^{(*,t)} - \bar{\gamma}_1) \xrightarrow{d} Z$ . Then for the setting (69), we have

$$\sqrt{n}(\hat{\gamma}_1 - \gamma_1) \xrightarrow{d} \max\{Z, -\mu_1\}. \quad (70)$$

In the following, we will show that  $\sqrt{m}(\hat{\gamma}_1^{(*,t)} - \hat{\gamma}_1)$  does not approximate the limiting distribution of  $\sqrt{n}(\hat{\gamma}_1 - \gamma_1)$  in (70) if  $\gamma_1 = \frac{\mu_1}{\sqrt{n}}$ . Note that

$$\begin{aligned} \sqrt{m}(\hat{\gamma}_1^{(*,t)} - \hat{\gamma}_1) &= \max\{\sqrt{m}[\bar{\gamma}_1^{(*,t)} - \hat{\gamma}_1], -\sqrt{m}\hat{\gamma}_1\} \\ &= \max\{\sqrt{m}[\bar{\gamma}_1^{(*,t)} - \bar{\gamma}_1] + \sqrt{m}[\bar{\gamma}_1 - \gamma_1], -\sqrt{m}\gamma_1\} - \sqrt{m}(\hat{\gamma}_1 - \gamma_1). \end{aligned}$$

When  $m \ll n$  and  $\gamma_1 \asymp 1/\sqrt{n}$ , then the following event happens with a probability larger than  $1 - n^{-c}$  for a small positive constant  $c > 0$ ,

$$\mathcal{A}_0 = \left\{ \max\{\sqrt{m}[\bar{\gamma}_1 - \gamma_1], \sqrt{m}\gamma_1, \sqrt{m}(\hat{\gamma}_1 - \gamma_1)\} \lesssim \sqrt{\frac{m \log n}{n}} \right\}.$$

Then conditioning on the event  $\mathcal{A}_0$ ,

$$\sqrt{m}(\hat{\gamma}_1^{(*,t)} - \hat{\gamma}_1) \xrightarrow{d} \max\{Z, 0\},$$

which is different from the limiting distribution of  $\sqrt{n}(\hat{\gamma}_1 - \gamma_1)$  in (70).

## C Proof of Theorem 1

We first introduce some notations. For  $L > 0$  and  $\alpha_0 \in (0, 0.01]$ , define

$$C^*(L, \alpha_0) = c^*(\alpha_0) \cdot \text{Vol} \left[ \frac{L(L+1)}{2} \right] \quad \text{with} \quad c^*(\alpha_0) = \frac{\exp \left( -\frac{L(L+1)}{3} \frac{z_{\alpha_0/[L(L+1)]}^2 (n \cdot \lambda_{\max}(\mathbf{V}) + \frac{4}{3}d_0)}{n \cdot \lambda_{\min}(\mathbf{V}) + \frac{2}{3}d_0} \right)}{\sqrt{2\pi} \prod_{i=1}^{\frac{L(L+1)}{2}} [n \cdot \lambda_i(\mathbf{V}) + 4d_0/3]^{1/2}}, \quad (71)$$

where  $\text{Vol} \left[ \frac{L(L+1)}{2} \right]$  denotes the volume of a unit ball in  $L(L+1)/2$  dimensions.

As a corollary of Proposition 5, we prove in Section D.2 that  $C^*(L, \alpha_0)$  and  $c^*(\alpha_0)$  are lower bounded by a positive constant with a high probability.

**Corollary 1.** *Consider the model (1). Suppose Condition (A1) holds,  $\frac{s \log p}{\min\{n, N_{\mathbb{Q}}\}} \rightarrow 0$  with  $n = \min_{1 \leq l \leq L} n_l$  and  $s = \max_{1 \leq l \leq L} \|b^{(l)}\|_0$ . If  $N_{\mathbb{Q}} \gtrsim \max\{n, p\}$ , then with probability larger than  $1 - \min\{n, p\}^{-c}$  for some positive constant  $c > 0$ , then  $\min\{C^*(L, \alpha_0), c^*(\alpha_0)\} \geq c'$  for a positive constant  $c' > 0$ .*

The following theorem is a generalization of Theorem 1 in the main paper, which implies Theorem 1 by setting  $\delta = 0$ .

**Theorem 3.** *Consider the model (1). Suppose Conditions (A1) and (A2) hold. If  $\text{err}_n(M)$  defined in (25) satisfies  $\text{err}_n(M) \ll \min\{1, c^*(\alpha_0), \lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta\}$  where  $c^*(\alpha_0)$  is defined in (71), then*

$$\liminf_{n, p \rightarrow \infty} \mathbf{P} \left( \min_{m \in \mathbb{M}} \|\hat{\gamma}_{\delta}^{[m]} - \gamma_{\delta}^*\|_2 \leq \frac{\sqrt{2} \text{err}_n(M)}{\lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta} \cdot \frac{1}{\sqrt{n}} \right) \geq 1 - \alpha_0, \quad (72)$$

where  $\alpha_0 \in (0, 0.01]$  is the pre-specified constant used in the construction of  $\mathbb{M}$  in (20).

In the following, we prove the more general Theorem 3.

### C.1 The Proof Plan of Theorem 3

The proof of Theorem 3 relies on Lemma 1, whose proof can be found in Section D.1.

**Lemma 1.** *Consider the model (1). Suppose that Conditions (A1) and (A2) hold, then the estimator  $\hat{\Gamma}^{\mathbb{Q}}$  in (23) satisfies*

$$\liminf_{n, p \rightarrow \infty} \mathbf{P} \left( \max_{1 \leq l, k \leq L} \frac{|\hat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}}|}{\sqrt{\hat{\mathbf{V}}_{\pi(l,k), \pi(l,k)} + d_0/n}} \leq 1.05 \cdot z_{\alpha_0/[L(L+1)]} \right) \geq 1 - \alpha_0, \quad (73)$$

for any  $\alpha_0 \in (0, 0.01]$ .

The proof of Theorem 3 consists of the following three steps.

1. In Section D.1, we establish that our proposed estimator  $\hat{\Gamma}^{\mathbb{Q}}$  satisfies (73) with a high probability.
2. In Section C.2, we prove that, if  $\hat{\Gamma}^{\mathbb{Q}}$  satisfies (73), then the following inequality holds,

$$\liminf_{n,p \rightarrow \infty} \mathbf{P} \left( \min_{m \in \mathbb{M}} \|\hat{\Gamma}^{[m]} - \Gamma^{\mathbb{Q}}\|_F \leq \sqrt{2} \text{err}_n(M) / \sqrt{n} \right) \geq 1 - \alpha_0, \quad (74)$$

where  $\mathbb{M}$  is defined in (20).

3. In Section C.3, we apply (74) to establish Theorem 3.

## C.2 Proof of (74)

Define the feasible set

$$\mathcal{F}' = \left\{ S \in \mathbb{R}^{L(L+1)/2} : \max_{1 \leq \pi(l,k) \leq L(L+1)/2} \left| \frac{S_{\pi(l,k)}}{\sqrt{\hat{\mathbf{V}}_{\pi(l,k),\pi(l,k)} + d_0/n}} \right| \leq 1.1 \cdot z_{\alpha_0/[L(L+1)]} \right\}, \quad (75)$$

and we have the equivalent expression of  $\mathbb{M}$  defined in (20) as

$$\mathbb{M} = \{1 \leq m \leq M : S^{[m]} \in \mathcal{F}'\}.$$

Denote the data by  $\mathcal{O}$ , that is,  $\mathcal{O} = \{X^{(l)}, Y^{(l)}\}_{1 \leq l \leq L} \cup \{X^{\mathbb{Q}}\}$ . Recall  $n = \min_{1 \leq l \leq L} n_l$  and write  $\hat{\Gamma} = \hat{\Gamma}^{\mathbb{Q}}$ . Define  $\hat{S} = \text{vecl}(\hat{\Gamma}) - \text{vecl}(\Gamma)$ , the rescaled difference  $\hat{Z} = \sqrt{n}\hat{S}$ , and  $Z^{[m]} = \sqrt{n}S^{[m]} = \sqrt{n}[\text{vecl}(\hat{\Gamma}) - \text{vecl}(\hat{\Gamma}^{[m]})]$  for  $1 \leq m \leq M$ . Note that

$$\hat{Z} - Z^{[m]} = \sqrt{n}[\text{vecl}(\hat{\Gamma}^{[m]}) - \text{vecl}(\Gamma)], \quad (76)$$

and for the given data  $\mathcal{O}$ ,  $\hat{Z}$  is fixed. We define the rescaled covariance matrices as

$$\mathbf{Cov} = n\mathbf{V} \quad \text{and} \quad \widehat{\mathbf{Cov}} = n\hat{\mathbf{V}}, \quad (77)$$

with  $\mathbf{V}$  and  $\hat{\mathbf{V}}$  defined in (50) and (24), respectively.

The density of the rescaled variable  $Z^{[m]} = \sqrt{n}S^{[m]}$  is

$$f(Z^{[m]} = Z \mid \mathcal{O}) = \frac{\exp\left(-\frac{1}{2}Z^\top(\widehat{\mathbf{Cov}} + d_0\mathbf{I})^{-1}Z\right)}{\sqrt{2\pi\det(\widehat{\mathbf{Cov}} + d_0\mathbf{I})}}.$$

We define the following function to facilitate the proof,

$$g(Z) = \frac{1}{\sqrt{2\pi\det(\mathbf{Cov} + \frac{4}{3}d_0\mathbf{I})}} \exp\left(-\frac{1}{2}Z^\top(\mathbf{Cov} + \frac{2}{3}d_0\mathbf{I})^{-1}Z\right). \quad (78)$$

We define the following events for the data  $\mathcal{O}$ ,

$$\begin{aligned} \mathcal{E}_1 &= \left\{ \|\widehat{\mathbf{Cov}} - \mathbf{Cov}\|_2 < d_0/3 \right\}, \\ \mathcal{E}_2 &= \left\{ \max_{1 \leq l, k \leq L} \frac{|\widehat{Z}_{\pi(l,k)}|}{\sqrt{\widehat{\mathbf{Cov}}_{\pi(l,k), \pi(l,k)} + d_0}} \leq 1.05 \cdot z_{\alpha_0/[L(L+1)]} \right\}, \end{aligned} \quad (79)$$

where  $\|\widehat{\mathbf{Cov}} - \mathbf{Cov}\|_2$  denotes the spectral norm of the matrix  $\widehat{\mathbf{Cov}} - \mathbf{Cov}$ . The following lemma shows that the event  $\mathcal{E}_1$  holds with a high probability, whose proof is presented in Section [G.3](#).

**Lemma 2.** *Suppose that the conditions of Theorem [1](#) hold, then we have*

$$\mathbf{P}(\mathcal{E}_1) \geq 1 - \min\{N_{\mathbb{Q}}, n, p\}^{-c} \quad (80)$$

for some positive constant  $c > 0$ .

Together with [\(73\)](#), we establish

$$\liminf_{n \rightarrow \infty} \mathbf{P}(\mathcal{E}_1 \cap \mathcal{E}_2) \geq 1 - \alpha_0. \quad (81)$$

On the event  $\mathcal{O} \in \mathcal{E}_1$ , we have

$$2(\mathbf{Cov} + \frac{2}{3}d_0\mathbf{I}) \succeq \mathbf{Cov} + \frac{4}{3}d_0\mathbf{I} \succ \widehat{\mathbf{Cov}} + d_0\mathbf{I} \succ \mathbf{Cov} + \frac{2}{3}d_0\mathbf{I}, \quad (82)$$

where  $A \succeq B$  and  $A \succ B$  denotes that the matrix  $A - B$  is positive semi-definite and positive definite, respectively. By (82), we have

$$f(Z^{[m]} = Z \mid \mathcal{O}) \cdot \mathbf{1}_{\{\mathcal{O} \in \mathcal{E}_1\}} \geq g(Z) \cdot \mathbf{1}_{\{\mathcal{O} \in \mathcal{E}_1\}}. \quad (83)$$

for any  $Z \in \mathbb{R}^{L(L+1)/2}$ . Furthermore, on the event  $\mathcal{E}_1 \cap \mathcal{E}_2$ , we have

$$\begin{aligned} \frac{1}{2} \widehat{Z}^\top (\mathbf{Cov} + \frac{2}{3} d_0 \mathbf{I})^{-1} \widehat{Z} &\leq \frac{L(L+1)}{4} \frac{\max_{1 \leq l, k \leq L} (\widehat{\mathbf{Cov}}_{\pi(l,k), \pi(l,k)} + d_0) \cdot (1.05 \cdot z_{\alpha_0/[L(L+1)]})^2}{\lambda_{\min}(\mathbf{Cov}) + \frac{2}{3} d_0} \\ &\leq \frac{L(L+1)}{3} \frac{z_{\alpha_0/[L(L+1)]}^2 (\lambda_{\max}(\mathbf{Cov}) + \frac{4}{3} d_0)}{\lambda_{\min}(\mathbf{Cov}) + \frac{2}{3} d_0}. \end{aligned}$$

Hence, on the event  $\mathcal{E}_1 \cap \mathcal{E}_2$ , we have

$$g(\widehat{Z}) \geq c^*(\alpha_0). \quad (84)$$

We further lower bound the targeted probability in (74) as

$$\begin{aligned} &\mathbf{P} \left( \min_{m \in \mathbb{M}} \|\widehat{\Gamma}^{[m]} - \Gamma^\mathbb{Q}\|_F \leq \sqrt{2} \text{err}_n(M) / \sqrt{n} \right) \\ &\geq \mathbf{P} \left( \min_{m \in \mathbb{M}} \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \right) \\ &\geq \mathbf{E}_\mathcal{O} \left[ \mathbf{P} \left( \min_{m \in \mathbb{M}} \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \right], \end{aligned} \quad (85)$$

where  $\mathbf{P}(\cdot \mid \mathcal{O})$  denotes the conditional probability given the observed data  $\mathcal{O}$  and  $\mathbf{E}_\mathcal{O}$  denotes the expectation taken with respect to the observed data  $\mathcal{O}$ .

For  $m \notin \mathbb{M}$ , the definition of  $\mathbb{M}$  implies that there exists  $1 \leq k_0 \leq l_0 \leq L$  such that

$$\frac{\left| Z_{\pi(l_0, k_0)}^{[m]} \right|}{\sqrt{\widehat{\mathbf{Cov}}_{\pi(l_0, k_0), \pi(l_0, k_0)} + d_0}} \geq 1.1 \cdot z_{\alpha_0/[L(L+1)]}.$$

Hence, on the event  $\mathcal{E}_1 \cap \mathcal{E}_2$ ,

$$\|Z^{[m]} - \widehat{Z}\|_2 \geq \left| Z_{\pi(l_0, k_0)}^{[m]} - \widehat{Z}_{\pi(l_0, k_0)} \right| \geq \sqrt{\frac{2d_0}{3}} \cdot 0.05 \cdot z_{\alpha_0/[L(L+1)]}.$$



If  $\text{err}_n(M) \leq \sqrt{\frac{2d_0}{3}} \cdot 0.05 \cdot z_{\alpha_0/[L(L+1)]}$ , then we have

$$\min_{m \notin \mathbb{M}} \|Z^{[m]} - \hat{Z}\|_2 \geq \text{err}_n(M).$$

As a consequence, for  $\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2$ ,

$$\mathbf{P} \left( \min_{m \in \mathbb{M}} \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) = \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right). \quad (86)$$

Together with (85), we have

$$\begin{aligned} & \mathbf{P} \left( \min_{m \in \mathbb{M}} \|\hat{\Gamma}^{[m]} - \Gamma^{\mathbb{Q}}\|_F \leq \sqrt{2} \text{err}_n(M) / \sqrt{n} \right) \\ & \geq \mathbf{E}_{\mathcal{O}} \left[ \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \right]. \end{aligned} \quad (87)$$

Note that

$$\begin{aligned} & \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \\ & = 1 - \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \hat{Z}\|_2 \geq \text{err}_n(M) \mid \mathcal{O} \right) \\ & = 1 - \prod_{1 \leq m \leq M} \left[ 1 - \mathbf{P} \left( \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \right] \end{aligned}$$

where the second equality follows from the conditional independence of  $\{Z^{[m]}\}_{1 \leq m \leq M}$  given the data  $\mathcal{O}$ .

Since  $1 - x \leq e^{-x}$ , we further lower bound the above expression as

$$\begin{aligned} & \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \\ & \geq 1 - \prod_{1 \leq m \leq M} \exp \left[ -\mathbf{P} \left( \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \right] \\ & = 1 - \exp \left[ -M \cdot \mathbf{P} \left( \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \right]. \end{aligned}$$

Hence, we have

$$\begin{aligned} & \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ & \geq \left( 1 - \exp \left[ -M \cdot \mathbf{P} \left( \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \right] \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ & = 1 - \exp \left[ -M \cdot \mathbf{P} \left( \|Z^{[m]} - \hat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \right]. \end{aligned} \quad (88)$$

Together with (87), it is sufficient to establish an lower bound for

$$\mathbf{P} \left( \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2}. \quad (89)$$

We apply the inequality (83) and further lower bound the targeted probability in (89) as

$$\begin{aligned} & \mathbf{P} \left( \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &= \int f(Z^{[m]} = Z \mid \mathcal{O}) \cdot \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &\geq \left[ \int g(Z) \cdot \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \right] \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &= \left[ \int g(\widehat{Z}) \cdot \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \right] \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &\quad + \left[ \int [g(Z) - g(\widehat{Z})] \cdot \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \right] \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2}. \end{aligned} \quad (90)$$

By (84), we have

$$\begin{aligned} & \int g(\widehat{Z}) \cdot \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &\geq c^*(\alpha_0) \cdot \int \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &\geq c^*(\alpha_0) \cdot \text{Vol}(L(L+1)/2) \cdot [\text{err}_n(M)]^{L(L+1)/2} \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2}, \end{aligned} \quad (91)$$

where  $\text{Vol}(L(L+1)/2)$  denotes the volume of the unit ball in  $\frac{L(L+1)}{2}$ -dimension.

Note that there exists  $t \in (0, 1)$  such that

$$g(Z) - g(\widehat{Z}) = [\nabla g(\widehat{Z} + t(Z - \widehat{Z}))]^\top (Z - \widehat{Z}),$$

with  $\nabla g(w) = \frac{\exp(-\frac{1}{2}w^\top(\mathbf{Cov} + \frac{2}{3}d_0\mathbf{I})^{-1}w)}{\sqrt{2\pi\det(\mathbf{Cov} + \frac{4}{3}d_0\mathbf{I})}}w^\top(\mathbf{Cov} + \frac{2}{3}d_0\mathbf{I})^{-1}w$ . Since  $\lambda_{\min}(\mathbf{Cov} + \frac{2}{3}d_0\mathbf{I}) \geq \frac{2}{3}d_0$ , then  $\nabla g$  is bounded from the above and  $|g(Z) - g(\widehat{Z})| \leq C\|Z - \widehat{Z}\|_2$  for a positive constant  $C > 0$ . Then we establish

$$\begin{aligned} & \left| \int [g(Z) - g(\widehat{Z})] \cdot \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \right| \\ &\leq C \cdot \text{err}_n(M) \cdot \int \mathbf{1}_{\{\|Z - \widehat{Z}\|_2 \leq \text{err}_n(M)\}} dZ \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ &= C \cdot \text{err}_n(M) \cdot \text{Vol}(L(L+1)/2) \cdot [\text{err}_n(M)]^{L(L+1)/2} \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2}. \end{aligned} \quad (92)$$

By assuming  $C \cdot \text{err}_n(M) \leq \frac{1}{2}c^*(\alpha_0)$ , we combine (90), (91) and (92) and obtain

$$\begin{aligned} & \mathbf{P} \left( \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ & \geq \frac{1}{2}c^*(\alpha_0) \cdot \text{Vol}(L(L+1)/2) \cdot [\text{err}_n(M)]^{L(L+1)/2} \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2}. \end{aligned}$$

Together with (88), we establish

$$\begin{aligned} & \mathbf{P} \left( \min_{1 \leq m \leq M} \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \mid \mathcal{O} \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \\ & \geq 1 - \exp \left[ -M \cdot \frac{1}{2}c^*(\alpha_0) \cdot \text{Vol}(L(L+1)/2) \cdot [\text{err}_n(M)]^{L(L+1)/2} \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2} \right] \\ & = \left( 1 - \exp \left[ -M \cdot \frac{1}{2}c^*(\alpha_0) \cdot \text{Vol}(L(L+1)/2) \cdot [\text{err}_n(M)]^{L(L+1)/2} \right] \right) \cdot \mathbf{1}_{\mathcal{O} \in \mathcal{E}_1 \cap \mathcal{E}_2}. \end{aligned} \tag{93}$$

Together with (87), we have

$$\begin{aligned} & \mathbf{P} \left( \min_{m \in \mathbb{M}} \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \right) \\ & \geq \left( 1 - \exp \left[ -M \cdot \frac{1}{2}c^*(\alpha_0) \cdot \text{Vol}(L(L+1)/2) \cdot [\text{err}_n(M)]^{L(L+1)/2} \right] \right) \mathbf{P}(\mathcal{E}_1 \cap \mathcal{E}_2). \end{aligned}$$

We choose

$$\text{err}_n(M) = \left\lceil \frac{4 \log n}{C^*(L, \alpha_0)M} \right\rceil^{\frac{2}{L(L+1)}}$$

with  $C^*(L, \alpha_0)$  defined in (71) and establish

$$\mathbf{P} \left( \min_{m \in \mathbb{M}} \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \right) \geq (1 - n^{-1}) \cdot \mathbf{P}(\mathcal{E}_1 \cap \mathcal{E}_2).$$

We further apply (81) and establish

$$\liminf_{n, p \rightarrow \infty} \mathbf{P} \left( \min_{M \in \mathbb{M}} \|Z^{[m]} - \widehat{Z}\|_2 \leq \text{err}_n(M) \right) \geq 1 - \alpha_0.$$

By the rescaling in (76), we establish (74).

### C.3 Proof of Theorem 3

The proof of Theorem 3 relies on (74) together with the following two lemmas, whose proofs are presented in Section C.4 and Section C.5.

