

# Projection-based Inference for High-dimensional Linear Models

Sangyoon Yi<sup>\*†</sup> and Xianyang Zhang<sup>†</sup>

Texas A&M University

**Abstract:** We develop a new method to estimate the projection direction in the debiased Lasso estimator. The basic idea is to decompose the overall bias into two terms corresponding to strong and weak signals respectively. We propose to estimate the projection direction by balancing the squared biases associated with the strong and weak signals as well as the variance of the projection-based estimator. Standard quadratic programming solver can efficiently solve the resulting optimization problem. In theory, we show that the unknown set of strong signals can be consistently estimated and the projection-based estimator enjoys the asymptotic normality under suitable assumptions. A slight modification of our procedure leads to an estimator with a potentially smaller order of bias comparing to the original debiased Lasso. We further generalize our method to conduct inference for a sparse linear combination of the regression coefficients. Numerical studies demonstrate the advantage of the proposed approach concerning coverage accuracy over some existing alternatives.

**Keywords:** Confidence interval, High-dimensional linear models, Lasso, Quadratic programming.

## 1 Introduction

Uncertainty quantification after model selection has been an active field of research in statistics for the past few years. The problem is challenging as the Lasso type estimator does not admit a tractable asymptotic limit due to its non-continuity at zero. Standard bootstrap and subsampling techniques cannot capture such non-continuity and thus fail for the Lasso estimator even in the low-dimensional regime. Several attempts have been made in the recent literature to tackle this challenge. For example, (Multi) sample-splitting and subsequent statistical inference procedures have been developed in Wasserman and Roeder (2009) and Meinshausen et al. (2009). Meinshausen and Bühlmann (2010) proposed the so-called stability selection method based on subsampling in combination with selection algorithms. Chatterjee and Lahiri (2011, 2013) have considered the bootstrap methods that can provide valid approximation to the limiting distributions of the Lasso and adaptive Lasso estimators, respectively.

---

<sup>\*</sup>Corresponding author. E-mail: syi@stat.tamu.edu.

<sup>†</sup>Department of Statistics, Texas A&M University, College Station, TX 77843, USA. The research was supported in part by NSF grants DMS-1607320 and DMS-1811747.

For statistical inference after model selection, Berk et al. (2013) developed a post-selection inference procedure by reducing the problem to one of simultaneous inference. Lockhart et al. (2014) constructed a statistic from the Lasso solution path and showed that it converges to a standard exponential distribution. To account for the effects of the selection, Lee et al. (2016) developed an exact post-selection inference procedure by characterizing the distribution of a post-selection estimator conditioned on the selection event. By leveraging the same core of statistical framework, Tibshirani et al. (2016) proposed a general scheme to derive post-selection hypothesis tests at any step of forward stepwise and least angle regression, or any step along the Lasso regularization path. Barber and Candès (2015) proposed an inferential procedure by adding knockoff variables to create certain symmetry among the original variables and their knockoff copies. By exploring such symmetry, they showed that the method provides finite sample false discovery rate control. The knockoff procedure has been extended to the high dimensional linear model in Barber and Candès (2019) and the settings in which the conditional distribution of the response is completely unknown in Candès et al. (2018).

Along with a different line that is more closely related to the current work, Zhang and Zhang (2014) first introduced the idea of regularized projection, which has been further explored and extended in van de Geer et al. (2014) and Javanmard and Montanari (2014). The common idea is to find a projection direction designed to remove the bias term in the Lasso estimator. The resulting debiased Lasso estimator which is no longer sparse was shown to admit an asymptotic normal limit. To find the projection direction, the nodewise Lasso regression by Meinshausen and Bühlmann (2006) was adopted in both Zhang and Zhang (2014) and van de Geer et al. (2014), while Javanmard and Montanari (2014) considered a convex optimization problem to approximate the precision matrix of the design. Zhang and Cheng (2017) and Dezeure et al. (2017) proposed bootstrap-assisted procedures to conduct simultaneous inference based on the debiased Lasso estimators. Belloni et al. (2014) developed a two-stage procedure with the so-called post-double-selection as first and least squares estimation as second stage. Ning and Liu (2017) proposed a decorrelated score test in a likelihood based framework. Zhu and Bradic (2018a, 2018b) developed projection-based methods that are robust to the lack of sparsity in the model parameter. More recent advances along this direction include Neykov et al. (2018) and Chang et al. (2019). Focusing on the theoretical aspects of debiased Lasso, Javanmard and Montanari (2018) studied the optimal sample size for debiased Lasso and Cai and Guo (2017) showed that the debiased estimator achieves the minimax rate. Although the methodology and theory for the debiased Lasso estimator are elegant, its empirical performance could be undesirable. For instance, the average coverage rate for active variables could be far lower than the nominal levels in finite sample [see, e.g., van de Geer et al. (2014)].

A natural question to ask is whether there exist alternative projection directions that can improve the finite sample performance in the original debiased Lasso estimator. In this paper, we propose a new method to estimate the projection direction and construct a novel Bias Reducing Projection (BRP) estimator, which is designed to further reduce the bias of the original debiased Lasso estimator. Different from the nodewise Lasso adopted in both Zhang and Zhang (2014)

and van de Geer et al. (2014), we propose a direct approach to estimate the projection direction. Our method is related to the procedure in Javanmard and Montanari (2014) but differs in the following aspects. (i) We formulate a different objective function which appropriately balances the squared bias and the variance of the BRP estimator; (ii) We decompose the bias term into two parts according to a preliminary estimate of the signal strength: one associated with the strong signals and the other one related to the weak signals and noise; (iii) We develop new methods to estimate the set of strong signals and to select the tuning parameters involved in the objective function.

Our approach relies crucially on the following observation in finite sample: the bias term associated with the strong signals contributes more to the overall bias. Motivated by this fact, we estimate the projection direction by minimizing an objective function that assigns different weights to the squared bias terms associated with the strong and weak signals. The set of strong signals is unknown but can be consistently estimated based on a preliminary debiased Lasso estimator. The resulting optimization problem can be cast into a quadratic programming problem which can be efficiently solved using a standard quadratic programming solver. We use residual bootstrap to estimate the coverage probabilities associated with different choices of weights and select the one that delivers the shortest interval width while ensuring that the bootstrap estimate of the coverage probability is close to the nominal level.

In theory, we show that the unknown set of strong signals can be consistently estimated by a surrogate set based on a preliminary projection-based Lasso estimator, where the projection direction is obtained using a novel formulation. The BRP estimator is shown to enjoy the asymptotic normality under suitable assumptions. As one of the main contributions, we prove that a slight modification of our BRP estimator leads to an estimator with a potentially smaller order of bias comparing to the original debiased Lasso. We further generalize our BRP estimator to conduct statistical inference for a sparse linear combination of the regression coefficients under suitable assumptions on a loading vector. We demonstrate the usefulness of the proposed approach by comparing it with the state-of-the-art approaches in simulations.

The rest of the paper is organized as follows. We introduce the projection-based estimator and develop a new formulation to find the projection direction in Section 2. We propose a method to estimate the set of strong signals and show its consistency in Section 3.1. We establish the asymptotic normality of the BRP estimator in Section 3.2 and the modified BRP estimator which could result in a potentially smaller order of bias compared to the original debiased Lasso is proposed in Section 3.3. Section 4 generalizes the method to conduct inference for a sparse linear combination of the regression coefficients. We develop a bootstrap-assisted procedure for choosing the tuning parameters in Section 5. Section 6 presents some numerical results. Section 7 concludes. Technical details and additional numerical results are gathered in Supplementary Material.

Throughout this paper, we use the following notations: For a matrix  $\mathbf{A} \in \mathbb{R}^{d \times d}$  and two sets  $I, J \subseteq [d] := \{1, 2, \dots, d\}$ , denote by  $\mathbf{A}_{I,J}$  ( $\mathbf{A}_{-I,-J}$ ) the submatrix of  $\mathbf{A}$  with (without) the rows in  $I$  and columns in  $J$ . Write  $\mathbf{A}_{[d],-I} = \mathbf{A}_{-I}$ . Similarly for a vector  $a \in \mathbb{R}^q$ , write  $a_I$  ( $a_{-I}$ ) the subvector of  $a$  with (without) the components in  $I$ . Let  $\|a\|_q$  with  $0 \leq q \leq \infty$  be the  $l_q$  norm of  $a$  and write  $\|a\| = \|a\|_2$ . For two sets  $\mathcal{S}_1, \mathcal{S}_2$ , let  $\mathcal{S}_1 \setminus \mathcal{S}_2$  be the set of elements in  $\mathcal{S}_1$  but not in  $\mathcal{S}_2$ .

Denote by  $|\mathcal{S}_1|$  the cardinality of  $\mathcal{S}_1$ . For a square matrix  $\mathbf{A}$ , let  $\lambda_{\max}(\mathbf{A})$  and  $\lambda_{\min}(\mathbf{A})$  be its largest and smallest eigenvalues respectively. Define  $\|\mathbf{A}\| = \|\mathbf{A}\|_{\text{op}} = \sup_{a \in \mathcal{S}^{d-1}} \|\mathbf{A}a\|$  as the operator norm of  $\mathbf{A}$ , where  $\mathcal{S}^{d-1}$  is the unit sphere in  $\mathbb{R}^d$ . The sub-gaussian norm of a random variable  $X$  which we denote by  $\|X\|_{\psi_2}$  is defined as  $\|X\|_{\psi_2} = \sup_{q \geq 1} q^{-1/2} (E|X|^q)^{1/q}$ . For a random vector  $X \in \mathbb{R}^d$ , its sub-gaussian norm can be defined as  $\|X\|_{\psi_2} = \sup_{a \in \mathcal{S}^{d-1}} \|a^\top X\|_{\psi_2}$ . The sub-exponential norm of a random variable  $X$  which we denote by  $\|X\|_{\psi_1}$  is defined as  $\|X\|_{\psi_1} = \sup_{q \geq 1} q^{-1} (E|X|^q)^{1/q}$ . For a random vector  $X \in \mathbb{R}^d$ , its sub-exponential norm can be defined as  $\|X\|_{\psi_1} = \sup_{a \in \mathcal{S}^{d-1}} \|a^\top X\|_{\psi_1}$ . Let  $(\mathcal{M}, \rho)$  be a metric space and let  $\varepsilon > 0$ . A subset  $\mathcal{N}_\varepsilon$  of  $\mathcal{M}$  is called an  $\varepsilon$ -net of  $\mathcal{M}$  if every point  $x \in \mathcal{M}$  can be approximated within  $\varepsilon$  by some point  $y \in \mathcal{N}_\varepsilon$ , i.e.,  $\rho(x, y) \leq \varepsilon$ . The minimal cardinality of an  $\varepsilon$ -net of  $\mathcal{M}$  is called the covering number of  $\mathcal{M}$ .

## 2 Projection-based estimator

To illustrate the idea, we shall focus on the high-dimensional linear model:

$$Y = \mathbf{X}\beta + \epsilon, \quad (1)$$

where  $Y = (y_1, \dots, y_n)^\top \in \mathbb{R}^{n \times 1}$  is the response vector,  $\mathbf{X} = (X_1, \dots, X_p) \in \mathbb{R}^{n \times p}$  is the design matrix,  $\beta = (\beta_1, \dots, \beta_p)^\top \in \mathbb{R}^{p \times 1}$  is the vector of unknown regression coefficients with  $\|\beta\|_0 = s_0$  and  $\epsilon = (\epsilon_1, \dots, \epsilon_n)^\top$  is the vector of independent errors with the common variance  $\sigma^2$ .

### 2.1 Motivation

Suppose we are interested in conducting inference for a single regression coefficient  $\beta_j$  for  $1 \leq j \leq p$ . We first rewrite model (1) as

$$\eta_j := Y - \mathbf{X}_{-j}\beta_{-j} = X_j\beta_j + \epsilon. \quad (2)$$

If the value of  $\eta_j$  is known, the problem would reduce to the inference about  $\beta_j$  in a simple linear regression model. As  $\eta_j$  is not directly observable, a natural idea is to replace  $\eta_j$  by a suitable estimator defined as

$$\hat{\eta}_j = Y - \mathbf{X}_{-j}\hat{\beta}_{-j} = X_j\beta_j + \epsilon + \mathbf{X}_{-j}(\beta_{-j} - \hat{\beta}_{-j}), \quad (3)$$

where  $\hat{\beta}$  is a preliminary estimator for  $\beta$ . Here (3) is an approximation to (2) with the extra term  $\mathbf{X}_{-j}(\beta_{-j} - \hat{\beta}_{-j})$  due to the estimation effect by replacing  $\beta_{-j}$  with  $\hat{\beta}_{-j}$ . In this paper, we focus on the Lasso estimator given by

$$\hat{\beta} = \underset{\tilde{\beta} \in \mathbb{R}^p}{\text{argmin}} \left\{ \frac{1}{2n} \|Y - \mathbf{X}\tilde{\beta}\|^2 + \lambda \|\tilde{\beta}\|_1 \right\}$$

whose properties have now been well understood [see e.g. Bühlmann and van de Geer (2011); Hastie et al. (2015)]. We also try the alternative Lasso formulation without penalizing  $\beta_j$  in our numerical studies and find that it does not improve the finite sample performance. Now given a projection vector  $v_j = (v_{j,1}, \dots, v_{j,n})^\top \in \mathbb{R}^{n \times 1}$  such that  $v_j^\top X_j = n$ , we define the projection-based estimator for  $\beta_j$  as

$$\tilde{\beta}_j(v_j) := \frac{1}{n} v_j^\top \hat{\eta}_j = \beta_j + \frac{1}{n} v_j^\top \epsilon + R(v_j, \beta_{-j}), \quad (4)$$

where  $R(v_j, \beta_{-j}) = n^{-1} v_j^\top \mathbf{X}_{-j} (\beta_{-j} - \hat{\beta}_{-j})$  is the bias term caused by the estimation effect. (4) implies that

$$\sqrt{n}(\tilde{\beta}_j(v_j) - \beta_j) = \frac{1}{\sqrt{n}} v_j^\top \epsilon + \sqrt{n} R(v_j, \beta_{-j}).$$

To ensure that  $\tilde{\beta}_j(v_j)$  has asymptotically tractable limiting distribution, we require the bias term  $\sqrt{n} R(v_j, \beta_{-j})$  to be dominated by the leading term  $n^{-1/2} v_j^\top \epsilon$ , which converges to a normal limit under suitable assumptions. In other words, the bias term  $\sqrt{n} R(v_j, \beta_{-j})$  controls the non-Gaussianity of  $\tilde{\beta}_j(v_j)$ . A practical challenge here is that the bias  $\sqrt{n} R(v_j, \beta_{-j})$  can be hardly estimated directly from the data. It is common in the literature to replace  $|\sqrt{n} R(v_j, \beta_{-j})|$  by a conservative estimator using the  $l_1 - l_\infty$  bound, i.e.,

$$\|\sqrt{n}(\beta_{-j} - \hat{\beta}_{-j})\|_1 \|n^{-1} v_j^\top \mathbf{X}_{-j}\|_\infty. \quad (5)$$

See Zhang and Zhang (2014), van de Geer et al. (2014), Javanmard and Montanari (2014). We note that the variance of  $n^{-1/2} v_j^\top \epsilon$  is equal to  $\sigma^2 n^{-1} \|v_j\|^2$ . To achieve efficiency, we shall also try to minimize  $\sigma^2 n^{-1} \|v_j\|^2$  given that the bias  $\sqrt{n} R(v_j, \beta_{-j})$  is properly controlled. Because the first term in (5) is independent of  $v_j$ , we can seek a projection direction to minimize a linear combination of  $\|n^{-1} v_j^\top \mathbf{X}_{-j}\|_\infty$  and the variance  $\sigma^2 n^{-1} \|v_j\|^2$ . However, the  $l_1 - l_\infty$  bound on the whole bias term could be conservative as it does not take into account the specific form of the bias term. We note that the bias term can be written as

$$\begin{aligned} \sqrt{n} R(v_j, \beta_{-j}) &= \frac{1}{\sqrt{n}} \sum_{k \neq j} v_j^\top X_k (\beta_k - \hat{\beta}_k) \\ &= \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{S}_j^{(1)}(\nu)} v_j^\top X_k (\beta_k - \hat{\beta}_k) + \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{S}_j^{(2)}(\nu)} v_j^\top X_k (\beta_k - \hat{\beta}_k) \\ &= \sqrt{n} R_{(1)}(v_j, \beta_{-j}) + \sqrt{n} R_{(2)}(v_j, \beta_{-j}), \end{aligned} \quad (6)$$

where  $\mathcal{S}_j^{(1)}(\nu) := \mathcal{S}(\nu) \setminus \{j\}$  and  $\mathcal{S}_j^{(2)}(\nu) := \mathcal{S}(\nu) \setminus \{j\}$  denote the index sets (except  $j$ ) associated with the strong and weak signals respectively for  $\mathcal{S}(\nu) := \{k : |\beta_k| \geq \nu\}$  and both  $R_{(1)}(v_j, \beta_{-j})$  and  $R_{(2)}(v_j, \beta_{-j})$  are defined accordingly. Here  $\nu$  is a threshold that separates the coefficients into two-groups namely the group with strong signals and the group with weak or zero signal. For

example, one can set  $\nu = c_0 \sqrt{\log(p)/n}$  for some large enough constant  $c_0$ , which is the minimax rate for support recovery.

The formulation (6) using the decomposition associated with signal strengths can be empirically motivated. Specifically, it generally provides a smaller bias than the one without such decomposition with the simulated data. Figure 4 illustrates one such representative case where we make a comparison of the biases for projection vectors calculated based on two different methods: the one solves (8) by using the estimated set of strong signals as in Section 3.1 (denoted by “With Decomposition”) and the other one solves the same problem but with  $\mathcal{A}_j^{(1)} = \emptyset$  (denoted by “Without Decomposition”). It can be seen that “With Decomposition” shows a smaller bias than “Without Decomposition.” Similar results could be observed in various simulation settings.

## 2.2 A new projection direction

In this subsection, we propose a novel formulation to find the projection direction. When  $|\mathcal{S}_j^{(1)}(\nu)| \leq n$ , we have the freedom to choose  $v_j$  to make the term  $\|n^{-1}v_j^\top \mathbf{X}_{\mathcal{S}_j^{(1)}(\nu)}\|_\infty$  arbitrarily small. In fact, we can always choose  $v_j$  such that it is orthogonal to all  $X_k$  with  $k \in \mathcal{S}_j^{(1)}(\nu)$ . The basic idea here is to find a projection direction  $v_j$  such that it is “more orthogonal” to the space spanned by  $\{X_k\}_{k \in \mathcal{S}_j^{(1)}(\nu)}$  as compared to the space spanned by  $\{X_k\}_{k \in \mathcal{S}_j^{(2)}(\nu)}$ . With this intuition in our mind and the goal to balance the squared bias with the variance, we formulate the following optimization problem

$$\begin{aligned} \min_{v_j} & \left( \gamma_1 \max_{k \in \mathcal{S}_j^{(1)}(\nu)} |n^{-1}v_j^\top X_k|^2 + \gamma_2 \max_{k \in \mathcal{S}_j^{(2)}(\nu)} |n^{-1}v_j^\top X_k|^2 + \sigma^2 n^{-1} \|v_j\|^2 \right), \\ \text{s.t.} & \quad v_j^\top X_j = n, \end{aligned} \tag{7}$$

where  $\gamma_1, \gamma_2 > 0$  are tuning parameters which control the trade-off between the squared bias and the variance. The term  $\gamma_1 \max_{k \in \mathcal{S}_j^{(1)}(\nu)} |n^{-1}v_j^\top X_k|^2$  ( $\gamma_2 \max_{k \in \mathcal{S}_j^{(2)}(\nu)} |n^{-1}v_j^\top X_k|^2$ ) corresponds to the  $l_1 - l_\infty$  bound for  $R_{(1)}^2$  ( $R_{(2)}^2$ ). By introducing two ancillary variables  $u_{j1}, u_{j2}$ , (7) can be cast into the following quadratic programming problem

$$\begin{aligned} \min_{u_{j1}, u_{j2}, v_j} & \quad (\gamma_1 u_{j1}^2 + \gamma_2 u_{j2}^2 + \sigma^2 n^{-1} \|v_j\|^2), \\ \text{s.t.} & \quad v_j^\top X_j = n, \\ & \quad -u_{j1} \leq n^{-1}v_j^\top X_k \leq u_{j1}, \quad k \in \mathcal{S}_j^{(1)}(\nu), \\ & \quad -u_{j2} \leq n^{-1}v_j^\top X_k \leq u_{j2}, \quad k \in \mathcal{S}_j^{(2)}(\nu), \end{aligned}$$

which can be solved efficiently using existing quadratic programming solver.

The set  $\mathcal{S}_j^{(1)}(\nu)$  is generally unknown and needs to be replaced by a surrogate set  $\mathcal{A}_j^{(1)}$  with  $|\mathcal{A}_j^{(1)}| \leq n$ . In Section 3.1, we describe a method to select  $\mathcal{A}_j^{(1)}$  based on a preliminary projection-

based estimators. We show that  $\mathcal{A}_j^{(1)}$  converges asymptotically to a nonrandom limit, i.e.,

$$P\left(\mathcal{A}_j^{(1)} = \mathcal{B}_j^{(1)}\right) \rightarrow 1,$$

for a nonrandom subset  $\mathcal{B}_j^{(1)}$  of  $[p]$ . We remark that  $\mathcal{B}_j^{(1)}$  does not need to agree with  $\mathcal{S}_j^{(1)}(\nu)$  for our procedure to be valid. To ensure that the remainder term is negligible, the theoretical analysis in Section 3.2 suggests that  $\gamma_1$  and  $\gamma_2$  should both be of the order  $O(\sigma^2 n / \log p)$ . Combining the above discussions, we now state the optimization problem for obtaining the optimal projection direction

$$\begin{aligned} \min_{u_{j1}, u_{j2}, v_j} & \left( C_1 \frac{n}{\log p} u_{j1}^2 + C_2 \frac{n}{\log p} u_{j2}^2 + n^{-1} \|v_j\|^2 \right), \\ \text{s.t.} & \quad v_j^\top X_j = n, \\ & \quad -u_{j1} \leq n^{-1} v_j^\top X_k \leq u_{j1}, \quad k \in \mathcal{A}_j^{(1)}, \\ & \quad -u_{j2} \leq n^{-1} v_j^\top X_k \leq u_{j2}, \quad k \in \mathcal{A}_j^{(2)}, \end{aligned} \tag{8}$$

where  $\mathcal{A}_j^{(2)} := \left(\mathcal{A}_j^{(1)}\right)^c \setminus \{j\}$  and  $C_1, C_2 > 0$  are tuning parameters whose choice will be discussed in Section 5.

**Remark 1.** A related method is the refitted Lasso by Liu and Yu (2013). The idea is to refit the model selected by the Lasso and conduct inference based on the refitted least squares estimator. Such an estimator fits into the framework of the projection-based estimators. To see this, let  $\hat{S}$  be the set of active variables selected by the Lasso and note that  $\hat{\beta}_k = 0$  for  $k \notin \hat{S}$ . For each  $j \in \hat{S}$ , let  $\hat{w}_j$  be the projection of  $X_j$  onto the orthogonal space of  $\mathbf{X}_{\hat{S} \setminus \{j\}}$ . Then the refitted least squares estimator is given by  $\hat{w}_j^\top (Y - X_{-j} \hat{\beta}_{-j}) / (\hat{w}_j^\top X_j)$ . It is easy to see that the bias for the refitted least squares estimator is proportional to  $\sum_{k \notin \hat{S}} \hat{w}_j^\top X_k \beta_k$ , which disappears when the selected model contains all significant variables. However, when the model selection consistency fails, such a procedure is no longer valid due to the nonnegligible bias.

## 3 Methodology

### 3.1 Surrogate set

We describe a procedure to estimate the set of strong signals based on a preliminary projection-based estimator. It should be noted that the estimator here is different from the original debiased Lasso because it is based on the novel formulation (8). Specifically, for some  $\tau > 0$ , we define our estimate for the set of strong signals as

$$\mathcal{A}(\tau) := \{l : |T_l| > \sqrt{\tau \log p}\} \quad \text{where} \quad T_l = \frac{\sqrt{n} \tilde{\beta}_l(\hat{v}_l)}{\hat{\sigma} n^{-1/2} \|\hat{v}_l\|} \tag{9}$$

where  $\hat{\sigma}$  is an estimator of the noise level  $\sigma$  and  $\tilde{\beta}_l(\hat{v}_l)$  is a projection-based estimator with  $\hat{v}_l$  being the solution to the following optimization problem

$$\begin{aligned} \min_{u_l, v_l} & \left( C_0 \frac{n}{\log p} u_l^2 + n^{-1} \|v_l\|^2 \right), \\ \text{s.t. } & v_l^\top X_l = n, \\ & -u_l \leq n^{-1} v_l^\top X_k \leq u_l, \quad k \neq l. \end{aligned} \tag{10}$$

In practice, both  $C_0$  and  $\tau$  need to be appropriately chosen. The details for the selection are discussed in Section 8.1. Note that (10) is a special case of (8) when we have no knowledge about the set of strong signals, that is,  $\mathcal{A}_l^{(1)} = \emptyset$ . We define the surrogate sets to be

$$\mathcal{A}_j^{(1)}(\tau) := \mathcal{A}(\tau) \setminus \{j\}, \quad \mathcal{A}_j^{(2)}(\tau) := \mathcal{A}(\tau) \setminus \{j\}. \tag{11}$$

Throughout the paper, we consider the variance estimator

$$\hat{\sigma}^2 = \frac{1}{n} \|Y - \mathbf{X}\hat{\beta}\|^2 \tag{12}$$

which appears to outperform an alternative estimator  $\|Y - \mathbf{X}\hat{\beta}\|^2 / (n - \|\hat{\beta}\|_0)$  studied in Reid et al. (2016), see Figure 22 in Supplementary Material for a comparison. Before presenting the main result of this subsection, we introduce some assumptions.

