

TESTABILITY OF HIGH-DIMENSIONAL LINEAR MODELS WITH NON-SPARSE STRUCTURES

BY JELENA BRADIC[§]
JIANQING FAN[¶] AND YINCHU ZHU[‡]

*University of California, San Diego[§], Princeton University[¶]
and University of Oregon[‡]*

Understanding statistical inference under possibly non-sparse high-dimensional models has gained much interest recently. For a given component of the regression coefficient, we show that the difficulty of the problem depends on the sparsity of the corresponding row of the precision matrix of the covariates, not the sparsity of the regression coefficients. We develop new concepts of uniform and essentially uniform non-testability that allow the study of limitations of tests across a broad set of alternatives. Uniform non-testability identifies a collection of alternatives such that the power of any test, against any alternative in the group, is asymptotically at most equal to the nominal size. Implications of the new constructions include new minimax testability results that, in sharp contrast to the current results, do not depend on the sparsity of the regression parameters. We identify new tradeoffs between testability and feature correlation. In particular, we show that, in models with weak feature correlations, minimax lower bound can be attained by a test whose power has the \sqrt{n} rate, regardless of the size of the model sparsity.

1. Introduction. Confidence intervals construction and hypothesis testing in high-dimensional studies arise in almost all modern application areas, ranging from biomedical imaging (Chalkidou et al., 2015) or disease tracking, to the discovery of genetic variants associated with normal and disorder-related phenotypic variance in brain function (Ganjgahi et al., 2018; Krishnan et al., 2016), to the evaluation of policy and marketing strategies (Verhoef et al., 2017), and many more. There has been considerable interest in developing valid statistical methods for the construction of confidence intervals in high-dimensional problems. Some notable recent advances include proposals based on the ridge estimate (Bühlmann, 2013; Nickl and van de Geer, 2013),

[§]Bradic’s research is supported by NSF Grant DMS-1712481.

[¶]Fan’s research is supported by NSF Grants DMS-1662139 and DMS-DMS-1712591 and NIH grants 5R01-GM072611-12.

MSC 2010 subject classifications: Primary 62C20, 62F03; Secondary 62F30, 62J07.

Keywords and phrases: Minimax theory, ℓ_2 -constraint, Confidence intervals, Uniform non-testability

on the lasso estimate (Van de Geer et al., 2014; Zhang and Zhang, 2014), score and orthogonal moments methods (Belloni et al., 2014a; Goeman et al., 2006), as well as combinations thereof (see for example Belloni et al. (2014b); Javanmard and Montanari (2014)).

This line of work has produced many promising methods. The literature, however, does not provide an answer as to how these methods should be adapted for the possible lack of sparse structures in the underlying models. First, there is no guidance on how to check whether a model is sparse or not. The majority of current approaches construct confidence intervals under a set of assumptions describing how sparse the underlying model is. The process of developing algorithms that detect model sparsity is still somewhat “unattainable”, therefore in practice effectively rendering a priori belief in the sparsity. Second, no formal guarantees have been provided, to either confirm or deny, the ability to perform a hypothesis test (or to construct optimal confidence intervals); not without imposing sparsity on model parameters.

In this paper, our primary goal is a theoretical understanding of the high-dimensional minimax theory that can address both of these concerns. Our framework allows for high dimensional linear models that are not necessarily sparse. We illustrate that moving away from assumptions on sparse parameters towards assumptions on the design matrix can allow for certain optimal inferences. Moreover, we show how the estimators and tests can be designed to achieve these new optimality results.

We formalize our results in terms of the high-dimensional linear regression:

$$(1) \quad \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}_n)$$

where $\mathbf{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$, \mathbf{X} is a collection of n i.i.d. vectors, $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)^\top \in \mathbb{R}^{n \times p}$ and $\boldsymbol{\beta} \in \mathbb{R}^p$ whose dimensionality p can be much larger than the sample size n . Here, the covariance matrix of \mathbf{X} is denoted by $\boldsymbol{\Sigma} = \mathbb{E}(\mathbf{X}_i \mathbf{X}_i^\top)$, whereas its precision matrix is denoted by $\boldsymbol{\Omega} = \boldsymbol{\Sigma}^{-1}$. We denote with $k = \|\boldsymbol{\beta}\|_0$. In this paper, our focus is on the problem of testing individual entries of $\boldsymbol{\beta}$. Without loss of generality, we consider the first entry and denote $\boldsymbol{\beta} = (\beta, \boldsymbol{\gamma}^\top)^\top \in \mathbb{R}^p$.

We provide a motivating result first. Note that β can be represented as a linear combination of easily estimable quantity, $\mathbb{E}(\mathbf{X}_i y_i)$, with the weights being the first row of $\boldsymbol{\Omega}$. We investigate if particular structures in $\boldsymbol{\Omega}$ can be leveraged to remove sparsity assumptions on k . When $\boldsymbol{\Omega}$ is known, we show that a simple plug-in estimate achieves the parametric rate for a full range $0 \leq k \leq p$. Hence, there is hope that strict sparsity requirements on k are not necessary for valid inference. However, there are significant hurdles that need to be cleared before minimax results can be directly developed for

inference on models that are not necessarily sparse.

An impediment to exploring high-dimensional models is the fear that the researchers will search for essential variables, and then report only the results for variables with extreme effects, which in turn, are dependent on the existence of only a small number of significant signals; highlighting thus the signal that may be purely spurious. For this reason, such practices must specify in advance that only a few signals are “real” and then they proceed to find them. However, such procedures can make it difficult to discover strong but unexpected signals. In this paper, we seek to address this challenge by developing a method that yields valid asymptotic confidence intervals for the real underlying signal by moving away from conditions on the conditional expectation of $\mathbf{y}|\mathbf{X}$ to exploring structures in the distribution of \mathbf{X} such as the sparsity of $\mathbf{\Omega}$. We showcase that sparsity in $\mathbf{\Omega}$ can allow for arbitrary growth of k .

In GWAS studies, an agnostic approach to the conditional distribution of the response is especially valuable. Since around 2006, the advent of GWAS, and more recently exome sequencing, has provided the first detailed understanding of the genetic basis of complex traits. To explain “missing heritability,” a new paradigm has emerged in which complex disease is driven by an accumulation of a large number of weak effects across all of the network of genetic pathways (Boyle et al., 2017; Chakravarti and Turner, 2016; Furlong, 2013). Similarly, it is deeply understood that microbial functional relationship to the host is highly complex, that microbial communities have highly complex structures and that small and numerous changes in the network affect the host adversely (Huttenhower et al., 2012). At the same time, it is widely believed that features in many studies have a sparse correlation structure (providing evidence of sparse $\mathbf{\Omega}$). For example, only a certain number of genes functionally depend on one-another, clump together. Similarly, far apart, SNPs are very nearly independent (Janson et al., 2017), so we may expect that the true $\mathbf{\Omega}$ has nearly banded structure.

Therefore, for many practically relevant examples, it is not necessary nor wise to impose a sparse structure on the conditional distribution of $\mathbf{y}|\mathbf{X}$; after all, if we are studying $\mathbf{y}|\mathbf{X}$, that typically means we do not know very much about it.

Our detection rates are stated in terms of s , the number of non-zero entries in the first row of $\mathbf{\Omega}$ as well as the size of the $\|\boldsymbol{\beta}\|_2$. Thanks to the newly defined optimality criterion, the rate $n^{-1/2} + sn^{-1} \log p$ is identified as the minimax rate of detection for the problem of identifying the null $\boldsymbol{\beta} = \boldsymbol{\beta}_0$ against the alternative $\boldsymbol{\beta} = \boldsymbol{\beta}_0 + h$, whenever $\|\boldsymbol{\beta}\|_2 \leq \kappa$ for a fixed $\kappa > 0$, regardless of the size of the model sparsity k . When $\mathbf{\Omega}$ ’s first row is sparse

enough, we provide a minimax optimal test and a confidence interval for β without assuming an upper bound on k . We identify as well that even with knowledge of sparsity of $\mathbf{\Omega}$, the detection rate will not tend to zero if $\kappa \gtrsim \sqrt{n}$ and no constraint is imposed on k .

We propose a novel framework to study the detection rates for β while allowing $k \lesssim p$. Impossibility results are established under the new concept of (essentially) uniform non-testability. We state that the null hypothesis is uniformly non-testable against the alternative if the power of any test of nominal level α against *any* point in the alternative is at most α . The proposed uniform non-testability results also provide new insights. Under uniform non-testability, testing the null hypothesis against one (arbitrary) point is impossible for *any* test. Since any test that has size control is powerless against every point in the alternative, our work indicates that the difficulty in these testing problems is quite fundamental. Besides, the new non-testability results allow for a characterization of non-adaptivity; in a certain sense, those two notions match. It will enable us to shed new light on the existing literature on the adaptivity of testing. Ideally, an adaptive confidence interval should have its length automatically adjusted to the actual sparsity of the unknown coefficient vector, while maintaining a pre-specified coverage probability. We showcase that with known $\mathbf{\Omega}$ this can be done while for the unknown $\mathbf{\Omega}$, adaptivity requires $s \ll \sqrt{n}/\log p$; both results do not depend on the size of k .

1.1. *Existing literature.* Under the linear model above, the parameter of interest can be written as

$$\beta = \mathbf{\Omega}_{,1} \mathbb{E}(\mathbf{X}_i y_i)$$

where $\mathbf{\Omega}_{,1} \in \mathbb{R}^p$ denotes the first row of $\mathbf{\Omega}$. As a consequence of this representation, it may be tempting to first estimate $\mathbf{\Omega}_{,1}$ as well as $\mathbb{E}(\mathbf{X}_i y_i)$ and then set $\tilde{\beta} = \hat{\mathbf{\Omega}}_{,1} \hat{\mu}$, where $\hat{\mu} = n^{-1} \sum_{i=1}^n \mathbf{X}_i y_i$. This simple approach, however, is often not optimal: Because $n^{-1} \sum_{i=1}^n \mathbf{X}_i y_i$ is a p -dimensional vector, that does not have to have any sparse structures, the product may be highly unstable. As an example, consider fitting the graphical lasso (Meinshausen and Bühlmann, 2006; Wainwright et al., 2007) to estimate $\mathbf{\Omega}_{,1}$. A naive approach would make a product of such an estimate and $\hat{\mu}$ to construct $\tilde{\beta}$. However, since $\hat{\mathbf{\Omega}}_{,1}$ is regularized towards zero, the bias at estimation will propagate in all elements of $\hat{\mu}$ and, therefore, the product.

The recent literature on high-dimensional inference has proposed several ideas on how to avoid such “regularization bias”. In particular, several recent papers have proposed structural changes to various regularized methods,

aimed at accurate estimation of β (Belloni et al., 2014a,b; Bühlmann, 2013; Goeman et al., 2006; Javanmard and Montanari, 2014; Nickl and van de Geer, 2013; Van de Geer et al., 2014; Zhang and Zhang, 2014). These approaches always correctly de-bias the estimates for valid high-dimensional inference. However, they assume various sparsity structures in their analysis without which no guarantees are provided for validity. In detail, their analysis relies on the assumption that the vector of the nuisance parameters belongs to the set of k -sparse regression vectors with $k \ll \sqrt{n}/\log p$. Such sparsity requirement recently raised considerable interest since it appears to be a much stronger condition than that needed for consistent estimation, which only imposes $k \ll n/\log p$; see, e.g., Negahban et al. (2009); Raskutti et al. (2011).

The natural question is whether the strong condition of $k \ll \sqrt{n}/\log p$ is needed. The pioneering work of Cai and Guo (2017) and Javanmard and Montanari (2018) aim to address this question, where the former derives the minimax rate for the expected length of confidence intervals assuming $k \lesssim n/\log p$ and the latter, in a different context, improves the condition $k \ll \sqrt{n}/\log p$ to $k \ll n/(\log p)^2$. This work provides a complementary study where we reveal an intricate relationship between sparsity (or the non-existence of thereof) and ℓ_2 -norm constraints.

Another line of work, closer to our paper, has focused on inference approaches not closely relying on sparsity assumptions; see e.g., Shah and Bühlmann (2017) and Janson et al. (2017). The work of Zhu and Bradic (2018) and Zhu and Bradic (2017) is particularly close to ours. There, the authors propose asymptotically exact confidence interval construction under no model sparsity assumption. However, therein, no formal optimality guarantees were derived beyond several specific examples in which the model parameters are restricted to be small or approximately sparse. Therefore, it is not apparent what the optimal detection rate is for general non-sparse models, and it is not expected that methods discussed therein can provide uniform guarantees for an ample parameter space. Inspired by those findings, we asked whether any formal, minimax guarantees can be provided for a class of dense models? If so, what kind of estimates would be able to achieve the fundamental limits of detection? We identify that sample-splitting helps guarantee uniform detection rates. We discuss in detail sparsity in the precision matrix as being sufficient and necessary tools for this purpose. We also showcase an increase in the minimax (testing) rates whenever ℓ_2 -norm of the model parameters is not bounded, and the model is not necessarily sparse.

1.2. *Organization of the paper.* The rest of the paper is organized as follows: After basic notation is introduced, Section 2 presents a precise formulation of the problem and some initial insights. Section 3 establishes two impossibility results under the lack of sparsity in the first row of $\mathbf{\Omega}$. These results provide a lower bound on the detection rates. Section 4 focuses on the upper bounds and the attainability of lower bounds. Section 5 discusses connections to the minimax rates of detection and adaptivity of the confidence intervals. Section 6 discusses minimax detection rates with growing ℓ_2 balls. The proofs of all of the results are presented in the Appendices: A-C are collected in the main document whereas D-L are presented in the Supplement.

2. Problem setup. We present in this section the framework for hypothesis in high-dimensional models that are not necessarily sparse. We begin with the notation that will be used throughout the manuscript.

2.1. *Notation.* For a matrix $\mathbf{X} \in \mathbb{R}^{n \times p}$, \mathbf{X}_i , \mathbf{X}_j and X_{ij} denote respectively the i -th row, j -th column and (i, j) entry of the matrix \mathbf{X} , $\mathbf{X}_{i,-j}$ denotes the i -th row of \mathbf{X} excluding the j -th coordinate, and \mathbf{X}_{-j} denotes the submatrix of \mathbf{X} excluding the j -th column. Let $[p] = \{1, 2, \dots, p\}$. For a subset $J \subseteq [p]$, \mathbf{X}_J denotes the submatrix of \mathbf{X} consisting of columns \mathbf{X}_j with $j \in J$ and for a vector $\mathbf{x} \in \mathbb{R}^p$, \mathbf{x}_J is the p -dimensional vector that has the same coordinates as \mathbf{x} on J and zero coordinates on the complement J^c of J . Let \mathbf{x}_{-J} denote the subvector with indices in J^c . For a set S , $|S|$ denotes its cardinality. For a vector $\mathbf{x} \in \mathbb{R}^p$, $\text{supp}(\mathbf{x})$ denotes the support of \mathbf{x} and the ℓ_q -norm of \mathbf{x} is defined as $\|\mathbf{x}\|_q^q = \sum_{j \in [p]} |x_j|^q$ for $q \geq 0$, with $\|\mathbf{x}\|_0 = |\text{supp}(\mathbf{x})|$ and $\|\mathbf{x}\|_\infty = \max_{j \in [p]} |x_j|$. For a matrix \mathbf{A} and $1 \leq q \leq \infty$, $\|\mathbf{A}\|_q = \sup_{\mathbf{x}: \|\mathbf{x}\|_q=1} \|\mathbf{A}\mathbf{x}\|_q$. For a symmetric matrix \mathbf{A} , $\lambda_{\min}(\mathbf{A})$ and $\lambda_{\max}(\mathbf{A})$ denote respectively the smallest and largest eigenvalue of \mathbf{A} . \mathbb{I}_q denotes the $q \times q$ identity matrix. For two positive sequences a_n and b_n , $a_n \lesssim b_n$ means $a_n \leq Cb_n$ for a positive constant C independent of n . Moreover, we use $a_n \asymp b_n$ if $b_n \lesssim a_n$ and $a_n \lesssim b_n$. Lastly, $a_n \ll b_n$ is used to denote that $\lim_{n \rightarrow \infty} a_n/b_n = 0$. For $a \in \mathbb{R}$, let $\lfloor a \rfloor$ denote the largest integer that is at most a .

2.2. *High-dimensional linear models that are not necessarily sparse.* We shall focus on the high-dimensional linear model (1) with the random design such that $\mathbf{X}_i \sim \mathcal{N}_p(0, \mathbf{\Sigma})$, $1 \leq i \leq n$ and are independent of the error ε . Note that both $\mathbf{\Sigma}$ and the noise level σ are considered as unknown. Since our problem is centered around the construction of confidence intervals for

the univariate parameter β , we re-parametrize model (1) as

$$(2) \quad \mathbf{y} = \mathbf{Z}\beta + \mathbf{W}\boldsymbol{\gamma} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \sigma^2 \mathbb{I}),$$

where $\beta \in \mathbb{R}$, $\boldsymbol{\gamma} \in \mathbb{R}^{p-1}$, $\mathbf{Z} = (Z_1, \dots, Z_n)^\top \in \mathbb{R}^n$ and $\mathbf{W} = (\mathbf{W}_1, \dots, \mathbf{W}_n)^\top \in \mathbb{R}^{n \times (p-1)}$. The distribution of the data is now indexed by the parameter

$$\theta = (\beta, \boldsymbol{\gamma}, \boldsymbol{\Sigma}, \sigma),$$

which consists of parameter of interest β , the nuisance parameters $\boldsymbol{\gamma}$, the covariance matrix $\boldsymbol{\Sigma} = \mathbb{E}[\mathbf{X}_i \mathbf{X}_i^\top]$ of the random design vector $\mathbf{X}_i = (Z_i, \mathbf{W}_i^\top)^\top$, and the variance of the noise σ . Observed data $\mathbf{D} = \{D_1, \dots, D_n\}$ consists of i.i.d. triplets $D_i = (y_i, Z_i, \mathbf{W}_i)$, for $i = 1, \dots, n$. Note that β in (2) can be represented as

$$(3) \quad \beta = \boldsymbol{\Omega}_1 \mathbb{E}(\mathbf{X}_i^\top y_i).$$

Since each element of $\mathbb{E}(\mathbf{X}_i^\top y_i)$ can be easily estimated at a root- n rate, the estimability of β depends on $\boldsymbol{\Omega}$ only through its first row. Hence, it seems prudent to define a parameter space that includes both the parameters of the model as well as the matrix $\boldsymbol{\Omega}$,

$$(4) \quad \tilde{\Theta} = \left\{ \theta = (\beta, \boldsymbol{\gamma}, \boldsymbol{\Sigma}, \sigma) : \beta \in \mathbb{R}, \right. \\ \left. M^{-1} \leq \lambda_{\min}(\boldsymbol{\Sigma}) \leq \lambda_{\max}(\boldsymbol{\Sigma}) \leq M, 0 \leq \sigma \leq M_1, \|\boldsymbol{\beta}\|_2 \leq M_2 \right\},$$

where $\boldsymbol{\beta} = (\beta, \boldsymbol{\gamma}^\top)^\top$ and $M > 1$, M_1 and M_2 are positive constants. Note that

$$\text{Var}(y_i) = \boldsymbol{\beta}^\top \boldsymbol{\Sigma} \boldsymbol{\beta} + \sigma^2 \geq \lambda_{\min}(\boldsymbol{\Sigma}) \|\boldsymbol{\beta}\|_2^2.$$

Thus, the constraint $\|\boldsymbol{\beta}\|_2 \leq M_2$ can be dropped in the above definition if we only consider bounded $\text{Var}(y_i)$.

Observe that whenever $\boldsymbol{\Omega}$ is known, a simple plug-in estimate

$$(5) \quad \hat{\beta} = \boldsymbol{\Omega}_1 \mathbf{X}^\top \mathbf{y} / n$$

achieves the parametric rate without any assumption on k . Namely, we provide the following result.

THEOREM 1. *For $\tilde{\Theta}$ defined in (4) we have*

$$(6) \quad \sup_{\theta \in \tilde{\Theta}} \mathbb{E}_\theta |\hat{\beta} - \beta| \asymp n^{-1/2}.$$

Theorem 1 is an oracle-like statement that holds for bounded constant M_2 . It indicates that a parametric rate of detection is possible for dense parameters (with $p \gg n$) with bounded ℓ_2 norm. Majority of the present paper focuses on the parameter space $\tilde{\Theta}$. In Section 6 we showcase minimax optimality rates that do not restrict the growth of M_2 .

Observe that the above result allows for $p \gg n$; in fact, it does not put any restrictions on the growth of p or k . Additionally, Theorem 1 identifies that the inference for non-sparse high-dimensional models is possible as long as the precision matrix is known. It indicates that the ability to decorrelate the features (i.e., to estimate $\mathbf{\Omega}_1$, well) is the key to efficient inference in high-dimensional non-sparse models.

To further study lower limits of detection of testing

$$H_0 : \beta = \beta_0$$

our focus is on the parameter spaces defined by

$$(7) \quad \Theta = \left\{ \theta = (\beta, \gamma, \mathbf{\Sigma}, \sigma) : \beta \in \mathbb{R}, M^{-1} \leq \lambda_{\min}(\mathbf{\Sigma}) \leq \lambda_{\max}(\mathbf{\Sigma}) \leq M, \right. \\ \left. \mathbf{\Sigma}_{(-1),(-1)} = \mathbb{I}_{p-1}, 0 \leq \sigma \leq M_1, \text{ and } \|\beta\|_2 \leq M_2 \right\},$$

where $M > 1$ and $M_1, M_2 > 0$ are some universal constants. To study upper limits of detection we still analyze $\tilde{\Theta}$ as defined in (4). It is worth pointing that none of the parameter spaces, $\tilde{\Theta}$ or Θ , restricts k , the number of non-zero elements in β of the linear model (1), or the ℓ_1 -norm of β (which can grow at a rate of \sqrt{p}). Our work is hence very different from existing minimax studies. We also define

$$\Theta(s_0, \beta_0) = \{ \theta = (\beta, \gamma, \mathbf{\Sigma}, \sigma) \in \Theta : \beta = \beta_0, \|\mathbf{\Omega}_1\|_0 \leq s_0 \}$$

and

$$\Theta(s_0) = \bigcup_{\beta \in \mathbb{R}} \Theta(s_0, \beta).$$

The main goal of this paper is to address the following questions:

1. *Is it possible to have accurate inference procedure about univariate parameters without requiring the model parameter β itself to be sparse?*
2. *Is the accuracy in terms of the detection rates uniform over the parameter space?*

3. Lower bound. For $0 < \alpha < 1$ and a given parameter space Θ_1 , the set of tests of nominal level $\alpha \in (0, 1)$ regarding the null hypothesis $\theta \in \Theta_1$ is denoted with

$$\Psi_\alpha(\Theta_1) = \left\{ \psi : \mathbf{D} \mapsto [0, 1] : \sup_{\theta \in \Theta_1} \mathbb{E}_\theta \psi \leq \alpha \right\},$$

see, e.g., [Lehmann and Romano \(2006\)](#). Here, we allow for both random and non-random tests.

DEFINITION 1 (Uniform non-testability). *Consider the hypothesis testing problem of $H_0 : \theta \in \Theta^{(1)}$ versus $H_1 : \theta \in \Theta^{(2)}$. We say that $\Theta^{(1)}$ is asymptotically uniformly non-testable against $\Theta^{(2)}$ at size $\alpha \in (0, 1)$ if $\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta^{(2)}} \mathbb{E}_\theta \psi \leq \alpha$ for any test $\psi \in \Psi_\alpha(\Theta^{(1)})$.*

Above Definition 1 introduces new concept of testability. Per Definition 1 there does not exist a test that is better than a simple coin toss. Since a simple coin toss is uniformly most powerful asymptotically, the data cannot provide sufficient statistical evidence to distinguish the null from the alternative hypothesis. This concept provides an alternative to the widely known minimax-type results which state that for any test, there is one “difficult” point in the alternative for which this test has no power; therefore, it is possible that beyond this “difficult” point, there might exist a test that has good power against all the other points. We could argue that we are proposing a different and not necessarily better characterization of optimality.

To characterize alternative hypothesis we introduce

$$\Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n) = \left\{ \theta \in \Theta(s/2, \beta_0 + a) : 0 \leq a \leq h_n, \|\beta\|_2 \leq \zeta M_2, \right. \\ \left. (\zeta M)^{-1} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq \zeta M, \kappa \leq \sigma \leq \zeta M_1 \right\},$$

where h_n is a sequence of positive numbers and $\zeta \in (M^{-1}, 1)$ and $\kappa \in (0, \zeta M_1)$ are constants.

THEOREM 2. *Suppose that $sn^{-1} \log p \leq 1/4$ and $2 \leq s \leq p^c$ for some constant $c < 1/2$. Then we have that for any β_0*

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi = \alpha,$$

where $h_n = \rho sn^{-1} \log p$ and

$$(8) \quad \rho = \min \left\{ 4, \frac{1/2 - c}{15(\kappa^{-2}M + 1)}, \frac{2(\zeta^{-1} - 1)^2}{M^3(2M + 1)}, \frac{2M(1 - \zeta)^2}{2M + 1}, \right. \\ \left. \frac{(1 - \zeta^2)M_2}{8\zeta\sqrt{M}}, \frac{\kappa^2(1 - \zeta^2)^2 M_2^2}{64\zeta^4 M M_1^2}, \frac{M_2\sqrt{1 - \zeta^2}}{2\sqrt{M}}, \frac{\kappa^2(1 - \zeta^2)M_2^2}{4\zeta^2 M_1^2 M} \right\}.$$

Theorem 2 establishes that $\Theta(s, \beta_0)$ is uniformly non-testable against all points in the alternative $\Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)$, i.e. *every* point in $\Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)$ is difficult for *every* test. The distance to the alternatives, $\rho sn^{-1} \log p$, depends on the unknown constants M , M_1 and M_2 characterizing invertability of the covariance matrix Σ , noise level σ and the norm $\|\beta\|_2$, respectively; see e.g., (8).

This result is unique in its treatment of nuisance parameters γ , which are allowed to be fully dense. In the case of dense models, the lower bound for detection depends on how sparse Ω_1 , is: it is impossible to have power in testing $\beta = \beta_0$ against $\beta = \beta_0 + h$ whenever $|h| \leq \rho sn^{-1} \log p$. One implication is that when Ω_1 is not sparse enough (i.e., $s \gtrsim n/\log p$), a detection of alternatives separated by a constant is not guaranteed; that is, even deviation of non-vanishing magnitude cannot be detected.

The proof of Theorem 2 is formulated in a novel way. For any point in the alternative hypothesis, we compute the χ^2 distance between that alternative and a large collection of points in the null hypothesis. Whenever this distance is small, it indicates that the average rejection probability for that particular alternative is close to the average rejection probability for many of the nulls—therefore indicating lack of power. Although uniform non-testability explores a class of α level tests (only), by inspecting the proof of Theorem 2, we see that the detection rate would not change even if we localize our problem and impose $k \lesssim n/\log(p)$. Therefore, the rate is not really driven by some ultra-dense (and hence seemingly hopeless) $k \gg n/\log(p)$ points in the parameter space.

Another novelty in the theoretical analysis is the construction of the prior. The prior used by Cai and Guo (2017) can be adapted to the case of sparse Ω_1 , (instead of sparse γ as in their paper). However, that adaption would assume $\Omega_{1,-1} = 0$ and thus would not be enough to show the *uniformity* of non-testability. We compare this adaption with our construction in Appendix M.

Next, we fine-tune the above result in search of a parametric rate of detection. In view of that fact, we introduce a slightly weaker notion of essentially uniform non-testability.

DEFINITION 2 (Essentially uniform non-testability). *Consider the hy-*

pothesis testing problem of $H_0 : \theta \in \Theta^{(1)}$ versus $H_1 : \theta \in \Theta^{(2)}$. We say that $\Theta^{(1)}$ is asymptotically essentially uniformly non-testable against $\Theta^{(2)}$ at size $\alpha \in (0, 1/2)$ if $\limsup_{n \rightarrow \infty} \sup_{\theta \in \Theta^{(2)}} \mathbb{E}_\theta \psi \leq 2\alpha$ for any test $\psi \in \Psi_\alpha(\Theta^{(1)})$.

Essentially uniform non-testability implies

$$\liminf_{n \rightarrow \infty} \left(\alpha + \inf_{\psi \in \Psi_\alpha(\Theta^{(1)})} \inf_{\theta \in \Theta^{(2)}} \mathbb{E}_\theta(1 - \psi) \right) \geq 1 - \alpha.$$

We note that this statement implies the following claim on the minimax total error probability (a notion discussed by [Ingster et al. \(2010\)](#))

$$\liminf_{n \rightarrow \infty} \left(\alpha + \inf_{\psi \in \Psi_\alpha(\Theta^{(1)})} \sup_{\theta \in \Theta^{(2)}} \mathbb{E}_\theta(1 - \psi) \right) \geq 1 - \alpha.$$

We also denote

$$\Theta_\kappa(s, \beta_0 + h_n) = \left\{ \theta \in \Theta(s, \beta_0 + a) : 0 \leq a \leq h_n, \|\beta\|_2 \leq M_2, \right. \\ \left. M^{-1} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq M, \kappa \leq \sigma \leq M_1 \right\}.$$

THEOREM 3. *Suppose that $sn^{-1} \log p \leq 1/4$ and $2 \leq s \leq p^c$ for some constant $c < 1/2$. Then for any constant $\kappa \in (0, M_1]$, we have that for any β_0 ,*

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_\kappa(s, \beta_0 + h_n)} \mathbb{E}_\theta \psi \leq 2\alpha,$$

where $h_n = n^{-1/2}\tau$ and $\tau = \kappa \sqrt{M^{-1} \log(1 + \alpha^2)}$.

This results implies that $\Theta(s, \beta_0)$ is essentially uniformly non-testable against $\Theta_\kappa(s, \beta_0 + n^{-1/2}\tau)$. This result confirms the intuition that parametric rate is a fundamental boundary for statistical inference, an insight from the classical results of, for example, [Lehmann and Romano \(2006\)](#); [Van der Vaart \(2000\)](#). Let $c_0 = \min\{\rho, \tau\}$. Then

$$\Theta_{\zeta, \kappa}(s/2, \beta_0 + c_0(n^{-1/2} + sn^{-1} \log p)) \\ \subset \Theta_{\zeta, \kappa}(s/2, \beta_0 + \rho sn^{-1} \log p) \cap \Theta_\kappa(s, \beta_0 + \tau n^{-1/2}).$$

Hence, Theorems 2 and 3 imply

$$(9) \quad \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + c_0(n^{-1/2} + sn^{-1} \log p))} \mathbb{E}_\theta \psi \leq 2\alpha.$$

Therefore, constructing a meaningful test with a detection rate smaller than that of $n^{-1/2} + sn^{-1} \log p$ is indeed impossible.

4. Upper bound. In this section, we show that the lower limit of detection matches the upper limit of detection. In this section we focus our analysis on the space $\tilde{\Theta}(s)$. Formally, we define $\tilde{\Theta}(s_0) = \bigcup_{\beta_0 \in \mathbb{R}} \tilde{\Theta}(s_0, \beta_0)$ and

$$\tilde{\Theta}(s_0, \beta_0) = \left\{ \theta = (\beta, \gamma, \Sigma, \sigma) \in \tilde{\Theta} : \beta = \beta_0, \|\Omega_1\|_0 \leq s_0 \right\}.$$

We propose a test that achieves the bounds of Section 3. The newly proposed estimator $\hat{\beta}$ of β utilizes the constants that define the parameter set of interest to us, $\tilde{\Theta}(s, \cdot)$, and is therefore of pure theoretical interest. It is based on delicately designed high-dimensional estimators of the nuisance parameters: both of the model as well as that of the partial correlations of the features; an ℓ_1 consistent in the big coordinates while ℓ_∞ consistent in the small coordinates. Lastly, the new estimates are based on cross-fitting concepts enabling adaptivity to the rates of Section 3.

