



Contents lists available at ScienceDirect

Journal of Econometrics

journal homepage: [www.elsevier.com/locate/jeconom](http://www.elsevier.com/locate/jeconom)

# Inference for high-dimensional instrumental variables regression<sup>☆</sup>

David Gold<sup>a</sup>, Johannes Lederer<sup>a</sup>, Jing Tao<sup>b,\*</sup><sup>a</sup> Department of Statistics, University of Washington, Seattle, WA, 98195, USA<sup>b</sup> Department of Economics, University of Washington, Seattle, WA, 98195, USA

## ARTICLE INFO

### Article history:

Received 31 July 2018

Received in revised form 19 June 2019

Accepted 29 September 2019

Available online xxxx

### JEL classification:

C14

C31

C36

### Keywords:

High-dimensional inference

Instrumental variables

De-biasing

## ABSTRACT

This paper concerns statistical inference for the components of a high-dimensional regression parameter despite possible endogeneity of each regressor. Given a first-stage linear model for the endogenous regressors and a second-stage linear model for the dependent variable, we develop a novel adaptation of the parametric one-step update to a generic second-stage estimator. We provide conditions under which the scaled update is asymptotically normal. We then introduce a two-stage Lasso procedure and show that the second-stage Lasso estimator satisfies the aforementioned conditions. Using these results, we construct asymptotically valid confidence intervals for the components of the second-stage regression coefficients. We complement our asymptotic theory with simulation studies, which demonstrate the performance of our method in finite samples.

© 2019 Elsevier B.V. All rights reserved.

## 1. Introduction

### 1.1. Overview

High-dimensional estimation has been extensively studied and is now ubiquitous in the data-intensive sciences (Bühlmann and van de Geer, 2011; Giraud, 2014; Hastie et al., 2015). High-dimensional inference, on the other hand, is much less developed. In particular, although considerable progress has been made for inference in standard high-dimensional regression (Javanmard and Montanari, 2014; van de Geer et al., 2014; Zhang and Zhang, 2014; Ning and Liu, 2017), much less is known for more complex models.

In this paper, we extend the study of high-dimensional inference to the *linear instrumental variables (IV) model*. To motivate the linear IV model, we consider the ordinary linear model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u},$$

where  $\mathbf{y}$  is the vector of responses,  $\mathbf{X}$  is the matrix of regressors,  $\boldsymbol{\beta}$  is the regression vector, and  $\mathbf{u}$  is a vector of random disturbances. Standard inference for  $\boldsymbol{\beta}$  using ordinary least squares is valid only if  $E[\mathbf{u}|\mathbf{X}] = 0$ . However, this assumption

<sup>☆</sup> We are grateful to the editors and anonymous referees for their careful review and valuable suggestions. We sincerely thank Alex Belloni, Mehmet Caner, Denis Chetverikov, Yanqin Fan, Amit Gandhi, Eric Gautier, Mohamed Hebiri, Joseph Salmon, Alexandre Tsybakov, and Jon Wellner for their insightful comments and Donghui Mai for her help in setting up the cloud computing. We acknowledge support through Cloud Credits for Research by Amazon and through the University of Washington Royalty Research Fund.

\* Corresponding author at: Department of Economics, University of Washington, Seattle, WA, 98195, USA.

E-mail addresses: [ledererj@uw.edu](mailto:ledererj@uw.edu) (J. Lederer), [jingtao@uw.edu](mailto:jingtao@uw.edu) (J. Tao).

is easily violated in practice in view of selection biases, omitted variables, measurement errors, and many other challenges common to data collection. Hence, it is often more reasonable to allow for  $E[\mathbf{u}|\mathbf{X}] \neq 0$  and instead assume that  $\mathbf{X}$  can be modeled based on observable variables  $\mathbf{Z}$  that satisfy  $E[\mathbf{u}|\mathbf{Z}] = 0$ . As standard in the econometric literature, we then call the regressors in  $\mathbf{X}$  *endogenous*, because they can be correlated with  $\mathbf{u}$ , and we call the instrumental variables  $\mathbf{Z}$  *exogenous*, because  $E[\mathbf{u}|\mathbf{Z}] = 0$ .

Inference for such models in low-dimensional settings, where the number of samples is much larger than both the number of regressors in  $\mathbf{X}$  and the number of regressors in  $\mathbf{Z}$ , has been extensively studied and put to use in economic applications and beyond (Angrist and Pischke, 2009). In the era of Big Data, however, low-dimensional settings often do not apply, because one wants to allow for flexible parameter combinations, or because many variables are measured in the first place. We are thus interested in inference for double high-dimensional settings with the number of samples dominated by both the number of regressors in  $\mathbf{X}$  and the number of regressors in  $\mathbf{Z}$ .

The method we propose is based on a novel adaptation of the parametric one-step update to a generic two-stage estimation procedure. In parametric models, the one-step update  $\hat{\beta}$  to an initial estimator  $\tilde{\beta}$  is one Newton–Raphson step in the direction of a solution to the empirical analogue of the score equations. This approach is similar to those of van de Geer et al. (2014) and Javanmard and Montanari (2014), who de-bias the Lasso, and Zhang and Zhang (2014), who use a low-dimensional projection technique, to obtain asymptotic pivots for the low-dimensional components of the high-dimensional linear regression models when endogeneity is absent. The present work extends this approach to two-stage estimation when both stages are high-dimensional and the regressors of interest may each be endogenous. To adapt the one-step update to handle endogeneity of  $\mathbf{X}$ , we (i) choose the update as a step towards the solution to the empirical analogue of a valid moment condition and (ii) apply the update to a generic second-stage estimator  $\tilde{\beta}$  that depends on the predicted conditional means  $\hat{E}[\mathbf{X}|\mathbf{Z}]$ . The resultant estimator decomposes into a main term and four remainder terms, which contrasts with the single remainder term in the case of the de-biased Lasso in the ordinary linear model.

We present high-level conditions under which the updated estimator yields asymptotic pivots for the components of  $\beta$ , and we show, as an example, how these conditions may be satisfied by a two-stage Lasso estimation routine. We assume a sub-Gaussian regime throughout for the noise elements and instrumental variables in order to support flexibility of distributional assumptions. The main challenges of establishing the example are due to (i) the involved structure of the remainder terms, whose control requires a variety of concentration bounds and lead to extensive proofs and (ii) the estimation of the population precision matrix of the conditional means  $E[\mathbf{X}|\mathbf{Z}]$ , since these are not observed directly.

## 1.2. Our contributions

Our primary contribution is to develop a method with which to conduct statistical inference for the components  $\beta_j, j = 1, \dots, p_x$  of a high-dimensional regression vector  $\beta$  despite endogeneity of the respective regressors. We develop a novel adaptation of the one-step update and high-level conditions under which the updated estimator yields asymptotically Gaussian pivots for each  $\beta_j$ . Our rigorous demonstration of conditions under which such inference is possible in the doubly high-dimensional setting differentiates the present paper from similar works such as Belloni et al. (2011), who develop inferential methods for IV models with high-dimensional instruments and low-dimensional endogenous regressors, and (Zhu, 2018), who work under a doubly high-dimensional regime but focus primarily on bounding the error of a two-stage Lasso estimator such as that of Section 4.

A related contribution concerns sparse inverse covariance matrix estimation. The updated estimator  $\tilde{\beta}$  depends on an estimate of the inverse covariance matrix  $\Theta$  of the conditional means  $E[\mathbf{x}_i|\mathbf{z}_i]$ . However, we do not observe these conditional means directly, and must base our estimate of  $\Theta$  on the predictions  $\hat{E}[\mathbf{x}_i|\mathbf{z}_i]$ . For this, we use a modification of the CLIME estimator  $\hat{\Theta}$  of Cai et al. (2011). Our paper is the first one to use such an estimator in the context of instrumental variable selection, and we account for the prediction step in deriving probabilistic guarantees for the estimator's performance.

Another contribution is to show that the updated second-stage Lasso estimator studied in Section 4 satisfies the high-level conditions in Section 3 and therefore supports inference for the  $\tilde{\beta}_j$ . To show as much, we develop probabilistic bounds for the second-stage  $\ell_1$  estimation error, and we use these bounds to show asymptotic negligibility of the four remainder terms described in the previous section. We also demonstrate the feasibility of the compatibility condition in the second-stage regression, thereby justifying the practical use of the second-stage rates.

A majority of the proofs factor into deterministic and stochastic components. This also allows future analysts easily to combine the generic bounds contained in Section B of the Appendix with concentration results for specific error and design matrix distribution regimes and thereby derive the growth conditions required for good asymptotic behavior of the updated second-stage estimator under a variety of models.

## 1.3. Related work

Our work is related to the recent research on inference for high-dimensional linear instrumental variables models such as Belloni et al. (2012, 2018), Fan and Liao (2014), Gautier and Tsybakov (2014), Cheng and Liao (2015) and Neykov et al. (2015). Belloni et al. (2012) use the Lasso to obtain representations of the optimal IVs of Amemiya (1974, 1977) and Hansen (1982) for models in which the conditional mean of the response is linear in a small and fixed number

of endogenous variables and show that the second-stage estimator is  $\sqrt{n}$ -consistent. Following the line of the seminal work by Belloni et al. (2012), some recent papers propose different novel procedures to select many instruments when the number of second-stage regressors remains to be fixed or low-dimensional (for instance, Hansen and Kozbur, 2014; Caner and Fan, 2015; Fan and Zhong, 2018, among others). Cheng and Liao (2015) propose a Lasso procedure to select valid and relevant moments for the GMM estimation when the number of moments increases with sample size. However, and in contrast with the present paper, the dimensions of moments and parameters of interest are both smaller than the sample size. Compared to Belloni et al. (2012) and Cheng and Liao (2015), we allow both the number of IVs and the number of endogenous regressors to be bigger than the sample size. Gautier and Tsybakov (2014, 2018) construct robust confidence sets based on their Self-Tuned Instrumental Variables (STIV) estimator and confidence bands after bias correction (requiring a type of strong instruments condition).

Much of the present work is devoted to solving the inference problem for parameters of interest in a high-dimensional linear IV model with homoscedasticity by accounting for the prediction error when the first- and second-stage regression models are both high-dimensional. This contrasts with the methods of Fan and Liao (2014), who do not account for the need to predict the optimal instruments (Amemiya, 1974, 1977; Newey, 1990; Imbens et al., 2003). To our knowledge, such analysis under an  $\ell_1$ -regularized estimation procedure is new in the literature.

Recent work by Lin et al. (2015) and Zhu (2018) also propose estimation methods for linear IV models when the regressors of both stages are high-dimensional but do not rigorously develop asymptotic methods for inference. We notice that Neykov et al. (2015) also provide an inferential method for high-dimensional linear IV models by using a Dantzig selector. However, their method requires that the number of instruments equals the number of endogenous regressors and they do not account for the first-stage of estimation. Another related work is concurrently developed by Belloni et al. (2018), who propose regularized estimation of nuisance parameters that appear in carefully constructed empirical orthogonality conditions. Our one-step update approach can be interpreted as an iteration of the Newton–Raphson method, while Belloni et al. (2018) build on the idea of Neyman orthogonality. Even though the two papers take different approaches, both show that  $\sqrt{n}$ -consistent estimators for low-dimensional parameters can be constructed in high-dimensional IV models.

