

# Robust estimation in high-dimensional sparse heteroscedastic linear models

Duzhe Wang      Po-Ling Loh  
dwang282@wisc.edu      ploh@stat.wisc.edu

Department of Statistics  
University of Wisconsin-Madison  
Madison, WI 53706

July 20, 2020

## Abstract

We study the problem of robust estimation for the regression parameter in high-dimensional sparse heteroscedastic linear models. In our model, the variance of the random error in each observation depends on both predictor variables and regression coefficients. Our first goal is to provide theoretical guarantees of the Lasso under the assumed type of heteroscedasticity. We show in both theory and experiments that the Lasso is estimation consistent and variable selection consistent under mild sufficient conditions, even in the presence of heteroscedastic noise. The second goal of this paper is to robustify the Lasso estimates given that the parametric form of the variance function is known. We propose a new one-step estimator and derive its asymptotical properties for uncertainty quantification and valid statistical inference. Our approach creates an efficient paradigm to treat heteroscedasticity for high-dimensional sparse linear models.

## 1 Introduction

In standard formulations of linear regression analysis, it is often assumed that the random error is independent and identically distributed, and the variance does not depend on both predictor variables and regression parameters. However, this *homoscedasticity* assumption is not sensible for many real-world applications. For instance, it has been shown that *heteroscedasticity*, a phenomenon of unequal noise levels across different observations, is not uncommon in a variety of scientific areas ranging from genomics [11], medical imaging [14] to econometrics [9]. Moreover, heteroscedastic noise naturally exists in large and complex datasets, where the observations may be collected from multiple sources or be contaminated by outliers. While modern data collection technology brings new opportunities for regression applications, the accompanying occurrence of heteroscedasticity may cause serious problems for statistical analysis. It is a classical result, implied by the well-known Cramér-Rao lower bound, that the ordinary least squares estimator is inefficient in the presence of heteroscedastic noise [18].

**Problem setup.** In this paper, our primary interest is in robust estimation and inference about the regression coefficients under heteroscedasticity. More specifically, we consider the linear model

$$y = X\beta^* + \varepsilon, \tag{1}$$

where  $X = (X_1, \dots, X_N)^T \in \mathbb{R}^{N \times n}$  is the design matrix with  $X_i \in \mathbb{R}^n$ ,  $y = (y_1, \dots, y_N)^T \in \mathbb{R}^N$  is the observed response,  $\beta^* = (\beta_1^*, \dots, \beta_n^*)^T \in \mathbb{R}^n$  is the unknown true regression parameter, and  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_N)^T \in \mathbb{R}^N$  is a vector of unobserved random errors with mean zero and diagonal conditional variance matrix  $W \in \mathbb{R}^{N \times N}$ . When the principal diagonal of  $W$  contains same constant elements, (1) reduces to the homoscedastic linear model, which has been studied extensively in statistics since as early as the 1800s. Otherwise, the matrix  $W$  expresses the heteroscedasticity of noise. We suspect a heteroscedastic model because the dispersion of the residuals is significantly affected by the magnitude of the fitted values [6]. In general, larger magnitude of observed value is accompanied by larger variance in practice. Therefore, it is usually common to assume that  $W_i$ , the  $i$ th diagonal component of  $W$ , is a given function of  $X_i$  and  $\beta^*$ . For example, Bickel [1] considered the log-linear form of  $W_i$ —namely,  $W_i = c \exp(c' X_i^T \beta^*)$ . On the other hand, Jia et al. [14] focused on the Poisson-like heteroscedasticity, where  $W_i$  can be written as  $W_i = c |X_i^T \beta^*|$ . Inspired by above settings, we are interested in the conditional heteroscedastic model, where  $W_i$  takes the following general form:

$$W_i = g(X_i, \beta^*). \quad (2)$$

Throughout the paper, we assume that the parametric function form of  $g$  is known. We will introduce more assumptions on  $g$  in Section 2 and Section 3. Furthermore, driven by a large amount of modern applications in various fields, we will work within the high-dimensional framework, which allows the number of predictors  $n$  to grow and substantially exceed the sample size  $N$ .

The bulk of statistics literature on heteroscedasticity—notably a series of influential work by Carroll and his collaborators [5, 3, 4, 6, 7, 8, 10]—dates back to the last century, when most of the theoretical work was devoted to studying the classical low-dimensional scenario that the sample size  $N$  was able to diverge with the ambient dimension  $n$  fixed. In the direction of hypothesis testing, Bickel [1], and Carroll and Ruppert [4] suggested various robust tests for heteroscedasticity. Under similar setup of the variance function to ours, Jobson and Fuller [15], and Carroll and Ruppert [6] studied different classes of generalized weighted least squares estimates for the regression parameter, and established the corresponding asymptotic properties of their estimators. On the other hand, Carroll [3] considered the case that parametric form of the variance function is unknown. He proved that if the variance function is smooth to the design or the mean response, then there still exist efficient estimates for the regression parameter. Unfortunately, none of these traditional methods are readily applied to the high-dimensional settings that are of interest in our paper.

**Our contributions.** Despite the substantial body of work on heteroscedastic linear regression in low dimensions as discussed earlier, however, very little research has been conducted in the high-dimensional regime. In this paper, we are interested in the following two key questions:

- (a) *What are theoretical guarantees of the Lasso under heteroscedasticity?* Although heteroscedasticity is ubiquitous in real-world high-dimensional problems, practitioners usually ignore its existence and choose the standard methods in applied linear regression. For example, to estimate the sparse regression parameter in high-dimensional heteroscedastic linear models, we often use the canonical Lasso in practice, which is defined by

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^n}{\operatorname{argmin}} \quad \frac{1}{2N} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1. \quad (3)$$

Here,  $\lambda > 0$  is a tuning parameter. However, from a theoretical perspective, the performance of the Lasso is only well understood under assumptions of the standard sparse linear model

with homoscedastic noise. Therefore, as an initial step to study heteroscedasticity in high dimensions, we are interested in performance of the Lasso for the conditional heteroscedastic model defined in (1).

- (b) *How can we robustify the Lasso estimates for the regression parameter when some side information on the variance function is available?* In particular, we focus on the scenario that the parametric form of the variance function  $g$  is known. In such a setting, we are interested in designing estimators with small variances and quantifying their uncertainties.

The current paper is well motivated by above questions. We now summarize the main contributions of our paper as follows:

- (a) Our first major contribution is to provide a general set of mild sufficient conditions under which the Lasso is statistically consistent, even in the presence of heteroscedastic noise. More specifically, in Section 2, we consider the sub-Gaussian random design, and derive upper bounds for the  $\ell_2$ -estimation error and the  $\ell_1$ -estimation error of the Lasso under the conditional heteroscedastic model (1). In Theorem 1, we show that for an appropriate choice of the tuning parameter  $\lambda$ , the Lasso is  $\ell_2$ -consistent if  $\frac{sL_2 \log n}{N} = o(1)$ , where  $s$  is the number of nonzero elements of  $\beta^*$  and  $L_2$  is the largest conditional error variance. Furthermore, in Theorem 2, we prove that the Lasso is variable selection consistent under similar mutual incoherence condition and beta-min condition as in the standard case.
- (b) Our second major contribution is to study robustification of the Lasso under heteroscedasticity. In Section 3, we develop a new type of one-step estimator for the regression parameter and derive its limiting distribution in Theorem 3. Asymptotical properties of the one-step estimator not only helps us to quantify its uncertainty, but also provides us theoretical foundations to make valid confidence intervals and statistical tests for the regression parameter.

**Related work.** Next, we highlight several recent papers on related topics. Jia et al. [14] examined performance of the Lasso under Poisson-like heteroscedasticity. They provided necessary and sufficient conditions for the sign consistency of the Lasso under a sparse Poisson-like model in the high-dimension regime. But the authors did not develop new estimation approaches to utilize the useful information given by Poisson-like heteroscedasticity. Another related work is that of Wagener and Dette [21]. In their paper, the authors proposed a weighted adaptive Lasso to estimate the regression parameter. They showed that the weighted adaptive Lasso performed consistent model selection and was asymptotically normal. Wagener and Dette [21] dealt with the deterministic design, while our paper focuses on the random design. Compared with theirs, the randomness in the design matrix requires a more involved analysis.

Finally, we discuss another two work by Daye et al. [11] and Sharpnack and Kolar [19]. Slightly different with our setup, both these two papers considered a log-linear heteroscedastic model with variances parameterized by an independent set of parameters rather than the regression vector. Daye et al. [11] focused exclusively on estimating the regression parameter. The authors discussed a likelihood-based approach, which optimized the regularized negative log-likelihood with  $\ell_1$ -penalties on both the regression and variance parameters simultaneously. They devised an efficient coordinate descent algorithm to minimize the non-convex objective function. On the other hand, the goal of Sharpnack and Kolar [19] was to estimate both the regression and variance parameters. They proposed a heteroscedastic iterative penalized pseudolikelihood optimizer, which employed non-convex penalties.

**Notation.** For functions  $f(n)$  and  $g(n)$ , we write  $f(n) \lesssim g(n)$  to mean that  $f(n) \leq cg(n)$  for some universal constant  $c \in (0, \infty)$ . Similarly, we write  $f(n) \gtrsim g(n)$  when  $f(n) \geq c'g(n)$  for some universal constant  $c' \in (0, \infty)$ , and write  $f(n) \asymp g(n)$  to mean that  $f(n) \lesssim g(n)$  and  $f(n) \gtrsim g(n)$  hold simultaneously. We write  $f(n) = o(g(n))$  to mean that  $\lim_{n \rightarrow \infty} f(n)/g(n) = 0$ . We write  $X_n = o_{\mathbb{P}}(a_n)$  to mean  $X_n/a_n$  converges to 0 in probability, and write  $X_n = \mathcal{O}_{\mathbb{P}}(a_n)$  to mean  $X_n/a_n$  is stochastically bounded. We use  $I_n$  to denote the identity matrix of size  $n \times n$ . For a matrix  $M \in \mathbb{R}^{m \times n}$ , we write  $\|M\|_{op}$  to denote the  $\ell_2$  operator norm,  $\|M\|_1 = \max_{j=1, \dots, n} \sum_{i=1}^m |M_{ij}|$  to denote the  $\ell_1$  operator norm,  $\|M\|_{\infty} = \max_{i=1, \dots, m} \sum_{j=1}^n |M_{ij}|$  to denote the  $\ell_{\infty}$  operator norm,  $\|M\|_F = \left( \sum_j \sum_i M_{ij}^2 \right)^{1/2}$  to denote the Frobenius norm, and write  $\|M\|_{1,1} = \sum_{i=1}^m \sum_{j=1}^n |M_{ij}|$  to denote the element-wise  $\ell_1$  norm. We write  $\|\cdot\|_{\infty}$  to denote the element-wise infinity norm for both vectors and matrices. We write  $\|\cdot\|_0$  to denote the number of non-zero elements in a vector. For squared matrices, we write  $\lambda_{\min}$  and  $\lambda_{\max}$  to denote the minimum and maximum eigenvalues respectively. For  $q, r > 0$ , we write  $\mathbb{B}_q(r)$  to denote the  $\ell_q$ -ball of radius  $r$  centered around 0. We use  $c, c', c'', C',$  and  $C''$  to denote positive constants, where we may use the same notation to refer to different constants as we move between results. For a constant  $\sigma > 0$ , a random variable  $X \in \mathbb{R}$  is said to be  $\sigma$ -sub-Gaussian if  $\mathbb{E}(X) = 0$  and its moment generating function satisfies  $\mathbb{E}[\exp(tX)] \leq \exp(\sigma^2 t^2/2)$  for all  $t \in \mathbb{R}$ . Similarly, a random vector  $X \in \mathbb{R}^n$  is said to be  $\sigma$ -sub-Gaussian if  $\mathbb{E}(X) = 0$  and  $u^T X$  is  $\sigma$ -sub-Gaussian for any unit vector  $u \in \mathbb{S}^{n-1}$ . Furthermore, if a random matrix  $X \in \mathbb{R}^{N \times n}$  is formed by drawing each row  $X_i \in \mathbb{R}^n$  in an i.i.d. manner from a  $\sigma$ -sub-Gaussian distribution with covariance matrix  $\Sigma$ , then we say  $X$  is a row-wise  $(\sigma, \Sigma)$ -sub-Gaussian random matrix.

## 2 The Lasso under heteroscedasticity

We start with studying theoretical properties of the Lasso for the conditional heteroscedastic linear model (1). In this section, we make use of the following assumptions:

**Assumption 1.** The true regression parameter  $\beta^*$  is  $s$ -sparse. That is,  $\|\beta^*\|_0 \leq s$ , where  $s < N \ll n$ . In what follows, let  $S$  denote the support set of  $\beta^*$ . Furthermore, the triple  $(s, N, n)$  is allowed to increase to  $\infty$ .

**Assumption 2.** (a) The design matrix  $X$  is a row-wise  $(\sigma_x, \Sigma_x)$ -sub-Gaussian random matrix. Furthermore, eigenvalues of  $\Sigma_x$  are bounded by dimension-free constants. That is,

$$c_{\min} \leq \lambda_{\min}(\Sigma_x) \leq \lambda_{\max}(\Sigma_x) \leq c_{\max},$$

where  $c_{\min} > 0$  and  $c_{\max} > 0$  are constants.

(b) The random errors  $\varepsilon_i$ ,  $i = 1, \dots, N$ , are independent, conditionally normal random variables with

$$\mathbb{E}(\varepsilon_i | X_i) = 0, \quad \mathbb{E}(\varepsilon_i^2 | X_i) = g(X_i, \beta^*),$$

where  $0 < L_1 \leq g(X_i, \beta^*) \leq L_2 < \infty$ .

**Assumption 3.** For covariance matrix  $\Sigma_x$ , there exists a constant  $\alpha \in (0, 1)$ , such that

$$\|(\Sigma_x)_{S^c S} (\Sigma_x)_{SS}^{-1}\|_{\infty} \leq \frac{\alpha}{2}.$$

Assumption 1 and Assumption 2 contain a standard set of conditions for elements of  $\beta^*$ ,  $X$ , and  $\varepsilon$ . Note that in Part (b) of Assumption 2, we assume that the variance function is lower and upper bounded, which is often imposed when analyzing heteroscedasticity [15, 21]. Assumption 3, namely the mutual incoherence condition, is usually assumed for variable selection of the Lasso in the standard high-dimensional sparse linear models [22]. With above assumptions at hand, we have the first main result, which concerns the estimation error of the Lasso under heteroscedasticity.

**Theorem 1.** Assume that Assumption 1 and Assumption 2 hold. If  $N \gtrsim s \log n$  and

$$\lambda \geq 4(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} \sqrt{\frac{L_2 \log n}{N}},$$

then with probability at least  $1 - c \exp(-c' \log n)$ , we have

$$\|\hat{\beta} - \beta^*\|_2 \leq 6\lambda_{\min}^{-1}(\Sigma_x) \lambda \sqrt{s}, \quad \|\hat{\beta} - \beta^*\|_1 \leq 24\lambda_{\min}^{-1}(\Sigma_x) \lambda s.$$

The proof of Theorem 1 is contained in Section 6.1.

**Remark 1.** We provide several comments on the upper bounds derived in Theorem 1. First, if  $\lambda \asymp \sqrt{\frac{L_2 \log n}{N}}$ , then we have

$$\|\hat{\beta} - \beta^*\|_2 \leq \mathcal{O}_{\mathbb{P}} \left( \sqrt{\frac{s L_2 \log n}{N}} \right) \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_1 \leq \mathcal{O}_{\mathbb{P}} \left( s \sqrt{\frac{L_2 \log n}{N}} \right). \quad (4)$$

Secondly, (4) implies that the Lasso is  $\ell_2$ -consistent if  $\frac{s L_2 \log n}{N} = o(1)$ . Lastly, if we consider the homoscedastic case—namely,  $W_i = \sigma_\varepsilon^2$  for  $i = 1, \dots, N$ , then inequalities in (4) reduce to

$$\|\hat{\beta} - \beta^*\|_2 \leq \mathcal{O}_{\mathbb{P}} \left( \sigma_\varepsilon \sqrt{\frac{s \log n}{N}} \right) \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_1 \leq \mathcal{O}_{\mathbb{P}} \left( \sigma_\varepsilon s \sqrt{\frac{\log n}{N}} \right),$$

which are standard results of the Lasso in statistics literature [17].