**Assumption 1.** *There exist a set  $\mathcal{B} \subseteq [p] = \{1, 2, \dots, p\}$  and  $0 \leq d_0 < d_1$  such that*

$$\begin{aligned} \max_{l \in \mathcal{B}^c} \frac{|\sqrt{n}\beta_l|}{\sigma} & \leq \sqrt{d_0 \log p}, \\ \min_{l \in \mathcal{B}} \frac{|\sqrt{n}\beta_l|}{\sigma} & \geq \sqrt{d_1 \log p}. \end{aligned}$$

**Assumption 2.** *The error  $\epsilon$  is a mean-zero sub-gaussian random vector with the sub-gaussian norm  $\kappa_\epsilon$ .*

**Assumption 3.** *The preliminary estimator satisfies that*

$$\sqrt{n} \|\hat{\beta} - \beta\|_1 = O_p(s_0 \sqrt{\log(p)}).$$

**Assumption 4.** *The variance estimator  $\hat{\sigma}^2$  is consistent in the sense that  $\hat{\sigma}/\sigma \xrightarrow{p} 1$ .*

**Assumption 5.** *Suppose the design matrix  $\mathbf{X} \in \mathbb{R}^{n \times p}$  has i.i.d. rows with zero population mean and covariance matrix  $\Sigma = (\Sigma_{i,j})_{i,j=1}^p$ . Assume that*

1.  $\max_j \Sigma_{j,j} < \infty$ ;
2.  $\lambda_{\min}(\Sigma) \geq \Lambda_{\min} > 0$ ;
3. *The rows of  $\mathbf{X}$  are sub-gaussian with the sub-gaussian norm  $\kappa < \infty$ .*



**Assumption 6.**  $n, p$  and  $s_0$  satisfy the rate condition  $s_0 \log p / \sqrt{n} = o(1)$ .

Assumption 1 allows the strengths of strong and weak signals to be the same order and thus is much weaker than the ‘‘beta-min’’ condition which requires the weak signals to be of smaller order. Assumptions 3 and 4 are satisfied for the Lasso estimator and the variance estimator  $\hat{\sigma}$  in (12) under suitable regularity conditions [Bühlmann and van de Geer (2011)]. Assumptions 2 and 5 require the error and design to be sub-gaussian. Similar assumptions have been made in van de Geer et al. (2014). Like Javanmard and Montanari (2014), the validity of our method does not rely on the sparsity of the precision matrix of the design, which is required in the nodewise Lasso regression for the original debiased Lasso. In view of Cai and Guo (2017), the rate condition in Assumption 6 cannot be relaxed without extra information. Zhu and Bradic (2018a, 2018b) proposed testing procedures in high-dimensional linear models which impose much weaker restrictions on model sparsity or the loading vector representing the hypothesis. However, their methods require certain auxiliary sparse models, which are not needed for our procedure.

Define  $\Sigma_{j \setminus -j} = \Sigma_{j,j} - \Sigma_{j,-j} \Sigma_{-j,-j}^{-1} \Sigma_{-j,j}$  and  $\kappa_{0j} = 2 \left( 1 + \sqrt{\Lambda_{\min}^{-1} \Sigma_{j,j}} \right) \kappa^2$  for  $1 \leq j \leq p$ . The following proposition shows that the surrogate set  $\mathcal{A}_j^{(1)}(\tau)$  with a properly chosen  $\tau$  converges to  $\mathcal{B} \setminus \{j\}$ .

**Proposition 1.** Define  $\mathcal{A}_j^{(1)}(\tau)$  and  $\mathcal{A}_j^{(2)}(\tau)$  as in (11) and let  $\hat{v}_l$  be the solution to (10) for  $l \neq j$ . Suppose  $d_0, d_1$  and  $\tau$  satisfy

$$\frac{\sigma^2}{32e\kappa_\epsilon^2} (\sqrt{\tau} - \sqrt{d_0 \max_l \Sigma_{l,l}})^2 > 1$$

and  $\sqrt{d_1/M} - \sqrt{\tau} > 0$  where

$$M = \left( \min_{1 \leq l \leq p} \Sigma_{l \setminus -l} \right)^2 \left( 2C_0 \left( \min_{1 \leq l \leq p} \frac{1}{8e^2 (\kappa_{0l})^2} \right)^{-1} + \max_{1 \leq l \leq p} \Sigma_{l \setminus -l} \right).$$

Then under Assumptions 1-6, we have

$$\begin{aligned} \mathbb{P} \left( \max_{l \in \mathcal{B}_j^{(2)}} |T_l| \leq \sqrt{\tau \log p} \right) &\rightarrow 1, \\ \mathbb{P} \left( \min_{l \in \mathcal{B}_j^{(1)}} |T_l| > \sqrt{\tau \log p} \right) &\rightarrow 1, \end{aligned}$$

where  $\mathcal{B}_j^{(1)} := \mathcal{B} \setminus \{j\}$  and  $\mathcal{B}_j^{(2)} := \left( \mathcal{B}_j^{(1)} \right)^c \setminus \{j\}$ . As a consequence,  $\mathbb{P} \left( \mathcal{A}_j^{(1)}(\tau) = \mathcal{B}_j^{(1)} \right) \rightarrow 1$ .

**Remark 2.** As shown in Proposition 1, the surrogate set in (11) has an asymptotic (nonrandom) limit, which implies that the projection direction obtained in (8) is asymptotically independent of the random error  $\epsilon$ . This fact is useful in the proof of Theorem 1 later. To ensure the independence between the projection direction and the random error, we can also employ the sample splitting strategy, i.e., we split the samples into two subsamples, estimate the set of strong signals based on

the first subsample and construct the projection-based estimator based on another subsample. As we use all samples in building the projection-based estimator, our method is more efficient than the sample splitting strategy.

**Remark 3.** When  $d_0 = 0$ ,  $\mathcal{B}$  coincides with the support of  $\beta$ . Proposition 1 suggests that one can consistently recover the support of  $\beta$  by thresholding the projection-based estimator.

### 3.2 Bias reducing projection (BRP) estimator

In this subsection, we introduce the bias reducing projection (BRP) estimator and study its asymptotic behavior. Let  $\tilde{v}_j$  be the solution to (8) based the surrogate sets in (11). Then the BRP estimator  $\tilde{\beta}_j(\tilde{v}_j)$  is defined as

$$\tilde{\beta}_j(\tilde{v}_j) = \frac{1}{n} \tilde{v}_j^\top \hat{\eta}_j = \frac{1}{n} \tilde{v}_j^\top (Y - \mathbf{X}_{-j} \hat{\beta}_{-j}).$$

In the following, we introduce the two asymptotic results depending on whether the surrogate set is estimated from the same data set used to find the projection direction. We first state the following theorem on the asymptotic normality when the surrogate set is estimated via (11).

**Theorem 1.** Denote by  $\tilde{v}_j$  the solution to (8) with  $\mathcal{A}_j^{(1)}(\tau)$  and  $\mathcal{A}_j^{(2)}(\tau)$  in (11). Suppose the assumptions in Proposition 1 hold and further assume that for some  $\delta > 0$ ,

$$\|\tilde{v}_j\|_{2+\delta} = o_{a.s.}(\|\tilde{v}_j\|). \quad (13)$$

Then we have

$$\frac{\sqrt{n} \left( \tilde{\beta}_j(\tilde{v}_j) - \beta_j \right)}{\hat{\sigma} n^{-1/2} \|\tilde{v}_j\|} \xrightarrow{d} N(0, 1). \quad (14)$$

Thus an asymptotic  $100(1 - \alpha)\%$  confidence interval for  $\beta_j$  is given by

$$\text{CI}(1 - \alpha) = \left\{ b \in \mathbb{R} : \left| \frac{\sqrt{n}(\tilde{\beta}_j(\tilde{v}_j) - b)}{\hat{\sigma} n^{-1/2} \|\tilde{v}_j\|} \right| \leq z_{1-\alpha/2} \right\}, \quad (15)$$

where  $z_{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of  $N(0, 1)$ .

(13) is a Lyapunov type condition which implies the central limit theorem. This type of assumption regarding the projection direction has also been imposed in Dezeure et al. (2017). It can be dropped under the Gaussian assumption on the errors. If the surrogate set is chosen based on prior knowledge or estimated from an independent data set (e.g., based on sample splitting), then Assumptions 1-2 can be relaxed and we have the following result.

**Corollary 1.** Suppose the surrogate set  $\mathcal{A}_j^{(1)}$  is independent of the data. Under Assumptions 3-6 and further assuming that for some  $\delta > 0$ ,  $E[|\epsilon_i|^{2+\delta}] < \infty$  and  $\|\tilde{v}_j\|_{2+\delta} = o_{a.s.}(\|\tilde{v}_j\|)$ , then (14) still holds.

### 3.3 Modified bias reducing projection (MBRP) estimator

We introduce a modified bias reducing projection (MBRP) estimator which is motivated by Proposition 1 and the refitted Lasso idea. This new estimator would lead to a potentially smaller order of bias compared to that of the original debiased Lasso estimator under suitable assumptions as shown in Proposition 2. Thus, it is expected to provide better empirical coverage probability. See more details in Section 6. To motivate the MBRP estimator, we note that the bias associated with the BRP estimator based on some estimator  $\check{\beta}$  for  $\beta$  can be written as

$$\begin{aligned}\sqrt{n}R(v_j, \beta_{-j}) &= \frac{1}{\sqrt{n}} \sum_{k \neq j} v_j^\top X_k (\beta_k - \check{\beta}_k) \\ &= \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{B}_j^{(1)}} v_j^\top X_k (\beta_k - \check{\beta}_k) + \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{B}_j^{(2)}} v_j^\top X_k (\beta_k - \check{\beta}_k)\end{aligned}$$

where  $\mathcal{B}_j^{(1)}, \mathcal{B}_j^{(2)}$  are the same as in Proposition 1. When  $|\mathcal{B}_j^{(1)}| \leq n$ , we can always require  $v_j$  to be exactly orthogonal to  $\mathbf{X}_{\mathcal{B}_j^{(1)}}$ . So, the bias associated with the set of strong signals becomes zero. Thus it suffices to control the bias term associated with  $\mathcal{B}_j^{(2)}$  by properly choosing  $v_j$  and  $\check{\beta}$ , which will be clarified below.

To find the projection direction for the MBRP estimator, we consider the optimization problem

$$\begin{aligned}\min_{u_{j2}, v_j} & \left( C_2 \frac{n}{\log p} u_{j2}^2 + n^{-1} \|v_j\|^2 \right), \\ \text{s.t.} & \quad v_j^\top X_j = n, \\ & \quad n^{-1} v_j^\top X_k = 0, \quad k \in \mathcal{A}_j^{(1)}, \\ & \quad -u_{j2} \leq n^{-1} v_j^\top X_k \leq u_{j2}, \quad k \in \mathcal{A}_j^{(2)}.\end{aligned}\tag{16}$$

Different from (8), we require the projection direction to be orthogonal to the column space of  $\mathbf{X}_{\mathcal{A}_j^{(1)}}$  in (16). Instead of using the Lasso estimator  $\hat{\beta}$ , we shall adopt the refitted least squares estimator  $\check{\beta}$  as our preliminary estimator, i.e.,

$$\check{\beta}_{\mathcal{A}_j^{(1)}} = \underset{\check{\beta}}{\operatorname{argmin}} \frac{1}{2n} \|Y - \mathbf{X}_{\mathcal{A}_j^{(1)}} \check{\beta}\|^2, \quad \check{\beta}_{\mathcal{A}_j^{(2)}} = 0.\tag{17}$$

The MBRP estimator is then defined as

$$\tilde{\beta}_j(\bar{v}_j) = \frac{1}{n} \bar{v}_j^\top (Y - \mathbf{X}_{-j} \check{\beta}_{-j}) = \beta_j + \frac{1}{n} \bar{v}_j^\top \epsilon + R(\bar{v}_j, \beta_{-j})\tag{18}$$

where  $R(\bar{v}_j, \beta_{-j}) = n^{-1} \bar{v}_j^\top \mathbf{X}_{-j} (\beta_{-j} - \check{\beta}_{-j})$  and  $\bar{v}_j$  is the solution to problem (16). The MBRP estimator can be viewed as an intermediate estimator between the refitted Lasso and the BRP estimator based on (8). While (16) is a variant of (8) seeking for a projection direction that is exactly orthogonal to the column space of  $\mathbf{X}_{\mathcal{A}_j^{(1)}}$ , the modified procedure uses the refitted estimator

for  $\beta$  as the refitted Lasso does as noted in Remark 1.

We argue that the bias term  $\sqrt{n}R(\bar{v}_j, \beta_{-j})$  which controls non-Gaussianity could have a potentially smaller order compared to that of the original debiased Lasso estimator in the following.

**Proposition 2.** Denote by  $\bar{v}_j$  the solution to (16) with  $\mathcal{A}_j^{(1)}(\tau)$  and  $\mathcal{A}_j^{(2)}(\tau)$  defined in (11). Let  $\check{\beta}$  be the refitted least square estimator in (17). Conditional on the event  $\{\mathcal{A}_j^{(2)} = \mathcal{B}_j^{(2)}\}$ , we have

$$|\sqrt{n}R(\bar{v}_j, \beta_{-j})| \leq O_p \left( \sqrt{d_0} \|\beta_{\mathcal{B}_j^{(2)}}\|_0 \frac{\log p}{\sqrt{n}} \right) \quad (19)$$

under Assumptions 1 and 5. If we further assume that

$$\sqrt{d_0} \|\beta_{\mathcal{B}_j^{(2)}}\|_0 = o(s_0), \quad (20)$$

the bias  $\sqrt{n}R(\bar{v}_j, \beta_{-j})$  is asymptotically negligible with smaller order than that of the original debiased Lasso given by  $O_p(s_0 \log p / \sqrt{n})$ .

In particular, (20) holds if  $d_0 = o(1)$  and  $d_1 = O(1)$ , i.e., the strength of weak signals is of smaller order compared to the strong signals. It is more stringent than Assumption 1 where the magnitudes of the set of strong signals and weak signals are allowed to be of the same order. However, it should be mentioned that Proposition 2 is not necessary for the asymptotic normality in Corollary 2 to be achieved. The following result shows the asymptotic normality of (18) which can be proved by using similar arguments as those for Theorem 1.

**Corollary 2.** Under the assumptions in Theorem 1, we have

$$\frac{\sqrt{n} \left( \tilde{\beta}_j(\bar{v}_j) - \beta_j \right)}{\hat{\sigma} n^{-1/2} \|\bar{v}_j\|} \xrightarrow{d} N(0, 1),$$

where  $\tilde{\beta}_j(\bar{v}_j)$  is defined in (18) and  $\bar{v}_j$  is the solution to (16).

## 4 Inference on a sparse linear combination of parameters

In some applications, one may be interested in conducting inference on  $a^\top \beta$  for a (sparse) loading vector  $a = (a_1, \dots, a_p)^\top \in \mathbb{R}^p$  with  $\|a\|_0 = s \ll n$ . Denote by  $S = S(a) = \{1 \leq j \leq p : a_j \neq 0\}$  the support set of  $a$ . Our method can be generalized to construct estimator and conduct inference for  $a^\top \beta = a_S^\top \beta_S$ . Recall that  $\hat{\beta}$  is the preliminary estimator of  $\beta$ . Define

$$\eta_S = Y - \mathbf{X}_{-S} \beta_{-S} = \mathbf{X}_S \beta_S + \epsilon$$

and

$$\hat{\eta}_S = Y - \mathbf{X}_{-S} \hat{\beta}_{-S} = \mathbf{X}_S \beta_S + \epsilon + \mathbf{X}_{-S} (\beta_{-S} - \hat{\beta}_{-S}).$$

We construct an estimator for  $a^\top \beta$  in the form of  $n^{-1}v_a^\top \hat{\eta}_S$ , where  $v_a = (v_{a,1}, \dots, v_{a,n})^\top$  is a projection direction such that  $n^{-1}v_a^\top \hat{\eta}_S$  has tractable asymptotic limit. Notice that

$$\begin{aligned} n^{-1}v_a^\top \hat{\eta}_S &= n^{-1}v_a^\top \mathbf{X}_S \beta_S + n^{-1}v_a^\top \epsilon + n^{-1}v_a^\top \mathbf{X}_{-S}(\beta_{-S} - \hat{\beta}_{-S}) \\ &= a_S^\top \beta_S + (n^{-1}v_a^\top \mathbf{X}_S - a_S^\top) \beta_S + n^{-1}v_a^\top \epsilon + n^{-1}v_a^\top \mathbf{X}_{-S}(\beta_{-S} - \hat{\beta}_{-S}). \end{aligned}$$

Under the equality constraint that  $n^{-1}v_a^\top \mathbf{X}_S - a_S^\top = 0$  and by rearranging the above terms, we have

$$\sqrt{n}(n^{-1}v_a^\top \hat{\eta}_S - a_S^\top \beta_S) = n^{-1/2}v_a^\top \epsilon + \sqrt{n}R(v_a, \beta_{-S}), \quad (21)$$

where  $R(v_a, \beta_{-S}) = n^{-1}v_a^\top \mathbf{X}_{-S}(\beta_{-S} - \hat{\beta}_{-S})$ . Similar to (6), the bias term can be decomposed into two parts corresponding to different strengths of the signals. Let  $\mathcal{A}_S^{(1)}$  be the surrogate set for the set of strong signals (excluding the elements in  $S$ ), which can be obtained in a similar way as described in Section 3.1. Following the derivations in Section 2, we can formulate the following optimization problem to find  $v_a$

$$\begin{aligned} \min_{u_{a1}, u_{a2}, v_a} & \left( C_1 \frac{n}{\log p} u_{a1}^2 + C_2 \frac{n}{\log p} u_{a2}^2 + n^{-1} \|v_a\|^2 \right), \\ \text{s.t. } & v_a^\top \mathbf{X}_S = n a_S^\top, \\ & -u_{a1} \leq n^{-1}v_a^\top \mathbf{X}_k \leq u_{a1}, \quad k \in \mathcal{A}_S^{(1)}, \\ & -u_{a2} \leq n^{-1}v_a^\top \mathbf{X}_k \leq u_{a2}, \quad k \in \mathcal{A}_S^{(2)}, \end{aligned} \quad (22)$$

where  $\mathcal{A}_S^{(2)} := (\mathcal{A}_S^{(1)} \cup S)^\complement$ . Denote by  $(\tilde{u}_{a1}, \tilde{u}_{a2}, \tilde{v}_a)$  the solution to (22). Our estimator for  $a^\top \beta$  is thus given by  $n^{-1}\tilde{v}_a^\top \hat{\eta}_S$  whose asymptotic normality is established in the following theorem.

**Theorem 2.** With  $\|a\|_0 = s \ll n$ , suppose the assumptions in Proposition 1 hold and  $\|\tilde{v}_a\|_{2+\delta} = o_{a.s.}(\|\tilde{v}_a\|)$  for some  $\delta > 0$ . Then, we have

$$\frac{\sqrt{n}(n^{-1}\tilde{v}_a^\top \hat{\eta}_S - a^\top \beta)}{\hat{\sigma} n^{-1/2} \|\tilde{v}_a\|} \xrightarrow{d} N(0, 1). \quad (23)$$

Thus an asymptotic  $100(1 - \alpha)\%$  confidence interval for  $a^\top \beta$  is given by

$$\text{CI}(1 - \alpha) = \left\{ b \in \mathbb{R} : \left| \frac{\sqrt{n}(n^{-1}\tilde{v}_a^\top \hat{\eta}_S - b)}{\hat{\sigma} n^{-1/2} \|\tilde{v}_a\|} \right| \leq z_{1-\alpha/2} \right\},$$

where  $z_{1-\alpha/2}$  is the  $1 - \alpha/2$  quantile of  $N(0, 1)$ .

We mention some existing works for inference on linear combinations of  $\beta$ . When the sparsity level  $s_0$  is known, Cai and Guo (2017) obtained the minimax expected length of confidence intervals for  $a^\top \beta$  in both the sparse and dense loading regions. They further showed that without the knowledge of  $s_0$ , rate-optimal adaptation in the sparse loading regime is only possible under Assumption

6 and in the dense loading regime, adaptation to  $s_0$  is impossible. In Zhu and Bradic (2018b), the authors proposed a test for linear hypothesis, which does not impose restriction on model sparsity or the loading vector representing the hypothesis. Nevertheless, compared to our method, the method by Zhu and Bradic (2018b) requires an additional sparse model to account for the dependence between the so-called synthesized feature and the stabilized feature.

Parallel to Corollary 1, if the surrogate set is estimated based on prior information or an independent data set, Assumptions 1-2 can be dropped and the asymptotic normality can be established as follows.

**Corollary 3.** Suppose the surrogate set  $\mathcal{A}_j^{(1)}$  is independent of the data. Under Assumptions 3-6 and further assuming that for some  $\delta > 0$ ,  $E[|\epsilon_i|^{2+\delta}] < \infty$  and  $\|\tilde{v}_a\|_{2+\delta} = o_{a.s.}(\|\tilde{v}_a\|)$ , then (23) still holds.

## 5 Selecting the tuning parameters

Bootstrap for debiased Lasso has been recently studied in both Zhang and Cheng (2017) and Dezeure et al. (2017) to approximate the sampling distribution of the debiased Lasso estimator. Here we propose a bootstrap-assisted approach for choosing the tuning parameters in (8), (10) and (16). Specifically, the residual bootstrap is used to obtain the empirical coverage rate and its standard error for selecting the optimal tuning parameters. We focus our discussions on (8) and remark that the procedure is applicable to (10) and (16) as well. Let

$$\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top = Y - \mathbf{X}\hat{\beta}$$

and  $\bar{\varepsilon}_i = \varepsilon_i - n^{-1} \sum_{j=1}^n \varepsilon_j$  be the centered residual where  $\hat{\beta}$  denotes the cross-validated Lasso estimator. Given a sequence of tuning parameters  $\{(c_{1,j,(k)}, c_{2,j,(k)})\}_{k=1}^K$ , we first calculate  $\tilde{v}_j(c_{1,j,(k)}, c_{2,j,(k)})$  which is the solution to (8) given  $(c_{1,j,(k)}, c_{2,j,(k)})$ . Note that the projection direction  $\tilde{v}_j$  only needs to be calculated once for each pair of tuning parameters. Given  $\{\tilde{v}_j(c_{1,j,(k)}, c_{2,j,(k)})\}_{k=1}^K$ , we do the following.

1. To generate the  $b$ -th bootstrap sample, we sample  $n$  residuals with replacement from  $\{\bar{\varepsilon}_i\}_{i=1}^n$  and denote the corresponding samples by  $\varepsilon_b^* = (\varepsilon_{b,1}^*, \dots, \varepsilon_{b,n}^*)^\top$ . Then, generate  $Y_b^*$  such that  $Y_b^* = \mathbf{X}\hat{\beta} + \varepsilon_b^*$ .
2. With  $(X, Y_b^*)$ , calculate the cross-validated Lasso estimator  $\hat{\beta}_b^*$  as well as the projection-based estimator

$$\tilde{\beta}_j(\tilde{v}_j(c_{1,j,(k)}, c_{2,j,(k)})) = \frac{\tilde{v}_j(c_{1,j,(k)}, c_{2,j,(k)})^\top (Y_b^* - \mathbf{X}_{-j}\hat{\beta}_{b,-j}^*)}{n},$$

where  $\hat{\beta}_{b,-j}^*$  denotes  $\hat{\beta}_b^*$  without the  $j$ -th component. We then calculate the  $100(1 - \alpha)\%$  confidence interval  $\text{CI}_{b,j,(k)}^*$  by using (15). For each  $j$ , calculate  $I(\hat{\beta}_j \in \text{CI}_{b,j,(k)}^*)$  which is 1 if  $\hat{\beta}_j$  is covered by  $\text{CI}_{b,j,(k)}^*$  and 0 otherwise. Also, calculate the length of  $\text{CI}_{b,j,(k)}^*$  and denote it as  $\text{Len}_{b,j,(k)}^*$ .

3. Repeat the above steps for  $B$  bootstrap samples. We choose the tuning parameters for  $\beta_j$  as

$$\begin{aligned} (c_{1,j,(k)}^*, c_{2,j,(k)}^*) &= \underset{k}{\operatorname{argmin}} \operatorname{AvgLen}_{j,(k)} \\ \text{s.t. } \check{\operatorname{Cover}}_{j,(k)} + \operatorname{SE}(\check{\operatorname{Cover}}_{j,(k)}) &\geq 1 - \alpha. \end{aligned}$$

where  $\operatorname{AvgLen}_{j,(k)} = B^{-1} \sum_{b=1}^B \operatorname{Len}_{b,j,(k)}^*$  and

$$\begin{aligned} \check{\operatorname{Cover}}_{j,(k)} &= \frac{\sum_{b=1}^B I(\hat{\beta}_j \in \operatorname{CI}_{b,j,(k)}^*)}{B}, \\ \operatorname{SE}(\check{\operatorname{Cover}}_{j,(k)}) &= \sqrt{\frac{\check{\operatorname{Cover}}_{j,(k)}(1 - \check{\operatorname{Cover}}_{j,(k)})}{B}}. \end{aligned}$$

In words, the optimal pair of tuning parameters is selected with the minimum average interval length among all the pairs whose empirical coverage rate increased by one standard error is at least the nominal level  $1 - \alpha$ .

## 6 Numerical results

### 6.1 Confidence interval for a single regression coefficient

We conduct simulations to evaluate the finite sample performance of the proposed BRP and MBRP estimators. We use the R package `quadprog` to solve the quadratic programming problems involved in our methods and the R package `doMC` with 5 cores for parallel computation. All the other implementation details are the same as described in Section 8.1. For comparison, we implement the debiased Lasso in van de Geer et al. (2014) (denoted by DB) using the R package `hdi` and the method in Javanmard and Montanari (2014) (denoted by JM) using the code posted on the authors' website. As we encounter some numerical issue when implementing JM's code for the equicorrelation covariance structure of  $\mathbf{X}$  in (ii). Therefore, we only report the results of JM for the toeplitz covariance structure of  $\mathbf{X}$ . In addition, we present the results of the double selection approach in Belloni et al. (2014) (denoted by BCH) using the R package `hdm`. Due to the high computational cost of BCH in the case of equicorrelation covariance, we only report the result for the active set. We also implement the method in Zhu and Bradic (2018b) (denoted by "ZB" and "ZB2"). The only difference between ZB and ZB2 lies on the choice of the constant  $c$  in the tuning parameter  $\eta = \sqrt{c(\log p)/n}$  in (12) of their paper. In ZB, we set  $c = 2$  as suggested by the authors while in ZB2, we let  $c = 10^{-3}$ .

In (1), the rows of  $\mathbf{X}$  are considered to be i.i.d realizations from  $N(0, \boldsymbol{\Sigma})$  with  $\Sigma_{jj} = 1$  under two scenarios: (i)  $\Sigma_{j,k} = 0.9^{|j-k|}$  (denoted as Tp); (ii)  $\Sigma_{j,k} = 0.8$  for all  $j \neq k$  (denoted as Eq). To generate  $\beta$ , we consider the following two cases,

Case 1:  $\beta_j \stackrel{i.i.d.}{\sim} U(0, 4)$  with  $s_0 = 3, 5, 10, 15$ .

Case 2: Half of the non-zero  $\beta_j$ 's are independently generated from  $U(0, 0.5)$  and the rest are

generated from  $U(2.5, 3)$  with  $s_0 = 4, 8, 12, 16$ .

The errors are independently generated from (a) the standard normal distribution; (b) the studentized  $t(4)$  distribution, i.e.,  $t(4)/\sqrt{2}$ ; (c) the centralized and studentized Gamma(4,1) distribution, i.e.,  $(\text{Gamma}(4, 1) - 4)/2$ . The simulation results for (b) and (c) are summarized in the supplementary material. To save space, we only included the results of BCH, ZB and ZB2 for case (a). Throughout the simulations, we set  $n = 100$ ,  $p = 500$  and the nominal level  $1 - \alpha = 0.95$ . All the simulation results are based on 100 independent simulation runs.

We summarize the empirical coverage probabilities, the corresponding confidence interval lengths and the absolute value of the overall normalized bias defined as

$$\text{Bias} = \frac{|\sqrt{n}R(v_j, \beta_{-j})|}{\sqrt{\hat{\sigma}^2 n^{-1} \|v_j\|^2}} \quad (24)$$

for both the active set and the inactive set in Figures 5-8. The R code of Javanmard and Montanari (2014) makes a finite sample adjustment. To avoid unfair comparison, we do not include their method in the bias comparison. As inverting the test statistic in Zhu and Bradic (2018b) doesn't provide a closed form of confidence interval, the interval lengths of ZB and ZB2 are numerically calculated by using the bisection-type method. To avoid computational burden therein, we only calculate the lengths of 5 confidence intervals of ZB and ZB2 for inactive set in each simulation runs.