In a very recent work Caner and Kock (2018b) introduce a de-sparsified,  $\ell_1$ -penalized two-stage GMM estimator. They develop estimation error bounds and inferential procedures based on this estimator in the doubly high-dimensional setting and extend our work in that they (i) allow conditional heteroskedasticity in the second-stage noise elements, (ii) specify random components in terms of moment conditions rather than sub-Gaussianity, and (iii) do not require  $\ell_0$ -sparsity of the first-stage regression vectors. Also noteworthy is that the authors do not predict the conditional means  $\mathbf{Z}\mathbf{A}$  for the weighting scheme of the second-stage. This leads to expressions for the scale factor of the asymptotic pivots in their Theorem 2.(i) that differ those of our Theorem 3.4. The scale factors of the latter such pivots is identical to the asymptotic variance of optimal estimators in low-dimensional IV models with homoscedastic structural errors (Chamberlain, 1987; Newey, 1990).<sup>1</sup>

#### 1.4. Organization

The rest of the paper is organized as follows. We introduce our model and a generic two-stage estimation procedure in Section 2. In Section 3, we propose the one-step update inference procedure and demonstrate conditions under which the update yields asymptotically Gaussian pivots. In Section 4, we introduce a two-stage Lasso estimator of the parameter  $\beta$  and show that it is suitable for use with the inference procedure developed in Section 3. Finally, in Section 5, we present the results of numerical studies that demonstrate the relevance of our theoretical results to finite samples. All proofs are contained in Appendix.

#### 1.5. Basic notation and preliminaries

We adopt the following general notational conventions. For  $p \in \mathbb{N}$ , we let  $[p] := \{1, \dots, p\}$ . We typically use bold and non-bold lowercase letters to denote vectors and scalars, respectively. We use bold uppercase letters to denote matrices. We typically denote the components of a vector (matrix) by the non-bold (lowercase) counterpart of the letter that denotes the vector (matrix). If  $\mathbf{M} \in \mathbb{R}^{n \times p}$ , with components  $m_{ij}$ , we use a superscript to refer to columns  $\mathbf{m}^j = (m_{1j}, \dots, m_{nj})^\top$  and a subscript to refer to rows  $\mathbf{m}_i = (m_{i1}, \dots, m_{ip})$ . We let  $\|\cdot\|_q$  and  $\langle \cdot, \cdot \rangle$  denote the usual  $\ell_q$  norm and inner product over Euclidean spaces, respectively.

For  $\mathbf{m} \in \mathbb{R}^p$ , we let  $\text{supp}(\mathbf{m}) := \{j \in [p] : m_j \neq 0\}$ ,  $\|\mathbf{m}\|_0 = |\text{supp}(\mathbf{m})|$ , and  $\|\mathbf{m}\|_\infty = \max_{j \in [p]} |m_j|$ . For matrices  $\mathbf{M} \in \mathbb{R}^{n \times p}$ , we let  $\|\mathbf{M}\|_\infty = \max_{i,j \in [n] \times [p]} |m_{ij}|$ ,  $\|\mathbf{M}\|_{L_1} = \max_{j \in [p]} \|\mathbf{m}^j\|_1$ , and  $\|\mathbf{M}\|_{L_2} = \max_{j \in [p]} \|\mathbf{m}^j\|_2$ . For matrices  $\mathbf{M}_1$  and  $\mathbf{M}_2 \in \mathbb{R}^{n \times p}$ , we write  $\mathbf{M}_1 \succ \mathbf{M}_2$  if  $\mathbf{M}_1 - \mathbf{M}_2$  is positive-definite.

For quantities  $x$  indexed by  $i \in [n]$ , we let  $\mathbb{E}_n[x_i] = n^{-1} \sum_{i=1}^n x_i$ . If  $X_n$  is a sequence of random variables, we write  $X_n \rightsquigarrow X$  if  $X_n$  converges weakly to  $X$ . For  $a, b \in \mathbb{R}$ , we let  $a \vee b = \max(a, b)$  and  $a \wedge b = \min(a, b)$ . We write  $a_n \lesssim b_n$  if  $a_n \leq C_n b_n$  for a  $C_n$  that is of constant order. We say that a sequence of events  $\mathcal{E} \equiv \mathcal{E}_n$  occurs with probability approaching one if  $\lim_{n \rightarrow \infty} \mathbb{P} \mathcal{E}_n = 1$ .

<sup>1</sup> In Caner and Kock (2018a), the same authors develop asymptotic theory to support both estimation and inference for the conservative Lasso under a wide variety (including heteroscedasticity and non-sub-Gaussianity) of error regimes, thereby laying foundations for future work considering two-stage estimation.

We recall the following definitions for sub-Gaussian and sub-exponential norms.

**Definition 1.1** (*Sub-Gaussian and Sub-exponential Norms*). For  $q \geq 1$  and a random variable  $X$ , we write

$$\|X\|_{\psi_q} := \inf\{t \in (0, \infty) : E[\exp(|X|^q/t^q) - 1] \leq 1\}.$$

if the infimum exists. The *sub-Gaussian norm* of a random variable  $X$  is given by  $\|X\|_{\psi_2}$ ; the *sub-exponential norm* of a random variable  $X$  is given by  $\|X\|_{\psi_1}$ . The corresponding norms for a random  $p$ -vector  $\mathbf{X}$  are given by

$$\|\mathbf{X}\|_{\psi_q} := \sup_{\mathbf{x} \in \mathbb{R}^p : \|\mathbf{x}\|_2=1} \|\langle \mathbf{X}, \mathbf{x} \rangle\|_{\psi_q}.$$

## 2. Two-stage estimation

To contend with endogeneity, the method of instrumental variables isolates variation in the endogenous regressors induced by the instrumental variables. In Section 2.1, we posit the two-stage linear IV model to describe this relationship. In Section 2.2, we discuss a generic two-stage estimation routine that respects the structure of the model.

### 2.1. Model

Our model of interest is

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta} + u_i, \quad (1)$$

$$x_{ij} = \mathbf{z}_i^\top \boldsymbol{\alpha}^j + v_{ij}, \quad (2)$$

where:  $i$  ranges from 1 to  $n$  (unless stated otherwise);  $j$  ranges from 1 to  $p_x$  (unless stated otherwise); the vectors  $\mathbf{x}_i \in \mathbb{R}^{p_x}$  consist of the *second-stage regressors*  $x_{i1}, \dots, x_{ip_x}$ ; the vector  $\boldsymbol{\beta} \in \mathbb{R}^{p_x}$  is the parameter of interest; the vectors  $\mathbf{z}_i \in \mathbb{R}^{p_z}$  consist of the *first-stage regressors*  $z_{i1}, \dots, z_{ip_z}$ ; the quantities  $u_i$  and  $\mathbf{v}_i := (v_{i1}, \dots, v_{ip_x})^\top$  are random noise elements that satisfy

$$E[u_i | \mathbf{z}_i] = 0, \quad E[\mathbf{v}_i | \mathbf{z}_i] = \mathbf{0}, \quad (3)$$

and the vectors  $\boldsymbol{\alpha}^j$  are regression parameters up to which the respective conditional means  $d_{ij} := E[x_{ij} | \mathbf{z}_i] = \mathbf{z}_i^\top \boldsymbol{\alpha}^j$  are specified. We call the models of (2) and (1) the first-stage and second-stage models, respectively. Note that the setup in (3) is similar to the one in Belloni et al. (2012) except for this significant difference: we consider both stages to be high-dimensional while the second stage parameters in Belloni et al. (2012) are low-dimensional. To focus on presenting the main idea and inference steps, we ignore approximation errors in the model but one can add approximation errors in both (1) and (2) at the expense of more tedious derivations and notation (see the discussion of Theorem 6.3 in Bühlmann and van de Geer, 2011).

In matrix notation, we write

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}$$

and

$$\mathbf{X} = \mathbf{D} + \mathbf{V} = \mathbf{Z}\mathbf{A} + \mathbf{V},$$

where the vectors  $\mathbf{y}, \mathbf{u} \in \mathbb{R}^n$  consist of the responses  $y_i$  and the noise components  $u_i$ , respectively; the matrix  $\mathbf{X} \in \mathbb{R}^{n \times p_x}$  has columns  $\mathbf{x}^j = (x_{1j}, \dots, x_{nj})^\top$  and rows  $\mathbf{x}_i = (x_{i1}, \dots, x_{ip_x})$ ; the matrix  $\mathbf{D} = E[\mathbf{X} | \mathbf{Z}] \in \mathbb{R}^{n \times p_x}$  has columns  $\mathbf{d}^j = (d_{1j}, \dots, d_{nj})^\top$  and rows  $\mathbf{d}_i = (d_{i1}, \dots, d_{ip_x})$ ; the matrix  $\mathbf{Z} \in \mathbb{R}^{n \times p_z}$  has columns  $\mathbf{z}^k = (z_{1k}, \dots, z_{nk})^\top$  and rows  $\mathbf{z}_i = (z_{i1}, \dots, z_{ip_z})$ ; and the matrix  $\mathbf{A} \in \mathbb{R}^{p_z \times p_x}$  has columns given by  $\boldsymbol{\alpha}^j$ . We make the following assumption concerning the  $n$ -indexed sequence of regression parameters  $\mathbf{A}, \boldsymbol{\beta}$ .

**Assumption 2.1** (*Regularity of  $\mathbf{A}, \boldsymbol{\beta}$* ). The quantities  $\|\mathbf{A}\|_{L_1}$  and  $\|\boldsymbol{\beta}\|_1$  are bounded above by universal constants  $m_A, m_\beta < \infty$ , respectively.

We let  $\widehat{\Sigma}_z = \mathbf{Z}^\top \mathbf{Z} / n$  denote the empirical Gram matrix of the instrumental variables. As remarked earlier, the linear IV model has been studied extensively in the low-dimensional setting, where the number  $p_x$  of endogenous variables  $\mathbf{x}^j$  is fixed. We are particularly concerned with the high-dimensional regime in which both  $p_x$  and the number  $p_z$  of instrumental variables  $\mathbf{z}^k$  increase with  $n$ . Our results generalize to the low-dimensional case in which  $p_z$  and  $p_x$  are held fixed with respect to  $n$ , but we do not treat this case explicitly in the present essay. Regardless of whether the model is high-dimensional, we require that  $p_x \leq p_z$  in order to maintain identifiability of  $E[\mathbf{x}_i | \mathbf{z}_i]$ .

We study a sub-Gaussian regime for the noise components as well as the instrumental variables, which we treat as random throughout and for which we give marginal results. This regime encompasses the typical Gaussian model considered in the high-dimensional literature and allows for flexibility in modeling assumptions.

**Assumption 2.2** (Specification of  $\mathbf{z}_i$ ). The instrumental variables  $\mathbf{z}_i$  are i.i.d. and sub-Gaussian with sub-Gaussian norm  $\tau_{\mathbf{z}} := \|\mathbf{z}_i\|_{\psi_2}$  and satisfy  $E[\mathbf{z}_i] = 0$  for each  $i \in [n]$ . Considered as components of an  $n$ -indexed sequence of models, the quantities  $\tau_{\mathbf{z}}$  are bounded away from zero and infinity.

**Remark 2.3** (Specification of  $\mathbf{z}_i$ ). Assumption [Assumption 2.2](#) is not restrictive. We require that the first-stage regressors  $\mathbf{z}_i$  have mean zero in order to simplify the following exposition and to apply concentration results under more specific distributional assumptions, such as in [Lemma 4.10](#). This assumption can be relaxed at the expense of brevity and given a sufficient reformulation of the required concentration results. Similarly, the condition that  $\tau_{\mathbf{z}} = O(1)$  can be relaxed at the expense of introducing more complex growth conditions in later results.

**Assumption 2.4** (Specification of  $\mathbf{v}^j$  and  $\mathbf{u}$ ). The noise vectors  $\mathbf{v}^j$  and  $\mathbf{u}$  are sub-Gaussian with sub-Gaussian norms  $\tau_{\mathbf{v}^j} := \|\mathbf{v}^j\|_{\psi_2}$  and  $\tau_{\mathbf{u}} := \|\mathbf{u}\|_{\psi_2}$ . Considered as components of an  $n$ -indexed sequence of models, the quantities  $\tau_{\mathbf{v}^j}$  and  $\tau_{\mathbf{u}}$  are bounded strictly away from zero and infinity.