We introduce notation that helps with our construction. The constructed method will utilize a sample-splitting scheme. Let $b_n = \lfloor n/4 \rfloor$. We consider four non-overlapping subsets of the original sample $H_1 = \{1, \dots, b_n\}$, $H_2 = \{b_n + 1, \dots, 2b_n\}$, $H_3 = \{2b_n + 1, \dots, 3b_n\}$ and $H_4 = \{3b_n + 1, \dots, 4b_n\}$.

Next, we observe that the first row Ω_1 , takes the form $(1, -\pi^\top)/\sigma_{\mathbf{V}}^2$, where π and $\sigma_{\mathbf{V}}^2$ are from the regression

$$(10) \quad \mathbf{Z} = \mathbf{W}\pi + \mathbf{V},$$

where the vector \mathbf{V} is independent of \mathbf{W} with $\sigma_{\mathbf{V}}^2 = \mathbb{E}(\mathbf{V}^\top \mathbf{V})/n$. Moreover, observe that

$$(11) \quad y_i = \mathbf{W}_i^\top (\pi\beta + \gamma) + \eta_i$$

for $\eta_i = \beta v_i + \varepsilon_i$. Then, we notice that the parameter of interest, β , can be defined through a moment condition

$$\mathbb{E}[v_i y_i] = \beta \sigma_{\mathbf{V}}^2.$$

Therefore, for a suitably chosen estimator $\check{\pi}$ of π , let

$$\hat{v}_i = Z_i - \mathbf{W}_i^\top \check{\pi}$$

denote the estimated residuals of the model (10) and consider a natural estimator of β arising from the above moment condition

$$(12) \quad \hat{\beta} = \frac{\sum_{i \in H_4} \hat{v}_i y_i}{\sum_{i \in H_4} \hat{v}_i^2}.$$

Observe that this estimator is computed on the last fold, H_4 of the data; the remaining three folds are used to construct the estimator $\check{\pi}$. Note that the numerator in (12) is estimating

$$\mathbb{E}[v_i y_i] = \mathbb{E}[(Z_i - \mathbf{W}_i^\top \boldsymbol{\pi}) y_i] = \mathbb{E}(Z_i y_i) - \boldsymbol{\pi}^\top \boldsymbol{\xi}, \quad \boldsymbol{\xi} = \mathbb{E}[\mathbf{W}_i y_i].$$

Although the estimation of $\boldsymbol{\pi}$ is a sparse high-dimensional regression problem, existing estimators, such as Lasso, Dantzig selector or their debiased version, do not possess the theoretical properties we need for inference on β . Therefore, we construct a new estimator that is suitable for the purpose of inference. This new projected de-biased estimator $\check{\pi}$ of $\boldsymbol{\pi}$ aims to balance the good properties of both Lasso as well as de-biased Lasso estimator; balancing ℓ_1 with ℓ_∞ estimation quality.

We use the second and fourth fold of the data to construct cross-validated de-biased estimator of $\boldsymbol{\pi}$ in the following way. On the second fold compute a simple ℓ_1 -regularized estimator $\hat{\boldsymbol{\pi}}$,

$$\hat{\boldsymbol{\pi}} = \arg \min_{\mathbf{q} \in \mathbb{R}^{p-1}} b_n^{-1} \sum_{i \in H_2} (Z_i - \mathbf{W}_i^\top \mathbf{q})^2 + \lambda_\pi \|\mathbf{q}\|_1,$$

with $\lambda_\pi = 24M \sqrt{b_n^{-1} \log p}$.

To shrink the bias in estimated large coefficients of $\boldsymbol{\pi}$, we define a cross-fitted estimator as

$$\tilde{\boldsymbol{\pi}} = \hat{\boldsymbol{\pi}} + b_n^{-1} \sum_{i \in H_4} \hat{\boldsymbol{\Omega}}_{\mathbf{W}} \mathbf{W}_i (Z_i - \mathbf{W}_i^\top \hat{\boldsymbol{\pi}}),$$

In the above, $\hat{\boldsymbol{\Omega}}_{\mathbf{W}}$ is a carefully designed candidate estimate of $\boldsymbol{\Omega}_{\mathbf{W}} = \boldsymbol{\Sigma}_{\mathbf{W}}^{-1}$, that utilizes model (11) while ensuring that $\hat{\boldsymbol{\Omega}}_{\mathbf{W}}$ is close to $\boldsymbol{\Sigma}_{\mathbf{W}}^{-1}$. We propose the following cross-fitted spectral estimate

$$(13) \quad \begin{aligned} \hat{\boldsymbol{\Omega}}_{\mathbf{W}} &= \arg \min_{\mathbf{Q} \in \mathbb{R}^{(p-1) \times (p-1)}} \lambda_{\max}(\mathbf{Q}) \\ \text{s.t.} \quad &\mathbf{Q} = \mathbf{Q}^\top \\ &\left\| \{\mathbb{I}_{p-1} - \hat{\boldsymbol{\Sigma}}_{\mathbf{W}} \mathbf{Q}\} \hat{\boldsymbol{\xi}}_A \right\|_\infty \leq \lambda_\Omega \\ &\hat{\boldsymbol{\xi}}_A^\top \{\mathbf{Q} \hat{\boldsymbol{\Sigma}}_{\mathbf{W}} \mathbf{Q}\} \hat{\boldsymbol{\xi}}_A \leq \eta_\Omega, \end{aligned}$$

for $\hat{\boldsymbol{\Sigma}}_{\mathbf{W}} = b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top$ as well as

$$\lambda_\Omega = 24 \sqrt{b_n^{-1} \log p} M^3 M_2, \quad \eta_\Omega = 32 M^5 M_2^2.$$

In the above $\widehat{\boldsymbol{\xi}}$, a thresholded, marginal, estimate, and $\widehat{\boldsymbol{\Sigma}}_{\mathbf{W}}$ are computed on different folds of the data. Correlation estimate, $\widehat{\boldsymbol{\xi}}_A$, is defined as a sparse vector containing the top largest elements of the empirical inner product $\langle \mathbf{W}, y \rangle$. Set A denotes the largest elements,

$$(14) \quad A = \left\{ j \in [p] : |\widetilde{\boldsymbol{\xi}}_j| > \tau_n \right\}, \quad \widetilde{\boldsymbol{\xi}} = b_n^{-1} \sum_{i \in H_1} \mathbf{W}_i y_i.$$

Here, $\tau_n = 4M b_n^{-1} \sqrt{n(\log p)(M_1^2 + M_2^2)}$. Then, $\{\widehat{\boldsymbol{\xi}}_A\}_j = 0$ for $j \notin A$ and $b_n^{-1} \sum_{i \in H_3} \mathbf{W}_{ij} y_i$ otherwise.

Finally, we construct the following projected de-biased estimator

$$(15) \quad \begin{aligned} \check{\boldsymbol{\pi}} &= \arg \min_{\mathbf{q} \in \mathbb{R}^{p-1}} \|\mathbf{q}\|_1 \\ \text{s.t.} \quad & \left| \widehat{\boldsymbol{\xi}}_A^\top (\mathbf{q}_A - \widetilde{\boldsymbol{\pi}}_A) \right| \leq \eta_\pi \\ & \left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i (Z_i - \mathbf{W}_i^\top \mathbf{q}) \right\|_\infty \leq \lambda_\pi / 4 \\ & b_n^{-1} \sum_{i \in H_4} (Z_i - \mathbf{W}_i^\top \mathbf{q})^2 \geq \frac{1}{2M}, \end{aligned}$$

where the last three lines define the constraint set and where the tuning parameter η_π satisfies

$$\eta_\pi = 6408 \sqrt{b_n^{-1} \log p} M^4 M_2 s \lambda_\pi + 8 b_n^{-1/2} M^2 M_2 \sqrt{M \log(100/\alpha)}.$$

The estimator $\check{\boldsymbol{\pi}}$ is carefully crafted in order to achieve the desirable bias-variance tradeoff: it has small bias for entries corresponding to “large” elements of $\boldsymbol{\xi}$ and has small variance on other entries. Here, sample splitting is helpful in providing several independence structures that we need for the theoretical analysis; for example, the set A that defines “large” and “small” components needs to be independent of the subsequent constructions. As a result, $\check{\boldsymbol{\pi}}$ is quite different from the debiased estimator $\widetilde{\boldsymbol{\pi}}$ and these two estimators only behave similarly on large elements, i.e., $\left| \widehat{\boldsymbol{\xi}}_A^\top (\check{\boldsymbol{\pi}}_A - \widetilde{\boldsymbol{\pi}}_A) \right|$ is small.

We propose the following test

$$\psi_* = \mathbf{1} \left\{ |\widehat{\beta} - \beta_0| > c_n \right\},$$

where $\widehat{\beta}$ is defined in (12) and

$$(16) \quad c_n = 2M \left(10b_n^{-1/2} \sqrt{M(4M_2^2 M^3 + M_1^2) \log(100/\alpha)} + \right. \\ \left. + 34M(1 + M_2)\lambda_\pi^2 s + 1608b_n^{-1} M^2 \sqrt{n(\log p)(M_1^2 + M_2^2)} \lambda_\pi s + 2\eta_\pi \right).$$

We now show that even on the larger parameter $\tilde{\Theta}(s)$ (compared to $\Theta(s)$), the test ψ_* is indeed valid and has the optimal detection rate.

THEOREM 4. *Suppose that $p \geq \max \{2(1 + 1764M^2)s, 360/\alpha\}$ and*

$$n \geq \max \{4 + 784 \log p, (5067 + 220M^2) \log(100/\alpha), \\ 4 + 4054 [1 + 1764M^2] s \log(16ep)\}.$$

Then $\psi_ \in \Psi_\alpha(\tilde{\Theta}(s, \beta_0))$, i.e.,*

$$\sup_{\theta \in \tilde{\Theta}(s, \beta_0)} \mathbb{E}_\theta \psi_* \leq \alpha.$$

Moreover, $c_n \asymp sn^{-1} \log p + n^{-1/2}$ and

$$\inf_{\theta \in \tilde{\Theta}(s, \beta_0 + 3c_n)} \mathbb{E}_\theta \psi_* \geq 1 - \alpha.$$

Theorem 4 demonstrates that lower bound in (9) is achievable by a test ψ_* as defined above. Notice that requirement on n and p in Theorem 4 is mild; the key requirement is $s \lesssim n/\log p$. The proposed uniform non-testability results indicate the new detection boundary of $n^{-1/2} + sn^{-1} \log p$. Theorem 4 establishes that deviations of magnitude $3c_n \asymp n^{-1/2} + sn^{-1} \log p$ are uniformly testable over $\tilde{\Theta}(s)$, whereas results in Section 3 imply that even on the smaller $\Theta(s)$, deviations smaller than this rate are (essentially) uniformly non-testable.

Moreover, the parametric rate can be attained whenever $s \lesssim \sqrt{n}/\log p$. The case of $\sqrt{n}/\log p \ll s \ll n/\log p$ is more difficult and our proposed test still achieves the optimal rate. Note our test ψ_* depends on the knowledge of s . It turns out that the uniform non-testability results in Section 3 imply that such knowledge is required to achieve the minimax rate, indicating lack of adaptivity to the precision matrix sparsity. We make this argument precise in Section 5.2 and in more generality in Section 5.3.

5. Connections to minimax rates and confidence intervals. In this section we highlight the implication of the obtained results on the minimax theory and adaptivity.

5.1. *Minimax rates.* The (essential) uniform non-testability leads to the following minimax lower bound.

COROLLARY 5. *If $sn^{-1} \log p \leq 1/4$ and $2 \leq s \leq p^c$ for some constant $c < 1/2$, then there exists a constant $h_0 > 0$ such that for any β_0*

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\tilde{\Theta}(s, \beta_0))} \inf_{\theta \in \tilde{\Theta}(s, \beta_0 + h_n)} \mathbb{E}_\theta \psi \leq 2\alpha,$$

with $h_n = h_0(n^{-1/2} + sn^{-1} \log p)$.

Observe that Corollary 5 establishes a minimax claim that spans the space of $\tilde{\Theta}(s, \beta_0 + h)$; it does not impose $\Sigma_{(-1), (-1)} = \mathbb{I}_{p-1}$ and does not restrict k (the sparsity of β). Therefore, Corollary 5 directly refines the existing results on minimax testing, which routinely assume $k \lesssim n/\log p$; see Cai and Guo (2017, 2018); Cai and Low (2004, 2006); Genovese and Wasserman (2008); Hoffmann and Nickl (2011); Nickl and van de Geer (2013); Robins and Van Der Vaart (2006). Corollary 5 establishes a lower bound for the minimax detection rate of the null $H_0 : \beta = \beta_0$ against the alternative

$$H_1 : \beta = \beta_0 + h_0(n^{-1/2} + sn^{-1} \log p)$$

regardless of the sparsity of the nuisance parameter γ in the regression model (2). As such this result is the first that derives the lower bound for the detection rate under fairly general model setting and in particular not requiring a model to be sparse. Theorem 4 entails that sparsity of the first row of the precision matrix (alone) is sufficient for minimax inference (per Corollary 5), and the sparsity on regression coefficients is not necessary.

When $s \geq c_0 n/\log p$, a direct consequence of Corollary 5 is that it is impossible to distinguish $H_0 : \beta = \beta_0$ and $H_1 : \beta = \beta_0 + c_0 h_0$ in a minimax sense; in other words, there is no power even against fixed alternatives. Whenever Ω_1 is sparse in that $\|\Omega_1\|_0 = o(n/\log p)$, the lower bound for minimax detection rate is of the order

$$n^{-1/2} + \|\Omega_1\|_0 n^{-1} \log p.$$

However, when Ω_1 is ultra sparse in that $\|\Omega_1\|_0 = o(\sqrt{n}/\log p)$, then this lower bound is the parametric rate, i.e.

$$1/\sqrt{n}.$$

5.2. *Confidence intervals.* The theoretical results in Sections 3 and 4 also imply that the expected length of confidence intervals cannot be adapted to s if $s \gg \sqrt{n}/\log p$.

We denote by $\mathcal{C}_\alpha(\Theta_1)$ the set of all $(1 - \alpha)$ level confidence intervals for β over the parameter space Θ_1 constructed from the observed data \mathbf{D} :

$$(17) \quad \mathcal{C}_\alpha(\Theta_1) = \left\{ [l(\mathbf{D}), u(\mathbf{D})] : \inf_{\theta \in \Theta_1} \mathbb{P}_\theta(l(\mathbf{D}) \leq \beta \leq u(\mathbf{D})) \geq 1 - \alpha \right\}.$$

The construction in Section 4 yields the following confidence interval

$$\mathcal{CI}_* = \left[\hat{\beta} - c_n, \hat{\beta} + c_n \right],$$

where c_n is defined in (16). Theorem 4 implies that

$$\inf_{\theta \in \tilde{\Theta}(s)} \mathbb{P}_\theta(\beta \in \mathcal{CI}_*) \geq 1 - \alpha.$$

Since the diameter of the confidence set, $\text{diam}(\mathcal{CI}_*) = 2c_n$, Theorem 4 states the rate for c_n and thus implies the following minimax upper bound for the expected length of the confidence intervals:

$$(18) \quad \inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\tilde{\Theta}(s))} \sup_{\theta \in \tilde{\Theta}(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) \lesssim n^{-1/2} + sn^{-1} \log p,$$

where $\text{diam}(\mathcal{CI})$ denotes the length of \mathcal{CI} .

Since confidence intervals can be used to construct tests, minimax results on tests have implications for the minimax length of confidence intervals.

COROLLARY 6. *Suppose that $sn^{-1} \log p \leq 1/4$ and $2 \leq s \leq p^c$ for some constant $c < 1/2$. Then for any $\alpha \in (0, 1/3)$, we have*

$$\inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \Theta(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) \gtrsim n^{-1/2} + sn^{-1} \log p.$$

See Theorem 1 and Equation (3.14) of [Cai and Guo \(2017\)](#) for quantification of optimal confidence interval width for sparse or moderately sparse models i.e., $k \lesssim n/\log p$. Complementary, we allow for non-sparse vectors β , i.e., $k \lesssim p$.

Moreover, since $\Theta(s) \subset \tilde{\Theta}(s)$ and $\mathcal{C}_\alpha(\tilde{\Theta}(s)) \subset \mathcal{C}_\alpha(\Theta(s))$, we have

$$\begin{aligned} \inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\tilde{\Theta}(s))} \sup_{\theta \in \tilde{\Theta}(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) &\geq \inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\tilde{\Theta}(s))} \sup_{\theta \in \Theta(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) \\ &\geq \inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \Theta(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha). \end{aligned}$$

Therefore, Corollary 6 still holds if we replace $\Theta(s)$ by $\tilde{\Theta}(s)$. Combining this with (18), we obtain the minimax optimal rate for the expected length of confidence intervals over $\Theta(s)$ and $\tilde{\Theta}(s)$:

$$\inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \Theta(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) \asymp n^{-1/2} + sn^{-1} \log p$$

and

$$\inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\tilde{\Theta}(s))} \sup_{\theta \in \tilde{\Theta}(s)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) \asymp n^{-1/2} + sn^{-1} \log p.$$

In Theorem 4, we have constructed a minimax rate-optimal confidence interval for β in the case that the sparsity s is assumed to be known. A significant drawback of the construction is that it requires prior knowledge of s , which is typically unavailable in practice. Is it possible to construct adaptive confidence intervals that have the guaranteed coverage and automatically adjust the length to s ? In other words, does there exist a confidence interval in $\mathcal{C}_\alpha(\tilde{\Theta}(s))$ that has expected length of the order $n^{-1/2} + s_1 n^{-1} \log p$ over all $\tilde{\Theta}(s_1)$ and any $s_1 \ll s$? One consequence of the uniform non-testability result is that such adaptivity is not possible.

THEOREM 7. *Suppose that $sn^{-1} \log p \leq 1/4$ and $2 \leq s \leq p^c$ for some constant $c < 1/2$. Then for any $\alpha \in (0, 1/4)$ and $s_1 \leq s/2$, we have*

$$\inf_{\mathcal{CI}_\alpha \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \Theta(s_1)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_\alpha) \asymp n^{-1/2} + sn^{-1} \log p.$$

Even for $s_1 \ll s$, the optimal rate over $\Theta(s_1)$ for all confidence intervals that do not take into account knowledge of s_1 , is larger than that with the knowledge of s_1 . Therefore, Theorem 7 implies that for dense models ($k \lesssim p$), adaptivity with respect to s is in general not possible if Ω_1 is in at the least moderately sparse regime ($\sqrt{n}/\log p \ll s \lesssim n/\log p$).

5.3. Characterization of uniform non-testability. Here, we showcase that uniform non-testability is equivalent to the lack of adaptivity in all subsets of the parameter space.

Let Θ be a parameter space for a general model. We are interested in confidence intervals of $g(\theta)$, where g is an arbitrary functional of the whole parameter space, characterized by θ . For any $\Theta_1 \subseteq \Theta$, define the set of valid confidence intervals on Θ_1 :

$$\mathcal{C}_\alpha(\Theta_1) = \left\{ CI : \inf_{\theta \in \Theta_1} \mathbb{P}_\theta(g(\theta) \in CI) \geq 1 - \alpha \right\}.$$

For $\Theta_1 \subseteq \Theta$, the minimax rate over Θ_1 confidence intervals valid over Θ can be defined as

$$(19) \quad L(\Theta_1, \Theta) = \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta_1} \mathbb{E}_\theta \text{diam}(CI).$$

For $\Theta_1 \subseteq \Theta$, we say that there is no adaptivity between Θ and Θ_1 if

$$L(\Theta_1, \Theta) \asymp L(\Theta, \Theta).$$

In other words, if we use a confidence interval that is valid over the larger set Θ , then even on the smaller set Θ_1 , the length of the confidence interval has no improvement. For confidence intervals, we say that points in Θ are uniformly non-testable if

$$(20) \quad \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \asymp \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

In other words, the minimax confidence intervals have the same order of magnitude in terms of length for all the points in Θ . The following result establishes the link between uniform non-testability and adaptivity.

THEOREM 8. *The uniform non-testability*

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \asymp \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI)$$

if and only if there exists a constant $c > 0$ such that

$$cL(\Theta, \Theta) \leq L(\Theta_1, \Theta) \leq L(\Theta, \Theta)$$

for any subset $\Theta_1 \subseteq \Theta$.

Theorem 8 establishes that uniform non-testability simply means that there is no adaptivity between Θ and any subset of Θ . Hence, uniform non-testability provides a way of looking at adaptivity. Intuitively, adaptivity means that a procedure can automatically adapt its efficiency to the parameter. Since uniform non-testability means that the minimax optimal procedure has the same efficiency at each point in the parameter space, this rules out the possibility that the efficiency of the optimal procedure can change from parameter to parameter.

In the above setup, consider the testing problem

$$(21) \quad H_0 : g(\theta) = \tau, \quad \text{vs} \quad H_1 : g(\theta) = \tau + c_1 h_n.$$

For any $\tau \in \mathbb{R}$, $\Theta(\tau) = \{\theta \in \Theta : g(\theta) = \tau\}$, i.e., $\Theta(\tau)$ is the set of parameters θ satisfying the null hypothesis $H_0 : g(\theta) = \tau$. For any $\Theta_1 \subseteq \Theta$, let the set of valid tests of size α over Θ_1 be denoted by $\Psi_\alpha(\Theta_1)$, i.e.,

$$\Psi_\alpha(\Theta_1) = \{\psi : \sup_{\theta \in \Theta_1} \mathbb{E}_\theta \psi \leq \alpha\}.$$

Next, we showcase that the result of Theorem 8 applies to the hypothesis testing problems studied in Section 3.

COROLLARY 9. *Suppose that there exist constants $c_1, c_2 > 0$ and a confidence interval $CI_* \in \mathcal{C}_\alpha(\Theta)$ such that*

(1) for any $\tau \in \mathbb{R}$

$$\sup_{\psi \in \Psi_\alpha(\Theta(\tau))} \sup_{\theta \in \Theta(\tau + c_1 h_n)} \mathbb{E}_\theta \psi \leq 2\alpha,$$

(2)

$$\sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI_*) \leq c_2 h_n.$$

Then,

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \asymp \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \asymp h_n.$$

Condition (1) states that any α level test about (21) has power at most 2α whereas Condition (2) assumes a valid confidence interval with expected length of the order h_n , therefore h_n is a detection boundary. Corollary 9 then states that we have uniform non-testability in the sense of (20) and the optimal rate is h_n which, by Sections 3 and 4 is $h_n = n^{-1/2} + sn^{-1} \log p$.

Theorem 8 suggests a much broader implication. Since the uniform non-testability implies lack of adaptivity with respect to any subset of the parameter space, our result indicates that it is impossible for a confidence interval to automatically exploit other structures of the model. In particular, if a confidence interval is valid on Θ , then it will have the same rate even at points with special structures, e.g., sparsity, homogeneity, etc. Hence, our result not only states that there is no adaptivity with respect to s , we show that there cannot be any adaptivity with respect to any structure.

6. Impact of an increasing $\|\beta\|_2$. We now discuss the case in which the ℓ_2 -norm of β for the model (2) is allowed to grow. To explicitly write out the dependence on $\|\beta\|_2$ i.e., M_2 , we introduce the notation

$$\tilde{\Theta}_{M_1, M_2}(s) = \left\{ \theta = (\beta, \gamma, \Sigma, \sigma) : M^{-1} \leq \lambda_{\min}(\Sigma) \leq \lambda_{\max}(\Sigma) \leq M, \right.$$

$$\left. \|\boldsymbol{\Omega}_1\|_0 \leq s, 0 \leq \sigma \leq M_1, \|\boldsymbol{\beta}\|_2 \leq M_2 \right\},$$

where $M > 1$ is a constant. Now define the minimax length

$$\mathbb{A}(s, M_1, M_2) = L\left(\tilde{\Theta}_{M_1, M_2}(s), \tilde{\Theta}_{M_1, M_2}(s)\right),$$

where $L(\cdot, \cdot)$ is defined in (19). The following result states a scaling property that allows us to derive the minimax result for dense models with growing ℓ_2 -norm of the parameter.

THEOREM 10. *For any constants $Q, M_1, M_2 > 0$,*

$$\mathbb{A}(s, QM_1, QM_2) = Q\mathbb{A}(s, M_1, M_2).$$

By Theorem 10, it suffices to derive $\mathbb{A}(s, M_0, 1)$ for all $M_0 > 0$. This is because $\mathbb{A}(s, M_1, M_2) = M_2\mathbb{A}(s, M_0, 1)$ with $M_0 = M_1M_2^{-1}$.

For that end we consider a specific asymptotic regime where M_2 is considered fixed while M_0 is allowed to grow to infinity or to shrink to zero. Hence, for an arbitrary constant $C > 0$, in the duration of this section, we assume that $M_2 > C$.

Upper bound is obtain as a corollary of Theorem 4, from which we can easily conclude

$$\mathbb{A}(s, M_0, 1) \leq C_1(M_0 + 1)(n^{-1/2} + sn^{-1} \log p),$$

where $C_1 > 0$ is a constant independent of M_1, M_2, n, s or p .

To establish a lower bound, we establish the following result.

THEOREM 11. *Assume that $p \geq 2n + 1$. If $\alpha \in (0, 1/3)$, then there exists a constant $C > 0$, depending only on α , such that*

$$\mathbb{A}(1, 0, 1) \geq Cn^{-1/2}.$$

Theorem 11 considers a particularly simple setting where the model has no noise and the sparsity of the precision matrix, $s = 1$. Even in this simple case, Theorem 11 suggests the following: Even if the noise level σ is zero, perfect recovery of β is not possible as long as the vector γ is allowed to be non-sparse with bounded ℓ_2 -norm.

First implication of this result is that imposing bounded ℓ_2 -norm is weaker than imposing sparsity. When the model parameter γ is assumed to be sparse $\|\gamma\|_0 \lesssim n/\log p$ and the noise level is zero, one can invoke classical results (e.g., Bickel et al. (2009); Raskutti et al. (2011)) and obtain that

exact recovery of the model parameter is achievable. However, Theorem 11 says that no estimator can exactly recover dense signals that are only known to have bounded ℓ_2 -norm. This is true even if the covariance matrix of the design is known to be diagonal; in fact, by inspecting the proof, the same result holds even if the covariance matrix is known to be \mathbb{I}_p . Therefore, the difficulty of identifying dense signals is quite fundamental (even if their ℓ_2 -norm is bounded) and is not due to noise in the response or to unknown distribution of the design.

Second implication of Theorem 11 is a lower bound of $\mathbb{A}(s, M_0, 1)$. Notice that for any non-singular covariance matrix, its inverse always has non-zero diagonal entries, which means $s \geq 1$. Hence, Theorem 11 implies that

$$\mathbb{A}(s, M_0, 1) \geq \mathbb{A}(1, 0, 1) \geq Cn^{-1/2}.$$

To sum up the above bounds, we have

COROLLARY 12. *Suppose that $p \geq 2n + 1$, $sn^{-1} \log p \leq 1/4$ and $1 \leq s \leq p^c$ for some constant $c < 1/2$. Then for any $\alpha \in (0, 1/4)$ and $M_2 \geq C$, we have*

$$C_1 M_2 n^{-1/2} + C_1 s n^{-1} \log p \leq \mathbb{A}(s, M_1, M_2) \leq C_2 M_2 (n^{-1/2} + s n^{-1} \log p),$$

where $C_1, C_2 > 0$ are constants depending only on M, M_1, α, C, c .

Corollary 12 outlines a unique phenomenon for dense high-dimensional models. Since efficiency for testing dense models depends on the ℓ_2 -norm of the model parameter, consistency is impossible if this magnitude in ℓ_2 -norm is of the order larger than \sqrt{n} . In contrast, for models with sparse parameters, results in (Cai and Guo, 2017; Javanmard and Montanari, 2018) show that ℓ_2 -norm requirements are not required.

An interesting, and yet challenging question arises from the above study: What is the exact minimax lower bound as a function of the ℓ_2 -norm for high-dimensional and dense models? We leave this investigation to future research.

7. Discussion. This paper establishes theoretical results for hypothesis testing problems in high-dimensional linear models. Our work pushes the frontier of high-dimensional inference by allowing the model sparsity to be arbitrary. We derive the optimal detection rates and show that the accuracy of statistical inference without imposing model sparsity depends on the ability to decorrelate the features. The sparsity of the first row of the precision matrix controls the optimal detection rate; for sparse enough

precision matrices, the parametric rate can be achieved. These results also provide additional insights into the adaptivity of optimal inference.

The theoretical development in this paper has potential implications beyond minimax detection rates. In particular, we show that the detection rate for *every* point in the alternative is the same, and thus the derived detection rate is uniform over the alternative, which indicates that a simple coin toss is a uniformly most powerful test asymptotically. For this reason, the detection rates established in this paper are driven by the fundamental difficulty that cannot be adequately described under the general minimax framework.

Some important extensions and refinements are left open. Our current results only provide confidence intervals and testability results regarding univariate parameters; extending our theory to the setting of global testing and especially multivariate testing, seems like a promising avenue for further work. Another challenge is that many new hypothesis testing problems have complex structures and some even non-convex boundaries. A systematic approach to studying such problems would improve and extend the current scope of inferential theoretical results. In general, work can be done to identify a subset of points for which attainable and optimal tests can be developed, in turn, paving the way for new inferential methods.

SUPPLEMENT contains the detailed proofs of all auxiliary lemmas as well as details of the proofs of Theorems 1, 4, 7, 8 and 10 as well as Corollaries 5, 6 and 9. Below we present proofs of Theorems 2, 3 and 11.

APPENDIX A: PROOF OF THEOREM 2

Proof of Theorem 2 has been split into a sequence of smaller results. First we present some notation, then auxiliary Lemmas 1-6 that are useful in the proof of Theorem 2 and lastly the proof of the result itself.

A.1. Notations. In the rest of Appendix A, we introduce the following notation. We utilize Lemma 6 below to pinpoint the structure of the covariance matrices Σ that are of interest to us, i.e., the lower-right corner is equal to \mathbb{I}_{p-1} .