We observe that (i) BRP and MBRP generally provide more accurate coverage for the active set in comparison to DB and JM. The coverage probability for the active set based on DB can be significantly lower than the nominal level. While BCH shows similar or slightly higher coverage rate than BRP for the Toeplitz covariance structure, its coverage rate is lower than the nominal level in the equicorrelation case; (ii) The interval length of BCH is generally similar or wider than the lengths of BRP and MBRP, which is in turn wider than that of DB for the active set. Both ZB and ZB2 tend to provide wider confidence intervals compared to the other methods. (iii) For the equicorrelation covariance structure and  $s_0 \geq 10$ , ZB2 delivers the most accurate coverage rate followed by MBRP. In contrast, the other methods significantly undercover in these cases. (iv) The better coverage of the active set for our method is closely related to the smaller bias. Interestingly, the coverage rate for the inactive set seems not sensitive to the bias; (v) The computation time of our method is between those of DB and ZB as shown in Table 1; (vi) The bias associated with the active set tends to be larger than that with the inactive set especially in the case of Toeplitz covariance. BRP seems to overallly reduce the bias associated with both the active and inactive sets in such case; (vii) The coverage rate for the inactive set is usually close or above the nominal level for all methods except for ZB. According to our extensive simulations, the over-coverage is partly caused by the overestimation of the noise level as illustrated in Figure 22 in the supplementary material. Overall, our proposed method appears to outperform DB, JM, BCH and ZB in terms of coverage accuracy.

Figures 9-10 plot the bias and length of BRP and MBRP against  $C_2$  selected by the procedure



in Section 5. It is interesting to note that for BRP, the interval width generally increases while the bias decreases with  $C_2$ . The pattern is less obvious for MBRP with most of the values of  $C_2$  concentrate around the lower end of the grid points in (25).

## 6.2 Confidence interval for a sparse linear combination of regression coefficients

In this subsection, we investigate the finite sample performance of the method in Section 4. We consider the case where a linear contrast for two coefficients is of interest. We set the true regression coefficient  $\beta = (b_1, b_1, b_2, b_3, 0, \dots, 0)^\top$ , where  $b_1, b_2, b_3$  are drawn independently from  $U(0, 4)$ . Depending on  $a$ , we consider the following two cases:

1. Contrast 1:  $a = (1, -1, 0, \dots, 0)^\top$  and  $a^\top \beta = b_1 - b_1 = 0$ ;
2. Contrast 2:  $a = (0, 0, 1, -1, 0, \dots, 0)^\top$  and  $a^\top \beta = b_2 - b_3 \neq 0$ .

We adopt the same procedures as before for choosing the surrogate set and the tuning parameters but the results are based on 300 independent simulation runs. The configuration for  $\epsilon$  is the same as in the previous subsection. The results for t-distributed and gamma errors are presented in the supplementary material.

Figure 11 shows the empirical coverage rates, the corresponding confidence interval widths as well as the bias for each contrast. For the Toeplitz covariance structure, BRP and MBRP provide closer coverage rate to the nominal level but with wider interval length than DB does. In particular, MBRP delivers the smallest bias. Thus, the better coverage for our method is again closely related to the smaller bias in the finite sample. For the equicorrelation covariance structure, the coverage rates of all the methods are close to the nominal level. We also note that ZB2 provides satisfactory coverage probabilities while ZB significantly undercovers in the case of Toeplitz covariance structure. Similar to the case for a single regression coefficient, the lengths of ZB and ZB2 are generally wider than those of the other methods.

## 6.3 Real data analysis

As a real data application, we consider a dataset of riboflavin (vitamin  $B_2$ ) production by *Bacillus subtilis*. The dataset is available in the R package `hdi` and has also been analyzed in van de Geer et al. (2014) and Javanmard and Montanari (2014). It contains  $n = 71$  observations of  $p = 4088$  covariates of gene expressions and a response of riboflavin production. We model the data using (1) and consider the following multiple hypothesis testing for the significance of each gene:

$$H_{j,0} : \beta_j = 0 \quad \text{for } j = 1, \dots, 4088.$$

We use Theorem 1 and Corollary 2 to calculate the p-values based on BRP and MBRP respectively. The Holm procedure is adopted for multiplicity adjustment with the 5% significance level. Neither of our methods finds any significant predictors, which is also the case for DB while there turn out to be two significant genes YXLD-at and YXLE-at identified by JM.

## 7 Concluding remark

We have proposed a new method to find the projection direction in the debiased Lasso estimator and demonstrated its advantage over the original debiased Lasso estimator in van de Geer et al. (2014) and the method in Javanmard and Montanari (2014). The main contributions of this work are summarized below.

- We propose a new formulation to estimate the projection direction by properly balancing the biases associated with the strong and weak signals respectively.
- We show that the set of strong signals can be consistently estimated and establish the asymptotic normality of the proposed estimator.
- We further propose a modified estimator which can lead to a smaller order of bias comparing to the original debiased Lasso both theoretically and empirically.
- We generalize our idea to conduct inference for a sparse linear combination of regression coefficients.

As for future research, we expect that our method can be extended to other settings such as the generalized linear models, the Cox proportional hazards model and nonparametric additive models.

## 8 Supplementary Material

We empirically investigate the sensitiveness of our method to the choice of tuning parameters in Section 8.1. Section 8.2 provides every figure and table for Sections 2 and 6 in the main paper. Technical details and additional numerical results are gathered in Sections 8.3 and 8.4, respectively.

### 8.1 Empirical analysis of tuning parameters

We empirically investigate the sensitiveness of our method to the choice of tuning parameters. Throughout this subsection, we suppose the rows of  $\mathbf{X} \in \mathbb{R}^{100 \times 500}$  are i.i.d realizations from  $N(0, \Sigma)$  with  $\Sigma_{j,k} = 0.9^{|j-k|}$  (Toeplitz) or  $\Sigma_{j,k} = 0.8$  (Equicorrelation) for  $j \neq k$  and  $\Sigma_{jj} = 1$ . Regression coefficients  $\beta_j$ 's are generated by either Case 1 with  $s_0 = 10$  or Case 2 with  $s_0 = 4$  as described in Section 6. The errors are independently generated from the standard normal distribution. The nominal level is 95% and results are based on 100 independent simulation runs.

We first explore the effect of  $C_0$  on the estimation of the surrogate set and the impact of  $C_1$  and  $C_2$  on the coverage rate and interval width of the BRP-based confidence interval. The results for  $\beta_j$  generated from Case 2 with  $s_0 = 4$  and Toeplitz covariance  $\Sigma$  are summarized in Figure 1. As seen from Panel A, the surrogate set  $\mathcal{A}(\tau)$  with  $\tau = 2$  correctly identifies the large coefficients when  $C_0 \geq 2$ . Panels B-D provide the average coverage rate, bias and length of the BRP-based confidence intervals for the active set over a prespecified set of grid points for  $(C_1, C_2)$ . The coverage probability and interval width both tend to increase with the values of  $C_1$  and  $C_2$ . These results

appear to suggest that fixing one parameter at a reasonably large value while choosing the other parameter to balance the coverage probability and interval width would generally deliver similar results as simultaneously selecting the two parameters.

To confirm this intuition, we set  $C_0 = 2$ ,  $C_1 = 8$  and use the procedure in Section 5 to select  $C_2$  over the following prespecified grid points

$$\{c_{2,j,(k)}\}_{k=1}^K = \{0.3, 0.6, \dots, 14.7, 15.0\}. \quad (25)$$

We denote the corresponding procedures by “Fix-BRP” and “Fix-MBRP” and compare their performance with the procedures that select all tuning parameters automatically using the method in Section 5. Notice that fixing  $C_0$  and  $C_1$  would significantly ease the computational burden. Figure 2 presents the empirical coverage probabilities and lengths of the 95% confidence intervals and the normalized overall bias as in (24). Fix-BRP and Fix-MBRP perform equally well in terms of the coverage accuracy and bias as compared to BRP and MBRP but with a much lower computational cost. Indeed similar results are observed for the other simulation setups in Section 6.1. For the rest of the paper, we shall adopt the above procedure by fixing  $C_0$  and  $C_1$  to implement the proposed method.

Finally, we study the impact of  $B$  and  $\tau$ . Figure 3 summarizes the performance of the BRP and MBRP-based confidence intervals with different values of  $B$  and  $\tau$ . The results are not sensitive to the bootstrap sample size  $B$ . We also observe that a larger  $\tau$  tends to deliver higher coverage for MBRP in the equicorrelation case. Unreported numerical studies show that similar phenomenon can be observed for the other simulation setups. In Section 6 below, we shall fix  $B = 200$  and  $\tau = 2$ .

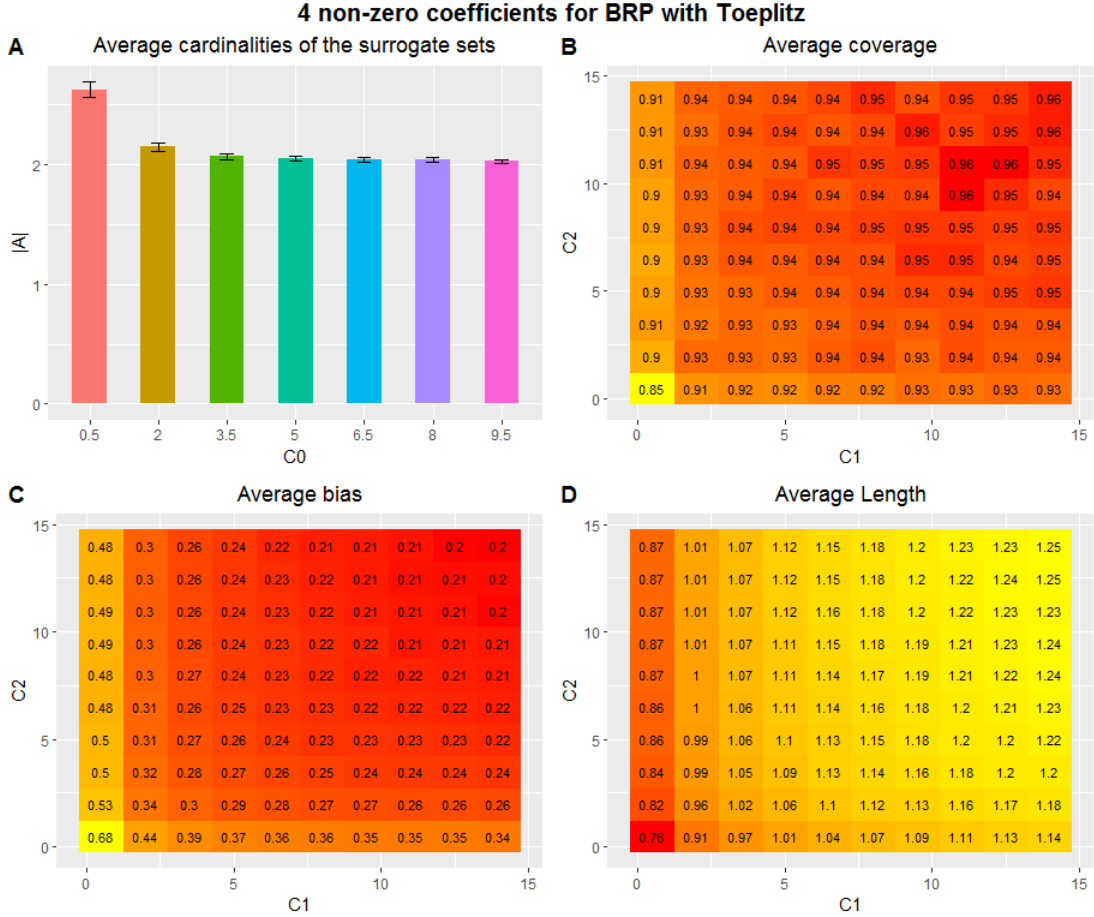


Figure 1: The first set of figures on empirical analysis of the effect of tuning parameters for BRP. Panel A shows the barplots of the average cardinality of  $\mathcal{A}(\tau)$  against  $C_0$ . Error bars in the barplots represent the interval within one standard error of the average value. Panel B (C or D) shows the heatmap of the average coverage rates (bias or length) by the BRP estimator over a prespecified grid points for  $(C_1, C_2)$ . The number represents the average coverage probability (bias or length) of the 95% confidence intervals for the active set.

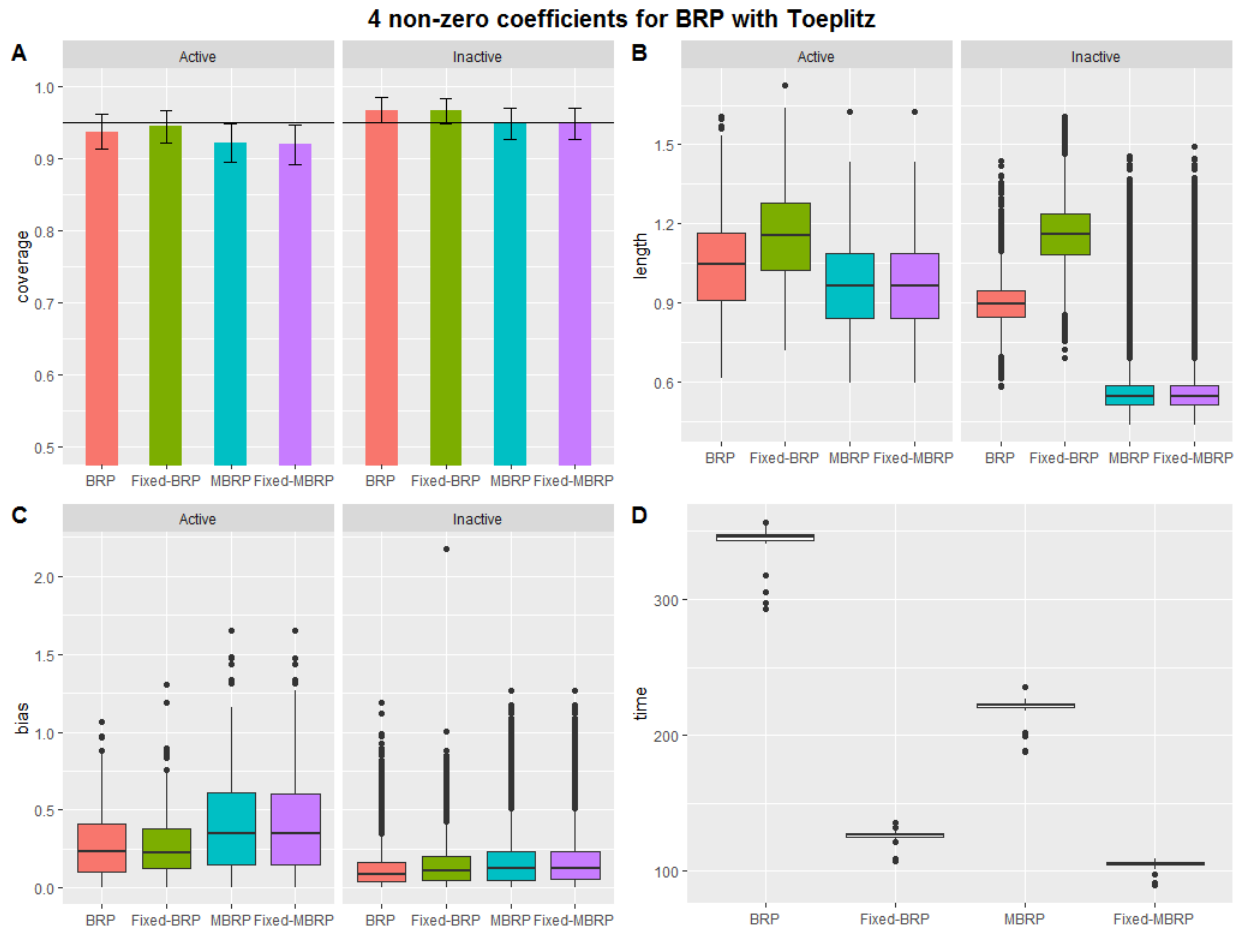


Figure 2: The second set of figures on empirical analysis of the effect of tuning parameters for BRP. Panel A shows the barplots of the empirical coverage and Panels B-C display the boxplots for the length and bias of the 95% confidence intervals of each method. In Panel A, the horizontal line indicates the nominal level and error bars represent the interval within one standard deviation of the empirical coverage. Panel D shows the boxplots of the computation time (in seconds) for each method.

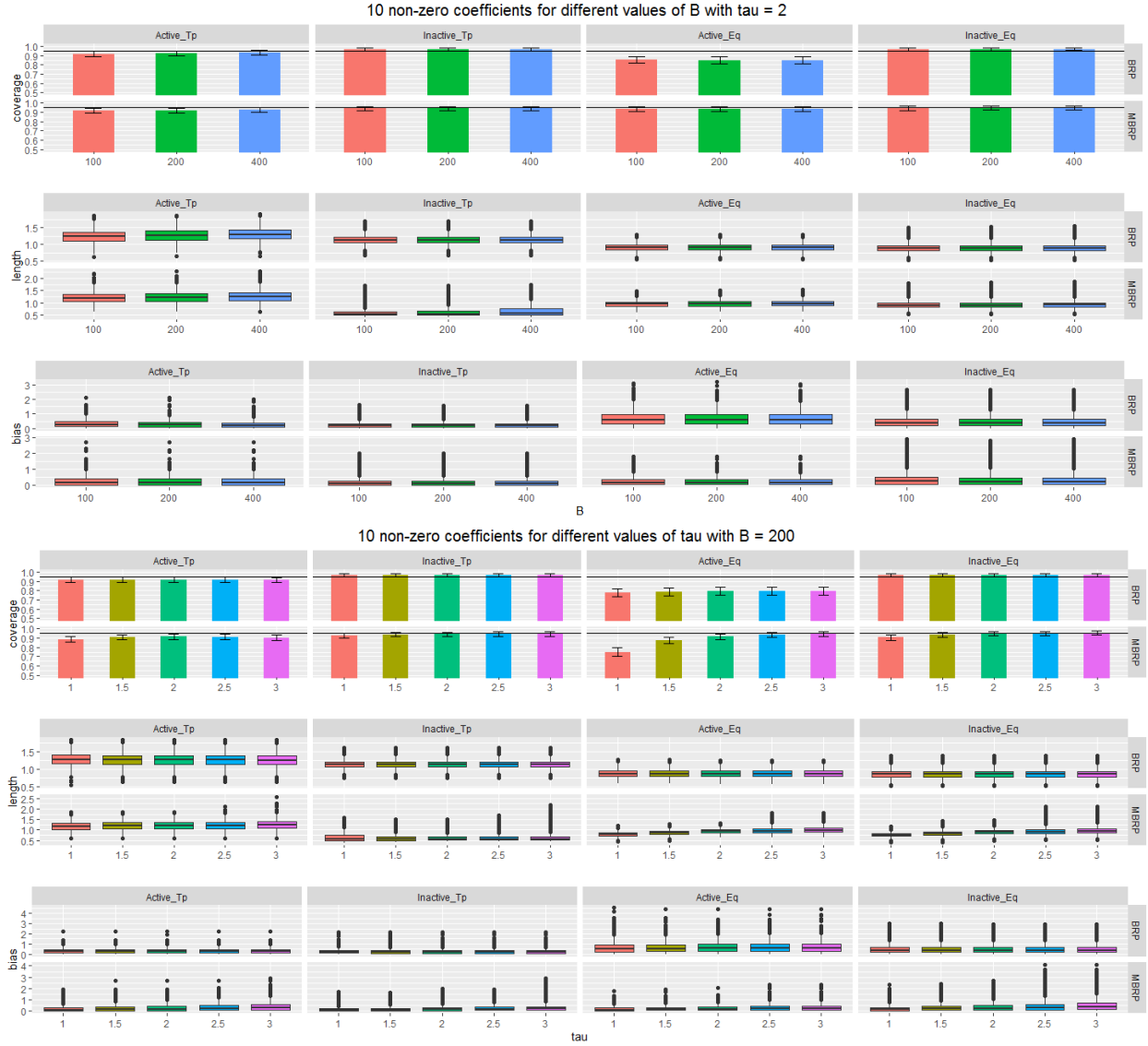


Figure 3: Set of figures on empirical analysis of the effect of tuning parameters  $B$  and  $\tau$ . Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals for both the active and inactive sets with different values of  $B$  and  $\tau$ . The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage. The data are independently generated from Case 1 with  $s_0 = 10$  and standard normal error as in Section 6.

## 8.2 Appendix for Sections 2 and 6

### 8.2.1 For Section 2

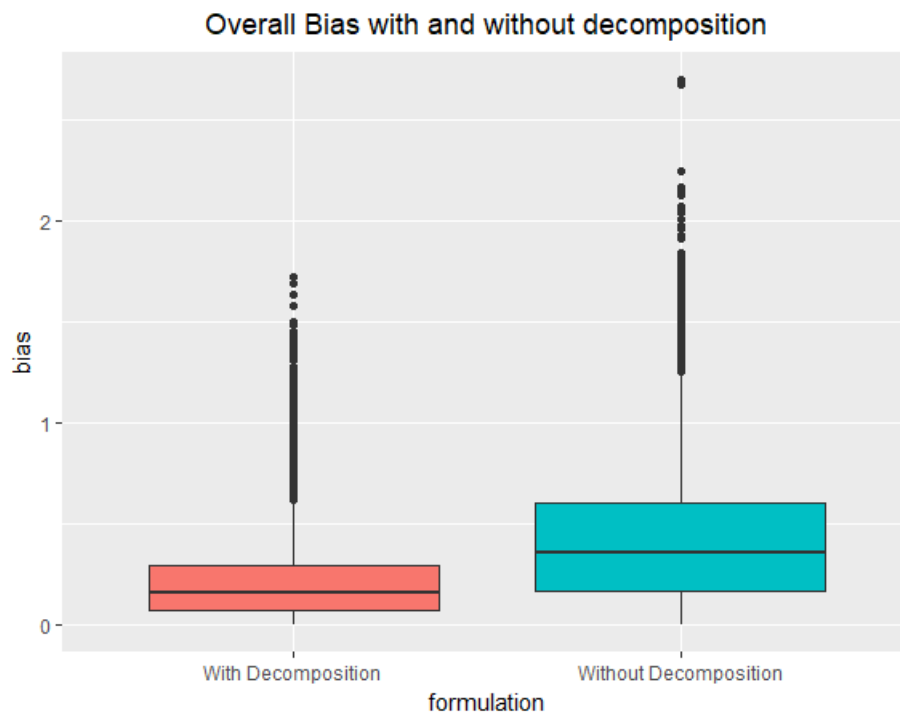


Figure 4: Boxplots of the absolute values of the normalized bias terms defined in (24) by “With Decomposition” and “Without Decomposition.” The non-zero  $\beta_j$ 's are independently generated from  $U(0, 4)$  with  $s_0 = 10$ . All the simulation settings are the same as the case with the Toeplitz covariance structure and standard normal error in Section 6. The results are based on 100 simulation runs.

## 8.2.2 For Section 6

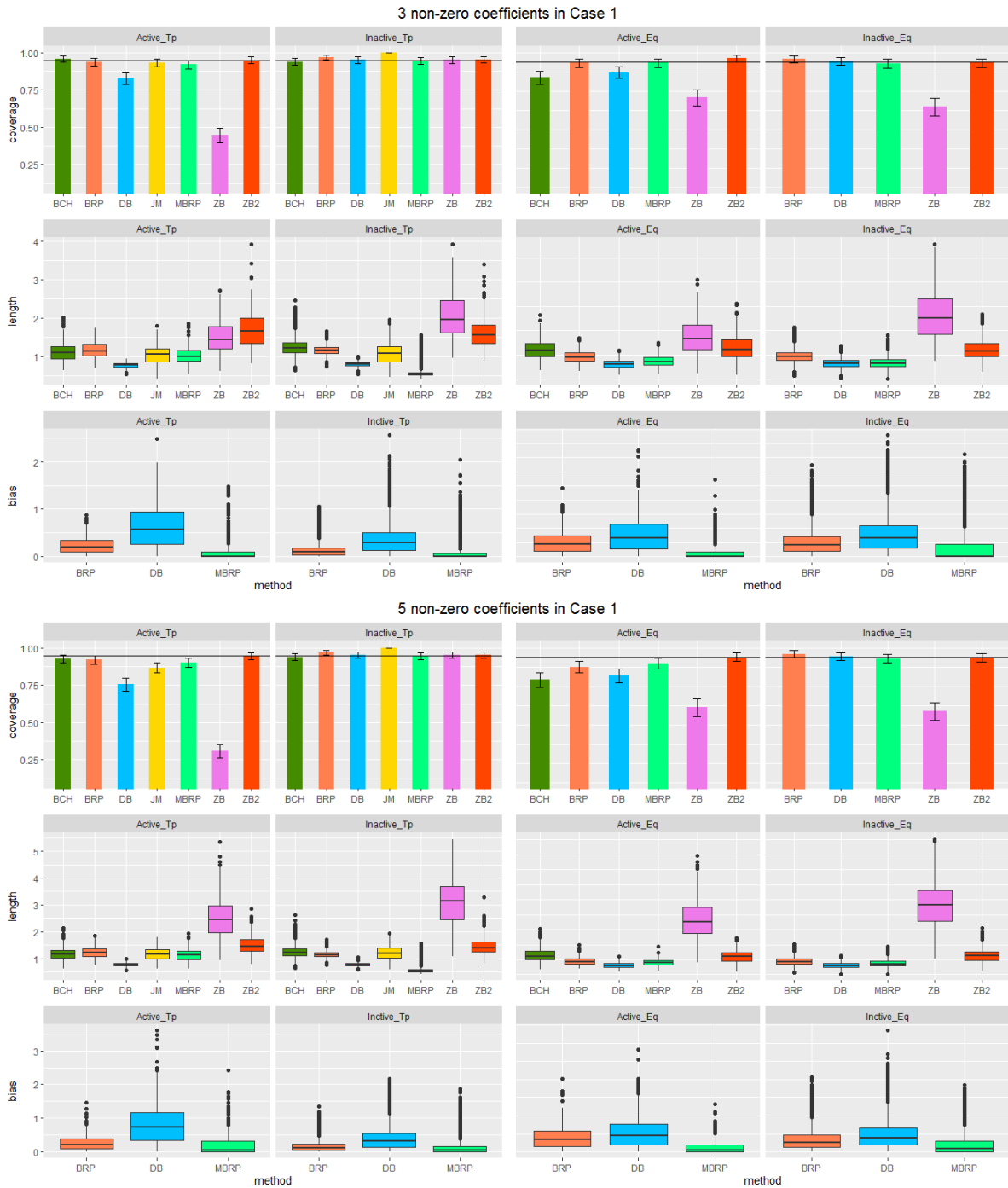


Figure 5: Simulation results for Case 1 with  $s_0 = 3, 5$  and standard normal random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.