Note that [Assumption 2.4](#) makes no stipulations concerning the joint covariance structure of the  $u_i$  and  $v_i$ . The assumption therefore allows for nontrivial covariance between the two stages of noise, which can be used to model endogeneity of the  $\mathbf{x}_i$ . Furthermore, [Assumption 2.4](#) allows for heteroscedasticity amongst the components of the  $\mathbf{v}^j$ . We require homoscedasticity of the second-stage noise elements  $u_i$  for [Theorem 3.4](#) and [Lemma 4.11](#); all other results hold in the presence of heteroscedasticity.

## 2.2. Generic two-stage estimators

We formulate our proposed method of inference for the components  $\beta_j$  of the second-stage regression parameter  $\boldsymbol{\beta}$  in terms of generic estimators that reflect the structure of the model described above. We now introduce notation that will be used in [Section 3](#).

For each  $j \in [p_x]$ , let  $\hat{\boldsymbol{\alpha}}^j \equiv \hat{\boldsymbol{\alpha}}^j(\mathbf{x}^j, \mathbf{Z})$  denote a generic *first-stage estimator* of the first-stage regression vector  $\boldsymbol{\alpha}^j$  based on the data  $\mathbf{x}^j$  and  $\mathbf{Z}$ . We write  $\hat{\mathbf{A}} := (\hat{\boldsymbol{\alpha}}^1, \dots, \hat{\boldsymbol{\alpha}}^{p_x})$  for the matrix of estimated regression vectors. From such an estimator  $\hat{\mathbf{A}}$  we may predict the conditional means  $\mathbf{d}_i = E[\mathbf{x}_i | \mathbf{z}_i]$  for  $i \in [n]$  with  $\hat{\mathbf{d}}_i := \mathbf{z}_i^\top \hat{\mathbf{A}}$ ; we write  $\hat{\mathbf{D}}$  for the predicted conditional mean matrix whose rows are given by the  $\hat{\mathbf{d}}_i$ , and we write  $\hat{\Sigma}_{\mathbf{d}} := \hat{\mathbf{D}}^\top \hat{\mathbf{D}}/n$  and  $\Sigma_{\mathbf{d}} := E[\hat{\Sigma}_{\mathbf{d}}]$ . Our choice of the notation  $\hat{\mathbf{D}}$  reflects the fact that this quantity predicts and, under suitable conditions, approaches in probability the conditional mean matrix  $\mathbf{D}$ . We write  $\hat{\boldsymbol{\beta}} \equiv \hat{\boldsymbol{\beta}}(\mathbf{y}, \hat{\mathbf{D}})$  for a generic *second-stage estimator* of the second-stage regression parameter  $\boldsymbol{\beta}$  based on the response  $\mathbf{y}$  and the predicted conditional means  $\hat{\mathbf{D}}$ .

## 3. Main proposal

Our main contribution is to develop a method for statistical inference for the components  $\beta_j$  of the second-stage regression vector  $\boldsymbol{\beta}$ . In general, statistical inference for high-dimensional regression parameters is a difficult problem. Regularized estimators, such as the Lasso and ridge regression, are often used for the purpose of high-dimensional parameter estimation but generally do not have asymptotic distributions suitable for inference ([Knight and Fu, 2000](#); [Pötscher and Leeb, 2009](#)). In studying the model of [Section 2.1](#), we must also account for the dependence of the second-stage estimator on the first-stage estimators.

The basis for our procedure is to adapt the parametric one-step update to the two-stage estimation procedure described in [Section 2.2](#). We first briefly review the use of the one-step estimator in parametric models and its application to high-dimensional inference for the ordinary linear model. Then we adapt the one-step update to the two-stage estimation procedure described in [Section 2.2](#). [Section 3.3](#) discusses high-level conditions under which the scaled updated estimator is asymptotically normal.

### 3.1. One-step with endogeneity

In this section, we develop an adaptation of the one-step update that, under suitable high-level conditions, yields asymptotic pivots for the second-stage components  $\beta_j$  of the two-stage model described in [Section 2.1](#). We note that the present development is valid for any initial second-stage estimator  $\hat{\boldsymbol{\beta}}$ . To demonstrate that the high-level conditions are satisfied requires consideration of particular estimators.

In summary, the one-step update is a general method for constructing efficient estimators for parametric and semiparametric models. Recall that the Newton–Raphson method for finding the root in  $\mathbf{b}$  to a *target system* of  $p_x$  equations

$$\mathbf{h}(\mathbf{y}_i, \mathbf{x}_i; \mathbf{b}) \equiv (h_1(\mathbf{y}_i, \mathbf{x}_i; \mathbf{b}), \dots, h_{p_x}(\mathbf{y}_i, \mathbf{x}_i; \mathbf{b}))^\top = \mathbf{0}$$

is to update an approximation  $\mathbf{b}^k$  by the rule

$$\mathbf{b}^{k+1} = \mathbf{b}^k - \left[ \frac{\partial \mathbf{h}}{\partial \mathbf{b}} \bigg|_{\mathbf{b}=\mathbf{b}^k} \right]^{-1} \mathbf{h}(\mathbf{y}_i, \mathbf{x}_i; \mathbf{b}^k),$$



where  $\frac{\partial \mathbf{h}}{\partial \mathbf{b}}|_{\mathbf{b}=\mathbf{b}^k}$  is the Jacobian matrix of  $\mathbf{h}$  with respect to  $\mathbf{b}$  evaluated at  $\mathbf{b}^k$ . In the ordinary least squares regression model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{u}, \quad (4)$$

the score function  $\mathbf{h}(\mathbf{y}_i, \mathbf{x}_i; \mathbf{b}) = -\mathbf{x}_i(\mathbf{y}_i - \mathbf{x}_i^\top \mathbf{b})$  satisfies  $E[\mathbf{h}(\mathbf{y}_i, \mathbf{x}_i; \boldsymbol{\beta})] = \mathbf{0}$  given the *orthogonality condition*

$$E[\mathbf{x}_i u_i] = \mathbf{0}. \quad (5)$$

The *one-step update*  $\tilde{\boldsymbol{\beta}}$  to an initial estimate  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  is given by

$$\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\Theta}} \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})/n,$$

where  $\hat{\boldsymbol{\Theta}}$  denotes the inverse of

$$\left. \frac{\partial (-\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{b})/n)}{\partial \mathbf{b}} \right|_{\mathbf{b}=\hat{\boldsymbol{\beta}}} = \mathbf{X}^\top \mathbf{X}/n = \hat{\boldsymbol{\Sigma}}_{\mathbf{x}}.$$

The case when the model is high-dimensional is less well studied. When  $p_{\mathbf{x}} > n$ , the empirical covariance matrix  $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}}$  is not invertible, and we have instead

$$\tilde{\boldsymbol{\beta}} = \boldsymbol{\beta} + \hat{\boldsymbol{\Theta}} \mathbf{X}^\top \mathbf{u}/n + \underbrace{(\hat{\boldsymbol{\Theta}} \hat{\boldsymbol{\Sigma}}_{\mathbf{x}} - \mathbf{I})(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})}_{\mathbf{f}/\sqrt{n}},$$

where  $\hat{\boldsymbol{\Theta}}$  denotes an approximate inverse of the Jacobian matrix. The latter term  $\mathbf{f}/\sqrt{n}$  in the above display is the “remainder” after incomplete inversion of  $\hat{\boldsymbol{\Sigma}}_{\mathbf{x}}$ . Thus, in the high-dimensional one-stage linear model, the one-step update satisfies

$$\sqrt{n}(\tilde{\beta}_j - \beta_j) = \frac{1}{\sqrt{n}} \sum_{i=1}^n \hat{\theta}_j^\top \mathbf{x}_i u_i + f_j, \quad (6)$$

where  $\hat{\theta}_j$  is the  $j$ th row of  $\hat{\boldsymbol{\Theta}}$ . The structure of the main term on the right-hand side above suggests to use  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\hat{\omega}_j$ , where  $\hat{\omega}_j$  is an appropriate estimate of  $\omega_j = (E[(\hat{\theta}_j, \mathbf{x}_i)^2 u_i^2])^{1/2}$ , as an asymptotic pivot for  $\beta_j$ .

When the initial estimator  $\hat{\boldsymbol{\beta}}$  is the Lasso, the updated estimator  $\tilde{\boldsymbol{\beta}}$  is sometimes called the desparsified (van de Geer et al., 2014) or de-biased (Javanmard and Montanari, 2014) Lasso, though these authors obtain the form of  $\tilde{\boldsymbol{\beta}}$  by means other than the one-step update. The general upshot of their results is that if  $\|\mathbf{f}\|_\infty = o_p(1)$ , and if  $\hat{\theta}_j$  and  $\mathbf{x}_i$  are independent of  $u_i$ , then the updated Lasso estimator yields asymptotically Gaussian pivots for the parameter components.

In contrast to the ordinary linear regression model, the challenge we face in the case of high-dimensional IV model is that the condition in (5) does not hold. Instead, the conditional moment restriction  $E[u_i | \mathbf{z}_i] = 0$  in (3) entails the orthogonality condition  $E[\mathbf{d}_i u_i] = \mathbf{0}$  for the conditional means  $\mathbf{d}_i = E[\mathbf{x}_i | \mathbf{z}_i]$ . This suggests that, to develop a one-step update for a generic second-stage estimator  $\hat{\boldsymbol{\beta}}$  of  $\boldsymbol{\beta}$  of the present model, we ought to take the empirical analogue

$$\mathbb{E}_n[\tilde{\mathbf{h}}(\mathbf{y}_i, \mathbf{x}_i, \hat{\mathbf{d}}_i; \mathbf{b})] := \mathbb{E}_n[-\hat{\mathbf{d}}_i(\mathbf{y}_i - \mathbf{x}_i^\top \mathbf{b})] = -\hat{\mathbf{D}}^\top (\mathbf{y} - \mathbf{X}\mathbf{b})/n = \mathbf{0},$$

of  $E[\mathbf{d}_i u_i] = \mathbf{0}$  as the target system for which the root is sought via a Newton–Raphson update. We have elected to base the target system on the moment condition  $E[\mathbf{d}_i u_i] = \mathbf{0}$  in accordance with optimal weighting regimes for generalized method of moments (GMM) estimators; see Amemiya (1974, 1977), Hansen (1982) and Newey (1990). Further, since the  $\mathbf{d}_i$  are generally unavailable, we instead use the predicted conditional mean matrix  $\hat{\mathbf{D}}$  in the target system above. The one-step update  $\tilde{\boldsymbol{\beta}}$  to a second-stage estimator  $\hat{\boldsymbol{\beta}}$  is then given by

$$\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - \hat{\boldsymbol{\Theta}} \mathbb{E}_n[\tilde{\mathbf{h}}(\mathbf{y}_i, \mathbf{x}_i, \hat{\mathbf{d}}_i; \mathbf{b})] = \hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\Theta}} \hat{\mathbf{D}}^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})/n, \quad (7)$$

where we continue to let  $\hat{\boldsymbol{\Theta}}$  denote an (approximate) inverse to the Jacobian matrix in  $\mathbf{b}$  of the score  $\tilde{\mathbf{h}}(\mathbf{y}_i, \mathbf{x}_i, \hat{\mathbf{d}}_i; \mathbf{b})$ .