Namely, we show that for any point $\theta = (\beta, \gamma, \Sigma, \sigma) \in \Theta$, we can write Σ as

$$\Sigma = \begin{pmatrix} \boldsymbol{\pi}^\top \boldsymbol{\pi} + \sigma_{\mathbf{V}}^2 & \boldsymbol{\pi}^\top \\ \boldsymbol{\pi} & \mathbb{I}_{p-1} \end{pmatrix},$$

where $\boldsymbol{\pi}$ is a suitably chosen vector and $\sigma_{\mathbf{V}}^2$ is a suitably chosen constant that is positive. This is equivalent to working with the vector $\boldsymbol{\pi} \in \mathbb{R}^{p-1}$ from

the following regression,

$$\mathbf{Z} = \mathbf{W}\boldsymbol{\pi} + \mathbf{V}$$

for a vector of residuals $\mathbf{V} \in \mathbb{R}^n$. Coincidentally, $\sigma_{\mathbf{V}}$ will be the standard deviation of the residuals \mathbf{V} . Recall that $\mathbb{E}_{\theta}[\mathbf{W}^{\top}\mathbf{W}]/n = \mathbb{I}_{p-1}$ for $\theta \in \Theta$. We also define a matrix L_{θ} as follows

$$(22) \quad L_{\theta} = L(\theta) = \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ \boldsymbol{\pi}^{\top} & \sigma_{\mathbf{V}} & 0 \\ (\boldsymbol{\pi}\beta + \boldsymbol{\gamma})^{\top} & \beta\sigma_{\mathbf{V}} & \sigma \end{pmatrix}.$$

From Lemma 6 below we know that the space of correlation matrices is spanned by the collection of $\boldsymbol{\Sigma}$'s as described above. Notice that under \mathbb{P}_{θ} , vector $(\mathbf{W}_i^{\top}, Z_i, y_i)^{\top} \in \mathbb{R}^{p+1}$ has gaussian distribution $\mathcal{N}(0, L_{\theta}L_{\theta}^{\top})$.

The plan of the proof proceeds as follows. We pick an arbitrary test $\psi_{*} \in \Psi_{\alpha}(\Theta(s, \beta_0))$ and an arbitrary point in the alternative

$$(23) \quad \theta_{*} = (\beta_{*}, \boldsymbol{\gamma}_{*}, \boldsymbol{\Sigma}_{*}, \sigma_{\varepsilon,*}) \in \Theta_{\zeta,\kappa}(s/2, \beta_0 + h_n),$$

where $h_n = \rho sn^{-1} \log p$ and

$$\boldsymbol{\Sigma}_{*} = \begin{pmatrix} \boldsymbol{\pi}_{*}^{\top} \boldsymbol{\pi}_{*} + \sigma_{\mathbf{V},*}^2 & \boldsymbol{\pi}_{*}^{\top} \\ \boldsymbol{\pi}_{*} & \mathbb{I}_{p-1} \end{pmatrix}.$$

We then construct points based on θ_{*} according to Definition 3 below. Observe that these points are chosen to be dependent on the alternative. Then we proceed to show that (1) these points are in the null space $\Theta(s, \beta_0)$ and (2) the average χ^2 -distance between these points and θ_{*} is small. Therefore, the power of ψ_{*} against θ_{*} is close to the average power against these points. Since these points are in the null space, the power against them is at most equal to the nominal level. As a result, the power against θ_{*} is also close to the nominal level.

In the rest of Appendix A, we denote $m = \lfloor s/2 \rfloor$ and define the set of all m -sparse vectors with entries taking values in $\{0, 1\}$ as \mathcal{M} , i.e.,

$$\mathcal{M} = \{\boldsymbol{\delta} \in \{0, 1\}^{p-1} : \|\boldsymbol{\delta}\|_0 = m\}$$

and let N denote the cardinality of \mathcal{M} . Clearly, $N = \binom{p-1}{m}$. We list \mathcal{M} as $\mathcal{M} = \{\boldsymbol{\delta}_{(1)}, \dots, \boldsymbol{\delta}_{(N)}\}$, i.e., $\boldsymbol{\delta}_{(j)}$ denotes the element j of the set \mathcal{M} .

DEFINITION 3. *Given $\theta_{*} \in \Theta_{\zeta,\kappa}(s/2, \beta_0 + h_n)$ as in (23), let $0 \leq d \leq \rho$ be such that $\beta_{*} = \beta_0 + h$ with $h = dsn^{-1} \log p$. Let $r = \sigma_{\mathbf{V},*}/\sigma_{\varepsilon,*} > 0$. For $j \in \{1, \dots, N\}$, define*

$$\theta_j = (\beta_0, \boldsymbol{\gamma}_{(j)}, \boldsymbol{\Sigma}_{(j)}, \sigma_{\varepsilon,0})$$

with

$$\begin{aligned}\beta_0 &= \beta_* - h \\ \boldsymbol{\gamma}_{(j)} &= \boldsymbol{\gamma}_* + h\boldsymbol{\pi}_{(j)} + r(1-h)\sigma_{\varepsilon,*}\sqrt{h/m}\boldsymbol{\delta}_{(j)} \\ \sigma_{\varepsilon,0} &= \sigma_{\varepsilon,*}\sqrt{1-hr^2+h^2r^2} \\ \boldsymbol{\Sigma}_{(j)} &= \begin{pmatrix} \boldsymbol{\pi}_{(j)}^\top \boldsymbol{\pi}_{(j)} + \sigma_{\mathbf{V},0}^2 & \boldsymbol{\pi}_{(j)}^\top \\ \boldsymbol{\pi}_{(j)} & \mathbb{I}_{p-1} \end{pmatrix}\end{aligned}$$

where

$$\boldsymbol{\pi}_{(j)} = \boldsymbol{\pi}_* + \sigma_{\mathbf{V},*}\sqrt{h/m}\boldsymbol{\delta}_{(j)}$$

and $\sigma_{\mathbf{V},0} = \sigma_{\mathbf{V},*}\sqrt{1-h}$.

A.2. Auxiliary results. Below we present useful auxiliary results.

LEMMA 1. For a constant $c \in (0, 1/2)$, let the sequence (m, n, p) be such that $1 \leq m \leq p^c$ as well as $mn^{-1} \log p \leq 1/4$ as $p \rightarrow \infty$. Then for any $a \in (0, (1-2c)/4)$,

$$\sum_{k=0}^m [1 - kan^{-1} \log p]^{-n} \frac{\binom{m}{k} \binom{p-m-1}{m-k}}{\binom{p-1}{m}} \leq 1 + o(1).$$

The next two results are useful for computing χ^2 -distance.

LEMMA 2. Let g_j denote the probability density function of $\mathcal{N}(0, \boldsymbol{\Sigma}_j)$ with nonsingular $\boldsymbol{\Sigma}_j \in \mathbb{R}^{k \times k}$ for $j = 0, 1, 2$. Suppose that $\boldsymbol{\Sigma}_j$ can be decomposed as $\boldsymbol{\Sigma}_j = L_j L_j^\top$. Then

$$\mathbb{E}_{g_0} \left(\frac{d\mathbb{P}_{g_1}}{d\mathbb{P}_{g_0}} \times \frac{d\mathbb{P}_{g_2}}{d\mathbb{P}_{g_0}} \right) = \frac{1}{\sqrt{\det(\mathbb{I}_k - [Q_1 Q_1^\top - \mathbb{I}_k] [Q_2 Q_2^\top - \mathbb{I}_k])}},$$

where $Q_j = L_0^{-1} L_j$ for $j = 1, 2$.

LEMMA 3. Consider the notations in Definition 3. Let $Q_j = L_{\theta_*}^{-1} L_{\theta_j}$. Then for any $j_1, j_2 \in \{1, \dots, N\}$,

$$\begin{aligned}\det \left[\mathbb{I}_{p+1} - \left(Q_{j_1} Q_{j_1}^\top - \mathbb{I}_{p+1} \right) \left(Q_{j_2} Q_{j_2}^\top - \mathbb{I}_{p+1} \right) \right] \\ = \left[1 - m^{-1} h [r^2 (1-h)^2 + 1] \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \right]^2.\end{aligned}$$

With the help of Lemmas 1, 2 and 3, we can provide the next result concerning the distance between the null and alternative hypothesis.

LEMMA 4. *Consider $\theta_* \in \Theta_{\zeta, \kappa}(m, \beta_0 + \rho sn^{-1} \log p)$ as defined in the proof of Theorem 2. Consider $\{\theta_j\}_{j=1}^N$ defined in Definition 3. Let ρ be defined as in (8). Then*

$$\limsup_{n \rightarrow \infty} \mathbb{E}_{\theta_*} \left(N^{-1} \sum_{j=1}^N \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - 1 \right)^2 = 0.$$

Note that θ_j is a function of ρ and γ . Next, we show that the designed points, θ_j belong to the the null parameter space.

LEMMA 5. *Consider $\theta_* \in \Theta_{\zeta, \kappa}(m, \beta_0 + h_n)$ with $h_n = \rho sn^{-1} \log p$ in the proof of Theorem 2. Consider $\{\theta_j\}_{j=1}^N$ defined in Definition 3. Suppose that the conditions in the statement of Theorem 2 hold. Then*

$$\{\theta_j : 1 \leq j \leq N\} \subset \Theta_{\zeta, \kappa}(2m, \beta_0).$$

Lastly, the following lemma describes the structure of the covariance matrices.

LEMMA 6. *Consider any $a > 0$ and $b \in \mathbb{R}^{p-1}$. Let Σ be a positive definite matrix. If all the eigenvalue of $\begin{pmatrix} a & b^\top \Sigma \\ \Sigma b & \Sigma \end{pmatrix}$ are positive, then $a > b^\top \Sigma b$.*

In particular, if all the eigenvalues of $\begin{pmatrix} a & b^\top \\ b & \mathbb{I}_{p-1} \end{pmatrix}$ are positive, then $a > b^\top b$.

Now, that all of the auxiliary results are established, we are ready to present the main proof.

A.3. Proof of Theorem 2. The proof methodology is novel in that for each possible candidate point in the alternative, we need to design a sequence of points in the null space and demonstrate that their χ^2 -distances to the candidate point in the alternative will be small therefore limiting the power of the test.

Proof of Theorem 2. Recall that m denotes the largest integer not exceeding $s/2$, i.e., $m = \lfloor s/2 \rfloor$. Fix any $\eta > 0$. Recall ρ defined in (8).

Observe that by the properties of the supremum, we can choose $\psi_* \in \Psi_\alpha(\Theta(s, \beta_0))$ and

$$\theta_* = (\beta_*, \gamma_*, \Sigma_*, \sigma_{\varepsilon,*}) \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)$$

with $h_n = \rho sn^{-1} \log p$ such that

$$(24) \quad \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi \leq \mathbb{E}_{\theta_*} \psi_* + \eta.$$

Since $\|\cdot\|_0$ can only take values in \mathbb{Z} , $\Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n) = \Theta_{\zeta, \kappa}(m, \beta_0 + h_n)$. By Lemma 6, there exist $\sigma_{\mathbf{V},*} > 0$ and $\boldsymbol{\pi}_* \in \mathbb{R}^{p-1}$ such that

$$\boldsymbol{\Sigma}_* = \begin{pmatrix} \boldsymbol{\pi}_*^\top \boldsymbol{\pi}_* + \sigma_{\mathbf{V},*}^2 & \boldsymbol{\pi}_*^\top \\ \boldsymbol{\pi}_* & \mathbb{I}_{p-1} \end{pmatrix}.$$

We construct $\{\theta_j\}_{j=1}^N$ as in the Definition 3.

Since $\mathbb{E}_{\theta_j} \psi_* = \mathbb{E}_{\theta_*} \left(\psi_* \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} \right)$, it follows that

$$\begin{aligned} & \left| N^{-1} \sum_{j=1}^N \mathbb{E}_{\theta_j} \psi_* - \mathbb{E}_{\theta_*} \psi_* \right| \\ &= \left| N^{-1} \sum_{j=1}^N \left(\mathbb{E}_{\theta_*} \psi_* \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - \mathbb{E}_{\theta_*} \psi_* \right) \right| = \left| \mathbb{E}_{\theta_*} \psi_* \left(N^{-1} \sum_{j=1}^N \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - 1 \right) \right| \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\theta_*} \left| N^{-1} \sum_{j=1}^N \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - 1 \right| \stackrel{(ii)}{\leq} \left[\mathbb{E}_{\theta_*} \left(N^{-1} \sum_{j=1}^N \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - 1 \right)^2 \right]^{1/2}, \end{aligned}$$

where (i) holds by $|\psi_*| \leq 1$ almost surely and (ii) holds by Lyapunov's inequality.

By Lemma 4 and the above display, we have

$$\limsup_{n \rightarrow \infty} \left| N^{-1} \sum_{j=1}^N \mathbb{E}_{\theta_j} \psi_* - \mathbb{E}_{\theta_*} \psi_* \right| = 0.$$

By Lemma 5, $\theta_j \in \Theta_{\zeta, \kappa}(2m, \beta_0) \subseteq \Theta(s, \beta_0)$ for all $j \in \{1, \dots, N\}$. This and the fact that $\psi_* \in \Psi_\alpha(\Theta(s, \beta_0))$ imply

$$N^{-1} \sum_{j=1}^N \mathbb{E}_{\theta_j} \psi_* \leq \alpha.$$

Hence, $\limsup_{n \rightarrow \infty} \mathbb{E}_{\theta_*} \psi_* \leq \alpha$. By (24), we have

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi \leq \limsup_{n \rightarrow \infty} \mathbb{E}_{\theta_*} \psi_* + \eta \leq \alpha + \eta.$$

Moreover, since $\eta > 0$ is arbitrary, we have

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi \leq \alpha.$$

Notice that a random test that rejects the hypothesis at random with probability α has power equal to α . Since $\Psi_\alpha(\Theta(s, \beta_0))$ includes such random tests, the above inequality holds with equality. The proof is complete. \square

APPENDIX B: PROOF OF THEOREM 3

PROOF OF THEOREM 3. Here we will show that no test can be better than the Likelihood Ratio test.

Recall $\tau = \kappa \sqrt{M^{-1} \log(1 + \alpha^2)}$. We choose an arbitrary test

$$\psi_{**} \in \Psi_\alpha(\Theta(s, \beta_0))$$

and an arbitrary point

$$\theta_{**} = (\beta_{**}, \gamma_{**}, \Sigma_{**}, \sigma_{\varepsilon, **}) \in \Theta_\kappa(s, \beta_0 + h_n)$$

with $h_n = \tau n^{-1/2}$. Throughout this proof we denote with θ_{**} the point in the alternative space. Notice that $\beta_{**} = \beta_0 + h$ with $0 \leq h \leq h_n$.

We define

$$\theta_0 = (\beta_0, \gamma_{**}, \Sigma_{**}, \sigma_{\varepsilon, **}).$$

Clearly,

$$\theta_0 \in \Theta(s, \beta_0) \quad \text{and thus} \quad E_{\theta_0} \psi_{**} \leq \alpha.$$

Recall the notation $\mathbf{X}_i = (Z_i, \mathbf{W}_i^\top)^\top \in \mathbb{R}^p$. Let $\sigma_z^2 = \mathbb{E}_{\theta_0} Z_i^2$. By the definition of $\Theta_\kappa(s, \beta_0)$,

$$(25) \quad \sigma_z^2 \leq \lambda_{\max}(\Sigma_{**}) \leq M \quad \text{and} \quad \sigma_{\varepsilon, **} \geq \kappa^2.$$

Then the likelihood of the data under $\mathbb{P}_{\theta_{**}}$ can be written as a product of the likelihood of \mathbf{y} given \mathbf{X} and the likelihood of \mathbf{X} :

$$\left[\frac{1}{(\sqrt{2\pi}\sigma_{\varepsilon, **})^n} \exp \left(-\frac{1}{2\sigma_{\varepsilon, **}^2} \sum_{i=1}^n (y_i - Z_i \beta_{**} - \mathbf{W}_i^\top \gamma_{**})^2 \right) \right] \\ \times \left[\frac{1}{(\sqrt{\det(2\pi\Sigma_{**})})^n} \exp \left(-\frac{1}{2} \sum_{i=1}^n \mathbf{X}_i^\top \Sigma_{**}^{-1} \mathbf{X}_i \right) \right].$$

Similarly, the likelihood of the data under \mathbb{P}_{θ_0} can be written as

$$\left[\frac{1}{(\sqrt{2\pi}\sigma_{\varepsilon,**})^n} \exp\left(-\frac{1}{2\sigma_{\varepsilon,**}^2} \sum_{i=1}^n (y_i - Z_i\beta_0 - \mathbf{W}_i^\top \boldsymbol{\gamma}_{**})^2\right) \right] \\ \times \left[\frac{1}{(\sqrt{\det(2\pi\boldsymbol{\Sigma}_{**})})^n} \exp\left(-\frac{1}{2} \sum_{i=1}^n \mathbf{X}_i^\top \boldsymbol{\Sigma}_{**}^{-1} \mathbf{X}_i\right) \right].$$

Hence, the likelihood ratio can be written as

$$\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} \\ = \exp\left(\frac{1}{2\sigma_{\varepsilon,**}^2} \sum_{i=1}^n \left[(y_i - Z_i\beta_0 - \mathbf{W}_i^\top \boldsymbol{\gamma}_{**})^2 - (y_i - Z_i\beta_{**} - \mathbf{W}_i^\top \boldsymbol{\gamma}_{**})^2 \right]\right) \\ (26) \\ \stackrel{(i)}{=} \exp\left(\frac{h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i \left[y_i - Z_i(\beta_0 + h/2) - \mathbf{W}_i^\top \boldsymbol{\gamma}_{**} \right]\right),$$

where (i) follows by $\beta_{**} = \beta_0 + h$. Thus,

$$\begin{aligned} |\mathbb{E}_{\theta_0} \psi_{**} - \mathbb{E}_{\theta_{**}} \psi_{**}| &= \left| \mathbb{E}_{\theta_0} \psi_{**} - \mathbb{E}_{\theta_0} \psi_{**} \frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} \right| \\ &= \left| \mathbb{E}_{\theta_0} \psi_{**} \left(\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} - 1 \right) \right| \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\theta_0} \left| \frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} - 1 \right| \\ (27) \quad &\leq \sqrt{\mathbb{E}_{\theta_0} \left(\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} - 1 \right)^2} = \sqrt{\mathbb{E}_{\theta_0} \left(\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} \right)^2} - 1, \end{aligned}$$

where (i) follows by $|\psi_{**}| \leq 1$. By (26), we have

$$\begin{aligned} \mathbb{E}_{\theta_0} \left(\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} \right)^2 &= \mathbb{E}_{\theta_0} \left[\exp\left(\frac{2h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i \left[y_i - Z_i(\beta_0 + h/2) - \mathbf{W}_i^\top \boldsymbol{\gamma}_{**} \right]\right) \right] \\ &\stackrel{(i)}{=} \mathbb{E}_{\theta_0} \left[\exp\left(\frac{2h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i [\varepsilon_i - Z_i h/2]\right) \right] \\ &= \mathbb{E}_{\theta_0} \left\{ \mathbb{E}_{\theta_0} \left[\exp\left(\frac{2h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i [\varepsilon_i - Z_i h/2]\right) \mid \mathbf{Z} \right] \right\} \\ (28) \quad &= \mathbb{E}_{\theta_0} \left\{ \mathbb{E}_{\theta_0} \left[\exp\left(\frac{2h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i \varepsilon_i\right) \mid \mathbf{Z} \right] \exp\left(-\frac{h^2}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i^2\right) \right\} \end{aligned}$$

where (i) follows by the fact that under \mathbb{P}_{θ_0} , $y_i = Z_i\beta_0 + \mathbf{W}_i^\top \boldsymbol{\gamma}_{**} + \varepsilon_i$.

Notice that under \mathbb{P}_{θ_0} , $\sum_{i=1}^n Z_i\varepsilon_i$ conditional on \mathbf{Z} has a Gaussian distribution with mean 0 and variance equal to $\sum_{i=1}^n Z_i^2\sigma_{\varepsilon,**}^2$. Hence, by the moment generating function of Gaussian distributions, it follows that

$$\mathbb{E}_{\theta_0} \left[\exp \left(\frac{2h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i\varepsilon_i \right) \mid \mathbf{Z} \right] = \exp \left(2\sigma_{\varepsilon,**}^{-2} h^2 \sum_{i=1}^n Z_i^2 \right).$$

Therefore, we can use the above display to continue (28) and obtain

$$\begin{aligned} \mathbb{E}_{\theta_0} \left(\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} \right)^2 &= \mathbb{E}_{\theta_0} \left\{ \mathbb{E}_{\theta_0} \left[\exp \left(\frac{2h}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i\varepsilon_i \right) \mid \mathbf{Z} \right] \exp \left(-\frac{h^2}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i^2 \right) \right\} \\ &= \mathbb{E}_{\theta_0} \left\{ \exp \left(2\sigma_{\varepsilon,**}^{-2} h^2 \sum_{i=1}^n Z_i^2 \right) \exp \left(-\frac{h^2}{\sigma_{\varepsilon,**}^2} \sum_{i=1}^n Z_i^2 \right) \right\} \\ &= \mathbb{E}_{\theta_0} \left[\exp \left(\sigma_{\varepsilon,**}^{-2} h^2 \sum_{i=1}^n Z_i^2 \right) \right] \\ &= \mathbb{E}_{\theta_0} \left[\exp \left(\sigma_z^2 \sigma_{\varepsilon,**}^{-2} h^2 \sum_{i=1}^n (Z_i^2 \sigma_z^{-2}) \right) \right] \\ &\stackrel{(i)}{\leq} \mathbb{E}_{\theta_0} \left[\exp \left([\log(1 + \alpha^2)] n^{-1} \sum_{i=1}^n (Z_i^2 \sigma_z^{-2}) \right) \right] \\ &\stackrel{(ii)}{=} (1 - 2n^{-1} \log(1 + \alpha^2))^{-n/2} \\ &\stackrel{(iii)}{\leq} \exp [\log(1 + \alpha^2)] = 1 + \alpha^2, \end{aligned}$$

where (i) follows by $0 \leq h \leq h_n = \tau n^{-1/2} = n^{-1/2} \kappa \sqrt{M^{-1} \log(1 + \alpha^2)}$ and (25), (ii) follows by the moment generating function of $\chi^2(n)$ (chi-squared distribution with n degrees of freedom) and the fact that

$$\sum_{i=1}^n Z_i^2 \sigma_z^{-2}$$

has a $\chi^2(n)$ distribution together with $n^{-1} \log(1 + \alpha^2) < 1/2$ (due to $\alpha^2 < 1/4$ and $\log(1.25) < 1/2$) and (iii) follows by the fact that

$$(1 - a/n)^{-n/2} \leq \exp(a/2)$$

for any $n \geq 1$ and $a \geq 0$.

Therefore, the above display and (27) imply that

$$|\mathbb{E}_{\theta_0} \psi_{**} - \mathbb{E}_{\theta_{**}} \psi_{**}| \leq \sqrt{\mathbb{E}_{\theta_0} \left(\frac{d\mathbb{P}_{\theta_{**}}}{d\mathbb{P}_{\theta_0}} \right)^2 - 1} = \sqrt{\alpha^2} = \alpha.$$

Since $\mathbb{E}_{\theta_0} \psi_{**} \leq \alpha$, it follows that $\mathbb{E}_{\theta_{**}} \psi_{**} \leq 2\alpha$. Since ψ_{**} and θ_{**} are chosen arbitrarily, the desired result follows. \square

APPENDIX C: PROOF OF THEOREM 11

Proof of Theorem 11 has been split into a sequence of smaller results. First we present some notation, then auxiliary Lemmas 7 - 10 that are useful in the proof of Theorem 11 and lastly the proof of the result itself.

We first recall the notions of total variation and KL divergence. Given two probability measures \mathbb{P}_0 and \mathbb{P}_1 that are absolutely continuous with each other, we define the total variation

$$\text{TV}(\mathbb{P}_0, \mathbb{P}_1) = \frac{1}{2} \int \left| \frac{d\mathbb{P}_1}{d\mathbb{P}_0} - 1 \right| d\mathbb{P}_0$$

and KL divergence:

$$\text{KL}(\mathbb{P}_0, \mathbb{P}_1) = \int \left(\log \frac{d\mathbb{P}_0}{d\mathbb{P}_1} \right) d\mathbb{P}_0.$$

C.1. Auxiliary results.

LEMMA 7. *Let \mathbb{P}_0 and \mathbb{P}_1 denote the probability measures for $\mathcal{N}(\mu_0, \Sigma_0)$ and $\mathcal{N}(\mu_1, \Sigma_1)$, respectively. Then*

$$\text{KL}(\mathbb{P}_0, \mathbb{P}_1) = \frac{1}{2} \left(\text{trace}(\Sigma_1^{-1}(\Sigma_0 - \Sigma_1)) + \log \left(\frac{\det(\Sigma_1)}{\det(\Sigma_0)} \right) + (\mu_1 - \mu_0)^\top \Sigma_1^{-1}(\mu_1 - \mu_0) \right).$$

The proof of Lemma 7 follows by straight-forward computation and is thus omitted. The next two results are useful bounding tools.

LEMMA 8. *Let $\mathbf{W} \in \mathbb{R}^{n \times 2n}$ and $\mathbf{Z} \in \mathbb{R}^n$ have entries being i.i.d standard normal random variables. Then for any $a > 0$*

$$\mathbb{P} \left(\mathbf{Z}^\top (\mathbf{W}\mathbf{W}^\top)^{-1} \mathbf{Z} > a \right) \leq 2 \exp(-0.005n) + 12/a.$$

LEMMA 9. *Let $\boldsymbol{\xi}$ be a random vector with distribution $\mathcal{N}(0, \Sigma)$. Then for any $x > 0$,*

$$\mathbb{P}(\|\boldsymbol{\xi}\|_2 > x) \leq x^{-2} \text{trace}(\Sigma).$$

The main lemma utilized in the proof is the following one.

LEMMA 10. *Assume that $p \geq 2n + 1$. For any $r \in \mathbb{R}$, we define*

$$\Theta_*(r) = \{\theta = (\beta, \gamma, \Sigma, \sigma) : \beta = r, \|\gamma\|_2 \leq 1, \Sigma = \mathbb{I}_p, \sigma = 0\}.$$

Then there exists a constant $K > 0$ depending only on α such that

$$\inf_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{E}_\theta \psi \leq 1 - 2\alpha$$

for any measurable function of the data $(\mathbf{y}, \mathbf{W}, \mathbf{Z})$ satisfying $|\psi(\mathbf{y}, \mathbf{W}, \mathbf{Z})| \leq 1$ and $\sup_{\theta \in \Theta_(0)} \mathbb{E}_\theta \psi \leq \alpha$.*

C.2. Proof of Theorem 11.

PROOF OF THEOREM 11. Let

$$CI(\mathbf{y}, \mathbf{Z}, \mathbf{W}) = [l(\mathbf{y}, \mathbf{Z}, \mathbf{W}), u(\mathbf{y}, \mathbf{Z}, \mathbf{W})]$$

be a confidence set for β with nominal coverage probability $1 - \alpha$ over $\tilde{\Theta}_{0,1}(1)$. Define $\psi(\mathbf{y}, \mathbf{Z}, \mathbf{W}) = \mathbf{1}\{0 \notin CI(\mathbf{y}, \mathbf{Z}, \mathbf{W})\}$.

From now on, we will write CI , ψ , u and l without $(\mathbf{y}, \mathbf{Z}, \mathbf{W})$ to simplify the notation.

Recall the notation $\Theta_*(r)$ from Lemma 10. Since $\Theta_*(0) \subset \tilde{\Theta}_{0,1}(1)$, we have

$$\sup_{\theta \in \Theta_*(0)} \mathbb{E}_\theta \psi \leq \alpha.$$

Moreover, by the same Lemma 10,

$$\inf_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{E}_\theta \psi \leq 2\alpha$$

for some constant $K > 0$ depending only on α . This means that

$$\inf_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{P}_\theta(l \leq 0 \leq u) \geq 1 - 2\alpha.$$

Since $\Theta_*(n^{-1/2}K) \subset \tilde{\Theta}_{0,1}(1)$, we have that

$$\inf_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{P}_\theta(l \leq n^{-1/2}K \leq u) \geq 1 - \alpha.$$

Therefore,

$$\inf_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{P}_\theta(l \leq 0 < n^{-1/2}K \leq u) \geq 1 - 3\alpha.$$

It follows that

$$\inf_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{P}_\theta \left(u - l \geq n^{-1/2}K \right) \geq 1 - 3\alpha$$

and thus

$$\begin{aligned} \sup_{\theta \in \tilde{\Theta}_{0,1}(1)} \mathbb{E}_\theta(u - l) &\geq \sup_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{E}_\theta(u - l) \\ &\geq \sup_{\theta \in \Theta_*(n^{-1/2}K)} \mathbb{E}_\theta \left[\mathbf{1} \left\{ u - l \geq n^{-1/2}K \right\} \times n^{-1/2}K \right] \\ &\geq n^{-1/2}K \times (1 - 3\alpha). \end{aligned}$$

Since the above bound holds for any confidence interval CI , the proof is complete. \square

REFERENCES

- Belloni, A., Chernozhukov, V., and Hansen, C. (2014a). Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650.
- Belloni, A., Chernozhukov, V., and Kato, K. (2014b). Uniform post-selection inference for least absolute deviation regression and other z-estimation problems. *Biometrika*, 102(1):77–94.
- Bickel, P. J., Ritov, Y., and Tsybakov, A. B. (2009). Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732.
- Boucheron, S., Lugosi, G., and Massart, P. (2013). *Concentration inequalities: A nonasymptotic theory of independence*. Oxford university press.
- Boyle, E. A., Li, Y. I., and Pritchard, J. K. (2017). An expanded view of complex traits: from polygenic to omnigenic. *Cell*, 169(7):1177–1186.
- Bühlmann, P. (2013). Statistical significance in high-dimensional linear models. *Bernoulli*, 19(4):1212–1242.
- Bühlmann, P. and Van De Geer, S. (2011). *Statistics for high-dimensional data: methods, theory and applications*. Springer Science & Business Media.
- 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.
- Cai, T. T. and Guo, Z. (2018). Accuracy assessment for high-dimensional linear regression. *Annals of Statistics*, 46(5):1807–1836.
- Cai, T. T. and Low, M. G. (2004). An adaptation theory for nonparametric confidence intervals. *The Annals of Statistics*, 32(5):1805–1840.
- Cai, T. T. and Low, M. G. (2006). Adaptive confidence balls. *The Annals of Statistics*, 34(1):202–228.
- Carpentier, A. and Verzelen, N. (2019). Optimal sparsity testing in linear regression model. *arXiv:1901.08802*.
- Chakravarti, A. and Turner, T. N. (2016). Revealing rate-limiting steps in complex disease biology: The crucial importance of studying rare, extreme-phenotype families. *BioEssays*, 38(6):578–586.
- Chalkidou, A., O’Doherty, M. J., and Marsden, P. K. (2015). False discovery rates in pet and ct studies with texture features: a systematic review. *PloS one*, 10(5):e0124165.