**Lemma 3.** *Define*

$$\hat{\gamma} = \arg \min_{\gamma \in \Delta^L} \gamma^\top \hat{\Gamma} \gamma \quad \text{and} \quad \gamma^* = \arg \min_{\gamma \in \Delta^L} \gamma^\top \Gamma \gamma. \quad (94)$$

If  $\lambda_{\min}(\Gamma) > 0$ , then

$$\|\hat{\gamma} - \gamma^*\|_2 \leq \frac{\|\hat{\Gamma} - \Gamma\|_2}{\lambda_{\min}(\Gamma)} \|\hat{\gamma}\|_2 \leq \frac{\|\hat{\Gamma} - \Gamma\|_F}{\lambda_{\min}(\Gamma)}. \quad (95)$$

**Lemma 4.** *Suppose that  $\Gamma$  is positive semi-definite, then we have*

$$\|\hat{\Gamma}_+ - \Gamma\|_F \leq \|\hat{\Gamma} - \Gamma\|_F.$$

We use  $m^* \in \mathbb{M}$  to denote one index such that  $\|\hat{\Gamma}^{[m^*]} - \Gamma^{\mathbb{Q}}\|_F = \min_{m \in \mathbb{M}} \|\hat{\Gamma}^{[m]} - \Gamma^{\mathbb{Q}}\|_F$ . By (74), with probability larger than  $1 - \alpha_0$ ,

$$\|\hat{\Gamma}^{[m^*]} - \Gamma^{\mathbb{Q}}\|_F \leq \sqrt{2} \text{err}(M) / \sqrt{n}.$$

Then we apply Lemma 3 with  $\hat{\Gamma} = (\hat{\Gamma}^{[m^*]} + \delta \cdot \mathbf{I})_+$  and  $\Gamma = \Gamma^{\mathbb{Q}} + \delta \cdot \mathbf{I}$  and establish

$$\|\hat{\gamma}_\delta^{[m^*]} - \gamma_\delta^*\|_2 \leq \frac{\|(\hat{\Gamma}^{[m^*]} + \delta \cdot \mathbf{I})_+ - (\Gamma^{\mathbb{Q}} + \delta \cdot \mathbf{I})\|_F}{\lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta} \leq \frac{\|\hat{\Gamma}^{[m^*]} - \Gamma^{\mathbb{Q}}\|_F}{\lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta}, \quad (96)$$

where the second inequality follows from Lemma 4. Together with (74), we establish (72).

#### C.4 Proof of Lemma 3

By the definition of  $\gamma^*$  in (94), for any  $t \in (0, 1)$ , we have

$$(\gamma^*)^\top \Gamma \gamma^* \leq [\gamma^* + t(\hat{\gamma} - \gamma^*)]^\top \Gamma [\gamma^* + t(\hat{\gamma} - \gamma^*)],$$

and hence

$$0 \leq 2t(\gamma^*)^\top \Gamma (\hat{\gamma} - \gamma^*) + t^2(\hat{\gamma} - \gamma^*)^\top \Gamma (\hat{\gamma} - \gamma^*).$$

By taking  $t \rightarrow 0+$ , we have

$$(\gamma^*)^\top \Gamma (\hat{\gamma} - \gamma^*) \geq 0. \quad (97)$$

By the definition of  $\hat{\gamma}$  in (94), for any  $t \in (0, 1)$ , we have

$$(\hat{\gamma})^\top \hat{\Gamma} \hat{\gamma} \leq [\hat{\gamma} + t(\gamma^* - \hat{\gamma})]^\top \hat{\Gamma} [\hat{\gamma} + t(\gamma^* - \hat{\gamma})].$$

This gives us

$$2(\gamma^*)^\top \hat{\Gamma} (\gamma^* - \hat{\gamma}) + (t - 2)(\gamma^* - \hat{\gamma})^\top \hat{\Gamma} (\gamma^* - \hat{\gamma}) \geq 0.$$

Since  $2 - t > 0$ , we have

$$(\gamma^* - \hat{\gamma})^\top \hat{\Gamma} (\gamma^* - \hat{\gamma}) \leq \frac{2}{2 - t} (\gamma^*)^\top \hat{\Gamma} (\gamma^* - \hat{\gamma}). \quad (98)$$

It follows from (97) that

$$(\gamma^*)^\top \hat{\Gamma} (\gamma^* - \hat{\gamma}) = (\gamma^*)^\top \Gamma (\gamma^* - \hat{\gamma}) + (\gamma^*)^\top (\hat{\Gamma} - \Gamma) (\gamma^* - \hat{\gamma}) \leq (\gamma^*)^\top (\hat{\Gamma} - \Gamma) (\gamma^* - \hat{\gamma}).$$

Combined with (98), we have

$$(\gamma^* - \hat{\gamma})^\top \hat{\Gamma} (\gamma^* - \hat{\gamma}) \leq \frac{2}{2 - t} (\gamma^*)^\top (\hat{\Gamma} - \Gamma) (\gamma^* - \hat{\gamma}) \leq \frac{2\|\gamma^*\|_2}{2 - t} \|\hat{\Gamma} - \Gamma\|_2 \|\gamma^* - \hat{\gamma}\|_2. \quad (99)$$

Note that the definitions of  $\hat{\gamma}$  and  $\gamma$  are symmetric. Specifically,  $\hat{\gamma}$  is defined as minimizing a quadratic form of  $\hat{\Gamma}$  while  $\gamma^*$  is defined with  $\Gamma$ . We switch the roles of  $\{\hat{\Gamma}, \hat{\gamma}\}$  and  $\{\Gamma, \gamma^*\}$  in (99) and establish

$$(\gamma^* - \hat{\gamma})^\top \Gamma (\gamma^* - \hat{\gamma}) \leq \frac{2\|\hat{\gamma}\|_2}{2 - t} \|\hat{\Gamma} - \Gamma\|_2 \|\gamma^* - \hat{\gamma}\|_2. \quad (100)$$

If  $\lambda_{\min}(\Gamma) > 0$ , we apply the above bound by taking  $t \rightarrow 0+$  and establish (95).

### C.5 Proof of Lemma 4

Write the eigenvalue decomposition of  $\hat{\Gamma}$  as  $\hat{\Gamma} = \sum_{l=1}^L \hat{A}_l u_l u_l^\top$ . Define

$$\hat{\Gamma}_+ = \sum_{l=1}^L \max\{\hat{A}_l, 0\} u_l u_l^\top \quad \text{and} \quad \hat{\Gamma}_- = \sum_{l=1}^L -\min\{\hat{A}_l, 0\} u_l u_l^\top.$$

We have

$$\hat{\Gamma} = \hat{\Gamma}_+ - \hat{\Gamma}_- \quad \text{with} \quad \text{Tr}(\hat{\Gamma}_+^\top \hat{\Gamma}_-) = 0. \quad (101)$$

Since  $\Gamma$  is positive semi-definite, we have

$$\text{Tr}(\Gamma^\top \hat{\Gamma}_-) = \sum_{l=1}^L -\min\{\hat{A}_{ll}, 0\} \text{Tr}(\Gamma^\top u_l u_l^\top) \geq 0. \quad (102)$$

We apply (101) and (102) and establish

$$\|\hat{\Gamma} - \Gamma\|_F^2 = \|\hat{\Gamma}_+ - \Gamma - \hat{\Gamma}_-\|_F^2 = \|\hat{\Gamma}_+ - \Gamma\|_F^2 + \|\hat{\Gamma}_-\|_F^2 - 2\text{Tr}(\hat{\Gamma}_+^\top \hat{\Gamma}_-) + 2\text{Tr}(\Gamma^\top \hat{\Gamma}_-) \geq \|\hat{\Gamma}_+ - \Gamma\|_F^2.$$

## D Proofs of Lemma 1, Corollary 1 and Propositions 5, 4, and 6

### D.1 Proof of Lemma 1

We start with the proof of Lemma 1 for the covariate shift setting, which relies on Propositions 4 and 5. On the event  $\mathcal{E}_1$  defined in (79), we apply the definitions in (77) and have

$$\|\hat{\mathbf{V}} - \mathbf{V}\|_2 \leq \frac{d_0}{3n}. \quad (103)$$

Note that

$$\begin{aligned} \frac{|\hat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}}|}{\sqrt{\hat{\mathbf{V}}_{\pi(l,k),\pi(l,k)} + d_0/n}} &\leq \frac{|\hat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}}|}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)} + 2d_0/3n}} \\ &\leq \frac{|D_{l,k}|}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)} + 2d_0/3n}} + \frac{|\text{Rem}_{l,k}|}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)} + 2d_0/3n}}, \end{aligned} \quad (104)$$

where  $D_{l,k}$  and  $\text{Rem}_{l,k}$  are defined in Proposition 4.

It follows from (53) that

$$\liminf_{n,p \rightarrow \infty} \mathbf{P} \left( \max_{1 \leq l,k \leq L} \frac{|D_{l,k}|}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)} + (2d_0/3n)}} \leq z_{\alpha_0/[L(L+1)]} \right) \geq 1 - \alpha_0. \quad (105)$$

Combining (54) and (55), we apply the boundedness on  $\max_{1 \leq l \leq L} \|b^{(l)}\|_2$  and establish that, with probability larger than  $1 - \min\{n, p\}^{-c}$ ,

$$\frac{|\text{Rem}_{l,k}|}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)} + 2d_0/(3n)}} \lesssim \frac{s \log p}{\sqrt{n}} + \sqrt{\frac{s(\log p)^2}{N_{\mathbb{Q}}}}.$$

By the condition  $\sqrt{n} \gg s \log p$  and  $N_{\mathbb{Q}} \gg s(\log p)^2$ , we then establish that, with probability larger than  $1 - \min\{n, p\}^{-c}$ ,

$$\max_{1 \leq l, k \leq L} \frac{|\text{Rem}_{l,k}|}{\sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)} + 3d_0/(3n)}} \leq 0.05 \cdot z_{\alpha_0/[L(L+1)]}$$

Combined with (104) and (105), we establish (73).

The proof of Lemma 1 for the no covariate shift setting relies on Proposition 6. The specific argument is the same as that for the covariate shift setting.

## D.2 Proof of Corollary 1

By Proposition 5, under Condition (A1),  $s \log p \ll n$  and  $N_{\mathbb{Q}} \gtrsim \max\{n, p\}$ , we show that  $\|n\mathbf{V}\|_{\infty}$  and  $d_0$  are bounded from above with probability larger than  $1 - \min\{n, p\}^{-c}$  for some positive constant  $c > 0$ . Furthermore, since  $L$  is finite,  $n\lambda_{\max}(V) \lesssim \|n\mathbf{V}\|_{\infty}$ . As a consequence,  $c^*(\alpha_0) \geq c'$  for a positive constant  $c' > 0$ . The lower bound for  $C^*(L, \alpha_0)$  follows from the boundedness on  $L$  and the lower bound for  $c^*(\alpha_0)$ .

## D.3 High probability events

We introduce the following events to facilitate the proofs of Propositions 4, 5 and 6.

$$\begin{aligned} \mathcal{G}_0 &= \left\{ \left\| \frac{1}{n_l} [X^{(l)}]^\top \epsilon^{(l)} \right\|_{\infty} \lesssim \sqrt{\frac{\log p}{n_l}} \quad \text{for } 1 \leq l \leq L \right\}, \\ \mathcal{G}_1 &= \left\{ \max \left\{ \|\widehat{b}_{init}^{(l)} - b^{(l)}\|_2, \frac{1}{\sqrt{n_l}} \|X^{(l)}(\widehat{b}_{init}^{(l)} - b^{(l)})\|_2 \right\} \lesssim \sqrt{\|b^{(l)}\|_0 \frac{\log p}{n_l}} \sigma_l \quad \text{for } 1 \leq l \leq L \right\}, \\ \mathcal{G}_2 &= \left\{ \|\widehat{b}_{init}^{(l)} - b^{(l)}\|_1 \lesssim \|b^{(l)}\|_0 \sqrt{\frac{\log p}{n_l}} \sigma_l, \|\widehat{b}_{init}^{(l)} - b^{(l)}\|_{\mathcal{S}_l^c} \leq C \|\widehat{b}_{init}^{(l)} - b^{(l)}\|_{\mathcal{S}_l} \quad \text{for } 1 \leq l \leq L \right\}, \\ \mathcal{G}_3 &= \left\{ |\widehat{\sigma}_l^2 - \sigma_l^2| \lesssim \|b^{(l)}\|_0 \frac{\log p}{n_l} + \sqrt{\frac{\log p}{n_l}} \quad \text{for } 1 \leq l \leq L \right\}, \end{aligned} \tag{106}$$

where  $\mathcal{S}_l \subset [p]$  denotes the support of  $b^{(l)}$  for  $1 \leq l \leq L$  and  $C > 0$  is a positive constant. Recall that, for  $1 \leq l \leq L$ ,  $\widehat{b}_{init}^{(l)}$  is the Lasso estimator defined in (43) with  $\lambda_l = \sqrt{(2+c) \log p / |A_l|} \sigma_l$  for some constant  $c > 0$ ;  $\widehat{\sigma}_l^2 = \|Y^{(l)} - X^{(l)} \widehat{b}^{(l)}\|_2^2 / n_l$  for  $1 \leq l \leq L$  with  $\widehat{b}^{(l)}$  denoting the Lasso estimator based on the non-split data.

We further define the following events,

$$\begin{aligned}
\mathcal{G}_4 &= \left\{ \|\tilde{\Sigma}^{\mathbb{Q}} - \Sigma^{\mathbb{Q}}\|_2 \lesssim \sqrt{\frac{p}{N_{\mathbb{Q}}}} + \frac{p}{N_{\mathbb{Q}}} \right\}, \\
\mathcal{G}_5 &= \left\{ \max_{\mathcal{S} \subset [p], |\mathcal{S}| \leq s} \max_{\|w_{\mathcal{S}^c}\|_1 \leq C\|w_{\mathcal{S}}\|_1} \left| \frac{w^{\top} \left( \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^{\top} \right) w}{w^{\top} E(X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]) w} - 1 \right| \lesssim \frac{s \log p}{N_{\mathbb{Q}}} \right\}, \\
\mathcal{G}_6(w, v, t) &= \left\{ \left| w^{\top} (\hat{\Sigma}^{\mathbb{Q}} - \Sigma^{\mathbb{Q}}) v \right| + \left| w^{\top} (\tilde{\Sigma}^{\mathbb{Q}} - \Sigma^{\mathbb{Q}}) v \right| \lesssim t \frac{\|(\Sigma^{\mathbb{Q}})^{1/2} w\|_2 \|(\Sigma^{\mathbb{Q}})^{1/2} v\|_2}{\sqrt{N_{\mathbb{Q}}}} \right\},
\end{aligned} \tag{107}$$

where  $\tilde{\Sigma}^{\mathbb{Q}} = \frac{1}{|A|} \sum_{i \in A} X_{i,\cdot}^{\mathbb{Q}} (X_{i,\cdot}^{\mathbb{Q}})^{\top}$ ,  $\hat{\Sigma}^{\mathbb{Q}} = \frac{1}{|B|} \sum_{i \in B} X_{i,\cdot}^{\mathbb{Q}} (X_{i,\cdot}^{\mathbb{Q}})^{\top}$  and  $t > 0$  is any positive constant and  $w, v \in \mathbb{R}^p$  are pre-specified vectors.

**Lemma 5.** *Suppose that Condition (A1) holds and  $s \lesssim n/\log p$ , then*

$$\mathbf{P}(\cap_{j=0}^3 \mathcal{G}_j) \geq 1 - \min\{n, p\}^{-c}, \tag{108}$$

$$\mathbf{P}(\mathcal{G}_4 \cap \mathcal{G}_5) \geq 1 - p^{-c}, \tag{109}$$

$$\mathbf{P}(\mathcal{G}_6(w, v, t)) \geq 1 - 2 \exp(-ct^2), \tag{110}$$

for some positive constant  $c > 0$ .

The above high-probability statement (108) follows from the existing literature results on the analysis of Lasso estimators and we shall point to the exact literature results. Specifically, the control of the probability of  $\mathcal{G}_0$  follows from Lemma 6.2 of [Bühlmann and van de Geer \(2011\)](#). Regarding the events  $\mathcal{G}_1$  and  $\mathcal{G}_2$ , the control of  $\|\hat{b}_{init}^{(l)} - b^{(l)}\|_1$ ,  $\|\hat{b}_{init}^{(l)} - b^{(l)}\|_2$  and  $\frac{1}{\sqrt{n_l}} \|X^{(l)}(\hat{b}_{init}^{(l)} - b^{(l)})\|_2$  can be found in Theorem 3 of [Ye and Zhang \(2010\)](#), Theorem 7.2 of [Bickel et al. \(2009\)](#) or Theorem 6.1 of [Bühlmann and van de Geer \(2011\)](#); the control of  $\|[\hat{b}_{init}^{(l)} - b^{(l)}]_{\mathcal{S}_l^c}\|_1 \leq C\|[\hat{b}_{init}^{(l)} - b^{(l)}]_{\mathcal{S}_l}\|_1$  can be found in Corollary B.2 of [Bickel et al. \(2009\)](#) or Lemma 6.3 of [Bühlmann and van de Geer \(2011\)](#). For the event  $\mathcal{G}_3$ , its probability can be controlled as Theorem 2 or (20) in [Sun and Zhang \(2012\)](#).

If  $X_i^{\mathbb{Q}}$  is sub-gaussian, it follows from equation (5.26) of [Vershynin \(2012\)](#) that the event  $\mathcal{G}_4$  holds with a probability larger than  $1 - \exp(-cp)$  for some positive constant  $c > 0$ ; it follows from Theorem 1.6 of [Zhou \(2009\)](#) that the event  $\mathcal{G}_5$  holds with a probability larger than  $1 - p^{-c}$  for some positive constant  $c > 0$ . The proof of (110) follows from Lemma 10 in the supplement of [Cai and Guo \(2020\)](#).



#### D.4 Proof of Proposition 5

We have the expression for the diagonal element of  $\mathbf{V}$  as

$$\begin{aligned} \mathbf{V}_{\pi(l,k),\pi(l,k)} &= \frac{\sigma_l^2}{|B_l|} (\widehat{u}^{(l,k)})^\top \widehat{\Sigma}^{(l)} [\widehat{u}^{(l,k)} + \widehat{u}^{(k,l)} \mathbf{1}(k=l)] + \frac{\sigma_k^2}{|B_k|} (\widehat{u}^{(k,l)})^\top \widehat{\Sigma}^{(k)} [\widehat{u}^{(l,k)} \mathbf{1}(l=k) + \widehat{u}^{(k,l)}] \\ &\quad + \frac{1}{|B|} (\mathbf{E}[b^{(l)}]^\top X_i^\mathbb{Q} [b^{(k)}]^\top X_i^\mathbb{Q} [b^{(l)}]^\top X_i^\mathbb{Q} [b^{(k)}]^\top X_i^\mathbb{Q} - (b^{(l)})^\top \Sigma^\mathbb{Q} b^{(k)} (b^{(l)})^\top \Sigma^\mathbb{Q} b^{(k)}). \end{aligned}$$

We introduce the following lemma, which restates Lemma 1 of [Cai et al. \(2021\)](#) in the current paper's terminology.

**Lemma 6.** *Suppose that Condition (A1) holds, then with probability larger than  $1 - p^{-c}$ ,*

$$\begin{aligned} c \frac{\|\omega^{(k)}\|_2^2}{n_l} &\leq \frac{1}{|B_l|} (\widehat{u}^{(l,k)})^\top \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)} \leq C \frac{\|\omega^{(k)}\|_2^2}{n_l}, \quad \text{for } 1 \leq l, k \leq L, \\ c \frac{\|\omega^{(l)}\|_2^2}{n_k} &\leq \frac{1}{|B_k|} (\widehat{u}^{(k,l)})^\top \widehat{\Sigma}^{(k)} \widehat{u}^{(k,l)} \leq C \frac{\|\omega^{(l)}\|_2^2}{n_k}, \quad \text{for } 1 \leq l, k \leq L, \end{aligned}$$

for some positive constants  $C > c > 0$ .

The bounds for  $\mathbf{V}_{\pi(l,k),\pi(l,k)}^{(a)}$  in (55) follow from Lemma 6. Since  $X_i^\mathbb{Q}$  is sub-gaussian, we have

$$|\mathbf{E}[b^{(l_1)}]^\top X_i^\mathbb{Q} [b^{(k_1)}]^\top X_i^\mathbb{Q} [b^{(l_2)}]^\top X_i^\mathbb{Q} [b^{(k_2)}]^\top X_i^\mathbb{Q}| \lesssim \|b^{(l_1)}\|_2 \|b^{(k_1)}\|_2 \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2, \quad (111)$$

and

$$(b^{(l_1)})^\top \Sigma^\mathbb{Q} b^{(k_1)} (b^{(l_2)})^\top \Sigma^\mathbb{Q} b^{(k_2)} \lesssim \|b^{(l_1)}\|_2 \|b^{(k_1)}\|_2 \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2. \quad (112)$$

We establish the upper bound for  $\mathbf{V}_{\pi(l,k),\pi(l,k)}^{(b)}$  in (55) by taking  $l_1 = l_2 = l$  and  $k_1 = k_2 = k$ .

For the setting of known  $\Sigma^\mathbb{Q}$ , on the event  $\mathcal{G}_2$  defined in (106), we establish

$$\|\omega^{(k)}\|_2 = \|\Sigma^\mathbb{Q} \widehat{b}_{init}^{(k)}\|_2 \lesssim \lambda_{\max}(\Sigma^\mathbb{Q}) \|\widehat{b}_{init}^{(k)}\|_2 \lesssim \lambda_{\max}(\Sigma^\mathbb{Q}) \left( \|b^{(k)}\|_2 + \sqrt{s \log p/n} \right).$$

To establish (56), we control  $\|\omega^{(k)}\|_2 = \|\widetilde{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)}\|_2$  as follows,

$$\|\widetilde{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)}\|_2 \leq \|\Sigma^\mathbb{Q} \widehat{b}_{init}^{(k)}\|_2 + \|(\widetilde{\Sigma}^\mathbb{Q} - \Sigma^\mathbb{Q}) \widehat{b}_{init}^{(k)}\|_2 \leq \lambda_{\max}(\Sigma^\mathbb{Q}) \|\widehat{b}_{init}^{(k)}\|_2 + \|\widetilde{\Sigma}^\mathbb{Q} - \Sigma^\mathbb{Q}\|_2 \|\widehat{b}_{init}^{(k)}\|_2.$$

On the event  $\mathcal{G}_4$ , we establish

$$\|\omega^{(k)}\|_2 \lesssim \lambda_{\max}(\Sigma^{\mathbb{Q}}) \left(1 + \sqrt{\frac{p}{N_{\mathbb{Q}}}} + \frac{p}{N_{\mathbb{Q}}}\right) \left(\|b^{(k)}\|_2 + \sqrt{s \log p/n}\right).$$

With a similar bound for  $\|\omega^{(l)}\|_2$ , we establish (56).