If one were to follow strictly the prescription of the Newton–Raphson method for selection of  $\hat{\boldsymbol{\Theta}}$  for the updated second-stage estimator  $\tilde{\boldsymbol{\beta}}$ , one would select  $\hat{\boldsymbol{\Theta}} \approx [\hat{\mathbf{D}}^\top \mathbf{X}/n]^{-1}$  to approximate the inverse of the Jacobian of  $\tilde{\mathbf{h}}$  evaluated at  $\hat{\boldsymbol{\beta}}$ . However, the decomposition obtained in the following lemma suggests that  $\hat{\boldsymbol{\Theta}}$  ought to control, say, the sup-norm of  $\hat{\boldsymbol{\Theta}} \hat{\boldsymbol{\Sigma}}_{\mathbf{d}} - \mathbf{I}$ , and hence aim to invert  $\hat{\boldsymbol{\Sigma}}_{\mathbf{d}} := \hat{\mathbf{D}}^\top \hat{\mathbf{D}}/n$  rather than  $\hat{\mathbf{D}}^\top \mathbf{X}/n$ . We emphasize that the one-step formulation, insofar as it follows the Newton–Raphson method, is merely a vehicle for producing an updated estimator  $\tilde{\boldsymbol{\beta}}$ . In particular, Lemma 3.1 is valid regardless of what convergence properties an actual Newton–Raphson algorithm incorporating a specific choice of  $\hat{\boldsymbol{\Theta}}$  may exhibit. We may choose  $\hat{\boldsymbol{\Theta}}$  in whatever manner is most appropriate for achieving our goal, which is to obtain a tractable limiting distribution for  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j$ , where  $\omega_j$  is an appropriate scale factor. That said, the two suggestions for how to choose  $\hat{\boldsymbol{\Theta}}$  may be reconciled somewhat by noting that both  $\hat{\mathbf{D}}^\top \mathbf{Z}/n$  and  $\hat{\mathbf{D}}^\top \hat{\mathbf{D}}/n$  are equal to the empirical Gram matrix  $\hat{\boldsymbol{\Sigma}}_{\mathbf{d}}$  modulo additional terms whose sup-norms can be controlled given a rate for  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$  and appropriate concentration results for  $\|\mathbf{Z}^\top \mathbf{v}/n\|_\infty$ . In turn, one finds  $\|\hat{\boldsymbol{\Sigma}}_{\mathbf{d}} - \boldsymbol{\Sigma}_{\mathbf{d}}\|_\infty = o_p(1)$  under appropriate growth restrictions on  $p_{\mathbf{x}}$ ; see Lemma C.5.

For our purposes, we consider the matrix  $\hat{\Theta}$  primarily as an estimator of the population quantity  $\Theta := E[\mathbf{d}_i \mathbf{d}_i^\top]^{-1}$ . In particular, we require good behavior of  $\hat{\Theta}$  as such an estimator to derive the asymptotic distribution of  $\sqrt{n}(\hat{\beta}_j - \beta_j)$ .

The following lemma characterizes a similar decomposition of the updated estimator  $\tilde{\beta}$  as in the one-stage model.

**Lemma 3.1** (Decomposition of One-step Second-stage Estimator). *Consider the two-stage linear model described in Section 2.1. Let  $\hat{\mathbf{D}}$  be a prediction of the conditional mean matrix  $\mathbf{D}$  from an estimate  $\hat{\mathbf{A}}$  of the first-stage regression matrix  $\mathbf{A}$ . Let  $\hat{\beta}$  be a second-stage estimator based on the predictions  $\hat{\mathbf{D}}$ . Let  $\hat{\Theta}$  denote an estimator of  $\Theta = E[\mathbf{d}_i \mathbf{d}_i^\top]^{-1}$ . The one-step second-stage estimator*

$$\tilde{\beta} = \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{y} - \mathbf{X} \hat{\beta}) / n$$

satisfies  $\sqrt{n}(\tilde{\beta} - \beta) = \Theta \mathbf{D}^\top \mathbf{u} / \sqrt{n} + \sum_{\ell=1}^4 \mathbf{f}_\ell$ , where

$$\begin{aligned} \mathbf{f}_1 &= (\hat{\Theta} - \Theta) \mathbf{D}^\top \mathbf{u} / \sqrt{n}, & \mathbf{f}_2 &= \hat{\Theta} (\hat{\mathbf{D}} - \mathbf{D})^\top \mathbf{u} / \sqrt{n}, \\ \mathbf{f}_3 &= \hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})(\beta - \hat{\beta}) / \sqrt{n}, & \mathbf{f}_4 &= \sqrt{n}(\hat{\Theta} \hat{\Sigma}_{\mathbf{d}} - \mathbf{I})(\beta - \hat{\beta}). \end{aligned}$$

As in the case of the main term for the ordinary one-step update discussed in (6), this observation is similar to that of Neykov et al. (2015), who derive a similar asymptotic linearization but do not account for prediction of the conditional means  $\mathbf{d}_i$ . Indeed, due to the two stages of estimation, the update incurs four remainder terms as opposed to the single term in Javanmard and Montanari (2014) and van de Geer et al. (2014). In Section 3.3, we show that the quantity  $\sqrt{n}(\hat{\beta}_j - \beta_j) / \omega_j$ , where  $\omega_j^2 := E[(\theta_j, \mathbf{d}_i)^2 u_i^2]$ , converges weakly to a  $\mathcal{N}(0, 1)$  random variable under high level conditions on the remainder terms  $\mathbf{f}_{\ell,j}$  and that the limit continues to hold if  $\omega_j$  is replaced by an appropriate estimator. From this result one may construct asymptotically valid confidence intervals for the regression components  $\beta_j$ .

We have described a strategy for inference for the components  $\tilde{\beta}_j$ . To implement the strategy for a specific choice of first- and second-stage estimators, one must identify the conditions under which the remainder terms  $\mathbf{f}_\ell$  vanish in probability. We demonstrate such an implementation in Section 4. The conditions in turn depend on the properties of the estimator  $\hat{\Theta}$ . In the following section, we introduce an estimator suitable for our purposes.

### 3.2. Estimating $\Theta$

The one-step second-stage estimator  $\tilde{\beta}$  depends on an estimator  $\hat{\Theta}$  of  $\Theta = \Sigma_{\mathbf{d}}^{-1}$ , the population precision matrix of the conditional means  $\mathbf{d}_i$ . In general, estimating the population precision matrix incurs two main difficulties in the high-dimensional setting. First, the empirical covariance matrix  $\hat{\Sigma}_{\mathbf{d}}$  is singular when  $p_{\mathbf{x}} > n$  and cannot be inverted to produce an estimator of  $\Theta$ . Second, even if an inverse were available, one cannot naively use the continuous mapping theorem to derive asymptotic guarantees if  $p_{\mathbf{x}} \rightarrow \infty$ , since the sequence of population covariance matrices  $\Sigma_{\mathbf{d}} \equiv \Sigma_{\mathbf{d},n}$  does not itself have a limit if  $p_{\mathbf{x}} \rightarrow \infty$ . In addition to these general difficulties, we must further contend with the fact that the conditional mean matrix  $\mathbf{D}$  is unknown. Hence any estimator of  $\hat{\Theta}$  will depend on the prediction  $\hat{\mathbf{D}}$ , and guarantees for such an estimator must account for such dependence.

We use a slight modification of the CLIME estimator of Cai et al. (2011) to contend with the challenges described above. The rows  $\hat{\theta}_j$  of the estimator  $\hat{\Theta}$  are obtained as solutions to the CLIME program codified below.

**Program 3.2** (Program for  $\hat{\theta}_j$ ).

$$\underset{\theta \in \mathbb{R}^{p_{\mathbf{x}}}}{\text{minimize:}} \quad Q(\theta) := \|\theta\|_1, \quad \text{subject to:} \quad \|\hat{\Sigma}_{\mathbf{d}} \theta - \mathbf{e}_j\|_\infty \leq \mu,$$

where  $\mathbf{e}_j$  denotes the  $j$ th canonical basis vector in  $p_{\mathbf{x}}$  dimensions and  $\mu > 0$  is a controlled tolerance.

The present estimator  $\hat{\Theta}$  differs in only one respect from that of the CLIME estimator of Cai et al. (2011). The latter authors symmetrize the matrix  $\Theta$  with rows obtained as solutions to the aforementioned optimization problem, whereas we use the raw solutions. We omit the symmetrization step for simplicity; the  $\ell_\infty$  and  $\ell_1$  guarantee that Cai et al. (2011) obtain for the estimation error of the CLIME estimator continue to hold. We include the requisite guarantees for the unsymmetrized estimator in Section A.2 of the Supplementary Materials.

The present estimator  $\hat{\Theta}$  also differs in an important respect from that of Javanmard and Montanari (2014). The latter also obtain an inverse Gram matrix approximation as a solution to a convex program with identical constraints as in Problem 3.2 but with objective function  $Q(\theta) = \mathbb{E}_n[(\theta, \mathbf{x}_i)^2]$ . To our knowledge, however, it is currently unknown whether the choice of  $Q$  in Javanmard and Montanari (2014) yields guarantees comparable to those of the CLIME estimator.

The  $\ell_1$  bound for  $\hat{\theta}_j - \theta_j$ , which we require for control of the remainder terms  $\mathbf{f}_\ell$ , depends on the following restriction on the class of population precision matrices  $\Theta$ .

**Definition 3.3** (Uniformity Class). Following Cai et al. (2011), we define the *uniformity class* of population precision matrices  $\Theta = \Sigma_{\mathbf{d}}^{-1}$  relative to the controlled tolerance  $q \in [0, 1]$  and the generalized sparsity level  $s_\Theta > 0$  by

$$\mathcal{U}(m_\Theta, q, s_\Theta) := \left\{ \Theta = (\theta_{jk})_{j,k=1}^{p_{\mathbf{x}}} > \mathbf{0} : \|\Theta\|_{L_1} \leq m_\Theta; \max_{j \in [p_{\mathbf{x}}]} \sum_{k \in [p_{\mathbf{x}}]} |\theta_{jk}|^q \leq s_\Theta \right\}.$$

In the sequel, we assume as part of high-level regularity conditions that  $\Theta \in \mathcal{U}(m_\Theta, q, s_\Theta)$  and that the model parameters  $m_\Theta$  and  $s_\Theta$  are well-behaved as functions of  $n$ . These parameters appear in the rates for the remainder terms in our analysis of the two-stage Lasso of Section 4. For high-level results, we also assume that the probability that the rows  $\theta_j$  of the population precision matrix are feasible for Problem 3.2 approaches one. To express this requirement formally, we define the event

$$\mathcal{T}_\Theta(\mu) := \{\|\Theta\widehat{\Sigma}_d - I\|_\infty \leq \mu\} \quad (8)$$

where  $\mu > 0$  is the tolerance of Problem 3.2, and require that  $P\mathcal{T}_\Theta(\mu) \rightarrow 1$  as  $n \rightarrow \infty$ . We identify a theoretical choice of  $\mu$  that satisfies the latter requirement in Lemma 4.10; we discuss a practical method for selecting  $\mu$  in Section 5. Note that, given the event  $\mathcal{T}_\Theta(\mu)$ , the rows  $\theta_j$  of  $\Theta_j$  are each feasible for the respective Problem 3.2.

Since the quantity  $\mu$  appears in the rates for the remainder terms  $\mathbf{f}_\ell$ , it must be chosen carefully so as to balance the growth of  $P\mathcal{T}_\Theta(\mu)$  with the decay of the  $\|\mathbf{f}_\ell\|_\infty$ . The appropriate choice of  $\mu$  depends on both the distribution of the  $\mathbf{z}_i$  as well as the rate for  $\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1}$  – see Lemma 4.10.