Next, we shift our attention to variable selection performance of the Lasso under heteroscedasticity. We say that the Lasso is variable selection consistent if there exists an appropriate choice of  $\lambda$ , such that  $\mathbb{P}(S = \hat{S}) \rightarrow 1$  as the triple  $(s, N, n)$  goes to  $\infty$ , where  $\hat{S}$  is the support set of  $\hat{\beta}$ . The following theorem, proved in Section 6.2, claims that the Lasso performs consistent variable selection, even in the presence of heteroscedasticity. The proof of Theorem 2 is based on the primal-dual witness (PDW) type arguments. See Appendix B for a detailed introduction to the PDW construction.

**Theorem 2.** Assume that Assumption 1, Assumption 2, and Assumption 3 hold. Let

$$\lambda \geq \frac{4}{1 - \alpha} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} \sqrt{\frac{L_2 \log(n - s)}{N}}.$$

(a) Under the sample size scaling

$$N^2 \gtrsim \frac{s^5 \log[s(n - s)]}{\alpha^2}, \quad N \gtrsim \frac{s^4 \log[s(n - s)]}{\alpha^2}, \quad \text{and} \quad N \gtrsim \frac{s^5}{\alpha^2},$$

with probability at least  $1 - c \exp(-c' \log s)$ , we have  $\hat{S} \subseteq S$ .

(b) In addition, if

$$\min_{i \in S} |\beta_i^*| > \sqrt{\frac{8}{\lambda_{\min}(\Sigma_x)}} \sqrt{\frac{L_2 \log s}{N}} + \lambda \left( \|(\Sigma_x)_{SS}^{-1}\|_{\infty} + cs \sqrt{\frac{s}{N}} \right),$$

then with probability at least  $1 - c' \exp(-c'' \log s)$ , we have  $S = \hat{S}$ .

Finally, results in Theorem 1 and Theorem 2 together lead to the following corollary, which is useful for the statistical analysis in the next section.

**Corollary 1.** Assume that Assumption 1, Assumption 2, and Assumption 3 hold. Let

$$\lambda = \frac{4}{1 - \alpha} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} \sqrt{\frac{L_2 \log n}{N}}. \quad (5)$$

Under the sample size scaling in Part (a) of Theorem 2, with probability at least  $1 - c \exp(-c' \log s)$ ,  $\hat{\beta}$  has the following properties:

- (a) No false inclusion:  $\hat{S} \subseteq S$ .
- (b)  $\ell_2$ -estimation error:

$$\|\hat{\beta} - \beta^*\|_2 \leq \frac{24}{1 - \alpha} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} \lambda_{\min}^{-1}(\Sigma_x) \sqrt{\frac{s L_2 \log n}{N}}.$$

- (c)  $\ell_1$ -estimation error:

$$\|\hat{\beta} - \beta^*\|_1 \leq \frac{96}{1 - \alpha} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} \lambda_{\min}^{-1}(\Sigma_x) s \sqrt{\frac{L_2 \log n}{N}}.$$

Theorem 1, Theorem 2, and Corollary 1 provide theoretical guarantees for the estimation and variable selection performance of the Lasso under heteroscedasticity. But as a starting point in the statistical analysis routine for high-dimensional heteroscedastic linear models, the Lasso has two drawbacks: (a) it is difficult to quantify uncertainty of the Lasso performance. In other words, the Lasso does not have a uniform tractable limiting distribution [16], which can be seen from Figure 5 in Section 4.2. (b) The standard procedure in (3) does not efficiently employ the given information of the variance function. We will discuss a mature treatment paradigm to heteroscedasticity in the next section.

### 3 Robustifying the Lasso under heteroscedasticity

We now address the problem of robustifying the Lasso estimates when the parametric form of the variance function  $g$  is known. As mentioned in previous sections, we have two main goals: (a) make use of the variance information to devise new estimators, and (b) quantify uncertainties of the proposed estimators.

Following recent advances in statistical inference for high-dimensional homoscedastic linear models [13, 20, 24], we adopt the general framework of one-step estimators: start with the Lasso estimates, and then take a single update towards a certain direction. In this paper, we carefully choose the desired direction in the heteroscedastic context. The whole procedure of our approach is described in Algorithm 1 below:

---

**Algorithm 1** One-step estimator for  $\beta^*$  in high-dimensional heteroscedastic linear models

---

**Input:**  $N$  pairs of  $(X_i, y_i)$ , variance function  $g$ , and tuning parameters  $\lambda$  and  $\mu$ .

**Output:** One-step estimator  $\tilde{\beta}$ .

1. Solve the following optimization problem to obtain a preliminary estimator  $\hat{\beta}$ :

$$\hat{\beta} = \underset{\beta \in \mathbb{R}^n}{\operatorname{argmin}} \quad \frac{1}{2N} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1.$$

2. Set  $\widehat{W}_i = g(X_i, \hat{\beta})$ ,  $\widehat{W} = \operatorname{diag}(\widehat{W}_1, \dots, \widehat{W}_N) \in \mathbb{R}^{N \times N}$ , and  $\widehat{\Sigma}_N = \frac{1}{N} X^T \widehat{W}^{-1} X \in \mathbb{R}^{n \times n}$ .

3. Solve the following optimization problem to obtain  $\widehat{\Theta}$ :

$$\begin{aligned} \widehat{\Theta} &= \underset{\Theta \in \mathbb{R}^{n \times n}}{\operatorname{argmin}} \quad \|\Theta\|_{1,1} \\ \text{subject to} \quad & \|\Theta \widehat{\Sigma}_N - I_n\|_\infty \leq \mu. \end{aligned} \tag{6}$$

The above optimization consists of  $n$  subproblems: for  $i = 1, 2, \dots, n$ ,

$$\begin{aligned} \hat{\theta}_i &= \underset{\theta \in \mathbb{R}^n}{\operatorname{argmin}} \quad \|\theta\|_1 \\ \text{subject to} \quad & \|\widehat{\Sigma}_N \theta - e_i\|_\infty \leq \mu. \end{aligned} \tag{7}$$

Then  $\widehat{\Theta} = (\hat{\theta}_1, \dots, \hat{\theta}_n)^T$ . (6) and (7) can be efficiently solved via the ADMM algorithm [2].

4. Construct the one-step estimator as

$$\tilde{\beta} = \hat{\beta} + \frac{1}{N} \widehat{\Theta} X^T \widehat{W}^{-1} (y - X\hat{\beta}).$$


---

Next, we present theoretical guarantees of the proposed one-step estimator. Our theory requires another set of assumptions:

**Assumption 4.** For all  $\beta$  and  $X_i$ , the first two partial derivatives of  $g(X_i, \beta)$  with respect to  $\beta$  are uniformly bounded functions of  $\beta$ . That is, there exists a constant  $L_g > 0$ , such that

$$\left\| \frac{\partial g(X_i, \beta)}{\partial \beta} \right\|_\infty \leq L_g \quad \text{and} \quad \left\| \frac{\partial^2 g(X_i, \beta)}{\partial \beta^2} \right\|_\infty \leq L_g.$$

**Assumption 5.** Expectation of  $\frac{X_i X_i^T}{W_i}$  exists. That is,  $\mathbb{E} \left( \frac{X_i X_i^T}{W_i} \right) = \Sigma^*$ , where  $\Sigma^* \in \mathbb{R}^{n \times n}$  is positive definite. In what follows, let  $\Theta^* \in \mathbb{R}^{n \times n}$  denote the inverse of  $\Sigma^*$ . Furthermore, we assume the following conditions for  $\Theta^*$ :

- (a) Eigenvalues of  $\Theta^*$  are bounded. That is,

$$c_{\min} \leq \lambda_{\min}(\Theta^*) \leq \lambda_{\max}(\Theta^*) \leq c_{\max},$$

where  $c_{\min} > 0$  and  $c_{\max} > 0$  are constants.

- (b) For  $1 \leq i \leq n$ ,  $\|\Theta_i^*\|_0 \leq k$ , where  $\Theta_i^*$  is the  $i$ th row of  $\Theta^*$  and  $k \ll n$ .
- (c) There exists a number  $M$ , such that  $\|\Theta^*\|_\infty \leq M$ .

Assumption 4 imposes similar conditions on the variance function  $g$  as in [15] and [21]. As mentioned in [15], this assumption is usually required in nonlinear least squares estimation. In Assumption 5, we define an appropriate inverse variance matrix  $\Theta^*$  in order to adapt the statistical analysis to the heteroscedastic setting. Note that  $\hat{\Theta}$  in step 3 of Algorithm 1 can be viewed as an estimator of  $\Theta^*$ . We now derive the limiting distribution of the one-step estimator in the following:

**Theorem 3.** Assume that Assumption 1 to Assumption 5 hold. In Algorithm 1, let

$$\lambda \asymp \sqrt{\frac{L_2 \log n}{N}} \quad \text{and} \quad \mu \asymp \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}}.$$

If  $N$  satisfies the sample size scaling in Part (a) of Theorem 2 and  $L_2 M s^2 (\log n)^2 / \sqrt{N} \rightarrow 0$ , then we have for  $1 \leq j \leq n$ ,

$$\sqrt{N} \left( \tilde{\beta}_j - \beta_j^* \right) \xrightarrow{d} N \left( 0, e_j^T \Theta^* e_j \right).$$

The proof of Theorem 3 is contained in Section 6.3.

**Remark 2.** Theorem 3 shows that  $\tilde{\beta}_j$  is asymptotically normal with mean  $\beta_j^*$  and variance  $\frac{e_j^T \Theta^* e_j}{N}$ . Furthermore, it has several consequences for statistical inference of regression parameter under heteroscedasticity. For example, similar arguments yield

$$\frac{\sqrt{N} \left( \tilde{\beta}_j - \beta_j^* \right)}{\sqrt{e_j^T \hat{\Theta} \hat{\Sigma}_N \hat{\Theta}^T e_j}} \rightarrow N(0, 1). \quad (8)$$

Let  $\Phi(x)$  be the cumulative distribution function of  $N(0, 1)$ . Then

$$\left[ \tilde{\beta}_j - \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right) \sqrt{\frac{e_j^T \hat{\Theta} \hat{\Sigma}_N \hat{\Theta}^T e_j}{N}}, \quad \tilde{\beta}_j + \Phi^{-1} \left( 1 - \frac{\alpha}{2} \right) \sqrt{\frac{e_j^T \hat{\Theta} \hat{\Sigma}_N \hat{\Theta}^T e_j}{N}} \right]$$

provides an asymptotically valid  $(1 - \alpha)$ -confidence interval for  $\beta_j$ . On the other hand, the pivot in (8) is a test statistic for testing whether  $\beta_j$  is equal to 0.

## 4 Simulations

In this section, we perform simulation studies to validate our theoretical results of this paper. Our focus of this section is two-fold. First, we study the finite sample performance of the Lasso under different kinds of heteroscedasticity. Second, we investigate the asymptotic behavior of the Lasso and one-step estimator. The code for replicating our experiments is available on GitHub <https://github.com/dzwang91/HHLM>.



#### 4.1 Performance of the Lasso under heteroscedasticity

In the first set of simulations, our goal is to verify statistical consistency (estimation consistency, prediction consistency, and variable selection consistency) of the Lasso when the underlying model is heteroscedastic. In all experiments, we set

$$\beta^* = (3, 4, 3, 1.5, 2, 1.5, 0, \dots, 0)^T.$$

We first generate predictor variables  $X_i$  from the normal distribution  $N(0, \Sigma_x)$ , where  $(\Sigma_x)_{S^c S} = 0$ , and  $(\Sigma_x)_{ij} = 0.5^{|i-j|}$  if both  $i$  and  $j$  are in  $S$  or both  $i$  and  $j$  are in  $S^c$ . Then the random error  $\varepsilon_i$  is drawn from the conditional distribution  $\varepsilon_i \mid X_i \sim N(0, W_i)$ . Finally, we generate the response  $y_i$  via the linear model  $y_i = X_i^T \beta^* + \varepsilon_i$ . We consider the following four different scenarios on  $W_i$ :

(a) Type 1:

$$W_i = \min \left( \frac{1}{25} \exp \left( \frac{1}{2} |X_i^T \beta^*| \right), 5 \right).$$

(b) Type 2:

$$W_i = \min \left( \frac{1}{50} \exp \left( \frac{1}{20} (X_i^T \beta^*)^2 \right), 5 \right).$$

(c) Type 3:

$$W_i = \begin{cases} \frac{1}{20} & \frac{1}{4} |X_i^T \beta^*| \leq \frac{1}{20} \\ \frac{1}{4} |X_i^T \beta^*| & \frac{1}{20} < \frac{1}{4} |X_i^T \beta^*| \leq 5 \\ 5 & \frac{1}{4} |X_i^T \beta^*| > 5 \end{cases}$$

(d) Type 4:

$$W_i = \begin{cases} \frac{1}{20} & \frac{1}{16} (X_i^T \beta^*)^2 \leq \frac{1}{20} \\ \frac{1}{16} (X_i^T \beta^*)^2 & \frac{1}{20} < \frac{1}{16} (X_i^T \beta^*)^2 \leq 5 \\ 5 & \frac{1}{16} (X_i^T \beta^*)^2 > 5 \end{cases}$$

These four types are similar to those in [21]. For each type of heteroscedastic models, we consider three ambient dimensions:  $n = 150, 300$ , and  $500$ . For each  $n$ , we run simulations for eight different sample sizes (see Figures 1, 2, 3, and 4). Furthermore, following our theory, we set the tuning parameter  $\lambda = 0.75 \sqrt{\lambda_{\max}(\Sigma_x) \frac{L_2 \log n}{N}}$  in each simulation.

Figure 1, Figure 2, Figure 3, and Figure 4 show the simulation results for those four heteroscedastic linear models listed above. Each point in these figures represents an average over 200 trials. Panel (a), Panel (b), and Panel (c) in each figure confirm that the  $\ell_2$ -error, the  $\ell_1$ -error, and the mean squared prediction error decrease to 0 as the sample size increases, showing the estimation and prediction consistency of the Lasso under different kinds of heteroscedasticity. In Panel (d), we see that the empirical probability of recovering the correct support set transitions sharply from 0 to 1, which implies that the Lasso is variable selection consistent, as stated in Theorem 2.