## D.5 Proof of Proposition 4

We decompose the error  $\widehat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}}$  as

$$\begin{aligned} \widehat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}} &= \frac{1}{|B_l|} (\widehat{u}^{(l,k)})^{\top} [X_{B_l}^{(l)}]^{\top} \epsilon_{B_l}^{(l)} + \frac{1}{|B_k|} (\widehat{u}^{(k,l)})^{\top} [X_{B_k}^{(k)}]^{\top} \epsilon_{B_k}^{(k)} \\ &\quad + (b^{(l)})^{\top} (\widehat{\Sigma}^{\mathbb{Q}} - \Sigma^{\mathbb{Q}}) b^{(k)} - (\widehat{b}_{init}^{(l)} - b^{(l)})^{\top} \widehat{\Sigma}^{\mathbb{Q}} (\widehat{b}_{init}^{(k)} - b^{(k)}) \\ &\quad + (\widehat{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^{\top} (\widehat{b}_{init}^{(l)} - b^{(l)}) + (\widehat{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(l)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(k,l)})^{\top} (\widehat{b}_{init}^{(k)} - b^{(k)}). \end{aligned} \quad (113)$$

We define  $D_{l,k} = D_{l,k}^{(a)} + D_{l,k}^{(b)}$  with

$$D_{l,k}^{(a)} = \frac{1}{|B_l|} (\widehat{u}^{(l,k)})^{\top} [X_{B_l}^{(l)}]^{\top} \epsilon_{B_l}^{(l)} + \frac{1}{|B_k|} (\widehat{u}^{(k,l)})^{\top} [X_{B_k}^{(k)}]^{\top} \epsilon_{B_k}^{(k)},$$

and

$$D_{l,k}^{(b)} = (b^{(l)})^{\top} (\widehat{\Sigma}^{\mathbb{Q}} - \Sigma^{\mathbb{Q}}) b^{(k)}.$$

Define  $\text{Rem}_{l,k}$  as

$$\begin{aligned} \text{Rem}_{l,k} &= -(\widehat{b}_{init}^{(l)} - b^{(l)})^{\top} \widehat{\Sigma}^{\mathbb{Q}} (\widehat{b}_{init}^{(k)} - b^{(k)}) + (\widehat{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^{\top} (\widehat{b}_{init}^{(l)} - b^{(l)}) \\ &\quad + (\widehat{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(l)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(k,l)})^{\top} (\widehat{b}_{init}^{(k)} - b^{(k)}). \end{aligned}$$

By (113), we have

$$\widehat{\Gamma}_{l,k}^{\mathbb{Q}} - \Gamma_{l,k}^{\mathbb{Q}} = D_{l,k} + \text{Rem}_{l,k}.$$

Note that  $D_{l,k}^{(a)}$  is a function of  $X_A^{\mathbb{Q}}$ ,  $\{X^{(l)}\}_{1 \leq l \leq L}$  and  $\{\epsilon^{(k)}\}_{1 \leq l \leq L}$  and  $D_{l,k}^{(b)}$  is a function of the sub-sample  $X_B^{\mathbb{Q}}$  and hence  $D_{l,k}^{(a)}$  is independent of  $D_{l,k}^{(b)}$ .

**Limiting distribution of  $D_{l,k}$ .** We check the Lindeberg's condition and establish the asymptotic normality of  $D_{l,k}$ . In the following, we focus on the setting  $l \neq k$  and the proof

can be extended to the setting  $l = k$ . We write

$$\begin{aligned}\frac{D_{l,k}^{(a)}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} &= \frac{1}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \left( (\hat{u}^{(l,k)})^\top \frac{1}{|B_l|} \sum_{i \in B_l} X_i^{(l)} \epsilon_i^{(l)} + (\hat{u}^{(k,l)})^\top \frac{1}{|B_k|} \sum_{i \in B_k} X_i^{(k)} \epsilon_i^{(k)} \right), \\ \frac{D_{l,k}^{(b)}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} &= \frac{1}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \frac{1}{|B|} \sum_{i \in B} [b^{(l)}]^\top (X_i^\mathbb{Q} (X_i^\mathbb{Q})^\top - \Sigma^\mathbb{Q}) b^{(k)}.\end{aligned}$$

Define

$$\begin{aligned}W_{l,i} &= \frac{1}{|B_l| \sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} (\hat{u}^{(l,k)})^\top X_i^{(l)} \epsilon_i^{(l)} \quad \text{for } i \in B_l, \\ W_{k,i} &= \frac{1}{|B_k| \sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} (\hat{u}^{(k,l)})^\top X_i^{(k)} \epsilon_i^{(k)} \quad \text{for } i \in B_k,\end{aligned}$$

and

$$W_{\mathbb{Q},i} = \frac{1}{|B| \sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} [b^{(l)}]^\top (X_i^\mathbb{Q} (X_i^\mathbb{Q})^\top - \Sigma^\mathbb{Q}) b^{(k)} \quad \text{for } i \in B.$$

Then we have

$$\frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} = \sum_{i \in B_l} W_{l,i} + \sum_{i \in B_k} W_{k,i} + \sum_{i \in B} W_{\mathbb{Q},i}.$$

We use  $\mathcal{O}_1 = \{X_i^\mathbb{Q}\}_{i \in A} \cup \{X_{A_l}^{(l)}, Y_{A_l}^{(l)}\}_{1 \leq l \leq L}$  to denote the subset of data used for computing  $\hat{u}^{(l,k)}$  and  $\hat{u}^{(k,l)}$ . Conditioning on  $\mathcal{O}_1$ , then  $\{W_{l,i}\}_{i \in B_l}$ ,  $\{W_{k,i}\}_{i \in B_k}$  and  $\{W_{\mathbb{Q},i}\}_{i \in B}$  are independent random variables. Note that  $\mathbf{E}(W_{l,i} \mid \mathcal{O}_1) = 0$  for  $i \in B_l$ ,  $\mathbf{E}(W_{k,i} \mid \mathcal{O}_1) = 0$  for  $i \in B_k$  and  $\mathbf{E}(W_{\mathbb{Q},i} \mid \mathcal{O}_1) = 0$  for  $i \in B$ . Furthermore, we have

$$\sum_{i \in B_l} \mathbf{E}(W_{l,i}^2 \mid \mathcal{O}_1) + \sum_{i \in B_k} \mathbf{E}(W_{k,i}^2 \mid \mathcal{O}_1) + \sum_{i \in B} \mathbf{E}(W_{\mathbb{Q},i}^2 \mid \mathcal{O}_1) = 1. \quad (114)$$

Define the event

$$\mathcal{E}_3 = \left\{ \frac{\|\omega^{(l)}\|_2^2}{n_k} + \frac{\|\omega^{(k)}\|_2^2}{n_l} \lesssim \mathbf{V}_{\pi(l,k),\pi(l,k)}^{(a)} \right\},$$

and it follows from Proposition 5 that

$$\mathbf{P}(\mathcal{O}_1 \in \mathcal{E}_3) \geq 1 - \min\{n, p\}^{-c}. \quad (115)$$

Let  $o_1$  denote one element of the event  $\mathcal{E}_3$ . To establish the asymptotic normality, it is sufficient to check the following Lindeberg's condition: for any constant  $c > 0$ ,

$$\begin{aligned} & \sum_{i \in B_l} \mathbf{E} (W_{l,i}^2 \mathbf{1}\{|W_{l,i}| \geq c\} \mid \mathcal{O}_1 = o_1) + \sum_{i \in B_k} \mathbf{E} (W_{k,i}^2 \mathbf{1}\{|W_{k,i}| \geq c\} \mid \mathcal{O}_1 = o_1) \\ & + \sum_{i \in B} \mathbf{E} (W_{\mathbb{Q},i}^2 \mathbf{1}\{|W_{\mathbb{Q},i}| \geq c\} \mid \mathcal{O}_1 = o_1) \rightarrow 0. \end{aligned} \quad (116)$$

We apply the optimization constraint in (46) and establish

$$\begin{aligned} & \sum_{i \in B_l} \mathbf{E} (W_{l,i}^2 \mathbf{1}\{|W_{l,i}| \geq c\} \mid \mathcal{O}_1) + \sum_{i \in B_k} \mathbf{E} (W_{k,i}^2 \mathbf{1}\{|W_{k,i}| \geq c\} \mid \mathcal{O}_1) \\ & \leq \sum_{i \in B_l} \frac{[\sigma_l^2 (\hat{u}^{(l,k)})^\top X_i^{(l)}]^2}{|B_l|^2 \mathbf{V}_{\pi(l,k), \pi(l,k)}} \mathbf{E} \left( \frac{(\epsilon_i^{(l)})^2}{\sigma_l^2} \mathbf{1} \left\{ \left| \epsilon_i^{(l)} \right| \geq \frac{c|B_l| \sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}}}{\|\omega^{(k)}\|_2 \sqrt{\log |B_l|}} \right\} \mid \mathcal{O}_1 \right) \\ & + \sum_{i \in B_k} \frac{[\sigma_k^2 (\hat{u}^{(k,l)})^\top X_i^{(k)}]^2}{|B_k|^2 \mathbf{V}_{\pi(l,k), \pi(l,k)}} \mathbf{E} \left( \frac{(\epsilon_i^{(k)})^2}{\sigma_k^2} \mathbf{1} \left\{ \left| \epsilon_i^{(k)} \right| \geq \frac{c|B_k| \sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}}}{\|\omega^{(l)}\|_2 \sqrt{\log |B_k|}} \right\} \mid \mathcal{O}_1 \right) \\ & \leq \sum_{i \in B_l} \frac{[\sigma_l^2 (\hat{u}^{(l,k)})^\top X_i^{(l)}]^2}{|B_l|^2 \mathbf{V}_{\pi(l,k), \pi(l,k)}} \mathbf{E} \left( \frac{(\epsilon_i^{(l)})^2}{\sigma_l^2} \mathbf{1} \left\{ \left| \epsilon_i^{(l)} \right| \geq \frac{c|B_l| \sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(a)}}}{\|\omega^{(k)}\|_2 \sqrt{\log |B_l|}} \right\} \mid \mathcal{O}_1 \right) \\ & + \sum_{i \in B_k} \frac{[\sigma_k^2 (\hat{u}^{(k,l)})^\top X_i^{(k)}]^2}{|B_k|^2 \mathbf{V}_{\pi(l,k), \pi(l,k)}} \mathbf{E} \left( \frac{(\epsilon_i^{(k)})^2}{\sigma_k^2} \mathbf{1} \left\{ \left| \epsilon_i^{(k)} \right| \geq \frac{c|B_k| \sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(a)}}}{\|\omega^{(l)}\|_2 \sqrt{\log |B_k|}} \right\} \mid \mathcal{O}_1 \right) \end{aligned}$$

where the last inequality follows from  $\mathbf{V}_{\pi(l,k), \pi(l,k)} \geq \mathbf{V}_{\pi(l,k), \pi(l,k)}^{(a)}$ .

By the definition of  $\mathcal{E}_3$ , we condition on  $\mathcal{O}_1 = o_1$  with  $o_1 \in \mathcal{E}_3$  and further upper bound the right hand side of the above inequality by

$$\begin{aligned} & \max_{1 \leq i \leq n_l} \mathbf{E} \left( \frac{(\epsilon_i^{(l)})^2}{\sigma_l^2} \mathbf{1} \left\{ \left| \epsilon_i^{(l)} \right| \geq \frac{c|B_l| \sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(a)}}}{\|\omega^{(k)}\|_2 \sqrt{\log |B_l|}} \right\} \mid \mathcal{O}_1 = o_1 \right) \\ & + \max_{1 \leq i \leq n_k} \mathbf{E} \left( \frac{(\epsilon_i^{(k)})^2}{\sigma_k^2} \mathbf{1} \left\{ \left| \epsilon_i^{(k)} \right| \geq \frac{c|B_k| \sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(a)}}}{\|\omega^{(l)}\|_2 \sqrt{\log |B_k|}} \right\} \mid \mathcal{O}_1 = o_1 \right) \\ & \leq \max_{1 \leq i \leq n_l} \mathbf{E} \left( \frac{(\epsilon_i^{(l)})^2}{\sigma_l^2} \mathbf{1} \left\{ \left| \epsilon_i^{(l)} \right| \geq \frac{c|B_l| \sqrt{\frac{\|\omega^{(l)}\|_2^2}{n_k} + \frac{\|\omega^{(k)}\|_2^2}{n_l}}}{\|\omega^{(k)}\|_2 \sqrt{\log |B_l|}} \right\} \mid \mathcal{O}_1 = o_1 \right) \\ & + \max_{1 \leq i \leq n_k} \mathbf{E} \left( \frac{(\epsilon_i^{(k)})^2}{\sigma_k^2} \mathbf{1} \left\{ \left| \epsilon_i^{(k)} \right| \geq \frac{c|B_k| \sqrt{\frac{\|\omega^{(l)}\|_2^2}{n_k} + \frac{\|\omega^{(k)}\|_2^2}{n_l}}}{\|\omega^{(l)}\|_2 \sqrt{\log |B_k|}} \right\} \mid \mathcal{O}_1 = o_1 \right) \lesssim \left( \frac{\log n}{n} \right)^{\frac{c}{2}}, \end{aligned} \quad (117)$$

where the last inequality follows from  $\mathbf{E}([\epsilon_i^{(l)}]_i^{2+c} \mid X_i^{(l)}) \leq C$  in Condition (A1).

Define

$$J_i = [b^{(l)}]^\top (X_i^\mathbb{Q} (X_i^\mathbb{Q})^\top - \Sigma^\mathbb{Q}) b^{(k)},$$

and

$$\bar{W}_{\mathbb{Q},i} = \frac{1}{|B| \sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}^{(b)}}} [b^{(l)}]^\top (X_i^\mathbb{Q} (X_i^\mathbb{Q})^\top - \Sigma^\mathbb{Q}) b^{(k)} = \frac{J_i}{\sqrt{|B| \text{Var}(J_i)}} \quad \text{for } i \in B.$$

Note that  $|W_{\mathbb{Q},i}| \leq |\bar{W}_{\mathbb{Q},i}|$ , and then we have

$$\begin{aligned} \sum_{i \in B} \mathbf{E} (W_{\mathbb{Q},i}^2 \mathbf{1}\{|W_{\mathbb{Q},i}| \geq c\} \mid \mathcal{O}_1) &= \sum_{i \in B} \mathbf{E} (W_{\mathbb{Q},i}^2 \mathbf{1}\{|\bar{W}_{\mathbb{Q},i}| \geq c\}) \\ &\leq \sum_{i \in B} \mathbf{E} (\bar{W}_{\mathbb{Q},i}^2 \mathbf{1}\{|\bar{W}_{\mathbb{Q},i}| \geq c\}) \\ &\leq \mathbf{E} (J_i^2 / \text{Var}(J_i)) \cdot \mathbf{1} \left( |J_i| / \sqrt{\text{Var}(J_i)} \geq c \sqrt{|B|} \right). \end{aligned}$$

Together with the dominated convergence theorem, we have

$$\sum_{i \in B} \mathbf{E} (W_{\mathbb{Q},i}^2 \mathbf{1}\{|W_{\mathbb{Q},i}| \geq c\} \mid \mathcal{O}_1) \rightarrow 0.$$

Combined with (117), we establish (116). Hence, for  $o_1 \in \mathcal{E}_3$ , we establish

$$\frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \mid \mathcal{O}_1 = o_1 \xrightarrow{d} \mathcal{N}(0, 1).$$

We calculate its characteristic function

$$\begin{aligned} &\mathbf{E} \exp \left( it \cdot \frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \right) - e^{-t^2/2} \\ &= \int \mathbf{E} \left( \left[ \exp \left( it \cdot \frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \right) \mid \mathcal{O}_1 = o_1 \right] - e^{-t^2/2} \right) \cdot \mathbf{1}_{o_1 \in \mathcal{E}_3} \cdot \mu(o_1) \\ &+ \int \mathbf{E} \left[ \exp \left( it \cdot \frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \right) \mid \mathcal{O}_1 = o_1 \right] \cdot \mathbf{1}_{o_1 \notin \mathcal{E}_3} \cdot \mu(o_1) - e^{-t^2/2} \cdot \mathbf{P}(\mathcal{E}_3^c). \end{aligned}$$

Combined with (115) and the bounded convergence theorem, we establish

$$\frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k),\pi(l,k)}}} \xrightarrow{d} \mathcal{N}(0, 1).$$

**Control of  $\text{Rem}_{l,k}$  in (54).** We introduce the following lemma, whose proof is presented in Section G.1.

**Lemma 7.** *Suppose that Condition (A1) holds, then with probability larger than  $1 - \min\{n, p\}^{-c}$ , we have*

$$\left| (\widehat{b}_{init}^{(l)} - b^{(l)})^\top \widehat{\Sigma}^\mathbb{Q} (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| \lesssim \sqrt{\frac{\|b^{(l)}\|_0 \|b^{(k)}\|_0 (\log p)^2}{n_l n_k}}, \quad (118)$$

$$\left| (\widehat{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| \lesssim \|\omega^{(k)}\|_2 \frac{\|b^{(l)}\|_0 \log p}{n_l} + \|\widehat{b}_{init}^{(k)}\|_2 \sqrt{\frac{\|b^{(l)}\|_0 (\log p)^2}{n_l N_\mathbb{Q}}}, \quad (119)$$

$$\left| (\widehat{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(l)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^\top \widehat{\Sigma}^\mathbb{Q} (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| \lesssim \|\omega^{(l)}\|_2 \frac{\|b^{(k)}\|_0 \log p}{n_k} + \|\widehat{b}_{init}^{(l)}\|_2 \sqrt{\frac{\|b^{(k)}\|_0 (\log p)^2}{n_k N_\mathbb{Q}}}. \quad (120)$$

On the event  $\mathcal{G}_1$ , we have

$$\|\widehat{b}_{init}^{(l)}\|_2 \leq \|b^{(l)}\|_2 + \sqrt{\frac{\|b^{(k)}\|_0 \log p}{n_k}}.$$

Combining this inequality with Lemma 7, we establish the upper bound for  $\text{Rem}_{l,k}$  in (54).

## D.6 Proof of Proposition 6

Define  $N = \sum_{l=1}^L n_l + N_\mathbb{Q}$ ,  $\widetilde{\Sigma}^{(l)} = \frac{1}{n_l} \sum_{i=1}^{n_l} X_i^{(l)} [X_i^{(l)}]^\top$  for  $1 \leq l \leq L$  and

$$\widehat{\Sigma} = \frac{1}{\sum_{l=1}^L n_l + N_\mathbb{Q}} \left( \sum_{l=1}^L \sum_{i=1}^{n_l} X_i^{(l)} [X_i^{(l)}]^\top + \sum_{i=1}^{N_\mathbb{Q}} X_i^{(l)} [X_i^{(l)}]^\top \right).$$

The difference between the proposed estimator  $\widehat{\Gamma}_{l,k}$  and  $\Gamma_{l,k}$  is

$$\widehat{\Gamma}_{l,k} - \Gamma_{l,k} = D_{l,k}^{(1)} + D_{l,k}^{(2)} + \text{Rem}_{l,k},$$

where

$$\begin{aligned} D_{l,k}^{(1)} &= \frac{1}{n_k} (b^{(l)})^\top [X^{(k)}]^\top \epsilon^{(k)} + \frac{n_k}{N} (b^{(l)})^\top \left( \widetilde{\Sigma}^{(k)} - \Sigma \right) b^{(k)} \\ &\quad + \frac{1}{n_l} (b^{(k)})^\top [X^{(l)}]^\top \epsilon^{(l)} + \frac{n_l}{N} (b^{(l)})^\top \left( \widetilde{\Sigma}^{(l)} - \Sigma \right) b^{(k)}, \\ D_{l,k}^{(2)} &= (b^{(l)})^\top \left( \widehat{\Sigma} - \frac{n_l}{N} \widetilde{\Sigma}^{(l)} - \frac{n_k}{N} \widetilde{\Sigma}^{(k)} - \frac{N - n_l - n_k}{N} \Sigma \right) b^{(k)}, \end{aligned}$$

and

$$\begin{aligned} \text{Rem}_{l,k} &= \frac{1}{n_l} (\widehat{b}_{init}^{(k)} - b^{(k)})^\top [X^{(l)}]^\top \epsilon^{(l)} + \frac{1}{n_k} (\widehat{b}_{init}^{(l)} - b^{(l)})^\top [X^{(k)}]^\top \epsilon^{(k)} \\ &\quad - (\widehat{b}_{init}^{(l)} - b^{(l)})^\top \widehat{\Sigma} (\widehat{b}_{init}^{(k)} - b^{(k)}) + [\widehat{b}_{init}^{(l)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(k)}) (\widehat{b}_{init}^{(k)} - b^{(k)}) + [\widehat{b}_{init}^{(k)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(l)}) (\widehat{b}_{init}^{(l)} - b^{(l)}). \end{aligned} \quad (121)$$

In the following, we shall control  $D^{(1)}$ ,  $D^{(2)}$  and  $\text{Rem}_{l,k}$  separately. Note that  $D_{l,k}^{(1)}$  and  $D_{l,k}^{(2)}$  are independent. We write

$$D_{l,k}^{(1)} = \sum_{i=1}^{n_k} W_{k,i}^{(1)} + \sum_{i=1}^{n_l} W_{l,i}^{(1)}$$

with

$$W_{k,i}^{(1)} = \frac{1}{n_k} \left[ (b^{(l)})^\top X_i^{(k)} \epsilon_i^{(k)} + \frac{n_k}{N} (b^{(l)})^\top \left( X_i^{(k)} (X_i^{(k)})^\top - \Sigma \right) b^{(k)} \right] \quad \text{for } 1 \leq i \leq n_k,$$

and

$$W_{l,i}^{(1)} = \frac{1}{n_l} \left[ (b^{(k)})^\top X_i^{(l)} \epsilon_i^{(l)} + \frac{n_l}{N} (b^{(k)})^\top \left( X_i^{(l)} (X_i^{(l)})^\top - \Sigma \right) b^{(k)} \right] \quad \text{for } 1 \leq i \leq n_l.$$

We write

$$D_{l,k}^{(2)} = \sum_{j \neq k, l} \sum_{i=1}^{n_j} W_{j,i}^{(2)} + \sum_{i=1}^{n_{\mathbb{Q}}} W_{\mathbb{Q},i}^{(2)},$$

where for  $1 \leq j \leq L$  and  $j \neq l, k$ ,

$$W_{j,i}^{(2)} = \frac{1}{N} X_i^{(j)} [X_i^{(j)}]^\top \quad \text{for } 1 \leq i \leq n_j, \quad \text{and} \quad W_{\mathbb{Q},i}^{(2)} = \frac{1}{N} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top \quad \text{for } 1 \leq i \leq N_{\mathbb{Q}}.$$

Since  $\{W_{l,i}^{(1)}\}_{1 \leq i \leq n_l}$  are i.i.d. random variables,  $\{W_{k,i}^{(1)}\}_{1 \leq i \leq n_k}$  are i.i.d. random variables,  $\{W_{j,i}^{(2)}\}_{n \neq l, k, 1 \leq i \leq n_j}$  and  $\{W_{\mathbb{Q},i}^{(2)}\}_{1 \leq i \leq N_{\mathbb{Q}}}$  are i.i.d. random variables and all of these random variables are independent, we apply central limit theorem and establish that

$$\frac{D_{l,k}}{\sqrt{\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(1)} + \mathbf{V}_{\pi(l,k), \pi(l,k)}^{(2)}}} \xrightarrow{d} \mathcal{N}(0, 1), \quad (122)$$

with  $\mathbf{V}^{(1)}$  and  $\mathbf{V}^{(2)}$  defined in (60) and (61), respectively.

We control  $\text{Rem}_{l,k}$  by the definition of  $\text{Rem}_{l,k}$  in (121) and the following Lemma.