- Furlong, L. I. (2013). Human diseases through the lens of network biology. *Trends in Genetics*, 29(3):150–159.
- Ganjgahi, H., Winkler, A. M., Glahn, D. C., Blangero, J., Donohue, B., Kochunov, P., and Nichols, T. E. (2018). Fast and powerful genome wide association of dense genetic data with high dimensional imaging phenotypes. *Nature communications*, 9(1):3254.
- Genovese, C. and Wasserman, L. (2008). Adaptive confidence bands. *The Annals of Statistics*, 36(2):875–905.
- Goeman, J. J., Van De Geer, S. A., and Van Houwelingen, H. C. (2006). Testing against a high dimensional alternative. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 68(3):477–493.
- Hoffmann, M. and Nickl, R. (2011). On adaptive inference and confidence bands. *The Annals of Statistics*, 39(5):2383–2409.
- Huttenhower, C., Gevers, D., Knight, R., Abubucker, S., Badger, J. H., Chinwalla, A. T., Creasy, H. H., Earl, A. M., FitzGerald, M. G., Fulton, R. S., et al. (2012). Structure, function and diversity of the healthy human microbiome. *Nature*, 486(7402):207.
- Ingster, Y. I., Tsybakov, A. B., and Verzelen, N. (2010). Detection boundary in sparse regression. *Electronic Journal of Statistics*, 4:1476–1526.
- Janson, L., Barber, R. F., and Candès, E. (2017). Eigenprism: inference for high dimensional signal-to-noise ratios. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(4):1037–1065.
- 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.
- Javanmard, A. and Montanari, A. (2018). De-biasing the lasso: Optimal sample size for gaussian designs. *to appear in the Annals of Statistics*, 46(6A):22593–2622.
- Krishnan, A., Zhang, R., Yao, V., Theesfeld, C. L., Wong, A. K., Tadych, A., Volfovsky, N., Packer, A., Lash, A., and Troyanskaya, O. G. (2016). Genome-wide prediction and functional characterization of the genetic basis of autism spectrum disorder. *Nature neuroscience*, 19(11):1454.
- Lehmann, E. L. and Romano, J. P. (2006). *Testing statistical hypotheses*. Springer.
- Meinshausen, N. and Bühlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *The Annals of Statistics*, 34(3):1436–1462.
- Negahban, S., Yu, B., Wainwright, M. J., and Ravikumar, P. K. (2009). A unified framework for high-dimensional analysis of m -estimators with decomposable regularizers. In *Advances in Neural Information Processing Systems*, pages 1348–1356.
- Nickl, R. and van de Geer, S. (2013). Confidence sets in sparse regression. *The Annals of Statistics*, 41(6):2852–2876.
- Raskutti, G., Wainwright, M. J., and Yu, B. (2011). Minimax rates of estimation for high-dimensional linear regression over l_q -balls. *IEEE transactions on information theory*, 57(10):6976–6994.
- Robins, J. and Van Der Vaart, A. (2006). Adaptive nonparametric confidence sets. *The Annals of Statistics*, 34(1):229–253.
- Rudelson, M. and Zhou, S. (2013). Reconstruction from anisotropic random measurements. *IEEE Transactions on Information Theory*, 59(6):3434–3447.
- Shah, R. D. and Bühlmann, P. (2017). Goodness-of-fit tests for high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- Tsybakov, A. B. (2008). *Introduction to Nonparametric Estimation*. Springer Publishing Company, Incorporated, 1st edition.
- 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.
- Van der Vaart, A. W. (2000). *Asymptotic statistics*, volume 3. Cambridge university press.
- Verhoef, P. C., Stephen, A. T., Kannan, P., Luo, X., Abhishek, V., Andrews, M., Bart, Y., Datta, H., Fong, N., Hoffman, D. L., et al. (2017). Consumer connectivity in a complex, technology-enabled, and mobile-oriented world with smart products. *Journal of Interactive Marketing*, 40:1–8.
- Vershynin, R. (2010). Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*.
- Wainwright, M. J., Lafferty, J. D., and Ravikumar, P. K. (2007). High-dimensional graphical model selection using l_1 -regularized logistic regression. In *Advances in neural information processing systems*, pages 1465–1472.
- Zhang, C.-H. and Zhang, S. S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242.
- Zhu, Y. and Bradic, J. (2017). Significance testing in non-sparse high-dimensional linear models. *ArXiv e-prints*.
- Zhu, Y. and Bradic, J. (2018). Linear hypothesis testing in dense high-dimensional linear models. *Journal of the American Statistical Association*, 113(524):1583–1600.

**SUPPLEMENT TO THE “TESTABILITY OF
HIGH-DIMENSIONAL LINEAR MODELS WITH
NON-SPARSE STRUCTURES”**

This document collects detailed proofs of Theorems 1, 4, 7, 8 and 10 as well as Corollaries 5, 6 and 9 of the main text, as well as detailed proofs of the twelve supplementary lemmas (alphabetically enumerated in this document): Lemma 1 - 12. In particular, Lemmas 1 - 7 are used for Theorem 4. Lemmas 8-9 are used for Theorem 8. Lemmas 10 - 12 are used for Theorem 10.

APPENDIX D: PROOF OF THEOREM 1

Proof of Theorem 1. To simply notation, we write \mathbb{E} instead of \mathbb{E}_θ . All the statements hold uniformly for any $\theta \in \Theta$. Let $\hat{\beta} = \Omega_1 \mathbf{X}^\top \mathbf{y}/n$. Then

$$\begin{aligned} \hat{\beta} - \beta &= \Omega_1 \mathbf{X}^\top (\mathbf{Z}\beta + \mathbf{W}\gamma + \varepsilon)/n - \beta \\ &= \left[\Omega_1 \mathbf{X}^\top \mathbf{Z}/n - 1 \right] \beta + \Omega_1 \mathbf{X}^\top (\mathbf{W}\gamma + \varepsilon)/n. \end{aligned}$$

Notice that

$$\Sigma = \begin{pmatrix} \gamma^\top \Sigma_{\mathbf{W}} \gamma + \sigma_{\mathbf{V}}^2 & \Sigma_{\mathbf{W}} \gamma \\ \gamma^\top \Sigma_{\mathbf{W}} & \Sigma_{\mathbf{W}} \end{pmatrix}$$

and $\Omega_1 = \sigma_{\mathbf{V}}^{-2} (1, -\boldsymbol{\pi}^\top)$, where $\Sigma_{\mathbf{W}} = \mathbb{E}(\mathbf{W}^\top \mathbf{W})/n$, $\sigma_{\mathbf{V}}^2 = \mathbb{E}(\mathbf{V}^\top \mathbf{V})/n$ and $\mathbf{V} = \mathbf{Z} - \mathbf{W}\boldsymbol{\pi}$. Then

$$\begin{aligned} &\left[\Omega_1 \mathbf{X}^\top \mathbf{Z}/n - 1 \right] \beta + \Omega_1 \mathbf{X}^\top (\mathbf{W}\gamma + \varepsilon)/n \\ &= n^{-1} \sum_{i=1}^n \left[(v_i Z_i \sigma_{\mathbf{V}}^{-2} - 1) + \sigma_{\mathbf{V}}^{-2} v_i (\mathbf{W}_i^\top \gamma + \varepsilon_i) \right]. \end{aligned}$$

where v_i , Z_i and $\mathbf{W}_i^\top \gamma$ denote the i -th entry of \mathbf{V} , \mathbf{Z} and $\mathbf{W}\gamma$, respectively.

Notice that

$$\left\{ (v_i Z_i \sigma_{\mathbf{V}}^{-2} - 1) + \sigma_{\mathbf{V}}^{-2} v_i (\mathbf{W}_i^\top \gamma + \varepsilon_i) \right\}_{i=1}^n$$

is an i.i.d sequence of random variables with bounded sub-exponential norms. Therefore,

$$\mathbb{E}(\hat{\beta} - \beta)^2 = n^{-2} \sum_{i=1}^n \left[(v_i Z_i \sigma_{\mathbf{V}}^{-2} - 1) + \sigma_{\mathbf{V}}^{-2} v_i (\mathbf{W}_i^\top \gamma + \varepsilon_i) \right]^2 \lesssim n^{-1}.$$

The desired result follows by noticing $\mathbb{E}|\hat{\beta} - \beta| \leq \sqrt{\mathbb{E}(\hat{\beta} - \beta)^2}$. \square

APPENDIX E: PROOF OF THEOREM 4

Before the main proof we establish a sequence of useful auxiliary results. Then we shall prove Theorem 4. To simplify notations, we write \mathbb{P} instead of \mathbb{P}_θ . Note that all the results here hold uniformly over $\theta \in \tilde{\Theta}(s)$ in finite samples. Therefore, we also omit $\sup_{\theta \in \tilde{\Theta}(s)}$ and $\inf_{\theta \in \tilde{\Theta}(s)}$ whenever possible.

E.1. Auxiliary results. The following result establishes a concentration result regarding the product of two Gaussian random variables that are allowed to be dependent. In particular, the result generalizes the concentration of measure of chi-squared random variables.

LEMMA 1. *Let $\{r_{i,1}\}_{i=1}^n$ and $\{r_{i,2}\}_{i=1}^n$ be sequences of i.i.d random variables with $\mathcal{N}(0, \sigma_1^2)$ and $\mathcal{N}(0, \sigma_2^2)$ distributions, respectively that are not necessarily independent from each other. Then for any $t > 0$,*

$$\mathbb{P}\left(\left|\sum_{i=1}^n (r_{i,1}r_{i,2} - \mathbb{E}r_{i,1}r_{i,2})\right| \geq t\sigma_1\sigma_2\right) \leq 2 \exp\left(-\frac{t^2}{2(2n+7t)}\right).$$

LEMMA 2. *Let the assumption of Theorem 4 hold. Then,*

- (1) *The population parameter $\boldsymbol{\xi}$ satisfies $\|\boldsymbol{\xi}\|_2 \leq 2M^2M_2$.*
- (2) *The estimator $\hat{\boldsymbol{\xi}}$ satisfies*

$$\mathbb{P}\left(\|\hat{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty > 2b_n^{-1}M\sqrt{n(\log p)(M_1^2 + M_2^2)}\right) \leq 2/p.$$

- (3) *Similarly,*

$$\mathbb{P}\left(\|\hat{\boldsymbol{\xi}} - b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i y_i\|_\infty > 4b_n^{-1}M\sqrt{n(\log p)(M_1^2 + M_2^2)}\right) \leq 4/p.$$

- (4) *Moreover,*

$$\mathbb{P}\left(\|\hat{\boldsymbol{\xi}}_{\mathcal{A}}\|_2 \leq 4M^2M_2\right) \geq 1 - 4/p.$$

- (5) *The ℓ_∞ -norm of estimation error of the thresholded estimator is*

$$\mathbb{P}\left(\|\hat{\boldsymbol{\xi}}_{\mathcal{A}^c}\|_\infty \leq 8b_n^{-1}M\sqrt{n(\log p)(M_1^2 + M_2^2)}\right) \geq 1 - 4/p.$$

(6) Lastly,

$$\begin{aligned} & \mathbb{P}\left(\left|b_n^{-1} \sum_{i \in H_4} v_i(\mathbf{W}_i^\top(\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \varepsilon_i)\right| \right. \\ & \quad \left. \leq 10b_n^{-1/2} \sqrt{M(4M_2^2M^3 + M_1^2) \log(100/\alpha)}\right) \geq 1 - 0.02\alpha, \end{aligned}$$

(7) and

$$\mathbb{P}\left(b_n^{-1} \sum_{i \in H_4} v_i^2 \geq (2M)^{-1}\right) \geq 1 - 2 \exp(-M^{-2}b_n/44).$$

We now discuss the estimation properties of the proposed regularized estimator $\widehat{\boldsymbol{\Omega}}_{\mathbf{W}}$.

LEMMA 3. *Let the assumption of Theorem 4 hold. Then $\boldsymbol{\Omega}_{\mathbf{W}}$ satisfies the constraint in (13) for $\widehat{\boldsymbol{\Omega}}_{\mathbf{W}}$ with probability at least $1 - 10/p - 2 \exp(-b_n/18)$.*

The next result establishes a lower bound on the restricted eigenvalue constant

$$\kappa(s) = \min_{|J| \subset \{1, \dots, p-1\}, |J| \leq s} \min_{\|\mathbf{q}_{J^c}\|_1 \leq 3\|\mathbf{q}_J\|_1} \frac{b_n^{-1} \sum_{i \in H_2} (\mathbf{W}_i^\top \mathbf{q})^2}{\|\mathbf{q}_J\|_2^2}.$$

LEMMA 4. *Let $\tau \in (0, 1)$ be an arbitrary constant. Whenever*

$$(1 + 36M^2(1 + \tau)^2(1 - \tau)^{-2}) s \leq p - 1,$$

and $b_n \geq 570 [1 + 36M^2(1 + \tau)^2(1 - \tau)^{-2}] \tau^{-2} s \log(12ep/\tau)$, then

$$\mathbb{P}(\kappa(s) > 0.24(1 - \tau)^2 M^{-1}) \geq 1 - 4 \exp(-\tau^2 b_n/570).$$

The following result establishes finite-sample properties of the Lasso estimator and follows by standard arguments. We include it here for completeness and clarity.

LEMMA 5. *Let the assumption of Theorem 4 hold. Then,*

$$\mathbb{P}(\|\widehat{\boldsymbol{\pi}} - \boldsymbol{\pi}\|_1 \leq 267s\lambda_\pi M) \geq 1 - 4 \exp(-3b_n/3040) - 2/p^2$$

and

$$\mathbb{P}\left(\left\|\sum_{i \in H_4} \mathbf{W}_i v_i\right\|_\infty / b_n \leq \lambda_\pi / 4\right) \geq 1 - 2/p^2.$$

The next two results establish the properties of the proposed regularized estimator $\check{\boldsymbol{\pi}}$.

LEMMA 6. *Let the assumption of Theorem 4 hold. Then $\boldsymbol{\pi}$ satisfies the constraints in (15) for $\check{\boldsymbol{\pi}}$ with probability at least $1 - 14/p - 0.02\alpha - 6 \exp(-3b_n/3040) - 2 \exp(-M^{-2}b_n/44)$.*

LEMMA 7. *Let the assumption of Theorem 4 hold. Then with probability at least $1 - 14/p - 0.02\alpha - 10 \exp(-3b_n/3040) - 2 \exp(-M^{-2}b_n/44)$,*

$$\|\check{\boldsymbol{\pi}} - \boldsymbol{\pi}\|_1 \leq 134M\lambda_{\boldsymbol{\pi}}s.$$

E.2. Proof of Theorem 4. Now we are ready to prove Theorem 4.

Proof of Theorem 4. Let $\boldsymbol{\delta} = \check{\boldsymbol{\pi}} - \boldsymbol{\pi}$. Notice that $\widehat{v}_i = v_i - \mathbf{W}_i^\top \boldsymbol{\delta}$. Then

$$(1) \quad \begin{aligned} \widehat{\beta} - \beta &= \frac{b_n^{-1} \sum_{i \in H_4} \widehat{v}_i (y_i - \beta \widehat{v}_i)}{b_n^{-1} \sum_{i \in H_4} \widehat{v}_i^2} \\ &= \underbrace{\frac{b_n^{-1} \sum_{i \in H_4} v_i (y_i - \beta \widehat{v}_i)}{b_n^{-1} \sum_{i \in H_4} \widehat{v}_i^2}}_{T_1} - \underbrace{\frac{b_n^{-1} \sum_{i \in H_4} \boldsymbol{\delta}^\top \mathbf{W}_i (y_i - \beta \widehat{v}_i)}{b_n^{-1} \sum_{i \in H_4} \widehat{v}_i^2}}_{T_2}. \end{aligned}$$

We now bound T_1 and T_2 in two steps. We first make the following observations. Notice that Lemma 7 implies

$$\mathbb{P}(\mathcal{M}_1) \geq 1 - 14/p - 0.02\alpha - 10 \exp(-3b_n/3040) - 2 \exp(-M^{-2}b_n/44),$$

where

$$\mathcal{M}_1 = \{\|\check{\boldsymbol{\pi}} - \boldsymbol{\pi}\|_1 \leq 134M\lambda_{\boldsymbol{\pi}}s\}.$$

Moreover, Lemma 2 implies that $\mathbb{P}(\mathcal{M}_2) \geq 1 - 8/p - 0.02\alpha$, where

$$\begin{aligned} \mathcal{M}_2 &= \left\{ \left\| \widehat{\boldsymbol{\xi}} - b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i y_i \right\|_\infty \leq 4b_n^{-1} M \sqrt{n(\log p)(M_1^2 + M_2^2)} \right\} \\ &\quad \cap \left\{ \|\widehat{\boldsymbol{\xi}}_{A^c}\|_\infty \leq 8b_n^{-1} M \sqrt{n(\log p)(M_1^2 + M_2^2)} \right\} \\ &\quad \cap \left\{ \left| b_n^{-1} \sum_{i \in H_4} v_i \left(\mathbf{W}_i^\top (\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \varepsilon_i \right) \right| \right. \\ &\quad \left. \leq 10b_n^{-1/2} \sqrt{M(4M_2^2 M^3 + M_1^2) \log(100/\alpha)} \right\}. \end{aligned}$$

Finally, Lemma 6 implies that

$$\mathbb{P}(\mathcal{M}_3) \geq 1 - 14/p - 0.02\alpha - 6 \exp(-3b_n/3040) - 2 \exp(-M^{-2}b_n/44),$$

where

$$\begin{aligned} \mathcal{M}_3 = & \left\{ \left| \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\pi}_A - \widehat{\boldsymbol{\xi}}_A^\top \widetilde{\boldsymbol{\pi}}_A \right| \leq \eta_\pi \right\} \cap \left\{ b_n^{-1} \sum_{i \in H_4} (Z_i - \mathbf{W}_i^\top \check{\boldsymbol{\pi}})^2 \geq \frac{1}{2M} \right\} \\ & \cap \left\{ \left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i (Z_i - \mathbf{W}_i^\top \boldsymbol{\pi}) \right\|_\infty \leq \lambda_\pi/4 \right\}. \end{aligned}$$

Define

$$\mathcal{M} = \mathcal{M}_1 \cap \mathcal{M}_2 \cap \mathcal{M}_3.$$

Since $b_n > n/4 - 1 > n/5$ (due to $n > 784$) and $p \geq 360/\alpha$, we have

$$\begin{aligned} \mathbb{P}(\mathcal{M}) & \geq 1 - 36/p - 0.06\alpha - 16 \exp(-3b_n/3040) - 4 \exp(-M^{-2}b_n/44) \\ & > 1 - 0.1\alpha - 0.06\alpha - 16 \exp(-3n/15200) - 4 \exp(-M^{-2}n/220) \\ (2) \quad & \stackrel{(i)}{\geq} 1 - 0.16\alpha - 16 \times 0.01\alpha - 4 \times 0.01\alpha > 1 - \alpha \end{aligned}$$

where (i) follows by the assumption of $n \geq 5067 \log(100/\alpha)$ and $n \geq 220M^2 \log(100/\alpha)$.

Since $\widehat{v}_i = Z_i - \mathbf{W}_i^\top \check{\boldsymbol{\pi}}$, we have that by definition, on the event \mathcal{M} ,

$$(3) \quad b_n^{-1} \sum_{i \in H_4} \widehat{v}_i^2 \geq \frac{1}{2M}.$$

Step 1: bound T_1 .

First observe that

$$y_i = Z_i \beta + \mathbf{W}_i^\top \boldsymbol{\gamma} + \varepsilon_i = \mathbf{W}_i^\top (\boldsymbol{\pi} \beta + \boldsymbol{\gamma}) + \beta v_i + \varepsilon_i.$$

Hence, $y_i - \beta \widehat{v}_i = \mathbf{W}_i^\top (\boldsymbol{\pi} \beta + \boldsymbol{\gamma}) + \mathbf{W}_i^\top \boldsymbol{\delta} + \varepsilon_i$. Therefore,

$$b_n^{-1} \sum_{i \in H_4} v_i (y_i - \beta \widehat{v}_i) = b_n^{-1} \underbrace{\sum_{i \in H_4} v_i \left(\mathbf{W}_i^\top (\boldsymbol{\pi} \beta + \boldsymbol{\gamma}) + \varepsilon_i \right)}_{T_{1,1}} + b_n^{-1} \underbrace{\sum_{i \in H_4} v_i \mathbf{W}_i^\top \boldsymbol{\delta}}_{T_{1,2}}.$$

By definition, on the event \mathcal{M} , we have

$$|T_{1,1}| \leq 10b_n^{-1/2} \sqrt{M (4M_2^2 M^3 + M_1^2) \log(100/\alpha)}.$$

Notice that $\mathbf{W}_i v_i = \mathbf{W}_i (Z_i - \mathbf{W}_i^\top \boldsymbol{\pi})$. Therefore, on the event \mathcal{M} ,

$$|T_{1,2}| \leq \|\boldsymbol{\delta}\|_1 \left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i v_i \right\|_\infty \leq (134M\lambda_\pi s) \times (\lambda_\pi/4) < 34M\lambda_\pi^2 s.$$

The above displays and (3) imply that on the event \mathcal{M} ,

$$(4) \quad |T_1| \leq 2M \left(10b_n^{-1/2} \sqrt{M(4M_2^2 M^3 + M_1^2) \log(100/\alpha)} + 34M\lambda_\pi^2 s \right).$$

Step 2: bound T_2 .

First notice that

$$(5) \quad \begin{aligned} & b_n^{-1} \sum_{i \in H_4} \boldsymbol{\delta}^\top \mathbf{W}_i (y_i - \beta \hat{v}_i) \\ &= \underbrace{b_n^{-1} \sum_{i \in H_4} \boldsymbol{\delta}^\top (\mathbf{W}_i y_i - \hat{\boldsymbol{\xi}})}_{T_{2,1}} + \underbrace{\boldsymbol{\delta}^\top \hat{\boldsymbol{\xi}}}_{T_{2,2}} - \underbrace{b_n^{-1} \sum_{i \in H_4} \boldsymbol{\delta}^\top \mathbf{W}_i \hat{v}_i \beta}_{T_{2,3}}. \end{aligned}$$

On the event \mathcal{M} , by Hölder's inequality, we have

$$(6) \quad \begin{aligned} |T_{2,1}| &\leq \|\boldsymbol{\delta}\|_1 \left\| b_n^{-1} \sum_{i \in H_4} (\mathbf{W}_i y_i - \hat{\boldsymbol{\xi}}) \right\|_\infty \\ &\leq (134M\lambda_\pi s) \times \left(4b_n^{-1} M \sqrt{n(\log p)(M_1^2 + M_2^2)} \right) \\ &= 536b_n^{-1} M^2 \sqrt{n(\log p)(M_1^2 + M_2^2)} \lambda_\pi s. \end{aligned}$$

To bound $T_{2,2}$, notice that on the event \mathcal{M} ,

$$\begin{aligned} |T_{2,2}| &= \left| \boldsymbol{\delta}_A^\top \hat{\boldsymbol{\xi}}_A + \boldsymbol{\delta}_{A^c}^\top \hat{\boldsymbol{\xi}}_{A^c} \right| \\ &\leq \left| (\tilde{\boldsymbol{\pi}}_A - \tilde{\boldsymbol{\pi}}_A)^\top \hat{\boldsymbol{\xi}}_A \right| + \left| (\tilde{\boldsymbol{\pi}}_A - \boldsymbol{\pi}_A)^\top \hat{\boldsymbol{\xi}}_A \right| + \left| \boldsymbol{\delta}_{A^c}^\top \hat{\boldsymbol{\xi}}_{A^c} \right| \\ &\stackrel{(i)}{\leq} \eta_\pi + \left| (\tilde{\boldsymbol{\pi}}_A - \boldsymbol{\pi}_A)^\top \hat{\boldsymbol{\xi}}_A \right| + \left| \boldsymbol{\delta}_{A^c}^\top \hat{\boldsymbol{\xi}}_{A^c} \right| \\ &\stackrel{(ii)}{\leq} \eta_\pi + \eta_\pi + \left| \boldsymbol{\delta}_{A^c}^\top \hat{\boldsymbol{\xi}}_{A^c} \right| \\ &\leq 2\eta_\pi + \|\boldsymbol{\delta}_{A^c}\|_1 \|\hat{\boldsymbol{\xi}}_{A^c}\|_\infty \\ &\leq 2\eta_\pi + \|\boldsymbol{\delta}\|_1 \|\hat{\boldsymbol{\xi}}_{A^c}\|_\infty \\ &\stackrel{(iii)}{\leq} 2\eta_\pi + (134M\lambda_\pi s) \times \left(8b_n^{-1} M \sqrt{n(\log p)(M_1^2 + M_2^2)} \right) \end{aligned}$$

$$(7) \quad = 2\eta_\pi + 1072b_n^{-1}M^2\sqrt{n(\log p)(M_1^2 + M_2^2)}\lambda_\pi s,$$

where (i) follows by the constraint (15) and (ii) and (iii) follow by the definition of \mathcal{M} .

To bound $T_{2,3}$, notice on the event \mathcal{M} , the constraint in (15) is satisfied by $\check{\boldsymbol{\pi}}$ and thus $\|b_n^{-1}\sum_{i \in H_4} \mathbf{W}_i(Z_i - \mathbf{W}_i^\top \check{\boldsymbol{\pi}})\|_\infty \leq \lambda_\pi/4$, which is

$$\left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \hat{v}_i \right\|_\infty \leq \lambda_\pi/4.$$

Therefore, on the event \mathcal{M} ,

$$(8) \quad \begin{aligned} |T_{2,3}| &\leq \|\boldsymbol{\delta}\|_1 \left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \hat{v}_i \right\|_\infty |\beta| \\ &\stackrel{(i)}{\leq} (134M\lambda_\pi s) \times (\lambda_\pi/4) \times M_2 < 34MM_2\lambda_\pi^2 s. \end{aligned}$$

where (i) follows by the definition of \mathcal{B} and the fact that $|\beta|^2 \leq \beta^2 + \|\boldsymbol{\gamma}\|_2^2 = \|\boldsymbol{\beta}\|_2^2 \leq M_2^2$.

In light of (5) and (3), we combine (6), (7) and (8), obtaining that on the event \mathcal{M} ,

$$(9) \quad |T_2| \leq 2M \left(1608b_n^{-1}M^2\sqrt{n(\log p)(M_1^2 + M_2^2)}\lambda_\pi s + 2\eta_\pi + 34MM_2\lambda_\pi^2 s \right).$$

By (1), (4) and (9), it follows that on the event \mathcal{M} ,

$$(10) \quad |\hat{\beta} - \beta| \leq c_n.$$

Therefore, by (2), for any $\theta \in \tilde{\Theta}(s, \beta_0)$, we have $\mathbb{E}_\theta \psi_* = \mathbb{P}_\theta(|\hat{\beta} - \beta_0| > c_n) = \mathbb{P}_\theta(|\hat{\beta} - \beta| > c_n) \leq \alpha$. This proves the first part of Theorem 4.

We now show the second part of Theorem 4. It is straight-forward to see that $b_n \asymp n$, $\lambda_\pi \asymp \sqrt{n^{-1} \log p}$ and $\eta_\pi \asymp sn^{-1} \log p + n^{-1/2}$. Therefore, $c_n \asymp n^{-1/2} + sn^{-1} \log p$.

Moreover, for any $\theta \in \tilde{\Theta}(s, \beta_0 + 3c_n)$, we have that on the event \mathcal{M} ,

$$|\hat{\beta} - \beta_0| \geq |\beta - \beta_0| - |\hat{\beta} - \beta| = 3c_n - |\hat{\beta} - \beta| \stackrel{(i)}{\geq} 2c_n > c_n,$$

where (i) follows by (10). Thus, for any $\theta \in \tilde{\Theta}(s, \beta_0 + 3c_n)$, we have

$$\mathbb{E}_\theta \psi_* = \mathbb{P}_\theta(|\hat{\beta} - \beta_0| > c_n) \geq \mathbb{P}_\theta(\mathcal{M}) \stackrel{(i)}{\geq} 1 - \alpha,$$

where (i) holds by (2). This proves the second part of Theorem 4. \square

APPENDIX F: PROOF OF COROLLARY 5

Proof of Corollary 5. Let $h_0 = \min\{\rho, \tau\}$, where ρ and τ are defined in Theorems 2 and 3, respectively. Notice that

$$\begin{aligned} \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_0(n^{-1/2} + sn^{-1} \log p)) \\ \subset \Theta_{\zeta, \kappa}(s/2, \beta_0 + \rho sn^{-1} \log p) \cap \Theta_{\kappa}(s, \beta_0 + \tau n^{-1/2}). \end{aligned}$$

Thus, Theorems 2 and 3 imply

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_{\alpha}(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_0(n^{-1/2} + sn^{-1} \log p))} \mathbb{E}_{\theta} \psi \leq 2\alpha.$$

Hence,

$$(11) \quad \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_{\alpha}(\Theta(s, \beta_0))} \inf_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_0(n^{-1/2} + sn^{-1} \log p))} \mathbb{E}_{\theta} \psi \leq 2\alpha.$$

The desired result follows by noticing that $\Psi_{\alpha}(\tilde{\Theta}(s, \beta_0)) \subset \Psi_{\alpha}(\Theta(s, \beta_0))$ and $\Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n) \subset \tilde{\Theta}(s, \beta_0 + h_n)$ with $h_n = h_0(n^{-1/2} + sn^{-1} \log p)$.

APPENDIX G: PROOF OF COROLLARY 6

Consider a sequence of $\mathcal{CI} = [l, u] \in \mathcal{C}_{\alpha}(\Theta(s))$ such that

$$\limsup_{n \rightarrow \infty} \inf_{\theta \in \Theta(s)} \mathbb{E}_{\theta} \text{diam}(\mathcal{CI}) = \limsup_{n \rightarrow \infty} \inf_{\mathcal{CI}' \in \mathcal{C}_{\alpha}(\Theta(s))} \sup_{\theta \in \Theta(s)} \mathbb{E}_{\theta} \text{diam}(\mathcal{CI}').$$

Consider the test

$$\psi = \mathbf{1}\{\beta_0 \notin \mathcal{CI}\}$$

for testing $\theta \in \Theta(s, \beta_0)$. Clearly, $\psi \in \Psi_{\alpha}(\Theta(s, \beta_0))$. Consider $\Theta(s, \beta_0 + h_n)$ with $h_n = h_0(n^{-1/2} + sn^{-1} \log p)$ defined in Corollary 5.