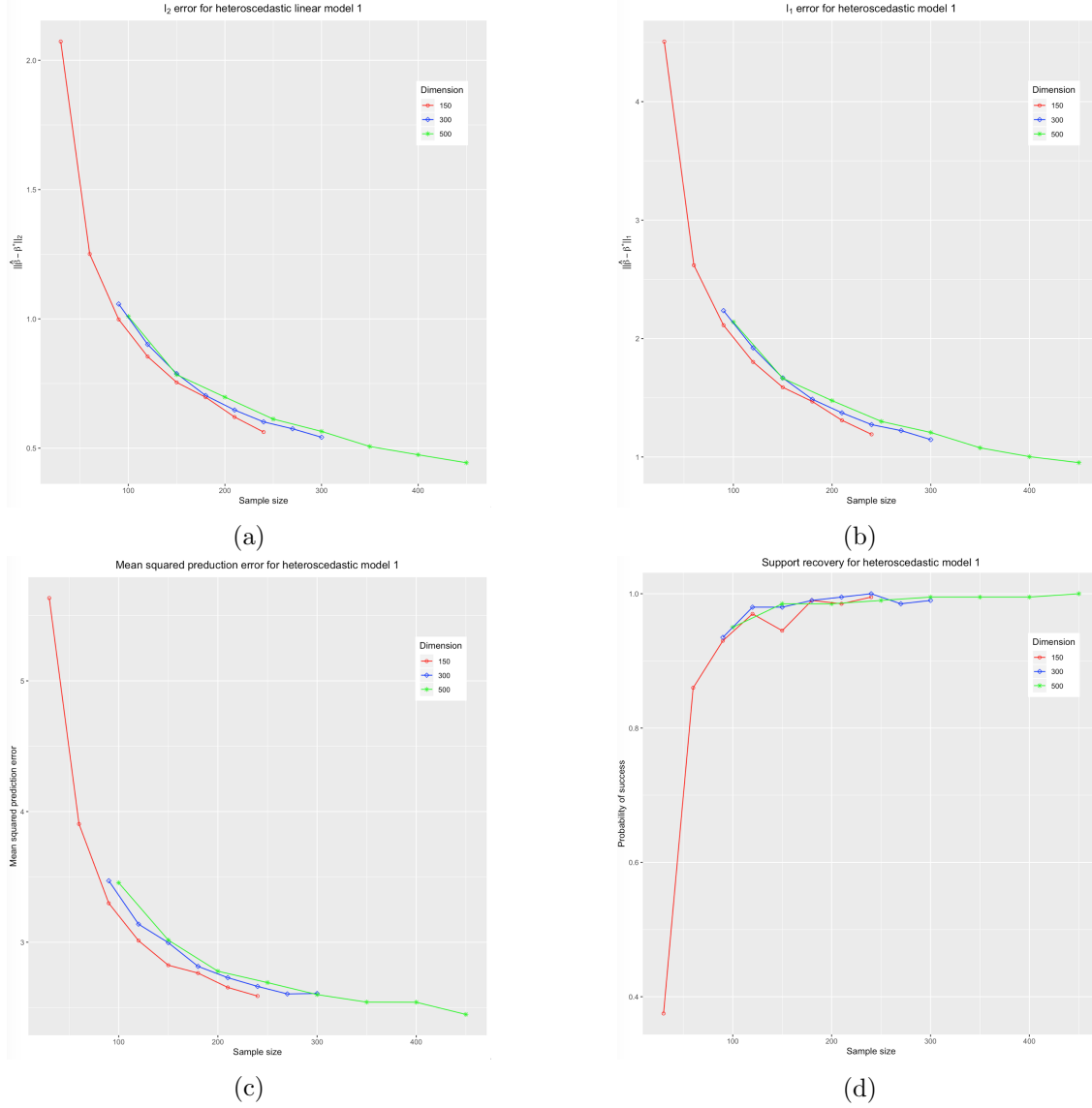


Figure 1: Plots showing simulation results of the Lasso for high-dimensional heteroscedastic linear model 1 with three problem sizes:  $n = 150$  (red),  $n = 300$  (blue), and  $n = 500$  (green). Each point represents an average over 200 trials. (a) Plot showing consistency of the  $\ell_2$ -estimation error  $\|\hat{\beta} - \beta^*\|_2$ . (b) Plot showing consistency of the  $\ell_1$ -estimation error  $\|\hat{\beta} - \beta^*\|_1$ . (c) Plot showing consistency of the mean squared prediction error  $\frac{1}{N} \|y - X\hat{\beta}\|_2^2$ . (d) Plot showing variable selection consistency, measured by the empirical success probability in recovering the support set of  $\beta^*$ .

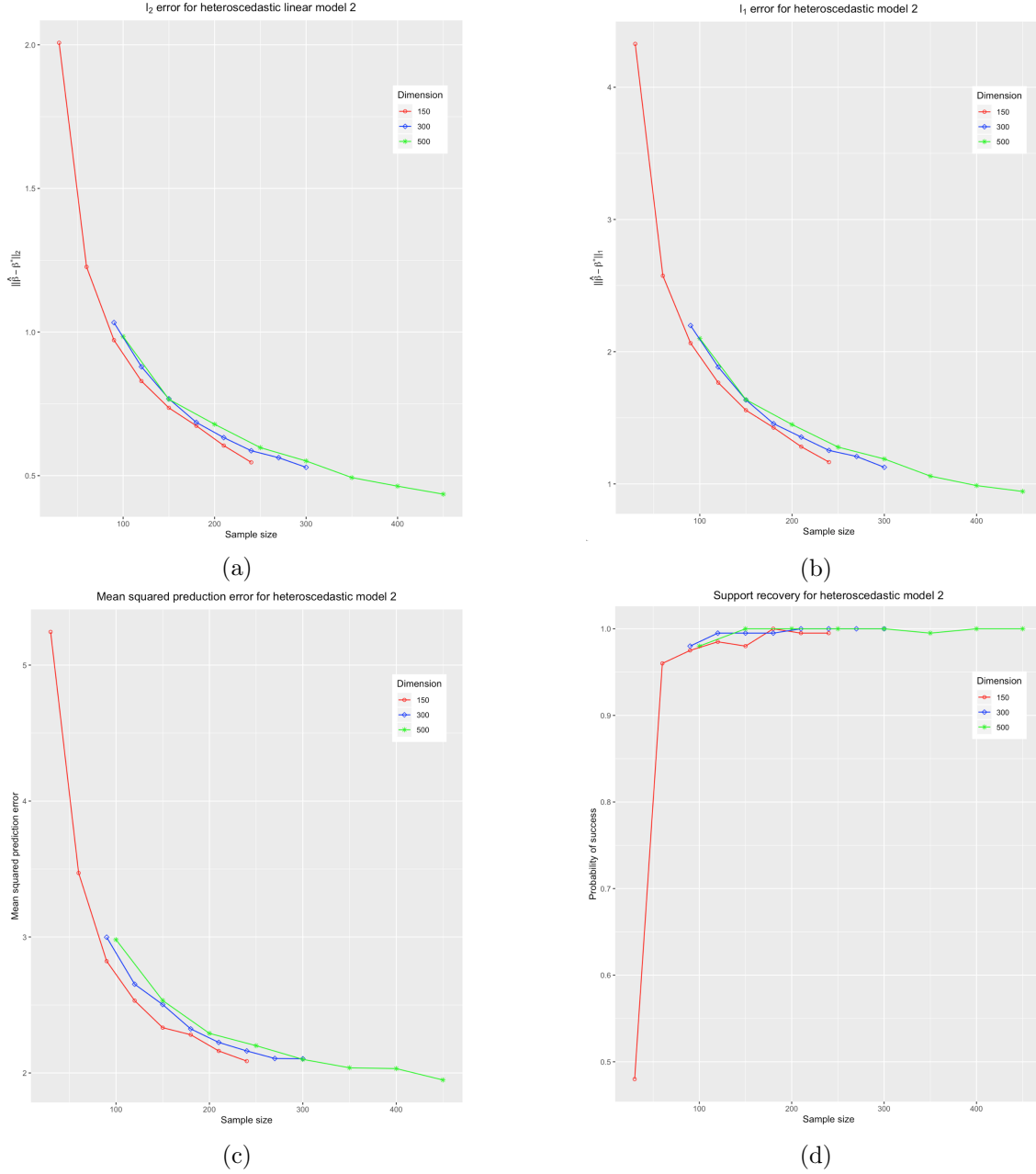


Figure 2: Plots showing simulation results of the Lasso for high-dimensional heteroscedastic linear model 2 with three problem sizes:  $n = 150$  (red),  $n = 300$  (blue), and  $n = 500$  (green). Each point represents an average over 200 trials. (a) Plot showing consistency of the  $\ell_2$ -estimation error  $\|\hat{\beta} - \beta^*\|_2$ . (b) Plot showing consistency of the  $\ell_1$ -estimation error  $\|\hat{\beta} - \beta^*\|_1$ . (c) Plot showing consistency of the mean squared prediction error  $\frac{1}{N} \|y - X\hat{\beta}\|_2^2$ . (d) Plot showing variable selection consistency, measured by the empirical success probability in recovering the support set of  $\beta^*$ .

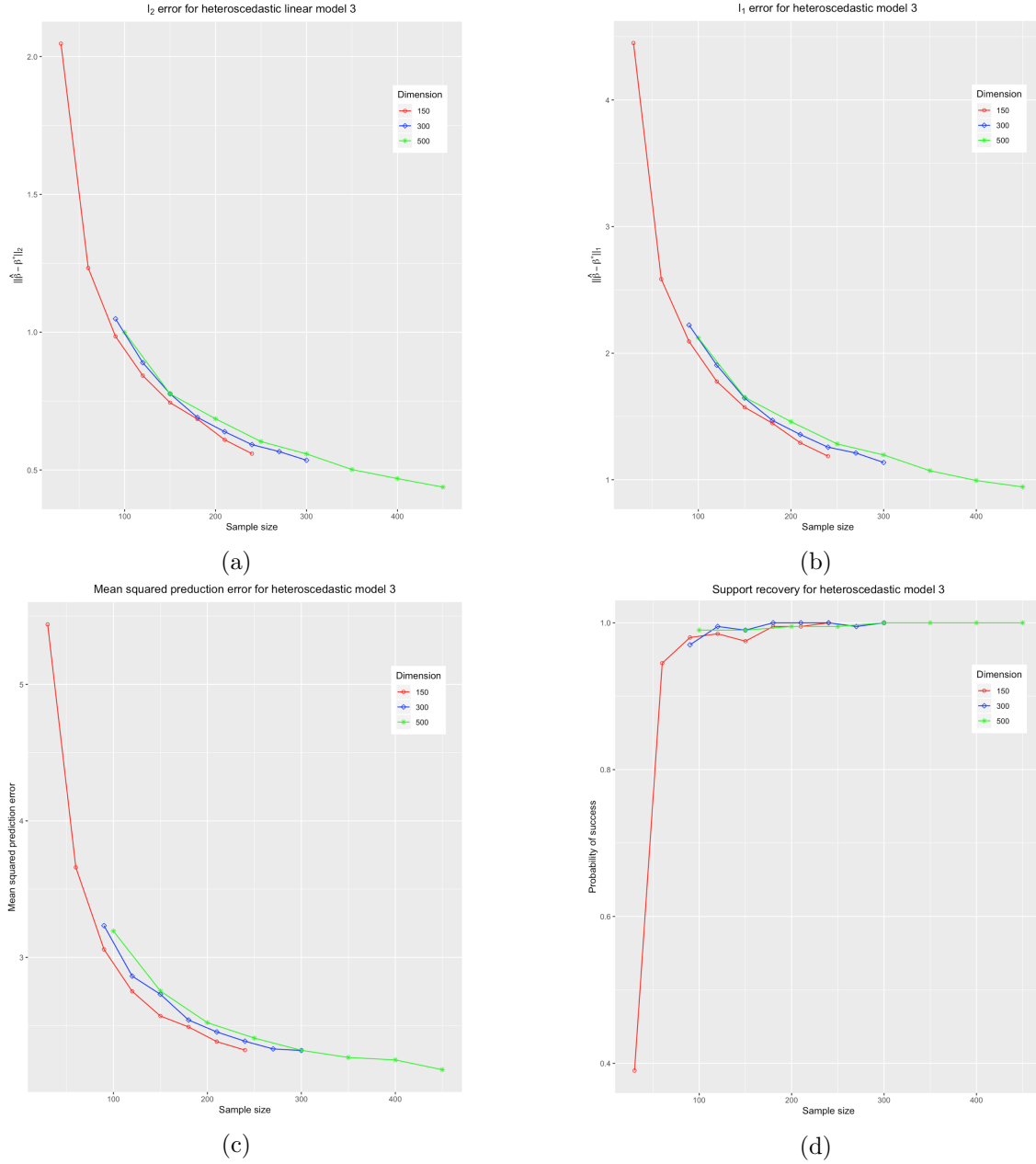


Figure 3: Plots showing simulation results of the Lasso for high-dimensional heteroscedastic linear model 3 with three problem sizes:  $n = 150$  (red),  $n = 300$  (blue), and  $n = 500$  (green). Each point represents an average over 200 trials. (a) Plot showing consistency of the  $\ell_2$ -estimation error  $\|\hat{\beta} - \beta^*\|_2$ . (b) Plot showing consistency of the  $\ell_1$ -estimation error  $\|\hat{\beta} - \beta^*\|_1$ . (c) Plot showing consistency of the mean squared prediction error  $\frac{1}{N} \|y - X\hat{\beta}\|_2^2$ . (d) Plot showing variable selection consistency, measured by the empirical success probability in recovering the support set of  $\beta^*$ .

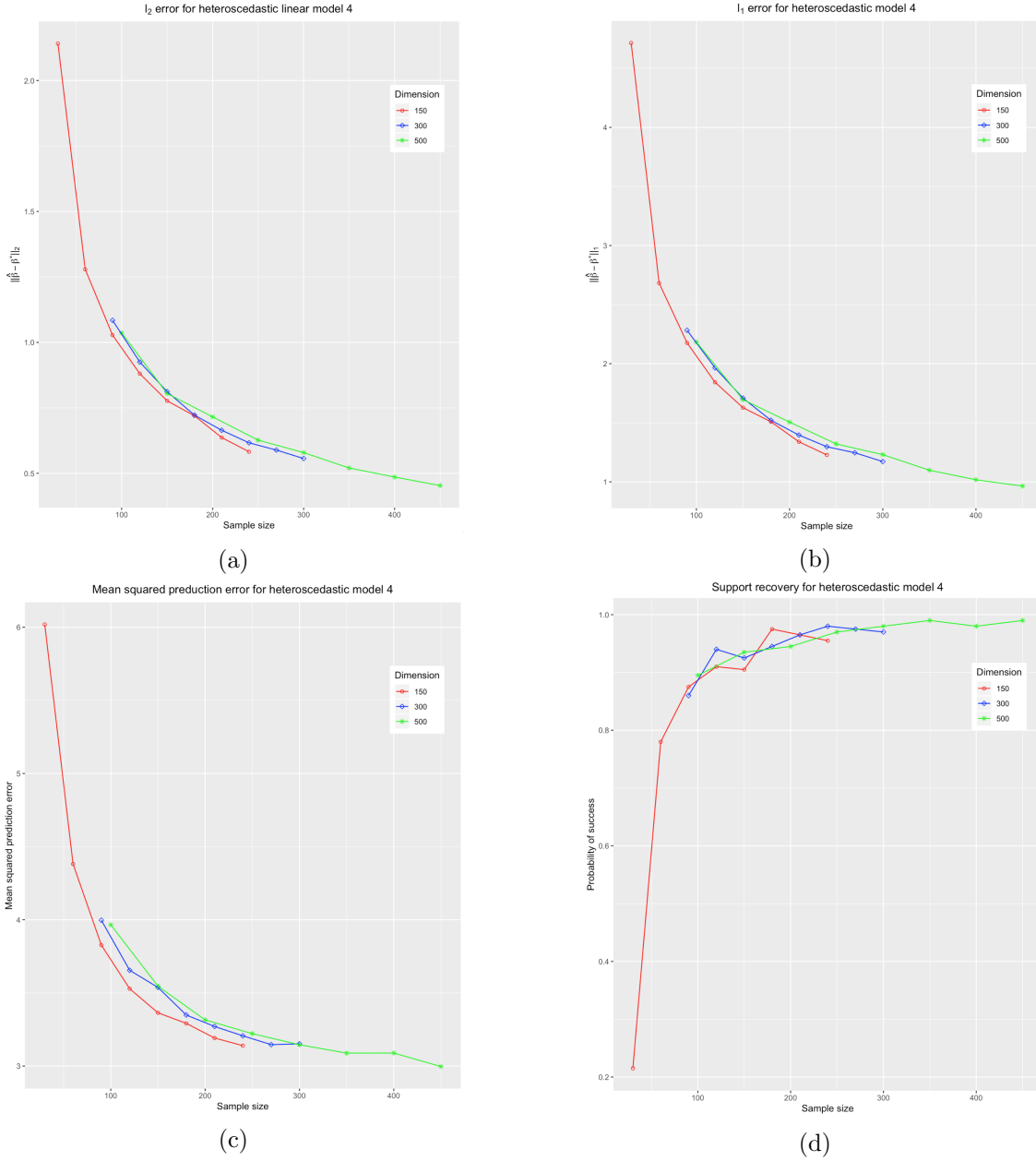


Figure 4: Plots showing simulation results of the Lasso for high-dimensional heteroscedastic linear model 4 with three problem sizes:  $n = 150$  (red),  $n = 300$  (blue), and  $n = 500$  (green). Each point represents an average over 200 trials. (a) Plot showing consistency of the  $\ell_2$ -estimation error  $\|\hat{\beta} - \beta^*\|_2$ . (b) Plot showing consistency of the  $\ell_1$ -estimation error  $\|\hat{\beta} - \beta^*\|_1$ . (c) Plot showing consistency of the mean squared prediction error  $\frac{1}{N} \|y - X\hat{\beta}\|_2^2$ . (d) Plot showing variable selection consistency, measured by the empirical success probability in recovering the support set of  $\beta^*$ .