**Lemma 8.** *With probability larger than  $1 - \min\{n, p\}^{-c}$ , we have*

$$\begin{aligned} \left| \frac{1}{n_l} (\widehat{b}_{init}^{(k)} - b^{(k)})^\top [X^{(l)}]^\top \epsilon^{(l)} \right| &\lesssim \frac{s \log p}{n}; \quad \left| \frac{1}{n_k} (\widehat{b}_{init}^{(l)} - b^{(l)})^\top [X^{(k)}]^\top \epsilon^{(k)} \right| \lesssim \frac{s \log p}{n}; \\ \left| (\widehat{b}_{init}^{(l)} - b^{(l)})^\top \widehat{\Sigma} (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| &\lesssim \frac{s \log p}{n}; \\ \left| [\widehat{b}_{init}^{(l)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(k)}) (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| &\lesssim (\|b^{(l)}\|_2 + 1) \frac{s \log p}{n}; \\ \left| [\widehat{b}_{init}^{(k)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(l)}) (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| &\lesssim (\|b^{(k)}\|_2 + 1) \frac{s \log p}{n}. \end{aligned}$$

**Proof of (63).** The control of  $\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(1)}$  follows from the definition (60). The control of  $\mathbf{V}_{\pi(l,k), \pi(l,k)}^{(2)}$  follows from (111) and (112).

## E Proof of Theorem 2 and (29)

We will prove the properties of the CI for the ridge-type maximin effect, which includes Theorem 2 as a special case. We first detail the generalized inference procedure in the following, which is a generalization of the inference method in Section 4.2. We construct the sampled weight as

$$\widehat{\gamma}_\delta^{[m]} = \arg \min_{\gamma \in \Delta^L} \gamma^\top (\widehat{\Gamma}^{[m]} + \delta \cdot \mathbf{I})_+ \gamma \quad \text{for } \delta \geq 0.$$

For  $1 \leq m \leq M$ , we compute  $\widehat{x_{\text{new}}^\top \beta}^{[m]} = \sum_{l=1}^L [\widehat{\gamma}_\delta^{[m]}]_l \cdot \widehat{x_{\text{new}}^\top b}^{(l)}$ , and construct the sampled interval as,

$$\text{Int}_\alpha^{[m]}(x_{\text{new}}) = \left( \widehat{x_{\text{new}}^\top \beta}^{[m]} - (1 + \eta_0) z_{\alpha/2} \widehat{\text{se}}^{[m]}(x_{\text{new}}), \widehat{x_{\text{new}}^\top \beta}^{[m]} + (1 + \eta_0) z_{\alpha/2} \widehat{\text{se}}^{[m]}(x_{\text{new}}) \right), \quad (123)$$

where  $\eta_0$  is any positive constant (with default value 0.01) and

$$\widehat{\text{se}}^{[m]}(x_{\text{new}}) = \sqrt{\sum_{l=1}^L (\widehat{\sigma}_l^2 / n_l^2) \cdot [\widehat{\gamma}_\delta^{[m]}]_l^2 [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}} \quad \text{with } \widehat{v}^{(l)} \text{ defined in (13).}$$

We slightly abuse the notation by using  $\text{Int}_\alpha^{[m]}(x_{\text{new}})$  to denote the sampled interval for both  $x_{\text{new}}^\top \beta_\delta^*$  and  $x_{\text{new}}^\top \beta^*$ . We construct the CI for  $x_{\text{new}}^\top \beta_\delta^*$  by aggregating the sampled intervals with  $m \in \mathbb{M}$ ,

$$\text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*) = \cup_{m \in \mathbb{M}} \text{Int}_\alpha^{[m]}(x_{\text{new}}), \quad (124)$$



where  $\mathbb{M}$  is defined in (20) and  $\text{Int}_\alpha^{[m]}(x_{\text{new}})$  is defined in (123).

In the following Theorem 4, we establish the coverage and precision properties of  $\text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)$  defined in (124). We apply Theorem 4 with  $\delta = 0$  and establish Theorem 2.

**Theorem 4.** *Suppose that the conditions of Theorem 3 hold. Then for any positive constant  $\eta_0 > 0$  used in (123), the confidence interval  $\text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)$  defined in (124) satisfies*

$$\lim_{n,p \rightarrow \infty} \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \in \text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) \geq 1 - \alpha - \alpha_0, \quad (125)$$

where  $\alpha \in (0, 1/2)$  is the pre-specified significance level and  $\alpha_0 \in (0, 0.01]$  is defined in (20). By further assuming  $N_\mathbb{Q} \gtrsim \max\{n, p\}$  and  $\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta \gg \sqrt{\log p / \min\{n, N_\mathbb{Q}\}}$ , then there exists some positive constant  $C > 0$  such that

$$\liminf_{n,p \rightarrow \infty} \mathbf{P}\left(\mathbf{Leng}(\text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) \leq C \max\left\{1, \frac{z_{\alpha_0/[L(L+1)]}}{\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta}\right\} \cdot \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}\right) = 1, \quad (126)$$

where  $\mathbf{Leng}(\text{CI}(x_{\text{new}}^\top \beta_\delta^*))$  denotes the interval length and  $z_{\alpha_0/[L(L+1)]}$  is the upper  $\alpha_0/[L(L+1)]$  quantile of the standard normal distribution.

In the following, we introduce the definitions of events, which are used to facilitate the proof of Theorem 4. We shall take  $m^*$  as any index such that  $\|\hat{\gamma}_\delta^{[m^*]} - \gamma_\delta^*\|_2 = \min_{m \in \mathbb{M}} \|\hat{\gamma}_\delta^{[m]} - \gamma_\delta^*\|_2$ . We introduce the following high-probability events to facilitate the discussion.

$$\begin{aligned} \mathcal{E}_4 &= \left\{ \frac{1}{n_l} [\hat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \hat{v}^{(l)} \asymp \|x_{\text{new}}\|_2 \quad \text{for } 1 \leq l \leq L \right\}, \\ \mathcal{E}_5 &= \left\{ \frac{\left| \sum_{l=1}^L \left( [\hat{\gamma}_\delta^{[m]}]_l - [\gamma_\delta^*]_l \right) \widehat{x_{\text{new}}^\top b^{(l)}} \right|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\hat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \hat{v}^{(l)}}} \lesssim \sqrt{n} \|\hat{\gamma}_\delta^{[m]} - \gamma_\delta^*\|_2, \text{ for } 1 \leq m \leq M \right\}, \\ \mathcal{E}_6 &= \left\{ \|\hat{\gamma}_\delta^{[m^*]} - \gamma_\delta^*\|_2 \leq \frac{\sqrt{2\text{err}_n(M)}}{\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta} \cdot \frac{1}{\sqrt{n}} \right\}, \\ \mathcal{E}_7 &= \left\{ \|\hat{\Gamma}^\mathbb{Q} - \Gamma^\mathbb{Q}\|_2 \lesssim \sqrt{\frac{\log p}{\min\{n, N_\mathbb{Q}\}}} + \frac{p\sqrt{\log p}}{\sqrt{n}N_\mathbb{Q}} \right\}, \\ \mathcal{E}_8 &= \left\{ \max_{1 \leq l \leq L} \frac{|\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)}|}{\|x_{\text{new}}\|_2} \lesssim \sqrt{\frac{\log n}{n}} \right\}, \\ \mathcal{E}_9 &= \left\{ \frac{\left| \sum_{l=1}^L (\hat{\gamma}_\delta^{[m^*]}]_l - [\gamma_\delta^*]_l \right) \cdot \widehat{x_{\text{new}}^\top b^{(l)}}}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\hat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \hat{v}^{(l)}}} \leq \eta_0 \cdot z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\hat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\sigma_l^2}{n_l^2} [\hat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \hat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\hat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \hat{v}^{(l)}}} \right\}. \end{aligned} \quad (127)$$

We apply Lemma 1 of [Cai et al. \(2021\)](#) and establish that  $\mathbf{P}(\mathcal{E}_4) \geq 1 - p^{-c}$ , for some positive constant  $c > 0$ . On the event  $\mathcal{E}_4$ , we have

$$\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}} \asymp \frac{\|\gamma_\delta^*\|_2 \|x_{\text{new}}\|_2}{\sqrt{n}} \asymp \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}, \quad (128)$$

where the last asymptotic equivalence holds since  $\frac{1}{\sqrt{L}} \leq \|\gamma_\delta^*\|_2 \leq 1$ . Similarly, on the event  $\mathcal{E}_4$ ,

$$\sqrt{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m]}]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}} \asymp \frac{\|\widehat{\gamma}_\delta^{[m]}\|_2 \|x_{\text{new}}\|_2}{\sqrt{n}} \asymp \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}. \quad (129)$$

We introduce the following lemma to justify the asymptotic normality of  $\widehat{x_{\text{new}}^\top b^{(l)}}$ , which follows the same proof as that of Proposition 1 in [Cai et al. \(2021\)](#).

**Lemma 9.** *Consider the model (1). Suppose Conditions (A1) and (A2) hold, then*

$$\frac{\sum_{l=1}^L c_l [\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)}]}{\sqrt{\sum_{l=1}^L c_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \xrightarrow{d} N(0, 1). \quad (130)$$

where  $|c_l| \leq 1$  for any  $1 \leq l \leq L$ , and  $\sum_{l=1}^L c_l = 1$ .

Note that

$$\begin{aligned} \left| \sum_{l=1}^L \left( [\widehat{\gamma}_\delta^{[m]}]_l - [\gamma_\delta^*]_l \right) \widehat{x_{\text{new}}^\top b^{(l)}} \right| &\leq \|\widehat{\gamma}_\delta^{[m]} - \gamma_\delta^*\|_2 \sqrt{\sum_{l=1}^L [\widehat{x_{\text{new}}^\top b^{(l)}}]^2} \\ &\lesssim \|\widehat{\gamma}_\delta^{[m]} - \gamma_\delta^*\|_2 \sqrt{\sum_{l=1}^L [\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)}]^2 + \sum_{l=1}^L [x_{\text{new}}^\top b^{(l)}]^2}. \end{aligned} \quad (131)$$

By Lemma 9 and (128), with probability larger than  $1 - \min\{n, p\}^{-c}$  for some positive constant  $c > 0$ ,

$$\begin{aligned} &\frac{\left| \sum_{l=1}^L \left( [\widehat{\gamma}_\delta^{[m]}]_l - [\gamma_\delta^*]_l \right) \widehat{x_{\text{new}}^\top b^{(l)}} \right|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \\ &\lesssim \|\widehat{\gamma}_\delta^{[m]} - \gamma_\delta^*\|_2 \cdot \frac{\sqrt{L} \cdot (\log n \cdot \frac{\|x_{\text{new}}\|_2}{\sqrt{n}} + \max_{1 \leq l \leq L} |x_{\text{new}}^\top b^{(l)}|)}{\frac{1}{\sqrt{L}} \cdot \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}} \lesssim \sqrt{n} \|\widehat{\gamma}_\delta^{[m]} - \gamma_\delta^*\|_2, \end{aligned} \quad (132)$$

where the last inequality follows from bounded  $\|b^{(l)}\|_2$  and finite  $L$ . This implies that  $\mathbf{P}(\mathcal{E}_5) \geq 1 - \min\{n, p\}^{-c}$  for some positive constant  $c > 0$ . It follows from (72) of Theorem 3 that  $\liminf_{n \rightarrow \infty} \mathbf{P}(\mathcal{E}_6) \geq 1 - \alpha_0$ . It follows from Propositions 4 and 5 that  $\liminf_{n \rightarrow \infty} \mathbf{P}(\mathcal{E}_7) = 1$ . We apply Lemma 9 with  $c_l = 1$  and  $c_j = 0$  for  $j \neq l$  and show that

$$\mathbf{P}(\mathcal{E}_8 \cap \mathcal{E}_4) \geq 1 - \min\{n, p\}^{-c},$$

for some positive constant  $c > 0$ . On the event  $\mathcal{E}_5 \cap \mathcal{E}_6$ , we have

$$\frac{|\sum_{l=1}^L (\widehat{\gamma}_\delta^{[m*]})_l - [\gamma_\delta^*]_l \cdot \widehat{x_{\text{new}}^\top b^{(l)}}|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \lesssim \frac{\sqrt{2\text{err}_n(M)}}{\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta}.$$

On the event  $\mathcal{G}_3 \cap \mathcal{E}_4 \cap \mathcal{E}_6$  with  $\mathcal{G}_3$  defined in (106) and  $\mathcal{E}_4$  and  $\mathcal{E}_6$  defined in (127), we apply (128) and (129) to establish  $z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m*]})_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \geq c$  for some small constant  $c > 0$ . If  $\eta_0$  used in (123) satisfies

$$\eta_0 \geq C \frac{\sqrt{2\text{err}_n(M)}}{\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta}, \quad (133)$$

for some positive constant  $C > 0$ , we establish  $\mathcal{E}_9$ , which implies

$$\liminf_{n, p \rightarrow \infty} \mathbf{P}(\mathcal{E}_9) \geq \liminf_{n, p \rightarrow \infty} \mathbf{P}(\mathcal{G}_3 \cap \mathcal{E}_4 \cap \mathcal{E}_5 \cap \mathcal{E}_6) \geq 1 - \alpha_0. \quad (134)$$

Since  $\text{err}_n(M) \ll \lambda_{\min}(\Gamma^\mathbb{Q}) + \delta$ , then any positive constant  $\eta_0 > 0$  will guarantee (134).

### E.1 Coverage Property: Proof of (125)

It follows from Lemma 9 that

$$\frac{\sum_{l=1}^L [\gamma_\delta^*]_l [\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)}]}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \xrightarrow{d} N(0, 1). \quad (135)$$

By the definition in (124), we have

$$\begin{aligned} \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) &\leq \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{Int}_\alpha^{[m^*]}(x_{\text{new}})) \\ &= \mathbf{P}\left(\frac{|\widehat{x_{\text{new}}^\top \beta}^{[m^*]} - x_{\text{new}}^\top \beta_\delta^*|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \geq (1 + \eta_0) \cdot z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}\right). \end{aligned}$$

By applying (26) with  $m = m^*$ , we further upper bound the above inequality as

$$\begin{aligned} \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) &\leq \mathbf{P}(\mathcal{E}_9^c) + \\ \mathbf{P}\left(\frac{|\sum_{l=1}^L [\gamma_\delta^*]_l \cdot (\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)})|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \geq z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}\right). \end{aligned} \quad (136)$$

Since  $\mathcal{E}_4 \cap \mathcal{E}_6$  holds with a high probability, we have

$$\sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \xrightarrow{p} 1.$$

Together with (135), we establish

$$\mathbf{P}\left(\frac{|\sum_{l=1}^L [\gamma_\delta^*]_l \cdot (\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)})|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \geq z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}\right) \rightarrow \alpha. \quad (137)$$

Combined with (136) and (134), we have

$$\limsup_{n, p \rightarrow \infty} \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) \leq \limsup_{n, p \rightarrow \infty} \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{Int}_\alpha^{[m^*]}(x_{\text{new}})) \leq \alpha + \alpha_0. \quad (138)$$

## E.2 Precision Property: Proof of (126)

Regarding the length of the confidence interval, we notice that

$$\text{Leng}(\text{CI}(x_{\text{new}}^\top \beta_\delta^*)) \leq 2 \max_{m \in \mathbb{M}} \left( \left| \widehat{x_{\text{new}}^\top \beta}^{[m]} - \widehat{x_{\text{new}}^\top \beta_\delta^*} \right| + \widehat{\text{se}}^{[m]}(x_{\text{new}}) \right), \quad (139)$$

where  $\widehat{x_{\text{new}}^\top \beta_\delta^*} = \sum_{l=1}^L [\widehat{\gamma}_\delta]_l \cdot \widehat{x_{\text{new}}^\top b^{(l)}}$  is defined in (31) and  $\widehat{x_{\text{new}}^\top \beta^{[m]}} = \sum_{l=1}^L [\widehat{\gamma}_\delta^{[m]}]_l \cdot \widehat{x_{\text{new}}^\top b^{(l)}}$ . Note that

$$\begin{aligned} \max_{m \in \mathbb{M}} \left| \widehat{x_{\text{new}}^\top \beta^{[m]}} - \widehat{x_{\text{new}}^\top \beta_\delta^*} \right| &= \max_{m \in \mathbb{M}} \left| \sum_{l=1}^L \left( [\widehat{\gamma}_\delta]_l - [\widehat{\gamma}_\delta^{[m]}]_l \right) \cdot \widehat{x_{\text{new}}^\top b^{(l)}} \right| \\ &\leq \max_{m \in \mathbb{M}} \|\widehat{\gamma}_\delta - \widehat{\gamma}_\delta^{[m]}\|_2 \cdot \sqrt{\sum_{l=1}^L (\widehat{x_{\text{new}}^\top b^{(l)}})^2}. \end{aligned} \quad (140)$$

We apply  $N_{\mathbb{Q}} \gtrsim \max\{n, p\}$  and  $\lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta \gg \sqrt{\log p / \min\{n, N_{\mathbb{Q}}\}}$  and establish

$$\lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta \gg \sqrt{\frac{\log p}{\min\{n, N_{\mathbb{Q}}\}}} + \frac{p\sqrt{\log p}}{\sqrt{n}N_{\mathbb{Q}}}.$$

On the event  $\mathcal{E}_7$ , we establish

$$\lambda_{\min}(\widehat{\Gamma}) + \delta \geq \frac{1}{2} (\lambda_{\min}(\Gamma^{\mathbb{Q}}) + \delta). \quad (141)$$

We now apply Lemma 3 with  $\widehat{\Gamma} = (\widehat{\Gamma}^{[m]} + \delta \cdot \mathbf{I})_+$  and  $\Gamma = \widehat{\Gamma} + \delta \cdot \mathbf{I}$  and establish that

$$\|\widehat{\gamma}_\delta - \widehat{\gamma}_\delta^{[m]}\|_2 \leq \frac{\|(\widehat{\Gamma}^{[m]} + \delta \cdot \mathbf{I})_+ - (\widehat{\Gamma} + \delta \cdot \mathbf{I})\|_F}{\lambda_{\min}(\widehat{\Gamma}) + \delta} \leq \frac{\|\widehat{\Gamma}^{[m]} - \widehat{\Gamma}\|_F}{\lambda_{\min}(\widehat{\Gamma}) + \delta},$$

where the last inequality follows from Lemma 4 together with  $\lambda_{\min}(\widehat{\Gamma}) + \delta > 0$  on the event  $\mathcal{E}_7$ . Together with (140), we establish

$$\max_{m \in \mathbb{M}} \left| \widehat{x_{\text{new}}^\top \beta^{[m]}} - \widehat{x_{\text{new}}^\top \beta_\delta^*} \right| \leq \max_{m \in \mathbb{M}} \frac{\|\widehat{\Gamma}^{[m]} - \widehat{\Gamma}\|_F}{\lambda_{\min}(\widehat{\Gamma}) + \delta} \cdot \sqrt{\sum_{l=1}^L (\widehat{x_{\text{new}}^\top b^{(l)}})^2}. \quad (142)$$

By the definition of  $\mathbb{M}$  in (20) and the definition of  $d_0$  in (17), we have

$$\max_{m \in \mathbb{M}} \|\widehat{\Gamma}^{[m]} - \widehat{\Gamma}\|_F \lesssim L \cdot \sqrt{d_0/n} \cdot 1.1 \cdot z_{\alpha_0/[L(L+1)]}.$$

Together with (142), we establish

$$\max_{m \in \mathbb{M}} \left| \widehat{x_{\text{new}}^\top \beta^{[m]}} - \widehat{x_{\text{new}}^\top \beta_\delta^*} \right| \lesssim \frac{1.1L \cdot \sqrt{d_0}}{\sqrt{n}[\lambda_{\min}(\widehat{\Gamma}) + \delta]} \cdot \sqrt{\sum_{l=1}^L (\widehat{x_{\text{new}}^\top b^{(l)}})^2 \cdot z_{\alpha_0/[L(L+1)]}}. \quad (143)$$

On the event  $\mathcal{E}_4$ , we apply (129) and establish

$$\max_{m \in \mathbb{M}} \widehat{\text{se}}^{[m]}(x_{\text{new}}) \lesssim \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}. \quad (144)$$

We combine (139), (143), (141) and (144) and establish that, on the event  $\mathcal{E}_4 \cap \mathcal{E}_7$ , we have

$$\mathbf{Leng}(\text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) \lesssim \frac{L \cdot \sqrt{d_0}}{\sqrt{n}[\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta]} \cdot \sqrt{\sum_{l=1}^L (\widehat{x_{\text{new}}^\top b^{(l)}})^2 \cdot z_{\alpha_0/[L(L+1)]}} + \frac{\|x_{\text{new}}\|_2}{\sqrt{n}}. \quad (145)$$

On the event  $\mathcal{E}_8$ , we apply the Condition (A2) and establish that

$$\frac{1}{\|x_{\text{new}}\|_2} \sqrt{\sum_{l=1}^L (\widehat{x_{\text{new}}^\top b^{(l)}})^2} \lesssim \sqrt{L} \left( \sqrt{\frac{\log n}{n}} + \frac{|x_{\text{new}}^\top b^{(l)}|}{\|x_{\text{new}}\|_2} \right) \leq C, \quad (146)$$

for some positive constant  $C > 0$ . For a finite  $L$ ,  $\text{Vol}(L(L+1)/2)$  and  $z_{\alpha_0/[L(L+1)]}$  are bounded from above. If  $N_\mathbb{Q} \gtrsim \max\{n, p\}$ ,  $s \log p/n \rightarrow 0$  and Condition (A2) holds, we apply (56) and show that

$$n \cdot \lambda_i(\mathbf{V}) \lesssim n \cdot \|\mathbf{V}\|_\infty \lesssim 1, \quad \text{and} \quad d_0 \lesssim 1. \quad (147)$$

Hence, if  $N_\mathbb{Q} \gtrsim \max\{n, p\}$  and  $\lambda_{\min}(\Gamma^\mathbb{Q}) + \delta \gg \sqrt{\log p / \min\{n, N_\mathbb{Q}\}}$ , we establish (126) by combining (145) with (141), (146) and (147).

### E.3 Proof of (29)

By (138), it is sufficient to establish

$$\liminf_{n, p \rightarrow \infty} \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{CI}_\alpha(x_{\text{new}}^\top \beta_\delta^*)) \geq \alpha - \alpha_0. \quad (148)$$

By applying (26) with  $m = m^*$ , we establish that, on the event  $\mathcal{E}_9$ ,

$$\begin{aligned} & \mathbf{P}(x_{\text{new}}^\top \beta_\delta^* \notin \text{Int}_\alpha^{[m^*]}(x_{\text{new}})) \geq \mathbf{P}(\{x_{\text{new}}^\top \beta_\delta^* \notin \text{Int}_\alpha^{[m^*]}(x_{\text{new}})\} \cap \mathcal{E}_9) \\ & \geq \mathbf{P}\left(\left\{\frac{|\sum_{l=1}^L [\gamma_\delta^*]_l \cdot (\widehat{x_{\text{new}}^\top b^{(l)}}) - x_{\text{new}}^\top b^{(l)}|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}} \geq (1 + 2\eta_0) z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}}\right\} \cap \mathcal{E}_9\right) \\ & \geq \mathbf{P}\left(\frac{|\sum_{l=1}^L [\gamma_\delta^*]_l \cdot (\widehat{x_{\text{new}}^\top b^{(l)}}) - x_{\text{new}}^\top b^{(l)}|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}} \geq (1 + 2\eta_0) z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}}}\right) - \mathbf{P}(\mathcal{E}_9^c), \end{aligned}$$

where the last inequality follows from the union bound. If  $\eta_0 \rightarrow 0$ , we apply the similar argument of (137) and establish

$$\mathbf{P} \left( \frac{|\sum_{l=1}^L [\gamma_\delta^*]_l \cdot (\widehat{x_{\text{new}}^\top b^{(l)}} - x_{\text{new}}^\top b^{(l)})|}{\sqrt{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \geq (1 + 2\eta_0) z_{\alpha/2} \sqrt{\frac{\sum_{l=1}^L [\widehat{\gamma}_\delta^{[m^*]}]_l^2 \frac{\widehat{\sigma}_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}{\sum_{l=1}^L [\gamma_\delta^*]_l^2 \frac{\sigma_l^2}{n_l^2} [\widehat{v}^{(l)}]^\top (X^{(l)})^\top X^{(l)} \widehat{v}^{(l)}}} \right) \rightarrow \alpha.$$

Combined with (134), we establish (148).