### 3.3. Asymptotic normality

We saw in Section 3.1 that the updated estimator  $\tilde{\beta}$  satisfies  $\sqrt{n}(\tilde{\beta}_j - \beta_j) = \sqrt{n}\mathbb{E}_n[\langle \theta_j, \mathbf{d}_i \rangle u_i] + \sum_{\ell=1}^4 f_{\ell,j}$ . If the remainder terms vanish in probability, then  $\sqrt{n}(\tilde{\beta}_j - \beta_j)$  shares the same weak limit, if it exists, as  $\sqrt{n}\mathbb{E}_n[\langle \theta_j, \mathbf{d}_i \rangle u_i]$ . If the model were fixed in  $n$ , the Central Limit Theorem would entail that the latter quantity converges weakly to a  $\mathcal{N}(0, \omega_j^2)$ , where  $\omega_j^2 = E[\langle \theta_j, \mathbf{d}_i \rangle^2 u_i^2]$ . In Theorem 3.4, we provide conditions under which  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j$  converges weakly to a standard Normal random variable when the model is not fixed in  $n$ . We also show that the limit continues to hold if  $\omega_j$  is replaced by an estimator  $\hat{\omega}_j$  that satisfies  $|\hat{\omega}_j - \omega_j| = o_p(1)$ . Note that Theorem 3.4 gives conditions under which the limit holds given homoscedastic Gaussian noise (Condition (5)) as well as conditions under which the limit holds given generic i.i.d. noise (Condition (6)). Condition (5) is unnecessarily restrictive in practice. Nonetheless, we include the result under Condition (5) because it serves as a benchmark and requires subsequently weaker assumptions concerning limiting rates. The latter result includes the case of Assumption 2.4, as well as any other i.i.d. second-stage noise regime.

**Theorem 3.4 (Weak Limits).** Suppose that

- (1) the quantity  $\|\bar{\Sigma}_d - \Sigma_d\|_\infty$  vanishes in probability, where  $\bar{\Sigma}_d := \mathbf{D}^\top \mathbf{D}/n$ ;
- (2) the remainder terms satisfy  $\|\mathbf{f}_\ell\|_\infty = o_p(1)$  for each  $1 \leq \ell \leq 4$ ;
- (3)  $\Theta_{jj} > \vartheta$  for some universal constant  $\vartheta > 0$  and each  $j \in [p_x]$ ;
- (4)  $\max_{j \in [p_x]} \|\theta_j\|_1 \leq m_\Theta$  for some universal constant  $m_\Theta < \infty$ .

If either

- (5) the noise elements  $u_i$  satisfy  $u_i \mid \mathbf{z}_i \sim_{i.i.d.} \mathcal{N}(0, \sigma_u)$ , where  $\sigma_u$  is bounded away from zero and infinity uniformly in  $n$ , or
- (6) the  $\mathbf{z}_i$  and  $u_i$  are i.i.d. with  $E[u_i^2 \mid \mathbf{z}_i] = \sigma_u^2$ , where  $\sigma_u$  is bounded away from zero and infinity uniformly in  $n$ , and there exist  $0 < \zeta < 1/2$  and  $\nu > 0$  such that
  - (a)  $P\{|\langle \theta_j, \mathbf{d}_i \rangle| > n^\zeta\} = o(1)$  and
  - (b)  $E[|u_i|^{2+\nu}] \lesssim \sigma_u^{2+\nu}$ ,

then

$$\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j \rightsquigarrow Z_j \sim \mathcal{N}(0, 1).$$

Furthermore, the limit continues to hold if  $\omega_j$  is replaced by an estimator  $\hat{\omega}_j$  that satisfies  $|\hat{\omega}_j - \omega_j| = o_p(1)$ .

The proof of Theorem 3.4 can be found in Section A.3 of the Supplementary Materials. We note that, unlike the setting of Bickel et al. (1998, Sections 2.5), we do not require  $\sqrt{n}$ -consistency of the initial estimator  $\tilde{\beta}$  in order for the updates  $\tilde{\beta}_j$  to be asymptotically Gaussian. Indeed, the rates of convergence required are dictated by the strategies used to bound the quantities  $\|\mathbf{f}_\ell\|_\infty$ . See Section 4.3 for an example of sufficient rates for the two-stage Lasso routine studied in Section 4.

The main application of Theorem 3.4 is the construction of asymptotically valid confidence intervals under a wide variety of noise regimes. Given  $j \in [p_x]$  and a confidence level  $\tau$ , an asymptotic  $100(1 - \tau)\%$  confidence interval  $\hat{\mathcal{I}}_{\tau,j}$  is given by

$$\hat{\mathcal{I}}_{\tau,j} := [\tilde{\beta}_j - z_\tau \hat{\omega}_j, \tilde{\beta}_j + z_\tau \hat{\omega}_j], \quad (9)$$

where  $z_\tau = \Phi^{-1}(1 - \tau/2)$  and  $\hat{\omega}_j$  satisfies the conditions of Theorem 3.4. We present a simulation study of the finite-sample properties of this procedure for the updated two-stage Lasso estimator in Section 5.

Similarly as for the de-biasing approach in linear regression (van de Geer et al., 2014, Corollary 2.1), the results in Theorem 3.4 hold uniformly over a class of parameters as long as the assumptions also hold uniformly over that class.



This requires in particular that the remainder terms are  $o_p(1)$  uniformly. In the Lasso case considered in Section 4, such uniformity is known for parameter sets with bounded sparsity – we refer again to van de Geer et al. (2014) for details.

Note of the scale factor  $\omega_j$  that

$$\omega_j^2 = E[(\theta_j, \mathbf{d}_i)^T \mathbf{u}_i^2] = \sigma_u^2 \theta_j^T E[\mathbf{d}_i \mathbf{d}_i^T] \theta_j = \sigma_u^2 \theta_j^T \Sigma_d \theta_j = \sigma_u^2 \Theta_{jj},$$

which, if the second-stage model were fixed in  $n$ , would equal the asymptotic variance of the optimal linear IV estimator (Amemiya, 1974; Chamberlain, 1987; Newey, 1990). Indeed, Belloni et al. (2012) show that, in such a case, the linear IV estimator still attains the semi-parametric efficiency bound if the first-stage mean model is unknown but well-approximated with a high-dimensional linear model. While the display above suggests a similar optimality result for our estimator as well, such a property is difficult to state formally in the present framework. The aforementioned authors study target model parameters assumed to be identical in  $n$  or that converge to some fixed limit. In the high-dimensional setting of the present paper, the number of model parameters is allowed to grow infinitely and thus cannot converge to a fixed limit within the parameter space for any given model. (Note that this is not the same issue as contending with an infinite-dimensional nuisance parameter.) Thus it is impossible to compare asymptotic covariance matrices directly. Indeed, we do not prove any results concerning the asymptotic variance of our estimator but rather that of the derived t-statistic in Theorem 3.4.

One might be tempted to impose regularity conditions that would allow for such comparisons. For instance, one might require that, for any  $p_0 \in \mathbb{N}$ , the  $n$ -indexed sequence of parameter subsets  $\beta_{1,n}, \dots, \beta_{p_0,n}$  converge to a well-defined  $p_0$ -dimensional vector. Then one might conjecture along the lines of: “For each such  $p_0$ , the respective  $p_0$ -length subset of  $\beta$  has asymptotic variance identical to the optimal variance in a fixed  $p_0$ -dimensional linear IV model”. However, the scope of such asymptotic statements is relative to a fixed  $p_0$ ; the approach does not describe the asymptotic behavior of  $\beta$  under joint growth of  $n$  and  $p_x$  when the latter may depend on the former, which is the setting we study. Hence we do not pursue this approach here. Instead, we believe that a proper study of optimal inference in high-dimensional linear IV models should follow along the lines of Jankova and van de Geer (2018). We conjecture that the scale factor in the asymptotic pivot of Theorem 3.4 is optimal in the sense of Jankova and van de Geer (2018, Theorem 3) and that the updated second-stage estimator achieves a similar efficiency bound. A more thorough investigation of this matter is required for future work.

We conclude the present section with a brief discussion of the feasibility of select conditions of Theorem 3.4. Condition (3) can be derived as a consequence of the standard assumption that the minimal and maximal eigenvalues of  $\Sigma_d$  be bounded strictly away from zero and infinity; see Proposition C.4. The feasibility of Conditions (1) and (6) depends on the distribution of the conditional means  $\mathbf{d}_i$  and hence of the instrumental variables  $\mathbf{z}_i$ . The following Lemma shows that both conditions are satisfied if the  $\mathbf{z}_i$  are sub-Gaussian.

**Lemma 3.5** (Feasibility of Conditions (1) and (6)a). Suppose that (i) the instrumental variables  $\mathbf{z}_i$  satisfy Assumption 2.2, (ii) we have  $\|\Sigma_z\|_\infty = O(1)$ , (iii) we have  $\|\Theta\|_{L_1} \leq m_\Theta$  for  $m_\Theta = O(1)$ , (iv)  $\sqrt{(\log p_x)/n} = o(1)$ , and (v)  $\exp(-n^{2\zeta}/\sqrt{\log p_x}) = o(1)$  for some  $0 < \zeta < 1/2$ . Then Conditions (1) and (6)a of Theorem 3.4 are satisfied.

The proof of Lemma 3.5 is found in Section A.3 of Appendix. The requirements that the quantities  $\tau$ ,  $\|\Sigma_z\|_\infty$ , and  $m_\Theta$  in the statement of Lemma 3.5 be of constant order can be relaxed at the cost of introducing more complex growth conditions. Condition (v) of Lemma 3.5 is satisfied under reasonable constraints on the growth of  $p_x$  – for instance, if  $\sqrt{\log p_x}/n^{2\zeta} \lesssim 1$ .

## 4. Two-stage Lasso

Theorem 3.4 depends on high-level assumptions that ensure good behavior of the remainder terms  $\mathbf{f}_\ell$  and standard error estimate  $\hat{\omega}_j$ . In this section, we demonstrate how such conditions may be satisfied in the high-dimensional setting. In particular, we introduce in Section 4.1 a two-stage Lasso estimation procedure, for which we provide theoretical bounds in Section 4.2. The rates for the second-stage estimation error are particularly involved due to the dependence on the predicted conditional means from the first-stage estimation. In Section 4.3, we identify conditions under which the remainder terms  $\mathbf{f}_\ell$  vanish in probability under the two-stage Lasso procedure. We also propose estimators of the standard errors and provide conditions under which these estimators converge to the true standard errors.

### 4.1. Two-stage estimator

For  $j \in [p_x]$ , we let  $\hat{\alpha}^j$  denote the first-stage Lasso estimator

$$\hat{\alpha}^j \in \arg \min_{\mathbf{a} \in \mathbb{R}^{p_z}} \{ \|\mathbf{x}^j - \mathbf{Z}\mathbf{a}\|_2^2 / (2n) + r_j \|\mathbf{a}\|_1 \}. \quad (10)$$

We let  $\hat{\mathbf{d}}_{ij} := \mathbf{z}_i^T \hat{\alpha}^j$  denote the predicted conditional mean of  $x_{ij}$  given  $\mathbf{z}_i$  based on the estimates  $\hat{\alpha}^j$  and write  $\hat{\mathbf{D}} = \mathbf{Z}\hat{\mathbf{A}}$ , where the matrix  $\hat{\mathbf{D}}$  has columns given by  $\hat{\mathbf{d}}^j := (\hat{d}_{1j}, \dots, \hat{d}_{nj})^T$  and the matrix  $\hat{\mathbf{A}} \in \mathbb{R}^{p_z \times p_x}$  has columns given by the  $\hat{\alpha}^j$ .