## 4.2 Performance of one-step estimators

Our second set of simulations investigates asymptotics of the proposed one-step estimator in Algorithm 1. We take the heteroscedastic linear model 1 with  $(N, n) = (120, 150)$  as an example. The data generating process in this set of simulations is the same as the previous set. Based on our theoretical results, in Algorithm 1, we set

$$\lambda = 0.75 \sqrt{\lambda_{\max}(\Sigma_x) \frac{L_2 \log n}{N}} \quad \text{and} \quad \mu = 0.0001 \left( \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}} \right).$$

Finally, we repeat the whole procedure of Algorithm 1 200 times.

In order to confirm the limiting distribution of one-step estimators stated in Theorem 3, we report simulation results in two cases: (a) the sampling distribution of  $\sqrt{N}(\tilde{\beta}_5 - \beta_5^*)$  where  $\beta_5^* \neq 0$ , and (b) the sampling distribution of  $\sqrt{N}(\tilde{\beta}_7 - \beta_7^*)$  where  $\beta_7^* = 0$ . In Figure 5, we present Q-Q plots of  $\sqrt{N}(\tilde{\beta}_5 - \beta_5^*)$  and  $\sqrt{N}(\tilde{\beta}_7 - \beta_7^*)$  in Panel (a) and Panel (b), respectively. It is very clear that both sampling distributions are normal with mean 0. We also perform the Shapiro-Wilk normality test for samples from  $\sqrt{N}(\tilde{\beta}_5 - \beta_5^*)$  and  $\sqrt{N}(\tilde{\beta}_7 - \beta_7^*)$ . The corresponding p-values are 0.30 and 0.95, which also certifies our conclusions. Furthermore, in Panel (c) and Panel (d) of Figure 5, we display Q-Q plots of  $\sqrt{N}(\hat{\beta}_5 - \beta_5^*)$  and  $\sqrt{N}(\hat{\beta}_7 - \beta_7^*)$ , where  $\hat{\beta}_5$  and  $\hat{\beta}_7$  are obtained in the first step of Algorithm 1. Obviously, we are not able to claim normality from these two Q-Q plots, compared with those in Panel (a) and Panel (b).

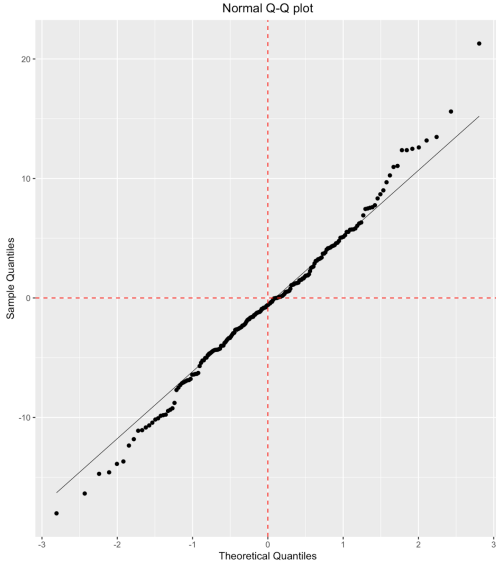
## 5 Discussion

Our work recommends a number of directions for future research. First, we have only studied the case that the variance is a function of regression parameter. We can extend the current model to a broader one, where variances also depend on parameters other than  $\beta^*$  as in [15]. Secondly, in this paper, we have access to partial information of the variance function—the parametric form of  $g$ , through domain knowledge. Another scenario of interest involves completely unknown variance function. As mentioned earlier in Section 1, when only assuming the variance function was smooth in the classical low-dimensional regime, Carroll [3] constructed an estimator of the regression parameter and showed that it was asymptotically equivalent to the weighted least squares estimator with known variances. Likewise, for the high-dimensional regime, it is worthwhile to explore whether there exists some nonparametric approach to recover the results in our paper.

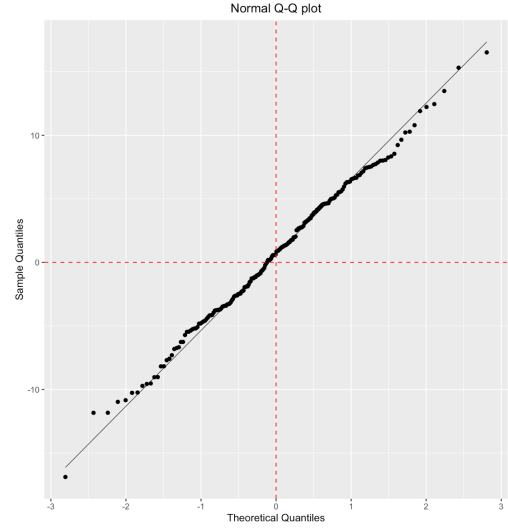
Finally, another line of future work is motivated by some interesting results of Jobson and Fuller [15], which are derived for the fixed design in the low-dimensional setting. Under a similar model to ours, the authors developed a weighted joint least squares estimator  $\hat{\beta}^{JF}$  which is asymptotically equivalent to the maximum likelihood estimator. More specifically, they proved

$$\sqrt{N}(\hat{\beta}^{JF} - \beta^*) \rightarrow N \left( 0, \left( \lim_{N \rightarrow \infty} \frac{1}{N} H^T F^{-1} H + \lim_{N \rightarrow \infty} X^T W^{-1} X \right)^{-1} \right), \quad (9)$$

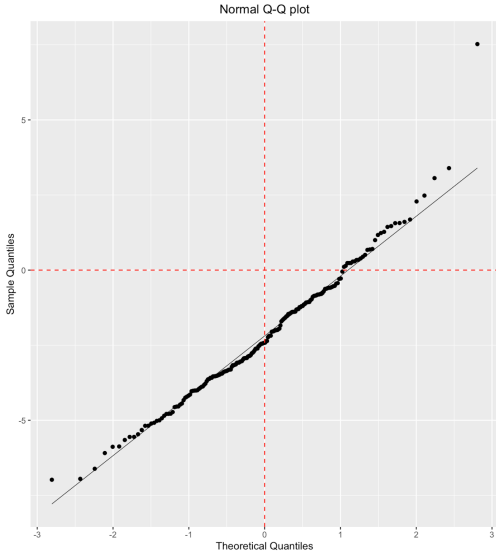
where  $F \in \mathbb{R}^{N \times N}$  is a diagonal matrix with  $i$ th diagonal element  $2g^2(X_i, \beta^*)$ , and  $H \in \mathbb{R}^{N \times n}$  with  $(i, j)$ th element  $H_{ij} = \frac{\partial g(X_i, \beta^*)}{\partial \beta_j}$ . It can be seen that the asymptotic variance in (9) is smaller than



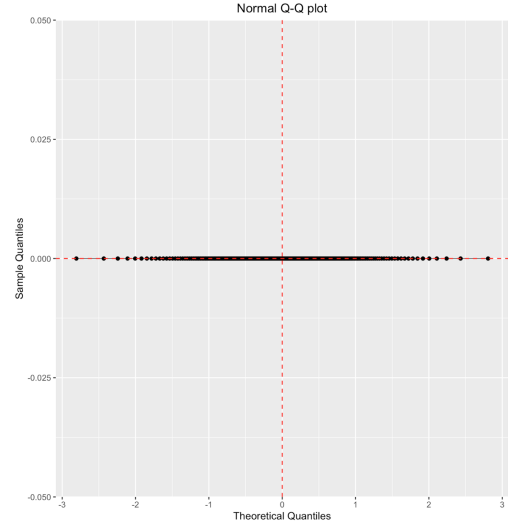
(a)



(b)



(c)



(d)

Figure 5: Plots showing simulation results for high-dimensional heteroscedastic linear model 1 with  $(N, n) = (120, 150)$ . (a) Q-Q plot of  $\sqrt{N}(\hat{\beta}_5 - \beta_5^*)$ . (b) Q-Q plot of  $\sqrt{N}(\hat{\beta}_7 - \beta_7^*)$ . (c) Q-Q plot of  $\sqrt{N}(\hat{\beta}_5 - \beta_5^*)$ . (d) Q-Q plot of  $\sqrt{N}(\hat{\beta}_7 - \beta_7^*)$ . Panel (a) and Panel (b) verify the asymptotic normality of one-step estimators, while Panel (c) and Panel (d) illustrate that the Lasso does not have a tractable limiting distribution.

$(\lim_{N \rightarrow \infty} \frac{1}{N} X^T W^{-1} X)^{-1}$ , the counterpart of which in the random design is  $\Theta^*$  defined in Assumption 5. Therefore, an important but challenging open question is whether our proposed one-step estimator is statistically efficient in some sense for the high-dimensional conditional heteroscedastic linear model (1). Of course, this question is somewhat vague in that the notion of efficiency has not been well established in the high-dimensional setting. From a semiparametric point of view, Janková and van de Geer [12] studied efficiency lower bounds on variance for high-dimensional homoscedastic linear models. It should be interesting to investigate the similar problem for high-dimensional heteroscedastic linear models to understand limits of the proposed one-step estimator.

## 6 Proofs of main results

In this section, we prove Theorem 1, Theorem 2, and Theorem 3. Proofs of the supporting lemmas are contained in Appendix A.

### 6.1 Proof of Theorem 1

First, we assume two conditions:

- The tuning parameter  $\lambda$  satisfies

$$\lambda > \frac{2}{N} \|\varepsilon^T X\|_\infty. \quad (10)$$

- There exists  $\eta > 0$ , such that

$$\frac{1}{N} \|Xv\|_2^2 \geq \eta \|v\|_2^2 \quad (11)$$

for  $v \in \mathbb{R}^n$  and  $\|v_{S^c}\|_1 \leq 3\|v_S\|_1$ .

By optimality of  $\hat{\beta}$ , we have the following basic inequality:

$$\frac{1}{2N} \|y - X\hat{\beta}\|_2^2 + \lambda \|\hat{\beta}\|_1 \leq \frac{1}{2N} \|y - X\beta^*\|_2^2 + \lambda \|\beta^*\|_1.$$

Rearranging terms and letting  $\Delta = \hat{\beta} - \beta^*$ , we have

$$\begin{aligned} \frac{1}{2N} \|X\Delta\|_2^2 &\leq \frac{1}{N} (y - X\beta^*)^T X\Delta + \lambda (\|\beta^*\|_1 - \|\hat{\beta}\|_1) \\ &= \frac{1}{N} \varepsilon^T X\Delta + \lambda (\|\beta_S^*\|_1 - \|\hat{\beta}_S\|_1 - \|\hat{\beta}_{S^c}\|_1) \\ &\leq \frac{1}{N} \varepsilon^T X\Delta + \lambda (\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1) \\ &\leq \frac{1}{N} \|\varepsilon^T X\|_\infty \|\Delta\|_1 + \lambda (\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1). \end{aligned}$$

When condition (10) holds, we have

$$0 \leq \frac{1}{2N} \|X\Delta\|_2^2 \leq \frac{\lambda}{2} (\|\Delta_S\|_1 + \|\Delta_{S^c}\|_1) + \lambda (\|\Delta_S\|_1 - \|\Delta_{S^c}\|_1).$$

Therefore, we have

$$\|\Delta_{S^c}\|_1 \leq 3\|\Delta_S\|_1.$$



Combining condition (11), we have

$$\eta \|\Delta\|_2^2 \leq \frac{1}{N} \|X\Delta\|_2^2.$$

Therefore, we have

$$\eta \|\Delta\|_2^2 \leq \frac{1}{N} \|X\Delta\|_2^2 \leq 3\lambda \|\Delta_S\|_1 \leq 3\lambda\sqrt{s} \|\Delta\|_2,$$

implying that

$$\|\Delta\|_2 \leq \frac{3\lambda\sqrt{s}}{\eta}. \quad (12)$$

Furthermore, we have

$$\|\Delta\|_1 \leq 4\|\Delta_S\|_1 \leq 4\sqrt{s} \|\Delta\|_2 \leq \frac{12\lambda s}{\eta}. \quad (13)$$

Next, to establish the results in the theorem, we use the following two lemmas, which are proved in Appendix A.1 and Appendix A.2, respectively.

**Lemma 1.** Assume that Assumption 2 holds. Let  $\lambda \geq 4(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} L_2^{1/2} (\log n/N)^{1/2}$ . Then we have

$$\mathbb{P}\left(\frac{2}{N} \|\varepsilon^T X\|_\infty > \lambda\right) \leq 2 \exp(-\log n) + \exp(\log n - cN),$$

where  $c > 0$  is a constant.

**Lemma 2.** Assume that Part (a) of Assumption 2 holds. Then we have

$$\frac{1}{N} \|Xv\|_2^2 \geq \frac{1}{2} \lambda_{\min}(\Sigma_x) \|v\|_2^2$$

for  $v \in \mathbb{R}^n$  and  $\|v_{S^c}\|_1 \leq 3\|v_S\|_1$  with probability at least  $1 - 2 \exp(cs \log n - c'N)$ , where  $s$  and  $S$  are defined in Assumption 1, and  $c > 0$  and  $c' > 0$  are constants.

Therefore, letting  $\lambda \geq 4(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} L_2^{1/2} (\log n/N)^{1/2}$  in (10) and  $\eta = \frac{1}{2} \lambda_{\min}(\Sigma_x)$  in (11), and combining Lemma 1, Lemma 2, (12), and (13), we have

$$\|\Delta\|_2 \leq 6\lambda_{\min}^{-1}(\Sigma_x) \lambda \sqrt{s}, \quad \|\Delta\|_1 \leq 24\lambda_{\min}^{-1}(\Sigma_x) \lambda s,$$

with probability at least  $1 - 2 \exp(-\log n) - c \exp(c' s \log n - c''N)$ . Finally, letting  $N \gtrsim s \log n$  yields the desired probability in the theorem. Hence the proof is complete.

## 6.2 Proof of Theorem 2

Results in Theorem 2 build upon the following lemma:

**Lemma 3.** Let  $S$  and  $\hat{S}$  be support sets of  $\beta^*$  and  $\hat{\beta}$ , respectively.

(a) Assume that the following conditions hold:

(C1) The smallest eigenvalue of  $\frac{1}{N} X_S^T X_S$  is bounded below. That is,

$$\lambda_{\min}\left(\frac{1}{N} X_S^T X_S\right) \geq c_{\min} > 0. \quad (14)$$

(C2) There exists a constant  $\alpha \in (0, 1)$  such that

$$\max_{j \in S^c} \|(X_S^T X_S)^{-1} X_S^T X_j\|_1 \leq \alpha. \quad (15)$$

(C3) The tuning parameter  $\lambda$  satisfies

$$\lambda \geq \frac{2}{1 - \alpha} \left\| X_{S^c}^T (I_N - X_S (X_S^T X_S)^{-1} X_S^T) \frac{\varepsilon}{N} \right\|_\infty. \quad (16)$$

Then we have  $\hat{S} \subseteq S$ .

(b) In addition, assume that the following condition holds:

(C4)  $\min_{i \in S} |\beta_i^*|$  satisfies

$$\min_{i \in S} |\beta_i^*| > B(\lambda, X_S, \varepsilon) = \left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} X_S^T \frac{\varepsilon}{N} \right\|_\infty + \lambda \left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} \right\|_\infty. \quad (17)$$

Then we have  $S = \hat{S}$ .

The proof of Lemma 3 is deferred to Appendix A.3. Therefore, it suffices to show that for the given assumptions in Theorem 2, conditions (14), (15), (16), and (17) hold with high probability. These probabilistic results are established in the following lemmas:

**Lemma 4.** Assume that Part (a) of Assumption 2 holds. Under the sample size scaling  $N \gtrsim s$ , we have

$$\lambda_{\min} \left( \frac{1}{N} X_S^T X_S \right) \geq \frac{1}{2} \lambda_{\min}(\Sigma_x),$$

with probability at least  $1 - 2 \exp(-cN)$ .

**Lemma 5.** Assume that Part (a) of Assumption 2 and Assumption 3 hold. When

$$N^2 \gtrsim \frac{s^5 \log[s(n-s)]}{\alpha^2}, \quad N \gtrsim \frac{s^4 \log[s(n-s)]}{\alpha^2}, \quad \text{and} \quad N \gtrsim \frac{s^5}{\alpha^2},$$

we have (15) holds with probability at least  $1 - c' \exp(-c''s) - C' \exp(-C'' \log[s(n-s)])$ .