Fix any $\theta \in \Theta(s, \beta_0 + h_n)$. We have that $\beta = \beta_0 + h'$ with $0 \leq h' \leq h_n$. Notice that

$$\begin{aligned} 1 - \mathbb{E}_{\theta} \psi &= \mathbb{P}_{\theta}(\beta_0 \in \mathcal{CI}) \\ &= \mathbb{P}_{\theta}(l \leq \beta_0 \leq u) \\ &= \mathbb{P}_{\theta}(l \leq \beta_0 \leq u \text{ and } \beta \in \mathcal{CI}) + \mathbb{P}_{\theta}(l \leq \beta_0 \leq u \text{ and } \beta \notin \mathcal{CI}) \\ &\stackrel{(i)}{=} \mathbb{P}_{\theta}(l \leq \beta_0 \leq u \text{ and } \beta_0 + h' \in \mathcal{CI}) + \mathbb{P}_{\theta}(l \leq \beta_0 \leq u \text{ and } \beta \notin \mathcal{CI}) \\ &= \mathbb{P}_{\theta}(\max\{l, l - h'\} \leq \beta_0 \leq \min\{u, u - h'\}) + \mathbb{P}_{\theta}(l \leq \beta_0 \leq u \text{ and } \beta \notin \mathcal{CI}) \\ &\leq \mathbb{P}_{\theta}(\max\{l, l - h'\} \leq \min\{u, u - h'\}) + \mathbb{P}_{\theta}(\beta \notin \mathcal{CI}) \\ &\leq \mathbb{P}_{\theta}(l \leq u - h') + \alpha \end{aligned}$$

$$= \mathbb{P}_\theta(\text{diam}(\mathcal{CI}) \geq h') + \alpha \leq \mathbb{P}_\theta(\text{diam}(\mathcal{CI}) \geq h_n) + \alpha$$

where (i) follows by $\beta = \beta_0 + h'$. Hence,

$$\inf_{\theta \in \Theta(s, \beta_0 + h_n)} \mathbb{E}_\theta \psi \geq 1 - \alpha - \sup_{\theta \in \Theta(s, \beta_0 + h_n)} \mathbb{P}_\theta(\text{diam}(\mathcal{CI}) \geq h_n).$$

By (11) in the proof of Corollary 5, we have that

$$\begin{aligned} & \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \inf_{\theta \in \Theta(s, \beta_0 + h_n)} \mathbb{E}_\theta \psi \\ & \leq \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \inf_{\theta \in \Theta(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi \\ & \leq \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \inf_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi \leq 2\alpha. \end{aligned}$$

The above two displays imply

$$\liminf_{n \rightarrow \infty} \sup_{\theta \in \Theta(s, \beta_0 + h)} \mathbb{P}_\theta(\text{diam}(\mathcal{CI}) \geq h_n) \geq 1 - 3\alpha.$$

The desired result follows by noticing that

$$\text{diam}(\mathcal{CI}) \geq \text{diam}(\mathcal{CI}) \mathbf{1}\{\text{diam}(\mathcal{CI}) \geq h_n\} \geq h_n \mathbf{1}\{\text{diam}(\mathcal{CI}) \geq h_n\}$$

and thus

$$\sup_{\theta \in \Theta(s, \beta_0 + h_n)} \mathbb{E}_\theta \text{diam}(\mathcal{CI}) \geq h_n \sup_{\theta \in \Theta(s, \beta_0 + h_n)} \mathbb{P}_\theta(\text{diam}(\mathcal{CI}) \geq h_n).$$

□

APPENDIX H: PROOF OF THEOREM 7

PROOF OF THEOREM 7. By Theorem 4, we have

$$\inf_{\mathcal{CI} \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI}) = O\left(n^{-1/2} + sn^{-1} \log p\right).$$

Hence, it suffices to show that

$$(12) \quad \liminf_{n \rightarrow \infty} \frac{\inf_{\mathcal{CI} \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI})}{(n^{-1/2} + sn^{-1} \log p)} > 0.$$

We proceed by contradiction. Let h_n be defined as in Theorem 2. Fix an arbitrary $\beta_0 \in \mathbb{R}$. Suppose that there exists $\mathcal{CI}_0 = [l_0, u_0] \in \mathcal{C}_\alpha(\Theta(s))$ such that

$$\sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_0) = \tilde{h}_n$$

with $\liminf_{n \rightarrow \infty} (\tilde{h}_n/h_n) = 0$. Define $\Delta = \alpha^{-1}\tilde{h}_n$. Consider

$$\psi_0 = \mathbf{1}\{\beta_0 \notin \mathcal{CI}_0\}$$

as the test for $H_0 : \beta = \beta_0$ vs $H_a : \beta = \beta_0 + \Delta$.

Notice that

$$\sup_{\theta \in \Theta(\beta_0, s)} \mathbb{E}_\theta \psi_0 = \sup_{\theta \in \Theta(\beta_0, s)} P_\theta(\beta_0 \notin \mathcal{CI}_0) \stackrel{(i)}{\leq} \alpha,$$

where (i) follows by $\mathcal{CI}_0 \in \mathcal{C}_\alpha(\Theta(s))$. Thus, $\psi_0 \in \Psi(\Theta(\beta_0, s))$.

Fix an arbitrary $\theta_1 \in \Theta_{\zeta, \kappa}(s_1, \beta_0 + \Delta)$. Notice that on the event

$$\{\beta_0 + \Delta \in \mathcal{CI}_0\} \cap \{u_0 - \Delta < l_0\},$$

we have $\beta_0 + \Delta \leq u_0$, which means $\beta_0 \leq u_0 - \Delta < l_0$ and thus $\beta_0 \notin \mathcal{CI}_0$. Hence,

$$\begin{aligned} \mathbb{E}_{\theta_1} \psi_0 &= \mathbb{P}_{\theta_1}(\beta_0 \notin \mathcal{CI}_0) \geq \mathbb{P}_{\theta_1}(\{\beta_0 + \Delta \in \mathcal{CI}_0\} \cap \{u_0 - \Delta < l_0\}) \\ &\geq \mathbb{P}_{\theta_1}(\beta_0 + \Delta \in \mathcal{CI}_0) - \mathbb{P}_{\theta_1}(u_0 - \Delta \geq l_0) \\ &\stackrel{(i)}{\geq} 1 - \alpha - \mathbb{P}_{\theta_1}(u_0 - l_0 \geq \Delta) \\ &\stackrel{(ii)}{\geq} 1 - \alpha - \frac{\mathbb{E}_{\theta_1}|u_0 - l_0|}{\Delta} \\ &\stackrel{(iii)}{\geq} 1 - \alpha - \frac{\tilde{h}_n}{\Delta} \\ &\stackrel{(iv)}{=} 1 - 2\alpha \end{aligned}$$

where (i) follows by $\mathcal{CI}_0 \in \mathcal{C}_\alpha(\Theta(s))$, (ii) follows by Markov's inequality, (iii) follows by the fact that $\theta_1 \in \Theta(\beta_0 + \Delta, s_1)$ and $\sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_0) \leq \tilde{h}_n$ and (iv) follows by $\Delta = \alpha^{-1}\tilde{h}_n$. Consequently, we obtain

$$\limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s_1, \beta_0 + \Delta)} \mathbb{E}_\theta \psi \geq 1 - 2\alpha.$$

Since $\Delta \asymp \tilde{h}_n = o(h_n)$ and $s_1 \leq s/2$, we have that $\Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)$ contains $\Theta_{\zeta, \kappa}(s_1, \beta_0 + \Delta)$ for large n and thus

$$\begin{aligned} \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s_1, \beta_0 + \Delta)} \mathbb{E}_\theta \psi \\ \leq \limsup_{n \rightarrow \infty} \sup_{\psi \in \Psi_\alpha(\Theta(s, \beta_0))} \sup_{\theta \in \Theta_{\zeta, \kappa}(s/2, \beta_0 + h_n)} \mathbb{E}_\theta \psi \stackrel{(i)}{\leq} \alpha, \end{aligned}$$

where (i) follows by Theorem 2. The above two displays imply that $\alpha \geq 1 - 2\alpha$. This is not possible since $\alpha < 1/3$. Hence, we have arrived at the contradiction.

Therefore, there does not exist $\mathcal{CI}_0 = [l_0, u_0] \in \mathcal{C}_\alpha(\Theta(s))$ such that $\sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI}_0) = O(\tilde{h}_n)$ with $\tilde{h}_n = o(h_n)$. Hence,

$$\liminf_{n \rightarrow \infty} \left(\inf_{\mathcal{CI} \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI}) \right) / h_n > 0.$$

Similarly using Theorem 3, we can show that

$$\liminf_{n \rightarrow \infty} \left(\inf_{\mathcal{CI} \in \mathcal{C}_\alpha(\Theta(s))} \sup_{\theta \in \mathcal{C}_\alpha(\Theta(s_1))} \mathbb{E}_\theta \text{diam}(\mathcal{CI}) \right) / (n^{-1/2}) > 0.$$

Therefore, we have proved that the claim in (12). The proof is complete. \square

APPENDIX I: PROOF OF THEOREM 8

We rely on the following two lemmas.

LEMMA 8. *Suppose that points in Θ are uniformly non-testable, i.e.,*

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \asymp \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

Then there exists a constant $c > 0$ such that $cL(\Theta, \Theta) \leq L(\Theta_1, \Theta) \leq L(\Theta, \Theta)$ for any $\Theta_1 \subseteq \Theta$.

LEMMA 9. *Suppose that there exists a constant $c > 0$ such that $cL(\Theta, \Theta) \leq L(\Theta_1, \Theta) \leq L(\Theta, \Theta)$ for any subset $\Theta_1 \subseteq \Theta$. Then*

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \asymp \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

Now we are ready to prove Theorem 8.

PROOF OF THEOREM 8. The result is simple consequence of the two Lemmas, Lemma 8 and 9 whose proofs can be found in Section L.3. \square

APPENDIX J: PROOF OF COROLLARY 9

PROOF OF COROLLARY 9. Clearly,

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \leq \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \leq c_2 h_n.$$

It remains to show that $\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \gtrsim h_n$. For that end, we fix an arbitrary $\tau \in \mathbb{R}$, an arbitrary $CI \in \mathcal{C}_\alpha(\Theta)$ as well as an arbitrary $\theta \in \Theta(\tau)$.

Define a test $\psi = \mathbf{1}\{\tau \notin CI\}$. Clearly, $\psi \in \Psi(\Theta(\tau))$. Let $[l, u] = CI$. Since $g(\theta) = \tau + c_1 h_n$ for $\theta \in \Theta(\tau + c_1 h_n)$ and $CI \in \mathcal{C}_\alpha(\Theta)$, we have that for any $\theta \in \Theta(\tau + c_1 h_n)$

$$\mathbb{P}_\theta(l \leq \tau + c_1 h_n \leq u) \geq 1 - \alpha.$$

By assumption,

$$\mathbb{P}_\theta(\{\tau < l\} \cup \{\tau > u\}) = \mathbb{E}_\theta \psi \leq 2\alpha.$$

Let $\mathcal{M} = \{l \leq \tau + c_1 h_n \leq u\} \cap \{l \leq \tau \leq u\}$. Clearly, $\mathbb{P}_\theta(\mathcal{M}) \geq 1 - 3\alpha$.

Notice that on the event \mathcal{M} , $l \leq \tau \leq u - c_1 h_n$, which means $u - l \geq c_1 h_n$. It follows that

$$\mathbb{E}_\theta \text{diam}(CI) \geq \mathbb{E}_\theta \text{diam}(CI) \times \mathbf{1}\{\mathcal{M}\} \geq c_1 h_n \mathbb{P}_\theta(\mathcal{M}) \geq (1 - 3\alpha)c_1 h_n.$$

Notice that the above bound holds for any $\theta \in \Theta(\tau)$ with any $\tau \in \mathbb{R}$. Hence,

$$\inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \geq (1 - 3\alpha)c_1 h_n.$$

Since the above bound holds for any $CI \in \mathcal{C}_\alpha(\Theta)$, we have

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \geq (1 - 3\alpha)c_1 h_n.$$

□

APPENDIX K: PROOF OF THEOREM 10

For $\theta = (\beta, \gamma, \Sigma, \sigma)$ and $Q > 0$, we denote $\theta \odot Q = (\beta Q, \gamma Q, \Sigma, \sigma Q)$. For any $C \subseteq \mathbb{R}$ and $Q > 0$, we define $Q \cdot C = \{Qx : x \in C\}$.

LEMMA 10. *For any $Q, N_1, N_2 > 0$,*

$$\tilde{\Theta}_{QN_1, QN_2}(s) = \{\theta \odot Q : \theta \in \tilde{\Theta}_{N_1, N_2}(s)\}.$$

LEMMA 11. For any $D, N_1, N_2 > 0$, let $\theta \in \tilde{\Theta}_{N_1, N_2}(s)$. Then $(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim (\theta \odot D)$ if and only if $(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W}) \sim \theta$.

LEMMA 12. For any $D, N_1, N_2 > 0$,

$$D\mathbb{A}(s, N_1, N_2) \geq \mathbb{A}(s, DN_1, DN_2).$$

PROOF OF THEOREM 10. By Lemma 12 with $(D, N_1, N_2) = (Q, M_1, M_2)$, we have that $Q\mathbb{A}(s, M_1, M_2) \geq \mathbb{A}(s, QM_1, QM_2)$.

We now apply Lemma 12 with $(D, N_1, N_2) = (Q^{-1}, QM_1, QM_2)$, obtaining $Q^{-1}\mathbb{A}(s, QM_1, QM_2) \geq \mathbb{A}(s, M_1, M_2)$. The desired result follows. \square

APPENDIX L: PROOF OF AUXILIARY LEMMAS

L.1. Proof of auxiliary lemmas used in proving Theorem 2.

PROOF OF LEMMA 1. Let

$$A_k = [1 - kan^{-1} \log p]^{-n} \frac{\binom{m}{k} \binom{p-m-1}{m-k}}{\binom{p-1}{m}}.$$

Notice that for $0 \leq k \leq m$,

$$\begin{aligned} \log \frac{A_{k+1}}{A_k} &= \log \left[\left(1 - \frac{an^{-1} \log p}{1 - kan^{-1} \log p} \right)^{-n} \frac{(m-k)^2}{(k+1)(p-2m+k)} \right] \\ &= -n \log \left(1 - \frac{an^{-1} \log p}{1 - kan^{-1} \log p} \right) + \log \frac{(m-k)^2}{(k+1)(p-2m+k)} \\ &\leq -n \log \left(1 - \frac{an^{-1} \log p}{1 - kan^{-1} \log p} \right) + \log \frac{(m-k)^2}{p-2m+k} \\ &\stackrel{(i)}{\leq} -n \log (1 - 2an^{-1} \log p) + \log \frac{p^{2c}}{p-2p^c} \\ &\stackrel{(ii)}{\leq} \frac{2a \log p}{1 - 2an^{-1} \log p} + \log \frac{p^{2c}}{p-2p^c} \\ &\stackrel{(iii)}{<} 4a \log p + \log \frac{p^{2c}}{p-2p^c} \\ (13) \quad &= \log \frac{p^{4a+2c-1}}{1 - 2p^{c-1}}, \end{aligned}$$

where (i) follows by the fact that

$$1 - kan^{-1} \log p \geq 1 - man^{-1} \log p \geq 1 - a/4 \geq 1/2,$$

(ii) follows by the fact that $\log(1-x) \geq x/(x-1)$ for any $x \in (0, 1)$ and $2an^{-1} \log p \in (0, 1)$ (due to $2an^{-1} \log p \leq 2a/(4m) \leq a/2 < 1/2$) and (iii) follows by $2an^{-1} \log p < 1/2$.

Notice that $4a+2c-1 < 0$ and $c-1 < 0$. Hence, for large p , $\log(A_{k+1}/A_k) \leq -\log 2$ for any $0 \leq k \leq m$. It follows that for large p ,

$$(14) \quad \sum_{k=0}^m A_k = A_0 + \sum_{k=1}^m A_k \leq A_0 + A_1 \sum_{k=1}^m 2^{-k} \leq A_0 + 2A_1$$

Notice that

$$A_0 = \frac{\binom{p-m-1}{m}}{\binom{p-1}{m}} = \prod_{j=0}^{m-1} \frac{p-2m+j}{p-m+j} = \prod_{j=0}^{m-1} \left(1 - \frac{m}{p-m+j}\right).$$

Hence,

$$\left(1 - \frac{m}{p-m}\right)^m \leq A_0 \leq \left(1 - \frac{m}{p}\right)^m$$

Since $m^2/p \leq p^{2c-1} \rightarrow 0$, both sides tend to 1 and thus $A_0 \rightarrow 1$. To bound A_1 , notice that (13) implies

$$A_1 \leq \frac{p^{4a+2c-1}}{1-2p^{c-1}} A_0 \stackrel{(i)}{=} o(A_0),$$

where (i) follows by $4a+2c-1 < 0$ and $c < 1$. Hence, $A_1 = o(1)$. In light of (14), the desired result follows. \square

PROOF OF LEMMA 2. Notice that

$$\mathbb{E}_{g_0} \left(\frac{d\mathbb{P}_{g_1}}{d\mathbb{P}_{g_0}} \times \frac{d\mathbb{P}_{g_2}}{d\mathbb{P}_{g_0}} \right) = \int_{\mathbb{R}^k} \frac{g_1(x)g_2(x)}{g_0(x)} dx.$$

By Lemma 11 in Cai and Guo (2017), we have

$$\begin{aligned} & \int_{\mathbb{R}^k} \frac{g_1(x)g_2(x)}{g_0(x)} dx \\ &= \frac{1}{\sqrt{\det(\mathbb{I}_k - \Sigma_0^{-1}[\Sigma_1 - \Sigma_0]\Sigma_0^{-1}[\Sigma_2 - \Sigma_0])}} \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\sqrt{\det (\mathbb{I}_k - (L_0^{-1})^\top L_0^{-1} [L_1 L_1^\top - L_0 L_0^\top] (L_0^{-1})^\top L_0^{-1} [L_2 L_2^\top - L_0 L_0^\top])}} \\
&= \frac{1}{\sqrt{\det (\mathbb{I}_k - (L_0^{-1})^\top [Q_1 Q_1^\top - \mathbb{I}_k] [Q_2 Q_2^\top - \mathbb{I}_k] L_0^\top)}} \\
&= \frac{1}{\sqrt{\det (\mathbb{I}_k - (L_0^{-1})^\top [Q_1 Q_1^\top - \mathbb{I}_k] [Q_2 Q_2^\top - \mathbb{I}_k] L_0^\top)}} \\
&= \frac{1}{\sqrt{\det \{ (L_0^{-1})^\top (\mathbb{I}_k - [Q_1 Q_1^\top - \mathbb{I}_k] [Q_2 Q_2^\top - \mathbb{I}_k]) L_0^\top \}}} \\
&= \frac{1}{\sqrt{\det [(L_0^{-1})^\top] \det (\mathbb{I}_k - [Q_1 Q_1^\top - \mathbb{I}_k] [Q_2 Q_2^\top - \mathbb{I}_k]) \det (L_0^\top)}} \\
&= \frac{1}{\sqrt{\det (\mathbb{I}_k - [Q_1 Q_1^\top - \mathbb{I}_k] [Q_2 Q_2^\top - \mathbb{I}_k])}}.
\end{aligned}$$

□

PROOF OF LEMMA 3. We first derive some preliminary results and then compute

$$\det \left(\mathbb{I}_{p+1} - [Q_{j_1} Q_{j_1}^\top - \mathbb{I}_{p+1}] [Q_{j_2} Q_{j_2}^\top - \mathbb{I}_{p+1}] \right).$$

Step 1: First we derive the form of the matrix $Q_j Q_j^\top - \mathbb{I}_p$ for $1 \leq j \leq N$. By straight-forward computation, we can verify that

$$L_{\theta_*}^{-1} = \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ -\sigma_{\mathbf{V},*}^{-1} \boldsymbol{\pi}_*^\top & \sigma_{\mathbf{V},*}^{-1} & 0 \\ -\sigma_{\varepsilon,*}^{-1} \boldsymbol{\gamma}_*^\top & -\beta_* \sigma_{\varepsilon,*}^{-1} & \sigma_{\varepsilon,*}^{-1} \end{pmatrix}.$$

Thus,

$$\begin{aligned}
Q_j &= L_{\theta_*}^{-1} L_{\theta_j} \\
&= \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ -\sigma_{\mathbf{V},*}^{-1} \boldsymbol{\pi}_*^\top & \sigma_{\mathbf{V},*}^{-1} & 0 \\ -\sigma_{\varepsilon,*}^{-1} \boldsymbol{\gamma}_*^\top & -\beta_* \sigma_{\varepsilon,*}^{-1} & \sigma_{\varepsilon,*}^{-1} \end{pmatrix} \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ \boldsymbol{\pi}_{(j)}^\top & \sigma_{\mathbf{V},0} & 0 \\ (\boldsymbol{\pi}_{(j)} \beta_0 + \boldsymbol{\gamma}_{(j)})^\top & \beta_0 \sigma_{\mathbf{V},0} & \sigma_{\varepsilon,0} \end{pmatrix} \\
&= \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ \sigma_{\mathbf{V},*}^{-1} (\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*)^\top & \sigma_{\mathbf{V},*}^{-1} \sigma_{\mathbf{V},0} & 0 \\ \sigma_{\varepsilon,*}^{-1} [\boldsymbol{\gamma}_{(j)} - \boldsymbol{\gamma}_* + (\beta_0 - \beta_*) \boldsymbol{\pi}_{(j)}]^\top & -h \sigma_{\varepsilon,*}^{-1} \sigma_{\mathbf{V},0} & \sigma_{\varepsilon,*}^{-1} \sigma_{\varepsilon,0} \end{pmatrix}
\end{aligned}$$

$$\stackrel{(i)}{=} \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ a_2 \boldsymbol{\delta}_{(j)}^\top & \sigma_{\mathbf{V},*}^{-1} \sigma_{\mathbf{V},0} & 0 \\ a_1 \boldsymbol{\delta}_{(j)}^\top & -h \sigma_{\varepsilon,*}^{-1} \sigma_{\mathbf{V},0} & \sigma_{\varepsilon,*}^{-1} \sigma_{\varepsilon,0} \end{pmatrix}$$

for $a_1 = r(1-h)\sqrt{h/m}$ and $a_2 = \sqrt{h/m}$, where (i) follows by Definition 3. Since $\boldsymbol{\delta}_{(j)}^\top \boldsymbol{\delta}_{(j)} = m$, we have

$$\begin{aligned} & Q_j Q_j^\top - \mathbb{I}_{p+1} \\ &= \begin{pmatrix} \mathbb{I}_{p-1} & 0 & 0 \\ a_2 \boldsymbol{\delta}_{(j)}^\top & \sigma_{\mathbf{V},*}^{-1} \sigma_{\mathbf{V},0} & 0 \\ a_1 \boldsymbol{\delta}_{(j)}^\top & -h \sigma_{\varepsilon,*}^{-1} \sigma_{\mathbf{V},0} & \sigma_{\varepsilon,*}^{-1} \sigma_{\varepsilon,0} \end{pmatrix} \begin{pmatrix} \mathbb{I}_{p-1} & a_2 \boldsymbol{\delta}_{(j)} & a_1 \boldsymbol{\delta}_{(j)} \\ 0 & \sigma_{\mathbf{V},*}^{-1} \sigma_{\mathbf{V},0} & -h \sigma_{\varepsilon,*}^{-1} \sigma_{\mathbf{V},0} \\ 0 & 0 & \sigma_{\varepsilon,*}^{-1} \sigma_{\varepsilon,0} \end{pmatrix} - \mathbb{I}_{p+1} \\ (15) \quad & \stackrel{(i)}{=} \begin{pmatrix} 0 & a_2 \boldsymbol{\delta}_{(j)} & a_1 \boldsymbol{\delta}_{(j)} \\ a_2 \boldsymbol{\delta}_{(j)}^\top & 0 & 0 \\ a_1 \boldsymbol{\delta}_{(j)}^\top & 0 & 0 \end{pmatrix}. \end{aligned}$$

where (i) follows by Definition 3 and the definitions of a_1 and a_2 .

Step 2: Compute $\det \left(\mathbb{I}_{p+1} - \left[Q_{j_1} Q_{j_1}^\top - \mathbb{I}_{p+1} \right] \left[Q_{j_2} Q_{j_2}^\top - \mathbb{I}_{p+1} \right] \right)$ for any $j_1, j_2 \in \{1, \dots, N\}$.

From Step 1, we have that for any $j_1, j_2 \in \{1, \dots, M\}$,

$$\begin{aligned} & \mathbb{I}_{p+1} - \left(Q_{j_1} Q_{j_1}^\top - \mathbb{I}_{p+1} \right) \left(Q_{j_2} Q_{j_2}^\top - \mathbb{I}_{p+1} \right) \\ &= \mathbb{I}_{p+1} - \begin{pmatrix} 0 & a_2 \boldsymbol{\delta}_{(j_1)} & a_1 \boldsymbol{\delta}_{(j_1)} \\ a_2 \boldsymbol{\delta}_{(j_1)}^\top & 0 & 0 \\ a_1 \boldsymbol{\delta}_{(j_1)}^\top & 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & a_2 \boldsymbol{\delta}_{(j_2)} & a_1 \boldsymbol{\delta}_{(j_2)} \\ a_2 \boldsymbol{\delta}_{(j_2)}^\top & 0 & 0 \\ a_1 \boldsymbol{\delta}_{(j_2)}^\top & 0 & 0 \end{pmatrix} \\ &= \begin{pmatrix} \mathbb{I}_{p-1} - (a_1^2 + a_2^2) \boldsymbol{\delta}_{(j_1)} \boldsymbol{\delta}_{(j_2)}^\top & 0 & 0 \\ 0 & 1 - a_2^2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} & -a_1 a_2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \\ 0 & -a_1 a_2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} & 1 - a_1^2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \end{pmatrix}. \end{aligned}$$

Since this is a block-diagonal matrix, the desired result follows by simple computation

$$\begin{aligned} & \det \left[\mathbb{I}_{p+1} - \left(Q_{j_1} Q_{j_1}^\top - \mathbb{I}_{p+1} \right) \left(Q_{j_2} Q_{j_2}^\top - \mathbb{I}_{p+1} \right) \right] \\ &= \det \left(\mathbb{I}_{p-1} - (a_1^2 + a_2^2) \boldsymbol{\delta}_{(j_1)} \boldsymbol{\delta}_{(j_2)}^\top \right) \det \begin{pmatrix} 1 - a_2^2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} & -a_1 a_2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \\ -a_1 a_2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} & 1 - a_1^2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \end{pmatrix} \\ &\stackrel{(i)}{=} \left[1 - (a_1^2 + a_2^2) \boldsymbol{\delta}_{(j_1)} \boldsymbol{\delta}_{(j_2)}^\top \right] \det \begin{pmatrix} 1 - a_2^2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} & -a_1 a_2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \\ -a_1 a_2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} & 1 - a_1^2 \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \end{pmatrix} \end{aligned}$$

$$= \left[1 - (a_1^2 + a_2^2) \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \right]^2,$$

where (i) follows by the Sylvester's determinant identity. The desired result follows by the definitions of a_1 and a_2 . \square

PROOF OF LEMMA 4. Recall all the notations in Lemma 3 and ρ defined in (8). Notice that

$$\begin{aligned} & \mathbb{E}_{\theta_*} \left(N^{-1} \sum_{j=1}^N \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - 1 \right)^2 \\ &= N^{-2} \sum_{j_2=1}^N \sum_{j_1=1}^N \mathbb{E}_{\theta_*} \left(\frac{dP_{\theta_{j_1}}}{dP_{\theta_*}} \times \frac{dP_{\theta_{j_2}}}{dP_{\theta_*}} \right) - 1 \\ &\stackrel{(i)}{=} N^{-2} \sum_{j_2=1}^N \sum_{j_1=1}^N \left[1 - m^{-1} h [r^2 (1-h)^2 + 1] \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \right]^{-n} - 1 \\ &\stackrel{(ii)}{=} N^{-1} \sum_{j=1}^N \left[1 - m^{-1} h [r^2 (1-h)^2 + 1] \boldsymbol{\delta}_{(1)}^\top \boldsymbol{\delta}_{(j)} \right]^{-n} - 1, \end{aligned}$$

where (i) follows by Lemmas 2 and 3 (since there are n i.i.d observations, likelihood is a simple product) and (ii) follows by observing that

$$\sum_{j_1=1}^N \left[1 - m^{-1} h [r^2 (1-h)^2 + 1] \boldsymbol{\delta}_{(j_1)}^\top \boldsymbol{\delta}_{(j_2)} \right]^{-n}$$

does not depend on j_2 . To see this, simply notice that $\{\boldsymbol{\delta}_{(j)}^\top \boldsymbol{\delta}_{(j_2)}\}_{1 \leq j \leq N}$ is a permutation of $\{\boldsymbol{\delta}_{(j)}^\top \boldsymbol{\delta}_{(1)}\}_{1 \leq j \leq N}$ for any $1 \leq j_2 \leq N$.

For $k \in \{0, 1, \dots, m\}$, let

$$S_k = \{j \in \{1, \dots, N\} : \boldsymbol{\delta}_{(1)}^\top \boldsymbol{\delta}_{(j)} = k\}.$$

Notice that the cardinality of S_k is $\binom{m}{k} \binom{p-m-1}{m-k}$. Recall that $N = \binom{p-1}{m}$. It follows that

$$\mathbb{E}_{\theta_*} \left(N^{-1} \sum_{j=1}^N \frac{d\mathbb{P}_{\theta_j}}{d\mathbb{P}_{\theta_*}} - 1 \right)^2$$

$$= \sum_{k=0}^m [1 - m^{-1}h[r^2(1-h)^2 + 1]k]^{-n} \frac{\binom{m}{k} \binom{p-m-1}{m-k}}{\binom{p-1}{m}} - 1.$$

By Lemma 1, it suffices to verify that we can choose $a \in (0, (1-2c)/4)$ such that $m^{-1}h[r^2(1-h)^2 + 1] \leq an^{-1} \log p$. We now verify the stronger condition of

$$\frac{m^{-1}h[r^2(1-h)^2 + 1]}{n^{-1} \log p} < (1-2c)/5.$$

To this end, we recall $h = dsn^{-1} \log p$, $0 \leq d \leq \rho$ and $s/2 - 1 < m \leq s/2$ from Definition 3. Since $m \geq 1$, we have $s/m \leq (2m+1)/m \leq 3$. Now we observe that

$$\begin{aligned} \frac{m^{-1}h[r^2(1-h)^2 + 1]}{n^{-1} \log p} &= \frac{(m^{-1}dsn^{-1} \log p) [r^2(1-h)^2 + 1]}{n^{-1} \log p} \\ &\leq m^{-1}\rho s[r^2(1-h)^2 + 1] \\ &\stackrel{(i)}{\leq} 3\rho[r^2(1-h)^2 + 1] \\ &= 3\rho[r^2(1-dsn^{-1} \log p)^2 + 1] \\ &\stackrel{(ii)}{\leq} 3\rho[r^2 + 1] \\ &\stackrel{(iii)}{\leq} 3\rho[\kappa^{-2}M + 1] \\ &\stackrel{(iv)}{\leq} (1/2 - c)/5, \end{aligned}$$

where (i) follows by $s/m \leq 3$, (ii) follows by $dsn^{-1} \log p \leq 1$ (due to $sn^{-1} \log p \leq 1/4$ and $0 \leq d \leq \rho \leq 4$), (iii) follows by $r \leq \sqrt{M}/\kappa$ (since $r = \sigma_{\mathbf{V},*}/\sigma_{\varepsilon,*}$, $\sigma_{\mathbf{V},*}^2 \leq M$ and $\sigma_{\varepsilon,*} \geq \kappa$) and (v) follows by the definition of ρ . The proof is complete. \square

PROOF OF LEMMA 5. Recall that from Lemma 6, we can write $\theta_* = (\beta_*, \gamma_*, \Sigma_*, \sigma_{\varepsilon,*}) \in \Theta$ using

$$\Sigma_* = \begin{pmatrix} \pi_*^\top \pi_* + \sigma_{\mathbf{V},*}^2 & \pi_*^\top \\ \pi_* & \mathbb{I}_{p-1} \end{pmatrix}.$$

Since $\theta_* \in \Theta_{\zeta, \kappa}(m, \beta_0 + h_n)$, we have (1) $\beta_* = \beta_0 + h$ with $h = dsn^{-1} \log p$ and $0 \leq d \leq \rho$ and (2) $\lambda_{\max}(\Sigma_*) \leq \zeta M < M$. Notice that $\pi_*^\top \pi_* + \sigma_{\mathbf{V},*}^2 \leq \lambda_{\max}(\Sigma_*)$. Hence,

$$(16) \quad \max \{ \|\pi_*\|_2, \sigma_{\mathbf{V},*} \} \leq \sqrt{M}.$$

Recall $r = \sigma_{\mathbf{V},*}/\sigma_{\varepsilon,*}$. By the definition of $\Theta_{\zeta,\kappa}(s, \beta_0 + h_n)$, we have

$$(17) \quad r \leq \sqrt{M}/\kappa.$$

The rest of the proof proceeds in four steps, where we verify that

- (1) $\sigma_{\varepsilon,0} \leq M_1$,
- (2) $\|(\boldsymbol{\Sigma}_{(j)}^{-1})_{,1}\|_0 \leq 2m$ and
- (3) $M^{-1} \leq \lambda_{\min}(\boldsymbol{\Sigma}_{(j)}) \leq \lambda_{\max}(\boldsymbol{\Sigma}_{(j)}) \leq M$.
- (4) $\beta_0^2 + \|\boldsymbol{\gamma}_{(j)}\|_2^2 \leq \zeta^2 M_2^2$.