## F Proofs of Propositions 1, 2, and 3

### F.1 Proofs of Proposition 2

For  $\mathbb{B} = \begin{pmatrix} b^{(1)} & b^{(2)} & \dots & b^{(L)} \end{pmatrix} \in \mathbb{R}^{p \times L}$ , we define its SVD as

$$\mathbb{B} = U \Lambda V^\top \quad \text{with} \quad U \in \mathbb{R}^{p \times L}, \Lambda \in \mathbb{R}^{L \times L}, \text{ and } V \in \mathbb{R}^{L \times L}$$

where  $\Lambda_{1,1} \geq \dots \geq \Lambda_{L,L} > 0$ . We define  $\Delta = \delta \cdot U \Lambda^{-2} U^\top \in \mathbb{R}^{p \times p}$ . Then

$$[\mathbb{B}]^\top \Delta \mathbb{B} = \delta \mathbf{I} \quad \text{and} \quad \mathbb{B}^\top (\Sigma^\mathbb{Q} + \Delta) \mathbb{B} = \Gamma^\mathbb{Q} + \delta \cdot \mathbf{I}. \quad (149)$$

It follows from Proposition 1 and the definition of  $\beta_\delta^*$  in (30) that

$$\beta_\delta^* = \max_{\beta \in \mathbb{R}^p} \min_{b \in \mathbb{B}} [2b^\top (\Sigma^\mathbb{Q} + \Delta) \beta - \beta^\top (\Sigma^\mathbb{Q} + \Delta) \beta]$$

and

$$\min_{b \in \mathbb{B}} [2b^\top (\Sigma^\mathbb{Q} + \Delta) \beta_\delta^* - [\beta_\delta^*]^\top (\Sigma^\mathbb{Q} + \Delta) \beta_\delta^*] = [\beta_\delta^*]^\top (\Sigma^\mathbb{Q} + \Delta) \beta_\delta^* = [\gamma_\delta^*]^\top (\Gamma^\mathbb{Q} + \delta \cdot \mathbf{I}) \gamma_\delta^* \quad (150)$$

Now we compute the lower bound for  $R_\mathbb{Q}(\beta_\delta^*) = \min_{b \in \mathbb{B}} [2b^\top \Sigma^\mathbb{Q} \beta_\delta^* - [\beta_\delta^*]^\top \Sigma^\mathbb{Q} \beta_\delta^*]$ .

We apply (149) and  $\beta_\delta^* = \mathbb{B} \gamma_\delta^*$  and establish

$$\begin{aligned} & \min_{b \in \mathbb{B}} [2b^\top (\Sigma^\mathbb{Q} + \Delta) \beta_\delta^* - [\beta_\delta^*]^\top (\Sigma^\mathbb{Q} + \Delta) \beta_\delta^*] \\ &= \min_{b \in \mathbb{B}} [2b^\top (\Sigma^\mathbb{Q} + \Delta) \beta_\delta^* - [\beta_\delta^*]^\top \Sigma^\mathbb{Q} \beta_\delta^*] - \delta \|\gamma_\delta^*\|_2^2 \\ &\leq \min_{b \in \mathbb{B}} [2b^\top \Sigma^\mathbb{Q} \beta_\delta^* - [\beta_\delta^*]^\top \Sigma^\mathbb{Q} \beta_\delta^*] + 2 \max_{b \in \mathbb{B}} b^\top \Delta \beta_\delta^* - \delta \|\gamma_\delta^*\|_2^2 \\ &= R_\mathbb{Q}(\beta_\delta^*) + 2\delta \max_{\gamma \in \Delta^L} \gamma^\top \gamma_\delta^* - \delta \|\gamma_\delta^*\|_2^2 \end{aligned} \quad (151)$$

where the last inequality follows from Proposition 1 and the function  $b^\top \Delta \beta_\delta^*$  is linear in  $b$ ,  $\beta_\delta^* = \mathbb{B} \gamma_\delta^*$  and (149). We combine (150) and (151) and establish

$$\begin{aligned}
R_{\mathbb{Q}}(\beta_\delta^*) &\geq [\gamma_\delta^*]^\top (\Gamma^{\mathbb{Q}} + \delta \cdot \mathbf{I}) \gamma_\delta^* - 2\delta \max_{\gamma \in \Delta^L} \gamma^\top \gamma_\delta^* + \delta \|\gamma_\delta^*\|_2^2 \\
&\geq [\gamma^*]^\top \Gamma^{\mathbb{Q}} \gamma^* + 2\delta \|\gamma_\delta^*\|_2^2 - 2\delta \max_{\gamma \in \Delta^L} \gamma^\top \gamma_\delta^* \\
&= R_{\mathbb{Q}}(\beta^*) - 2\delta \left( \max_{\gamma \in \Delta^L} \gamma^\top \gamma_\delta^* - \|\gamma_\delta^*\|_2^2 \right)
\end{aligned} \tag{152}$$

where the second inequality follows from the definition of  $\gamma^*$ . Note that  $\max_{\gamma \in \Delta^L} \gamma^\top \gamma_\delta^* - \|\gamma_\delta^*\|_2^2 \geq 0$  and  $\max_{\gamma \in \Delta^L} \gamma^\top \gamma_\delta^* - \|\gamma_\delta^*\|_2^2 = \|\gamma_\delta^*\|_\infty - \|\gamma_\delta^*\|_2^2$ . We use  $j^* \in [L]$  to denote the index such that  $[\gamma_\delta^*]_{j^*} = \|\gamma_\delta^*\|_\infty$ . Then we have

$$\begin{aligned}
\|\gamma_\delta^*\|_\infty - \|\gamma_\delta^*\|_2^2 &= [\gamma_\delta^*]_{j^*} - [\gamma_\delta^*]_{j^*}^2 - \sum_{l \neq j^*} [\gamma_\delta^*]_l^2 \\
&\leq [\gamma_\delta^*]_{j^*} - [\gamma_\delta^*]_{j^*}^2 - \frac{1}{L-1} \left( \sum_{l \neq j^*} [\gamma_\delta^*]_l \right)^2 \\
&= [\gamma_\delta^*]_{j^*} - [\gamma_\delta^*]_{j^*}^2 - \frac{1}{L-1} (1 - [\gamma_\delta^*]_{j^*})^2
\end{aligned} \tag{153}$$

We take the maximum value of the right hand side with respect to  $[\gamma_\delta^*]_{j^*}$  over the domain  $[1/L, 1]$ . Then we obtain

$$\max_{\frac{1}{L} \leq [\gamma_\delta^*]_{j^*} \leq 1} [\gamma_\delta^*]_{j^*} - [\gamma_\delta^*]_{j^*}^2 - \frac{1}{L-1} (1 - [\gamma_\delta^*]_{j^*})^2 = \frac{1}{4} \left( 1 - \frac{1}{L} \right)$$

where the maximum value is achieved at  $[\gamma_\delta^*]_{j^*} = \frac{1+\frac{1}{L}}{2}$ . Combined with (152) and (153), we establish

$$R_{\mathbb{Q}}(\beta_\delta^*) \geq R_{\mathbb{Q}}(\beta^*) - \frac{\delta}{2} \cdot \left( 1 - \frac{1}{L} \right).$$

## F.2 Proof of Proposition 1

We now supply a proof of Proposition 1, which essentially follows from the same argument as that of Theorem 1 in Meinshausen and Bühlmann (2015). Let  $\mathbb{H}$  denote the convex hull



of the set  $\mathbb{B} = \{b^{(1)}, \dots, b^{(L)}\}$ . By the linear form of  $b$ , we have

$$\begin{aligned}\beta^* &= \arg \max_{\beta \in \mathbb{R}^p} \min_{b \in \mathbb{B}} [2b^\top \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta] \\ &= \arg \max_{\beta \in \mathbb{R}^p} \min_{b \in \mathbb{H}} [2b^\top \Sigma^{\mathbb{Q}} \beta - \beta^\top \Sigma^{\mathbb{Q}} \beta]\end{aligned}$$

We decompose  $\Sigma^{\mathbb{Q}} = C^\top C$  such that  $C$  is invertible. Define  $\tilde{\mathbb{H}} = C^{-1}\mathbb{H}$ . Then we have  $\beta^* = C^{-1}\xi^*$  with

$$\xi^* = \arg \max_{\xi \in \mathbb{R}^p} \min_{u \in \tilde{\mathbb{H}}} [2u^\top \xi - \xi^\top \xi] \quad (154)$$

If we interchange min and max in the above equation, then we have

$$\xi^* = \arg \min_{\xi \in \tilde{\mathbb{H}}} \xi^\top \xi \quad (155)$$

We will justify this inter-change by showing that the solution  $\xi^*$  defined in (155) is the solution to (154). For any  $\nu \in [0, 1]$  and  $\mu \in \tilde{\mathbb{H}}$ , we use the fact  $\xi^* + \nu(\mu - \xi^*) \in \tilde{\mathbb{H}}$  and obtain

$$\|\xi^* + \nu(\mu - \xi^*)\|_2^2 \geq \|\xi^*\|_2^2$$

This leads to

$$(\xi^*)^\top \mu - (\xi^*)^\top \xi^* \geq 0$$

and hence

$$2(\xi^*)^\top \mu - (\xi^*)^\top \xi^* \geq (\xi^*)^\top \xi^*, \quad \text{for any } \mu \in \tilde{\mathbb{H}}$$

By taking  $\xi$  as  $\xi^*$  in the optimization problem (154), we have

$$\max_{\xi \in \mathbb{R}^p} \min_{u \in \tilde{\mathbb{H}}} [2u^\top \xi - \xi^\top \xi] \geq \min_{u \in \tilde{\mathbb{H}}} [2u^\top \xi^* - [\xi^*]^\top \xi^*] \geq (\xi^*)^\top \xi^*.$$

In (154), we take  $u = \xi^*$ , then we have

$$\max_{\xi \in \mathbb{R}^p} \min_{u \in \tilde{\mathbb{H}}} [2u^\top \xi - \xi^\top \xi] \leq \max_{\xi \in \mathbb{R}^p} [2[\xi^*]^\top \xi - \xi^\top \xi] = (\xi^*)^\top \xi^*$$

By matching the above two bounds,  $\xi^*$  is the optimal solution to (154) and

$$\max_{\xi \in \mathbb{R}^p} \min_{u \in \tilde{\mathbb{H}}} [2u^\top \xi - \xi^\top \xi] = [\xi^*]^\top \xi^*.$$

Since  $\beta^* = C^{-1}\xi^*$  and  $\Sigma^{\mathbb{Q}} = C^{\top}C$ , we have

$$\beta^* = \arg \min_{\beta \in \mathbb{H}} \beta^{\top} \Sigma^{\mathbb{Q}} \beta \quad (156)$$

and

$$\max_{\beta \in \mathbb{R}^p} \min_{b \in \mathbb{H}} [2b^{\top} \Sigma^{\mathbb{Q}} \beta - \beta^{\top} \Sigma^{\mathbb{Q}} \beta] = [\beta^*]^{\top} \Sigma^{\mathbb{Q}} \beta^*.$$

We establish (10) by combining (156) and the fact that  $\beta \in \mathbb{H}$  can be expressed as  $\beta = \mathbb{B}\gamma$  for  $\gamma \in \Delta^L$ .

### F.3 Proof of Proposition 3

We can write the maximin definition in the following form

$$\beta^{*,\text{MP}} = \arg \max_{\|\beta\|_2 \leq 1} \min_{b \in \mathbb{B}} \beta^{\top} b \quad (157)$$

where  $\mathbb{B} = \{b^{(1)}, \dots, b^{(L)}\}$ . Since  $b^{\top} \beta$  is linear in  $b$ , we can replace  $\mathbb{B}$  with its convex hull  $\mathbb{H}$ ,

$$\beta^{*,\text{MP}} = \arg \max_{\|\beta\|_2 \leq 1} \min_{b \in \mathbb{H}} b^{\top} \beta$$

We exchange the max and min in the above equation and have

$$\min_{b \in \mathbb{H}} \max_{\|\beta\|_2 \leq 1} b^{\top} \beta = \min_{b \in \mathbb{H}} \|b\|_2$$

We define

$$\xi = \arg \min_{b \in \mathbb{H}} \|b\|_2.$$

We claim that  $\xi^* = \xi/\|\xi\|$  is the optimal solution of (157). For any  $\mu \in \mathbb{H}$ , we have  $\xi + \nu(\mu - \xi) \in \mathbb{H}$  for  $\nu \in [0, 1]$  and have

$$\|\xi + \nu(\mu - \xi)\|_2^2 \geq \|\xi\|_2^2 \quad \text{for any } \nu \in [0, 1]$$

By taking  $\nu \rightarrow 0$ , we have

$$\mu^{\top} \xi - \|\xi\|_2^2 \geq 0$$

By dividing both sides by  $\|\xi\|_2$ , we have

$$\mu^\top \xi^* \geq \|\xi\|_2 \quad \text{for any } \mu \in \mathbb{H}. \quad (158)$$

In the definition of (157), we take  $\beta = \xi^*$  and have

$$\max_{\|\beta\|_2 \leq 1} \min_{b \in \mathbb{H}} b^\top \beta \geq \min_{b \in \mathbb{H}} b^\top \xi^* \geq \|\xi\|_2 \quad (159)$$

where the last inequality follows from (158). Additionally, we take  $b = \xi$  in the definition of (157) and have

$$\max_{\|\beta\|_2 \leq 1} \min_{b \in \mathbb{H}} b^\top \beta \leq \max_{\|\beta\|_2 \leq 1} \xi^\top \beta = \|\xi\|_2$$

Combined with (159), we have shown that

$$\xi^* = \arg \max_{\|\beta\|_2 \leq 1} \min_{b \in \mathbb{H}} b^\top \beta$$

that is,  $\beta^{*,\text{MP}} = \xi^*$ .

## G Proofs of Extra Lemmas

### G.1 Proof of Lemma 7

On the event  $\mathcal{G}_1 \cap \mathcal{G}_6(\widehat{b}_{init}^{(l)} - b^{(l)}, \widehat{b}_{init}^{(l)} - b^{(l)}, \sqrt{\log p})$ , we have

$$\frac{1}{|B|} \sum_{i \in B} [(X_i^\mathbb{Q})^\top (\widehat{b}_{init}^{(l)} - b^{(l)})]^2 \lesssim \frac{\|b^{(l)}\|_0 \log p}{n_l} \sigma_l^2.$$

Then we have

$$\begin{aligned} \left| (\widehat{b}_{init}^{(l)} - b^{(l)})^\top \widehat{\Sigma}^\mathbb{Q} (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| &\leq \frac{1}{|B|} \|X_B^\mathbb{Q} (\widehat{b}_{init}^{(l)} - b^{(l)})\|_2 \|X_B^\mathbb{Q} (\widehat{b}_{init}^{(k)} - b^{(k)})\|_2 \\ &\lesssim \sqrt{\frac{\|b^{(l)}\|_0 \|b^{(k)}\|_0 (\log p)^2}{n_l n_k}} \end{aligned}$$

and establish (118). We decompose

$$\begin{aligned} &(\widehat{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) \\ &= (\widetilde{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) + [\widehat{b}_{init}^{(k)}]^\top (\widehat{\Sigma}^\mathbb{Q} - \widetilde{\Sigma}^\mathbb{Q})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}). \end{aligned} \quad (160)$$

Regarding the first term of (160), we apply Hölder's inequality and establish

$$\left| (\tilde{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| \leq \|\tilde{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)}\|_\infty \|\widehat{b}_{init}^{(l)} - b^{(l)}\|_1.$$

By the optimization constraint (44), on the event  $\mathcal{G}_2$ , we have

$$\left| (\tilde{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k)} - \widehat{\Sigma}^{(l)} \widehat{u}^{(l,k)})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| \lesssim \|\omega^{(k)}\|_2 \sqrt{\frac{\log p}{n_l}} \cdot \|b^{(l)}\|_0 \sqrt{\frac{\log p}{n_l}}. \quad (161)$$

Regarding the second term of (160), conditioning on  $\widehat{b}_{init}^{(k)}$  and  $\widehat{b}_{init}^{(l)}$ , on the event  $\mathcal{G}_6(\widehat{b}_{init}^{(k)}, \widehat{b}_{init}^{(l)} - b^{(l)}, \sqrt{\log p})$ ,

$$\left| [\widehat{b}_{init}^{(k)}]^\top (\widehat{\Sigma}^{\mathbb{Q}} - \tilde{\Sigma}^{\mathbb{Q}})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| \lesssim \frac{\sqrt{\log p}}{\sqrt{N_{\mathbb{Q}}}} \|\widehat{b}_{init}^{(k)}\|_2 \|\widehat{b}_{init}^{(l)} - b^{(l)}\|_2.$$

On the event  $\mathcal{G}_1$ , we further have

$$\left| [\widehat{b}_{init}^{(k)}]^\top (\widehat{\Sigma}^{\mathbb{Q}} - \tilde{\Sigma}^{\mathbb{Q}})^\top (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| \lesssim \|\widehat{b}_{init}^{(k)}\|_2 \sqrt{\frac{\|b^{(l)}\|_0 (\log p)^2}{n_l N_{\mathbb{Q}}}}.$$

Combined with (161), we establish (119). We establish (120) through applying the similar argument for (119) by exchanging the role of  $l$  and  $k$ . Together with (108), (109) and (110) with  $t = \sqrt{\log p}$ , we establish the lemma.

## G.2 Proof of Lemma 8

On the event  $\mathcal{G}_0 \cap \mathcal{G}_2$ , we apply the Hölder's inequality and establish

$$\left| \frac{1}{n_l} (\widehat{b}_{init}^{(k)} - b^{(k)})^\top [X^{(l)}]^\top \epsilon^{(l)} \right| \leq \|\widehat{b}_{init}^{(k)} - b^{(k)}\|_1 \left\| \frac{1}{n_l} [X^{(l)}]^\top \epsilon^{(l)} \right\|_\infty \lesssim \frac{s \log p}{n}.$$

Similarly, we establish  $\left| \frac{1}{n_k} (\widehat{b}_{init}^{(l)} - b^{(l)})^\top [X^{(k)}]^\top \epsilon^{(k)} \right| \lesssim s \log p / n$ . We define the event

$$\mathcal{G}'_5 = \left\{ \max_{\substack{\mathcal{S} \subset [p], |\mathcal{S}| \leq s \\ \mathcal{T} \subset [p], |\mathcal{T}| \leq s}} \max_{\substack{\|w_{\mathcal{S}^c}\|_1 \leq C \\ \|v_{\mathcal{T}^c}\|_1 \leq C}} \frac{v^\top \left( \frac{1}{N} \sum_{i=1}^N X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top \right) w}{\sqrt{v^\top E(X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]) v} \sqrt{w^\top E(X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]) w}} \leq 1 + C \sqrt{\frac{s \log p}{N}} \right\}$$

It follows from Theorem 1.6 of Zhou (2009) that the event  $\mathcal{G}'_5$  holds with probability larger than  $1 - p^{-c}$  for some positive constant  $c > 0$ . On the event  $\mathcal{G}_1 \cap \mathcal{G}'_5$  with  $N = \sum_{l=1}^L n_l + N_{\mathbb{Q}}$ ,

we have

$$\left| (\widehat{b}_{init}^{(l)} - b^{(l)})^\top \widehat{\Sigma} (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| \lesssim \frac{s \log p}{n} \cdot \left( 1 + \sqrt{\frac{s \log p}{\sum_{l=1}^L n_l + N_{\mathbb{Q}}}} \right). \quad (162)$$

Note that

$$\begin{aligned} & [\widehat{b}_{init}^{(l)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(k)}) (\widehat{b}_{init}^{(k)} - b^{(k)}) \\ &= [\widehat{b}_{init}^{(l)} - b^{(l)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(k)}) (\widehat{b}_{init}^{(k)} - b^{(k)}) + [b^{(l)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(k)}) (\widehat{b}_{init}^{(k)} - b^{(k)}) \\ &= [\widehat{b}_{init}^{(l)} - b^{(l)}]^\top \widehat{\Sigma} (\widehat{b}_{init}^{(k)} - b^{(k)}) + [b^{(l)}]^\top (\widehat{\Sigma} - \Sigma) (\widehat{b}_{init}^{(k)} - b^{(k)}) \\ &\quad - [\widehat{b}_{init}^{(l)} - b^{(l)}]^\top \widetilde{\Sigma}^{(k)} (\widehat{b}_{init}^{(k)} - b^{(k)}) - [b^{(l)}]^\top (\widetilde{\Sigma}^{(k)} - \Sigma) (\widehat{b}_{init}^{(k)} - b^{(k)}) \end{aligned} \quad (163)$$

With a similar proof for (162), we show that, on the event  $\mathcal{G}_1 \cap \mathcal{G}'_5$ ,

$$\left| [\widehat{b}_{init}^{(l)} - b^{(l)}]^\top \widetilde{\Sigma}^{(k)} (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| \lesssim \frac{s \log p}{n} \cdot \left( 1 + \sqrt{\frac{s \log p}{n}} \right). \quad (164)$$

By Hölder's inequality, we have

$$\begin{aligned} \left| [b^{(l)}]^\top (\widehat{\Sigma} - \Sigma) (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| &\leq \|(\widehat{\Sigma} - \Sigma) b^{(l)}\|_\infty \|\widehat{b}_{init}^{(k)} - b^{(k)}\|_1 \\ \left| [b^{(l)}]^\top (\widetilde{\Sigma}^{(k)} - \Sigma) (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| &\leq \|(\widetilde{\Sigma}^{(k)} - \Sigma) b^{(l)}\|_\infty \|\widehat{b}_{init}^{(k)} - b^{(k)}\|_1 \end{aligned}$$

We define

$$\mathcal{G}'_6(w, v, t) = \left\{ \left| w^\top (\widehat{\Sigma} - \Sigma) v \right| \lesssim t \frac{\|\Sigma^{1/2} w\|_2 \|\Sigma^{1/2} v\|_2}{\sqrt{\sum_{l=1}^L n_l + N_{\mathbb{Q}}}}, \left| w^\top (\widetilde{\Sigma}^{(k)} - \Sigma) v \right| \lesssim t \frac{\|\Sigma^{1/2} w\|_2 \|\Sigma^{1/2} v\|_2}{\sqrt{n_k}} \right\}$$

and it follows from Lemma 10 in the supplement of [Cai and Guo \(2020\)](#) that

$$\mathbf{P}(\mathcal{G}'_6(w, v, t)) \geq 1 - \exp(-t^2).$$

On the event  $\cap_{j=1}^p \mathcal{G}'_6(b^{(l)}, e_j, \sqrt{\log p}) \cap \mathcal{G}_1$ , we have

$$\max \left\{ \|(\widehat{\Sigma} - \Sigma) b^{(l)}\|_\infty, \|(\widetilde{\Sigma}^{(k)} - \Sigma) b^{(l)}\|_\infty \right\} \lesssim \sqrt{\frac{\log p}{n}} \|b^{(l)}\|_2$$

and

$$\max \left\{ \left| [b^{(l)}]^\top (\widehat{\Sigma} - \Sigma) (\widehat{b}_{init}^{(k)} - b^{(k)}) \right|, \left| [b^{(l)}]^\top (\widetilde{\Sigma}^{(k)} - \Sigma) (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| \right\} \lesssim \|b^{(l)}\|_2 \cdot \frac{s \log p}{n}.$$

Together with (162), (163), (164), (108) and (109), we establish

$$\mathbf{P} \left( \left| [\widehat{b}_{init}^{(l)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(k)}) (\widehat{b}_{init}^{(k)} - b^{(k)}) \right| \lesssim (\|b^{(l)}\|_2 + 1) \frac{s \log p}{n} \right) \geq 1 - \min\{n, p\}^{-c}.$$

We apply a similar argument to show that

$$\mathbf{P} \left( \left| [\widehat{b}_{init}^{(k)}]^\top (\widehat{\Sigma} - \widetilde{\Sigma}^{(l)}) (\widehat{b}_{init}^{(l)} - b^{(l)}) \right| \lesssim (\|b^{(k)}\|_2 + 1) \frac{s \log p}{n} \right) \geq 1 - \min\{n, p\}^{-c}.$$

### G.3 Proof of Lemma 2

For  $\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}$  defined in (50), we express it as

$$\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)} = \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} + \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} \quad (165)$$

where  $\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)}$  and  $\mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)}$  defined in (51) and (52), respectively.