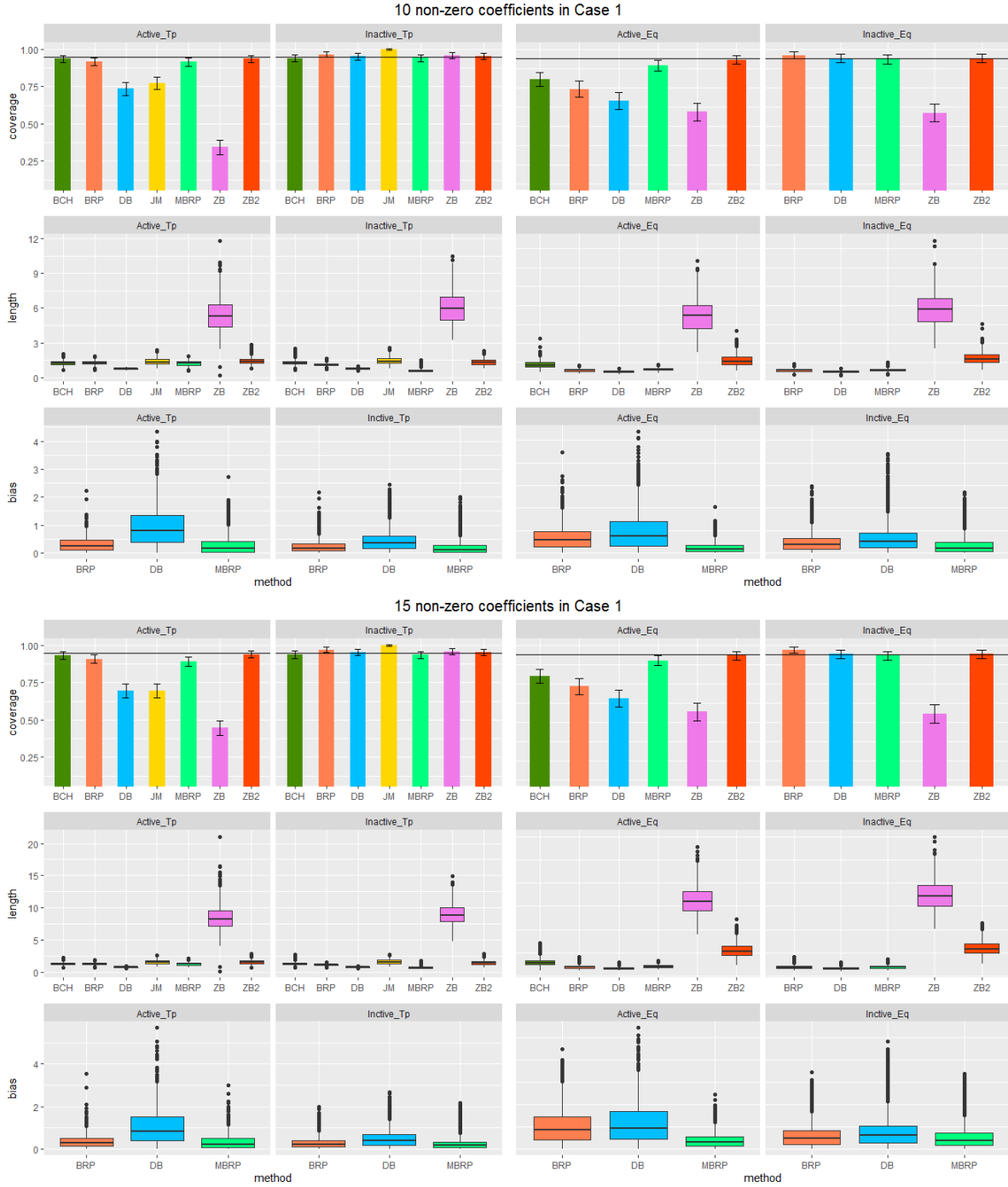


Figure 6: Simulation results for Case 1 with  $s_0 = 10, 15$  and standard normal random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

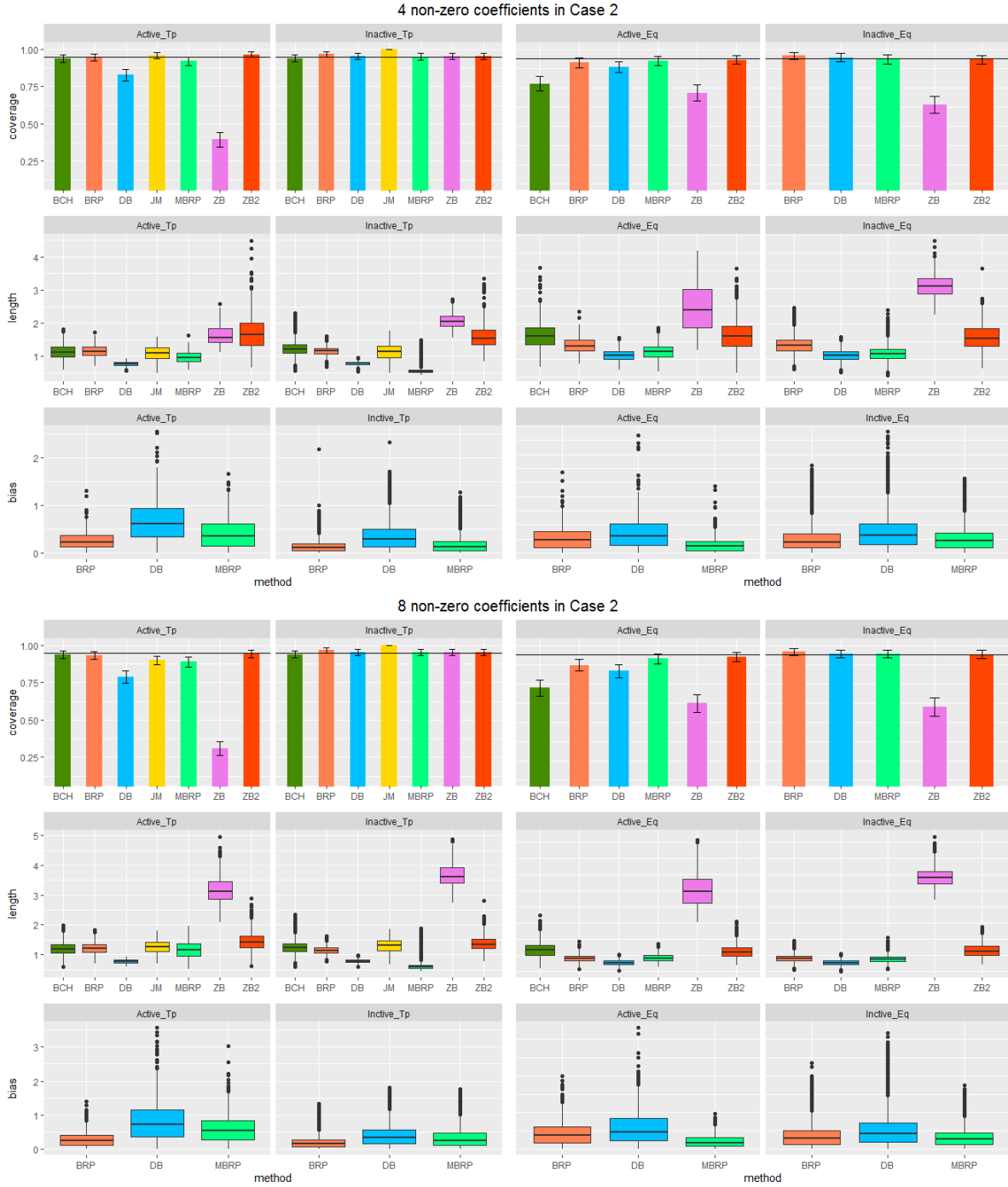


Figure 7: Simulation results for Case 2 with  $s_0 = 4, 8$  and standard normal random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

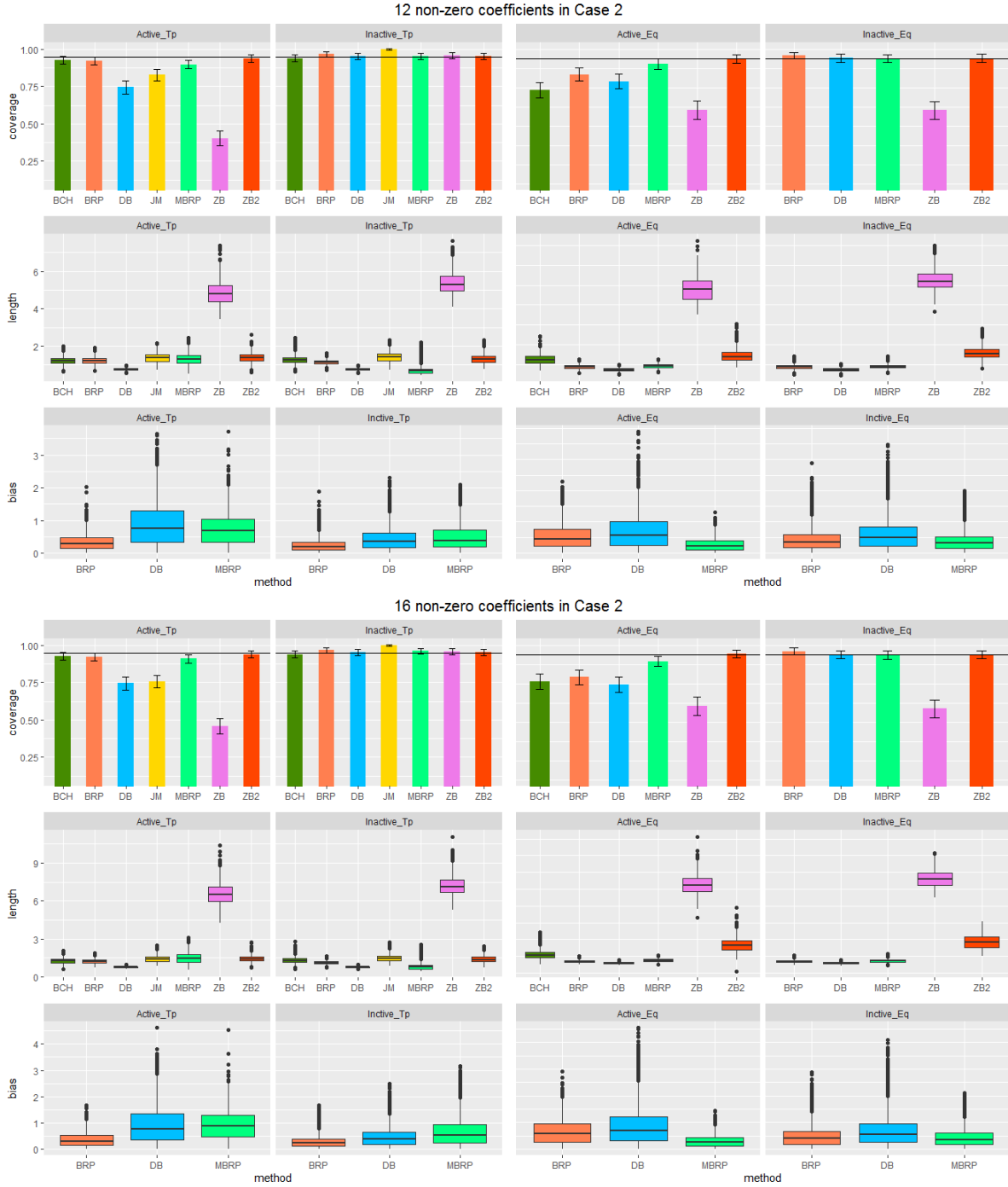


Figure 8: Simulation results for Case 2 with  $s_0 = 12, 16$  and standard normal random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

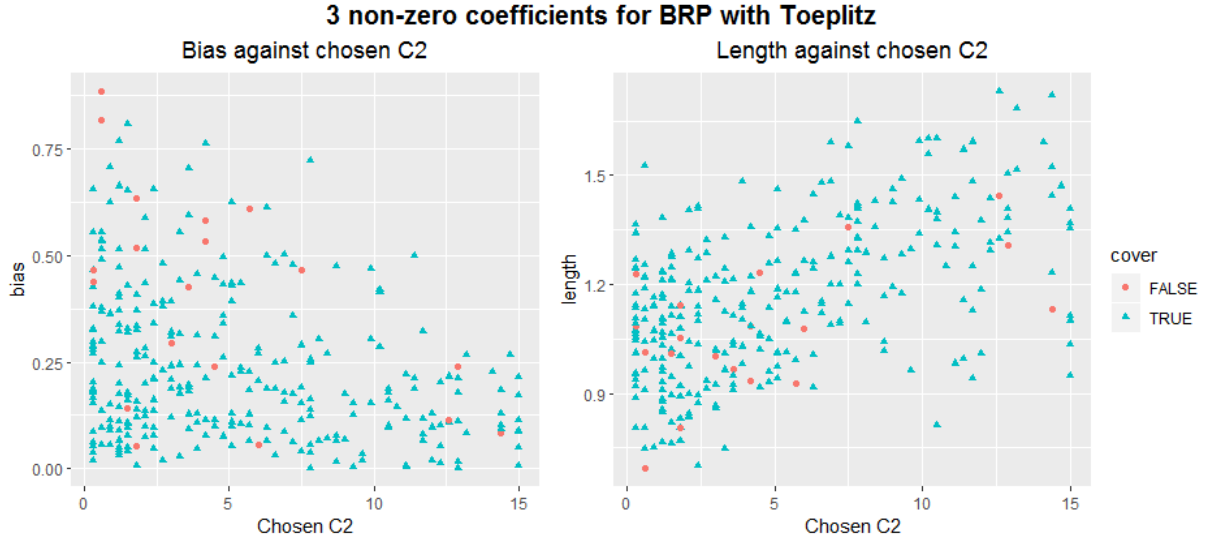


Figure 9: Scatterplots of the bias and length of the BRP-based confidence interval for the active set with  $s_0 = 3$  and Toeplitz covariance structure for  $\mathbf{X}$  against the selected  $C_2$ . The point shapes and colors indicate whether the constructed confidence intervals include the true parameter or not.

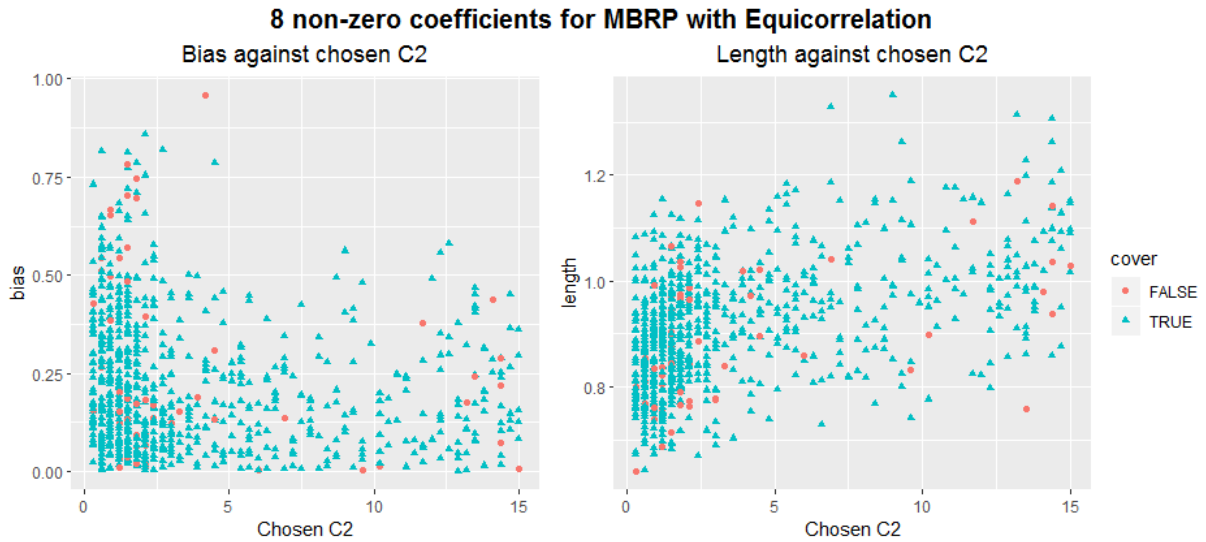


Figure 10: Scatterplots of the bias and length of the MBRP-based confidence interval for the active set with  $s_0 = 8$  and equicorrelation covariance structure for  $\mathbf{X}$  against the selected  $C_2$ . The point shapes and colors indicate whether the constructed confidence intervals include the true parameter or not.

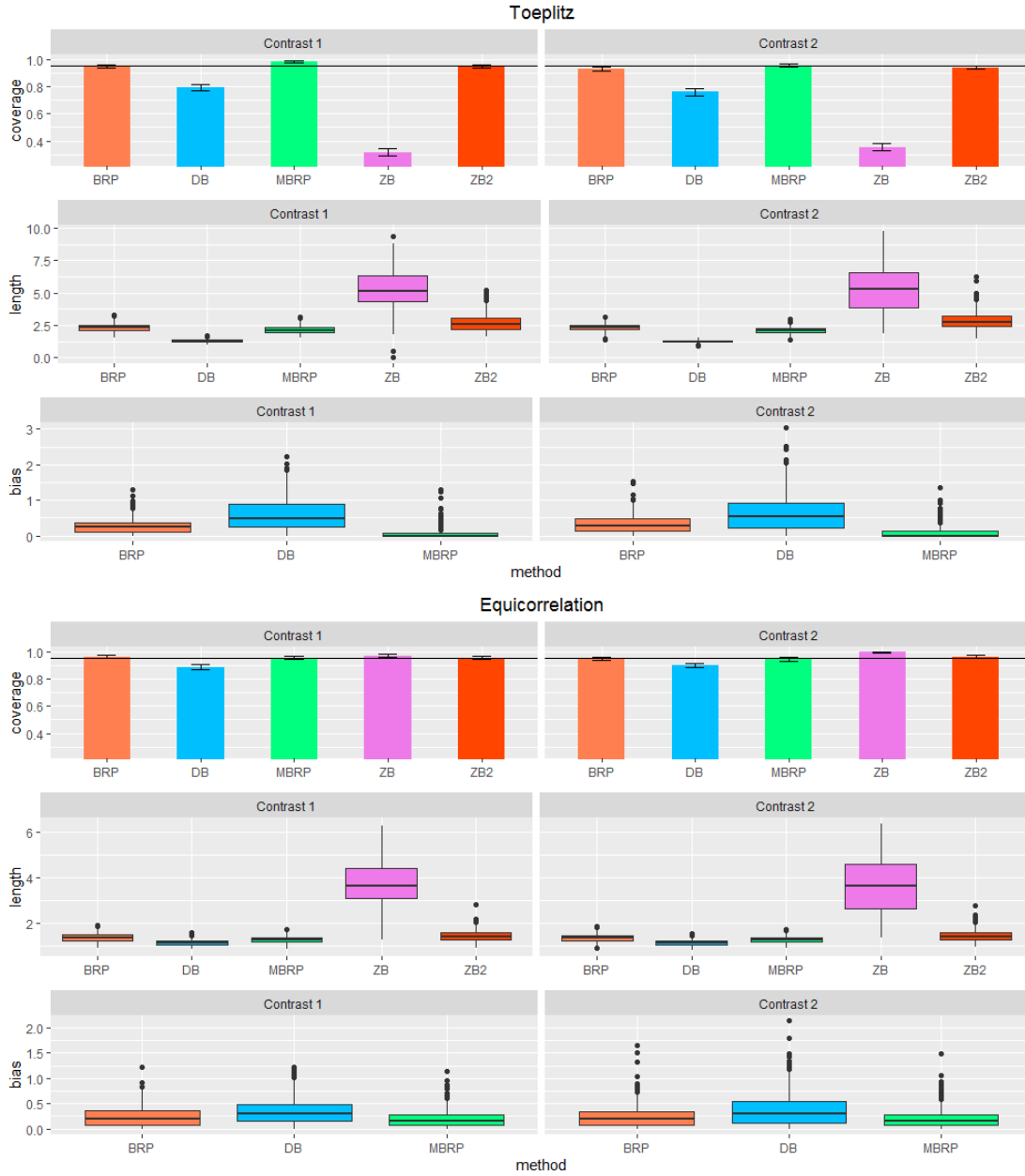


Figure 11: Simulation results for a sparse linear combination of  $\beta$  and standard normal random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals for each contrast. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

	Min	Q1	Median	Q3	Max
BRP	107.60	125.30	126.70	127.80	135.20
MBRP	89.65	104.34	105.21	106.35	109.03
DB	26.29	33.45	34.41	35.66	38.20
ZB	457.90	471.30	476.50	483.20	499.50

Table 1: Computation time (in seconds) of each method for constructing 500 confidence intervals calculated by the R package `mi` `crobenchmark`. The five number summaries are obtained based on 100 independent simulation runs.

### 8.3 Technical Details

#### 8.3.1 Concentration Inequalities

We first define several quantities which will appear throughout the supplementary material. Let  $\theta_j = X_j - \mathbf{X}_{-j}b_{-j}$  and

$$b_{-j} = \operatorname{argmin}_{\tilde{b} \in \mathbb{R}^{p-1}} E \|\mathbf{X}_j - \mathbf{X}_{-j}\tilde{b}\|_2^2 = \Sigma_{-j,-j}^{-1} \Sigma_{-j,j}.$$

Define  $\kappa_1 = 2\kappa^2$ ,  $\kappa_{2j} = 2\kappa^2 \sqrt{\Lambda_{\min}^{-1} \Sigma_{j,j}}$  and  $\kappa_{3j} = 2\kappa^2 \Lambda_{\min}^{-1} \Sigma_{j,j}$ .

The following lemmas shows the concentration inequalities for sub-exponential and sub-gaussian random variables which are motivated by Lemmas 5.5, 5.15 and Propositions 5.10, 5.16 in [26].

**Lemma 1.** *Let  $X_1, \dots, X_N$  be i.i.d. mean-zero sub-exponential random variables with  $\|X_i\|_{\psi_1} = K_1$ . Then, for every  $a = (a_1, \dots, a_N)^\top \in \mathbb{R}^{N \times 1}$  and any  $t \geq 0$ , we have*

$$\mathbb{P} \left( \left| \sum_{i=1}^N a_i X_i \right| \geq t \right) \leq 2 \exp \left\{ - \min \left( \frac{t^2}{8e^2 \|a\|^2 K_1^2}, \frac{t}{4e K_1 \|a\|_\infty} \right) \right\}.$$

**Proof of Lemma 1.** We first derive an upper bound of the moment generating function of  $X_i$ . By expanding the exponential function in the Taylor series, we have

$$\begin{aligned} E[\exp(\lambda X_i)] &= E \left[ 1 + \lambda X_i + \sum_{p=2}^{\infty} \frac{(\lambda X_i)^p}{p!} \right] = 1 + \sum_{p=2}^{\infty} \frac{\lambda^p E[X_i^p]}{p!} \\ &\leq 1 + \sum_{p=2}^{\infty} \frac{\lambda^p (K_1 p)^p}{(p/e)^p} = 1 + \sum_{p=2}^{\infty} (e\lambda K_1)^p = 1 + \frac{(e\lambda K_1)^2}{1 - (e\lambda K_1)} \end{aligned}$$

provided that  $|e\lambda K_1| < 1$ . The inequality follows by the definition of sub-exponential norm

$$E[X_i^p] \leq (K_1 p)^p$$

and Stirling's approximation  $p! \geq (p/e)^p$ . In addition, if  $|e\lambda K_1| < 0.5$ , the quantity on the right

hand side can be bounded above by

$$1 + 2(e\lambda K_1)^2 \leq \exp(2(e\lambda K_1)^2).$$

Thus, combining all of the above implies

$$E[\exp(\lambda X_i)] \leq \exp(2(e\lambda K_1)^2) \quad \text{for } |\lambda| < \frac{1}{2eK_1}. \quad (26)$$

Next, for  $\lambda > 0$ , we have

$$\begin{aligned} \mathbb{P}\left(\sum_{i=1}^N a_i X_i \geq t\right) &= \mathbb{P}\left(\exp\left(\lambda \sum_{i=1}^N a_i X_i\right) \geq \exp(\lambda t)\right) \\ &\leq \exp(-\lambda t) E\left[\exp\left(\lambda \sum_{i=1}^N a_i X_i\right)\right] = \exp(-\lambda t) \prod_{i=1}^N E[\exp(\lambda a_i X_i)] \end{aligned}$$

by the exponential Markov inequality for  $\sum_{i=1}^N a_i X_i$ . If  $\lambda$  is small enough so that  $|\lambda| < (2eK_1\|a\|_\infty)^{-1}$ , (26) gives

$$\mathbb{P}\left(\sum_{i=1}^N a_i X_i \geq t\right) \leq \exp(-\lambda t) \prod_{i=1}^N \exp(2(e\lambda a_i K_1)^2) = \exp(-\lambda t + 2e^2\lambda^2\|a\|^2 K_1^2).$$

By choosing  $\lambda = \min(t(4e^2\|a\|^2 K_1^2)^{-1}, (2eK_1\|a\|_\infty)^{-1})$ , we obtain

$$\mathbb{P}\left(\sum_{i=1}^N a_i X_i \geq t\right) \leq \exp\left\{-\min\left(\frac{t^2}{8e^2\|a\|^2 K_1^2}, \frac{t}{4eK_1\|a\|_\infty}\right)\right\}.$$

The second term in min can be obtained as follows. When  $\lambda = (2eK_1\|a\|_\infty)^{-1}$ , we have

$$-\lambda t + 2e^2\lambda^2\|a\|^2 K_1^2 = -\frac{t}{2eK_1\|a\|_\infty} + \frac{\|a\|^2}{2\|a\|_\infty^2} \leq -\frac{t}{4eK_1\|a\|_\infty}$$

where the last inequality follows as

$$\lambda = \frac{1}{2eK_1\|a\|_\infty} \leq \frac{t}{(4e^2\|a\|^2 K_1^2)}$$

which implies

$$\frac{\|a\|^2}{\|a\|_\infty} \leq \frac{t}{2eK_1}.$$

By repeating the same argument for  $-X_i$ , we get the same bound for  $\mathbb{P}(-\sum_{i=1}^N a_i X_i \geq t)$ , which completes the proof.  $\diamond$

**Lemma 2.** *Let  $X_1, \dots, X_N$  be i.i.d. mean-zero sub-gaussian random variables with  $\|X_i\|_{\psi_2} = K_2$ . Then, we have the following results.*

1. For any  $|\omega_1| \leq 1$ ,

$$E \left[ \exp \left( \omega_1^2 \frac{X_i^2}{4eK_2^2} \right) \right] \leq \exp(\omega_1^2). \quad (27)$$

2. For  $\omega_2 \in \mathbb{R}$ ,

$$E[\exp(\omega_2 X_i)] \leq \exp(8eK_2^2 \omega_2^2). \quad (28)$$

3. For every  $a = (a_1, \dots, a_N) \in \mathbb{R}^N$  and any  $t \geq 0$ ,

$$\mathbb{P} \left( \left| \sum_{i=1}^N a_i X_i \right| \geq t \right) \leq 2 \exp \left( -\frac{t^2}{32eK_2^2 \|a\|^2} \right). \quad (29)$$

**Proof of Lemma 2.** Let  $Y_i = X_i/(2\sqrt{e}K_2)$ . We note that, for  $|\omega_1^2/2| < 1$ ,

$$\begin{aligned} E[\exp(\omega_1^2 Y_i^2)] &= 1 + \sum_{k=1}^{\infty} \frac{\omega_1^{2k} E[Y_i^{2k}]}{k!} \\ &\leq 1 + \sum_{k=1}^{\infty} \frac{1}{(4e)^k} \frac{(2\omega_1^2 k)^k}{(k/e)^k} = \sum_{k=0}^{\infty} \left( \frac{\omega_1^2}{2} \right)^k = \left( 1 - \frac{\omega_1^2}{2} \right)^{-1} \end{aligned}$$

by the Taylor series expansion of the exponential function and Stirling's approximation. We can further bound

$$E[\exp(\omega_1^2 Y_i^2)] \leq \exp(\omega_1^2) \quad \text{for } |\omega_1| \leq 1$$

by using the inequality  $(1-x)^{-1} \leq \exp(2x)$  for  $0 \leq x \leq 0.5$ , which completes (27).