Step 1: Show $\sigma_{\varepsilon,0} \leq M_1$.

Notice that

$$\sigma_{\varepsilon,0} = \sigma_{\varepsilon,*} \sqrt{1 - hr^2 + h^2 r^2} \leq \zeta M_1 \sqrt{1 - hr^2 + h^2 r^2} \stackrel{(i)}{\leq} \zeta M_1 < M_1,$$

where (i) $-hr^2 + h^2 r^2 \leq 0$ (since $0 \leq h \leq \rho s n^{-1} \log p \leq \rho/4 \leq 1$).

Step 2: Show $\|(\boldsymbol{\Sigma}_{(j)}^{-1})_{,1}\|_0 \leq 2m$.

Observe that $(\boldsymbol{\Sigma}_{(j)}^{-1})_{,1} = \begin{pmatrix} 1 \\ -\boldsymbol{\pi}_{(j)} \end{pmatrix} \sigma_{\mathbf{V},0}^{-2}$ and $(\boldsymbol{\Sigma}_*^{-1})_{,1} = \begin{pmatrix} 1 \\ -\boldsymbol{\pi}_* \end{pmatrix} \sigma_{\mathbf{V},*}^{-2}$. Hence,

$$\|(\boldsymbol{\Sigma}_{(j)}^{-1})_{,1}\|_0 = \|\boldsymbol{\pi}_{(j)}\|_0 + 1$$

and $\|(\boldsymbol{\Sigma}_*^{-1})_{,1}\|_0 = \|\boldsymbol{\pi}_*\|_0 + 1$. Since

$$\|\boldsymbol{\pi}_{(j)}\|_0 \leq \|\boldsymbol{\pi}_*\|_0 + \|\boldsymbol{\delta}_{(j)}\|_0 = \|\boldsymbol{\pi}_*\|_0 + m$$

and $\theta_* \in \Theta_{\zeta,\kappa}(m, \beta_0 + h_n)$, we have

$$\|(\boldsymbol{\Sigma}_{(j)}^{-1})_{,1}\|_0 \leq \|(\boldsymbol{\Sigma}_*^{-1})_{,1}\|_0 + m \leq 2m.$$

Step 3: Show $M^{-1} \leq \lambda_{\min}(\boldsymbol{\Sigma}_{(j)}) \leq \lambda_{\max}(\boldsymbol{\Sigma}_{(j)}) \leq M$.

Since $2m \leq s \leq 2m + 1$ and $s \geq 2$, we have $m \geq 1$ and

$$2 \leq s/m \leq 2 + 1/m \leq 3.$$

Notice that $\|\boldsymbol{\delta}_{(j)}\|_2 = \sqrt{m}$ and

$$(18) \quad \|\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*\|_2 = \sigma_{\mathbf{V},*} \sqrt{h/m} \|\boldsymbol{\delta}_{(j)}\|_2 = \sigma_{\mathbf{V},*} \sqrt{h}.$$

Let $\|\cdot\|_F$ denote the Frobenius norm and observe that

$$\|\boldsymbol{\Sigma}_{(j)} - \boldsymbol{\Sigma}_*\|_F^2 = \left(\boldsymbol{\pi}_{(j)}^\top \boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*^\top \boldsymbol{\pi}_* + \sigma_{\mathbf{V},0}^2 - \sigma_{\mathbf{V},*}^2 \right)^2 + 2\|\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*\|_2^2$$

$$\begin{aligned}
&\stackrel{(i)}{=} \left(\|\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*\|_2^2 + 2(\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*)^\top \boldsymbol{\pi}_* - h\sigma_{\mathbf{V},*}^2 \right)^2 \\
&\quad + 2\|\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*\|_2^2 \\
&\stackrel{(ii)}{=} \left(2(\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*)^\top \boldsymbol{\pi}_* \right)^2 + 2\sigma_{\mathbf{V},*}^2 h \\
&\leq \left(2\|\boldsymbol{\pi}_{(j)} - \boldsymbol{\pi}_*\|_2 \times \|\boldsymbol{\pi}_*\|_2 \right)^2 + 2\sigma_{\mathbf{V},*}^2 h \\
&\stackrel{(iii)}{\leq} \left(2\sigma_{\mathbf{V},*} \sqrt{hM} \right)^2 + 2\sigma_{\mathbf{V},*}^2 h \\
&\stackrel{(iv)}{\leq} \left(2\sqrt{hM} \right)^2 + 2Mh \\
&\stackrel{(v)}{\leq} M(2M+1)\rho/2 \\
&\stackrel{(vi)}{\leq} \min \left\{ \frac{1}{M^2} \left(\frac{1}{\zeta} - 1 \right)^2, M^2(1-\zeta)^2 \right\},
\end{aligned}$$

where (i) follows by $\sigma_{\mathbf{V},0}^2 - \sigma_{\mathbf{V},*}^2 = -\sigma_{\mathbf{V},*}^2 h$ (due to Definition 3), (ii) follows by (18), (iii) follows by (18), (iv) follows by (16), (v) follows by $h \leq \rho/4$ (due to $h = ds n^{-1} \log p$ with $0 \leq d \leq \rho$ and $sn^{-1} \log p \leq 1/4$) and (vi) follows by $0 \leq d \leq \rho$ and the definition of ρ in (8).

Let $\|\cdot\|$ denote the spectral norm of a matrix (i.e., $\|A\| = \sqrt{\lambda_{\max}(A^\top A)}$). Notice that

$$\lambda_{\min}(\boldsymbol{\Sigma}_{(j)}) \geq \lambda_{\min}(\boldsymbol{\Sigma}_*) - \|\boldsymbol{\Sigma}_{(j)} - \boldsymbol{\Sigma}_*\|$$

and $\lambda_{\max}(\boldsymbol{\Sigma}_{(j)}) \leq \lambda_{\max}(\boldsymbol{\Sigma}_*) + \|\boldsymbol{\Sigma}_{(j)} - \boldsymbol{\Sigma}_*\|$. Since $\|\boldsymbol{\Sigma}_{(j)} - \boldsymbol{\Sigma}_*\| \leq \|\boldsymbol{\Sigma}_{(j)} - \boldsymbol{\Sigma}_*\|_F$, the above display implies that

$$\begin{aligned}
\lambda_{\min}(\boldsymbol{\Sigma}_{(j)}) &\geq \lambda_{\min}(\boldsymbol{\Sigma}_*) - \min \left\{ \frac{1}{M} \left(\frac{1}{\zeta} - 1 \right), M(1-\zeta) \right\} \\
&\geq \lambda_{\min}(\boldsymbol{\Sigma}_*) - \frac{1}{M} \left(\frac{1}{\zeta} - 1 \right)
\end{aligned}$$

and similarly

$$\begin{aligned}
\lambda_{\max}(\boldsymbol{\Sigma}_{(j)}) &\leq \lambda_{\max}(\boldsymbol{\Sigma}_*) + \min \left\{ \frac{1}{M} \left(\frac{1}{\zeta} - 1 \right), M(1-\zeta) \right\} \\
&\leq \lambda_{\max}(\boldsymbol{\Sigma}_*) + M(1-\zeta).
\end{aligned}$$

Since $(\zeta M)^{-1} \leq \lambda_{\min}(\boldsymbol{\Sigma}_*) \leq \lambda_{\max}(\boldsymbol{\Sigma}_*) \leq \zeta M$, we obtain

$$M^{-1} \leq \lambda_{\min}(\boldsymbol{\Sigma}_{(j)}) \leq \lambda_{\max}(\boldsymbol{\Sigma}_{(j)}) \leq M$$

Step 4: Show $\beta_0^2 + \|\gamma_{(j)}\|_2^2 \leq \zeta^2 M_2^2$.

Since $\theta_* \in \Theta_{\zeta, \kappa}(m, \beta_0 + h_n)$, we have that

$$(19) \quad (\beta_0 + h)^2 + \|\gamma_*\|_2^2 \leq \zeta^2 M_2^2.$$

Therefore, we need to show that

$$(20) \quad [\beta_0^2 + \|\gamma_{(j)}\|_2^2] - [(\beta_0 + h)^2 + \|\gamma_*\|_2^2] \leq (1 - \zeta^2)M_2^2.$$

Let $\delta_{\gamma, j} = \gamma_{(j)} - \gamma_*$. Notice that

$$(21) \quad \begin{aligned} & [\beta_0^2 + \|\gamma_{(j)}\|_2^2] - [(\beta_0 + h)^2 + \|\gamma_*\|_2^2] \\ &= -2(\beta_0 + h)h + h^2 + \|\delta_{\gamma, j}\|_2^2 + 2\gamma_*^\top \delta_{\gamma, j} \\ &\leq 2|\beta_0 + h|h + h^2 + \|\delta_{\gamma, j}\|_2^2 + 2\|\gamma_*\|_2 \cdot \|\delta_{\gamma, j}\|_2 \\ &\stackrel{(i)}{\leq} 2\zeta M_2 h + h^2 + \|\delta_{\gamma, j}\|_2^2 + 2\zeta M_2 \|\delta_{\gamma, j}\|_2 \\ &\leq 2\zeta M_2 h_n + h_n^2 + \|\delta_{\gamma, j}\|_2^2 + 2\zeta M_2 \|\delta_{\gamma, j}\|_2, \end{aligned}$$

where (i) follows by $\|\gamma_*\|_2 \leq \zeta M_2$ and $|\beta_0 + h| \leq \zeta M_2$ (due to (19)).

By the assumption of $sn^{-1} \log p \leq 1/4$, $M > 1$ and the definition of ρ in (8) we have that

$$(22) \quad h_n^2 = \rho^2 (sn^{-1} \log p)^2 \leq \rho^2 / 16 \leq \frac{(1 - \zeta^2)M_2^2}{64M} < (1 - \zeta^2)M_2^2 / 4$$

and

$$(23) \quad \begin{aligned} 2\zeta M_2 h_n &= 2\zeta M_2 \rho sn^{-1} \log p \leq \zeta M_2 \rho / 2 \\ &\leq \frac{(1 - \zeta^2)M_2^2}{16\sqrt{M}} < (1 - \zeta^2)M_2^2 / 4. \end{aligned}$$

By Definition 3, we have

$$\|\delta_{\gamma, j}\|_2 \leq h\|\pi_{(j)}\|_2 + r\sigma_{\varepsilon, *}\sqrt{h/m}\|\delta_{(j)}\|_2 = h\|\pi_{(j)}\|_2 + r\sigma_{\varepsilon, *}\sqrt{h}.$$

By (16) and (18), $\|\pi_{(j)}\|_2 \leq \|\pi_*\|_2 + \|\pi_{(j)} - \pi_*\|_2 \leq \sqrt{M} + \sigma_{\mathbf{V}, *}\sqrt{h}$. Since $h \leq h_n = \rho sn^{-1} \log p \leq \rho/4$, we have that

$$\begin{aligned} \|\delta_{\gamma, j}\|_2 &\leq \frac{1}{4}\rho \left(\sqrt{M} + \sigma_{\mathbf{V}, *}\sqrt{\rho/4} \right) + r\sigma_{\varepsilon, *}\sqrt{\rho/4} \\ &\stackrel{(i)}{\leq} \frac{1}{4}\rho \left(1 + \sqrt{\rho/4} \right) \sqrt{M} + \kappa^{-1}\sqrt{M}\zeta M_1 \sqrt{\rho/4} \end{aligned}$$

$$\stackrel{(ii)}{\leq} \frac{1}{2}\rho\sqrt{M} + \kappa^{-1}\sqrt{M}\zeta M_1\sqrt{\rho/4},$$

where (i) follows by $\sigma_{\mathbf{V},*} \leq \sqrt{M}$ (due to (16)), $\sigma_{\varepsilon,*} \leq \zeta M_1$ (due to the definition of $\Theta_{\zeta,\kappa}(s)$) and $r \leq \sqrt{M}/\kappa$ (due to (17)) and (ii) follows by $\rho \leq 4$. By the definition of ρ in (8), we have

$$\begin{aligned} 2\zeta M_2 \|\delta_{\gamma,j}\|_2 &\leq \zeta M_2 \sqrt{M} \rho + \frac{\sqrt{M}}{\kappa} \zeta^2 M_1 M_2 \sqrt{\rho} \\ (24) \quad &\leq \frac{(1-\zeta^2)M_2^2}{8} + \frac{(1-\zeta^2)M_2^2}{8} \leq \frac{(1-\zeta^2)M_2^2}{4}. \end{aligned}$$

By the elementary inequality of $(a+b)^2 \leq 2a^2 + 2b^2$, we also have

$$\begin{aligned} \|\delta_{\gamma,j}\|_2^2 &\leq \left(\frac{1}{2}\rho\sqrt{M} + \kappa^{-1}\sqrt{M}\zeta M_1\sqrt{\rho/4} \right)^2 \\ &\leq \frac{1}{2}\rho^2 M + \frac{M}{2\kappa^2} \zeta^2 M_1^2 \rho \\ (25) \quad &\stackrel{(i)}{\leq} \frac{(1-\zeta^2)M_2^2}{8} + \frac{(1-\zeta^2)M_2^2}{8} \leq \frac{(1-\zeta^2)M_2^2}{4}, \end{aligned}$$

where (i) follows by the definition of ρ in (8).

In light of (21), we obtain (20) by combining (22), (23), (24) and (25). The proof is complete. \square

PROOF OF LEMMA 6. Notice that

$$\begin{aligned} &\begin{pmatrix} a & b^\top \Sigma \\ \Sigma b & \Sigma \end{pmatrix}^{-1} \\ &= \begin{pmatrix} a^{-1} + a^{-2} b^\top \Sigma (\Sigma - a^{-1} \Sigma b b^\top \Sigma)^{-1} \Sigma b & -b^\top \Sigma (\Sigma - a^{-1} \Sigma b b^\top \Sigma)^{-1} \\ -(\Sigma - a^{-1} \Sigma b b^\top \Sigma)^{-1} \Sigma b & (\Sigma - a^{-1} \Sigma b b^\top \Sigma)^{-1} \end{pmatrix}. \end{aligned}$$

Since all the eigenvalues of the above matrix are positive, the eigenvalues of the blocks on the diagonal are also positive. This means that the eigenvalues of $\Sigma - a^{-1} \Sigma b b^\top \Sigma$ are positive. Notice that

$$\Sigma - a^{-1} \Sigma b b^\top \Sigma = \Sigma^{1/2} (\mathbb{I} - a^{-1} \Sigma^{1/2} b b^\top \Sigma^{1/2}) \Sigma^{1/2}.$$

Since $\Sigma^{1/2}$ is positive definite, we have that all the eigenvalues of $\mathbb{I} - a^{-1} \Sigma^{1/2} b b^\top \Sigma^{1/2}$ is positive. It follows that

$$\det(\mathbb{I} - a^{-1} \Sigma^{1/2} b b^\top \Sigma^{1/2}) > 0.$$

By Sylvester's determinant identity, we have $\det(\mathbb{I} - a^{-1} \Sigma^{1/2} b b^\top \Sigma^{1/2}) = 1 - a^{-1} b^\top \Sigma b$. The desired result follows. \square

L.2. Proof of auxiliary lemmas used in proving Theorem 4.

PROOF OF LEMMA 1. We first prove the result assuming $\sigma_1 = \sigma_2 = 1$. Let $r_i = r_{i,1}r_{i,2}$. Then for any $m \geq 3$,

$$|r_i|^m = |r_{i,1}r_{i,2}|^m \stackrel{(i)}{\leq} 2^{-m}(r_{i,1}^2 + r_{i,2}^2)^m \stackrel{(ii)}{\leq} \frac{1}{2}(|r_{i,1}|^{2m} + |r_{i,2}|^{2m}),$$

where (i) follows by $|r_{i,1}r_{i,2}| \leq (r_{i,1}^2 + r_{i,2}^2)/2$, (ii) follows by the elementary inequality $(a + b)^m \leq 2^{m-1}(a^m + b^m)$ for $a, b \geq 0$ and $m \geq 2$. Hence,

$$\sum_{i=1}^n \mathbb{E}|r_i|^m \leq \frac{n}{2} (\mathbb{E}|r_{1,1}|^{2m} + \mathbb{E}|r_{1,2}|^{2m}).$$

Since $r_{1,1} \sim \mathcal{N}(0, 1)$, we have that $r_{1,1}^2 \sim \chi^2(1)$. The moment generating function of χ^2 distributions implies

$$\mathbb{E} \exp(r_{1,1}^2/3) = (1 - 2/3)^{-1} = 3.$$

Notice that by Taylor's series,

$$\mathbb{E} \exp(r_{1,1}^2/3) = 1 + \sum_{j=1}^{\infty} \frac{3^{-j} \mathbb{E} \exp(|r_{1,1}|^{2j})}{j!}.$$

Therefore, for any $j \geq 1$,

$$\frac{3^{-j} \mathbb{E} \exp(|r_{1,1}|^{2j})}{j!} < 3.$$

Similarly, we can show that for any $j \geq 1$,

$$\frac{3^{-j} \mathbb{E} \exp(|r_{1,2}|^{2j})}{j!} < 3.$$

Let $\nu = 2n$. Hence, for $m \geq 6$,

$$\begin{aligned} \sum_{i=1}^n \mathbb{E}|r_i|^m &\leq \frac{n}{2} (\mathbb{E}|r_{1,1}|^{2m} + \mathbb{E}|r_{1,2}|^{2m}) \\ &\leq \frac{n}{2} (3^{m+1}m! + 3^{m+1}m!) = n3^{m+1}m! < \frac{m!}{2} \nu 7^{m-2}. \end{aligned}$$

Since both $r_{1,1}$ and $r_{1,2}$ are standard normal, we can easily compute for $m = 3, 4, 5$

$$\sum_{i=1}^n \mathbb{E}|r_i|^m \leq \frac{n}{2} (\mathbb{E}|r_{1,1}|^{2m} + \mathbb{E}|r_{1,2}|^{2m}) = \begin{cases} 15n & m = 3 \\ 105n & m = 4 \\ 945n & m = 5. \end{cases}$$

Thus, $\sum_{i=1}^n \mathbb{E}|r_i|^m \leq \frac{m!}{2} \nu 7^{m-2}$ for $m \geq 3$. Clearly, $\sum_{i=1}^n \mathbb{E}(r_i^2) = n < \nu$. Therefore, by Corollary 2.11 of [Boucheron et al. \(2013\)](#), we have that for any $t > 0$,

$$\mathbb{P}\left(\sum_{i=1}^n (r_i - \mathbb{E}r_i) \geq t\right) \leq \exp\left(-\frac{t^2}{2(2n+7t)}\right).$$

Similarly, we can show the same result for $-r_i$: for any $t > 0$,

$$\mathbb{P}\left(-\sum_{i=1}^n (r_i - \mathbb{E}r_i) \geq t\right) \leq \exp\left(-\frac{t^2}{2(2n+7t)}\right).$$

Hence,

$$\mathbb{P}\left(\left|\sum_{i=1}^n (r_i - \mathbb{E}r_i)\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(2n+7t)}\right).$$

We have proved the result for $\sigma_1 = \sigma_2 = 1$. In the general case, notice that $r_{i,1}\sigma_1^{-1} \sim \mathcal{N}(0, 1)$ and $r_{i,2}\sigma_2^{-1} \sim \mathcal{N}(0, 1)$. Hence, the above display implies

$$\mathbb{P}\left(\left|\sum_{i=1}^n (r_{i,1}r_{i,2}\sigma_1^{-1}\sigma_2^{-1} - \mathbb{E}r_{i,1}r_{i,2}\sigma_1^{-1}\sigma_2^{-1})\right| \geq t\right) \leq 2 \exp\left(-\frac{t^2}{2(2n+7t)}\right).$$

The desired result follows. \square

PROOF OF LEMMA 2. By the definition of $\tilde{\Theta}(s)$, we have that $\beta^2 + \|\boldsymbol{\gamma}\|_2^2 \leq M_2^2$. Notice that the first row of $\boldsymbol{\Sigma}^{-1}$ is $(1, -\boldsymbol{\pi}^\top)\sigma_{\mathbf{V}}^{-2}$. Therefore,

$$M^{-1} \leq \lambda_{\min}(\boldsymbol{\Sigma}^{-1}) \leq \|\boldsymbol{\pi}\|_2^2 \sigma_{\mathbf{V}}^{-2} + \sigma_{\mathbf{V}}^{-2} \leq \lambda_{\max}(\boldsymbol{\Sigma}^{-1}) \leq M.$$

This means that $M^{-1/2} \leq \sigma_{\mathbf{V}} \leq M^{1/2}$ and $\|\boldsymbol{\pi}\|_2 \leq M$. Since $M > 1$, it follows that

$$\|\boldsymbol{\xi}\|_2 \leq \lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{W}}) (|\beta| \cdot \|\boldsymbol{\pi}\|_2 + \|\boldsymbol{\gamma}\|_2)$$

$$\leq M(M_2M + M_2) = M_2M(M + 1) < 2M^2M_2.$$

This proves part (1).

Since $\mathbb{E}\mathbf{W}_{i,j}y_i = \mathbf{\Sigma}_{\mathbf{W}}(\boldsymbol{\pi}\beta + \boldsymbol{\gamma})$, we have that for each $1 \leq j \leq p-1$,

$$\widehat{\boldsymbol{\xi}}_j - \boldsymbol{\xi}_j = b_n^{-1} \sum_{i \in H_3} [\mathbf{W}_{i,j}y_i - \mathbb{E}\mathbf{W}_{i,j}y_i].$$

Notice that both $\mathbf{W}_{i,j}$ and y_i are normal random variables with mean zero. Moreover,

$$\mathbb{E}\mathbf{W}_{i,j}^2 \leq \lambda_{\max}(\mathbf{\Sigma}) \leq M$$

and

$$\mathbb{E}y_i^2 = \sigma^2 + \boldsymbol{\beta}^\top \mathbf{\Sigma} \boldsymbol{\beta} \leq \sigma^2 + \lambda_{\max}(\mathbf{\Sigma}_{\mathbf{W}}) \|\boldsymbol{\beta}\|_2^2 \leq M_1^2 + MM_2^2,$$

where we recall $\boldsymbol{\beta} = (\beta, \boldsymbol{\gamma}^\top)^\top \in \mathbb{R}^p$.

It follows by Lemma 1 that $\forall t > 0$,

$$\mathbb{P} \left(b_n |\widehat{\boldsymbol{\xi}}_j - \boldsymbol{\xi}_j| > t \sqrt{M(M_1^2 + MM_2^2)} \right) \leq 2 \exp \left(-\frac{t^2}{2(2b_n + 7t)} \right).$$

We set $t = 2\sqrt{n \log p}$. Since $n/4 - 1 < b_n \leq n/4$ and $n/\log p \geq 784 = 28^2$, the union bound implies

$$\begin{aligned} & \mathbb{P} \left(\|\widehat{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty > 2b_n^{-1} \sqrt{n(\log p)M(M_1^2 + MM_2^2)} \right) \\ & \leq 2p \exp \left(-\frac{t^2}{2(2b_n + 7t)} \right) \\ & \leq 2p \exp \left(-\frac{4n \log p}{2(n/2 + 14\sqrt{n \log p})} \right) \\ & = 2 \exp \left(\left(1 - \frac{4}{1 + 28\sqrt{n^{-1} \log p}} \right) \log p \right) \\ (26) \quad & \leq 2 \exp \left(\left(1 - \frac{4}{1 + 1} \right) \log p \right) = 2/p. \end{aligned}$$

Since $M > 1$, we have proved part (2).

By the same argument,

$$(27) \quad \mathbb{P} \left(\|\widetilde{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty > 2b_n^{-1} M \sqrt{n(\log p)(M_1^2 + M_2^2)} \right) \leq 2/p$$

and

$$\mathbb{P} \left(\left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_{i,j} y_i - \boldsymbol{\xi} \right\|_\infty > 2b_n^{-1} M \sqrt{n(\log p)(M_1^2 + M_2^2)} \right) \leq 2/p.$$

Part (3) follows.

Now we prove part (4).

Denote $\tau = 2b_n^{-1}M\sqrt{n(\log p)(M_1^2 + M_2^2)}$ and the event $\mathcal{B} = \{\|\tilde{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty \leq \tau\}$. Notice that $A = \{j : |\tilde{\boldsymbol{\xi}}_j| \geq 2\tau\}$ by the definition in (14). Define $A_\tau = \{j : |\boldsymbol{\xi}_j| \geq \tau\}$.

Since $|\boldsymbol{\xi}_j| \geq |\tilde{\boldsymbol{\xi}}_j| - |\tilde{\boldsymbol{\xi}}_j - \boldsymbol{\xi}_j|$, we have that $|\boldsymbol{\xi}_j| \geq |\tilde{\boldsymbol{\xi}}_j| - \|\tilde{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty$. Therefore, on the event \mathcal{B} , $|\boldsymbol{\xi}_j| \geq \tau$ for any $j \in A$. In other words, on the event \mathcal{B} , $A \subseteq A_\tau$ and thus $|A| \leq |A_\tau|$. To bound $|A_\tau|$, notice that $\tau^2|A_\tau| \leq \|\boldsymbol{\xi}\|_2^2$.

Define the event $\mathcal{B}' = \{\|\hat{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty \leq \tau\}$. On the event $\mathcal{B} \cap \mathcal{B}'$,

$$\begin{aligned} \|\hat{\boldsymbol{\xi}}_A\|_2 &\leq \|\boldsymbol{\xi}_A\|_2 + \|\hat{\boldsymbol{\xi}}_A - \boldsymbol{\xi}_A\|_2 \\ &\leq \|\boldsymbol{\xi}\|_2 + \sqrt{|A|}\|\hat{\boldsymbol{\xi}}_A - \boldsymbol{\xi}_A\|_\infty \\ &\leq \|\boldsymbol{\xi}\|_2 + \sqrt{|A_\tau|}\|\hat{\boldsymbol{\xi}}_A - \boldsymbol{\xi}_A\|_\infty \\ &\leq \|\boldsymbol{\xi}\|_2 + \sqrt{\|\boldsymbol{\xi}\|_2^2 \tau^{-2} \tau} \\ &= 2\|\boldsymbol{\xi}\|_2 \leq 4M^2M_2. \end{aligned}$$

Part (4) follows because (27) and part (2) imply $\mathbb{P}(\mathcal{B} \cap \mathcal{B}') \geq 1 - 4/p$.

To see part (5), notice that for any $j \in A^c$,

$$|\hat{\boldsymbol{\xi}}_j| \leq \|\hat{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty + \|\tilde{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty + |\tilde{\boldsymbol{\xi}}_j| \leq \|\hat{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty + \|\tilde{\boldsymbol{\xi}} - \boldsymbol{\xi}\|_\infty + 2\tau.$$

Therefore, on the event $\mathcal{B} \cap \mathcal{B}'$, $|\hat{\boldsymbol{\xi}}_j| \leq 4\tau$ for any $j \in A^c$. Part (5) follows.

Now we show part (6). The argument is similar to the proof of part (2). Notice that $v_i \sim \mathcal{N}(0, \sigma_{\mathbf{V}}^2)$ and

$$\mathbf{W}_i^\top(\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \varepsilon_i \sim \mathcal{N}(0, (\boldsymbol{\pi}\beta + \boldsymbol{\gamma})^\top \boldsymbol{\Sigma}_{\mathbf{W}}(\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \sigma^2).$$

Also notice that $\sigma_{\mathbf{V}}^2 \leq M$ and

$$\begin{aligned} (\boldsymbol{\pi}\beta + \boldsymbol{\gamma})^\top \boldsymbol{\Sigma}_{\mathbf{W}}(\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \sigma^2 &\leq \lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{W}})\|\boldsymbol{\pi}\beta + \boldsymbol{\gamma}\|_2^2 + M_1^2 \\ &\leq \lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{W}})(\|\boldsymbol{\pi}\|_2 \cdot |\beta| + \|\boldsymbol{\gamma}\|_2)^2 + M_1^2 \\ &\leq M(M_2M + M_2)^2 + M_1^2 \stackrel{(i)}{<} 4M_2^2M^3 + M_1^2, \end{aligned}$$

where (i) follows by $M > 1$.

Since $\mathbb{E}v_i[\mathbf{W}_i^\top(\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \varepsilon_i] = 0$, it follows by Lemma 1 that for any $t > 0$,

$$\mathbb{P}\left(\left|\sum_{i \in H_4} v_i \left(\mathbf{W}_i^\top(\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \varepsilon_i\right)\right| > tb_n^{1/2} \sqrt{M(4M_2^2M^3 + M_1^2)}\right)$$

$$(28) \quad \leq 2 \exp\left(-\frac{t^2 b_n}{2(2b_n + 7tb_n^{1/2})}\right) = 2 \exp\left(-\frac{t^2}{2(2 + 7tb_n^{-1/2})}\right).$$

Now we take $t = 10\sqrt{\log(100/\alpha)}$. The assumption of $(n-4)/\log p \geq 784$ implies that $n > 784$. Hence, $b_n > n/4 - 1 > n/5$, which means $b_n^{-1/2} < \sqrt{5/n}$. Thus, the assumption of Theorem 4 implies that $n > 500 \log(100/\alpha)$ and thus $tb_n^{-1/2} \leq 10\sqrt{5n^{-1} \log(100/\alpha)} < 1$. The above display implies

$$\begin{aligned} & \mathbb{P}\left(\left|\sum_{i \in H_4} v_i (\mathbf{W}_i^\top (\boldsymbol{\pi}\beta + \boldsymbol{\gamma}) + \varepsilon_i)\right| > 10b_n^{1/2} \sqrt{M(4M_2^2 M^3 + M_1^2) \log(100/\alpha)}\right) \\ & \leq 2 \exp\left(-\frac{100 \log(100/\alpha)}{2(2+7)}\right) = 2 \exp\left(-\frac{50}{9} \log(100/\alpha)\right) \\ & < 2 \exp(-\log(100/\alpha)) = \alpha/50. \end{aligned}$$

This proves part (6).