We define the second-stage Lasso estimator to be

$$\hat{\beta} \in \arg \min_{\mathbf{b} \in \mathbb{R}^{p_x}} \{ \|\mathbf{y} - \hat{\mathbf{D}}\mathbf{b}\|_2^2 / (2n) + r_\beta \|\mathbf{b}\|_1 \}. \quad (11)$$

In Section 4.2.2, we develop sparsity-based results that require the following quantities. We write  $S_j = \text{supp } \alpha^j$  for the active sets of the first-stage regression parameters, and we write  $s_{\alpha^j} := |S_j|$  and  $s_{\mathbf{A}} := \max_{j \in [p_x]} s_{\alpha^j}$ ; we write  $S_{\beta} = \text{supp } \beta$  for the active set of the second-stage regression parameter and  $s_{\beta} := |S_{\beta}|$ . We note that  $\ell_0$  sparsity is not a limitation in principle and that more general regression vectors may be considered at the price of additional complexity (Bühlmann and van de Geer, 2011, Sections 6.2.3–4), (Belloni et al., 2012).

**Remark 4.1 (Model Identifiability).** Let  $\alpha_{S_j}^j := \{\alpha_k^j : \alpha_k^j \neq 0\}$  denote the restriction of  $\alpha^j$  to  $S_j$ , and let  $S = \cup_j S_j$ . Write  $\mathbf{A}_S$  for the  $|S| \times p_x$  matrix with columns  $\alpha_{S_j}^j$  and  $\mathbf{Z}_S$  for the  $n \times |S|$  matrix with columns (re-indexed as necessary) corresponding to the columns of  $\mathbf{A}_S$ . Note that  $\mathbf{D} = \mathbf{Z}_S^T \mathbf{A}_S$  and hence that  $\Sigma_{\mathbf{d}} = \text{cov}(\mathbf{D}) = \mathbf{A}_S^T \text{cov}(\mathbf{Z}_S) \mathbf{A}_S$ . We require that  $\Sigma_{\mathbf{d}}$  be invertible and hence of full rank – for model identifiability, probabilistic guarantees we discuss in the following sections, and because we work directly with the inverse covariance matrix  $\Theta$ . Thus we must have  $p_x = \text{rank } \Sigma_{\mathbf{d}} \leq \min(\text{rank } \mathbf{A}_S, \text{rank } \text{cov}(\mathbf{Z}_S)) \leq |S|$ . This is a relatively strong assumption, which we implicitly require throughout the sequel. However, we emphasize that this requirement is an artifact of the  $\ell_0$ -sparsity-based methods by which we derive the bounds for the first- and second-stage estimation errors. As noted above, we chose these methods to simplify our exposition and that one can obtain morally similar but more complex bounds even when the  $\alpha^j$  are not sparse (Bühlmann and van de Geer, 2011, Sections 6.2.3–4). Thus, in general, this restriction is not impractical.

## 4.2. Estimation error bounds

In this section, we present estimation error bounds for the first- and second-stage estimators described in Section 4.1. Both such bounds depend on the same fundamental strategy for proving finite-sample guarantees for  $\ell_1$ -regularized estimators. This strategy consists of two parts. The first part is the *oracle inequality*, which establishes a deterministic bound for the estimation and prediction performance of the Lasso on a particular set of interest. The second part is the control of the *empirical process term*, which defines the set of interest. We include such prerequisites in Section B.1 of Appendix.

Before we present the estimation error bounds for the first- and second-stage Lasso estimators, we first discuss the compatibility condition, which is required in the proof of the oracle inequality.

### 4.2.1. Compatibility condition

The oracle bounds rely on the good behavior of certain moduli of continuity of the empirical Gram matrices  $\widehat{\Sigma}_{\mathbf{z}} = \mathbf{Z}^T \mathbf{Z}/n$  and  $\widehat{\Sigma}_{\mathbf{d}} = \widehat{\mathbf{D}}^T \widehat{\mathbf{D}}/n$ . We codify this requirement in the following definition.

**Definition 4.2 (Compatibility Condition).** For a given index set  $S \subseteq [p]$ ,  $p \in \mathbb{N}$ , define the double-cone

$$\mathcal{C}(S) := \{\delta \in \mathbb{R}^p \setminus \mathbf{0} : \|\delta_{S^c}\|_1 \leq 3\|\delta_S\|_1\}. \quad (12)$$

We say that the *compatibility condition* holds for the matrix  $\mathbf{M} \in \mathbb{R}^{n \times p}$  relative to the index set  $S$  and the constant  $\phi^2 > 0$  defined as

$$\phi^2 = \inf_{\delta \in \mathcal{C}(S)} \frac{|S| \|\mathbf{M}\delta\|_2^2}{n \|\delta_S\|_1^2} \quad (13)$$

if the latter is greater than zero. We call the quantity  $\phi^2$  the *compatibility constant*.

The compatibility condition is so named because it interfaces between the  $\ell_1$  norm of the estimation error and the  $\ell_2$  prediction error of the Lasso estimator. It is a standard tool in the  $\ell_1$ -regularized estimation literature to ensure identifiability: it limits the correlations among the predictors such that the estimator can discriminate between the “relevant” parameters with index in  $S$  and the remaining parameters. For this purpose, the index set  $S$  is taken to be the active set of the target regression parameter (Bühlmann and van de Geer, 2011, Chapter 6).

The set  $\mathcal{C}(S)$  increases with increasing active set  $S$ ; hence, the larger  $S$ , the more restrictive the compatibility condition becomes. This means that in such theories, identifiability and sparsity are closely intertwined.

The constant 3 is arbitrary. Alternative choices require adjustment of other constants that appear in the bounds (Bickel et al., 2009). A related, slightly stronger condition known as the *restricted eigenvalue condition* is elsewhere used for the same end (Bickel et al., 2009; Bühlmann and van de Geer, 2011; van de Geer and Bühlmann, 2009). The only known task where such conditions can be avoided is prediction (Dalalyan et al., 2017; Hebiri and Lederer, 2013; Lederer et al., 2019; van de Geer and Lederer, 2013; Zhuang and Lederer, 2018).

The compatibility and restricted eigenvalue conditions are sometimes defined more generally in terms of the cardinality  $s$  of the index set  $S$  rather than a specific index set. For instance, (Bickel et al., 2009, Assumption RE( $s, c_0$ )) require for their restricted eigenvalue condition that the quantity

$$\kappa(s, c_0) := \min_{S \subseteq [p] : |S| \leq s} \min_{\substack{\delta \neq \mathbf{0} \\ \|\delta_{S^c}\|_1 \leq c_0 \|\delta_S\|_1}} \|\mathbf{M}\delta\|_2 / (\sqrt{n} \|\delta_S\|_2) \quad (14)$$

be bounded away from zero. The rationale for taking the minimum over all such index sets  $S$  is that the true support of  $\beta$  is unknown. See also the discussion of Rudelson and Zhou (2012). We note also that the compatibility and restricted eigenvalue conditions can be replaced by slightly weaker assumptions at the cost of more involved definitions (Dalalyan et al., 2017).

Practical use of the first- and second-stage Lasso estimators requires selection of the tuning parameter  $r_j$  and  $r_\beta$  for  $j \in [p_x]$ . A number of proposals for theoretical choices of tuning parameters (Belloni and Chernozhukov, 2013; Belloni et al., 2011; Bickel et al., 2009; Bühlmann and van de Geer, 2011) for the ordinary (one-stage) linear model exist in the regularized regression literature. The bounds of Section 4.2.2 are based on oracle choices of the tuning parameters  $r_j, r_\beta$ , which depend explicitly on inestimable quantities. It would be preferable to give results for *data-adaptive* tuning parameters such as Chételat et al. (2017) and Chichignoud et al. (2016) for which the respective Lasso problems can feasibly be implemented. For the present work, we are content to demonstrate that there exist sequences of oracle tuning parameters that tend to zero sufficiently fast to ensure that the remainder terms  $f_\ell$  are asymptotically negligible. In practice, cross-validated choices of Lasso tuning parameters and  $\mu$  chosen according to the scheme described in Section 3.2 suffice in favorable parameter configurations. We provide evidence for this claim in Section 5.

#### 4.2.2. Estimation error bounds

Simultaneous control of the first-stage estimation errors  $\hat{\alpha}^j - \alpha^j$  is a straightforward consequence of the standard theoretical results for the Lasso. We include these bounds in Section B.2 of the Supplementary Materials. On the other hand, the bounds for the second-stage estimation error  $\hat{\beta} - \beta$ , which we study in the present section, are more involved due to the dependence of  $\hat{\beta}$  on the predicted conditional means  $\hat{d}_i$ . Our strategy is to write  $y = \hat{D}\beta + \tilde{u}$ , where  $\tilde{u} := u + [(D - \hat{D}) + V]\beta$  and apply concentration results to bound the probability of the event  $\{4\|\hat{D}^\top \tilde{u}\|_\infty \leq r_\beta\}$ , allowing us to adapt oracle inequality arguments for the Lasso to the present case.

We require the following assumption for the bounds of this section.

#### Assumption 4.3 (Compatibility Conditions).

- (1) There exists for each active set  $S_j$  of the first-stage model a constant  $\phi_j^2 > 0$  such that  $Z$  satisfies the compatibility condition with respect to  $S_j$  and  $\phi_j^2$ . We assume that the  $n$ -indexed sequences of such constants  $\phi_j^2$  are bounded strictly away from zero uniformly in  $n$ . We write  $\phi_A^2 := \max_{j \in [p_x]} \phi_j^2$ .
- (2) There exists a constant  $\phi_\beta$  such that  $\hat{D}$  satisfies the compatibility condition with respect to  $S_\beta$  and  $\phi_\beta$ . We assume that the  $n$ -indexed sequence of such constants is bounded strictly away from zero uniformly in  $n$ .

Assumption 4.3 imposes the compatibility condition on the random object  $\hat{D}$ . In Lemma 4.6, we provide sufficient conditions under which Assumption 4.3 holds with high probability.

Note that whether  $Z$  satisfies the compatibility condition with respect to one active set  $S_{j_1}$  does not bear directly on whether it satisfies the compatibility condition with respect to another active set  $S_{j_2}$  for  $j_1, j_2 \in [p_x]$ . As such, it is non-trivial to assume that the compatibility condition as specified in Definition 4.2 holds for each active set  $S_j$  for  $j \in [p_x]$  when  $p_x$  tends to infinity. However, the condition that  $Z$  satisfies the compatibility condition with respect to each active set  $S_j$  is entailed by requiring that  $\kappa(s, c_0)$  of (14) for  $s = \max_{j \in [p_x]} s_{\phi_j}$  and  $c_0 = 3$  be bounded away from 0. Thus, the need to accommodate multiple active sets does not thereby significantly alter the treatment of the compatibility condition.

We will refer to the following choices of tuning parameters throughout the sequel.

**Definition 4.4 (Tuning Parameters).** (1) For the first-stage Lasso estimator, set  $r_j := c_j(\|\hat{\Sigma}_Z\|_\infty(\log p_Z)/n)^{1/2}$ , where  $c_j > 0$  is a controlled quantity. We let  $r := (r_1, \dots, r_{p_x})$  denote the tuple of first-stage tuning parameters, and we write  $r_A := \max_{j \in [p_x]} r_j = c_V(\|\hat{\Sigma}_Z\|_\infty(\log p_Z)/n)^{1/2}$ , where  $c_V = \max_{j \in [p_x]} c_j$ . (2) For the second-stage Lasso estimator, set

$$r_\beta = 16 \frac{s_A}{\phi_A^2} r_A \|\hat{\Sigma}_Z\|_\infty (4m_\beta \frac{s_A}{\phi_A^2} r_A + m_A) + (4 \frac{s_A}{\phi_A^2} r_A + m_A)(m_\beta \lambda_V + \lambda_u),$$

where  $m_A, m_\beta$  are as defined in Assumption 2.1 and

$$\lambda_V := c_V(\|\hat{\Sigma}_Z\|_\infty(\log p_Z)/n)^{1/2}, \quad \lambda_u := c_u(\|\hat{\Sigma}_Z\|_\infty(\log p_Z)/n)^{1/2},$$

where  $c_u > 0$  is a controlled quantity.