**Lemma 6.** Assume that Assumption 2 holds. Let

$$\lambda \geq \frac{4}{1 - \alpha} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} \sqrt{\frac{L_2 \log(n-s)}{N}}.$$

Under the sample size scaling  $N \gtrsim \log(n-s)$ , we have

$$\lambda \geq \frac{2}{1 - \alpha} \left\| X_{S^c}^T (I_N - X_S (X_S^T X_S)^{-1} X_S^T) \frac{\varepsilon}{N} \right\|_\infty,$$

with probability at least  $1 - c \exp(-c' \log(n-s))$ .

**Lemma 7.** Assume that Assumption 2 holds. Under the sample size scaling  $N \gtrsim s$ , we have

$$B(\lambda, X_S, \varepsilon) \leq \sqrt{\frac{8}{\lambda_{\min}(\Sigma_x)}} \sqrt{\frac{L_2 \log s}{N}} + \lambda \left( \left\| (\Sigma_x)_{SS}^{-1} \right\|_\infty + cs \sqrt{\frac{s}{N}} \right),$$

with probability at least  $1 - c' \exp(-c'' \log s)$ .

Proofs of Lemma 4, Lemma 5, Lemma 6, and Lemma 7 are deferred to Appendix A.4, Appendix A.5, Appendix A.6, and Appendix A.7, respectively. Therefore, the proof of Theorem 2 is complete.

### 6.3 Proof of Theorem 3

Let  $\check{\beta} = \hat{\beta} + \frac{1}{N}\Theta^*X^TW^{-1}(y - X\hat{\beta})$  and  $\Sigma_N = \frac{1}{N}X^TW^{-1}X$ . We divide the proof of Theorem 3 into three main steps: first, we show that

$$\max_j \frac{|\sqrt{N}(\check{\beta}_j - \check{\beta}_j^*)|}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} = o_{\mathbb{P}}(1) \quad (18)$$

under given assumptions in Section 6.3.1. Next, in Section 6.3.2, we show that for  $1 \leq j \leq n$ ,

$$\frac{\sqrt{N}(\check{\beta}_j - \check{\beta}_j^*)}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} \xrightarrow{d} N(0, 1). \quad (19)$$

Finally, we combine these previous results in Section 6.3.3 to obtain the desired result in the theorem.

#### 6.3.1 Step 1 of proof

For  $1 \leq j \leq n$ , we have

$$\begin{aligned} \sqrt{N}(\check{\beta}_j - \check{\beta}_j^*) &= \sqrt{N}e_j^T (\Theta^* \Sigma_N - I_n) (\hat{\beta} - \beta^*) - \sqrt{N}e_j^T (\hat{\Theta} \hat{\Sigma}_N - I_n) (\hat{\beta} - \beta^*) \\ &\quad + \frac{1}{\sqrt{N}}e_j^T (\hat{\Theta} - \Theta^*) X^T \widehat{W}^{-1} \varepsilon + \frac{1}{\sqrt{N}}e_j^T \Theta^* X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon. \end{aligned}$$

Therefore, we have

$$\begin{aligned} \max_j |\sqrt{N}(\check{\beta}_j - \check{\beta}_j^*)| &\leq \max_j |\sqrt{N}e_j^T (\Theta^* \Sigma_N - I_n) (\hat{\beta} - \beta^*)| + \max_j |\sqrt{N}e_j^T (\hat{\Theta} \hat{\Sigma}_N - I_n) (\hat{\beta} - \beta^*)| \\ &\quad + \max_j \left| \frac{1}{\sqrt{N}}e_j^T (\hat{\Theta} - \Theta^*) X^T \widehat{W}^{-1} \varepsilon \right| + \max_j \left| \frac{1}{\sqrt{N}}e_j^T \Theta^* X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right| \\ &\leq \sqrt{N}(\|\Theta^* \Sigma_N - I_n\|_{\infty} + \mu) \|\hat{\beta} - \beta^*\|_1 + \|\hat{\Theta} - \Theta^*\|_{\infty} \left\| \frac{1}{\sqrt{N}} X^T \widehat{W}^{-1} \varepsilon \right\|_{\infty} + \|\Theta^*\|_{\infty} \left\| \frac{1}{\sqrt{N}} X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right\|_{\infty}. \end{aligned}$$

We've bounded  $\|\hat{\beta} - \beta^*\|_1$  in Corollary 1. Therefore, it remains to bound  $\|\Theta^* \Sigma_N - I_n\|_{\infty}$ ,  $\|\hat{\Theta} - \Theta^*\|_{\infty}$ ,  $\left\| \frac{1}{\sqrt{N}} X^T \widehat{W}^{-1} \varepsilon \right\|_{\infty}$ , and  $\left\| \frac{1}{\sqrt{N}} X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right\|_{\infty}$ . We establish these bounds in the following lemmas and their proofs are contained in Appendix A.8, Appendix A.9, and Appendix A.10, respectively.

**Lemma 8.** Assume that Part (a) of Assumption 2 holds. We have

$$\|\Theta^* \Sigma_N - I_n\|_{\infty} \leq \frac{c}{L_1} \sqrt{\frac{\log n}{N}},$$

with probability at least  $1 - \frac{2}{n}$ .

**Lemma 9.** Assume that same conditions as Corollary 1 and Assumption 4 hold. Let

$$\mu \asymp \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}}$$

in Algorithm 1 and  $N \geq s^2 L_2 \log n$ . Then we have

$$\|\hat{\Theta} - \Theta^*\|_\infty \leq ckM \left( \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}} \right),$$

with probability at least  $1 - c' \exp(-c'' \log s)$ .

**Lemma 10.** Assume that same conditions as Corollary 1 and Assumption 4 hold.

(a) We have

$$\left\| \frac{1}{\sqrt{N}} X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right\|_\infty \leq c \frac{L_2 s^2 (\log n)^2}{\sqrt{N}},$$

with probability at least  $1 - c' \exp(-c'' \log s)$ .

(b) We have

$$\left\| \frac{1}{\sqrt{N}} X^T \widehat{W}^{-1} \varepsilon \right\|_\infty \leq c \left( \frac{L_2 s^2 (\log n)^2}{\sqrt{N}} + \sqrt{\log n} \right),$$

with probability at least  $1 - c' \exp(-c'' \log s)$ .

Therefore, by above lemmas, we have

$$\max_j \left| \sqrt{N} (\tilde{\beta}_j - \check{\beta}_j) \right| \leq c \left[ \frac{s \log n}{L_1} \sqrt{\frac{L_2}{N}} + \frac{L_2 s^2 \log n}{\sqrt{N}} + kM \left( \frac{\log n}{L_1 \sqrt{N}} + s \log n \sqrt{\frac{L_2}{N}} \right) + \frac{L_2 M s^2 (\log n)^2}{\sqrt{N}} \right],$$

with probability at least  $1 - c' \exp(-c'' \log s)$ . Hence, under the given assumption that

$$\frac{L_2 M s^2 (\log n)^2}{\sqrt{N}} \rightarrow 0, \tag{20}$$

we have

$$\max_j \left| \sqrt{N} (\tilde{\beta}_j - \check{\beta}_j) \right| = o_{\mathbb{P}}(1). \tag{21}$$

On the other hand, we have

$$\max_j |e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j - e_j^T \Theta^* e_j| = \max_j |e_j^T (\Theta^* \Sigma_N - I_n) \Theta^* e_j| \leq \|\Theta^* \Sigma_N - I_n\|_\infty \|\Theta^*\|_\infty.$$

Therefore, by Lemma 8 and Assumption 5, we have

$$\mathbb{P} \left( \max_j |e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j - e_j^T \Theta^* e_j| \leq c \frac{M}{L_1} \sqrt{\frac{\log n}{N}} \right) \geq 1 - \frac{2}{n}.$$

Since (20) implies that

$$\frac{M}{L_1} \sqrt{\frac{\log n}{N}} \rightarrow 0,$$

so we have

$$\max_j |e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j - e_j^T \Theta^* e_j| = o_{\mathbb{P}}(1). \tag{22}$$

Hence, combining (21) and (22) yields the desired result in (18).

### 6.3.2 Step 2 of proof

By definition of  $\check{\beta}$ , we have

$$\frac{\sqrt{N}(\check{\beta}_j - \beta_j^*)}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} = \frac{\frac{1}{\sqrt{N}} e_j^T \Theta^* X^T W^{-1} \varepsilon}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} - \frac{\sqrt{N} e_j^T (\Theta^* \Sigma_N - I_n) (\hat{\beta} - \beta^*)}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}}.$$

First, let

$$Z = \frac{\frac{1}{\sqrt{N}} e_j^T \Theta^* X^T W^{-1} \varepsilon}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}}.$$

We show that  $Z$  is a standard normal random variable. Note that the characteristic function of  $Z$  is

$$\begin{aligned} \mathbb{E}(e^{itZ}) &= \mathbb{E} \left[ \exp \left( it \frac{\frac{1}{\sqrt{N}} e_j^T \Theta^* X^T W^{-1} \varepsilon}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} \right) \right] \\ &= \mathbb{E}_X \left[ \mathbb{E}_\varepsilon \left[ \exp \left( it \frac{\frac{1}{\sqrt{N}} e_j^T \Theta^* X^T W^{-1} \varepsilon}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} \right) \middle| X \right] \right] \\ &= \mathbb{E}_X \left[ \exp \left( -\frac{t^2}{2} \right) \middle| X \right] = \exp \left( -\frac{t^2}{2} \right). \end{aligned}$$

Therefore, we obtain the desired result. Next, let

$$\Psi_j = \frac{\sqrt{N} e_j^T (\Theta^* \Sigma_N - I_n) (\hat{\beta} - \beta^*)}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}}.$$

Our goal is to show that  $\max_j |\Psi_j| = o_{\mathbb{P}}(1)$  under given assumptions. We have

$$\max_j \left| \sqrt{N} e_j^T (\Theta^* \Sigma_N - I_n) (\hat{\beta} - \beta^*) \right| \leq \sqrt{N} \|\Theta^* \Sigma_N - I_n\|_\infty \|\hat{\beta} - \beta^*\|_1.$$

Therefore, by Lemma 8 and Corollary 1, we have

$$\max_j \left| \sqrt{N} e_j^T (\Theta^* \Sigma_N - I_n) (\hat{\beta} - \beta^*) \right| \leq c \frac{s \log n}{L_1} \sqrt{\frac{L_2}{N}}, \quad (23)$$

with probability at least  $1 - c' \exp(-c'' \log s)$ . Hence, combining (22) and (23), we have  $\max_j |\Psi_j| = o_{\mathbb{P}}(1)$  under the given high-dimensional regime. Finally, applying Slutsky's theorem yields the desired result in (19).

### 6.3.3 Conclusion of proof

We have

$$\sqrt{N}(\tilde{\beta}_j - \beta_j^*) = \sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j} \left[ \frac{\sqrt{N}(\check{\beta}_j - \beta_j^*)}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} + \frac{\sqrt{N}(\tilde{\beta}_j - \check{\beta}_j)}{\sqrt{e_j^T \Theta^* \Sigma_N (\Theta^*)^T e_j}} \right].$$

By (18), (19), and Slutsky's theorem, under given assumptions, we have

$$\frac{\sqrt{N}(\check{\beta}_j - \beta_j^*)}{\sqrt{e_j^T \Theta^* \Sigma_N(\Theta^*)^T e_j}} + \frac{\sqrt{N}(\tilde{\beta}_j - \check{\beta}_j)}{\sqrt{e_j^T \Theta^* \Sigma_N(\Theta^*)^T e_j}} \xrightarrow{d} N(0, 1). \quad (24)$$

Then combining (22) and (24), and applying Slutsky's theorem again, we have

$$\sqrt{N}(\tilde{\beta}_j - \beta_j^*) \xrightarrow{d} N(0, e_j^T \Theta^* e_j).$$

Therefore, the proof is complete.

## References

- [1] P. Bickel. Using residuals robustly I: tests for heteroscedasticity, nonlinearity. *The Annals of Statistics*, 6(2):266–291, 03 1978.
- [2] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. *Distributed Optimization and Statistical Learning via the Alternating Direction Method of Multipliers*. Now Publishers Inc, 2011.
- [3] R. Carroll. Adapting for heteroscedasticity in linear models. *Annals of Statistics*, 10(4):1224–1233, 12 1982.
- [4] R. Carroll and D. Ruppert. On robust tests for heteroscedasticity. *Annals of Statistics*, 9(1):206–210, 01 1981.
- [5] R. Carroll and D. Ruppert. A comparison between maximum likelihood and generalized least squares in a heteroscedastic linear model. *Journal of the American Statistical Association*, 77(380):878–882, 1982.
- [6] R. Carroll and D. Ruppert. Robust estimation in heteroscedastic linear models. *The Annals of Statistics*, 10(2):429–441, 06 1982.
- [7] R. Carroll and D. Ruppert. *Transformation and Weighting in Regression*, volume 30. CRC Press, 1988.
- [8] R. Carroll, J. Wu, and D. Ruppert. The effect of estimating weights in weighted least squares. *Journal of the American Statistical Association*, 83(404):1045–1054, 1988.
- [9] M. Cattaneo, M. Jansson, and W. Newey. Inference in linear regression models with many covariates and heteroscedasticity. *Journal of the American Statistical Association*, 113(523):1350–1361, 2018.
- [10] M. Davidian and R. Carroll. Variance function estimation. *Journal of the American Statistical Association*, 82(400):1079–1091, 1987.
- [11] Z. Daye, J. Chen, and H. Li. High-dimensional heteroscedastic regression with an application to eQTL data analysis. *Biometrics*, 68(1):316–326, 2012.

- [12] J. Janková and S. van de Geer. Semiparametric efficiency bounds for high-dimensional models. *Annals of Statistics*, 46(5):2336–2359, 10 2018.
- [13] A. Javanmard and A. Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research*, 15:2869–2909, 2014.
- [14] J. Jia, K. Rohe, and B. Yu. The Lasso under Poisson-like heteroscedasticity. *Statistica Sinica*, pages 99–118, 2013.
- [15] J. Jobson and W. Fuller. Least squares estimation when the covariance matrix and parameter vector are functionally related. *Journal of the American Statistical Association*, 75(369):176–181, 1980.
- [16] K. Knight and W. Fu. Asymptotics for Lasso-type estimators. *Annals of Statistics*, 28(5):1356–1378, 10 2000.
- [17] G. Raskutti, M. Wainwright, and B. Yu. Minimax rates of estimation for high-dimensional linear regression over  $\ell_q$ -balls. *IEEE transactions on information theory*, 57(10):6976–6994, 2011.
- [18] G. Seber and A. Lee. *Linear Regression Analysis*, volume 329. John Wiley & Sons, 2012.
- [19] J. Sharpnack and M. Kolar. Mean and variance estimation in high-dimensional heteroscedastic models with non-convex penalties. *arXiv preprint arXiv:1410.7874*, 2014.
- [20] S. van de Geer, P. Bühlmann, Y. Ritov, and R. Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 06 2014.
- [21] J. Wagener and H. Dette. The adaptive Lasso in high-dimensional sparse heteroscedastic models. *Mathematical Methods of Statistics*, 22(2):137–154, 2013.
- [22] M. Wainwright. Sharp thresholds for high-dimensional and noisy sparsity recovery using  $\ell_1$ -constrained quadratic programming (Lasso). *IEEE transactions on information theory*, 55(5):2183–2202, 2009.
- [23] D. Wang and P. Loh. Adaptive estimation and statistical inference for high-dimensional graph-based linear models. *arXiv preprint arXiv:2001.10679*, 2020.
- [24] C. Zhang and S. Zhang. Confidence intervals for low-dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.

## A Proofs

In this appendix, we provide proofs of various lemmas stated in Section 6 of the paper. Furthermore, see Figure A.1 for a flow diagram which illustrates the dependence structure among lemmas, theorems, and corollaries in the paper.