It remains to show part (7). Notice that $v_i \sim \mathcal{N}(0, \sigma_{\mathbf{V}}^2)$ and $M^{-1} \leq \sigma_{\mathbf{V}}^2 \leq M$. By an argument similar to (28), we have that for any $t > 0$,

$$\mathbb{P}\left(\left|\sum_{i \in H_4} (v_i^2 - \mathbb{E}v_i^2)\right| > tM\right) \leq 2 \exp\left(-\frac{t^2}{2(2b_n + 7t)}\right).$$

Now we take $t = b_n/(2M^2)$. Hence,

$$\begin{aligned} & \mathbb{P}\left(b_n^{-1} \sum_{i \in H_4} v_i^2 < \frac{1}{2M}\right) \leq \mathbb{P}\left(b_n^{-1} \sum_{i \in H_4} (v_i^2 - \mathbb{E}v_i^2) < \frac{1}{2M} - \mathbb{E}v_i^2\right) \\ & \leq \mathbb{P}\left(b_n^{-1} \sum_{i \in H_4} (v_i^2 - \mathbb{E}v_i^2) < \frac{1}{2M} - \frac{1}{M}\right) \\ & \leq \mathbb{P}\left(b_n^{-1} \left|\sum_{i \in H_4} (v_i^2 - \mathbb{E}v_i^2)\right| > \frac{1}{2M}\right) \\ & = \mathbb{P}\left(\left|\sum_{i \in H_4} (v_i^2 - \mathbb{E}v_i^2)\right| > tM\right) \\ & \leq 2 \exp\left(-\frac{b_n^2/(4M^4)}{2(2b_n + 7b_n/(2M^2))}\right) \\ & = 2 \exp\left(-\frac{b_n/M^4}{(16 + 28/M^2)}\right) \stackrel{(i)}{<} 2 \exp(-M^{-2}b_n/44), \end{aligned}$$

where (i) follows by $28/M^2 < 28$ (since $M > 1$). The proof is complete. \square

PROOF OF LEMMA 3. We need to show

$$(29) \quad \mathbb{P} \left(\left\| \left(\mathbb{I}_{p-1} - b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \Omega_{\mathbf{W}} \right) \widehat{\boldsymbol{\xi}}_A \right\|_\infty > 24 \sqrt{b_n^{-1} \log p M^3 M_2} \right) < 6/p$$

and

$$(30) \quad \mathbb{P} \left(\widehat{\boldsymbol{\xi}}_A^\top \Omega_{\mathbf{W}} \left(b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A \leq 32M^5 M_2^2 \right) \leq 2 \exp(-b_n/18) + 4/p.$$

We prove these two claims in two steps.

Step 1: show (29).

Define $q_i = \mathbf{W}_i \mathbf{W}_i^\top \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A$ and $q_{i,j} = \mathbf{W}_{i,j} \mathbf{W}_i^\top \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A$ for $1 \leq j \leq p-1$. Let \mathcal{F} denote the σ -algebra generated by $\{(\mathbf{W}_i, y_i, Z_i)\}_{i \in H_1 \cup H_3}$. Notice that $\widehat{\boldsymbol{\xi}}_A$ is \mathcal{F} -measurable and $\{\mathbf{W}_i\}_{i \in H_4}$ is independent of \mathcal{F} due to the sample splitting. Therefore, for $i \in H_4$, conditional on \mathcal{F} , $\mathbf{W}_{i,j}$ and $\mathbf{W}_i^\top \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A$ are both Gaussian with mean zero.

Also observe that for $i \in H_4$, $\mathbb{E}(\mathbf{W}_{i,j}^2 | \mathcal{F}) \leq \lambda_{\max}(\boldsymbol{\Sigma}_{\mathbf{W}}) \leq M$ and

$$\begin{aligned} \mathbb{E}[(\mathbf{W}_i^\top \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A)^2 | \mathcal{F}] &\leq \widehat{\boldsymbol{\xi}}_A^\top \Omega_{\mathbf{W}} \boldsymbol{\Sigma}_{\mathbf{W}} \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A = \widehat{\boldsymbol{\xi}}_A^\top \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A \\ &\leq \lambda_{\max}(\Omega_{\mathbf{W}}) \|\widehat{\boldsymbol{\xi}}_A\|_2^2 \leq \frac{\|\widehat{\boldsymbol{\xi}}_A\|_2^2}{\lambda_{\min}(\boldsymbol{\Sigma}_{\mathbf{W}})} \leq M \|\widehat{\boldsymbol{\xi}}_A\|_2^2. \end{aligned}$$

Therefore, Lemma 1 implies that for any $t > 0$,

$$\mathbb{P} \left(\left| \sum_{i \in H_4} [q_{i,j} - \mathbb{E}(q_{i,j} | \mathcal{F})] \right| > tM \|\widehat{\boldsymbol{\xi}}_A\|_2 | \mathcal{F} \right) \leq 2 \exp \left(-\frac{t^2}{2(2b_n + 7t)} \right).$$

Since $\mathbb{E}(q_i | \mathcal{F}) = \mathbb{E}(\mathbf{W}_i \mathbf{W}_i^\top \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A | \mathcal{F}) = \boldsymbol{\Sigma}_{\mathbf{W}} \Omega_{\mathbf{W}} \widehat{\boldsymbol{\xi}}_A = \widehat{\boldsymbol{\xi}}_A$, we apply the union bound and obtain that $\forall t > 0$,

$$\begin{aligned} &\mathbb{P} \left(\left\| \left(b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \Omega_{\mathbf{W}} - \mathbb{I}_{p-1} \right) \widehat{\boldsymbol{\xi}}_A \right\|_\infty > tM \|\widehat{\boldsymbol{\xi}}_A\|_2 | \mathcal{F} \right) \\ &= \mathbb{P} \left(\max_{1 \leq j \leq p-1} \left| \sum_{i \in H_4} [q_{i,j} - \mathbb{E}(q_{i,j} | \mathcal{F})] \right| > tb_n M \|\widehat{\boldsymbol{\xi}}_A\|_2 | \mathcal{F} \right) \end{aligned}$$

$$\leq 2p \exp\left(-\frac{t^2 b_n^2}{2(2b_n + 7tb_n)}\right) = 2p \exp\left(-\frac{t^2 b_n}{2(2 + 7t)}\right).$$

By choosing $t = 6\sqrt{b_n^{-1} \log p}$, it follows that

$$\begin{aligned} & \mathbb{P}\left(\left\|\left(b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \boldsymbol{\Omega} \mathbf{w} - \mathbb{I}_{p-1}\right) \widehat{\boldsymbol{\xi}}_A\right\|_\infty > 6\sqrt{b_n^{-1} \log p} M \|\widehat{\boldsymbol{\xi}}_A\|_2\right) \\ & \leq 2p \exp\left(-\frac{36 \log p}{4 + 14 \times 8\sqrt{b_n^{-1} \log p}}\right) \\ & \stackrel{(i)}{\leq} 2p \exp\left(-\frac{36 \log p}{4 + 14 \times 6/14}\right) = 2p^{-2.6} < 2p^{-2}, \end{aligned}$$

where (i) follows by the fact that $b_n > n/4 - 1$ and the assumption $(n - 4)/\log p \geq 784 = 28^2$. By Lemma 2, $\mathbb{P}\left(\|\widehat{\boldsymbol{\xi}}_A\|_2 \leq 4M^2 M_2\right) \geq 1 - 4/p$. Therefore, we have

$$\begin{aligned} & \mathbb{P}\left(\left\|\left(\mathbb{I}_{p-1} - b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \boldsymbol{\Omega} \mathbf{w}\right) \widehat{\boldsymbol{\xi}}_A\right\|_\infty > 24\sqrt{b_n^{-1} \log p} M^3 M_2\right) \\ & \leq 4/p + 2p^{-2} < 6/p. \end{aligned}$$

We have proved (29).

Step 2: show (30).

Let $r_i = \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \mathbf{W}_i$. For $i \in H_4$, notice that conditional on \mathcal{F} , r_i is Gaussian with mean zero and variance $\widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \boldsymbol{\Sigma} \mathbf{w} \boldsymbol{\Omega} \mathbf{w} \widehat{\boldsymbol{\xi}}_A = \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \widehat{\boldsymbol{\xi}}_A$. It follows by Lemma 1 that

$$\mathbb{P}\left(\left|\sum_{i \in H_4} [r_i^2 - \mathbb{E}(r_i^2 | \mathcal{F})]\right| > t \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \widehat{\boldsymbol{\xi}}_A \mid \mathcal{F}\right) \leq 2 \exp\left(-\frac{t^2}{2(2b_n + 7t)}\right).$$

Since $\mathbb{E}(r_i^2 | \mathcal{F}) = \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \widehat{\boldsymbol{\xi}}_A$, we have

$$\begin{aligned} & \mathbb{P}\left(b_n^{-1} \sum_{i \in H_4} r_i^2 > (1 + b_n^{-1}t) \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \widehat{\boldsymbol{\xi}}_A\right) \\ & = \mathbb{P}\left(\sum_{i \in H_4} [r_i^2 - \mathbb{E}(r_i^2 | \mathcal{F})] > t \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{w} \widehat{\boldsymbol{\xi}}_A\right) \leq 2 \exp\left(-\frac{t^2}{2(2b_n + 7t)}\right). \end{aligned}$$

Notice that $\widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{W} \widehat{\boldsymbol{\xi}}_A \leq \lambda_{\max}(\boldsymbol{\Omega} \mathbf{W}) \|\widehat{\boldsymbol{\xi}}_A\|_2^2 = \|\widehat{\boldsymbol{\xi}}_A\|_2^2 / \lambda_{\max}(\boldsymbol{\Sigma} \mathbf{W}) \leq M \|\widehat{\boldsymbol{\xi}}_A\|_2^2$. By Lemma 2, $\|\widehat{\boldsymbol{\xi}}_A\|_2 \leq 4M^2 M_2$ with probability at least $1 - 4/p$. Therefore, we have that

$$\mathbb{P} \left(b_n^{-1} \sum_{i \in H_4} r_i^2 > 16 (1 + b_n^{-1} t) M^5 M_2^2 \right) \leq 2 \exp \left(-\frac{t^2}{2(2b_n + 7t)} \right) + 4/p.$$

Since $b_n^{-1} \sum_{i \in H_4} r_i^2 = \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\Omega} \mathbf{W} (b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top) \boldsymbol{\Omega} \mathbf{W} \widehat{\boldsymbol{\xi}}_A$, we choose $t = b_n$ and obtain (30). \square

PROOF OF LEMMA 4. We invoke Corollary 18 of Rudelson and Zhou (2013) and Lemma 4.1 of Bickel et al. (2009).

For any k between 1 and p , we define the sparse eigenvalues

$$\phi_{\min}(k) = \min_{1 \leq \|\mathbf{q}\|_0 \leq k} \frac{b_n^{-1} \sum_{i \in H_2} (\mathbf{W}_i^\top \mathbf{q})^2}{\|\mathbf{q}\|_2^2}$$

and

$$\phi_{\max}(k) = \max_{1 \leq \|\mathbf{q}\|_0 \leq k} \frac{b_n^{-1} \sum_{i \in H_2} (\mathbf{W}_i^\top \mathbf{q})^2}{\|\mathbf{q}\|_2^2}.$$

The proof proceeds in two steps. We first verify a sufficient condition for the sparse eigenvalue condition and then derive the desired result.

Step 1: Show that rows of $\boldsymbol{\Sigma}_{\mathbf{W}}^{-1/2} \mathbf{W}$ are isotropic and ψ_2 with constant $\sqrt{8/3}$.

Notice that $\boldsymbol{\Sigma}_{\mathbf{W}}^{-1/2} \mathbf{W}$ is a matrix whose entries are i.i.d $\mathcal{N}(0, 1)$. Let \mathbf{r}^\top denote the first row of $\boldsymbol{\Sigma}_{\mathbf{W}}^{-1/2} \mathbf{W}$. For any nonzero vector $\mathbf{q} \in \mathbb{R}^{p-1}$, $(\mathbf{r}^\top \mathbf{q})^2 / \|\mathbf{q}\|_2^2$ has a chi-squared distribution with one degree of freedom. By the moment generating function of chi-squared distributions, we have that for any $t > \sqrt{2} \|\mathbf{q}\|_2$,

$$\mathbb{E} \left[\exp \left((\mathbf{r}^\top \mathbf{q})^2 / t^2 \right) \right] = \mathbb{E} \left[\exp \left(\frac{(\mathbf{r}^\top \mathbf{q})^2}{\|\mathbf{q}\|_2^2} \times \frac{\|\mathbf{q}\|_2^2}{t^2} \right) \right] = \left(1 - \frac{2\|\mathbf{q}\|_2^2}{t^2} \right)^{-1/2}.$$

Thus,

$$\inf \left\{ t : \mathbb{E} \left[\exp \left((\mathbf{r}^\top \mathbf{q})^2 / t^2 \right) \right] \right\} \leq \sqrt{8/3} \|\mathbf{q}\|_2.$$

In other words, \mathbf{r} is isotropic and ψ_2 with constant $\sqrt{8/3}$; see Definition 5 of Rudelson and Zhou (2013).

Step 2: Show the desired result.

By Corollary 18 of [Rudelson and Zhou \(2013\)](#), we have that with probability at least $1 - 2 \exp(-\tau^2 b_n / 570)$,

$$(1 - \tau)^2 M^{-1} \leq \phi_{\min}(k) \leq \phi_{\max}(k) \leq (1 + \tau)^2 M$$

if $b_n \geq 570\tau^{-2}k \log(12ep/\tau)$. Let m be the smallest integer satisfying

$$m \geq 36M^2(1 + \tau)^2(1 - \tau)^{-2}s.$$

This means that if $b_n \geq 570\tau^{-2}(s + m) \log(12ep/\tau)$, then $\mathbb{P}(\mathcal{B}) \geq 1 - 4 \exp(-\tau^2 b_n / 570)$, where the event \mathcal{B} is defined as

$$\mathcal{B} = \{ \phi_{\min}(s + m) \geq (1 - \tau)^2 M^{-1} \text{ and } \phi_{\max}(m) \leq (1 + \tau)^2 M \}.$$

Notice that on the event \mathcal{B} , $m\phi_{\min}(m + s) > c_0^2 s \phi_{\max}(m)$ with $c_0 = 3$. By Lemma 4.1(ii) of [Bickel et al. \(2009\)](#), on the event \mathcal{B}

$$\begin{aligned} \sqrt{\kappa(s)} &= \sqrt{\phi_{\min}(m + s)} \left(1 - c_0 \sqrt{\frac{s \phi_{\max}(m)}{m \phi_{\min}(m + s)}} \right) \\ &= \sqrt{\phi_{\min}(m + s)} - c_0 \sqrt{\frac{s}{m} \phi_{\max}(s)} \\ &\geq (1 - \tau)M^{-1/2} - 3 \times \sqrt{\frac{s}{36M^2(1 + \tau)^2(1 - \tau)^{-2}s}} \times (1 + \tau)^2 M \\ &= 0.5(1 - \tau)M^{-1/2}. \end{aligned}$$

The desired result follows. \square

PROOF OF LEMMA 5. We invoke Theorem 6.1 of [Bühlmann and Van De Geer \(2011\)](#). We first show a concentration result for $\|\sum_{i \in H_2} \mathbf{W}_i v_i\|_\infty$.

For $1 \leq j \leq p - 1$, $\mathbf{W}_{i,j} \sim \mathcal{N}(0, \mathbb{E}(\mathbf{W}_{i,j}^2))$ with $\mathbb{E}(\mathbf{W}_{i,j}^2) \leq \lambda_{\max}(\boldsymbol{\Sigma}) \leq M$. Also observe that $v_i \sim \mathcal{N}(0, \sigma_V^2)$ with $\sigma_V^2 \leq \lambda_{\max}(\boldsymbol{\Sigma}) \leq M$. Since $\mathbb{E}(\mathbf{W}_{i,j} v_i) = 0$, it follows by Lemma 1 that $\forall t > 0$,

$$\mathbb{P} \left(\left| \sum_{i \in H_2} \mathbf{W}_{i,j} v_i \right| > tM \right) \leq 2 \exp \left(-\frac{t^2}{2(2b_n + 7t)} \right).$$

By the union bound, we have

$$\mathbb{P} \left(\left\| \sum_{i \in H_2} \mathbf{W}_i v_i \right\|_\infty > tM \right) \leq 2p \exp \left(-\frac{t^2}{2(2b_n + 7t)} \right).$$

Taking $t = 6\sqrt{b_n \log p}$, we have that

$$\begin{aligned} \mathbb{P} \left(\left\| \sum_{i \in H_2} \mathbf{W}_i v_i \right\|_{\infty} > 6M \sqrt{b_n \log p} \right) &\leq 2p \exp \left(- \frac{36 \log p}{4 + 14 \times 6 \sqrt{b_n^{-1} \log p}} \right) \\ &\stackrel{(i)}{\leq} 2p \exp \left(- \frac{36 \log p}{4 + 6} \right) = 2p^{-2.6} < 2/p^2, \end{aligned}$$

where (i) follows by $b_n > n/4 - 1$ and the assumption $(n - 4)/\log p \geq 784 = 28^2$. In other words,

$$(31) \quad \mathbb{P} \left(2 \left\| \sum_{i \in H_2} \mathbf{W}_i v_i \right\|_{\infty} / b_n \leq \lambda_{\pi} / 2 \right) \geq 1 - 2/p^2.$$

By the assumptions of Theorem 4 and $b_n > n/4 - 1$, we can easily verify the assumption of Lemma 4 with $\tau = 3/4$. Thus, we apply Lemma 4 with $\tau = 3/4$ and obtain the restricted eigenvalue condition

$$(32) \quad \mathbb{P} \left(\kappa(s) > 0.015M^{-1} \right) \geq 1 - 4 \exp(-3b_n/3040),$$

where $\kappa(s)$ is defined in Lemma 4. Notice that due to Hölder's inequality, $\kappa(s)$ is smaller than the compatibility constant in Equation (6.4) of [Bühlmann and Van De Geer \(2011\)](#):

$$\begin{aligned} \kappa(s) &= \min_{|J| \subset \{1, \dots, p-1\}, |J| \leq s} \min_{\|q_{J^c}\|_1 \leq 3\|q_J\|_1} \frac{b_n^{-1} \sum_{i \in H_2} (\mathbf{W}_i^{\top} q)^2}{\|q_J\|_2^2} \\ &\leq \min_{|J| \subset \{1, \dots, p-1\}, |J| \leq s} \min_{\|q_{J^c}\|_1 \leq 3\|q_J\|_1} \frac{b_n^{-1} \sum_{i \in H_2} (\mathbf{W}_i^{\top} q)^2}{\|q_J\|_1^2 / s}. \end{aligned}$$

By (31) and (32), together with Theorem 6.1 of [Bühlmann and Van De Geer \(2011\)](#), we have that

$$\mathbb{P} \left(\|\widehat{\boldsymbol{\pi}} - \boldsymbol{\pi}\|_1 \leq 267s\lambda_{\pi}M \right) \geq 1 - 4 \exp(-3b_n/3040) - 2/p^2.$$

This proves the first claim. For the second claim, we simply follow the same argument as in (31) with H_2 replaced by H_4 . \square

PROOF OF LEMMA 6. We need to show that with high probability,

$$(33) \quad \left| \widehat{\boldsymbol{\xi}}_A^{\top} \boldsymbol{\pi}_A - \widehat{\boldsymbol{\xi}}_A^{\top} \widetilde{\boldsymbol{\pi}}_A \right| \leq \eta_{\pi},$$

and

$$\left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i (Z_i - \mathbf{W}_i^\top \boldsymbol{\pi}) \right\|_\infty \leq \lambda_\pi / 4$$

as well as

$$b_n^{-1} \sum_{i \in H_4} (Z_i - \mathbf{W}_i^\top \boldsymbol{\pi})^2 \geq \frac{1}{2M}.$$

Since $Z_i - \mathbf{W}_i^\top \boldsymbol{\pi} = v_i$, Lemmas 5 and 2 imply that

$$(34) \quad \mathbb{P} \left(\left\| b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i (Z_i - \mathbf{W}_i^\top \boldsymbol{\pi}) \right\|_\infty \leq \lambda_\pi / 4 \right) \geq 1 - 2/p^2 > 1 - 2/p$$

and

$$(35) \quad \mathbb{P} \left(b_n^{-1} \sum_{i \in H_4} (Z_i - \mathbf{W}_i^\top \boldsymbol{\pi})^2 \geq \frac{1}{2M} \right) \geq 1 - 2 \exp(-M^{-2} b_n / 44).$$

It remains to show (33). Notice that

$$\tilde{\boldsymbol{\pi}} - \boldsymbol{\pi} = \left(\mathbb{I}_{p-1} - \widehat{\boldsymbol{\Omega}}_{\mathbf{W}} b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) (\widehat{\boldsymbol{\pi}} - \boldsymbol{\pi}) + b_n^{-1} \sum_{i \in H_4} \widehat{\boldsymbol{\Omega}}_{\mathbf{W}} \mathbf{W}_i v_i$$

and thus

$$\begin{aligned} & \widehat{\boldsymbol{\xi}}_A^\top \tilde{\boldsymbol{\pi}}_A - \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\pi}_A \\ &= \underbrace{\widehat{\boldsymbol{\xi}}_A^\top \left(\mathbb{I}_{p-1} - \widehat{\boldsymbol{\Omega}}_{\mathbf{W}} b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) (\widehat{\boldsymbol{\pi}} - \boldsymbol{\pi})}_{T_1} + \underbrace{b_n^{-1} \sum_{i \in H_4} \widehat{\boldsymbol{\xi}}_A^\top \widehat{\boldsymbol{\Omega}}_{\mathbf{W}} \mathbf{W}_i v_i}_{T_2}. \end{aligned}$$

We proceed in two steps. We first bound T_1 and then T_2 .

Let \mathcal{B} denote the event that $\boldsymbol{\Omega}_{\mathbf{W}}$ satisfies the constraint in (13) for $\widehat{\boldsymbol{\Omega}}_{\mathbf{W}}$. By Lemma 3,

$$(36) \quad \mathbb{P}(\mathcal{B}) \geq 1 - 10/p - 2 \exp(-b_n/18).$$

Step 1: bound T_1

Notice that on the event \mathcal{B} , $\widehat{\Omega}_{\mathbf{W}}$ satisfies the constraint in (13) and therefore,

$$\begin{aligned} |T_1| &\leq \left\| \widehat{\xi}_A^\top \left(\mathbb{I}_{p-1} - \widehat{\Omega}_{\mathbf{W}} b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) \right\|_\infty \|\widehat{\pi} - \pi\|_1 \\ &\stackrel{(i)}{\leq} 24 \sqrt{b_n^{-1} \log p M^3 M_2} \|\widehat{\pi} - \pi\|_1, \end{aligned}$$

where (i) follows by the constraint in (13). By the bound in Lemma 5, we have that

$$(37) \quad \mathbb{P} \left(|T_1| > 6408 \sqrt{b_n^{-1} \log p M^4 M_2 s \lambda_\pi} \text{ and } \mathcal{B} \right) \leq 4 \exp(-3b_n/3040) + 2/p^2.$$

Step 2: bound T_2

Let \mathcal{F} be the σ -algebra generated by $\{(y_i, \mathbf{W}_i, Z_i)\}_{i \in H_1 \cup H_3}$ and $\{\mathbf{W}_i\}_{i \in H_4}$. Notice that $\{v_i\}_{i \in H_4}$ is independent of both $\{\mathbf{W}_i\}_{i \in H_4}$ and $\{(y_i, \mathbf{W}_i, Z_i)\}_{i \in H_1 \cup H_3}$. Hence, $\{v_i\}_{i \in H_4}$ is independent of \mathcal{F} . On the other hand, notice that $\{\widehat{\xi}_A^\top \widehat{\Omega}_{\mathbf{W}} \mathbf{W}_i\}_{i \in H_4}$ is \mathcal{F} -measurable. Since $\{v_i\}_{i \in H_4}$ is i.i.d $\mathcal{N}(0, \sigma_{\mathbf{V}}^2)$, we have that conditional on \mathcal{F} , T_2 is Gaussian with mean zero and variance

$$\widehat{\xi}_A^\top \widehat{\Omega}_{\mathbf{W}} \left(b_n^{-2} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) \widehat{\Omega}_{\mathbf{W}} \widehat{\xi}_A.$$

By the elementary bound of $\mathbb{P}(|X| > t\sigma) \leq 2 \exp(-t^2/2)$ for $X \sim \mathcal{N}(0, \sigma^2)$, we have that for any $t > 0$,

$$\mathbb{P} \left(|T_2| > t \sqrt{\widehat{\xi}_A^\top \widehat{\Omega}_{\mathbf{W}} \left(b_n^{-2} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) \widehat{\Omega}_{\mathbf{W}} \widehat{\xi}_A} \mid \mathcal{F} \right) \leq 2 \exp(-t^2/2).$$

We notice that, on the event \mathcal{B} , $\widehat{\Omega}_{\mathbf{W}}$ satisfies the constraint in (13) and thus

$$\widehat{\xi}_A^\top \widehat{\Omega}_{\mathbf{W}} \left(b_n^{-2} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top \right) \widehat{\Omega}_{\mathbf{W}} \widehat{\xi}_A \leq 32M^5 M_2^2 b_n^{-1}.$$

It follows that for any $t > 0$,

$$\mathbb{P} \left(|T_2| > 4t b_n^{-1/2} M^2 M_2 \sqrt{2M} \mid \mathcal{F} \right) \leq 2 \exp(-t^2/2).$$

We take $t = \sqrt{2 \log(100/\alpha)}$ and obtain

$$(38) \quad \mathbb{P} \left(|T_2| > 8b_n^{-1/2} M^2 M_2 \sqrt{M \log(100/\alpha)} \text{ and } \mathcal{B} \right) \leq 0.02\alpha.$$

Now we combine (36), (37) and (38), obtaining

$$\begin{aligned} \mathbb{P} (|T_1| + |T_2| > \eta\pi) &\leq 10/p + 2 \exp(-b_n/18) + 0.02\alpha + 4 \exp(-3b_n/3040) + 2/p^2 \\ &< 12/p + 0.02\alpha + 6 \exp(-3b_n/3040). \end{aligned}$$

Since $\widehat{\boldsymbol{\xi}}_A^\top \tilde{\boldsymbol{\pi}}_A - \widehat{\boldsymbol{\xi}}_A^\top \boldsymbol{\pi}_A = T_1 + T_2$, we have proved that (33) holds with probability at least $1 - 12/p - 0.02\alpha - 6 \exp(-3b_n/3040)$. By recalling (34) and (35), we complete the proof. \square

PROOF OF LEMMA 7. Let $\boldsymbol{\delta} = \check{\boldsymbol{\pi}} - \boldsymbol{\pi}$ and $\widehat{\boldsymbol{\Sigma}}_{\mathbf{W}} = b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i \mathbf{W}_i^\top$. Let $J_0 = \text{supp}(\boldsymbol{\pi})$. Define \mathcal{B} to be the event that $\boldsymbol{\pi}$ satisfies the constraint in (15) and $\kappa(s) \geq 0.24(1 - \tau)^2 M^{-1}$, where $\kappa(s)$ is defined in Lemma 4 and $\tau \in (0, 1)$ is a constant to be determined later.

On the event \mathcal{B} , we have that $\|\check{\boldsymbol{\pi}}\|_1 \leq \|\boldsymbol{\pi}\|_1$, which means $\|\boldsymbol{\pi} + \boldsymbol{\delta}_{J_0}\|_1 + \|\boldsymbol{\delta}_{J_0^c}\|_1 \leq \|\boldsymbol{\pi}\|_1$. Hence, on the event \mathcal{B} ,

$$(39) \quad \|\boldsymbol{\delta}_{J_0^c}\|_1 \leq \|\boldsymbol{\delta}_{J_0}\|_1.$$

Also observe that on the event \mathcal{B} , $\|b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i Z_i - \widehat{\boldsymbol{\Sigma}}_{\mathbf{W}} \boldsymbol{\pi}\|_\infty \leq \lambda_\pi/4$ and $\|b_n^{-1} \sum_{i \in H_4} \mathbf{W}_i Z_i - \widehat{\boldsymbol{\Sigma}}_{\mathbf{W}} \check{\boldsymbol{\pi}}\|_\infty \leq \lambda_\pi/4$, which means

$$\|\widehat{\boldsymbol{\Sigma}} \boldsymbol{\delta}\|_\infty \leq \lambda_\pi/2.$$

Therefore, on the event \mathcal{B} ,

$$\begin{aligned} \boldsymbol{\delta}^\top \widehat{\boldsymbol{\Sigma}} \boldsymbol{\delta} &\leq \|\boldsymbol{\delta}\|_1 \|\widehat{\boldsymbol{\Sigma}}_{\mathbf{W}} \boldsymbol{\delta}\|_\infty \leq 0.5 \lambda_\pi \|\boldsymbol{\delta}\|_1 \\ &= 0.5 \lambda_\pi (\|\boldsymbol{\delta}_{J_0}\|_1 + \|\boldsymbol{\delta}_{J_0^c}\|_1) \stackrel{(i)}{\leq} \lambda_\pi \|\boldsymbol{\delta}_{J_0}\|_1 \leq \lambda_\pi \sqrt{s} \|\boldsymbol{\delta}_{J_0}\|_2, \end{aligned}$$

where (i) follows by (39).

On the other hand, we can lower bound $\boldsymbol{\delta}^\top \widehat{\boldsymbol{\Sigma}} \boldsymbol{\delta}$ via the restricted eigenvalue condition. By (39), we have that on the event \mathcal{B} , $\|\boldsymbol{\delta}_{J_0^c}\|_1 \leq \|\boldsymbol{\delta}_{J_0}\|_1 \leq 3 \|\boldsymbol{\delta}_{J_0}\|_1$. Thus, we have that

$$\boldsymbol{\delta}^\top \widehat{\boldsymbol{\Sigma}} \boldsymbol{\delta} \geq \kappa(s) \|\boldsymbol{\delta}_{J_0}\|_2^2 \geq 0.24(1 - \tau)^2 M^{-1} \|\boldsymbol{\delta}_{J_0}\|_2^2.$$

Now we combine the above two displays and obtain that on the event \mathcal{B} ,

$$\|\boldsymbol{\delta}_{J_0}\|_2 \leq \frac{M \lambda_\pi \sqrt{s}}{0.24(1 - \tau)^2}.$$

Therefore, (39) implies that on the event \mathcal{B} ,

$$\|\delta\|_1 \leq 2\|\delta_{J_0}\|_1 \leq 2\sqrt{s}\|\delta_{J_0}\|_2 \leq \frac{2M\lambda_\pi s}{0.24(1-\tau)^2}.$$

Notice that by Lemmas 4 and 6,

$$\begin{aligned} \mathbb{P}(\mathcal{B}) &\geq 1 - 14/p - 0.02\alpha - 6 \exp(-3b_n/3040) \\ &\quad - 2 \exp(-M^{-2}b_n/44) - 4 \exp(-\tau^2 b_n/570). \end{aligned}$$

Hence, the desired result follows by choosing $\tau = 3/4$. \square

L.3. Proof of auxiliary results used in proving Theorem 8.

PROOF OF LEMMA 8. Clearly, we always have $L(\Theta_1, \Theta) \leq L(\Theta, \Theta)$. We only need to show the other direction. Let $c > 0$ be a constant such that

$$c \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \leq \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

Notice that

$$\begin{aligned} L(\Theta_1, \Theta) &= \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta_1} \mathbb{E}_\theta \text{diam}(CI) \\ &\geq \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta_1} \mathbb{E}_\theta \text{diam}(CI) \\ &\geq \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \\ &\geq c \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) = cL(\Theta, \Theta). \end{aligned}$$

The proof is complete. \square

PROOF OF LEMMA 9. Clearly, we have

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \geq \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

We only need to show the other direction. Let $CI_* \in \mathcal{C}_\alpha(\Theta)$ and $\theta_* \in \Theta$ be such that

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \geq 0.9 \mathbb{E}_{\theta_*} \text{diam}(CI_*).$$

Now define $\Theta_1 = \{\theta_*\}$. Clearly,

$$\mathbb{E}_{\theta_*} \text{diam}(CI_*) = \sup_{\theta \in \Theta_1} \mathbb{E}_\theta \text{diam}(CI_*) \geq \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta_1} \mathbb{E}_\theta \text{diam}(CI).$$

By the assumption of $cL(\Theta, \Theta) \leq L(\Theta_1, \Theta)$, we have

$$\mathbb{E}_{\theta_*} \text{diam}(CI_*) \geq c \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

Hence,

$$\inf_{CI \in \mathcal{C}_\alpha(\Theta)} \inf_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI) \geq 0.9c \inf_{CI \in \mathcal{C}_\alpha(\Theta)} \sup_{\theta \in \Theta} \mathbb{E}_\theta \text{diam}(CI).$$

The proof is complete. \square

L.4. Proof of auxiliary results used in proving Theorem 10.