For (28), we notice that

$$E[\exp(\omega Y_i)] \leq E[\omega Y_i + \exp(\omega^2 Y_i^2)] \leq \exp(\omega^2) \quad (30)$$

for  $|\omega| \leq 1$  where the first inequality follows by  $e^x \leq x + e^{x^2}$  for any  $x \in \mathbb{R}$  and the second one does by (27). If  $|\omega| \geq 1$ , we have

$$E[\exp(\omega Y_i)] \leq \exp(\omega^2) E[\exp(Y_i^2)] \leq \exp(\omega^2 + 1) \leq \exp(2\omega^2) \quad (31)$$

due to  $\omega Y_i \leq \omega^2 + Y_i^2$  for any  $\omega, Y_i$  and (27). Thus, combining (30) with (31) gives

$$E[\exp(\omega Y_i)] \leq \exp(2\omega^2).$$

for any  $\omega \in \mathbb{R}$ . Letting  $\omega_2 = \omega/(2\sqrt{e}K_2)$  completes (28).



For (29), notice that

$$\begin{aligned} E \left[ \exp \left( \omega_2 \sum_{i=1}^N a_i X_i \right) \right] &= \prod_{i=1}^N E \left[ \exp(\omega_2 a_i X_i) \right] \\ &\leq \prod_{i=1}^N \exp(8eK_2^2 \omega_2^2 a_i^2) = \exp(8eK_2^2 \omega_2^2 \|a\|^2). \end{aligned}$$

For  $\omega_2 \geq 0$ , we have

$$\begin{aligned} \mathbb{P} \left( \sum_{i=1}^N a_i X_i \geq t \right) &= \mathbb{P} \left( \exp \left( \omega_2 \sum_{i=1}^N a_i X_i \right) \geq \exp(\omega_2 t) \right) \\ &\leq \exp(-\omega_2 t) E \left[ \exp \left( \omega_2 \sum_{i=1}^N a_i X_i \right) \right] \\ &\leq \exp(-\omega_2 t + 8e\omega_2^2 K_2^2 \|a\|^2) \\ &\leq \exp \left( -\frac{t^2}{32eK_2^2 \|a\|^2} \right) \end{aligned}$$

and the same bound can be obtained for  $\mathbb{P} \left( -\sum_{i=1}^N a_i X_i \geq t \right)$ . Thus, combining those bounds gives (29).  $\diamond$

### 8.3.2 Technical details in Section 3

**Lemma 3.** *Under Assumption 5,*

$$\mathbb{P} \left( n^{-1} \|\theta_j^\top \mathbf{X}_{-j}\|_\infty \geq \varepsilon_{0j} \sqrt{\frac{\log p}{n}} \right) \leq 2 \exp \left\{ \left( 1 - \frac{1}{8e^2} \frac{\varepsilon_{0j}^2}{(\kappa_{0j})^2} \right) \log p \right\}$$

for  $0 < \varepsilon_{0j} \leq \kappa_{0j} \sqrt{n(\log p)^{-1}}$ .

**Proof of Lemma 3.** Let  $Z = (Z_1 \cdots Z_{p-1}) = n^{-1} (X_j^\top \mathbf{X}_{-j} - b_{-j}^\top \mathbf{X}_{-j}^\top \mathbf{X}_{-j})$ . Then we have

$$Z = \frac{1}{n} \sum_{i=1}^n (X_{i,j} - b_{-j}^\top X_{i,-j}) X_{i,-j}^\top$$

where  $X_{i,j}$  is the value of the  $j$ th predictor of the  $i$ th observation and

$$X_{i,-j}^\top = (X_{i,1} \cdots X_{i,j-1}, X_{i,j+1} \cdots X_{i,p}).$$

Fix some  $k \in \{1, 2, \dots, p\} \setminus \{j\}$  and let  $Z_{i,j}^{(k)} = (X_{i,j} - b_{-j}^\top X_{i,-j}) X_{i,-j}^{(k)}$ , where  $X_{i,-j}^{(k)}$  denotes the  $k$ th element of  $X_{i,-j}$ . Then  $Z_k = n^{-1} \sum_{i=1}^n Z_{i,j}^{(k)}$ , where  $E[Z_{i,j}^{(k)}] = 0$  and  $Z_{i,j}^{(k)}$ 's are independent across  $1 \leq i \leq n$ .

We derive an upper bound for  $\|Z_{i,j}^{(k)}\|_{\psi_1}$ . Notice that

$$\begin{aligned}\|Z_{i,j}^{(k)}\|_{\psi_1} &= \|(X_{i,j} - b_{-j}^\top X_{i,-j})X_{i,-j}^{(k)}\|_{\psi_1} \leq 2\|X_{i,j} - b_{-j}^\top X_{i,-j}\|_{\psi_2}\|X_{i,-j}^{(k)}\|_{\psi_2} \\ &= 2\|X_{i,\cdot}^\top \gamma_{-j}\|_{\psi_2}\|X_{i,-j}^{(k)}\|_{\psi_2} \\ &\leq 2\kappa^2\|\gamma_{-j}\|_2 \\ &\leq 2(1 + \sqrt{\Lambda_{\min}^{-1}\Sigma_{j,j}})\kappa^2,\end{aligned}$$

where  $X_{i,\cdot}^\top = (X_{i,j}, X_{i,-j}^\top)$  and  $\gamma_{-j}^\top = (1, -b_{-j}^\top)$ . Here, the first inequality holds from the fact that  $\|XY\|_{\psi_1} \leq 2\|X\|_{\psi_2}\|Y\|_{\psi_2}$  for any two random variables  $X, Y$ ; the second inequality comes from

$$q^{-1/2}(E|X_{i,\cdot}^\top \gamma_{-j}|^q)^{1/q} = \|\gamma_{-j}\|_2 q^{-1/2}\{E|X_{i,\cdot}^\top (\gamma_{-j}/\|\gamma_{-j}\|)|^q\}^{1/q} \leq \|\gamma_{-j}\|_2 \kappa$$

and the third inequality follows from

$$\|\gamma_{-j}\|_2 = \sqrt{1 + \|b_{-j}\|^2} \leq 1 + \|b_{-j}\| \leq 1 + \sqrt{\lambda_{\max}(\Sigma_{-j,-j}^{-1})\Sigma_{j,j}}.$$

By Lemma 1, for any  $\varepsilon > 0$ , we have

$$\mathbb{P}\left(\frac{1}{n}\left|\sum_{i=1}^n Z_{i,j}^{(k)}\right| \geq \varepsilon\right) \leq 2 \exp\left\{-n \min\left(\frac{1}{8e^2}\left(\frac{\varepsilon}{\kappa_{0j}}\right)^2, \frac{1}{4e}\frac{\varepsilon}{\kappa_{0j}}\right)\right\}.$$

Choosing  $\varepsilon = \varepsilon_{0j}\sqrt{n^{-1}\log p}$  and assuming that  $n \geq \varepsilon_{0j}^2(\kappa_{0j})^{-2}\log p$ , then

$$\mathbb{P}\left(\frac{1}{n}\left|\sum_{i=1}^n Z_{i,j}^{(k)}\right| \geq \varepsilon_{0j}\sqrt{\frac{\log p}{n}}\right) \leq 2 \exp\left\{-\frac{1}{8e^2}\frac{\varepsilon_{0j}^2}{(\kappa_{0j})^2}\log p\right\}.$$

The result follows from the union bound over  $k \in \{1, 2, \dots, p-1\}$ .  $\diamond$

An implication of Lemma 3 is that

$$n^{-1}\|\theta_j^\top \mathbf{X}_{-j}\|_\infty \leq \varepsilon_{0j}\sqrt{\frac{\log p}{n}} \quad (32)$$

with probability tending to 1 for a fixed  $\varepsilon_{0j}$  such that  $\varepsilon_0^2 > (\kappa_{0j})^2 8e^2$ . We introduce an additional result below for a later use.

**Lemma 4.** *Under Assumption 5, we have*

$$\begin{aligned}\mathbb{P}\left(\left|\frac{n}{\theta_j^\top X_j} - \frac{1}{\Sigma_{j \setminus -j}}\right| \leq \varepsilon_{1j}\right) &\geq 1 - 2 \exp\left\{-\frac{1}{8e^2}\left(\frac{\Sigma_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_1}\right)^2 n\right\} \\ &\quad - 2 \exp\left\{-\frac{1}{8e^2}\left(\frac{\Sigma_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_{2j}}\right)^2 n\right\}\end{aligned}$$

for  $0 < \varepsilon_{1j} \leq \min\{(\boldsymbol{\Sigma}_{j \setminus -j})^{-1}, 4 \min(\kappa_1, \kappa_{2j})(\boldsymbol{\Sigma}_{j \setminus -j})^{-2}\}$  and

$$\begin{aligned} \mathbb{P}\left(\left|\frac{\|\theta_j\|^2}{n} - \boldsymbol{\Sigma}_{j \setminus -j}\right| \leq \varepsilon_{2j}\right) &\geq 1 - 2 \exp\left\{-\frac{1}{8e^2} \left(\frac{\varepsilon_{2j}}{3\kappa_1}\right)^2 n\right\} - 2 \exp\left\{-\frac{1}{8e^2} \left(\frac{\varepsilon_{2j}}{6\kappa_{2j}}\right)^2 n\right\} \\ &\quad - 2 \exp\left\{-\frac{1}{8e^2} \left(\frac{\varepsilon_{2j}}{3\kappa_{3j}}\right)^2 n\right\} \end{aligned}$$

for  $0 < \varepsilon_{2j} \leq 3 \min(\kappa_1, 2\kappa_{2j}, \kappa_{3j})$ .

**Proof of Lemma 4.** We notice that

$$\theta_j^\top X_j = X_j^\top X_j - \sum_{i=1}^n \sum_{k=1}^{p-1} b_{-j,k} X_{i,-j}^{(k)} X_{i,j},$$

where  $b_{-j,k}$  is the  $k$ th element of  $b_{-j}$  and  $X_{i,-j}^{(k)}$  is the  $k$ th element of  $\mathbf{X}_{-j}$ . Then, we see that

$$\frac{\theta_j^\top X_j}{n} - \boldsymbol{\Sigma}_{j \setminus -j} = \frac{1}{n} \sum_{i=1}^n (X_{i,j}^2 - \boldsymbol{\Sigma}_{j,j}) - \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=1}^{p-1} b_{-j,k} X_{i,-j}^{(k)} X_{i,j} - \boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j} \right).$$

By Lemma 1,

$$\begin{aligned} \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n (X_{i,j}^2 - \boldsymbol{\Sigma}_{j,j})\right| \leq \delta_j\right) &\geq 1 - 2 \exp\left\{-\frac{1}{8e^2} \left(\frac{\delta_j}{\kappa_1}\right)^2 n\right\} \\ \mathbb{P}\left(\left|\frac{1}{n} \sum_{i=1}^n \left(\sum_{k=1}^{p-1} b_{-j,k} X_{i,-j}^{(k)} X_{i,j} - \boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j}\right)\right| \leq \delta_j\right) &\geq 1 - 2 \exp\left\{-\frac{1}{8e^2} \left(\frac{\delta_j}{\kappa_{2j}}\right)^2 n\right\} \end{aligned}$$

for  $0 < \delta_j \leq \min(\kappa_1, \kappa_{2j})$ . Also, for  $\varepsilon_{1j} \leq (\boldsymbol{\Sigma}_{j \setminus -j})^{-1}$ , we have

$$\begin{aligned} &\left\{ \left| \frac{n}{\theta_j^\top X_j} - \frac{1}{\boldsymbol{\Sigma}_{j \setminus -j}} \right| \geq \varepsilon_{1j} \right\} \\ &\subset \left[ \left\{ \left| \frac{1}{n} \sum_{i=1}^n (X_{i,j}^2 - \boldsymbol{\Sigma}_{j,j}) \right| \geq \frac{\boldsymbol{\Sigma}_{j \setminus -j}^2}{4} \varepsilon_{1j} \right\} \right. \\ &\quad \left. \cup \left\{ \left| \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=1}^{p-1} b_{-j,k} X_{i,-j}^{(k)} X_{i,j} - \boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j} \right) \right| \geq \frac{\boldsymbol{\Sigma}_{j \setminus -j}^2}{4} \varepsilon_{1j} \right\} \right]. \end{aligned}$$

Thus, for  $\varepsilon_{1j} \leq \min\{(\boldsymbol{\Sigma}_{j \setminus -j})^{-1}, 4 \min(\kappa_1, \kappa_{2j})(\boldsymbol{\Sigma}_{j \setminus -j})^{-2}\}$ , we have

$$\begin{aligned} \mathbb{P} \left( \left| \frac{n}{\boldsymbol{\theta}_j^\top \mathbf{X}_j} - \frac{1}{\boldsymbol{\Sigma}_{j \setminus -j}} \right| \leq \varepsilon_{1j} \right) &\geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\boldsymbol{\Sigma}_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_1} \right)^2 n \right\} \\ &\quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\boldsymbol{\Sigma}_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_{2j}} \right)^2 n \right\} \end{aligned}$$

which proves the first inequality. Next, we note that

$$\begin{aligned} \frac{\|\boldsymbol{\theta}_j\|^2}{n} - \boldsymbol{\Sigma}_{j \setminus -j} &= \underbrace{\left( \frac{\mathbf{X}_j^\top \mathbf{X}_j}{n} - \boldsymbol{\Sigma}_{j,j} \right)}_{(*)} - 2 \underbrace{\left( \frac{\mathbf{X}_j^\top \mathbf{X}_{-j}}{n} - \boldsymbol{\Sigma}_{j,-j} \right) \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j}}_{(**)} \\ &\quad + \underbrace{\boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \left( \frac{\mathbf{X}_{-j}^\top \mathbf{X}_{-j}}{n} - \boldsymbol{\Sigma}_{-j,-j} \right) \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j}}_{(***)}. \end{aligned}$$

The concentration inequalities for (\*) and (\*\*) are given respectively as

$$\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n (X_{i,j}^2 - \boldsymbol{\Sigma}_{j,j}) \right| \leq \frac{\varepsilon_{2j}}{3} \right) \geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\varepsilon_{2j}}{3\kappa_1} \right)^2 n \right\}, \quad (33)$$

and

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=1}^{p-1} b_{-j,k} X_{i,-j}^{(k)} X_{i,j} - \boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j} \right) \right| \leq \frac{\varepsilon_{2j}}{6} \right) \\ &\geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\varepsilon_{2j}}{6\kappa_{2j}} \right)^2 n \right\}, \end{aligned} \quad (34)$$

for  $0 < \varepsilon_{2j} \leq \min(3\kappa_1, 6\kappa_{2j})$ . Also, we notice that

$$(***) = \frac{1}{n} \sum_{i=1}^n \left( \sum_{k=1}^{p-1} X_{i,-j}^{(k)} b_{-j,k} \right)^2 - \boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j}.$$

Lemma 1 gives us

$$\begin{aligned} &\mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n \left( \left( \sum_{k=1}^{p-1} X_{i,-j}^{(k)} b_{-j,k} \right)^2 - \boldsymbol{\Sigma}_{j,-j} \boldsymbol{\Sigma}_{-j,-j}^{-1} \boldsymbol{\Sigma}_{-j,j} \right) \right| \leq \frac{\varepsilon_{2j}}{3} \right) \\ &\geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\varepsilon_{2j}}{3\kappa_{3j}} \right)^2 n \right\} \end{aligned} \quad (35)$$

for  $0 < \varepsilon_{2j} \leq 3\kappa_{3j}$ . Combining (33), (34) and (35) finishes the proof.  $\diamond$

The following result directly follows by Lemmas 3 and 4.

**Corollary 4.** *Let  $\check{v}_j = n\theta_j/(\theta_j^\top X_j)$  and  $\check{u}_j = n^{-1}\|\check{v}_j^\top \mathbf{X}_{-j}\|_\infty$ . Under Assumption 5,  $\check{v}_j$  satisfies  $\check{v}_j^\top X_j = n$ ,*

$$\begin{aligned} & \mathbb{P} \left( n^{-1}\|\check{v}_j\|^2 \leq \left( \frac{1}{\Sigma_{j \setminus -j}} + \varepsilon_{1j} \right)^2 (\Sigma_{j \setminus -j} + \varepsilon_{2j}) \right) \\ & \geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\varepsilon_{2j}}{3\kappa_{1j}} \right)^2 n \right\} - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\varepsilon_{2j}}{6\kappa_{2j}} \right)^2 n \right\} - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\varepsilon_{2j}}{3\kappa_{3j}} \right)^2 n \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\Sigma_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_{1j}} \right)^2 n \right\} - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\Sigma_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_{2j}} \right)^2 n \right\}, \end{aligned}$$

and

$$\begin{aligned} & \mathbb{P} \left( \check{u}_j \leq \varepsilon_{0j} \sqrt{\frac{\log p}{n}} \left( \frac{1}{\Sigma_{j \setminus -j}} + \varepsilon_{1j} \right) \right) \\ & \geq 1 - 2 \exp \left\{ \left( 1 - \frac{1}{8e^2} \frac{\varepsilon_{0j}^2}{(\kappa_{0j})^2} \right) \log p \right\} - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\Sigma_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_1} \right)^2 n \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \frac{\Sigma_{j \setminus -j}^2 \varepsilon_{1j}}{4\kappa_{2j}} \right)^2 n \right\}, \end{aligned}$$

for  $\varepsilon_{0j}, \varepsilon_{1j}, \varepsilon_{2j}$  given in Lemmas 3 and 4.

**Lemma 5.** *Let  $R_l = n^{-1}\hat{v}_l^\top \mathbf{X}_{-l}(\beta_{-l} - \hat{\beta}_{-l})$  where  $\hat{v}_l$  and  $\hat{u}_l$  denote the solution to (10). Then,*

$$\max_l R_l = o_p(1).$$

**Proof of Lemma 5.** According to the definition of  $\hat{v}_l$ ,

$$C_2 \frac{n}{\log p} \hat{u}_l^2 + n^{-1}\|\hat{v}_l\|^2 \leq C_2 \frac{n}{\log p} \check{u}_l^2 + n^{-1}\|\check{v}_l\|^2$$

where  $\hat{u}_l = n^{-1}\|\hat{v}_l^\top \mathbf{X}_{-l}\|_\infty$ ,  $\check{v}_l = n\theta_l/\theta_l^\top X_l$  and  $\check{u}_l = n^{-1}\|\check{v}_l^\top \mathbf{X}_{-l}\|_\infty$ . Then, we have

$$\sqrt{C_2 \frac{n}{\log p}} \hat{u}_l \leq C_2 \frac{n}{\log p} \check{u}_l^2 + n^{-1}\|\check{v}_l\|^2$$

which implies

$$\sqrt{C_2 \frac{n}{\log p}} \max_l n^{-1}\|\hat{v}_l^\top \mathbf{X}_{-l}\|_\infty \leq \left( \max_l (\Sigma_{l \setminus -l})^{-1} + \varepsilon'_1 \right)^2 \left( C_2 (\varepsilon'_0)^2 + \max_l \Sigma_{l \setminus -l} + \varepsilon'_2 \right)$$

with probability tending to 1 by (37). Therefore,

$$\max_l R_l \leq \max_l n^{-1} \|\hat{v}_l^\top \mathbf{X}_{-l}\|_\infty \sqrt{n} \|\hat{\beta} - \beta\|_1 = o_p(1)$$

by Assumptions 3 and 6. ◇

The following inequalities are direct consequences of Lemmas 3-4 and the definition of  $\hat{v}_l$ .

**Corollary 5.** *Let  $\hat{v}_l$  be the solution to (10). Then, we have*

$$\begin{aligned} & \mathbb{P} \left( \max_l n^{-1} \|\theta_l^\top \mathbf{X}_{-l}\|_\infty \geq \varepsilon'_0 \sqrt{\frac{\log p}{n}} \right) \leq 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{1}{\kappa_{0l}^2} \right) (\varepsilon'_0)^2 \log p + 2 \log p \right\}, \\ & \mathbb{P} \left( \max_l \frac{n}{|\theta_l^\top X_l|} \leq \max_l (\boldsymbol{\Sigma}_{l \setminus -l})^{-1} + \varepsilon'_1 \right) \geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{(\boldsymbol{\Sigma}_{l \setminus -l})^4}{(4\kappa_1)^2} \right) (\varepsilon'_1)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{(\boldsymbol{\Sigma}_{l \setminus -l})^4}{(4\kappa_{2l})^2} \right) (\varepsilon'_1)^2 n + \log p \right\}, \\ & \mathbb{P} \left( \max_l \frac{\|\theta_l\|^2}{n} \leq \max_l \boldsymbol{\Sigma}_{l \setminus -l} + \varepsilon'_2 \right) \geq 1 - 2 \exp \left\{ -\frac{1}{8e^2 (3\kappa_1)^2} (\varepsilon'_2)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{1}{(6\kappa_{2l})^2} \right) (\varepsilon'_2)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{1}{(3\kappa_{3l})^2} \right) (\varepsilon'_2)^2 n + \log p \right\}, \end{aligned}$$

and

$$\begin{aligned} & \mathbb{P} \left( \max_l n^{-1} \|\hat{v}_l\|^2 \leq M' \right) \geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{1}{(\kappa_{0l})^2} \right) (\varepsilon'_0)^2 \log p + 2 \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{(\boldsymbol{\Sigma}_{l \setminus -l})^4}{(4\kappa_1)^2} \right) (\varepsilon'_1)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{(\boldsymbol{\Sigma}_{l \setminus -l})^4}{(4\kappa_{2l})^2} \right) (\varepsilon'_1)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2 (3\kappa_1)^2} (\varepsilon'_2)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{1}{(6\kappa_{2l})^2} \right) (\varepsilon'_2)^2 n + \log p \right\} \\ & \quad - 2 \exp \left\{ -\frac{1}{8e^2} \left( \min_l \frac{1}{(3\kappa_{3l})^2} \right) (\varepsilon'_2)^2 n + \log p \right\}, \end{aligned}$$

where

$$M' = (\max_l (\boldsymbol{\Sigma}_{l \setminus -l})^{-1} + \varepsilon'_1)^2 (C_2 \varepsilon_0'^2 + \max_l \boldsymbol{\Sigma}_{l \setminus -l} + \varepsilon'_2) \quad (36)$$

for  $0 < \varepsilon'_0 \leq (\min_l \kappa_{0l}) \sqrt{n(\log p)^{-1}}$ ,  $0 < \varepsilon'_1 \leq \min_l \{ \min((\boldsymbol{\Sigma}_{l \setminus -l})^{-1}, 4 \min(\kappa_1, \kappa_{2l})(\boldsymbol{\Sigma}_{l \setminus -l})^{-2}) \}$  and  $0 < \varepsilon'_2 \leq \min_l (\min(3\kappa_1, 6\kappa_{2l}, 3\kappa_{3l}))$ .

Under Assumption 6, Corollary 5 implies that

$$\begin{aligned}
\max_l n^{-1} \|\theta_l^\top \mathbf{X}_{-l}\|_\infty &\leq \varepsilon'_0 \sqrt{\frac{\log p}{n}}, \\
\max_l \frac{n}{|\theta_l^\top X_l|} &\leq \max_l (\boldsymbol{\Sigma}_{l \setminus -l})^{-1} + \varepsilon'_1, \\
\max_l n^{-1} \|\theta_l\|^2 &\leq \max_l \boldsymbol{\Sigma}_{l \setminus -l} + \varepsilon'_2, \\
\max_l n^{-1} \|\hat{v}_l\|^2 &\leq M',
\end{aligned} \tag{37}$$

with probability tending to 1 for a fixed  $(\varepsilon'_0)^2 \min_l (\kappa_{0l})^{-2} > 16e^2$  and fixed  $\varepsilon'_1, \varepsilon'_2$  as in Corollary 5.

**Proof of Proposition 1.** Noting that  $(n^{-1/2} \|\hat{v}_l\|)^{-1} \leq \|X_l\|/\sqrt{n}$ , we have

$$\begin{aligned}
|T_l| &= \frac{\sigma}{\hat{\sigma}} \left| \frac{\sqrt{n}(\tilde{\beta}_l(\hat{v}_l) - \beta_l)}{\sigma n^{-1/2} \|\hat{v}_l\|} + \frac{\sqrt{n}\beta_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| \\
&= \frac{\sigma}{\hat{\sigma}} \left| \frac{1}{\sigma n^{-1/2} \|\hat{v}_l\|} \left( \frac{1}{\sqrt{n}} \hat{v}_l^\top \epsilon + \sqrt{n} R_l \right) + \frac{\sqrt{n}\beta_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| \\
&\leq \frac{\sigma}{\hat{\sigma}} \left( |Z_l| + \left| \frac{\sqrt{n} R_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| + \left| \frac{\sqrt{n}\beta_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| \right) \\
&\leq \frac{\sigma}{\hat{\sigma}} \left( |Z_l| + \frac{\|X_l\|}{\sqrt{n}} \left| \frac{\sqrt{n} R_l}{\sigma} \right| + \frac{\|X_l\|}{\sqrt{n}} \left| \frac{\sqrt{n}\beta_l}{\sigma} \right| \right),
\end{aligned}$$

and

$$\begin{aligned}
|T_l| &\geq \frac{\sigma}{\hat{\sigma}} \left( \left| \frac{\sqrt{n}\beta_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| - |Z_l| - \left| \frac{\sqrt{n} R_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| \right) \\
&\geq \frac{\sigma}{\hat{\sigma}} \left( \left| \frac{\sqrt{n}\beta_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| - |Z_l| - \frac{\|X_l\|}{\sqrt{n}} \left| \frac{\sqrt{n} R_l}{\sigma} \right| \right),
\end{aligned}$$

where  $R_l = n^{-1} \hat{v}_l^\top \mathbf{X}_{-l} (\beta_{-l} - \hat{\beta}_{-l})$ . We also observe that

$$\begin{aligned}
&\left[ \left\{ \left| \frac{\hat{\sigma}}{\sigma} - 1 \right| \leq \varepsilon \right\} \cap \left\{ \max_{l \in \mathcal{B}_j^{(2)}} |Z_l| + D' \max_{l \in \mathcal{B}_j^{(2)}} \left| \frac{\sqrt{n} R_l}{\sigma} \right| + D' \max_{l \in \mathcal{B}_j^{(2)}} \left| \frac{\sqrt{n}\beta_l}{\sigma} \right| \leq \sqrt{\tau \log p} \right\} \right. \\
&\quad \left. \cap \left\{ \max_l \frac{\|X_l\|}{\sqrt{n}} \leq D' \right\} \right] \subset \left\{ \max_{l \in \mathcal{B}_j^{(2)}} |T_l| \leq \sqrt{\tau \log p} \right\},
\end{aligned}$$

and

$$\begin{aligned}
&\left[ \left\{ \min_{l \in \mathcal{B}_j^{(1)}} \left| \frac{\sqrt{n}\beta_l}{\sigma n^{-1/2} \|\hat{v}_l\|} \right| - \max_{l \in \mathcal{B}_j^{(1)}} |Z_l| - D' \max_{l \in \mathcal{B}_j^{(1)}} \left| \frac{\sqrt{n} R_l}{\sigma} \right| > \sqrt{\tau \log p} \right\} \right. \\
&\quad \left. \cap \left\{ \left| \frac{\hat{\sigma}}{\sigma} - 1 \right| \leq \varepsilon \right\} \cap \left\{ \max_l \frac{\|X_l\|}{\sqrt{n}} \leq D' \right\} \right] \subset \left\{ \min_{l \in \mathcal{B}_j^{(1)}} |T_l| > \sqrt{\tau \log p} \right\}
\end{aligned}$$

and where  $D' = \sqrt{\max_l \Sigma_{l,l} + \varepsilon'}$  for  $0 < \varepsilon' \leq 2\kappa^2$ . Note that

$$\mathbb{P} \left\{ \max_l \frac{\|X_l\|}{\sqrt{n}} \leq D' \right\} \geq 1 - 2 \exp \left\{ -\frac{1}{8e^2} \frac{(\varepsilon')^2}{4\kappa^4} n + \log p \right\}. \quad (38)$$

We prove Proposition 1 in the following two steps.