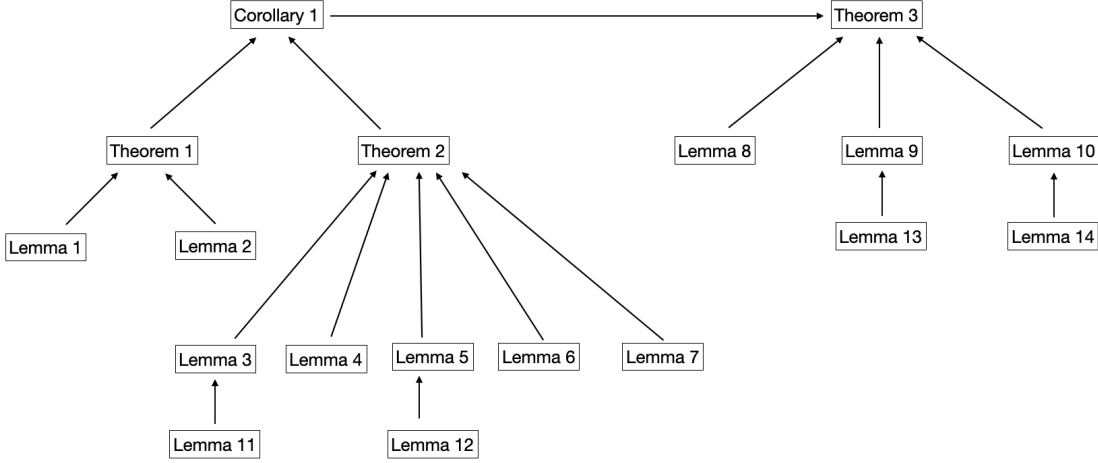


Figure A.1: Flowchart proofs.  $A \rightarrow B$  reads as “result B is partially dependent on result A”.

### A.1 Proof of Lemma 1

First, we write

$$\|\varepsilon^T X\|_\infty = \max_{i=1,\dots,n} |\varepsilon^T X e_i|.$$

Note that the conditional distribution of  $\varepsilon^T X e_i$  given  $X$  is

$$\varepsilon^T X e_i \mid X \sim N(0, e_i^T X^T W X e_i).$$

Furthermore, by our assumption on  $W$ , we have

$$e_i^T X^T W X e_i \leq L_2 \|X e_i\|_2^2. \quad (\text{A.1})$$

Applying Lemma 13 in [23], we have

$$\mathbb{P}(\|X e_i\|_2^2 \leq N(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)) \geq 1 - \exp(-Nc), \quad (\text{A.2})$$

where  $c > 0$  is a constant. Combining (A.1) and (A.2), we have

$$\mathbb{P}(e_i^T X^T W X e_i \leq L_2 N(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)) \geq 1 - \exp(-Nc). \quad (\text{A.3})$$

Next, for  $i = 1, \dots, n$ , we define the following event

$$\mathcal{E}_i = \{X \in \mathbb{R}^{N \times n} : e_i^T X^T W X e_i \leq L_2 N(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)\}.$$

Conditioning on  $\mathcal{E}_i$  and applying the standard tail bound of Gaussian random variable, for  $t > 0$ , we have

$$\mathbb{P}(|\varepsilon^T X e_i| > t \mid \mathcal{E}_i) \leq 2 \exp \left\{ -\frac{t^2}{2L_2 N(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)} \right\}. \quad (\text{A.4})$$

Therefore, combining (A.3) and (A.4), we have

$$\mathbb{P}(|\varepsilon^T X e_i| > t) \leq 2 \exp \left\{ -\frac{t^2}{2L_2 N(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)} \right\} + \exp(-Nc).$$



Finally, applying a union bound and letting  $t \geq 2(\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)^{1/2} L_2^{1/2} (N \log n)^{1/2}$  yield

$$\mathbb{P}(\|\varepsilon^T X\|_\infty > t) \leq 2 \exp(-\log n) + \exp(\log n - cN),$$

implying the desired result in the lemma.

## A.2 Proof of Lemma 2

First, we assume the following condition: Let  $\mathbb{L}_0(s) = \mathbb{B}_0(s) \cap \mathbb{B}_2(1)$ . For  $v \in \mathbb{L}_0(2s)$ , we have

$$\left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right| \leq \frac{\lambda_{\min}(\Sigma_x)}{150}. \quad (\text{A.5})$$

Note that for  $v \in \mathbb{R}^n$ , we have  $v^T \Sigma_x v \geq \lambda_{\min}(\Sigma_x) \|v\|_2^2$ . Therefore, when condition (A.5) holds, by Lemma 11 in [23], we have

$$\frac{1}{N} \|Xv\|_2^2 \geq \frac{\lambda_{\min}(\Sigma_x)}{2} \|v\|_2^2 \quad (\text{A.6})$$

for  $v \in \mathbb{R}^n$  and  $\|v_{S^c}\|_1 \leq 3\|v_S\|_1$ .

Next, we use a discretization argument to verify that condition (A.5) holds with high probability. For an index subset  $I \subset [n]$  with  $|I| \leq 2s$ , we define

$$S_I = \{v \in \mathbb{R}^n; \|v\|_2 \leq 1, \text{support}(v) \subset I\}.$$

Then we have  $\mathbb{L}_0(2s) = \cup_{|I| \leq 2s} S_I$ . For a fixed  $S_I$ , let  $\mathcal{A}_I$  be a  $\frac{1}{3}$ -cover of  $S_I$  with  $|\mathcal{A}_I| \leq 9^{2s}$ . Therefore, for  $v \in S_I$ , there exists  $a_v \in \mathcal{A}_I$ , such that  $\|\Delta_v\|_2 = \|v - a_v\|_2 \leq \frac{1}{3}$ . Hence,

$$\begin{aligned} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right| &= \left| (\Delta_v + a_v)^T \left( \frac{X^T X}{N} - \Sigma_x \right) (\Delta_v + a_v) \right| \\ &\leq \left| \Delta_v^T \left( \frac{X^T X}{N} - \Sigma_x \right) \Delta_v \right| + 2 \left| a_v^T \left( \frac{X^T X}{N} - \Sigma_x \right) \Delta_v \right| + \left| a_v^T \left( \frac{X^T X}{N} - \Sigma_x \right) a_v \right| \\ &\leq \frac{1}{9} \sup_{v \in S_I} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right| + \frac{2}{3} \sup_{v \in S_I} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right| + \sup_{v \in \mathcal{A}_I} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right|. \end{aligned}$$

Rearranging terms, we have

$$\sup_{v \in S_I} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right| \leq \frac{9}{2} \sup_{v \in \mathcal{A}_I} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right|.$$

Applying Lemma 13 in [23] and taking union bounds, we have

$$\mathbb{P} \left( \sup_{v \in \mathbb{L}_0(2s)} \left| v^T \left( \frac{X^T X}{N} - \Sigma_x \right) v \right| \geq \frac{\lambda_{\min}(\Sigma_x)}{150} \right) \leq 2 \exp(cs \log n - c'N),$$

where  $c > 0$  and  $c' > 0$  are constants.

Therefore, we obtain that inequality (A.6) holds at least with the same probability.

### A.3 Proof of Lemma 3

First, we show Part (a). By Lemma 11, it suffices to show that  $\|\hat{z}_{S^c}\|_\infty < 1$ . Note that we have

$$\hat{z}_{S^c} = X_{S^c}^T X_S (X_S^T X_S)^{-1} \hat{z}_S + \frac{1}{\lambda N} X_{S^c}^T (I_N - X_S (X_S^T X_S)^{-1} X_S^T) \varepsilon,$$

which implies

$$\|\hat{z}_{S^c}\|_\infty \leq \|X_{S^c}^T X_S (X_S^T X_S)^{-1} \hat{z}_S\|_\infty + \left\| \frac{1}{\lambda N} X_{S^c}^T (I_N - X_S (X_S^T X_S)^{-1} X_S^T) \varepsilon \right\|_\infty.$$

Therefore, it suffices to bound  $\|X_{S^c}^T X_S (X_S^T X_S)^{-1} \hat{z}_S\|_\infty$  and  $\left\| \frac{1}{\lambda N} X_{S^c}^T (I_N - X_S (X_S^T X_S)^{-1} X_S^T) \varepsilon \right\|_\infty$ . For the first term, we have

$$\begin{aligned} \|X_{S^c}^T X_S (X_S^T X_S)^{-1} \hat{z}_S\|_\infty &= \max_j |e_j^T X_{S^c}^T X_S (X_S^T X_S)^{-1} \hat{z}_S| \\ &\leq \max_j \|e_j^T X_{S^c}^T X_S (X_S^T X_S)^{-1}\|_1 \|\hat{z}_S\|_\infty \leq \alpha, \end{aligned}$$

where the last inequality follows from (15) and  $\|\hat{z}_S\|_\infty \leq 1$ . For the second term, we have

$$\left\| \frac{1}{\lambda N} X_{S^c}^T (I_N - X_S (X_S^T X_S)^{-1} X_S^T) \varepsilon \right\|_\infty \leq \frac{1 - \alpha}{2},$$

where the inequality follows from (16). Therefore, we have

$$\|\hat{z}_{S^c}\|_\infty \leq \alpha + \frac{1 - \alpha}{2} = \frac{1 + \alpha}{2} < 1.$$

Hence, Lemma 11 implies that there is a unique solution to (3) and its support  $\hat{S}$  is contained in  $S$ .

Next, we show Part (b). It remains to show that the true support  $S$  is contained in the set  $\hat{S}$ . Note that we have

$$\hat{\beta}_S - \beta_S^* = (X_S^T X_S)^{-1} X_S^T \varepsilon - \lambda N (X_S^T X_S)^{-1} \hat{z}_S.$$

Therefore, we have

$$\begin{aligned} \|\hat{\beta}_S - \beta_S^*\|_\infty &\leq \|(X_S^T X_S)^{-1} X_S^T \varepsilon\|_\infty + \|\lambda N (X_S^T X_S)^{-1} \hat{z}_S\|_\infty \\ &\leq \left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} X_S^T \frac{\varepsilon}{N} \right\|_\infty + \lambda \left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} \right\|_\infty = B(\lambda, X_S, \varepsilon). \end{aligned}$$

By (17), we have  $\hat{\beta}_i \neq 0$  for  $i \in S$ , which implies  $S \subseteq \hat{S}$ . Putting together the pieces, we conclude that  $S = \hat{S}$ .

Therefore, the proof is complete.

#### A.4 Proof of Lemma 4

Since  $X \in \mathbb{R}^{N \times n}$  is a row-wise  $(\sigma_x, \Sigma_x)$ -sub-Gaussian random matrix, so  $X_S \in \mathbb{R}^{N \times s}$  is also a row-wise  $(\sigma_x, (\Sigma_x)_{SS})$ -sub-Gaussian random matrix, where  $(\Sigma_x)_{SS} \in \mathbb{R}^{s \times s}$  is a submatrix of  $\Sigma_x$  with rows and columns are restricted to S. By Lemma 13 in [23], we have

$$\mathbb{P} \left( \lambda_{\min} \left( \frac{1}{N} X_S^T X_S \right) \geq \frac{1}{2} \lambda_{\min}((\Sigma_x)_{SS}) \right) \geq 1 - 2 \exp(cs - c'N).$$

Furthermore, we have  $\lambda_{\min}((\Sigma_x)_{SS}) \geq \lambda_{\min}(\Sigma_x)$ . Hence, we have

$$\lambda_{\min} \left( \frac{1}{N} X_S^T X_S \right) \geq \frac{1}{2} \lambda_{\min}(\Sigma_x),$$

with probability at least  $1 - 2 \exp(-cN)$  assuming the scaling  $N \gtrsim s$ . Therefore, the proof is complete.

#### A.5 Proof of Lemma 5

Let  $\hat{\Gamma} = \frac{1}{N} X^T X$ . Then we have

$$\begin{aligned} \left\| \hat{\Gamma}_{S^c S} \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{S^c S} (\Sigma_x)_{SS}^{-1} \right\|_{\infty} &\leq \left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty} \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty} \\ &\quad + \left\| \left[ \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right] (\Sigma_x)_{SS}^{-1} \right\|_{\infty} + \left\| (\Sigma_x)_{SS} \left[ \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right] \right\|_{\infty} \\ &\leq \left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty} \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty} \\ &\quad + \left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty} \left\| (\Sigma_x)_{SS}^{-1} \right\|_{\infty} + \left\| (\Sigma_x)_{SS} \right\|_{\infty} \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty}. \end{aligned}$$

Therefore, it suffices to bound  $\left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty}$ ,  $\left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty}$ ,  $\left\| (\Sigma_x)_{SS}^{-1} \right\|_{\infty}$ , and  $\left\| (\Sigma_x)_{SS} \right\|_{\infty}$ .

First, we bound the term  $\left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty}$ . We have

$$\begin{aligned} \left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty} &= \max_{j \in S^c} \left\| \hat{\Gamma}_{jS} - (\Sigma_x)_{jS} \right\|_1 \\ &\leq s \times \max_{j \in S^c} \left\| \hat{\Gamma}_{jS} - (\Sigma_x)_{jS} \right\|_{\infty} = s \times \max_{j \in S^c, k \in S} \left| \hat{\Gamma}_{jk} - (\Sigma_x)_{jk} \right|. \end{aligned}$$

By Part (a) of Lemma 12, we have

$$\mathbb{P} \left( \left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_{\infty} \leq cs \sqrt{\frac{\log[s(n-s)]}{N}} \right) \geq 1 - c' \exp(-c'' \log[s(n-s)]).$$

Next, we bound the second term  $\left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty}$ . We have

$$\begin{aligned} \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty} &= \max_{j \in S} \left\| \left( \hat{\Gamma}_{SS}^{-1} \right)_j - \left( (\Sigma_x)_{SS}^{-1} \right)_j \right\|_1 \\ &\leq s \times \max_{j \in S} \left\| \left( \hat{\Gamma}_{SS}^{-1} \right)_j - \left( (\Sigma_x)_{SS}^{-1} \right)_j \right\|_{\infty} = s \times \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty}, \end{aligned}$$

where  $(\hat{\Gamma}_{SS}^{-1})_j$  and  $(\Sigma_x)_{SS}^{-1}$  are  $j$ th row of  $\hat{\Gamma}_{SS}^{-1}$  and  $(\Sigma_x)_{SS}^{-1}$ , respectively. Therefore, by Part (b) of Lemma 12, we have

$$\mathbb{P}\left(\left\|\hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1}\right\|_{\infty} \leq cs\sqrt{\frac{s}{N}}\right) \geq 1 - c' \exp(-c''s).$$

Furthermore, for the last two terms, we have

$$\left\|(\Sigma_x)_{SS}^{-1}\right\|_{\infty} \leq s \left\|(\Sigma_x)_{SS}^{-1}\right\|_{\infty} \leq s \left\|(\Sigma_x)_{SS}^{-1}\right\|_{op} \leq \frac{s}{\lambda_{\min}(\Sigma_x)},$$

and

$$\left\|(\Sigma_x)_{SS}\right\|_{\infty} \leq s \left\|(\Sigma_x)_{SS}\right\|_{\infty} \leq s \left\|(\Sigma_x)_{SS}\right\|_{op} \leq \lambda_{\max}(\Sigma_x)s.$$

Putting together the pieces, we conclude that

$$\left\|\hat{\Gamma}_{S^cS}\hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{S^cS}(\Sigma_x)_{SS}^{-1}\right\|_{\infty} \leq c \left( \sqrt{\frac{s^5 \log[s(n-s)]}{N^2}} + s^2 \sqrt{\frac{\log[s(n-s)]}{N}} + s^2 \sqrt{\frac{s}{N}} \right),$$

with probability at least  $1 - c' \exp(-c''s) - C' \exp(-C'' \log[s(n-s)])$ . Therefore, when

$$N^2 \gtrsim \frac{s^5 \log[s(n-s)]}{\alpha^2}, \quad N \gtrsim \frac{s^4 \log[s(n-s)]}{\alpha^2}, \quad \text{and} \quad N \gtrsim \frac{s^5}{\alpha^2},$$

we have

$$\left\|\hat{\Gamma}_{S^cS}\hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{S^cS}(\Sigma_x)_{SS}^{-1}\right\|_{\infty} \leq \frac{\alpha}{2},$$

with at least the same probability. Hence, we have

$$\left\|\hat{\Gamma}_{S^cS}\hat{\Gamma}_{SS}^{-1}\right\|_{\infty} \leq \left\|(\Sigma_x)_{S^cS}(\Sigma_x)_{SS}^{-1}\right\|_{\infty} + \left\|\hat{\Gamma}_{S^cS}\hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{S^cS}(\Sigma_x)_{SS}^{-1}\right\|_{\infty} \leq \frac{\alpha}{2} + \frac{\alpha}{2} = \alpha,$$

with at least the same probability. Therefore, the proof is complete.