For  $\widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}$  defined in (24), we express it as

$$\widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)} = \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} + \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} \quad (166)$$

with

$$\begin{aligned} \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} &= \frac{\widehat{\sigma}_{l_1}^2}{|B_{l_1}|} (\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} [\widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = l_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = l_1)] \\ &\quad + \frac{\widehat{\sigma}_{k_1}^2}{|B_{k_1}|} (\widehat{u}^{(k_1, l_1)})^\top \widehat{\Sigma}^{(k_1)} [\widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = k_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = k_1)] \end{aligned} \quad (167)$$

$$\widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} = \frac{\sum_{i=1}^{N_{\mathbb{Q}}} \left( (\widehat{b}_{init}^{(l_1)})^\top X_i^{\mathbb{Q}} (\widehat{b}_{init}^{(k_1)})^\top X_i^{\mathbb{Q}} (\widehat{b}_{init}^{(l_2)})^\top X_i^{\mathbb{Q}} (\widehat{b}_{init}^{(k_2)})^\top X_i^{\mathbb{Q}} - (\widehat{b}_{init}^{(l_1)})^\top \bar{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k_1)} (\widehat{b}_{init}^{(l_2)})^\top \bar{\Sigma}^{\mathbb{Q}} \widehat{b}_{init}^{(k_2)} \right)}{|B| N_{\mathbb{Q}}} \quad (168)$$

where  $\bar{\Sigma}^{\mathbb{Q}} = \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} (X_i^{\mathbb{Q}})^\top$  and  $\widehat{\sigma}_l^2 = \|Y^{(l)} - X^{(l)} \widehat{b}^{(l)}\|_2^2 / n_l$  for  $1 \leq l \leq L$ .

The control of the event  $\mathcal{E}_1$  follows from the following high probability inequalities: with probability larger than  $1 - \exp(-cn) - \min\{N_{\mathbb{Q}}, p\}^{-c}$  for some positive constant  $c > 0$ ,

$$n \cdot \left| \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} - \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} \right| \leq C d_0 \left( \frac{s \log p}{n} + \sqrt{\frac{\log p}{n}} \right) \leq \frac{d_0}{4}. \quad (169)$$

$$N_{\mathbb{Q}} \cdot \left| \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} - \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} \right| \lesssim \log \max\{N_{\mathbb{Q}}, p\} \sqrt{\frac{s \log p \log N_{\mathbb{Q}}}{n}} + \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}}. \quad (170)$$

The proofs of (169) and (170) are presented in Sections G.3.1 and G.3.2, respectively.

We combine (169) and (170) and establish

$$\begin{aligned} \|\widehat{\mathbf{Cov}} - \mathbf{Cov}\|_2 &\lesssim \max_{(l_1, k_1), (l_2, k_2) \in \mathcal{I}_L} \left| \widehat{\mathbf{Cov}}_{\pi(l_1, k_1), \pi(l_2, k_2)} - \mathbf{Cov}_{\pi(l_1, k_1), \pi(l_2, k_2)} \right| \\ &\leq n \cdot \max_{(l_1, k_1), (l_2, k_2) \in \mathcal{I}_L} \left| \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} - \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} \right| \\ &\quad + n \cdot \max_{(l_1, k_1), (l_2, k_2) \in \mathcal{I}_L} \left| \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} - \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(b)} \right| \\ &\leq \frac{d_0}{4} + \frac{\sqrt{n \cdot s} [\log \max\{N_{\mathbb{Q}}, p\}]^2}{N_{\mathbb{Q}}} + \frac{n \cdot (\log N_{\mathbb{Q}})^{5/2}}{N_{\mathbb{Q}}^{3/2}} \leq d_0/2, \end{aligned}$$

where the first inequality holds for a finite  $L$  and the last inequality follows from Condition (A2).

### G.3.1 Proof of (169)

$$\begin{aligned} n \cdot \left| \widehat{\mathbf{V}}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} - \mathbf{V}_{\pi(l_1, k_1), \pi(l_2, k_2)}^{(a)} \right| &\lesssim \left| \widehat{\sigma}_{l_1}^2 - \sigma_{l_1}^2 \right| (\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} \left[ \widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = l_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = l_1) \right] \\ &\quad + \left| \widehat{\sigma}_{k_1}^2 - \sigma_{k_1}^2 \right| (\widehat{u}^{(k_1, l_1)})^\top \widehat{\Sigma}^{(k_1)} \left[ \widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = k_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = k_1) \right] \end{aligned} \quad (171)$$

Since

$$\begin{aligned} &\left| (\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} \left[ \widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = l_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = l_1) \right] \right| \\ &\leq \sqrt{(\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} \widehat{u}^{(l_1, k_1)} \cdot (\widehat{u}^{(l_1, k_2)})^\top \widehat{\Sigma}^{(l_1)} \widehat{u}^{(l_1, k_2)}} \\ &\quad + \sqrt{(\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} \widehat{u}^{(l_1, k_1)} \cdot (\widehat{u}^{(l_1, l_2)})^\top \widehat{\Sigma}^{(l_1)} \widehat{u}^{(l_1, l_2)}} \end{aligned} \quad (172)$$

we have

$$\left| (\widehat{u}^{(l_1, k_1)})^\top \widehat{\Sigma}^{(l_1)} \left[ \widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = l_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = l_1) \right] \right| \lesssim n \max_{(l, k) \in \mathcal{I}_L} \mathbf{V}_{\pi(l, k), \pi(l, k)}^{(a)} \lesssim d_0$$

Similarly, we have  $\left| (\widehat{u}^{(k_1, l_1)})^\top \widehat{\Sigma}^{(k_1)} \left[ \widehat{u}^{(l_2, k_2)} \mathbf{1}(l_2 = k_1) + \widehat{u}^{(k_2, l_2)} \mathbf{1}(k_2 = k_1) \right] \right| \lesssim d_0$ . Hence, on the event  $\mathcal{G}_3$ , we establish (169).

### G.3.2 Proof of (170)

Define

$$W_{i,1} = [b^{(l_1)}]^\top X_i^\mathbb{Q}, \quad W_{i,2} = [b^{(k_1)}]^\top X_i^\mathbb{Q}, \quad W_{i,3} = [b^{(l_2)}]^\top X_i^\mathbb{Q}, \quad W_{i,4} = [b^{(k_2)}]^\top X_i^\mathbb{Q};$$

and

$$\widehat{W}_{i,1} = (\widehat{b}_{init}^{(l_1)})^\top X_i^\mathbb{Q}, \quad \widehat{W}_{i,2} = (\widehat{b}_{init}^{(k_1)})^\top X_i^\mathbb{Q}, \quad \widehat{W}_{i,3} = (\widehat{b}_{init}^{(l_2)})^\top X_i^\mathbb{Q}, \quad \widehat{W}_{i,4} = (\widehat{b}_{init}^{(k_2)})^\top X_i^\mathbb{Q}.$$

With the above definitions, we have

$$\begin{aligned} & \mathbf{E}[b^{(l_1)}]^\top X_i^\mathbb{Q} [b^{(k_1)}]^\top X_i^\mathbb{Q} [b^{(l_2)}]^\top X_i^\mathbb{Q} [b^{(k_2)}]^\top X_i^\mathbb{Q} - (b^{(l_1)})^\top \Sigma^\mathbb{Q} b^{(k_1)} (b^{(l_2)})^\top \Sigma^\mathbb{Q} b^{(k_2)} \\ &= \mathbf{E} \prod_{t=1}^4 W_{i,t} - \mathbf{E} W_{i,1} W_{i,2} \cdot \mathbf{E} W_{i,3} W_{i,4} \end{aligned} \quad (173)$$

and

$$\begin{aligned} & \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \left( (\widehat{b}_{init}^{(l_1)})^\top X_i^\mathbb{Q} (\widehat{b}_{init}^{(k_1)})^\top X_i^\mathbb{Q} (\widehat{b}_{init}^{(l_2)})^\top X_i^\mathbb{Q} (\widehat{b}_{init}^{(k_2)})^\top X_i^\mathbb{Q} - (\widehat{b}_{init}^{(l_1)})^\top \bar{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k_1)} (\widehat{b}_{init}^{(l_2)})^\top \bar{\Sigma}^\mathbb{Q} \widehat{b}_{init}^{(k_2)} \right) \\ &= \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \prod_{t=1}^4 \widehat{W}_{i,t} - \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \widehat{W}_{i,1} \widehat{W}_{i,2} \cdot \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \widehat{W}_{i,3} \widehat{W}_{i,4} \end{aligned} \quad (174)$$

Hence, it is sufficient to control the following terms.

$$\begin{aligned} & \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \prod_{t=1}^4 \widehat{W}_{i,t} - \mathbf{E} \prod_{t=1}^4 W_{i,t} = \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \prod_{t=1}^4 \widehat{W}_{i,t} - \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \prod_{t=1}^4 W_{i,t} + \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \prod_{t=1}^4 W_{i,t} - \mathbf{E} \prod_{t=1}^4 W_{i,t} \\ & \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \widehat{W}_{i,1} \widehat{W}_{i,2} - \mathbf{E} W_{i,1} W_{i,2} = \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \widehat{W}_{i,1} \widehat{W}_{i,2} - \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} W_{i,1} W_{i,2} + \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} W_{i,1} W_{i,2} - \mathbf{E} W_{i,1} W_{i,2} \\ & \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \widehat{W}_{i,3} \widehat{W}_{i,4} - \mathbf{E} W_{i,3} W_{i,4} = \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \widehat{W}_{i,3} \widehat{W}_{i,4} - \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} W_{i,3} W_{i,4} + \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} W_{i,3} W_{i,4} - \mathbf{E} W_{i,3} W_{i,4} \end{aligned}$$

Specifically, we will show that, with probability larger than  $1 - \min\{N_\mathbb{Q}, p\}^{-c}$ ,

$$\left| \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} W_{i,1} W_{i,2} - \mathbf{E} W_{i,1} W_{i,2} \right| \lesssim \|b^{(l_1)}\|_2 \|b^{(k_1)}\|_2 \sqrt{\frac{\log N_\mathbb{Q}}{N_\mathbb{Q}}}, \quad (175)$$



$$\left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} W_{i,3} W_{i,4} - \mathbf{E} W_{i,3} W_{i,4} \right| \lesssim \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2 \sqrt{\frac{\log N_{\mathbb{Q}}}{N_{\mathbb{Q}}}}, \quad (176)$$

$$\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left( \prod_{t=1}^4 W_{i,t} - \mathbf{E} \prod_{t=1}^4 W_{i,t} \right) \lesssim \|b^{(l_1)}\|_2 \|b^{(k_1)}\|_2 \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2 \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}}, \quad (177)$$

$$\left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \widehat{W}_{i,1} \widehat{W}_{i,2} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} W_{i,1} W_{i,2} \right| \lesssim \sqrt{\frac{s \log p}{n}} \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_1)}\|_2 + \|b^{(k_1)}\|_2) + \sqrt{\frac{s \log p}{n}} \right), \quad (178)$$

$$\left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \widehat{W}_{i,3} \widehat{W}_{i,4} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} W_{i,3} W_{i,4} \right| \lesssim \sqrt{\frac{s \log p}{n}} \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_2)}\|_2 + \|b^{(k_2)}\|_2) + \sqrt{\frac{s \log p}{n}} \right). \quad (179)$$

If we further assume that  $\|b^{(l)}\|_2 \leq C$  for  $1 \leq l \leq L$  and  $s^2(\log p)^2/n \leq c$  for some positive constants  $C > 0$  and  $c > 0$ , then with probability larger than  $1 - \min\{N_{\mathbb{Q}}, p\}^{-c}$ ,

$$\left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \prod_{t=1}^4 \widehat{W}_{i,t} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \prod_{t=1}^4 W_{i,t} \right| \lesssim \log \max\{N_{\mathbb{Q}}, p\} \sqrt{\frac{s \log p \log N_{\mathbb{Q}}}{n}}. \quad (180)$$

By the expression (173) and (174), we establish (170) by applying (175), (176), (177), (178), (179), (180). In the following, we prove (175), (176) and (177). Then we will present the proofs of (178), (179), (180).

**Proofs of (175), (176) and (177).** We shall apply the following lemma to control the above terms, which re-states the Lemma 1 in Cai and Liu (2011).

**Lemma 10.** *Let  $\xi_1, \dots, \xi_n$  be independent random variables with mean 0. Suppose that there exists some  $c > 0$  and  $U_n$  such that  $\sum_{i=1}^n \mathbf{E} \xi_i^2 \exp(c|\xi_i|) \leq U_n^2$ . Then for  $0 < t \leq U_n$ ,*

$$\mathbf{P} \left( \sum_{i=1}^n \xi_i \geq C U_n t \right) \leq \exp(-t^2), \quad (181)$$

where  $C = c + c^{-1}$ .

Define

$$W_{i,1}^0 = \frac{[b^{(l_1)}]^\top X_i^{\mathbb{Q}}}{\sqrt{[b^{(l_1)}]^\top \Sigma^{\mathbb{Q}} b^{(l_1)}}}, \quad W_{i,2}^0 = \frac{[b^{(k_1)}]^\top X_i^{\mathbb{Q}}}{\sqrt{[b^{(k_1)}]^\top \Sigma^{\mathbb{Q}} b^{(k_1)}}}$$

and

$$W_{i,3}^0 = \frac{[b^{(l_2)}]^\top X_i^\mathbb{Q}}{\sqrt{[b^{(l_2)}]^\top \Sigma^\mathbb{Q} b^{(l_2)}}} \quad W_{i,4}^0 = \frac{[b^{(k_2)}]^\top X_i^\mathbb{Q}}{\sqrt{[b^{(k_2)}]^\top \Sigma^\mathbb{Q} b^{(k_2)}}}$$

Since  $X_i^\mathbb{Q}$  is sub-gaussian,  $W_{i,t}^0$  is sub-gaussian and both  $W_{i,1}^0 W_{i,2}^0$  and  $W_{i,3}^0 W_{i,4}^0$  are sub-exponential random variables, which follows from Remark 5.18 in [Vershynin \(2012\)](#). By Corollary 5.17 in [Vershynin \(2012\)](#), we have

$$\mathbf{P} \left( \left| \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} (W_{i,1}^0 W_{i,2}^0 - \mathbf{E} W_{i,1}^0 W_{i,2}^0) \right| \geq C \sqrt{\frac{\log N_\mathbb{Q}}{N_\mathbb{Q}}} \right) \leq 2N_\mathbb{Q}^{-c}$$

and

$$\mathbf{P} \left( \left| \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} (W_{i,3}^0 W_{i,4}^0 - \mathbf{E} W_{i,3}^0 W_{i,4}^0) \right| \geq C \sqrt{\frac{\log N_\mathbb{Q}}{N_\mathbb{Q}}} \right) \leq 2N_\mathbb{Q}^{-c}$$

where  $c$  and  $C$  are positive constants. The above inequalities imply (175) and (176) after rescaling.

For  $1 \leq t \leq 4$ , since  $W_{i,t}^0$  is a sub-gaussian random variable, there exist positive constants  $C_1 > 0$  and  $c > 2$  such that the following concentration inequality holds,

$$\sum_{i=1}^{N_\mathbb{Q}} \mathbf{P} \left( \max_{1 \leq t \leq 4} |W_{i,t}^0| \geq C_1 \sqrt{\log N_\mathbb{Q}} \right) \leq N_\mathbb{Q} \max_{1 \leq i \leq N_\mathbb{Q}} \mathbf{P} \left( \max_{1 \leq t \leq 4} |W_{i,t}^0| \geq C_1 \sqrt{\log N_\mathbb{Q}} \right) \lesssim N_\mathbb{Q}^{-c} \quad (182)$$

Define

$$H_{i,a} = \prod_{t=1}^4 W_{i,t}^0 \cdot \mathbf{1} \left( \max_{1 \leq t \leq 4} |W_{i,t}^0| \leq C_1 \sqrt{\log N_\mathbb{Q}} \right) \quad \text{for } 1 \leq t \leq 4,$$

and

$$H_{i,b} = \prod_{t=1}^4 W_{i,t}^0 \cdot \mathbf{1} \left( \max_{1 \leq t \leq 4} |W_{i,t}^0| \geq C_1 \sqrt{\log N_\mathbb{Q}} \right) \quad \text{for } 1 \leq t \leq 4.$$

Then we have

$$\frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} \prod_{t=1}^4 W_{i,t}^0 - \mathbf{E} \prod_{t=1}^4 W_{i,t}^0 = \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} (H_{i,a} - \mathbf{E} H_{i,a}) + \frac{1}{N_\mathbb{Q}} \sum_{i=1}^{N_\mathbb{Q}} (H_{i,b} - \mathbf{E} H_{i,b}) \quad (183)$$

By applying the Cauchy-Schwarz inequality, we bound  $\mathbf{E}H_{i,b}$  as

$$\begin{aligned} |\mathbf{E}H_{i,b}| &\leq \sqrt{\mathbf{E} \left( \prod_{t=1}^4 W_{i,t}^0 \right)^2 \mathbf{P} \left( \max_{1 \leq t \leq 4} |W_{i,t}^0| \geq C_1 \sqrt{\log N_{\mathbb{Q}}} \right)} \\ &\lesssim \mathbf{P} \left( |W_{i,t}^0| \geq C_1 \sqrt{\log N_{\mathbb{Q}}} \right)^{1/2} \lesssim N_{\mathbb{Q}}^{-1/2}, \end{aligned} \quad (184)$$

where the second and the last inequalities follow from the fact that  $W_{i,t}^0$  is a sub-gaussian random variable. Now we apply Lemma 10 to bound  $\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} (H_{i,a} - \mathbf{E}H_{i,a})$ . By taking  $c = c_1/(C_1^2 \log N_{\mathbb{Q}})^2$  for some small positive constant  $c_1 > 0$ , we have

$$\sum_{i=1}^{N_{\mathbb{Q}}} \mathbf{E} (H_{i,a} - \mathbf{E}H_{i,a})^2 \exp(c |H_{i,a} - \mathbf{E}H_{i,a}|) \leq C \sum_{i=1}^{N_{\mathbb{Q}}} \mathbf{E} (H_{i,a} - \mathbf{E}H_{i,a})^2 \leq C_2 N_{\mathbb{Q}}.$$

By applying Lemma 10 with  $U_n = \sqrt{C_2 N_{\mathbb{Q}}}$ ,  $c = c_1/(C_1^2 \log N_{\mathbb{Q}})^2$  and  $t = \sqrt{\log N_{\mathbb{Q}}}$ , then we have

$$\mathbf{P} \left( \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} (H_{i,a} - \mathbf{E}H_{i,a}) \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) \lesssim N_{\mathbb{Q}}^{-c}. \quad (185)$$

Note that

$$\begin{aligned} \mathbf{P} \left( \left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} H_{i,b} \right| \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) &\leq \mathbf{P} \left( \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} |H_{i,b}| \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) \\ &\leq \sum_{i=1}^{N_{\mathbb{Q}}} \mathbf{P} \left( |H_{i,b}| \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) \\ &\leq \sum_{i=1}^{N_{\mathbb{Q}}} \mathbf{P} \left( \max_{1 \leq t \leq 4} |W_{i,t}^0| \geq C_1 \sqrt{\log N_{\mathbb{Q}}} \right) \lesssim N_{\mathbb{Q}}^{-c} \end{aligned} \quad (186)$$

where the last inequality follows from (182).

By the decomposition (183), we have

$$\begin{aligned}
& \mathbf{P} \left( \left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left( \prod_{t=1}^4 W_{i,t}^0 - \mathbf{E} \prod_{t=1}^4 W_{i,t}^0 \right) \right| \geq 3C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) \\
& \leq \mathbf{P} \left( \left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} (H_{i,a} - \mathbf{E} H_{i,a}) \right| \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) \\
& + \mathbf{P} \left( |\mathbf{E} H_{i,b}| \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) + \mathbf{P} \left( \left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} H_{i,b} \right| \geq C \frac{(\log N_{\mathbb{Q}})^{5/2}}{\sqrt{N_{\mathbb{Q}}}} \right) \lesssim N_{\mathbb{Q}}^{-c}.
\end{aligned}$$

where the final upper bound follows from (184), (185) and (186). Hence, we establish that (177) holds with probability larger than  $1 - N_{\mathbb{Q}}^{-c}$ .

**Proofs (178), (179) and (180).** It follows from the definitions of  $\widehat{W}_{i,t}$  and  $W_{i,t}$  that

$$\begin{aligned}
& \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \widehat{W}_{i,1} \widehat{W}_{i,2} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} W_{i,1} W_{i,2} = [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^\top \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \\
& + [b^{(l_1)}]^\top \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] + [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^\top \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top b^{(k_1)}
\end{aligned} \tag{187}$$

On the event  $\mathcal{G}_2 \cap \mathcal{G}_5$  with  $\mathcal{G}_2$  defined in (106) and  $\mathcal{G}_5$  defined in (107), we establish (178).