1. Under Assumption 1, it suffices to show that

$$\mathbb{P} \left( \max_{l \in \mathcal{B}_j^{(2)}} |Z_l| + \frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(2)}} |\sqrt{n}R_l| \leq c_1 \sqrt{\log p} \right) \rightarrow 1,$$

where  $c_1 = \sqrt{\tau} - D' \sqrt{d_0}$ . We have, for  $\varepsilon'' > 0$ ,

$$\begin{aligned} & \mathbb{P} \left( \max_{l \in \mathcal{B}_j^{(2)}} |Z_l| + \frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(2)}} |\sqrt{n}R_l| \leq c_1 \sqrt{\log p} \right) \\ & \geq \mathbb{P} \left( \left\{ \max_{l \in \mathcal{B}_j^{(2)}} |Z_l| \leq c_1 \sqrt{\log p} - \varepsilon'' \right\} \cap \left\{ \frac{D'}{\sigma} \max_{k \in \mathcal{B}_j^{(2)}} |\sqrt{n}R_k| \leq \varepsilon'' \right\} \right) \\ & \geq \mathbb{P} \left( \max_{l \in \mathcal{B}_j^{(2)}} |Z_l| \leq c_1 \sqrt{\log p} - \varepsilon'' \right) + \mathbb{P} \left( \frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(2)}} |\sqrt{n}R_l| \leq \varepsilon'' \right) - 1 \\ & \geq \mathbb{P} \left( \frac{D'}{\sigma} \max_{k \in \mathcal{B}_j^{(2)}} |\sqrt{n}R_k| \leq \varepsilon'' \right) - 2p \exp \left\{ -\frac{\sigma^2 (c_1 \sqrt{\log p} - \varepsilon'')^2}{32e\kappa_\epsilon^2} \right\}. \end{aligned}$$

Here the last inequality follows by Lemma 2 under Assumption 2, i.e.,

$$\begin{aligned} & \mathbb{P} \left( \max_{l \in \mathcal{B}_j^{(2)}} |Z_l| \geq c_1 \sqrt{\log p} - \varepsilon'' \right) \\ & \leq \mathbb{P} \left( \bigcup_{l \in \mathcal{B}_j^{(2)}} \{ |Z_l| \geq c_1 \sqrt{\log p} - \varepsilon'' \} \right) \\ & \leq |\mathcal{B}_j^{(2)}| \times \mathbb{P} \left( \left| \frac{1}{\sigma \|\hat{v}_l\|} \sum_{i=1}^n \hat{v}_{li} \epsilon_i \right| \geq c_1 \sqrt{\log p} - \varepsilon'' \right) \\ & \leq |\mathcal{B}_j^{(2)}| \times 2 \times \exp \left\{ -\frac{\sigma^2 (c_1 \sqrt{\log p} - \varepsilon'')^2}{32e\kappa_\epsilon^2} \right\} \end{aligned} \quad (39)$$

conditional on  $\hat{v}_l$ . By the assumption

$$\frac{\sigma^2}{32e\kappa_\epsilon^2} (\sqrt{\tau} - \sqrt{d_0 \max_l \Sigma_{l,l}})^2 > 1,$$



we have

$$\frac{\sigma^2}{32e\kappa_\epsilon^2}c_1^2 > 1$$

for small enough  $\epsilon'$ . Together with (38), Lemma 5 and Assumption 4, we obtain

$$\mathbb{P}\left(\max_{l \in \mathcal{B}_j^{(2)}} |T_l| \leq \sqrt{\tau \log p}\right) \rightarrow 1.$$

2. We define  $c_2 = \sqrt{d_1/M''} - \sqrt{\tau}$ , where

$$M'' = \left(\min_l \Sigma_{l \setminus -l} + \epsilon'_1\right)^2 \left(\frac{2C_2}{8e^2} \left(\min_l \frac{1}{(\kappa_{0l})^2}\right)^{-1} + \max_l \Sigma_{l \setminus -l} + 2\epsilon'_2\right)$$

by letting  $(\epsilon'_0)^2 = 2((8e^2)^{-1} \min_l (\kappa_{0l})^{-2})^{-1} + \epsilon'_2$  in (36). We have, for  $\epsilon'' > 0$ ,

$$\begin{aligned} & \mathbb{P}\left(\min_{l \in \mathcal{B}_j^{(1)}} \left|\frac{\sqrt{n}\beta_l}{\sigma n^{-1/2}\|\hat{v}_l\|}\right| - \max_{l \in \mathcal{B}_j^{(1)}} |Z_l| - D' \max_{l \in \mathcal{B}_j^{(1)}} \left|\frac{\sqrt{n}R_l}{\sigma}\right| > \sqrt{\tau \log p}\right) \\ & \geq \mathbb{P}\left(\left\{\min_{l \in \mathcal{B}_j^{(1)}} \left|\frac{\sqrt{n}\beta_l}{\sigma n^{-1/2}\|\hat{v}_l\|}\right| - \max_{l \in \mathcal{B}_j^{(1)}} |Z_l| > \sqrt{\tau \log p} + \epsilon''\right\} \cap \left\{\frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(1)}} |\sqrt{n}R_l| \leq \epsilon''\right\}\right) \\ & \geq \mathbb{P}\left(\left\{\min_{l \in \mathcal{B}_j^{(1)}} \left|\frac{\sqrt{n}\beta_l}{\sigma\sqrt{M''}}\right| - \max_{l \in \mathcal{B}_j^{(1)}} |Z_l| > \sqrt{\tau \log p} + \epsilon''\right\} \cap \left\{\min_{l \in \mathcal{B}_j^{(1)}} \frac{1}{n^{-1/2}\|\hat{v}_l\|} \geq \frac{1}{\sqrt{M''}}\right\}\right) \\ & \quad + \mathbb{P}\left(\frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(1)}} |\sqrt{n}R_l| \leq \epsilon''\right) - 1 \\ & \geq \mathbb{P}\left(\min_{l \in \mathcal{B}_j^{(1)}} \left|\frac{\sqrt{n}\beta_l}{\sigma\sqrt{M''}}\right| - \max_{l \in \mathcal{B}_j^{(1)}} |Z_l| > \sqrt{\tau \log p} + \epsilon''\right) + \mathbb{P}\left(\max_{l \in \mathcal{B}_j^{(1)}} n^{-1}\|\hat{v}_l\|^2 \leq M''\right) \\ & \quad + \mathbb{P}\left(\frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(1)}} |\sqrt{n}R_l| \leq \epsilon''\right) - 2 \\ & \geq \mathbb{P}\left(\max_{l \in \mathcal{B}_j^{(1)}} |Z_l| < c_2\sqrt{\log p} - \epsilon''\right) + \mathbb{P}\left(\max_{l \in \mathcal{B}_j^{(1)}} n^{-1}\|\hat{v}_l\|^2 \leq M''\right) + \mathbb{P}\left(\frac{D'}{\sigma} \max_{l \in \mathcal{B}_j^{(1)}} |\sqrt{n}R_l| \leq \epsilon''\right) - 2 \\ & \geq 1 - 2|\mathcal{B}_j^{(1)}| \exp\left\{-\frac{\sigma^2}{32e\kappa_\epsilon^2}(b\sqrt{\log p} - \epsilon'')^2\right\} + \mathbb{P}\left(\max_{k \in \mathcal{B}_j^{(1)}} n^{-1}\|\check{v}_k\|^2 \leq M''\right) \\ & \quad + \mathbb{P}\left(\frac{D'}{\sigma} \max_{k \in \mathcal{B}_j^{(1)}} |\sqrt{n}R_k| \leq \epsilon''\right) - 2 \end{aligned}$$

where the last inequality follows from (39). By the assumption  $\sqrt{d_1/M} - \sqrt{\tau} > 0$ , we have

$c_2 = \sqrt{d_1/M''} - \sqrt{\tau} > 0$  for small enough  $\varepsilon'_1, \varepsilon'_2$ . Since  $|\mathcal{B}_j^{(1)}| \leq s_0 \ll p$ , by (37), (38), Lemma 5 and Assumption 4, we get  $\mathbb{P}\left(\min_{l \in \mathcal{B}_j^{(1)}} |T_l| > \sqrt{\tau \log p}\right) \rightarrow 1$ .

◇

**Proof of Theorem 1.** The argument below is conditional on the event  $\{\mathcal{A}_j^{(k)}(\tau) = \mathcal{B}_j^{(k)} \text{ for } k = 1, 2\}$  which occurs almost surely by Proposition 1. Let  $\check{u}_{j1} = \max_{k \in \mathcal{A}_j^{(1)}(\tau)} n^{-1} |\check{v}_j^\top X_k|$  and  $\check{u}_{j1} = \max_{k \in \mathcal{A}_j^{(2)}(\tau)} n^{-1} |\check{v}_j^\top X_k|$  where  $\check{v}_j$  is as in Corollary 4. Then,  $(\check{u}_{j1}, \check{u}_{j1}, \check{v}_j)$  is a feasible point to problem (8). By the definition of  $\tilde{v}_j$ ,

$$C_1 \frac{n}{\log p} \tilde{u}_{j1}^2 + C_2 \frac{n}{\log p} \tilde{u}_{j2}^2 + n^{-1} \|\tilde{v}_j\|^2 \leq C_1 \frac{n}{\log p} \check{u}_{j1}^2 + C_2 \frac{n}{\log p} \check{u}_{j2}^2 + n^{-1} \|\check{v}_j\|^2,$$

where  $\tilde{u}_{j1} = \max_{k \in \mathcal{A}_j^{(1)}} n^{-1} |\tilde{v}_j^\top X_k|$  and  $\tilde{u}_{j2} = \max_{k \in \mathcal{A}_j^{(2)}} n^{-1} |\tilde{v}_j^\top X_k|$ . Then, for  $i = 1, 2$ , we must have

$$\sqrt{C_i \frac{n}{\log p}} \tilde{u}_{ji} \leq \max\{C_1, C_2\} \varepsilon_{0j}^2 \left( \frac{1}{\Sigma_{j \setminus -j}} + \varepsilon_{1j} \right)^2 + \left( \frac{1}{\Sigma_{j \setminus -j}} + \varepsilon_{1j} \right)^2 (\Sigma_{j \setminus -j} + \varepsilon_{2j}),$$

with probability tending to 1 by Corollary 4. Then, by Assumptions 3 and 6,

$$|\sqrt{n}R(\tilde{v}_j, \beta_{-j})| = n^{-1/2} |\tilde{v}_j^\top \mathbf{X}_{-j}(\beta_{-j} - \hat{\beta}_{-j})| \leq n^{-1} \max_{k \neq j} |\tilde{v}_j^\top X_k| \sqrt{n} \|\hat{\beta}_{-j} - \beta_{-j}\|_1 = o_p(1).$$

Hence, we obtain

$$\sqrt{n}(\tilde{\beta}_j(\tilde{v}_j) - \beta_j) = \frac{1}{\sqrt{n}} \tilde{v}_j^\top \epsilon + o_p(1). \quad (40)$$

Note that

$$\frac{\sum_{i=1}^n E[(\tilde{v}_{j,i} \epsilon_i)^{2+\delta} | \tilde{v}_j]}{\sigma^{2+\delta} \|\tilde{v}_j\|^{2+\delta}} = \frac{E \epsilon_1^{2+\delta} \|\tilde{v}_j\|_{2+\delta}^{2+\delta}}{\sigma^{2+\delta} \|\tilde{v}_j\|^{2+\delta}} = o_{a.s.}(1).$$

Conditional on the event that  $\{\|\tilde{v}_j\|_{2+\delta}/\|\tilde{v}_j\| \rightarrow 0\}$ , the Lyapunov condition is satisfied and thus  $\tilde{v}_j^\top \epsilon / \{\sigma \|\tilde{v}_j\|\}$  converges to  $N(0, 1)$ . If  $\epsilon \sim N(0, \sigma^2 \mathbf{I})$ ,  $\tilde{v}_j^\top \epsilon / \{\sigma \|\tilde{v}_j\|\} \sim N(0, 1)$  conditional on  $\tilde{v}_j$ . The conclusion thus follows from (40) and Assumption 4 by the Slutsky's theorem. ◇

**Proof of Proposition 2.** All the arguments below are conditional on the event  $\{\mathcal{A}_j^{(2)} = \mathcal{B}_j^{(2)}\}$  which occurs almost surely by Proposition 1. With the projection direction  $\bar{v}_j$  from (16) and the

refitted least square estimator  $\check{\beta}$ , the bias (6) reduces to

$$\begin{aligned}
\sqrt{n}R(\bar{v}_j, \beta_{-j}) &= \frac{1}{\sqrt{n}} \sum_{k \neq j} \bar{v}_j^\top X_k (\beta_k - \check{\beta}_k) \\
&= \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{B}_j^{(1)}} \bar{v}_j^\top X_k (\beta_k - \check{\beta}_k) + \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{B}_j^{(2)}} \bar{v}_j^\top X_k (\beta_k - \check{\beta}_k) \\
&= \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{A}_j^{(1)}} \bar{v}_j^\top X_k (\beta_k - \check{\beta}_k) + \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{A}_j^{(2)}} \bar{v}_j^\top X_k (\beta_k - \check{\beta}_k) \\
&= \frac{1}{\sqrt{n}} \sum_{k \in \mathcal{A}_j^{(2)}} \bar{v}_j^\top X_k \beta_k,
\end{aligned}$$

where we have used the fact that  $\bar{v}_j^\top X_k = 0$  for  $k \in \mathcal{A}_j^{(1)}$  from (16) and  $\check{\beta}_{\mathcal{A}_j^{(2)}} = 0$  by (17). Thus, we have

$$\begin{aligned}
|\sqrt{n}R(\bar{v}_j, \beta_{-j})| &\leq \|n^{-1} \bar{v}_j^\top \mathbf{X}_{-j}\|_\infty \sqrt{n} \|\beta_{\mathcal{A}_j^{(2)}}\|_1 \\
&\leq \|n^{-1} \bar{v}_j^\top \mathbf{X}_{-j}\|_\infty \sigma \sqrt{d_0 \log p} \|\beta_{\mathcal{B}_j^{(2)}}\|_0 \\
&\leq O_p \left( \sqrt{\frac{\log p}{n}} \right) \sigma \sqrt{d_0 \log p} \|\beta_{\mathcal{B}_j^{(2)}}\|_0
\end{aligned}$$

where the second inequality holds by Assumption 1 under the event  $\{\mathcal{A}_j^{(2)} = \mathcal{B}_j^{(2)}\}$ . The last inequality follows from the fact that  $\|n^{-1} \bar{v}_j^\top \mathbf{X}_{-j}\|_\infty = O_p(\sqrt{\log p/n})$ , which can be verified by using similar arguments as in the proof of Corollary 1 together with the definition of  $\bar{v}_j$  under Assumption 5. The last statement follows immediately from condition (20).  $\diamond$

### 8.3.3 Technical details in Section 4

We first state the following results which are parallel to Lemma 3 and the first inequality in Corollary 5. As the proof is similar to the one in Lemma 3, we omit the details.

**Corollary 6.** *Let  $\theta_l = X_l - \mathbf{X}_{-S} b_l$  with  $b_l = \operatorname{argmin}_{\tilde{b}} E \|X_l - \mathbf{X}_{-S} \tilde{b}\|^2$  for  $l \in S$ . Under Assumption 5,*

$$\mathbb{P} \left( n^{-1} \|\theta_l^\top \mathbf{X}_{-S}\|_\infty \geq \xi_{0l} \sqrt{\frac{\log p}{sn}} \right) \leq 2 \exp \left\{ \left( 1 - c_{l,S} \frac{\delta_{0l}^2}{s(\xi_{0l})^2} \right) \log p \right\}$$

for  $0 < \xi_{0l} \leq \kappa_{0l} \sqrt{sn(\log p)^{-1}}$  where  $c_{l,S} > 0$  is an absolute constant and  $\kappa_{0l} = 2 \left( 1 + \sqrt{\Lambda_{\min}^{-1} \Sigma_{l,l}} \right) \kappa^2$ . As a consequence, we have

$$\mathbb{P} \left( \max_{l \in S} n^{-1} \|\theta_l^\top \mathbf{X}_{-S}\|_\infty \geq \xi'_0 \sqrt{\frac{\log p}{sn}} \right) \leq 2 \exp \left\{ - \left( \min_l \frac{c_{l,S}}{s(\kappa_{0l})^2} \right) (\xi'_0)^2 \log p + 2 \log p \right\}.$$

for  $0 < \xi'_0 \leq \min_l \kappa_{0l} \sqrt{sn(\log p)^{-1}}$ .

The following results are introduced for the proof of Theorem 2 which follows from a direct application of Proposition 2.1 in [27].

**Lemma 6.** For every  $\delta > 0$ , we have

$$\mathbb{P} \left( \|n^{-1} \mathbf{X}_S^\top \mathbf{X}_S - \boldsymbol{\Sigma}_{S,S}\| \leq \sqrt{\frac{4}{C_\kappa} \frac{s}{n} \log \frac{2}{\delta}} \right) \geq 1 - \delta,$$

where  $C_\kappa > 0$  is an absolute constant which only depends on  $\delta$  and  $\kappa$ .

We next introduce the following lemma which provides an upper bound for the operator norm of a matrix.

**Lemma 7.** Let  $\mathbf{B}$  be a  $m \times m$  matrix and  $\mathcal{N}_\varepsilon$  be an  $\varepsilon$ -net of the unit sphere  $\mathcal{S}^{m-1}$  for some  $\varepsilon \in (0, 1/2)$ . Then

$$\|\mathbf{B}\| \leq (1 - 2\varepsilon)^{-1} \sup_{c,d \in \mathcal{N}_\varepsilon} |c^\top \mathbf{B}d|.$$

**Proof of Lemma 7.** For any  $c, d \in \mathcal{S}^{m-1}$ , we can choose  $c_{\mathcal{N}}, d_{\mathcal{N}} \in \mathcal{N}_\varepsilon$  such that  $\max\{\|c - c_{\mathcal{N}}\|, \|d - d_{\mathcal{N}}\|\} \leq \varepsilon$ . Some algebra gives us

$$c^\top \mathbf{B}d = c_{\mathcal{N}}^\top \mathbf{B}d_{\mathcal{N}} + (c - c_{\mathcal{N}})^\top \mathbf{B}d + c_{\mathcal{N}}^\top \mathbf{B}(d - d_{\mathcal{N}}),$$

which implies that

$$|c^\top \mathbf{B}d| \leq 2\varepsilon \|\mathbf{B}\| + \sup_{c_{\mathcal{N}}, d_{\mathcal{N}} \in \mathcal{N}_\varepsilon} |c_{\mathcal{N}}^\top \mathbf{B}d_{\mathcal{N}}|.$$

Taking supremum over all  $c, d \in \mathcal{S}^{m-1}$  and rearranging terms give us the desired result.  $\diamond$

**Lemma 8.** For every  $\delta > 0$ , we have

$$\mathbb{P} \left( \|(n^{-1} \mathbf{X}_S^\top \mathbf{X}_{-S} - \boldsymbol{\Sigma}_{S,-S}) \boldsymbol{\Sigma}_{-S,-S}^{-1} \boldsymbol{\Sigma}_{-S,S}\| \leq 3 \sqrt{\frac{8}{C_{\kappa'}} \log \frac{2}{\delta} \frac{s}{n}} \right) \geq 1 - \delta$$

where  $C_{\kappa'} > 0$  denotes an absolute constant which only depends on  $\kappa' = 2\kappa^2 \sqrt{\Lambda_{\min}^{-1} D^2}$ .

**Proof of Lemma 8.** We prove the result in several steps. First, we bound the operator norm by using the so-called  $\varepsilon$ -net argument. Then we apply the concentration inequality for sub-exponential random variables and finally use the union bound to finish the proof. For two vectors  $a, b \in \mathbb{R}^{q \times 1}$ , write  $\langle a, b \rangle = a^\top b$ . By Lemma 7 and Lemma 5.2 in [27], we have

$$\begin{aligned} & \|(n^{-1} \mathbf{X}_S^\top \mathbf{X}_{-S} - \boldsymbol{\Sigma}_{S,-S}) \boldsymbol{\Sigma}_{-S,-S}^{-1} \boldsymbol{\Sigma}_{-S,S}\| \\ &= \sup_{c,d \in \mathcal{S}^{s-1}} \left| \frac{1}{n} \sum_{i=1}^n \left( \langle X_{i,S}, c \rangle \langle \boldsymbol{\Sigma}_{S,-S} \boldsymbol{\Sigma}_{-S,-S}^{-1} X_{i,-S}, d \rangle - c^\top \boldsymbol{\Sigma}_{S,-S} \boldsymbol{\Sigma}_{-S,-S}^{-1} \boldsymbol{\Sigma}_{-S,S} d \right) \right| \\ &\leq 3 \sup_{c,d \in \mathcal{N}_{1/3}} \left| \frac{1}{n} \sum_{i=1}^n \left( \langle X_{i,S}, c \rangle \langle \boldsymbol{\Sigma}_{S,-S} \boldsymbol{\Sigma}_{-S,-S}^{-1} X_{i,-S}, d \rangle - c^\top \boldsymbol{\Sigma}_{S,-S} \boldsymbol{\Sigma}_{-S,-S}^{-1} \boldsymbol{\Sigma}_{-S,S} d \right) \right| \end{aligned}$$

where  $\mathcal{N}_{1/3}$  denotes a  $1/3$ -net of  $\mathcal{S}^{s-1}$  with the covering number  $|\mathcal{N}_{1/3}| \leq 7^s$ .