PROOF OF LEMMA 10. Due to length of the work we comment that the result above is quite easy to verify. We leave the details to the reader. \square

PROOF OF LEMMA 11. If $(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim (\theta \odot D)$ with

$$\theta = (\beta, \gamma, \Sigma, \sigma) \in \tilde{\Theta}_{N_1, N_2}(s),$$

then

$$\mathbf{y} = \mathbf{Z}\beta D + \mathbf{W}\gamma D + \varepsilon$$

with $\varepsilon \sim \mathcal{N}_n(0, \mathbb{I}_n(\sigma D)^2)$ and rows of $[\mathbf{Z}, \mathbf{W}]$ being i.i.d $N(0, \Sigma)$. Now we divide both sides by D , obtaining

$$\mathbf{y}D^{-1} = \mathbf{Z}\beta + \mathbf{W}\gamma + \tilde{\varepsilon}$$

with $\tilde{\varepsilon} = \varepsilon D^{-1}$. Notice that $\tilde{\varepsilon} \sim \mathcal{N}_n(0, \mathbb{I}_n \sigma^2)$ and rows of $[\mathbf{Z}, \mathbf{W}]$ being i.i.d $N(0, \Sigma)$. Thus, $(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta$. This shows the ‘‘only if’’ direction. The ‘‘if’’ direction follows by an analogous argument. \square

PROOF OF LEMMA 12. Here, for notational simplicity, we use $|\cdot|$ to denote $\text{diam}(\cdot)$. Fix any $\eta > 0$. By the definition of infimum, there exists $T_* \in \mathcal{C}_\alpha(\tilde{\Theta}_{N_1, N_2}(s))$ satisfying

$$\begin{aligned} \mathbb{A}(s, N_1, N_2) &= \inf_{T \in \mathcal{C}_\alpha(\tilde{\Theta}_{N_1, N_2}(s))} \sup_{\theta \in \tilde{\Theta}_{N_1, N_2}(s)} \mathbb{E}_\theta |T(\mathbf{y}, \mathbf{Z}, \mathbf{W})| \\ (40) \quad &\geq \sup_{\theta \in \tilde{\Theta}_{N_1, N_2}(s)} \mathbb{E}_\theta |T_*(\mathbf{y}, \mathbf{Z}, \mathbf{W})| - \eta. \end{aligned}$$

Define \tilde{T} by

$$\tilde{T}(\mathbf{y}, \mathbf{Z}, \mathbf{W}) = D \cdot T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W}).$$

For an arbitrary $\theta_0 = (\beta_0, \gamma_0, \boldsymbol{\Sigma}, \sigma_0) \in \tilde{\Theta}_{DN_1, DN_2}(s)$, we define $\theta_1 = (\beta_1, \gamma_1, \boldsymbol{\Sigma}, \sigma_1)$ by

$$\beta_1 = \beta_0 D^{-1}, \quad \gamma_1 = \gamma_0 D^{-1}, \quad \text{and} \quad \sigma_1 = \sigma_0 D^{-1}.$$

Notice that $\theta_0 = \theta_1 \odot D$. Notice that

$$|\tilde{T}(\mathbf{y}, \mathbf{Z}, \mathbf{W})| = D |T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})|.$$

Therefore,

$$\begin{aligned} & \sup_{\theta_1 \in \Theta(s, DN_1, DN_2)} \mathbb{E}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta_1} |\tilde{T}(\mathbf{y}, \mathbf{Z}, \mathbf{W})| \\ &= D \sup_{\theta_1 \in \tilde{\Theta}_{DN_1, DN_2}(s)} \mathbb{E}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta_1} |T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})| \\ &\stackrel{(i)}{=} D \sup_{\theta \in \tilde{\Theta}_{N_1, N_2}(s)} \mathbb{E}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim (\theta \odot D)} |T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})| \\ &\stackrel{(ii)}{=} D \sup_{\theta \in \tilde{\Theta}_{N_1, N_2}(s)} \mathbb{E}_{(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W}) \sim \theta} |T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})| \\ (41) \quad &\stackrel{(iii)}{\leq} D(\mathbb{A}(s, N_1, N_2) + \eta), \end{aligned}$$

where (i) follows by Lemma 10, (ii) follows by Lemma 11 and (iii) follows by (40).

Now we show $\tilde{T} \in \mathcal{C}_\alpha(\Theta(s, DN_1, DN_2))$. Notice that

$$\begin{aligned} \mathbb{P}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta_0}(\beta_0 \in \tilde{T}(\mathbf{y}, \mathbf{Z}, \mathbf{W})) &= \mathbb{P}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta_0}(\beta_1 D \in D \cdot T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})) \\ &= \mathbb{P}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta_0}(\beta_1 \in T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})) \\ &= \mathbb{P}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim (\theta_1 \odot D)}(\beta_1 \in T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})) \\ &\stackrel{(i)}{=} \mathbb{P}_{(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W}) \sim \theta_1}(\beta_1 \in T_*(\mathbf{y}D^{-1}, \mathbf{Z}, \mathbf{W})) \\ &\stackrel{(ii)}{\geq} 1 - \alpha, \end{aligned}$$

where (i) follows by Lemma 11 and (ii) follows by $T_* \in \mathcal{C}_\alpha(\tilde{\Theta}_{N_1, N_2}(s))$ and $\theta_1 \in \tilde{\Theta}_{N_1, N_2}(s)$. Hence, $\tilde{T} \in \mathcal{C}_\alpha(\Theta(s, DN_1, DN_2))$ and

$$\sup_{\theta_1 \in \Theta(s, DN_1, DN_2)} \mathbb{E}_{(\mathbf{y}, \mathbf{Z}, \mathbf{W}) \sim \theta_1} |\tilde{T}(\mathbf{y}, \mathbf{Z}, \mathbf{W})| \geq \mathbb{A}(s, DN_1, DN_2).$$

By (41), it follows that

$$D(\mathbb{A}(s, N_1, N_2) + \eta) \geq \mathbb{A}(s, DN_1, DN_2).$$

Since $\eta > 0$ is arbitrary, we have $D\mathbb{A}(s, N_1, N_2) \geq \mathbb{A}(s, DN_1, DN_2)$. \square

L.5. Proof of auxiliary results used in proving Theorem 11.

PROOF OF LEMMA 8. Let $\lambda_{\min}(\cdot)$ denote the minimal eigenvalue. Then

$$\mathbb{P}\left(\mathbf{Z}^\top (\mathbf{W}\mathbf{W}^\top)^{-1} \mathbf{Z} > a\right) \leq \mathbb{P}\left(\lambda_{\max}[(\mathbf{W}\mathbf{W}^\top)^{-1}] \|\mathbf{Z}\|_2^2 > a\right) = \mathbb{P}\left(\|\mathbf{Z}\|_2^2 > \lambda_{\min}(\mathbf{W}\mathbf{W}^\top) a\right).$$

By Corollary 5.35 of Vershynin (2010), we have that

$$\mathbb{P}\left(\sqrt{\lambda_{\min}(\mathbf{W}\mathbf{W}^\top)} < \sqrt{2n} - \sqrt{n} - 0.1\sqrt{n}\right) \leq 2\exp(-0.01n/2).$$

Since $\sqrt{2} - 1 - 0.1 > 0.3$, we have that

$$\begin{aligned} \mathbb{P}\left(\mathbf{Z}^\top (\mathbf{W}\mathbf{W}^\top)^{-1} \mathbf{Z} > a\right) &\leq \mathbb{P}\left(\|\mathbf{Z}\|_2^2 > \lambda_{\min}(\mathbf{W}\mathbf{W}^\top) a\right) \\ &\leq 2\exp(-0.01n/2) + \mathbb{P}\left(\|\mathbf{Z}\|_2^2 > 0.09na\right) \\ &\leq 2\exp(-0.01n/2) + \frac{\mathbb{E}\|\mathbf{Z}\|_2^2}{0.09na} \\ &= 2\exp(-0.01n/2) + \frac{n}{0.09na} \\ &< 2\exp(-0.01n/2) + \frac{12}{a}. \end{aligned}$$

\square

PROOF OF LEMMA 9. We first notice that

$$\mathbb{E}\|\boldsymbol{\xi}\|_2^2 = \mathbb{E}\boldsymbol{\xi}^\top \boldsymbol{\xi} = \mathbb{E}\text{trace}(\boldsymbol{\xi}\boldsymbol{\xi}^\top) = \text{trace}(\mathbb{E}\boldsymbol{\xi}\boldsymbol{\xi}^\top) = \text{trace}(\boldsymbol{\Sigma}).$$

Then the desired result follows by Markov's inequality

$$\mathbb{P}(\|\boldsymbol{\xi}\|_2 > x) = \mathbb{P}(\|\boldsymbol{\xi}\|_2^2 > x^2) \leq x^{-2} \mathbb{E}\|\boldsymbol{\xi}\|_2^2.$$

\square

PROOF OF LEMMA 10. We use an argument that is inspired by the proof of Proposition 1 of Carpentier and Verzelen (2019). Let $C_1 > 0$ be a constant to be chosen later. Let $\mu_n(\cdot)$ denote the probability measure of the

Gaussian distribution $\mathcal{N}(0, \mathbb{I}_{p-1} C_1^2 (p-1)^{-1})$. Recall for the parameter $\theta = (\beta, \gamma, \Sigma, \sigma) \in \Theta_*(r)$, we have $\sigma = 0$ and $\Sigma = \mathbb{I}_p$. Thus, we can write $\mathbf{y} = \mathbf{Z}\beta + \mathbf{W}\gamma$ with $\gamma \in \mathbb{R}^{p-1}$, where entries of $\mathbf{Z} \in \mathbb{R}^n$ and $\mathbf{W} \in \mathbb{R}^{n \times (p-1)}$ are i.i.d standard normal random variables.

Since $p \geq 2n+1$, we can without loss of generality set $p = 2n+1$ and hence $\mathbf{W} \in \mathbb{R}^{n \times 2n}$. If $p > 2n+1$, then we can simply apply this distribution to the first $(2n+1)$ elements of γ and leave the other $(p-2n-1)$ elements to be zero; since doing so would create additional unnecessary notations without really changing the argument, we work with $p = 2n+1$ for notational simplicity. Define two probability measures

$$\mathbb{P}_{[A]} = \int_{\mathbb{R}^{p-1}} \mathbb{P}_{(0, \gamma, \mathbb{I}_p, 0)} d\mu_n(\gamma)$$

and

$$\mathbb{P}_{[B]} = \int_{\mathbb{R}^{p-1}} \mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)} d\mu_n(\gamma),$$

where $r_n > 0$ is a sequence to be determined. We define the event

$$\mathcal{A} = \left\{ \mathbf{Z}^\top (\mathbf{W}\mathbf{W}^\top)^{-1} \mathbf{Z} \leq C_2 \right\},$$

where $C_2 > 0$ is a constant to be determined. For a fixed (\mathbf{W}, \mathbf{Z}) , \mathbf{y} follows $\mathcal{N}(0, \mathbf{W}\mathbf{W}^\top C_1^2 (p-1)^{-1})$ under $\mathbb{P}_{[A]}$ and follows $\mathcal{N}(\mathbf{Z}r_n, \mathbf{W}\mathbf{W}^\top C_1^2 (p-1)^{-1})$ under $\mathbb{P}_{[B]}$.

Let $\mathcal{B} = \{v \in \mathbb{R}^{p-1} : \|v\|_2 \leq 1\}$. Let $\tilde{\mu}_n(\cdot)$ be the truncated Gaussian measure on \mathcal{B} , i.e., $\tilde{\mu}_n(C) = \mu_n(C \cap \mathcal{B}) / \mu_n(\mathcal{B})$ for any set C . Define

$$\mathbb{P}_{\tilde{A}} = \int_{\mathbb{R}^{p-1}} \mathbb{P}_{(0, \gamma, \mathbb{I}_p, 0)} d\tilde{\mu}_n(\gamma).$$

and

$$\mathbb{P}_{\tilde{B}} = \int_{\mathbb{R}^{p-1}} \mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)} d\tilde{\mu}_n(\gamma).$$

The rest of proof proceeds in three steps in which we bound (1) difference between $\mathbb{P}_{[A]}$ and $\mathbb{P}_{[B]}$, (2) difference between $\mathbb{P}_{[A]}$ and $\mathbb{P}_{\tilde{A}}$ and (3) difference between $\mathbb{P}_{[B]}$ and $\mathbb{P}_{\tilde{B}}$.

Step 1: bound the difference between $\mathbb{P}_{[A]}$ and $\mathbb{P}_{[B]}$

Let $\tilde{\gamma}$ be a random vector that is independent of $(\mathbf{y}, \mathbf{W}, \mathbf{Z})$ and has the distribution μ_n . Let $\mathcal{Y} \times \mathcal{W} \times \mathcal{Z}$ be the support of $(\mathbf{y}, \mathbf{W}, \mathbf{Z})$. We notice that

$$\mathbb{E}_{\mathbb{P}_{[B]}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) = \int_{\mathcal{Y} \times \mathcal{W} \times \mathcal{Z}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) d\mathbb{P}_{[B]}(\mathbf{y}, \mathbf{W}, \mathbf{Z})$$

$$\begin{aligned}
&= \int_{\mathcal{Y} \times \mathcal{W} \times \mathcal{Z}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \left(\int_{\mathbb{R}^n} d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{y}, \mathbf{W}, \mathbf{Z}) d\mu_n(\gamma) \right) \\
&\stackrel{(i)}{=} \int_{\mathbb{R}^{p-1}} \int_{\mathcal{Y} \times \mathcal{W} \times \mathcal{Z}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{y}, \mathbf{W}, \mathbf{Z}) d\mu_n(\gamma) \\
&= \int_{\mathbb{R}^{p-1}} \int_{\mathcal{W} \times \mathcal{Z}} \psi(\mathbf{Z}r_n + \mathbf{W}\gamma, \mathbf{W}, \mathbf{Z}) d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{W}, \mathbf{Z}) d\mu_n(\gamma) \\
&= \int_{\mathcal{W} \times \mathcal{Z}} \int_{\mathbb{R}^{p-1}} \psi(\mathbf{Z}r_n + \mathbf{W}\gamma, \mathbf{W}, \mathbf{Z}) d\mu_n(\gamma) d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{W}, \mathbf{Z}) \\
&= \mathbb{E}\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z})
\end{aligned}$$

with \mathbb{E} being expectation over random elements \mathbf{W} , \mathbf{Z} and $\tilde{\gamma}$, where (i) and (ii) follow by Fubini's theorem (since $\psi(\mathbf{Z}r_n + \mathbf{W}\gamma, \mathbf{W}, \mathbf{Z}) d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{W}, \mathbf{Z})$ is integrable). Here, notice that \mathbf{W} , \mathbf{Z} and $\tilde{\gamma}$ are mutually independent, where entries of \mathbf{W} and \mathbf{Z} follow the standard normal distribution.

Similarly, we have

$$\mathbb{E}_{\mathbb{P}_{[\mathcal{A}]}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) = \mathbb{E}\psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}).$$

Let $Q_{(w, z; r_n)}(\cdot)$ denote the distribution

$$\mathcal{N}(zr_n, ww^\top C_1^2(p-1)^{-1}).$$

Then we have

$$\begin{aligned}
\mathbb{E}\psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) &= \mathbb{E} \left(\int_{\mathcal{Y}} \psi(y, \mathbf{W}, \mathbf{Z}) Q_{(\mathbf{W}, \mathbf{Z}; 0)}(dy) \right) \\
&= \mathbb{E} \left(\int_{\mathcal{Y}} \psi(y, \mathbf{W}, \mathbf{Z}) Q_{(\mathbf{W}, \mathbf{Z}; r_n)}(dy) \frac{Q_{(\mathbf{W}, \mathbf{Z}; 0)}(dy)}{Q_{(\mathbf{W}, \mathbf{Z}; r_n)}(dy)} \right).
\end{aligned}$$

Moreover,

$$\begin{aligned}
&\left| \mathbb{E}_{\mathbb{P}_{[\mathcal{A}]}} \psi - \mathbb{E}_{\mathbb{P}_{[\mathcal{B}]}} \psi \right| \\
&= |\mathbb{E}[\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) - \psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z})]| \\
&= |\mathbb{E} \{ \mathbb{E}[\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) - \psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) \mid \mathbf{W}, \mathbf{Z}] \}| \\
&= |\mathbb{E} \{ \mathbf{1}_{\mathcal{A}} \times \mathbb{E}[\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) - \psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) \mid \mathbf{W}, \mathbf{Z}] \}| \\
&\quad + |\mathbb{E} \{ \mathbf{1}_{\mathcal{A}^c} \times \mathbb{E}[\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) - \psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) \mid \mathbf{W}, \mathbf{Z}] \}| \\
&\leq |\mathbb{E} \{ \mathbf{1}_{\mathcal{A}} \times \mathbb{E}[\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) - \psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) \mid \mathbf{W}, \mathbf{Z}] \}| + \mathbb{P}(\mathcal{A}^c) \\
&\leq \mathbb{E} |\mathbf{1}_{\mathcal{A}} \times \mathbb{E}[\psi(\mathbf{Z}r_n + \mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) - \psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) \mid \mathbf{W}, \mathbf{Z}]| + \mathbb{P}(\mathcal{A}^c)
\end{aligned}$$

$$\begin{aligned}
&= \mathbb{E} \left| \mathbf{1}_{\mathcal{A}} \times \mathbb{E} \left[\psi(\mathbf{W}\tilde{\gamma}, \mathbf{W}, \mathbf{Z}) \left(\frac{dQ_{(\mathbf{W}, \mathbf{Z}; r_n)}(\mathbf{W}\tilde{\gamma})}{dQ_{(\mathbf{W}, \mathbf{Z}; 0)}}(\mathbf{W}\tilde{\gamma}) - 1 \right) \mid \mathbf{W}, \mathbf{Z} \right] \right| + \mathbb{P}(\mathcal{A}^c) \\
&\leq \mathbb{E} \left| \mathbf{1}_{\mathcal{A}} \times \mathbb{E} \left[\left| \frac{dQ_{(\mathbf{W}, \mathbf{Z}; r_n)}(\mathbf{W}\tilde{\gamma})}{dQ_{(\mathbf{W}, \mathbf{Z}; 0)}}(\mathbf{W}\tilde{\gamma}) - 1 \right| \mid \mathbf{W}, \mathbf{Z} \right] \right| + \mathbb{P}(\mathcal{A}^c) \\
&= \mathbb{E} (\mathbf{1}_{\mathcal{A}} \times \text{TV}(Q_{(\mathbf{W}, \mathbf{Z}; r_n)}, Q_{(\mathbf{W}, \mathbf{Z}; 0)})) + \mathbb{P}(\mathcal{A}^c) \\
&\stackrel{(i)}{\leq} \mathbb{E} \left(\mathbf{1}_{\mathcal{A}} \times \sqrt{\text{KL}(Q_{(\mathbf{W}, \mathbf{Z}; r_n)}, Q_{(\mathbf{W}, \mathbf{Z}; 0)})/2} \right) + \mathbb{P}(\mathcal{A}^c),
\end{aligned}$$

where (i) follows by the first Pinsker's inequality (Lemma 2.5 of [Tsybakov \(2008\)](#)). By Lemma 7, we have

$$\text{KL}(Q_{(\mathbf{W}, \mathbf{Z}; r_n)}, Q_{(\mathbf{W}, \mathbf{Z}; 0)}) = \frac{1}{2} r_n^2 \mathbf{Z}^\top \left[\mathbf{W}\mathbf{W}^\top C_1^2 n^{-1} \right]^{-1} \mathbf{Z} = \frac{nr_n^2}{2C_1^2} \mathbf{Z}^\top (\mathbf{W}\mathbf{W}^\top)^{-1} \mathbf{Z}.$$

Thus,

$$\mathbf{1}_{\mathcal{A}} \times \text{KL}(Q_{(\mathbf{W}, \mathbf{Z}; r_n)}, Q_{(\mathbf{W}, \mathbf{Z}; 0)}) = \frac{nr_n^2}{2C_1^2} \mathbf{Z}^\top (\mathbf{W}\mathbf{W}^\top)^{-1} \mathbf{Z} \times \mathbf{1}_{\mathcal{A}} \leq \frac{nr_n^2 C_2}{2C_1^2}$$

and

$$\left| \mathbb{E}_{\mathbb{P}_{[\mathcal{A}]}} \psi - \mathbb{E}_{\mathbb{P}_{[\mathcal{B}]}} \psi \right| \leq \frac{n^{1/2} r_n \sqrt{C_2}}{2C_1} + \mathbb{P}(\mathcal{A}^c).$$

Fix an arbitrary $\alpha > 0$. By Lemma 8, there exists a constant C_2 depending only on α such that $\mathbb{P}(\mathcal{A}^c) \leq \alpha/4$. Then we take $r_n = n^{-1/2} C_2^{-1/2} C_1 \alpha/2$ and obtain that

$$(42) \quad \left| \mathbb{E}_{\mathbb{P}_{[\mathcal{A}]}} \psi - \mathbb{E}_{\mathbb{P}_{[\mathcal{B}]}} \psi \right| \leq \alpha/2.$$

Step 2: bound the difference between $\mathbb{P}_{[\mathcal{B}]}$ and $\mathbb{P}_{\tilde{\mathcal{B}}}$.

Recall from Step 1 that

$$\begin{aligned}
\mathbb{E}_{\mathbb{P}_{[\mathcal{B}]}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) &= \int_{\mathcal{W} \times \mathcal{Z}} \int_{\mathbb{R}^{p-1}} \psi(\mathbf{Z}r_n + \mathbf{W}\gamma, \mathbf{W}, \mathbf{Z}) d\mu_n(\gamma) d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{W}, \mathbf{Z}) \\
&= \int_{\mathbb{R}^{p-1}} \phi(\gamma) d\mu_n(\gamma),
\end{aligned}$$

where $\phi(\gamma) = \int_{\mathcal{W} \times \mathcal{Z}} \psi(\mathbf{Z}r_n + \mathbf{W}\gamma, \mathbf{W}, \mathbf{Z}) d\mathbb{P}_{(r_n, \gamma, \mathbb{I}_p, 0)}(\mathbf{W}, \mathbf{Z})$. Similarly, we have

$$\mathbb{E}_{\mathbb{P}_{\tilde{\mathcal{B}}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) = \int_{\mathbb{R}^{p-1}} \phi(\gamma) d\tilde{\mu}_n(\gamma) = \frac{1}{\mu_n(\mathcal{B})} \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma).$$

We observe

$$\begin{aligned}
& \left| \mathbb{E}_{\mathbb{P}_{[\mathcal{B}]}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) - \mathbb{E}_{\mathbb{P}_{\tilde{\mathcal{B}}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \right| \\
&= \left| \int_{\mathbb{R}^{p-1}} \phi(\gamma) d\mu_n(\gamma) - \frac{1}{\mu_n(\mathcal{B})} \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma) \right| \\
&= \left| \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma) + \int_{\mathcal{B}^c} \phi(\gamma) d\mu_n(\gamma) - \frac{1}{\mu_n(\mathcal{B})} \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma) \right| \\
&\leq \left| \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma) - \frac{1}{\mu_n(\mathcal{B})} \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma) \right| + \left| \int_{\mathcal{B}^c} \phi(\gamma) d\mu_n(\gamma) \right| \\
&= \left| 1 - \frac{1}{\mu_n(\mathcal{B})} \right| \times \left| \int_{\mathcal{B}} \phi(\gamma) d\mu_n(\gamma) \right| + \left| \int_{\mathcal{B}^c} \phi(\gamma) d\mu_n(\gamma) \right| \\
&\stackrel{(i)}{\leq} \left| 1 - \frac{1}{\mu_n(\mathcal{B})} \right| + \mu_n(\mathcal{B}^c) = \frac{\mu_n(\mathcal{B}^c)}{1 - \mu_n(\mathcal{B}^c)} + \mu_n(\mathcal{B}^c),
\end{aligned}$$

where (i) follows by $|\phi(\gamma)| \leq 1$ (since $|\psi| \leq 1$). By Lemma 9,

$$\mu_n(\mathcal{B}^c) \leq \text{trace}(\mathbb{I}_{p-1} C_1^2 (p-1)^{-1}) = C_1^2.$$

Now we choose $C_1 = \sqrt{\alpha/12}$. This means that $\mu_n(\mathcal{B}^c) \leq \alpha/12$. Hence, $\mu_n(\mathcal{B}^c) < 1/2$.

$$\begin{aligned}
& \left| \mathbb{E}_{\mathbb{P}_{[\mathcal{B}]}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) - \mathbb{E}_{\mathbb{P}_{\tilde{\mathcal{B}}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \right| \\
(43) \quad & \leq \frac{\mu_n(\mathcal{B}^c)}{1 - \mu_n(\mathcal{B}^c)} + \mu_n(\mathcal{B}^c) \leq 2\mu_n(\mathcal{B}^c) + \mu_n(\mathcal{B}^c) \leq \alpha/4.
\end{aligned}$$

Step 3: bound the difference between $\mathbb{P}_{[A]}$ and $\mathbb{P}_{\tilde{A}}$.

Similarly to Step 2, we can show that

$$(44) \quad \left| \mathbb{E}_{\mathbb{P}_{[A]}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) - \mathbb{E}_{\mathbb{P}_{\tilde{A}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \right| \leq \alpha/4.$$

Now we combine (42), (43) and (44), obtaining

$$\left| \mathbb{E}_{\mathbb{P}_{\tilde{A}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) - \mathbb{E}_{\mathbb{P}_{\tilde{\mathcal{B}}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \right| \leq \alpha.$$

Since $\sup_{\theta \in \Theta_*(0)} \mathbb{E}_\theta \psi \leq \alpha$ and $\mathbb{P}_{\tilde{A}}$ is by definition a mixture of distributions in $\Theta_*(0)$, we have $\mathbb{E}_{\mathbb{P}_{\tilde{A}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \leq \alpha$, which means

$$\mathbb{E}_{\mathbb{P}_{\tilde{\mathcal{B}}}} \psi(\mathbf{y}, \mathbf{W}, \mathbf{Z}) \leq 2\alpha.$$

Notice that $\mathbb{P}_{\tilde{\mathcal{B}}}$ is a mixture of distributions in $\Theta_*(r_n)$, we have that

$$\inf_{\theta \in \Theta_*(r_n)} \mathbb{E}_\theta \psi \leq 2\alpha.$$

The proof is complete since

$$r_n = n^{-1/2} C_2^{-1/2} C_1 \alpha / 2$$

with C_1, C_2 depending only on α . \square

APPENDIX M: COMPARISON OF PRIORS

To provide a comparison of the priors, we outline an adaptation of the prior from [Cai and Guo \(2017\)](#) and compare with our prior for the proof of minimax lower bound. This comparison illustrates the main differences. (We thank an anonymous reviewer for suggesting this.)

A simple adaptation of the prior considered in [Cai and Guo \(2017\)](#) under our notation: $\mathbf{y} = \mathbf{Z}\beta + \mathbf{W}\gamma + \varepsilon$ and $\Sigma = \begin{pmatrix} \boldsymbol{\pi}^\top \boldsymbol{\pi} + \sigma_{\mathbf{V}}^2 & \boldsymbol{\pi}^\top \\ \boldsymbol{\pi} & I_{p-1} \end{pmatrix}$. Let the parameter be indexed by $(\beta, \gamma, \boldsymbol{\pi}, \sigma_{\mathbf{V}}, \sigma_\varepsilon)$.

The priors used by [Cai and Guo \(2017\)](#) in Equation (7.13) on page 636 therein can be adapted (switching) as follows. Given $(\beta_*, \gamma_*, 0, 1, \sigma_0)$, their prior is

$$\begin{aligned} \beta &= \beta_* \\ \gamma &= \gamma_* + c_1 \sqrt{\frac{\log(p/m^2)}{n}} \boldsymbol{\delta} \\ \boldsymbol{\pi} &= c_2 \sqrt{\frac{\log(p/m^2)}{n}} \boldsymbol{\delta} \\ \sigma_{\mathbf{V}} &= \sqrt{1 - c_2^2 \frac{m \log(p/m^2)}{n}} \\ \sigma_\varepsilon &= \sigma_0 \end{aligned}$$

where $\boldsymbol{\delta}$ is from the uniform distribution from the set $\mathcal{M} = \{\mathbf{v} \in \{0, 1\}^{p-1} : \|\mathbf{v}\|_0 = m\}$. Here, $\|(\beta_*, \gamma_*^\top)^\top\|_0 = m$ and $c_1 > 0$ is a constant.

Here is our prior in Definition 3 from Appendix A. Given $(\beta_*, \gamma_*, \boldsymbol{\pi}_*, \sigma_{\mathbf{V},*}, \sigma_{\varepsilon,*})$ with $\|\boldsymbol{\pi}_*\|_0 = s$ (with $m \asymp s$), we define

$$\begin{aligned} \beta &= \beta_* - h \quad \text{with} \quad h = \frac{ds \log p}{n} \\ \gamma &= \gamma_* + h\boldsymbol{\pi} + r(1-h)\sigma_{\varepsilon,*} \sqrt{2d \log(p)/n} \boldsymbol{\delta} = \gamma_* + \frac{ds \log p}{n} \boldsymbol{\pi}_* + \sigma_{\mathbf{V},*} \sqrt{\frac{2d \log p}{n}} \boldsymbol{\delta} \\ \boldsymbol{\pi} &= \boldsymbol{\pi}_* + \sigma_{\mathbf{V},*} \sqrt{\frac{2d \log p}{n}} \boldsymbol{\delta} \\ \sigma_{\mathbf{V}} &= \sigma_{\mathbf{V},*} \sqrt{1 - \frac{ds \log p}{n}} \end{aligned}$$

$$\sigma_\varepsilon = \sigma_{\varepsilon,*}$$

where $d > 0$ is a constant and $r = \sigma_{\mathbf{V},*}/\sigma_{\varepsilon,*}$.

From the above comparison, the difference between our prior and that in [Cai and Guo \(2017\)](#) is not simply that γ and π are switched. Notice that in our prior, the construction of γ depends on π_* , whereas in [Cai and Guo \(2017\)](#), π_* is set to be zero. A priori, it is not obvious whether there exists a construction of γ under nonzero π_* such that the calculation in our Appendix A would go through. From this perspective, the prior of [Cai and Guo \(2017\)](#) is just a special case of our construction. For the *uniform* non-testability result to hold, we need to build the prior around a general point $(\beta_*, \gamma_*, \pi_*, \sigma_{\mathbf{V},*}, \sigma_{\varepsilon,*})$.

DEPARTMENT OF MATHEMATICS,
HALICIOĞLU DATA SCIENCE INSTITUTE,
UNIVERSITY OF CALIFORNIA, SAN DIEGO,
E-MAIL: jbradic@ucsd.edu

DEPARTMENT OF OPERATIONS RESEARCH
AND FINANCIAL ENGINEERING,
PRINCETON UNIVERSITY,
E-MAIL: jqfan@princeton.edu

LUNDQUIST COLLEGE OF BUSINESS,
UNIVERSITY OF OREGON,
E-MAIL: yzhu6@oregon.edu