By a similar argument, we establish (179). Furthermore, we define the event

$$\begin{aligned}
\mathcal{G}_7 &= \left\{ \max_{1 \leq l \leq L} \max_{1 \leq i \leq N_{\mathbb{Q}}} |X_i^{\mathbb{Q}} b^{(l)}| \lesssim (\sqrt{C_0} + \sqrt{\log N_{\mathbb{Q}}}) \|b^{(l)}\|_2 \right\} \\
\mathcal{G}_8 &= \left\{ \max_{1 \leq i \leq N_{\mathbb{Q}}} \|X_i^{\mathbb{Q}}\|_{\infty} \lesssim (\sqrt{C_0} + \sqrt{\log N_{\mathbb{Q}} + \log p}) \right\}
\end{aligned} \tag{188}$$

It follows from the assumption (A1) that  $\mathbf{P}(\mathcal{G}_7) \geq 1 - N_{\mathbb{Q}}^{-c}$  and  $\mathbf{P}(\mathcal{G}_8) \geq 1 - \min\{N_{\mathbb{Q}}, p\}^{-c}$  for some positive constant  $c > 0$ . Note that

$$\begin{aligned}
& \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| \widehat{W}_{i,1} \widehat{W}_{i,2} - W_{i,1} W_{i,2} \right| \leq \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^\top X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| \\
& + \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [b^{(l_1)}]^\top X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| + \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^\top X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top b^{(k_1)} \right|
\end{aligned} \tag{189}$$

By the Cauchy-Schwarz inequality, we have

$$\begin{aligned} & \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| \\ & \leq \frac{1}{N_{\mathbb{Q}}} \sqrt{\sum_{i=1}^{N_{\mathbb{Q}}} \left( [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} \right)^2 \sum_{i=1}^{N_{\mathbb{Q}}} \left( [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right)^2} \end{aligned}$$

Hence, on the event  $\mathcal{G}_1 \cap \mathcal{G}_6(\widehat{b}_{init}^{(k_1)} - b^{(k_1)}, \widehat{b}_{init}^{(k_1)} - b^{(k_1)}, \sqrt{\log p}) \cap \mathcal{G}_6(\widehat{b}_{init}^{(l_1)} - b^{(l_1)}, \widehat{b}_{init}^{(l_1)} - b^{(l_1)}, \sqrt{\log p})$ ,

$$\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| \lesssim \frac{s \log p}{n} \quad (190)$$

On the event  $\mathcal{G}_7$ , we have

$$\begin{aligned} & \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| \\ & \lesssim (\sqrt{C_0} + \sqrt{\log N_{\mathbb{Q}}}) \|b^{(l_1)}\|_2 \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| \\ & \leq (\sqrt{C_0} + \sqrt{\log N_{\mathbb{Q}}}) \|b^{(l_1)}\|_2 \frac{1}{\sqrt{N_{\mathbb{Q}}}} \sqrt{\sum_{i=1}^{N_{\mathbb{Q}}} \left( [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right)^2} \end{aligned}$$

where the last inequality follows from the Cauchy-Schwarz inequality. Hence, on the event  $\mathcal{G}_1 \cap \mathcal{G}_7 \cap \mathcal{G}_6(\widehat{b}_{init}^{(k_1)} - b^{(k_1)}, \widehat{b}_{init}^{(k_1)} - b^{(k_1)}, \sqrt{\log p})$ , we establish

$$\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^{\top} [\widehat{b}_{init}^{(k_1)} - b^{(k_1)}] \right| \lesssim (\sqrt{C_0} + \sqrt{\log N_{\mathbb{Q}}}) \|b^{(l_1)}\|_2 \sqrt{\frac{s \log p}{n}}. \quad (191)$$

Similarly, we establish

$$\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| [\widehat{b}_{init}^{(l_1)} - b^{(l_1)}]^{\top} X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^{\top} b^{(k_1)} \right| \lesssim (\sqrt{C_0} + \sqrt{\log N_{\mathbb{Q}}}) \|b^{(k_1)}\|_2 \sqrt{\frac{s \log p}{n}}.$$

Combined with (189), (190) and (191), we establish

$$\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| \widehat{W}_{i,1} \widehat{W}_{i,2} - W_{i,1} W_{i,2} \right| \lesssim \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_1)}\|_2 + \|b^{(k_1)}\|_2) + \sqrt{\frac{s \log p}{n}} \right) \sqrt{\frac{s \log p}{n}} \quad (192)$$

Similarly, we establish

$$\frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \left| \widehat{W}_{i,3} \widehat{W}_{i,4} - W_{i,3} W_{i,4} \right| \lesssim \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_2)}\|_2 + \|b^{(k_2)}\|_2) + \sqrt{\frac{s \log p}{n}} \right) \sqrt{\frac{s \log p}{n}} \quad (193)$$

Define  $H_{i,1} = W_{i,1} W_{i,2}$ ,  $H_{i,2} = W_{i,3} W_{i,4}$ ,  $\widehat{H}_{i,1} = \widehat{W}_{i,1} \widehat{W}_{i,2}$  and  $\widehat{H}_{i,2} = \widehat{W}_{i,3} \widehat{W}_{i,4}$ . Then we have

$$\begin{aligned} & \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \prod_{t=1}^4 \widehat{W}_{i,t} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \prod_{t=1}^4 W_{i,t} = \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \widehat{H}_{i,1} \widehat{H}_{i,2} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} H_{i,1} H_{i,2} \\ &= \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} (\widehat{H}_{i,1} - H_{i,1}) H_{i,2} + \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} (\widehat{H}_{i,2} - H_{i,2}) H_{i,1} + \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} (\widehat{H}_{i,1} - H_{i,1}) (\widehat{H}_{i,2} - H_{i,2}) \end{aligned} \quad (194)$$

On the event  $\mathcal{G}_7$ , we have

$$|H_{i,1}| \lesssim (C_0 + \log N_{\mathbb{Q}}) \|b^{(l_1)}\|_2 \|b^{(k_1)}\|_2 \quad \text{and} \quad |H_{i,2}| \lesssim (C_0 + \log N_{\mathbb{Q}}) \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2 \quad (195)$$

On the event  $\mathcal{G}_7 \cap \mathcal{G}_8$ , we have

$$\begin{aligned} & \left| \widehat{H}_{i,2} - H_{i,2} \right| \leq \left| [\widehat{b}_{init}^{(l_2)} - b^{(l_2)}]^\top X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top [\widehat{b}_{init}^{(k_2)} - b^{(k_2)}] \right| \\ &+ \left| [b^{(l_2)}]^\top X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top [\widehat{b}_{init}^{(k_2)} - b^{(k_2)}] \right| + \left| [\widehat{b}_{init}^{(l_2)} - b^{(l_2)}]^\top X_i^{\mathbb{Q}} [X_i^{\mathbb{Q}}]^\top b^{(k_2)} \right| \\ &\lesssim (C_0 + \log N_{\mathbb{Q}} + \log p) \left( s^2 \frac{\log p}{n} + s \sqrt{\frac{\log p}{n}} \|b^{(k_2)}\|_2 + s \sqrt{\frac{\log p}{n}} \|b^{(l_2)}\|_2 \right) \end{aligned} \quad (196)$$

By the decomposition (194), we combine (195), (196), (192) and (193) and establish

$$\begin{aligned}
& \left| \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \prod_{t=1}^4 \widehat{W}_{i,t} - \frac{1}{N_{\mathbb{Q}}} \sum_{i=1}^{N_{\mathbb{Q}}} \prod_{t=1}^4 W_{i,t} \right| \\
& \leq (C_0 + \log N_{\mathbb{Q}}) \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2 \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_1)}\|_2 + \|b^{(k_1)}\|_2) + \sqrt{\frac{s \log p}{n}} \right) \sqrt{\frac{s \log p}{n}} \\
& + (C_0 + \log N_{\mathbb{Q}}) \|b^{(l_1)}\|_2 \|b^{(k_1)}\|_2 \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_2)}\|_2 + \|b^{(k_2)}\|_2) + \sqrt{\frac{s \log p}{n}} \right) \sqrt{\frac{s \log p}{n}} \\
& + \left( \sqrt{\log N_{\mathbb{Q}}} (\|b^{(l_1)}\|_2 + \|b^{(k_1)}\|_2) + \sqrt{\frac{s \log p}{n}} \right) \sqrt{\frac{s \log p}{n}} \\
& \cdot (C_0 + \log N_{\mathbb{Q}} + \log p) \left( s^2 \frac{\log p}{n} + s \sqrt{\frac{\log p}{n}} \|b^{(k_2)}\|_2 + s \sqrt{\frac{\log p}{n}} \|b^{(l_2)}\|_2 + \|b^{(l_2)}\|_2 \|b^{(k_2)}\|_2 \right)
\end{aligned} \tag{197}$$

If we further assume that  $\|b^{(l)}\|_2 \leq C$  for  $1 \leq l \leq L$  and  $s^2(\log p)^2/n \leq c$  for some positive constants  $C > 0$  and  $c > 0$ , then we establish (180).

## H Additional Simulation

### H.1 Additional Simulation settings

In addition to the simulation settings in Section 6.2 in the main paper, we consider the following simulation settings to evaluate the finite-sample performance of our proposed method.

**Setting 2 ( $L = 2$  with covariate shift).**  $b^{(1)}$  and  $b^{(2)}$  are the same as setting 1, except for  $b_{498}^{(1)} = 0.5$ ,  $b_j^{(1)} = -0.5$  for  $j = 499, 500$ , and  $b_{500}^{(2)} = 1$ .  $[x_{\text{new}}]_j = 1$  for  $498 \leq j \leq 500$ , and  $[x_{\text{new}}]_j = 0$  otherwise.  $\Sigma_{i,i}^{\mathbb{Q}} = 1.5$  for  $1 \leq i \leq 500$ ,  $\Sigma_{i,j}^{\mathbb{Q}} = 0.9$  for  $1 \leq i \neq j \leq 5$ ,  $\Sigma_{i,j}^{\mathbb{Q}} = 0.9$  for  $499 \leq i \neq j \leq 500$  and  $\Sigma_{i,j}^{\mathbb{Q}} = \Sigma_{i,j}$  otherwise.

**Setting 5 (perturbation setting, no covariate shift with  $L = 2$ ):**  $b_j^{(1)} = j/40$  for  $1 \leq j \leq 10$ ,  $b_j^{(1)} = (10 - j)/40$  for  $11 \leq j \leq 20$ ,  $b_j^{(1)} = 0.2$  for  $j = 21$ ,  $b_j^{(1)} = 1$  for  $j = 22, 23$ ;  $b_j^{(2)} = b_j^{(1)} + \text{perb}/\sqrt{300}$  for  $1 \leq j \leq 10$ ,  $b_j^{(2)} = 0$  for  $11 \leq j \leq 20$ ,  $b_j^{(2)} = 0.5$  for  $j = 21$ ,  $b_j^{(2)} = 0.2$  for  $j = 22, 23$ . We vary the values of perb across  $\{1, 1.125, 1.25, 1.5, 3.75, 4, 5, 7, 10, 12\}$ .  $[x_{\text{new}}]_j = j/5$  for  $1 \leq j \leq 5$ .

**Setting 6 (opposite effect, no covariate shift with  $L = 2$ ):**

(6a)  $b^{(l)}$  for  $1 \leq l \leq 2$  are almost the same as setting 1, except for  $b_j^{(1)} = 0$  for  $j = 499$ ,  $b_j^{(1)} = 0.2$  for  $j = 500$ ,  $b_j^{(2)} = -0.2$  for  $j = 500$ .  $[x_{\text{new}}]_j = 1$  for  $j = 500$ .

(6b) Same as 6(a) except for  $b_j^{(2)} = -0.4$  for  $j = 500$ .

**Setting 7 (covariate shift with  $L = 5$ ):**

(7a)  $b_j^{(1)} = j/10$  for  $1 \leq j \leq 10$ ,  $b_j^{(1)} = (10 - j)/10$  for  $11 \leq j \leq 20$ ,  $b_j^{(1)} = 1/5$  for  $j = 21$ ,  $b_j^{(1)} = 1$  for  $j = 22, 23$ ; For  $2 \leq l \leq L$ ,  $b_j^{(l)} = b_j^{(1)} + 0.1 \cdot (l - 1)/\sqrt{300}$  for  $1 \leq j \leq 10$ ,  $b_j^{(2)} = -0.3 \cdot (l - 1)/\sqrt{300}$  for  $11 \leq j \leq 20$ ,  $b_j^{(l)} = 0.5 \cdot (l - 1)$  for  $j = 21$ ,  $b_j^{(l)} = 0.2 \cdot (j - 1)$  for  $j = 22, 23$ ;  $[x_{\text{new}}]_j = 1$  for  $21 \leq j \leq 23$ ;  $\Sigma_{i,i}^{\mathbb{Q}} = 1.1$  for  $1 \leq i \leq 500$ ,  $\Sigma_{i,j}^{\mathbb{Q}} = 0.75$  for  $1 \leq i \neq j \leq 6$  and  $\Sigma_{i,j}^{\mathbb{Q}} = \Sigma_{i,j}$  otherwise.

(7b)  $b^{(l)}$  for  $1 \leq l \leq 2$  and  $\Sigma^{\mathbb{Q}}$  are the same as (a); for  $l \geq 3$ ,  $\{b_j^{(l)}\}_{1 \leq j \leq 6}$  are independently generated following standard normal and  $b_j^{(l)} = 0$  for  $7 \leq j \leq 500$ ;  $x_{\text{new}} \sim \mathcal{N}(\mathbf{0}, \Sigma^{\text{new}})$  with  $\Sigma_{i,j}^{\text{new}} = 0.5^{1+|i-j|}/25$  for  $1 \leq i, j \leq 500$ .

## H.2 Choice of $\delta$ : Instability Measure and Reward Plots

In Table S2, we report the instability measure  $\mathbb{I}(\delta)$  for settings 2, 3, 6, and 7 detailed in Section 6.2 and Section H.1 in the supplement. setting 5 is not included since it consists of more than 10 subsettings and all other settings have been reported in Table 1 in the main paper.

setting	$\delta = 0$	$\delta = 0.1$	$\delta = 0.5$	$\delta = 1$	$\delta = 2$	$\Gamma_{11}^{\mathbb{Q}} + \Gamma_{22}^{\mathbb{Q}} - 2\Gamma_{12}^{\mathbb{Q}}$
2	0.023	0.021	0.014	0.010	0.006	4.635
3(a)	0.094	0.082	0.044	0.023	0.010	1.935
3(b)	0.058	0.050	0.029	0.017	0.008	2.810
4(a)	0.108	0.087	0.044	0.024	0.011	2.007
4(b)	0.076	0.059	0.027	0.014	0.006	$L = 5$
4(c)	0.052	0.039	0.016	0.008	0.003	$L = 10$
6(a)	3.305	1.449	0.221	0.065	0.018	0.160
6(b)	1.451	0.816	0.168	0.056	0.017	0.360
7(a)	0.007	0.007	0.007	0.005	0.003	$L = 5$
7(b)	0.005	0.005	0.004	0.003	0.001	$L = 5$

Table S2: The instability measure  $\mathbb{I}(\delta)$  for settings 2, 3, 4, 6, and 7. The reported values are averaged over 100 repeated simulations.



In Figure S1, we demonstrate the effect of  $\delta$  on the reward across simulation settings with large  $\mathbb{I}(0)$  values. For these settings, if the stable aggregation is an important goal, our recommended value is  $\delta = 2$  as  $R_{\mathbb{Q}}(\hat{\gamma}_{\delta})/R_{\mathbb{Q}}(\hat{\gamma}_{\delta=0})$  is above 95% even for  $\delta = 2$ .

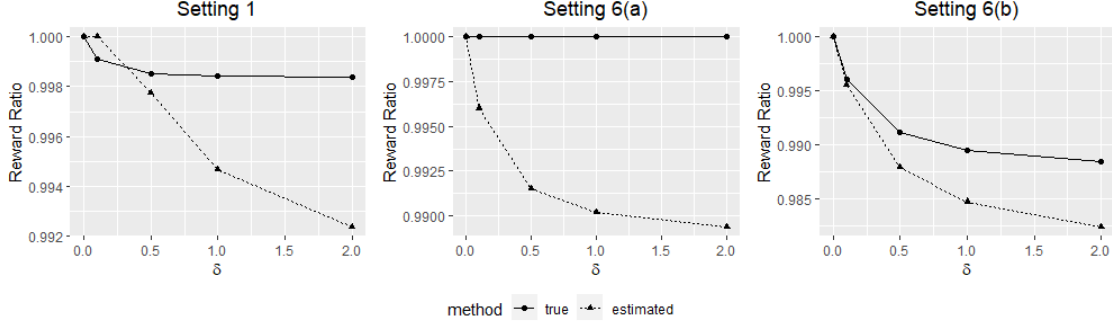


Figure S1: The reward ratio with respect to different  $\delta$  values. The true and estimated reward ratios respectively denote  $R_{\mathbb{Q}}(\hat{\gamma}_{\delta})/R_{\mathbb{Q}}(\hat{\gamma}_{\delta=0})$  and  $\hat{R}(\hat{\gamma}_{\delta})/\hat{R}(\hat{\gamma}_{\delta=0})$ , where  $\hat{\gamma}_{\delta}$  and  $\hat{R}(\hat{\gamma}_{\delta})$  are estimated based on a single data set with  $n = 500$  and  $N_{\mathbb{Q}} = 2000$ .

In Figure S2, we report the reward ratios for extra settings detailed in Section 6.2 and Section H.1. Since  $\mathbb{I}(0)$  is small in these settings, the maximin effect without the ridge penalty is already a stable aggregation.

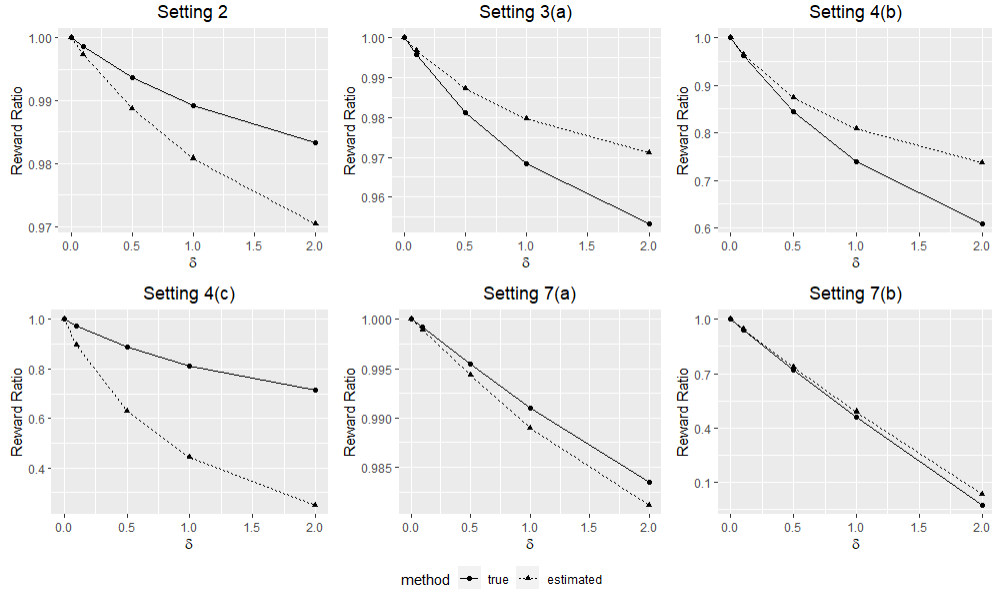


Figure S2: The reward ratio with respect to different  $\delta$  values. The true and estimated reward ratios respectively denote  $R_{\mathbb{Q}}(\hat{\gamma}_{\delta})/R_{\mathbb{Q}}(\hat{\gamma}_{\delta=0})$  and  $\hat{R}(\hat{\gamma}_{\delta})/\hat{R}(\hat{\gamma}_{\delta=0})$ , where  $\hat{\gamma}_{\delta}$  and  $\hat{R}(\hat{\gamma}_{\delta})$  are estimated based on a single data set with  $n = 500$  and  $N_{\mathbb{Q}} = 2,000$ .

### H.3 Simulation Results for Settings 2, 4, 5, 6 and 7

The results for settings 2, 4, 5, 6 and 7 are similar to those presented in Section 6.2 in the main paper. Our proposed CIs achieve the desired coverage level and the intervals become shorter with a larger  $n$  or  $\delta$ . We mainly highlight three interesting observations and then present other simulation results.

**Opposite effect setting.** In Figure S3, we investigate the opposite effect setting, which is setting 6. Note that the maximin effect tends to shrink the opposite effects to zero (Meinshausen and Bühlmann, 2015). For example, in setting 6,  $b_{500}^{(1)}$  and  $b_{500}^{(2)}$  have opposite signs and the maximin effect is zero for  $\delta = 0$ . We use the Empirical Rejection Rate (ERR) to denote the proportion of rejecting the null hypothesis out of 500 simulations. In Table S3, we observe that ERR is below 5% for  $\delta = 0$ , which indicate that the corresponding maximin effect is not significant.

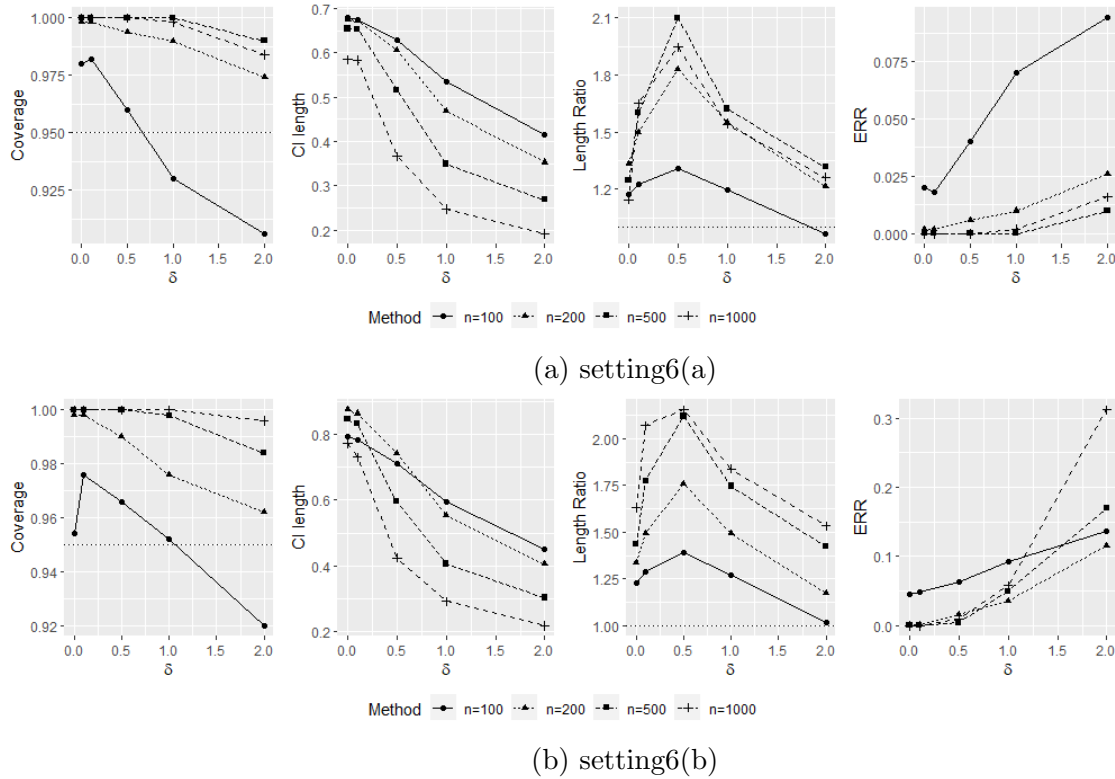


Figure S3: setting 6 with opposite effects. “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34); “ERR” represents the empirical rejection rate out of 500 simulations.