Let us fix  $c, d \in \mathcal{N}_{1/3}$ . Because each row of  $\mathbf{X}_S$  and  $\mathbf{X}_{-S}$  is independent sub-gaussian random vector, we can apply the concentration inequality in Corollary 5.17 of [26]. Specifically, we have

$$\begin{aligned} & \mathbb{P} \left( \left| \frac{1}{n} \sum_{i=1}^n \left( \langle X_{i,S}, c \rangle \langle \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} X_{i,-S}, d \rangle - c^\top \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S} d \right) \right| \geq \epsilon \right) \\ & \leq 2 \exp \left( -cn \frac{\epsilon^2}{(\kappa')^2} \right) \end{aligned}$$

provided that  $\epsilon^2 \leq (\kappa')^2$ , where  $\|\langle X_{i,S}, c \rangle \langle \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} X_{i,-S}, d \rangle\|_{\psi_1} \leq \kappa'$  and  $c > 0$  is an absolute constant. Applying the union bound over  $c, d \in \mathcal{N}_{1/3}$ , we have

$$\sup_{c,d \in \mathcal{N}_{1/3}} \left| \frac{1}{n} \sum_{i=1}^n \left( \langle X_{i,S}, c \rangle \langle \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} X_{i,-S}, d \rangle - c^\top (\Sigma_{S,-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}) d \right) \right| \geq \epsilon$$

with probability at most  $2|\mathcal{N}_{1/3}|^2 \exp[-cn\epsilon^2/(\kappa')^2]$ , which implies

$$\begin{aligned} \mathbb{P} \left( \|(n^{-1} \mathbf{X}_S^\top \mathbf{X}_{-S} - \Sigma_{S,-S}) \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}\| < 3\epsilon \right) & \geq 1 - 2|\mathcal{N}_{1/3}|^2 \exp \left[ -cn \left( \frac{\epsilon}{\kappa'} \right)^2 \right] \\ & \geq 1 - 2 \exp [4s - n\epsilon^2 C_{\kappa'}] \end{aligned}$$

where  $C_{\kappa'} = c/(\kappa')^2$ . Then by letting  $\epsilon^2 = (8/C_{\kappa'}) \log(2/\delta)(s/n)$ , we have

$$\mathbb{P} \left( \|(n^{-1} \mathbf{X}_S^\top \mathbf{X}_{-S} - \Sigma_{S,-S}) \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}\| \leq 3 \sqrt{\frac{8}{C_{\kappa'}} \log \frac{2s}{\delta n}} \right) \geq 1 - \delta$$

which completes the proof.  $\diamond$

**Lemma 9.** *Let  $\hat{\mathbf{A}} = n^{-1} \mathbf{X}_S^\top \boldsymbol{\Theta}$  and  $\mathbf{A} = \Sigma_{S,S} - \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}$ . Under the assumption that  $s/n = o(1)$  and  $\|\mathbf{A}^{-1}\| \leq B$  for some constant  $B > 0$ , we have  $\|w\| = O_p(\|a_S\|)$ .*

**Proof of Lemma 9.** Note that

$$\|w\| = \|(n^{-1} \mathbf{X}_S^\top \boldsymbol{\Theta})^{-1} a_S\| \leq \|\hat{\mathbf{A}}^{-1}\| \|a_S\|.$$

We want to bound  $\|\hat{\mathbf{A}}^{-1}\|$ . Using the properties of operator norm, we have

$$\|\hat{\mathbf{A}}^{-1}\| \leq \|\hat{\mathbf{A}}^{-1} - \mathbf{A}^{-1}\| + \|\mathbf{A}^{-1}\| \leq \|\hat{\mathbf{A}}^{-1}\| \|\mathbf{A}^{-1}\| \|\mathbf{A} - \hat{\mathbf{A}}\| + \|\mathbf{A}^{-1}\|.$$

Rearranging the terms, we obtain

$$\|\hat{\mathbf{A}}^{-1}\| (1 - \|\mathbf{A}^{-1}\| \|\mathbf{A} - \hat{\mathbf{A}}\|) \leq \|\mathbf{A}^{-1}\|.$$

With the assumption  $\|\mathbf{A}^{-1}\| \leq B$ , we have

$$1 - \|\mathbf{A}^{-1}\| \|\mathbf{A} - \hat{\mathbf{A}}\| \geq 1 - B \|\mathbf{A} - \hat{\mathbf{A}}\|.$$

Under the assumption  $s/n = o(1)$ , by Lemmas 6 and 8, we have  $\|\mathbf{A} - \hat{\mathbf{A}}\| = o_p(1)$ . Thus  $1 - \|\mathbf{A}^{-1}\| \|\mathbf{A} - \hat{\mathbf{A}}\|$  is bounded from below by a positive constant with probability tending to one. Thus

$$\|\hat{\mathbf{A}}^{-1}\| \leq (1 - \|\mathbf{A}^{-1}\| \|\mathbf{A} - \hat{\mathbf{A}}\|)^{-1} \|\mathbf{A}^{-1}\| \leq (1 - B \|\mathbf{A} - \hat{\mathbf{A}}\|)^{-1} \|\mathbf{A}^{-1}\|$$

which implies that  $\|\hat{\mathbf{A}}^{-1}\| = O_p(1)$ . The conclusion follows directly.  $\diamond$

**Lemma 10.** *Let  $\Theta \in \mathbb{R}^{n \times s}$  where the  $l$ -th column vector is  $\theta_l$  for  $l \in S$  as in Corollary 6 and  $\check{v}_a = \Theta w$  where  $w = (n^{-1} \mathbf{X}_S^\top \Theta)^{-1} a_S$ . Then, under Assumption 5 and  $\|a_S\| = O(1)$ , we have  $n^{-1} \|\check{v}_a\|^2 = O_p(1)$ .*

*Proof of Lemma 10.* We note that

$$\|\Theta w\|^2 \leq \|\mathbf{X}_S - \mathbf{X}_{-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}\|^2 \|w\|_2^2 \leq 2 \underbrace{\{\|\mathbf{X}_S\|^2 + \|\mathbf{X}_{-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}\|^2\}}_I \|w\|_2^2.$$

We shall control  $I$  below. Lemma 5.3 in [27] gives us

$$\begin{aligned} \|\mathbf{X}_S\|^2 &\leq 4 \max_{c \in \mathcal{N}_{1/2}} c^\top \mathbf{X}'_S \mathbf{X}_S c, \\ \|\mathbf{X}_{-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}\|^2 &\leq 4 \max_{d \in \mathcal{N}_{1/2}} d^\top \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} \mathbf{X}'_{-S} \mathbf{X}_{-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S} d. \end{aligned}$$

Let  $\mathbf{Q} = \Sigma_{S,-S} \Sigma_{-S,-S}^{-1} (n^{-1} \mathbf{X}_{-S}^\top \mathbf{X}_{-S} - \Sigma_{-S,-S}) \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}$ . Since the elements of the terms inside the maximization can be expressed as a sum of independent sub-exponential random variables, we can use similar arguments as in the proof of Lemma 3 to show that for every  $\delta > 0$ ,

$$\mathbb{P} \left( \|\mathbf{Q}\| \leq \sqrt{\frac{4}{C_{\kappa'}} \frac{s}{n} \log \frac{2}{\delta}} \right) \leq 1 - \delta$$

where  $C_{\kappa'} > 0$  is an absolute constant which only depends on  $\kappa' = 2\kappa^2 \sqrt{\Lambda_{\min}^{-1} D^2}$ . Together with Lemma 6, we have

$$\begin{aligned} n^{-1} \{\|\mathbf{X}_S\|^2 + \|\mathbf{X}_{-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}\|^2\} &\leq C_0 \left\{ o_p(1) + \lambda_{\max}(\Sigma_{S,S}) + \lambda_{\max}(\Sigma_{S,-S} \Sigma_{-S,-S}^{-1} \Sigma_{-S,S}) \right\} \\ &\leq C_0 \{o_p(1) + 2\lambda_{\max}(\Sigma_{S,S})\}, \end{aligned}$$

for some constant  $C_0$ . Therefore, we have

$$n^{-1} \|\Theta w\|_2^2 \leq 2C_0 (o_p(1) + 2\Lambda_{\min}^{-2}) O_p(\|a_S\|^2) = O_p(1).$$

$\diamond$

**Proof of Theorem 2.** The arguments below are conditional on the sets  $\mathcal{A}_S^{(1)}$  and  $\mathcal{A}_S^{(2)}$  which have nonrandom limits by Proposition 1. Let  $\check{u}_{a1} = \max_{k \in \mathcal{A}_j^{(1)}} n^{-1} |\check{v}_a^\top X_k|$  and  $\check{u}_{a2} = \max_{k \in \mathcal{A}_j^{(2)}} n^{-1} |\check{v}_a^\top X_k|$ , where  $\check{v}_a$  is as in Lemma 10. Then,  $(\check{u}_{a1}, \check{u}_{a2}, \check{v}_a)$  is a feasible point to problem (22). By the definition of  $\tilde{v}_a$ ,

$$C_1 \frac{n}{\log p} \tilde{u}_{a1}^2 + C_2 \frac{n}{\log p} \tilde{u}_{a2}^2 + n^{-1} \|\tilde{v}_a\|^2 \leq C_1 \frac{n}{\log p} \check{u}_{a1}^2 + C_2 \frac{n}{\log p} \check{u}_{a2}^2 + n^{-1} \|\check{v}_a\|^2.$$

Then, for  $i = 1, 2$ , we must have

$$\begin{aligned} \sqrt{C_i \frac{n}{\log p}} \tilde{u}_{ai} &\leq \max\{C_1, C_2\} \frac{n}{\log p} \max_{k \notin S} n^{-1} |w^\top \Theta^\top X_k| + n^{-1} \|\check{v}_a\|^2 \\ &\leq \max\{C_1, C_2\} \|w\| (\xi'_0)^2 + M_a \end{aligned}$$

with probability tending to 1 for  $0 < \xi'_0 \leq \min_l \kappa_{0l} \sqrt{sn(\log p)^{-1}}$  and some constant  $M_a$  according to Corollary 6 and Lemma 10. Then, by Assumptions 3 and 6,

$$|\sqrt{n}R(\tilde{v}_a, \beta_{-S})| = n^{-1/2} |\tilde{v}_a^\top \mathbf{X}_{-S}(\beta_{-S} - \hat{\beta}_{-S})| \leq n^{-1} \max_{k \notin S} |\tilde{v}_a^\top X_k| \sqrt{n} \|\hat{\beta}_{-S} - \beta_{-S}\|_1 = o_p(1).$$

Hence, we obtain

$$\sqrt{n}(\tilde{\beta}_S(\tilde{v}_a) - a_S^\top \beta_S) = \frac{1}{\sqrt{n}} \tilde{v}_a^\top \epsilon + o_p(1). \quad (\text{S1})$$

Finally we can apply the central limit theorem as in the proof of Theorem 1, which completes the proof.  $\diamond$

## 8.4 Additional numerical results

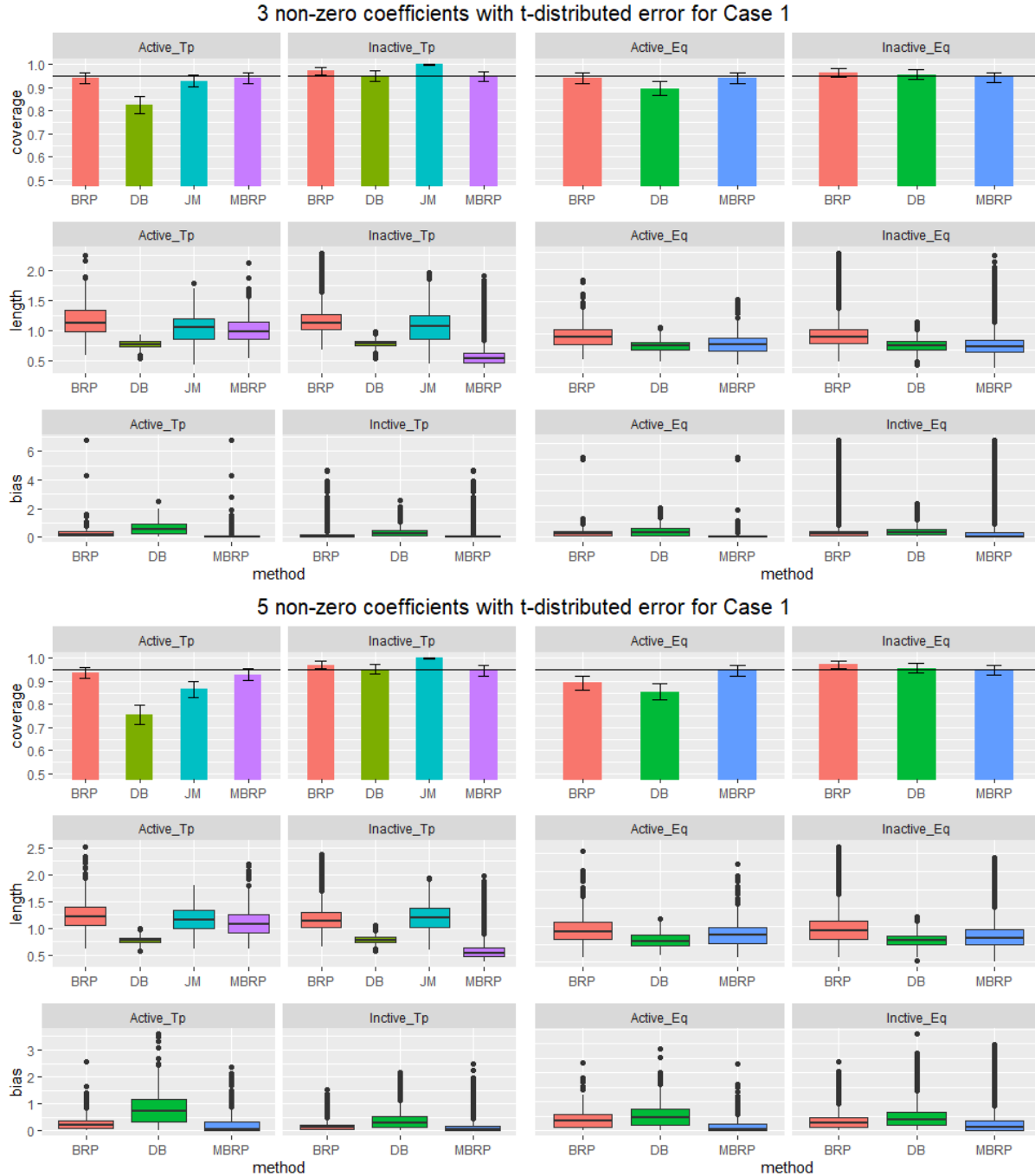


Figure 12: Simulation results for Case 1 with  $s_0 = 3, 5$  and  $t$ -distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.



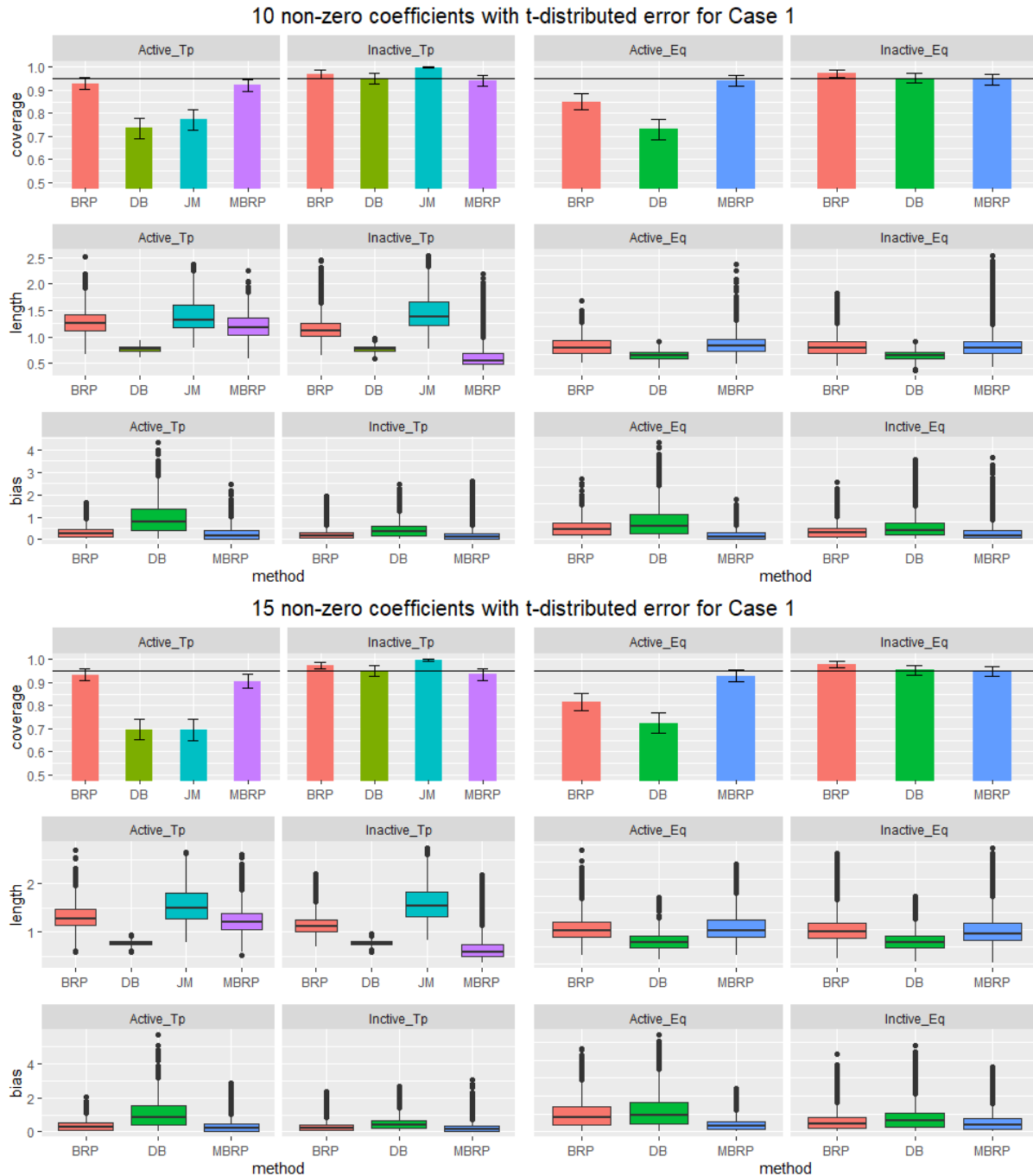


Figure 13: Simulation results for Case 1 with  $s_0 = 10, 15$  and  $t$ -distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

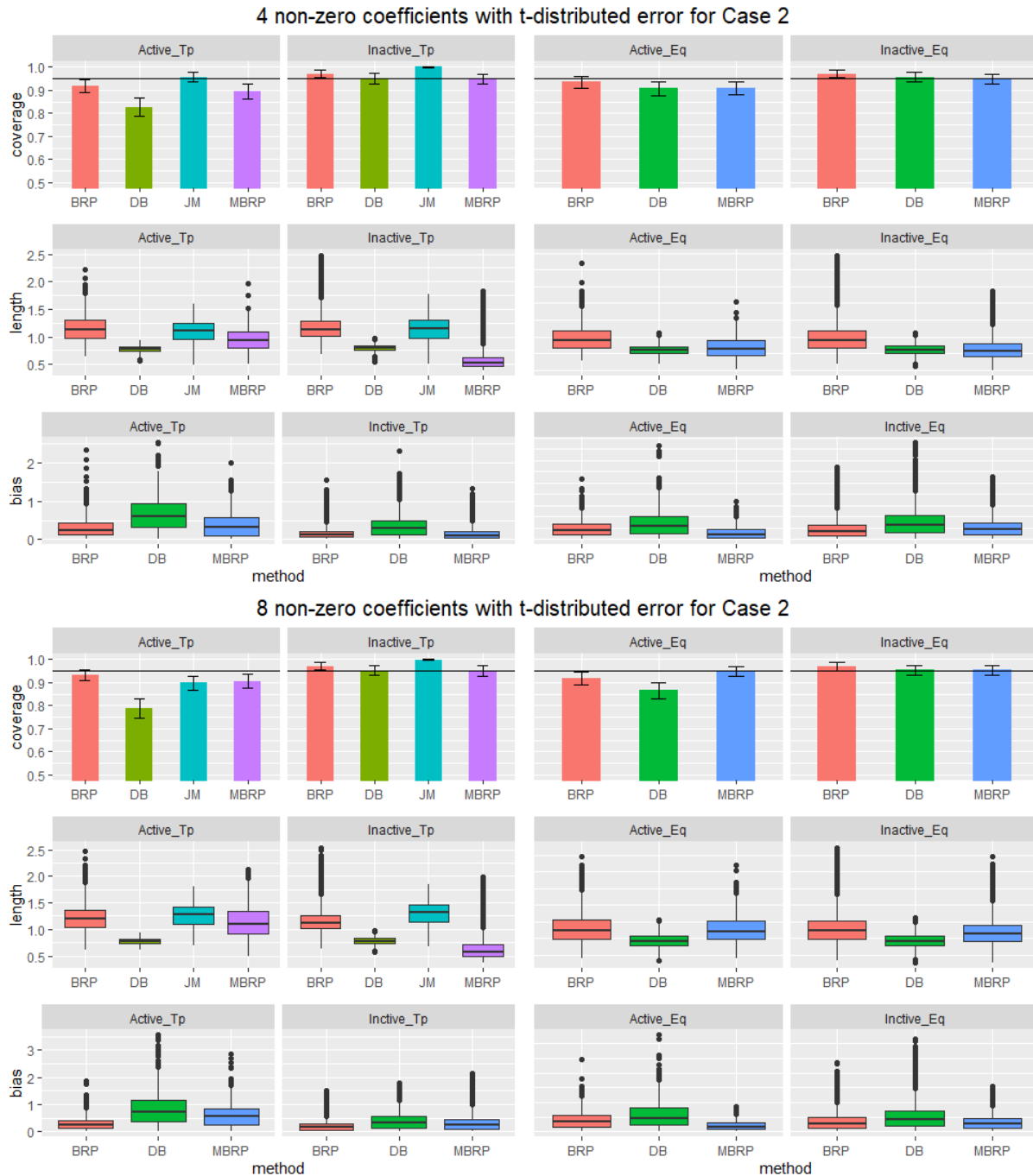


Figure 14: Simulation results for Case 1 with  $s_0 = 4, 8$  and  $t$ -distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

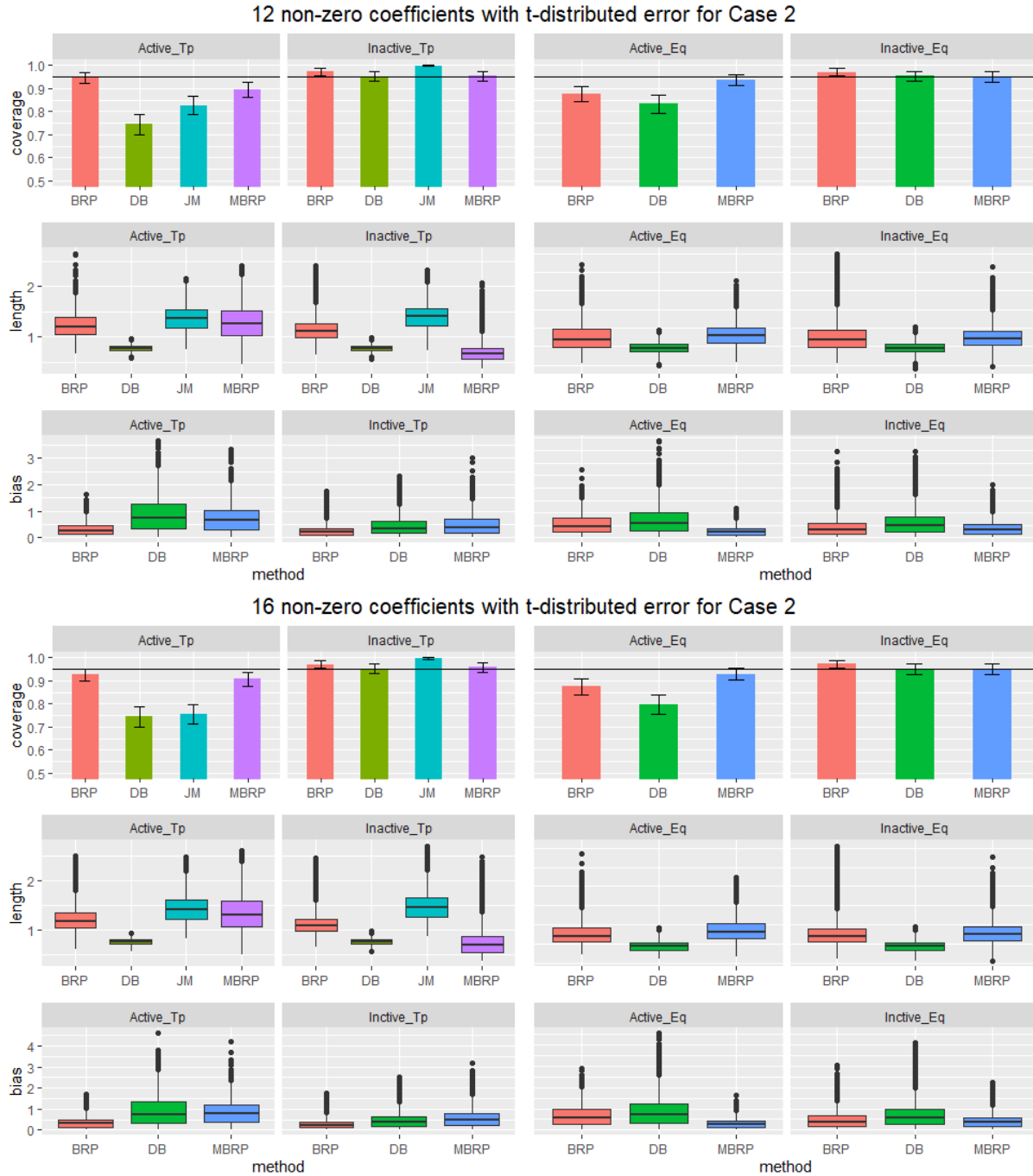


Figure 15: Simulation results for Case 1 with  $s_0 = 12, 16$  and  $t$ -distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

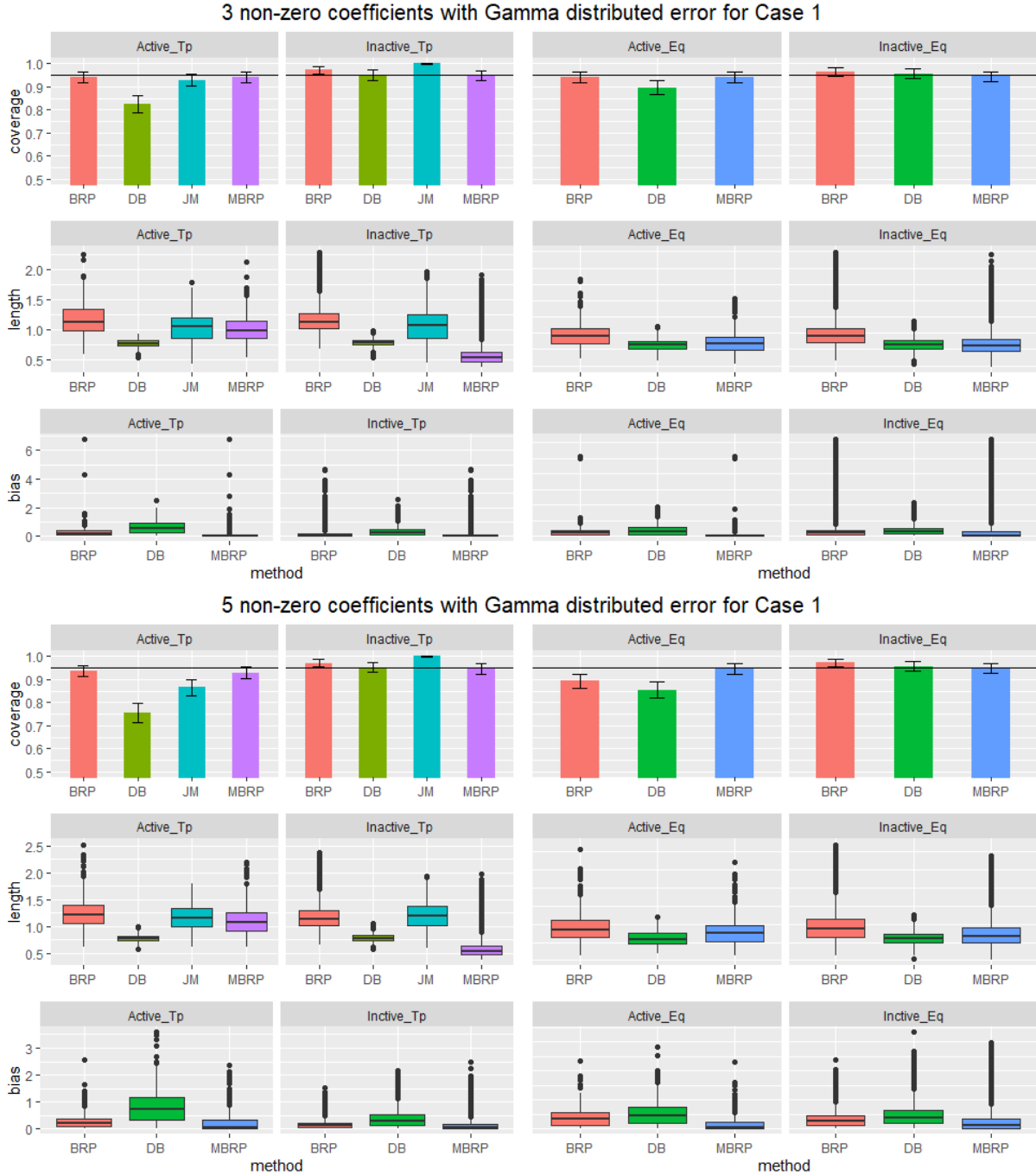


Figure 16: Simulation results for Case 1 with  $s_0 = 3, 5$  and Gamma-distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