## A.6 Proof of Lemma 6

First, we write

$$\left\|X_{S^c}^T (I_N - X_S(X_S^T X_S)^{-1} X_S^T) \frac{\varepsilon}{N}\right\|_{\infty} = \max_{1 \leq j \leq n-s} \left| e_j^T X_{S^c}^T (I_N - X_S(X_S^T X_S)^{-1} X_S^T) \frac{\varepsilon}{N} \right|.$$

Let  $\Pi_{S^{\perp}}(X) = I_N - X_S(X_S^T X_S)^{-1} X_S^T$ . Then we have

$$e_j^T X_{S^c}^T \Pi_{S^{\perp}}(X) \frac{\varepsilon}{N} \mid X \sim N\left(0, \frac{1}{N^2} e_j^T X_{S^c}^T \Pi_{S^{\perp}}(X) W \Pi_{S^{\perp}}(X) X_{S^c} e_j\right).$$

By our assumption on  $W$  and the fact that  $\Pi_{S^{\perp}}(X)$  is an orthogonal projection matrix, we have

$$\frac{1}{N^2} e_j^T X_{S^c}^T \Pi_{S^{\perp}}(X) W \Pi_{S^{\perp}}(X) X_{S^c} e_j \leq \frac{1}{N^2} \|W\|_{op} \|\Pi_{S^{\perp}}(X)\|_{op}^2 \|X_{S^c} e_j\|_2^2 \leq \frac{1}{N^2} L_2 \|X_{S^c} e_j\|_2^2.$$

Applying Lemma 13 in [23], we have

$$\mathbb{P} \left( \frac{L_2}{N^2} \|X_{S^c} e_j\|_2^2 \leq \frac{L_2}{N} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2) \right) \geq 1 - \exp(-Nc),$$

where  $c > 0$  is a constant.

Next, for  $j = 1, \dots, n - s$ , we define the following event

$$\mathcal{E}_j = \left\{ X \in \mathbb{R}^{N \times n} : \frac{1}{N^2} e_j^T X_{S^c}^T \Pi_{S^\perp}(X) W \Pi_{S^\perp}(X) X_{S^c} e_j \leq \frac{L_2}{N} (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2) \right\}.$$

Conditioning on  $\mathcal{E}_j$  and applying the standard tail bound of Gaussian random variable, for  $t > 0$ , we have

$$\mathbb{P} \left( \left| e_j^T X_{S^c}^T \Pi_{S^\perp}(X) \frac{\varepsilon}{N} \right| > t \mid \mathcal{E}_j \right) \leq 2 \exp \left\{ -\frac{Nt^2}{2L_2 (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)} \right\}.$$

Therefore, we have

$$\mathbb{P} \left( \left| e_j^T X_{S^c}^T \Pi_{S^\perp}(X) \frac{\varepsilon}{N} \right| > t \right) \leq 2 \exp \left\{ -\frac{Nt^2}{2L_2 (\lambda_{\max}(\Sigma_x) + 4\sigma_x^2)} \right\} + \exp(-Nc).$$

Finally, applying a union bound and letting  $t \geq 2\sqrt{\lambda_{\max}(\Sigma_x) + 4\sigma_x^2} \sqrt{\frac{L_2 \log(n-s)}{N}}$  yields the desired result in the lemma. Therefore, the proof is complete.

## A.7 Proof of Lemma 7

It suffices to bound  $\left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} X_S^T \frac{\varepsilon}{N} \right\|_\infty$  and  $\lambda \left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} \right\|_\infty$ . We use a similar argument with Lemma 6 to bound the first term, and conclude that

$$\mathbb{P} \left( \left\| \hat{\Gamma}_{SS}^{-1} X_S^T \frac{\varepsilon}{N} \right\|_\infty > \sqrt{\frac{8}{\lambda_{\min}(\Sigma_x)}} \sqrt{\frac{L_2 \log s}{N}} \right) \leq 2 \exp(-\log s) + 2 \exp(\log s + cs - c'N).$$

For the second term, we use a similar argument with Lemma 5, and conclude that

$$\lambda \left\| \left( \frac{X_S^T X_S}{N} \right)^{-1} \right\|_\infty \leq \lambda \left( \left\| (\Sigma_x)_{SS}^{-1} \right\|_\infty + cs \sqrt{\frac{s}{N}} \right),$$

with probability greater than  $1 - 2 \exp(-c's) - 2 \exp(c''s - C'N)$ . Therefore, the proof is complete.

## A.8 Proof of Lemma 8

Let  $\Gamma = \Theta^* \Sigma_N - I_n$ . Then the  $(j, k)$ -th entry of  $\Gamma$  is

$$\Gamma^{jk} = \frac{1}{N} \sum_{i=1}^N e_j^T \left( \Theta^* \frac{X_i X_i^T}{W_i} - I_n \right) e_k = \frac{1}{N} \sum_{i=1}^N \Gamma_i^{jk}.$$

Note that  $\mathbb{E}(\Gamma_i^{jk}) = 0$ , and

$$\|\Gamma_i^{jk}\|_{\psi_1} \leq 2 \left\| e_j^T \Theta^* \frac{X_i X_i^T}{W_i} e_k \right\|_{\psi_1} \leq \frac{4}{L_1} \|e_j^T \Theta^* X_i\|_{\psi_2} \|X_i^T e_k\|_{\psi_2} \leq \frac{4c}{L_1} \lambda_{\max}(\Theta^*) \sigma_x^2 = K.$$

Applying Lemma 12 in [23], we have for  $t > 0$ ,

$$\mathbb{P}\left(\left|\Gamma^{jk}\right| \geq t\right) = \mathbb{P}\left(\left|\frac{1}{N} \sum_{i=1}^N \Gamma_i^{jk}\right| \geq t\right) \leq 2 \exp\left[-C_b N \min\left(\frac{t^2}{K^2}, \frac{t}{K}\right)\right].$$

Taking a union bound, we have

$$\|\Gamma\|_\infty \leq c \frac{\lambda_{\max}(\Theta^*) \sigma_x^2}{L_1} \sqrt{\frac{\log n}{N}}, \quad (\text{A.7})$$

with probability at least  $1 - \frac{2}{n}$ . Therefore, the proof is complete.

## A.9 Proof of Lemma 9

Let  $\hat{\Theta}_j$  and  $\Theta_j^*$  be the  $j$ th row of  $\hat{\Theta}$  and  $\Theta^*$ , respectively. First, we show that on the event

$$\mathcal{E} = \left\{ \left\| \hat{\Theta}_j \right\|_1 \leq \left\| \Theta_j^* \right\|_1, 1 \leq j \leq n \right\},$$

we have

$$\left\| \hat{\Theta} - \Theta^* \right\|_\infty \leq 8k \left\| \hat{\Theta} - \Theta^* \right\|_\infty, \quad (\text{A.8})$$

where  $k$  is defined in Assumption 5. Let  $\xi = \left\| \hat{\Theta} - \Theta^* \right\|_\infty$  and  $r_j = \hat{\Theta}_j - \Theta_j^*$ . Then we have

$$r_j = \left[ \left( \hat{\Theta}_{ji} 1_{\{|\hat{\Theta}_{ji}| \geq 2\xi\}}; 1 \leq i \leq n \right) - \Theta_j^* \right] + \left( \hat{\Theta}_{ji} 1_{\{|\hat{\Theta}_{ji}| < 2\xi\}}; 1 \leq i \leq n \right) = r_j^1 + r_j^2.$$

Furthermore, on the event  $\mathcal{E}$ , we have

$$\left\| \Theta_j^* \right\|_1 - \left\| r_j^1 \right\|_1 + \left\| r_j^2 \right\|_1 \leq \left\| \Theta_j^* + r_j^1 \right\|_1 + \left\| r_j^2 \right\|_1 = \left\| \hat{\Theta}_j \right\|_1 \leq \left\| \Theta_j^* \right\|_1,$$

where the first inequality follows from the triangle inequality. Therefore we have

$$\left\| r_j \right\|_1 \leq \left\| r_j^1 \right\|_1 + \left\| r_j^2 \right\|_1 \leq 2 \left\| r_j^1 \right\|_1. \quad (\text{A.9})$$

Furthermore, we have

$$\begin{aligned} \left\| r_j^1 \right\|_1 &\leq \sum_{i=1}^n |\hat{\Theta}_{ji} - \Theta_{ji}^*| 1_{\{|\hat{\Theta}_{ji}| \geq 2\xi\}} + \sum_{i=1}^n |\Theta_{ji}^*| 1_{\{|\hat{\Theta}_{ji}| < 2\xi\}} \\ &\leq \sum_{i=1}^n \xi 1_{\{|\Theta_{ji}^*| \geq \xi\}} + \sum_{i=1}^n |\Theta_{ji}^*| 1_{\{|\Theta_{ji}^*| < 3\xi\}} \leq \xi k + 3\xi k = 4\xi k, \end{aligned} \quad (\text{A.10})$$

where the second inequality follows from the triangle inequality. Therefore, combining (A.9) and (A.9) yields the result in (A.8).

Next, we have

$$\begin{aligned} \left\| \hat{\Theta} - \Theta^* \right\|_\infty &= \left\| \hat{\Theta}(I_n - \hat{\Sigma}_N \Theta^*) + (\hat{\Theta} \hat{\Sigma}_N - I_n) \Theta^* \right\|_\infty \\ &\leq \left\| \hat{\Theta} \right\|_\infty \left\| I_n - \hat{\Sigma}_N \Theta^* \right\|_\infty + \left\| \hat{\Theta} \hat{\Sigma}_N - I_n \right\|_\infty \left\| \Theta^* \right\|_1 \\ &\leq \left\| \hat{\Theta} \right\|_\infty \left\| I_n - \hat{\Sigma}_N \Theta^* \right\|_\infty + \mu M. \end{aligned}$$

We define the event

$$\mathcal{E}' = \left\{ X : \left\| \Theta^* \hat{\Sigma}_N - I_n \right\|_\infty \leq \mu \right\}.$$

Then on the event  $\mathcal{E}'$ , we have  $\|\hat{\Theta}_j\|_1 \leq \|\Theta_j^*\|_1$  for  $1 \leq j \leq n$ , which implies  $\|\hat{\Theta}\|_\infty \leq \|\Theta^*\|_\infty$ . Therefore on the event  $\mathcal{E}'$ , we have

$$\|\hat{\Theta} - \Theta^*\|_\infty \leq 2\mu M.$$

Furthermore, since event  $\mathcal{E}'$  implies event  $\mathcal{E}$ , so on the event  $\mathcal{E}'$ , we have

$$\|\hat{\Theta} - \Theta^*\|_\infty \leq 8k \|\hat{\Theta} - \Theta^*\|_\infty \leq 16k\mu M.$$

Finally, letting

$$\mu = c \left( \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}} \right)$$

and applying Lemma 13, we have

$$\|\hat{\Theta} - \Theta^*\|_\infty \leq ckM \left( \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}} \right),$$

with probability at least  $1 - c' \exp(-c'' \log s)$ . Therefore, the proof is complete.

## A.10 Proof of Lemma 10

First, we show Part (a). We have

$$\left\| \frac{1}{\sqrt{N}} X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right\|_\infty = \max_{1 \leq j \leq n} \left| \frac{1}{\sqrt{N}} e_j^T X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right| = \max_{1 \leq j \leq n} \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\widehat{W}_i^{-1} - W_i^{-1}) \right|.$$

For a fixed  $j$ , using Taylor's theorem for  $\widehat{W}_i^{-1} - W_i^{-1}$ , we have

$$\begin{aligned} \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\widehat{W}_i^{-1} - W_i^{-1}) &= \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\hat{\beta} - \beta^*)^T d_i(\beta^*) \\ &+ \frac{1}{2\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\hat{\beta} - \beta^*)^T H_i(\beta^* + c_i(\hat{\beta} - \beta^*)) (\hat{\beta} - \beta^*), \end{aligned}$$

where  $d_i(\beta)$ ,  $H_i(\beta)$ , and  $c_i$  are defined in the proof of Lemma 13. Therefore, we obtain

$$\begin{aligned} \max_j \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\widehat{W}_i^{-1} - W_i^{-1}) \right| &\leq \max_j \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\hat{\beta} - \beta^*)^T d_i(\beta^*) \right| \\ &+ \max_j \left| \frac{1}{2\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\hat{\beta} - \beta^*)^T H_i(\beta^* + c_i(\hat{\beta} - \beta^*)) (\hat{\beta} - \beta^*) \right| \end{aligned}$$

Let  $I_j = \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\hat{\beta} - \beta^*)^T d_i(\beta^*)$ . We have

$$\max_j |I_j| \leq \|\hat{\beta} - \beta^*\|_1 \max_j \left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i d_i(\beta^*) \right\|_\infty.$$

By Corollary 1 and Part (a) of Lemma 14, we have

$$\mathbb{P} \left( \max_j |I_j| \leq cs \log n \sqrt{\frac{L_2}{N}} \right) \geq 1 - c' \exp(-c'' \log s). \quad (\text{A.11})$$

Next, let  $II_j = \frac{1}{2\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i (\hat{\beta} - \beta^*)^T H_i(\beta^* + c_i(\hat{\beta} - \beta^*)) (\hat{\beta} - \beta^*)$ . We have

$$\max_j |II_j| \leq \frac{1}{2} \left( \max_j \frac{1}{\sqrt{N}} \sum_{i=1}^N |X_{ij} \varepsilon_i| \right) \max_i \left| (\hat{\beta} - \beta^*)^T H_i(\beta^* + c_i(\hat{\beta} - \beta^*)) (\hat{\beta} - \beta^*) \right|.$$

If

$$\hat{S} \subseteq S \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_2 \leq c \sqrt{\frac{s L_2 \log n}{N}},$$

then we have

$$\begin{aligned} \max_j |II_j| &\leq \frac{1}{2} \left( \max_j \frac{1}{\sqrt{N}} \sum_{i=1}^N |X_{ij} \varepsilon_i| \right) \|\hat{\beta} - \beta^*\|_2^2 \|H_i(\beta^* + c_i(\hat{\beta} - \beta^*))\|_{SS} \\ &\leq \frac{1}{2} \left( \max_j \frac{1}{\sqrt{N}} \sum_{i=1}^N |X_{ij} \varepsilon_i| \right) \|\hat{\beta} - \beta^*\|_2^2 \|H_i(\beta^* + c_i(\hat{\beta} - \beta^*))\|_{SS} \\ &\leq \frac{c}{2} \left( \max_j \frac{1}{\sqrt{N}} \sum_{i=1}^N |X_{ij} \varepsilon_i| \right) \frac{s^2 L_2 \log n}{N}, \end{aligned}$$

where the last inequality follows from Part (b) of Assumption 2 and Assumption 4. Therefore, by Part (b) of Lemma 14 and Corollary 1, we have

$$\mathbb{P} \left( \max_j |II_j| \leq c \frac{L_2 s^2 (\log n)^2}{\sqrt{N}} \right) \geq 1 - c' \exp(-c'' \log s). \quad (\text{A.12})$$

Finally, combining (A.11) and (A.12), we have

$$\left\| \frac{1}{\sqrt{N}} X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon \right\|_\infty \leq c \frac{L_2 s^2 (\log n)^2}{\sqrt{N}},$$

with probability at least  $1 - c' \exp(-c'' \log s)$ .

Next, we show Part (b). We have

$$\frac{1}{\sqrt{N}} X^T \widehat{W}^{-1} \varepsilon = \frac{1}{\sqrt{N}} X^T (\widehat{W}^{-1} - W^{-1}) \varepsilon + \frac{1}{\sqrt{N}} X^T W^{-1} \varepsilon,$$



which implies that

$$\left\| \frac{1}{\sqrt{N}} X^T \widehat{W}^{-1} \varepsilon \right\|_{\infty} \leq \left\| \frac{1}{\sqrt{N}} X^T \left( \widehat{W}^{-1} - W^{-1} \right) \varepsilon \right\|_{\infty} + \left\| \frac{1}{\sqrt{N}} X^T W^{-1} \varepsilon \right\|_{\infty}.$$

Therefore, by Part (a) of this lemma and Part (c) of Lemma 14, we have

$$\left\| \frac{1}{\sqrt{N}} X^T \widehat{W}^{-1} \varepsilon \right\|_{\infty} \leq c \left( \frac{L_2 s^2 (\log n)^2}{\sqrt{N}} + \sqrt{\log n} \right),$$

with probability at least  $1 - c' \exp(-c'' \log s)$ . Therefore, the proof is complete.