**Dependence of RMSE on  $\delta$  and  $n$ .** For settings 1 and 2, the Root Mean Square Error (RMSE) of the point estimator  $\widehat{x_{\text{new}}^T \beta_\delta^*}$  in (31) is reported in Table S3.

	setting 1					setting 2				
n	$\delta = 0$	$\delta = 0.1$	$\delta = 0.5$	$\delta = 1$	$\delta = 2$	$\delta = 0$	$\delta = 0.1$	$\delta = 0.5$	$\delta = 1$	$\delta = 2$
100	0.137	0.136	0.134	0.131	0.129	0.271	0.260	0.227	0.202	0.178
200	0.129	0.104	0.089	0.086	0.084	0.183	0.177	0.157	0.142	0.126
300	0.125	0.088	0.072	0.069	0.068	0.137	0.132	0.118	0.107	0.094
500	0.120	0.077	0.058	0.056	0.056	0.108	0.104	0.093	0.084	0.074

Table S3: Root Mean Square Error of  $\widehat{x_{\text{new}}^T \beta_\delta^*}$  for settings 1 and 2.

For setting 1, the RMSE decreases with an increasing  $n$  or a larger  $\delta$  value. For  $n = 300, 500$ , the RMSE for  $\delta = 2$  is round half of that for  $\delta = 0$ . For setting 2, the decrease of the RMSE with respect to  $\delta$  is not as much as that of setting 1, since the weight optimization in setting 2 is much more stable than that in setting 1.

**Higher dimension  $p$ .** In Figure S4, we have explored our proposed method for a larger  $p$  value and our proposed CIs are still valid for  $p = 1000, 2000$  and  $3000$ .

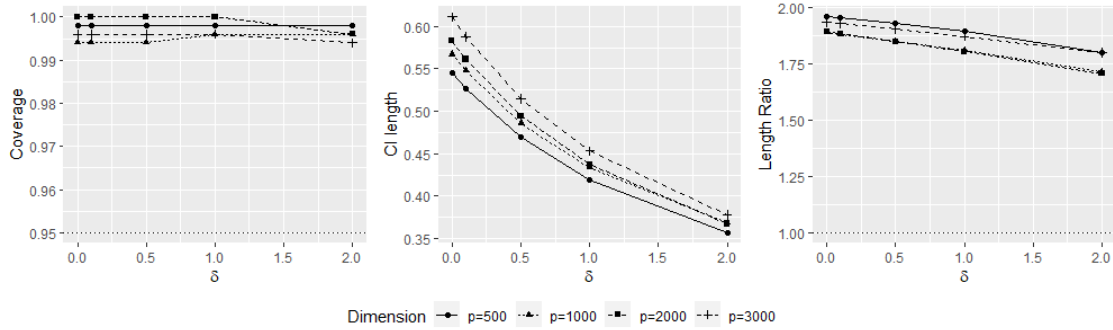
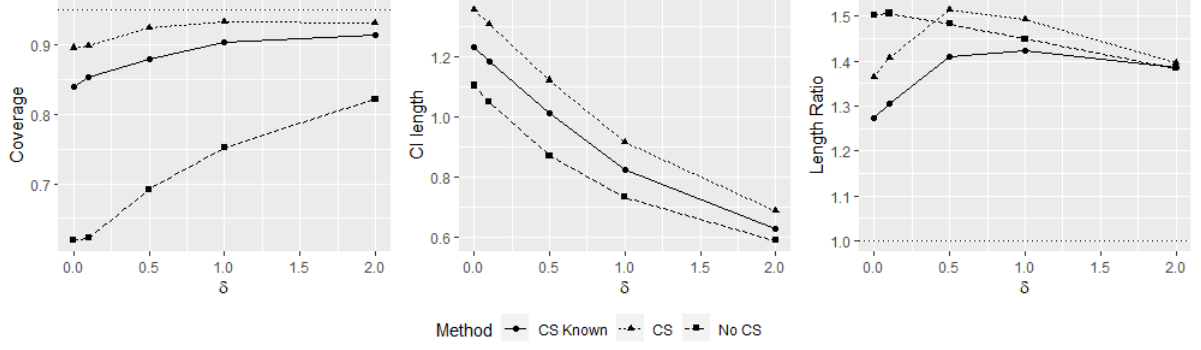
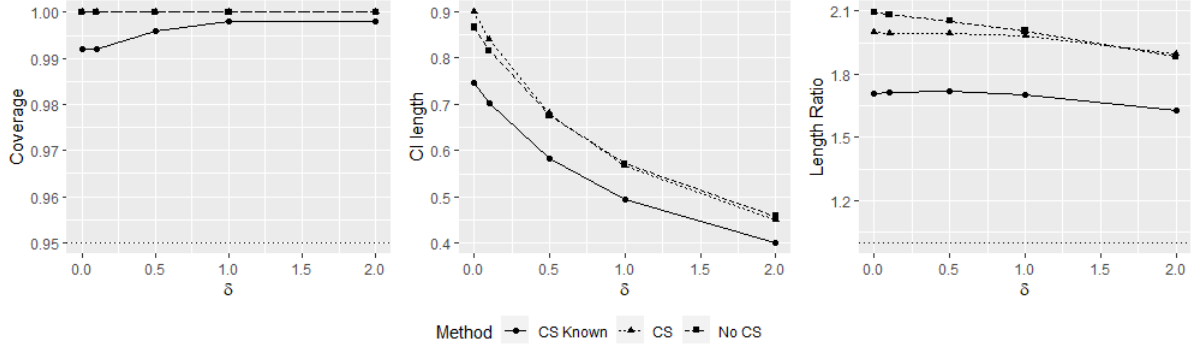


Figure S4: Dependence on  $\delta$  and  $p$ : setting 2 (covariate shift) with  $n = 500$ . “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).

**Additional results for settings 3, 4, 5 and 7.** We report the results for setting 3 with  $n = 200$  in Figure S5. We report the results for settings 4(a) with  $L = 2$  and 4(c) with  $L = 10$  in Figure S6. We report the results for settings 5 and 7 in Figures S7 and S8, respectively.

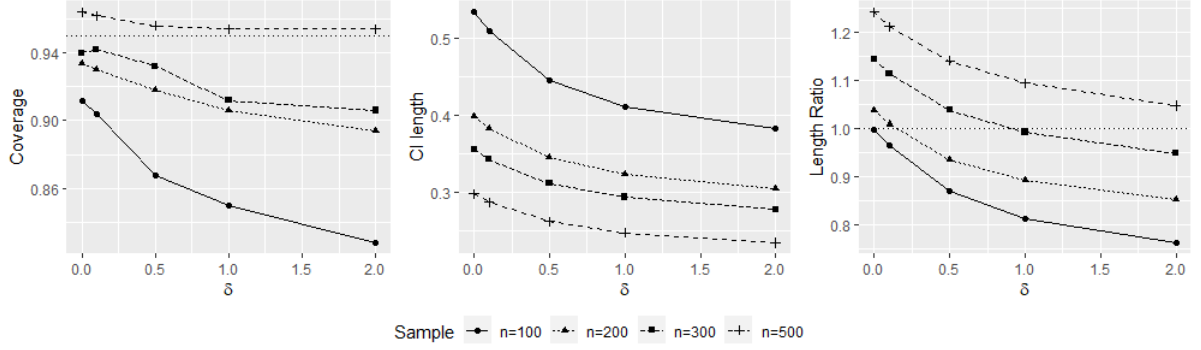


(a) setting 3(a) with covariate shift

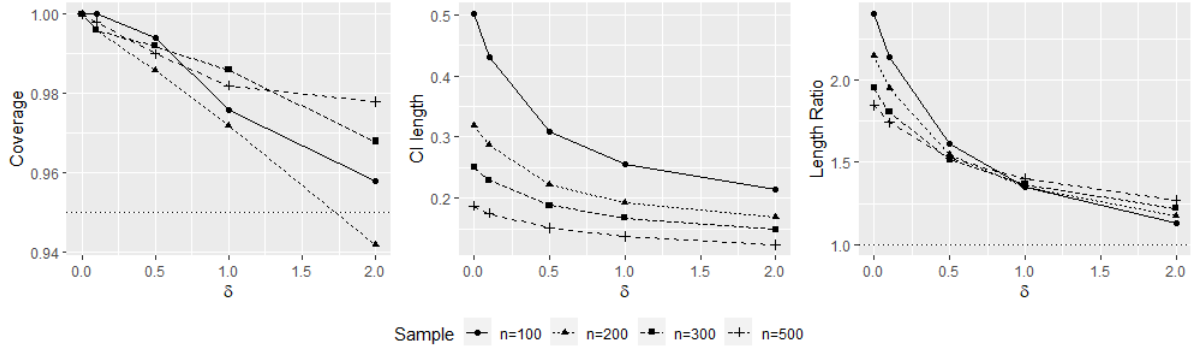


(b) setting 3(b) with no covariate shift

Figure S5: Comparison of covariate shift and no covariate shift algorithms with  $n = 200$ . “CS Known”, “CS” and “No CS” represent Algorithm 1 with known  $\Sigma^Q$ , Algorithm 1 with covariate shift but unknown  $\Sigma^Q$ , and Algorithm 1 with no covariate shift, respectively. “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).



(a) setting 4(a) with  $L = 2$



(b) setting 4(c) with  $L = 10$

Figure S6: Dependence on  $\delta$  and  $n$ . “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).

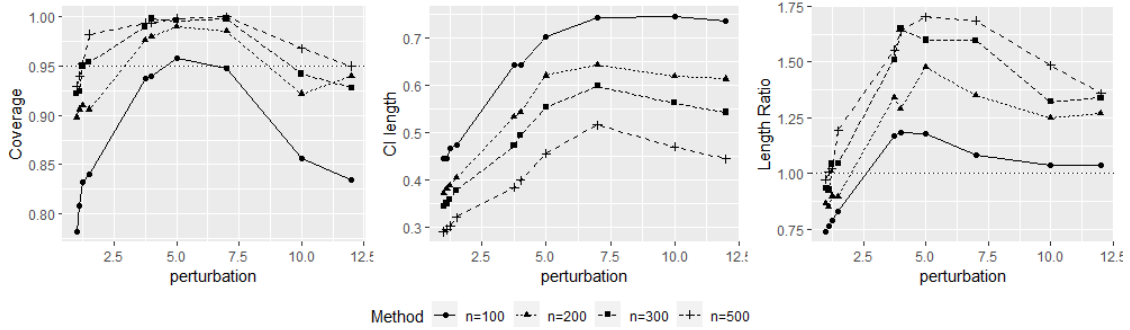


Figure S7: Dependence on the value perb for setting 5 (the perturbation setting). “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).

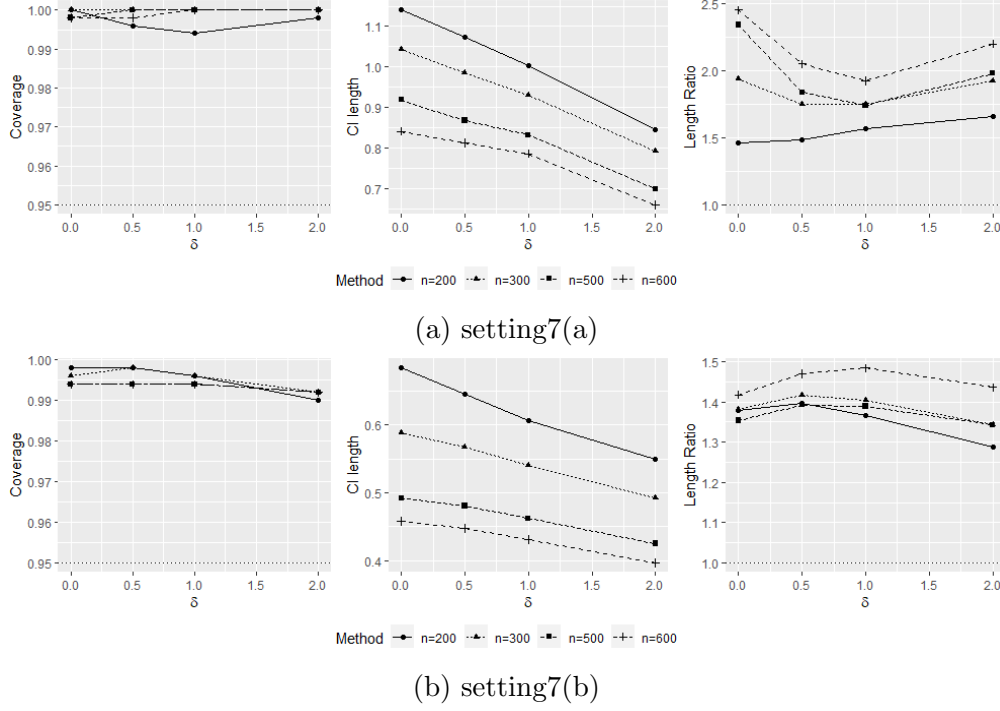


Figure S8: Dependence on  $\delta$  and  $n$ : setting 7 (covariate shift with  $L=5$ ). “Coverage” and “CI Length” stand for the empirical coverage and the average length of our proposed CI, respectively; “Length Ratio” represents the ratio of the average length of our proposed CI to the normality CI in (34).

#### H.4 Comparison of Different Threshold Levels

Our proposed CI in (21) relies on the index set  $\mathbb{M}$  in (20). As an alternative, we define

$$\text{CI}'_{\alpha}(x_{\text{new}}^{\top}\beta_{\delta}^{*}) = \cup_{m \in \mathbb{M}'} \text{Int}_{\alpha}^{[m]}(x_{\text{new}}), \quad (198)$$

where  $\mathbb{M}'$  is defined in (22). The CI in (198) uses a different index set  $\mathbb{M}'$  in (22). Note that both  $\mathbb{M}$  and  $\mathbb{M}'$  depend on a pre-specified value  $\alpha_0$  with the default value  $\alpha_0 = 0.01$ . Furthermore, we explore the performance of constructed CIs with  $\alpha_0 = 0.05$ . We compare the following CIs,

(V0) We do not screen out any sampled data:

$$\text{CI}_{\alpha}(x_{\text{new}}^{\top}\beta_{\delta}^{*}) = \cup_{m=1}^M \text{Int}_{\alpha}^{[m]}(x_{\text{new}}). \quad (199)$$

(V1-1) CI in (21) where  $\mathbb{M}$  is defined in (20) with  $\alpha_0 = 0.01$ .

(V1-5) CI in (21) where  $\mathbb{M}$  is defined in (20) with  $\alpha_0 = 0.05$ .

(V2-1) CI in (198) where  $\mathbb{M}'$  is defined in (22) with  $\alpha_0 = 0.01$ .

(V2-5) CI in (198) where  $\mathbb{M}'$  is defined in (22) with  $\alpha_0 = 0.05$ .

We compare the above five methods and the normality CI in (34) over settings 1 and 2. The results for settings 1 and 2 are reported in Tables S4 and S5, respectively. We report the empirical coverage and the length ratio for  $n \in \{100, 200, 300, 500\}$  and  $\delta \in \{0, 0.1, 0.5, 1, 2\}$ , where all length ratios are defined as the ratios with respect to the length of the normality CI. For different choices of the index sets, our proposed CIs have a similar performance in terms of both empirical coverage and length ratios. In general, our proposed CIs with the index set  $\mathbb{M}'$  are slightly shorter than those with the index set  $\mathbb{M}$ . Furthermore, our proposed CIs get shorter with a larger value of  $\alpha_0$ .

$n$	$\delta$	Coverage						Length Ratio				
		normality	V0	V1-1	V1-5	V2-1	V2-5	V0	V1-1	V1-5	V2-1	V2-5
100	0.0	0.948	0.980	0.980	0.980	0.978	0.974	1.189	1.188	1.188	1.186	1.176
	0.1	0.932	0.962	0.962	0.962	0.962	0.958	1.231	1.230	1.228	1.224	1.208
	0.5	0.920	0.944	0.942	0.940	0.928	0.920	1.275	1.270	1.257	1.244	1.196
	1.0	0.920	0.902	0.898	0.896	0.892	0.878	1.159	1.148	1.128	1.118	1.057
	2.0	0.918	0.838	0.836	0.832	0.830	0.824	0.931	0.924	0.912	0.905	0.876
200	0.0	0.914	0.980	0.980	0.980	0.980	0.980	1.248	1.248	1.248	1.248	1.245
	0.1	0.946	0.996	0.996	0.996	0.996	0.996	1.345	1.345	1.344	1.343	1.336
	0.5	0.942	0.986	0.986	0.984	0.986	0.982	1.511	1.501	1.480	1.463	1.388
	1.0	0.928	0.960	0.960	0.952	0.950	0.942	1.273	1.260	1.237	1.216	1.152
	2.0	0.926	0.926	0.924	0.924	0.920	0.914	1.019	1.013	1.002	0.997	0.970
300	0.0	0.886	0.972	0.972	0.972	0.972	0.970	1.259	1.259	1.259	1.259	1.258
	0.1	0.946	1.000	1.000	1.000	1.000	1.000	1.416	1.416	1.415	1.414	1.410
	0.5	0.948	0.996	0.996	0.994	0.994	0.988	1.600	1.592	1.570	1.555	1.458
	1.0	0.942	0.968	0.968	0.968	0.964	0.958	1.280	1.270	1.249	1.233	1.167
	2.0	0.938	0.940	0.940	0.940	0.938	0.934	1.042	1.038	1.029	1.025	0.999
500	0.0	0.860	0.966	0.966	0.966	0.966	0.966	1.294	1.294	1.294	1.294	1.293
	0.1	0.934	0.998	0.998	0.998	0.998	0.998	1.528	1.528	1.528	1.527	1.521
	0.5	0.958	0.996	0.996	0.996	0.996	0.996	1.723	1.709	1.678	1.655	1.541
	1.0	0.962	0.984	0.984	0.984	0.982	0.980	1.297	1.287	1.266	1.256	1.201
	2.0	0.962	0.974	0.974	0.974	0.974	0.972	1.074	1.070	1.061	1.058	1.036

Table S4: Comparison of our proposed CIs with different index sets for setting 1.

$n$	$\delta$	Coverage						Length Ratio				
		normality	V0	V1-1	V1-5	V2-1	V2-5	V0	V1-1	V1-5	V2-1	V2-5
100	0.0	0.910	0.928	0.920	0.904	0.910	0.890	1.194	1.165	1.114	1.124	1.032
	0.1	0.908	0.930	0.922	0.906	0.914	0.892	1.200	1.171	1.120	1.132	1.037
	0.5	0.912	0.922	0.918	0.904	0.914	0.890	1.195	1.165	1.114	1.132	1.041
	1.0	0.920	0.922	0.914	0.898	0.906	0.880	1.150	1.122	1.075	1.094	1.014
	2.0	0.930	0.906	0.902	0.894	0.898	0.884	1.054	1.031	0.995	1.011	0.951
200	0.0	0.906	0.976	0.972	0.966	0.972	0.956	1.554	1.503	1.433	1.452	1.318
	0.1	0.904	0.976	0.972	0.966	0.972	0.958	1.546	1.495	1.425	1.444	1.313
	0.5	0.902	0.976	0.972	0.968	0.972	0.960	1.494	1.447	1.382	1.404	1.289
	1.0	0.902	0.970	0.968	0.962	0.968	0.960	1.435	1.394	1.335	1.358	1.259
	2.0	0.912	0.974	0.968	0.958	0.966	0.950	1.335	1.302	1.253	1.275	1.196
300	0.0	0.898	0.994	0.990	0.990	0.992	0.982	1.786	1.733	1.645	1.667	1.520
	0.1	0.900	0.996	0.990	0.990	0.992	0.982	1.774	1.723	1.636	1.659	1.516
	0.5	0.910	0.996	0.992	0.992	0.994	0.984	1.729	1.682	1.603	1.628	1.500
	1.0	0.910	0.994	0.992	0.990	0.990	0.984	1.675	1.632	1.562	1.586	1.473
	2.0	0.914	0.988	0.988	0.986	0.988	0.980	1.569	1.534	1.476	1.498	1.407
500	0.0	0.916	0.992	0.992	0.990	0.992	0.990	1.809	1.767	1.686	1.706	1.564
	0.1	0.916	0.992	0.992	0.990	0.992	0.990	1.804	1.762	1.683	1.704	1.566
	0.5	0.920	0.994	0.994	0.992	0.992	0.992	1.785	1.746	1.672	1.694	1.568
	1.0	0.918	0.994	0.994	0.994	0.994	0.992	1.757	1.721	1.653	1.676	1.561
	2.0	0.916	0.992	0.992	0.992	0.992	0.992	1.688	1.658	1.600	1.621	1.525

Table S5: Comparison of our proposed CIs with different index sets for setting 2.

## H.5 Sample Splitting Comparison

We compare the algorithm with and without sample splitting and report the results in Table S6. For the sample splitting algorithm, we split the samples into two equal size sub-samples. For  $n = 100$ , no sample splitting algorithm is slightly under-coverage (the empirical coverage level is still above 90%). When  $n \geq 200$ , both the algorithm with and without sample splitting achieve the desired coverage levels. As expected, the CIs with sample splitting are longer than those without sample splitting. In Table S6, under the column indexed with “Length ratio”, we report the ratio of the average length of CI with sample splitting to that without sample splitting.



		Coverage		Length		
$\delta$	$n$	Splitting	No Splitting	Splitting	No Splitting	Length ratio
0.0	100	0.978	0.920	1.921	1.062	1.809
	200	0.994	0.972	1.618	0.898	1.802
	300	0.994	0.990	1.336	0.789	1.693
	500	1.000	0.992	0.999	0.651	1.534
0.1	100	0.980	0.922	1.904	1.023	1.861
	200	0.994	0.972	1.567	0.863	1.816
	300	0.994	0.990	1.286	0.759	1.695
	500	1.000	0.992	0.952	0.629	1.514
0.5	100	0.982	0.918	1.814	0.890	2.038
	200	0.996	0.972	1.373	0.748	1.836
	300	0.996	0.992	1.099	0.666	1.651
	500	1.000	0.994	0.800	0.559	1.433
1.0	100	0.986	0.914	1.639	0.769	2.131
	200	0.998	0.968	1.155	0.655	1.764
	300	0.996	0.992	0.911	0.589	1.546
	500	1.000	0.994	0.681	0.499	1.364
2.0	100	0.986	0.902	1.291	0.630	2.049
	200	0.996	0.968	0.871	0.549	1.587
	300	0.998	0.988	0.706	0.499	1.415
	500	0.998	0.992	0.549	0.426	1.290

Table S6: Comparison of algorithm with/without sample-splitting in setting 2 with  $p = 500$ .