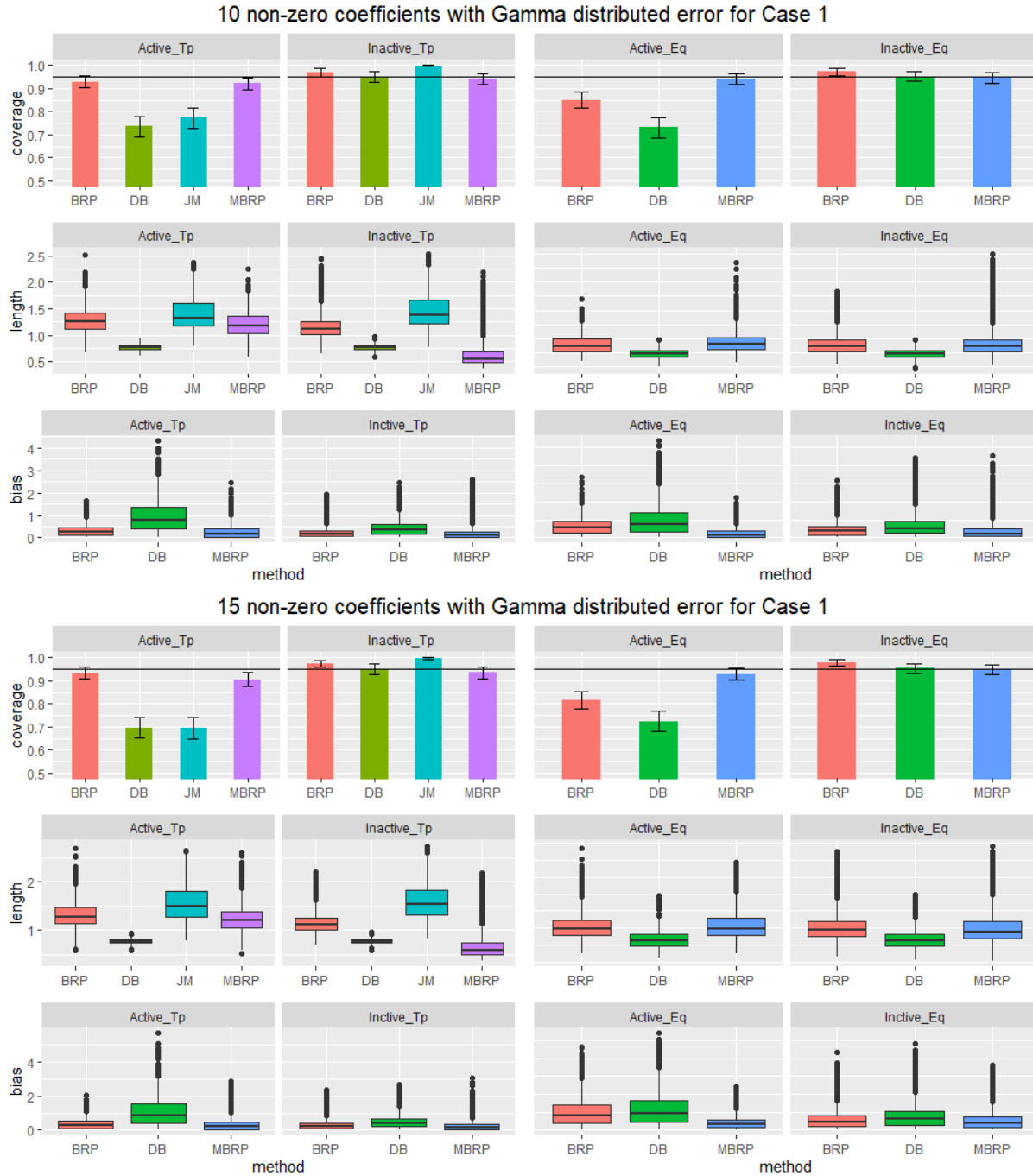


Figure 17: Simulation results for Case 1 with  $s_0 = 10, 15$  and Gamma-distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

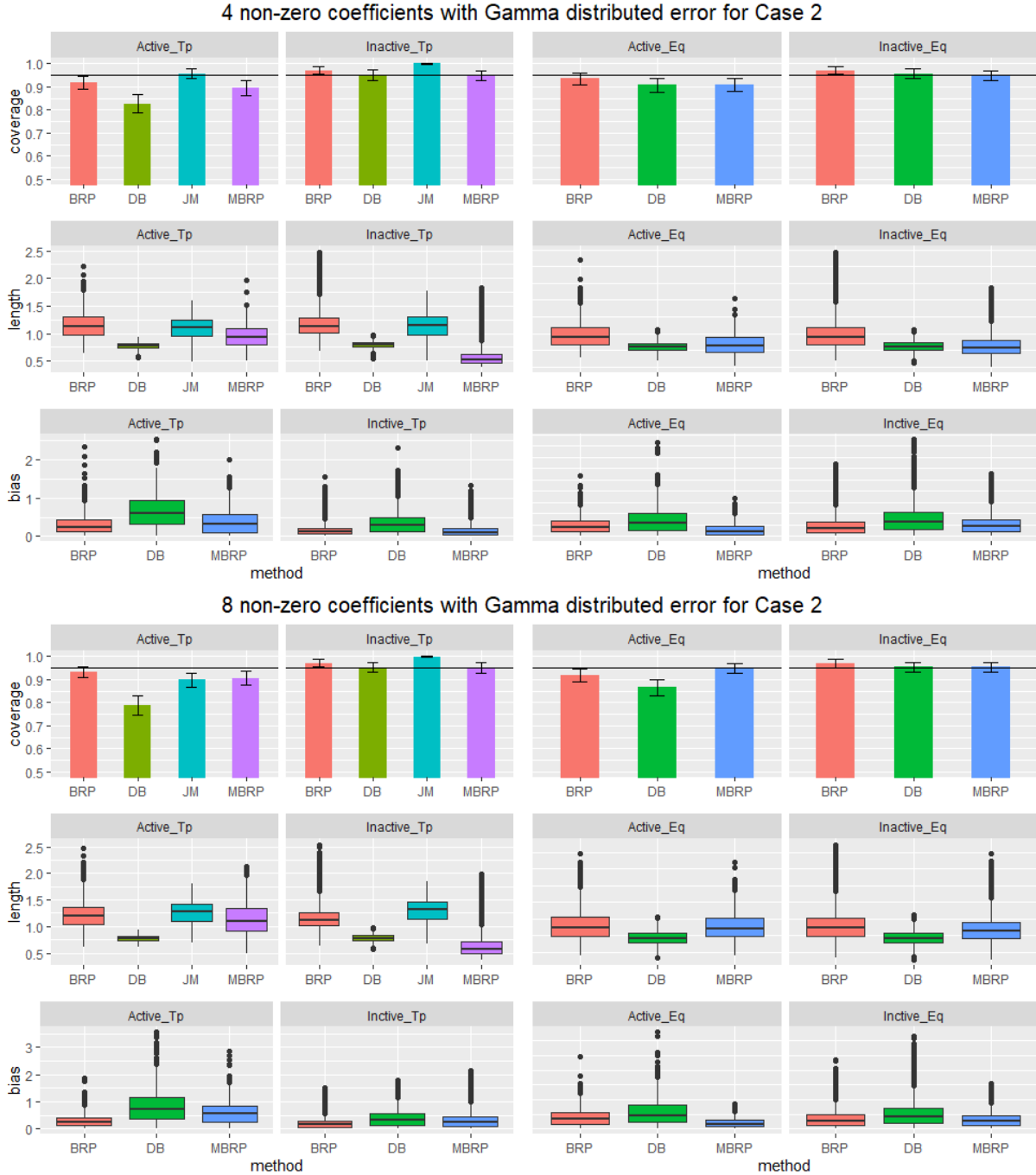


Figure 18: Simulation results for Case 2 with  $s_0 = 4, 8$  and Gamma-distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

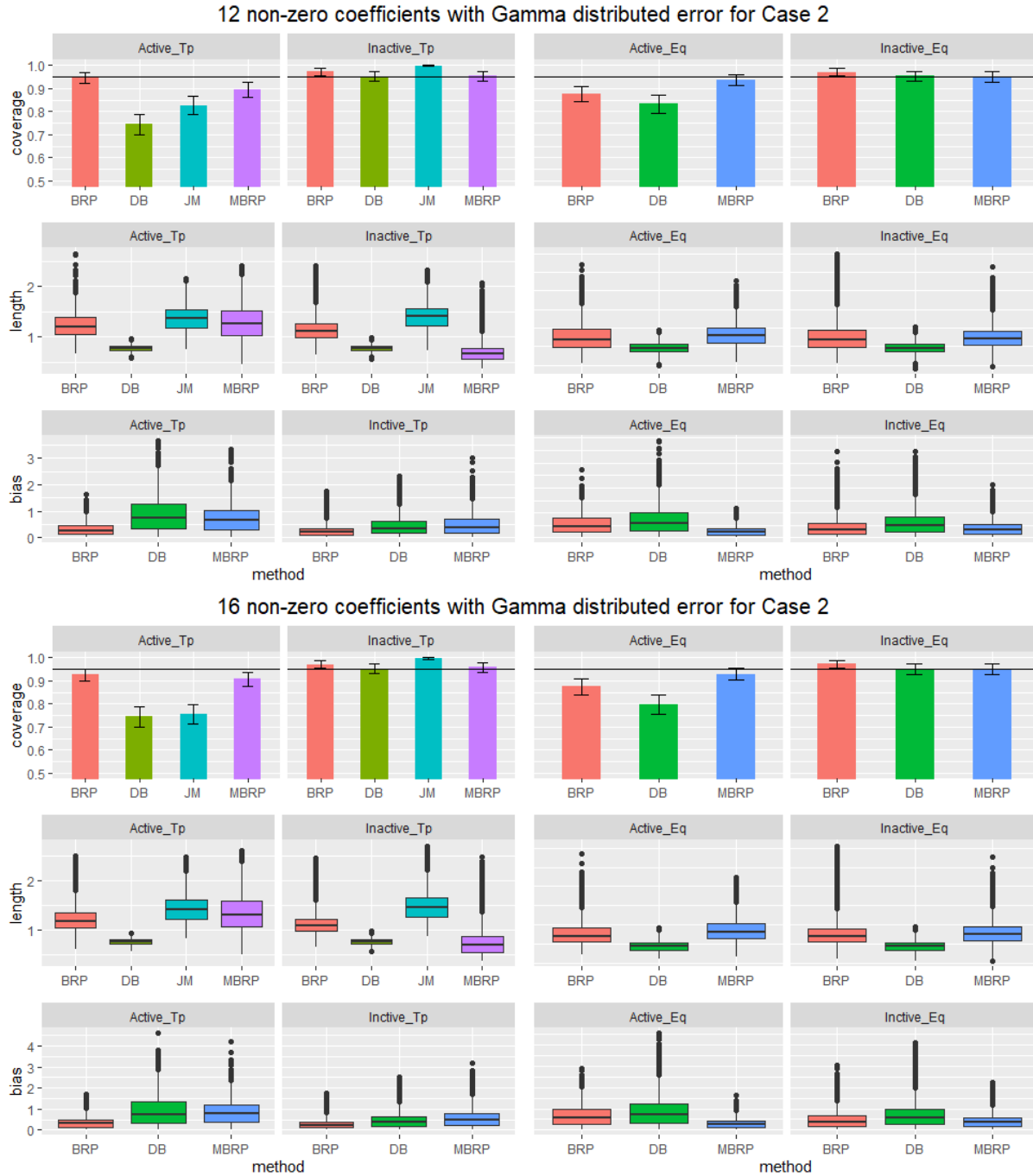


Figure 19: Simulation results for Case 2 with  $s_0 = 12, 16$  and Gamma-distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

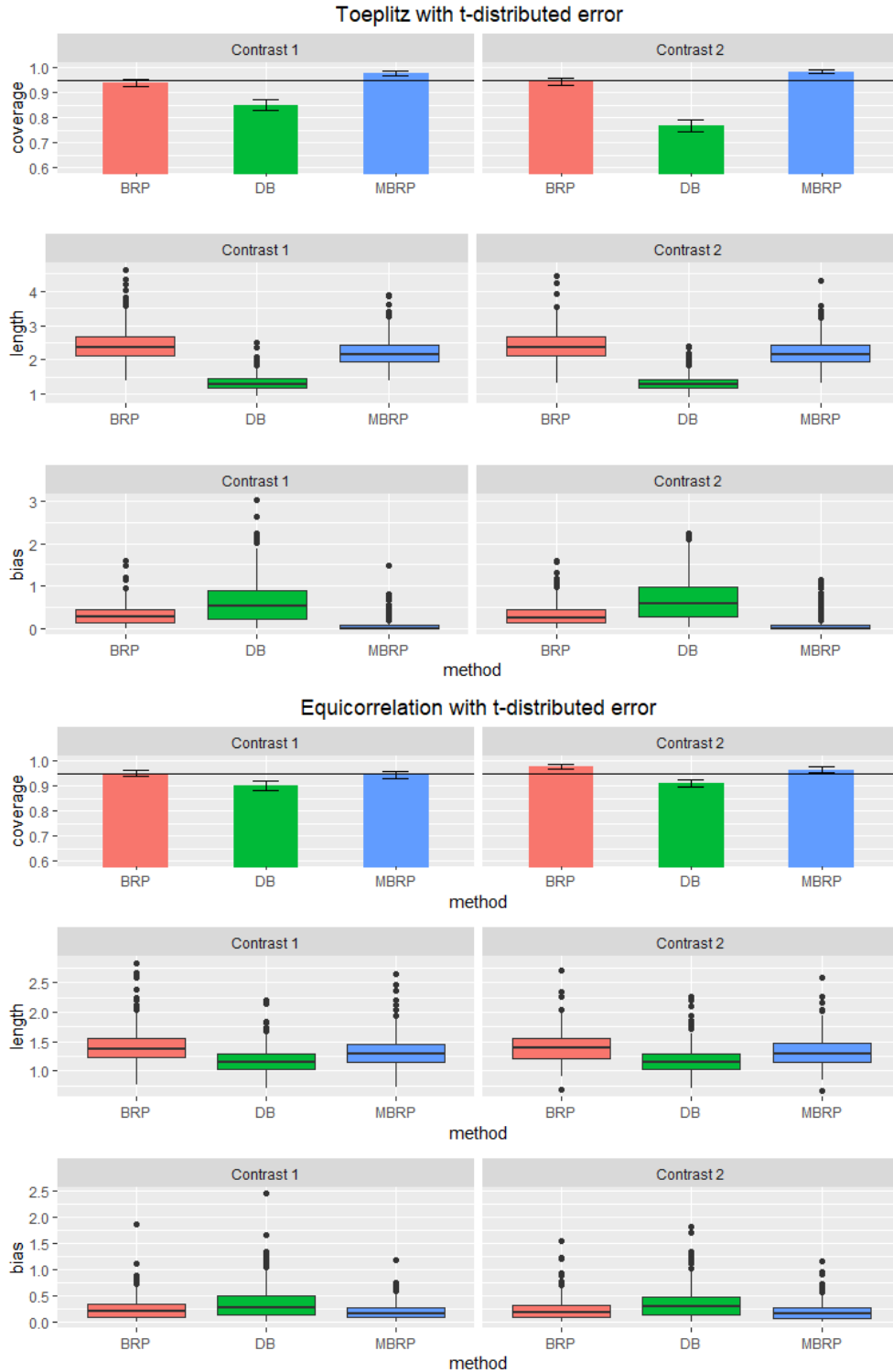


Figure 20: Simulation results for a sparse linear combination of  $\beta$  and  $t$ -distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals for each contrast. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.



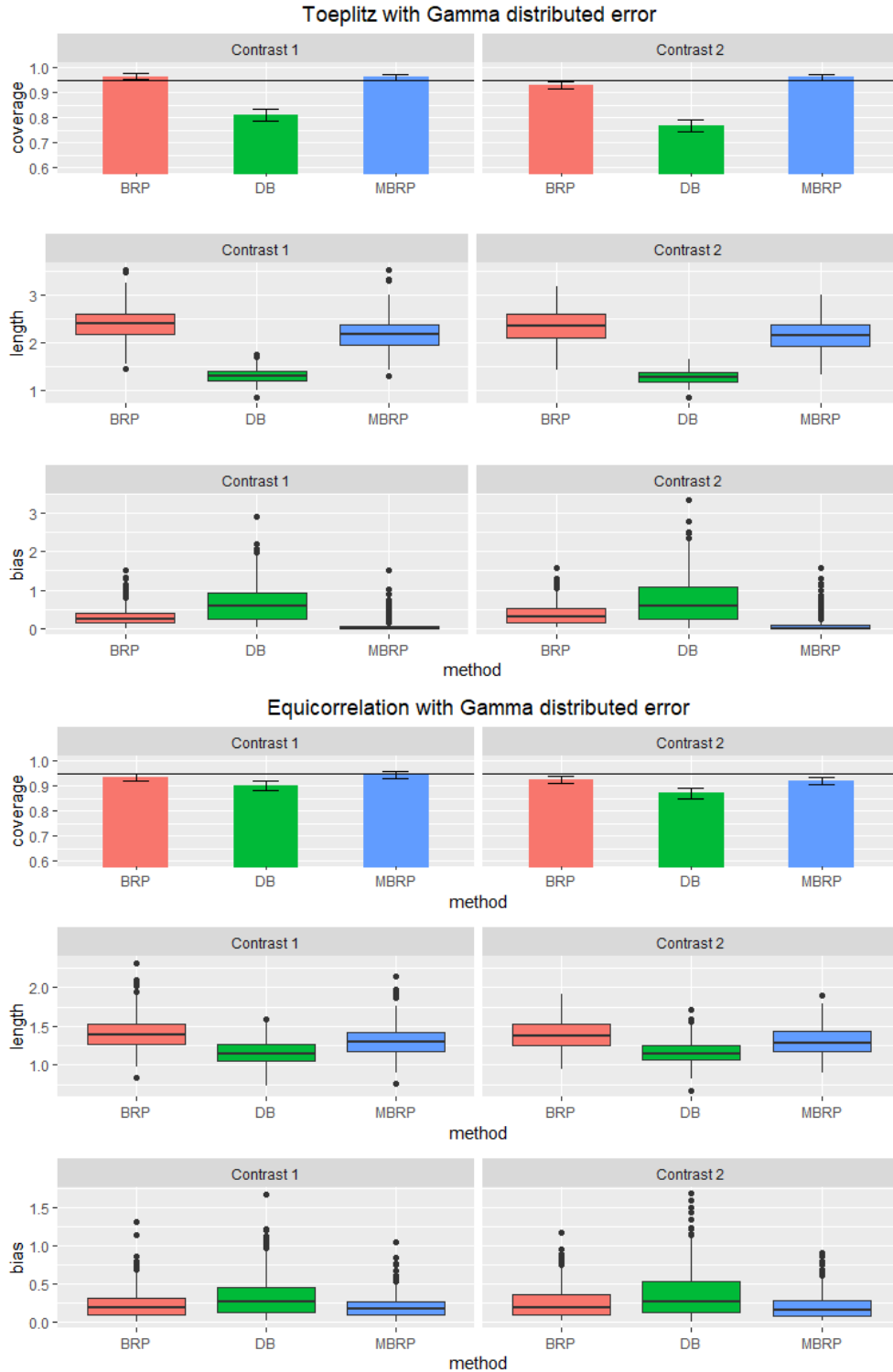


Figure 21: Simulation results for a sparse linear combination of  $\beta$  and Gamma-distributed random error. Barplots for the empirical coverage and boxplots for the length and bias of the 95% confidence intervals for each contrast. The horizontal line in the barplots indicates the nominal level. Error bars in the barplots represent the interval within one standard deviation of the empirical coverage.

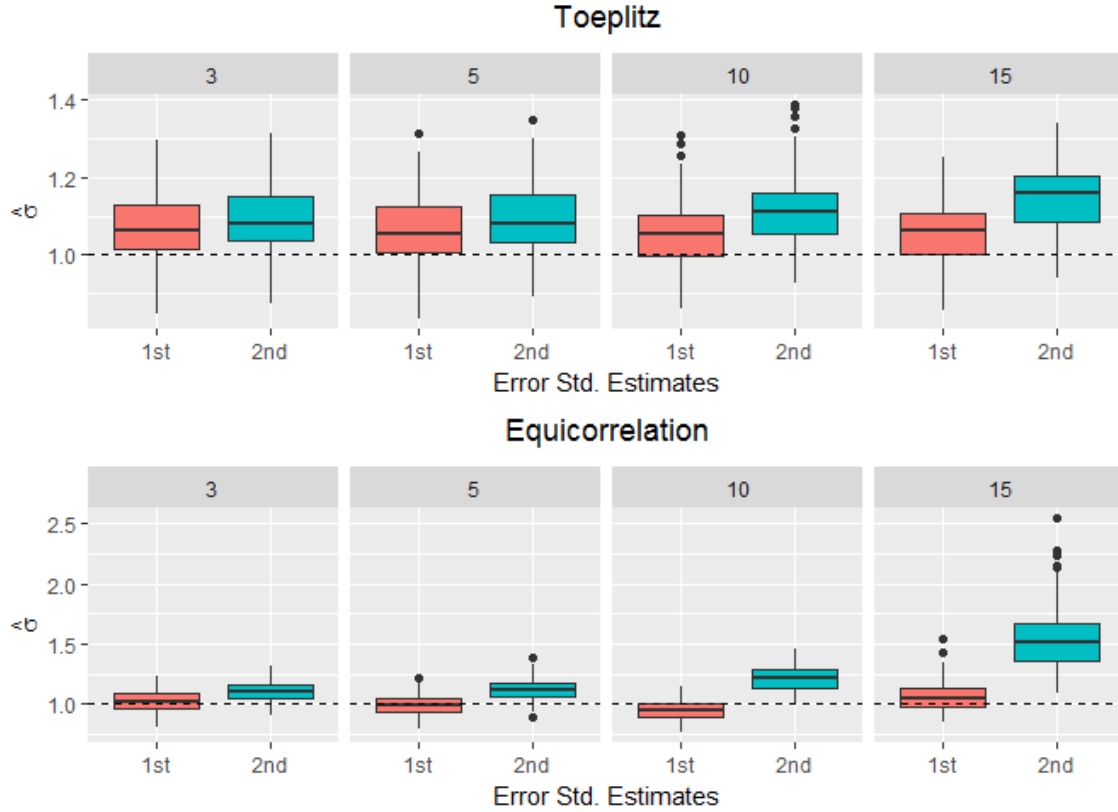


Figure 22: Boxplots of the two different error variance estimators. Data is generated by Case 1 with  $s_0 = 3, 5, 10$  and  $15$ . “1st” denotes the estimator  $\|Y - X\hat{\beta}\|^2/n$  and “2nd” denotes the estimator  $\|Y - X\hat{\beta}\|^2/(n - \|\hat{\beta}\|_0)$ . The number on the top of each panel denotes the number of non-zero coefficients. The horizontal dashed line corresponds to the true error variance.

## References

- [1] Barber, R. F. and Candès, E. J. (2015). Controlling the false discovery rate via knockoffs. *The Annals of Statistics*, 43(5), 2055-2085.
- [2] Barber, R. F. and Candès, E. J. (2019). A knockoff filter for high-dimensional selective inference. *The Annals of Statistics*, 47(5), 2504-2537.
- [3] Belloni, A., Chernozhukov, V. and Hansen, C. (2014). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2), 608-650.
- [4] Berk, R., Brown, L., Buja, A., Zhang, K. and Zhao, L. (2013). Valid post-selection inference. *The Annals of Statistics*, 41(2), 802-837.
- [5] Bühlmann, P. and van de Geer, S. (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.

- [6] Cai, T. T. and Guo, Z. (2017). Confidence intervals for high-dimensional linear regression: Minimax rates and adaptivity. *The Annals of Statistics*, 45(2), 615-646.
- [7] Candès, E., Fan, Y., Janson, L. and Lv, J. (2018). Panning for gold: ‘model-X’ knockoffs for high dimensional controlled variable selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 80(3), 551-577.
- [8] Chang, J., Chen, S. X., Tang, C. Y. and Wu, T. T. (2019). High-dimensional empirical likelihood inference. *arXiv preprint arXiv:1805.10742*.
- [9] Chatterjee, A. and Lahiri, S. N. (2011). Bootstrapping lasso estimators. *Journal of the American Statistical Association*, 106(494), 608-625.
- [10] Chatterjee, A. and Lahiri, S. N. (2013). Rates of convergence of the adaptive LASSO estimators to the oracle distribution and higher order refinements by the bootstrap. *The Annals of Statistics*, 41(3), 1232-1259.
- [11] Dezeure, R., Bühlmann, P. and Zhang, C.-H. (2017). High-dimensional simultaneous inference with the bootstrap. *Test*, 26(4), 685-719
- [12] Hastie, T., Tibshirani, R. and Wainwright, M. (2015). *Statistical learning with sparsity: the lasso and generalizations*. CRC press.
- [13] Javanmard, A. and Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research*, 15(1), 2869-2909.
- [14] Javanmard, A. and Montanari, A. (2018). Debiasing the lasso: Optimal sample size for Gaussian designs. *The Annals of Statistics*, 46(6A), 2593-2622.
- [15] Lee, J. D., Sun, D. L., Sun, Y. and Taylor, J. E. (2016). Exact post-selection inference, with application to the lasso. *The Annals of Statistics*, (44)(3), 907-927.
- [16] Liu, H. and Yu, B. (2013). Asymptotic properties of Lasso+mLS and Lasso+Ridge in sparse high-dimensional linear regression. *Electronic Journal of Statistics*, 7, 3124-3169.
- [17] Lockhart, R., Taylor, J., Tibshirani, R. and Tibshirani, R. (2014). A significance test for the lasso. *The Annals of Statistics*, 42(2), 413-468.
- [18] Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics*, 34(3), 1436-1462.
- [19] Meinshausen, N., Meier, L. and Bühlmann, P. (2009). P-values for high-dimensional regression. *Journal of the American Statistical Association*, 104(488), 1671-1681.
- [20] Meinshausen, N. and Bühlmann, P. (2010). Stability selection. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(4), 417-473.

- [21] Neykov, M., Ning, Y., Liu, J. S. and Liu, H. (2018). A unified theory of confidence regions and testing for high-dimensional estimating equations. *Statistical Science*, 33(3), 427-443.
- [22] Ning, Y. and Liu, H. (2017). A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *Annals of Statistics*, 45, 158-195.
- [23] Reid, S., Tibshirani, R. and Friedman, J. (2016). A study of error variance estimation in lasso regression. *Statistica Sinica*, 35-67.
- [24] Tibshirani, R. J., Taylor, J., Lockhart, R. and Tibshirani, R. (2016). Exact post-selection inference for sequential regression procedures. *Journal of the American Statistical Association*, 111(514), 600-620.
- [25] van de Geer, S., Bühlmann, P., Ritov, Y. and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high dimensional models. *The Annals of Statistics*, 42(3), 1166-1202.
- [26] Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
- [27] Vershynin, R. (2012). How close is the sample covariance matrix to the actual covariance matrix? *Journal of Theoretical Probability*, 25(3), 655-686.
- [28] Wasserman, L. and Roeder, K. (2009). High dimensional variable selection. *The Annals of Statistics*, 37(5A), 2178-2201.
- [29] Zhang, C. H. and Zhang, S. S. (2014). Confidence intervals for low-dimensional parameters with high-dimensional data. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1), 217-242.
- [30] Zhang, X. and Cheng, G. (2017). Simultaneous inference for high-dimensional linear models. *Journal of the American Statistical Association*, 112(518), 757-768.
- [31] Zhu, Y. and Bradic, J. (2018a). Significance testing in non-sparse high-dimensional linear models. *Electronic Journal of Statistics*, 12(2), 3312-3364.
- [32] Zhu, Y. and Bradic, J. (2018b). Linear hypothesis testing in dense high-dimensional linear models. *Journal of the American Statistical Association*, 113(524), 1583-1600.