## B Primal-dual witness construction

In this appendix, we introduce the primal-dual witness construction [22], which is used for establishing the variable selection consistency. PDW includes the following three steps:

- (1) Set  $\hat{\beta}_{S^c} = 0$ .
- (2) Determine  $\hat{\beta}_S \in \mathbb{R}^s$  by solving

$$\hat{\beta}_S \in \operatorname{argmin}_{\beta_S \in \mathbb{R}^s} \frac{1}{2N} \|y - X_S \beta_S\|_2^2 + \lambda \|\beta_S\|_1. \quad (\text{B.1})$$

Then choose  $\hat{z}_S \in \partial \|\hat{\beta}_S\|_1$  such that

$$\frac{1}{N} X_S^T (X_S \hat{\beta}_S - y) + \lambda \hat{z}_S = 0.$$

- (3) Solve for  $\hat{z}_{S^c} \in \mathbb{R}^{n-s}$  via the equation

$$\frac{1}{N} X^T (X \hat{\beta} - y) + \lambda \hat{z} = 0,$$

and check whether or not the condition  $\|\hat{z}_{S^c}\|_{\infty} < 1$  holds.

The PDW construction has the following important consequence:

**Lemma 11.** If there exists a constant  $c_{\min} > 0$  such that  $\lambda_{\min}(\frac{1}{N} X_S^T X_S) \geq c_{\min}$ , then the success of the PDW construction implies that  $(\hat{\beta}_S^T, 0)^T \in \mathbb{R}^n$  is the unique optimal solution to (3).

*Proof.* Note if the PDW construction succeeds, then by zero subgradient condition,  $\hat{\beta} = (\hat{\beta}_S, 0)$  is an optimal solution to (3) with associated subgradient vector  $\hat{z}$  satisfying  $\|\hat{z}_{S^c}\|_{\infty} < 1$  and  $\hat{z}^T \hat{\beta} = \|\hat{\beta}\|_1$ .

Let  $\bar{\beta}$  be any other optimal solution to (3). Then we have

$$\frac{1}{2N} \|y - X \hat{\beta}\|_2^2 + \lambda \hat{z}^T \hat{\beta} = \frac{1}{2N} \|y - X \bar{\beta}\|_2^2 + \lambda \|\bar{\beta}\|_1.$$

Therefore, we have

$$\frac{1}{2N} \|y - X \hat{\beta}\|_2^2 - \lambda \hat{z}^T (\bar{\beta} - \hat{\beta}) = \frac{1}{2N} \|y - X \bar{\beta}\|_2^2 + \lambda (\|\bar{\beta}\|_1 - \hat{z}^T \bar{\beta}).$$

Furthermore, by zero subgradient condition, we have

$$\lambda \hat{z} = -\frac{1}{N} X^T (X \hat{\beta} - y),$$

which implies that

$$\frac{1}{2N} \|y - X \hat{\beta}\|_2^2 + \frac{1}{N} (X \hat{\beta} - y)^T X (\bar{\beta} - \hat{\beta}) - \frac{1}{2N} \|y - X \bar{\beta}\|_2^2 = \lambda (\|\bar{\beta}\|_1 - \hat{z}^T \bar{\beta}) \leq 0.$$

Therefore, we have  $\|\bar{\beta}\|_1 \leq \hat{z}^T \bar{\beta}$ . On the other hand, we have

$$\hat{z}^T \bar{\beta} \leq \|\hat{z}\|_\infty \|\bar{\beta}\|_1 \leq \|\bar{\beta}\|_1.$$

Combining above arguments, we have  $\|\bar{\beta}\|_1 = \hat{z}^T \bar{\beta}$ . Since  $\|\hat{z}_{S^c}\|_\infty < 1$ , so we have  $\bar{\beta}_j = 0$  for  $j \in S^c$ . Therefore, all optimal solution are supported only on the subset  $S$ , and hence can be obtained by solving the problem (B.1). Given the condition stated in the lemma, (B.1) is strictly convex, and so has a unique minimizer. Therefore, the proof is complete.  $\square$

## C Additional supporting lemmas

**Lemma 12.** Assume that Part (a) of Assumption 2 holds. Let  $\hat{\Gamma} = \frac{1}{N} X^T X$ .

(a) We have

$$\left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_\infty \leq c \sqrt{\frac{\log[s(n-s)]}{N}},$$

with probability at least  $1 - c' \exp(-c'' \log[s(n-s)])$ .

(b) Under the sample size scaling  $N \gtrsim s$ , we have

$$\left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_\infty \leq \frac{2c}{\lambda_{\min}^2(\Sigma_x)} \sqrt{\frac{s}{N}},$$

with probability at least  $1 - c' \exp(-c'' s)$ .

*Proof.* First, we show Part (a). For a fixed  $j \in S^c$  and a fixed  $k \in S$ , we have

$$\left| \hat{\Gamma}_{jk} - (\Sigma_x)_{jk} \right| = \left| \frac{1}{N} \sum_{i=1}^N [X_{ij} X_{ik} - (\Sigma_x)_{jk}] \right|,$$

and

$$\|X_{ij} X_{ik} - (\Sigma_x)_{jk}\|_{\psi_1} \leq 2 \|X_{ij} X_{ik}\|_{\psi_1} \leq 4 \|X_{ij}\|_{\psi_2} \|X_{ik}\|_{\psi_2} \leq 4 c \sigma_x^2.$$

Then by Lemma 12 in [23], for  $0 < t \leq K = 4c\sigma_x^2$ , we have

$$\mathbb{P} \left( \left| \hat{\Gamma}_{jk} - (\Sigma_x)_{jk} \right| \geq t \right) \leq 2 \exp \left( -C_b N \frac{t^2}{K^2} \right).$$

Applying a union bound, we further have

$$\mathbb{P} \left( \left\| \hat{\Gamma}_{S^c S} - (\Sigma_x)_{S^c S} \right\|_\infty \geq t \right) \leq 2s(n-s) \exp \left( -C_b N \frac{t^2}{K^2} \right).$$

Therefore, letting  $t = c\sqrt{\frac{\log[s(n-s)]}{N}}$  yields the desired result in Part (a).

Next, we show Part (b). Since

$$\hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} = \hat{\Gamma}_{SS}^{-1} \left[ (\Sigma_x)_{SS} - \hat{\Gamma}_{SS} \right] (\Sigma_x)_{SS}^{-1},$$

so we have

$$\begin{aligned} \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty} &= \max_{j,k} \left| e_j^T \hat{\Gamma}_{SS}^{-1} \left( (\Sigma_x)_{SS} - \hat{\Gamma}_{SS} \right) (\Sigma_x)_{SS}^{-1} e_k \right| \\ &\leq \left\| \hat{\Gamma}_{SS}^{-1} \right\|_{op} \left\| (\Sigma_x)_{SS} - \hat{\Gamma}_{SS} \right\|_{op} \left\| (\Sigma_x)_{SS}^{-1} \right\|_{op} = \frac{1}{\lambda_{\min}(\hat{\Gamma}_{SS})} \left\| (\Sigma_x)_{SS} - \hat{\Gamma}_{SS} \right\|_{op} \frac{1}{\lambda_{\min}((\Sigma_x)_{SS})}. \end{aligned}$$

By Lemma 13 in [23], we have

$$\mathbb{P} \left( \left\| (\Sigma_x)_{SS} - \hat{\Gamma}_{SS} \right\|_{op} \geq c\sqrt{\frac{s}{N}} \right) \leq 2 \exp(-c's).$$

Then combining Lemma 4 and the fact that  $\lambda_{\min}((\Sigma_x)_{SS}) \geq \lambda_{\min}(\Sigma_x)$ , we have

$$\mathbb{P} \left( \left\| \hat{\Gamma}_{SS}^{-1} - (\Sigma_x)_{SS}^{-1} \right\|_{\infty} \leq \frac{2c}{\lambda_{\min}^2(\Sigma_x)} \sqrt{\frac{s}{N}} \right) \geq 1 - c' \exp(-c''s).$$

Hence, we obtain the result in Part (b). □

**Lemma 13.** Assume that same conditions as Corollary 1 and Assumption 4 hold. Under the sample size scaling  $N \geq s^2 L_2 \log n$ , we have

$$\left\| \Theta^* \hat{\Sigma}_N - I_n \right\|_{\infty} \leq c \left( \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}} \right),$$

with probability at least  $1 - c' \exp(-c'' \log s)$ .

*Proof.* We have

$$\begin{aligned} \left\| \Theta^* \hat{\Sigma}_N - I_n \right\|_{\infty} &= \left\| (\Theta^* \Sigma_N - I_n) + \Theta^* (\hat{\Sigma}_N - \Sigma_N) \right\|_{\infty} \leq \left\| \Theta^* \Sigma_N - I_n \right\|_{\infty} + \left\| \Theta^* (\hat{\Sigma}_N - \Sigma_N) \right\|_{\infty} \\ &\leq \left\| \Theta^* \Sigma_N - I_n \right\|_{\infty} + \left\| \frac{1}{N} \sum_{i=1}^N (\hat{W}_i^{-1} - W_i^{-1}) \Theta^* (X_i X_i^T - \Sigma_x) \right\|_{\infty} + \left\| \frac{1}{N} \sum_{i=1}^N (\hat{W}_i^{-1} - W_i^{-1}) \Theta^* \Sigma_x \right\|_{\infty} \\ &\leq \left\| \Theta^* \Sigma_N - I_n \right\|_{\infty} + \left\| \hat{W}^{-1} - W^{-1} \right\|_{\infty} \left[ \max_{j,k} \left( \frac{1}{N} \sum_{i=1}^N |e_j^T \Theta^* (X_i X_i^T - \Sigma_x) e_k| \right) + \lambda_{\max}(\Theta^*) \lambda_{\max}(\Sigma_x) \right]. \end{aligned}$$

We've bounded  $\left\| \Theta^* \Sigma_N - I_n \right\|_{\infty}$  in Lemma 8. Furthermore, using the Bernstein-type inequality and taking a union bound, we have

$$\mathbb{P} \left( \max_{j,k} \left( \frac{1}{N} \sum_{i=1}^N |e_j^T \Theta^* (X_i X_i^T - \Sigma_x) e_k| \right) > 4\lambda_{\max}(\Theta^*) \sigma_x^2 + c\sqrt{\frac{\log n}{N}} \right) \leq \frac{1}{n}. \quad (\text{C.1})$$

Therefore, it remains to bound  $\|\widehat{W}^{-1} - W^{-1}\|_\infty$ . First, by Taylor's theorem, we have

$$\frac{1}{g(X_i, \hat{\beta})} - \frac{1}{g(X_i, \beta^*)} = (\hat{\beta} - \beta^*)^T d_i(\beta^*) + \frac{1}{2} (\hat{\beta} - \beta^*)^T H_i(\beta^* + c_i(\hat{\beta} - \beta^*)) (\hat{\beta} - \beta^*),$$

where

$$c_i \in (0, 1), \quad d_i(\beta) = \frac{-1}{g^2(X_i, \beta)} \frac{\partial g(X_i, \beta)}{\partial \beta},$$

and

$$H_i(\beta) = \frac{2}{g^3(X_i, \beta)} \frac{\partial g(X_i, \beta)}{\partial \beta} \left( \frac{\partial g(X_i, \beta)}{\partial \beta} \right)^T - \frac{1}{g^2(X_i, \beta)} \frac{\partial^2 g(X_i, \beta)}{\partial \beta^2}.$$

Next, if

$$\hat{S} \subseteq S \quad \text{and} \quad \|\hat{\beta} - \beta^*\|_2 \leq c \sqrt{\frac{s L_2 \log n}{N}},$$

then we have

$$\begin{aligned} \|\widehat{W}^{-1} - W^{-1}\|_\infty &= \max_i \left| \frac{1}{g(X_i, \hat{\beta})} - \frac{1}{g(X_i, \beta^*)} \right| \leq \|\hat{\beta} - \beta^*\|_2 \max_i \|d_i(\beta^*)_S\|_2 \\ &+ \frac{1}{2} \|\hat{\beta} - \beta^*\|_2^2 \max_i \|H_i(\beta^* + c_i(\hat{\beta} - \beta^*))_{SS}\|_{op} \leq c \left( s \sqrt{\frac{L_2 \log n}{N}} + \frac{s^2 L_2 \log n}{N} \right), \end{aligned}$$

where  $d_i(\beta^*)_S$  is subvector of  $d_i(\beta^*)$  with elements restricted to  $S$ ,  $H_i(\beta^* + c_i(\hat{\beta} - \beta^*))_{SS}$  is submatrix of  $H_i(\beta^* + c_i(\hat{\beta} - \beta^*))$  with rows and columns restricted to  $S$ , and the last inequality follows from

$$\max_i \|d_i(\beta^*)_S\|_2 \leq c \sqrt{s},$$

and

$$\max_i \|H_i(\beta^* + c_i(\hat{\beta} - \beta^*))_{SS}\|_{op} \leq \|H_i(\beta^* + c_i(\hat{\beta} - \beta^*))_{SS}\|_F \leq c' s.$$

Applying Corollary 1 and assuming  $N \geq s^2 L_2 \log n$ , we have

$$\mathbb{P} \left( \|\widehat{W}^{-1} - W^{-1}\|_\infty \leq c s \sqrt{\frac{L_2 \log n}{N}} \right) \geq 1 - c' \exp(-c'' \log s). \quad (\text{C.2})$$

Finally, combining Lemma 8, (C.1), and (C.2), we have

$$\mathbb{P} \left( \|\Theta^* \widehat{\Sigma}_N - I_n\|_\infty \leq c \left( \frac{1}{L_1} \sqrt{\frac{\log n}{N}} + s \sqrt{\frac{L_2 \log n}{N}} \right) \right) \geq 1 - c' \exp(-c'' \log s).$$

Therefore, the proof is complete.  $\square$

**Lemma 14.** (a) Assume that Assumption 2 and Assumption 4 hold. Under the scaling  $N \gtrsim \log n$ , we have

$$\max_j \left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i d_i(\beta^*) \right\|_\infty \leq c \sqrt{\log n},$$

with probability at least  $1 - c' \exp(-c'' \log n)$ .

(b) Assume that Assumption 2 holds. We have

$$\max_j \frac{1}{\sqrt{N}} \sum_{i=1}^N |X_{ij} \varepsilon_i| \leq c\sqrt{N} \log n,$$

with probability at least  $1 - c' \exp(-c'' \log n)$ .

(c) Assume that Assumption 2 holds. Under the scaling  $N \gtrsim \log n$ , we have

$$\left\| \frac{1}{\sqrt{N}} X^T W^{-1} \varepsilon \right\|_{\infty} \leq c\sqrt{\log n},$$

with probability at least  $1 - c' \exp(-c'' \log n)$ .

*Proof.* First, we show Part (a). We have

$$\max_j \left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i d_i(\beta^*) \right\|_{\infty} = \max_{1 \leq j \leq n, 1 \leq k \leq n} \left| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i d_{ik}(\beta^*) \right|,$$

where  $d_{ik}(\beta^*)$  is the  $k$ th element of  $d_i(\beta^*)$ . Note that

$$\frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i d_{ik}(\beta^*) \mid X \sim N \left( 0, \frac{1}{N} \sum_{i=1}^N X_{ij}^2 d_{ik}^2 W_i \right),$$

where

$$\frac{1}{N} \sum_{i=1}^N X_{ij}^2 d_{ik}^2 W_i \leq \frac{c}{N} \sum_{i=1}^N X_{ij}^2.$$

By Lemma 13 in [23], we have

$$\mathbb{P} \left( \frac{1}{N} \sum_{i=1}^N X_{ij}^2 \leq (\Sigma_x)_{jj} + c \right) \geq 1 - \exp(-c' N).$$

Therefore, using a similar argument with Lemma 1 and assuming  $N \gtrsim \log n$ , we have

$$\max_j \left\| \frac{1}{\sqrt{N}} \sum_{i=1}^N e_j^T X_i \varepsilon_i d_i(\beta^*) \right\|_{\infty} \leq c\sqrt{\log n},$$

with probability at least  $1 - c' \exp(-c'' \log n)$ .

Next, using similar arguments with Part (a) yields results in Part (b) and Part (c). Therefore, the proof is complete.  $\square$