Note that, unlike as in much of the related literature, the tuning parameters we identify above are random, in particular due to the term  $\|\hat{\Sigma}_Z\|_\infty$ . Our theory handles the consequences of this allowance in Condition (5) of Assumption 4.7, for which we provide subsequent justification. As we note above, our practical choice of tuning parameters is guided by cross-validation.

We now present probabilistic bounds for the  $\ell_1$  estimation error for the second-stage Lasso estimator  $\hat{\beta}$ .

**Lemma 4.5 (Bound for  $\|\hat{\beta} - \beta\|_1$ ).** Suppose that Assumption 2.4 holds and that the compatibility conditions 4.3 are satisfied with probability at least  $t_n = o(1)$ . For each  $j \in [p_x]$ , set  $r_j$  according to Definition 4.4; set  $r_\beta$  according to Definition 4.4.

Then,

$$\begin{aligned} & \mathbb{P}\left\{\|\hat{\beta} - \beta\|_1 > 4 \frac{s_\beta}{\phi_\beta^2} \left( 4 \frac{s_A}{\phi_A^2} c_V (m_\beta c_V [16 \frac{s_A}{\phi_A^2} \|\widehat{\Sigma}_Z\|_\infty + 1] + c_U) \|\widehat{\Sigma}_Z\|_\infty (\log p_Z)/n \right. \right. \\ & \quad \left. \left. + m_A (16 \frac{s_A}{\phi_A^2} \|\widehat{\Sigma}_Z\|_\infty c_V + c_U [m_\beta + 1]) (\|\widehat{\Sigma}_Z\|_\infty (\log p_Z)/n)^{1/2} \right) \right\} \\ & \leq e p_Z^{1 - c_U^2 c_0 / \tau_U^2} + e p_Z^{2 - C_0 \min_{j \in [p_X]} \{c_j^2 / \|\mathbf{v}^j\|_{\psi_2}^2\}} + t_n, \end{aligned}$$

where  $C_0$  is as specified in Lemma C.1.

**Lemma 4.5** entails that  $\|\hat{\beta} - \beta\|_1 = O_p(s_\beta s_A^2 (\log p_Z)/n + s_\beta s_A \sqrt{(\log p_Z)/n})$ . Thus, we see that the convergence rate of the second-stage Lasso estimator is slower than the typical rate of  $s_\beta \sqrt{\log(p_X)/n}$  in the ordinary (sub-)Gaussian linear model. [Zhu \(2018\)](#) provides  $L_1$  improved convergence rate of the second stage estimator with different set of assumptions. Whether the rate can be improved under the assumptions of the present paper is a direction for future work.

Since [Lemma 4.5](#) requires that [Assumption 4.3](#) holds, we need to demonstrate the latter's feasibility. The following Lemma provides such a guarantee. For other approaches to studying the empirical compatibility constants and restricted eigenvalues of random matrices, see [Rudelson and Zhou \(2012\)](#) and [van de Geer and Muro \(2014\)](#). Unlike the extant literature on the compatibility condition, however, we must account for the prediction error of  $\widehat{\mathbf{D}}$ .

**Lemma 4.6** (Second-stage Compatibility Constant). Suppose that the  $\mathbf{z}_i$  and  $\mathbf{v}^j$  satisfy [Assumptions 2.2](#) and [2.4](#), respectively. Set  $r_j$  according to [Definition 4.4](#) for each  $j \in [p_X]$ ; set  $r_A = \max_{j \in [p_X]} r_j$ . Let  $\sqrt{(\log p_Z)/n} = o(1)$ . Then, for  $n$  sufficiently large,

$$\begin{aligned} & \mathbb{P}\left\{\phi^2(\widehat{\mathbf{D}}, S_\beta) < \Lambda_{\min}(\Sigma_{\mathbf{d}}) - (a + 384 m_A s_\beta \frac{s_A}{\phi_A^2} c_V \|\widehat{\Sigma}_Z\|_\infty^{3/2}) \sqrt{(\log p_Z)/n} \right\} \\ & \leq e p_Z^{2 - C_0 \min_{j \in [p_X]} \{c_j^2 / \|\mathbf{v}^j\|_{\psi_2}^2\}} + 2 p_X^{2 - a^2 / (6e^2 \kappa^2)}, \end{aligned}$$

where  $a > 0$  is a controlled quantity,  $\Lambda_{\min}(\Sigma_{\mathbf{d}})$  denotes the minimal eigenvalue of  $\Sigma_{\mathbf{d}}$ , and  $\kappa = m_A^2 (2\tau_Z^2 + \|\Sigma_Z\|_\infty / \log 2)$ .

[Lemma 4.6](#) entails that, under mild growth conditions  $\phi^2(\widehat{\mathbf{D}}, S_\beta)$  is bounded below by a sequence of quantities approaching  $\Lambda_{\min}(\Sigma_{\mathbf{d}})$ , and hence that  $\widehat{\mathbf{D}}$  satisfies the compatibility condition with probability approaching one. Note that we must also choose the controlled quantity  $a$  above so that the exponent  $2 - a^2 / (6e^2 \kappa^2)$  is negative. The sub-Gaussian regime of the present paper entails that we may choose such an  $a$  of constant order. Throughout the present essay we use the symbol  $a$  in various bounds to denote a controlled quantity that plays this role as above.

#### 4.3. Remainder terms and scale factors

The asymptotic results of [Section 3.3](#) depend on the high-level assumption that the remainder terms  $\mathbf{f}_\ell$  and satisfy  $\|\mathbf{f}_\ell\|_\infty = o_p(1)$ . In this Section, we identify the specific conditions under which this assumption is satisfied for the two-stage Lasso. The primary goal of these conditions, which we present in [Assumption 4.7](#), is to ensure the  $\ell_1$  consistency of the first- and second-stage estimators and of the estimator  $\widehat{\Theta}$  specified in [Section 3.2](#). We implicitly refer to  $n$ -indexed sequences of all quantities mentioned below.

**Assumption 4.7** (Model Regularity for Inference of  $\beta_j$ ).

- (1) [Assumption 4.3](#) holds;
- (2) The growth condition  $\max_{j \in [p_X]} s_{\omega^j} r_j = o(1)$  holds;
- (3) The sequence of population quantities  $\Theta$  satisfies  $\Theta \in \mathcal{U}(m_\Theta, q, s_\Theta)$  for a universal constant  $m_\Theta, s_\Theta > 0$  and controlled  $q \in [0, 1]$ ;
- (4) The condition  $\mathbb{P} \mathcal{T}_\Theta(\mu)^c = o(1)$  holds;
- (5) The quantity  $\|\widehat{\Sigma}_Z\|_\infty = \|\mathbb{E}_n[\mathbf{z}_i \mathbf{z}_i^\top]\|_\infty$  satisfies  $\lim_{n \rightarrow \infty} \mathbb{P}\{\|\widehat{\Sigma}_Z\|_\infty > m_Z\} = 0$  for a universal constant  $m_Z$ ;

and

- (6) The following growth conditions hold:
  - (a)  $\mu^{1-q} s_\Theta \sqrt{\log p_Z} = o(1)$ ;
  - (b)  $s_A^3 s_\beta (\log p_Z)^{3/2} / n + s_A^2 s_\beta (\log p_Z) / \sqrt{n} = o(1)$
  - (c)  $\mu s_\beta (s_A^2 \log p_Z / \sqrt{n} + s_A \sqrt{\log p_Z}) = o(1)$ .

Condition (1) is a prerequisite for the bounds on the first- and second-stage estimation errors; we discuss the feasibility of these assumptions in [Section 4.2.2](#). Condition (2) is required for asymptotic control of the remainder terms and is comparable to typical growth rates required for Lasso consistency. Condition (3) is required to control  $\theta_j - \theta_j$  under  $\ell_\infty$  and  $\ell_1$  norms as discussed in [Section 3.2](#). Condition (4) is a high-level requirement for asymptotic negligibility of the

remainder terms  $\mathbf{f}_\ell$ ; it can be obtained as a consequence of specific model assumptions as in [Lemma 4.10](#). Condition (5) is similarly a high-level condition required for asymptotic negligibility of the remainder terms: it ensures that the empirical quantity  $\|\widehat{\Sigma}_z\|_\infty = \|\mathbf{Z}^\top \mathbf{Z}/n\|_\infty$  behaves in probability as of constant order. It can be derived as a consequence of the standard requirement that the minimal and maximal eigenvalues of  $\Sigma_z$  be bounded away from zero and infinity uniformly in  $n$  if  $\|\widehat{\Sigma}_z - \Sigma_z\|_\infty = o_p(1)$ ; the latter condition can in turn be derived from distributional assumptions on the  $\mathbf{z}_i$  using the tools of, say, [Vershynin \(2012\)](#). Condition (6) lists the model parameter growth conditions required for asymptotic negligibility of the remainder terms under the sub-Gaussian noise regime of [Assumption 2.4](#). We can compare these conditions with the requirement  $s \log p/\sqrt{n} = o(1)$  in [Javanmard and Montanari \(2014\)](#) and [van de Geer et al. \(2014\)](#) for negligibility of the single remainder term that occurs under the ordinary linear model. The conditions on the sparsity here are generally similar to those of the ordinary linear model, yet slightly more strict when comparing the powers at which the sparsity factors enter. It is not clear if our conditions can be relaxed further, or whether there are more fundamental reasons for the differences.

Note that while the quantity  $m_z$  of Condition (5) appears in the bounds of Lemmas B.8–B.14 of the Appendix, which give the rates for the remainder terms, we do not include it in the growth conditions of Condition (6). This is because, under the presently studied regime,  $m_z$  is assumed of constant order. One could consider more general scenarios where the maximum entry of  $\widehat{\Sigma}_z$  is not bounded in probability and include the quantity  $m_z$  in the aforementioned growth conditions. Doing so would in turn affect the rate at which  $s_A$ ,  $s_\beta$ , and  $p_z$  may be allowed to grow with  $n$  while maintaining asymptotic negligibility of the remainder terms  $\mathbf{f}_\ell$ .

In addition to the model regularity conditions of [Assumption 4.7](#), we require appropriate choices of the first- and second-stage Lasso tuning parameters and the estimator  $\widehat{\Theta}$ . We codify such choices in the following Assumption and then conclude the asymptotic negligibility of the remainder terms.

**Assumption 4.8** (Specification of Estimators). Let  $\widehat{\mathbf{A}}$  and  $\widehat{\beta}$  be the first- and second-stage Lasso estimators, respectively. The tuning parameters under the sub-Gaussian noise regime of [Assumption 2.4](#) are chosen according to

- (i) [Definition 4.4](#) for the first-stage tuning parameters  $\mathbf{r} = (r_j)_{j=1}^{p_x}$  and the quantity  $r_A$ ;
- (ii) [Definition 4.4](#) for the second-stage tuning parameter  $r_\beta$  and quantities  $\lambda_u$  and  $\lambda_v$ ;
- (iii) let  $\widehat{\Theta}$  be an estimator of  $\Theta$  with rows  $\widehat{\theta}_j$  given by solutions to [Problem 3.2](#).

**Lemma 4.9** (Negligibility of Remainders for [Theorem 3.4](#)). Suppose that [Assumption 2.4](#) and Conditions (1)–(6) of [Assumption 4.7](#) hold and that the estimators  $\widehat{\mathbf{A}}$ ,  $\widehat{\beta}$ , and  $\widehat{\Theta}$  are chosen according to [Assumption 4.8](#)(i)–(iii). Then,  $\|\mathbf{f}_\ell\|_\infty = o_p(1)$  for  $\ell \in [4]$ .

The primary use of [Lemma 4.9](#) is to verify Condition (2) of [Theorem 3.4](#). Indeed, the result justifies the use of the one-step update to the second-stage Lasso estimator to construct asymptotically valid confidence intervals for the components  $\beta_j$  according to (9).

We note that the quantity  $\mu$ , which we recall is the tolerance parameter for [Problem 3.2](#), must be given careful consideration. Conditions (6)a and (6)b of [Assumption 4.7](#) require  $\mu$  to be of small order  $(s_\Theta \sqrt{\log p_z})^{\frac{1}{q-1}}$  and  $(s_\beta s_A \log p_z)^{-1}$ , respectively. However,  $\mu$  must not tend to zero so fast that the probability that  $\Theta$  is feasible for [Problem 3.2](#), which we recall is formally denoted by  $P \mathcal{T}_\Theta(\mu)$ , becomes bounded away from zero. The following Lemma identifies a choice of  $\mu$  that satisfies these competing objectives for sub-Gaussian  $\mathbf{z}_i$  and first-stage noise elements.

**Lemma 4.10** (Probability of  $\mathcal{T}_\Theta(\mu)$ ). Suppose that (i) the  $\mathbf{z}_i$  and  $\mathbf{v}^j$  satisfy [Assumptions 2.2](#) and [2.4](#), respectively; (ii)  $\widehat{\mathbf{A}}$  consists of first-stage Lasso estimates of the  $\alpha^j$  with tuning parameters  $r_j$  chosen according to [Definition 4.4](#). Set

$$\mu = \frac{m_\Theta}{\sqrt{n}} \left( a \sqrt{\log p_x} + 12 m_A c_V \|\widehat{\Sigma}_z\|_\infty^{\frac{3}{2}} \frac{s_A}{\phi_A^2} \sqrt{\log p_z} \right),$$

where  $c_V$  is as in [Definition 4.4](#). Then, for  $n$  sufficiently large,

$$P\{\|\Theta \widehat{\Sigma}_d - \mathbf{I}\|_\infty > \mu\} \leq 2p_x^{2-a^2/(6e^2\eta^2)} + ep_z^{2-C_0 \min_{j \in [p_x]} \{c_j^2/\|\mathbf{v}^j\|_{\psi_2}^2\}},$$

where  $\eta = m_A^2(2\tau_z^2 + \|\Sigma_z\|_\infty/\log 2)$  and  $C_0$  is as defined in [Lemma C.1](#).

If there exists  $m_z = O(1)$  that satisfies  $P\{\|\widehat{\Sigma}_z\|_\infty > m_z\} = o(1)$ , we may substitute the former quantity into the specification of  $\mu$  in [Lemma 4.10](#) to obtain

$$P\{\|\Theta \widehat{\Sigma}_d - \mathbf{I}\|_\infty > \mu\} \leq 2p_x^{2-a^2/(6e^2\eta^2)} + 2p_z^{1-c_{ep}} + P\{\|\widehat{\Sigma}_z\|_\infty > m_z\},$$

which tends to zero under appropriate specification of the controlled quantities  $a$  and  $c$ . Note that  $\mu \lesssim s_A \sqrt{\log(p_z)/n}$ ; Condition (6)a of [Assumption 4.7](#) then becomes

$$s_\Theta s_A^{1-q} (\log p_z)^{1-\frac{q}{2}} / n^{\frac{1-q}{2}} = o(1), \quad (15)$$



and Condition (6)c becomes

$$\mathbf{s}_\beta \mathbf{s}_\mathbf{A}^3 (\log p_z)^{3/2} / n + \mathbf{s}_\beta \mathbf{s}_\mathbf{A}^2 \log(p_z) / \sqrt{n} = o(1), \quad (16)$$

which is identical to Condition (6)b.

Recall that Theorem 3.4 specifies the conditions under which  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j$  converges weakly to a  $\mathcal{N}(0, 1)$  random variable and demonstrates that the limit continues to hold if  $\omega_j$  is replaced with an estimator  $\hat{\omega}_j$  that satisfies  $|\hat{\omega}_j - \omega_j| = o_p(1)$ . In practice,  $\omega_j$  is not available, hence it is crucial to demonstrate the existence of an estimator  $\hat{\omega}_j$  that satisfies the foregoing condition. The following Lemma identifies such an estimator  $\hat{\omega}_j$  and the conditions under which it is appropriate for use with Theorem 3.4.

**Lemma 4.11** (Existence of Appropriate  $\hat{\omega}_j$ ). Suppose that (i) (a) the second-stage noise elements  $u_i$  satisfy  $E[u_i^2 | \mathbf{z}_i] = \sigma_u^2$  and  $\|\mathbf{u}\|_2^2/n - \sigma_u^2 = o_p(1)$ , (b) the  $\mathbf{z}_i$  and  $\mathbf{v}^j$  satisfy Assumptions 2.2 and 2.4, (c) we have that  $\mathbf{s}_\beta \mathbf{s}_\mathbf{A}^2 \log(p_z)/n + \mathbf{s}_\beta \mathbf{s}_\mathbf{A} \sqrt{\log(p_z)/n} = o(1)$  and  $\max_{j \in [p_x]} \max_{i \in [n]} E[x_{ij}^2] = O(1)$ . Let  $\hat{\mathbf{A}}$ ,  $\hat{\beta}$ , and  $\hat{\Theta}$  be as specified in Assumption 4.8. Define the estimator  $\hat{\sigma}_u$  of the second-stage noise level  $\sigma_u$  by

$$\hat{\sigma}_u^2 := \|\mathbf{y} - \mathbf{X}\hat{\beta}\|_2^2/n. \quad (17)$$

Then  $\hat{\sigma}_u - \sigma_u = o_p(1)$ . If, in addition to the conditions in (i), we have (ii) (a) Condition (4) of Assumption 4.7 holds, (b) the sequence of minimal and maximal eigenvalues of  $\Sigma_d$ , denoted respectively by  $\Lambda_{\min}(\Sigma_d)$  and  $\Lambda_{\max}(\Sigma_d)$ , are bounded away from zero and infinity uniformly in  $n$ ; and (c) the tolerance  $\mu$  satisfies  $\mu = o(1)$ , then  $\hat{\omega}_j$  defined by

$$\hat{\omega}_j^2 := \hat{\sigma}_u^2 \hat{\Theta}_{jj} \quad (18)$$

satisfies  $\hat{\omega}_j - \omega_j = o_p(1)$ .

## 5. Numerical experiments

In this section, we present a Monte Carlo simulation study of the finite-sample properties of the inferential procedure developed in Section 3 using the two-stage Lasso studied in Section 4. Our objective is to test this method under a variety of parameter configurations chosen to reflect settings of practical interest. In Section 5.1, we describe the general scheme according to which the data for each trial are generated and the metrics gathered for each configuration. In Section 5.2, we enumerate the specific parameter configurations studied and discuss the results.

### 5.1. General experimental design

Each trial contains a data-generation step and an estimation step. We specify the regression parameters  $\beta$  and  $\mathbf{A}$  for the data-generation step as follows. For each configuration, we set the second-stage regression parameter  $\beta$  according to  $\beta_j = 1$  for  $j \in S_\beta$  and  $\beta_j = 0$  otherwise, where  $S_\beta \subset [p_x]$  is a random set of  $s_\beta$  generated by uniformly random draws from  $[p_x]$  without replacement. Similarly, we set the first-stage regression parameters  $\alpha^j$  for  $j \in [p_x]$  according to  $\alpha_k^j = 1$  for  $k \in S_j$  and  $\alpha_k = 0$  otherwise, where  $S_j \subset [p_z]$  is a random set of  $s_\mathbf{A}$  generated by uniformly random draws from  $[p_z]$  without replacement. We let  $s_\beta, s_\mathbf{A}$  vary over configurations.

Having specified the regression parameters, we then draw  $n$  i.i.d. observations  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  according to

$$\begin{aligned} \mathbf{z}_i &\sim \mathcal{N}_{p_z}(\mathbf{0}, \Sigma_z), & x_{ij} &= \langle \mathbf{z}_i, \alpha^j \rangle + v_{ij}, \\ (u_i, \mathbf{v}_i) | \mathbf{z}_i &\sim \mathcal{N}_{1+p_x}(\mathbf{0}, \Sigma_{uv}) & y_i &= \langle \mathbf{x}_i, \beta \rangle + u_i, \end{aligned}$$

where  $n, p_x, p_z, \beta, \{\alpha^j\}_{j \in [p_x]}$ , and the structure of  $\Sigma_z$  vary amongst configurations. For all configurations, we set

$$\Sigma_{uv} = \begin{pmatrix} \sigma_u & \sigma_{uv}^\top \\ \sigma_{uv} & \sigma_v^\top \mathbf{I} \end{pmatrix}$$

where  $\sigma_u, \sigma_v = \sqrt{.7}$  are held fixed and  $\sigma_{uv} = (\sigma_{uv^1}, \dots, \sigma_{uv^{p_x}})$  is given as follows. For each configuration, we set one  $\sigma_{uv^j}$  chosen at random equal to .5, nine  $\sigma_{uv^j}$  chosen at random equal to .25, and the remaining  $\sigma_{uv^j}$  equal to .05. The present covariance structure for the noise reflects the constraint that  $\Sigma_{uv}$  be positive-definite; our choices of  $\sigma_{uv^j}$  are an attempt to balance this requirement with the goal of studying non-trivial correlations between the first- and second-stage noise elements.

We consider two forms for the covariance matrix  $\Sigma_z$ . The first form is a Toeplitz (TZ) structure given by

$$\Sigma_z^{\text{TZ}}|_{jk} = \rho^{|j-k|}, \quad j, k \in [p_z], \quad \rho = 0.8.$$

The second is a circulant-symmetric (CS) structure given for  $j \leq k$  by

$$\Sigma_z^{\text{CS}}|_{jk} = \begin{cases} 1 & k = j, \\ 0.1 & k \in \{j+1, \dots, j+5\} \cup \{j+p_z-5, \dots, j+p_z-1\}, \\ 0 & \text{otherwise.} \end{cases}$$

**Table 1**  
Simulation results.

$(n, p_x, p_z)$	$(s_\beta, s_{\alpha^j})$	$\Sigma_z^{CS}$			$\Sigma_z^{TZ}$		
		$\widehat{cv\hat{g}}$	$\widehat{len}$	$MSE(\hat{\beta})$	$\widehat{cv\hat{g}}$	$\widehat{len}$	$MSE(\hat{\beta})$
(100, 125, 150)	(3, 5)	0.942	0.225	0.004	0.895	0.201	0.005
	(5, 10)	0.941	0.211	0.004	0.672	0.212	0.014
	(10, 15)	0.930	0.190	0.003	0.545	0.219	0.030
(200, 250, 275)	(3, 5)	0.947	0.157	0.002	0.942	0.140	$\leq 0.001$
	(5, 10)	0.941	0.171	0.002	0.673	0.192	0.011
	(10, 15)	0.930	0.190	0.003	0.545	0.219	0.030
(300, 400, 500)	(3, 5)	0.947	0.094	$\leq 0.001$	0.952	0.092	$\leq 0.001$
	(5, 10)	0.955	0.085	$\leq 0.001$	0.945	0.082	$\leq 0.001$
	(10, 15)	0.961	0.067	$\leq 0.001$	0.927	0.064	$\leq 0.001$
(500, 600, 700)	(3, 5)	0.947	0.094	$\leq 0.001$	0.952	0.092	$\leq 0.001$
	(5, 10)	0.951	0.082	$\leq 0.001$	0.950	0.088	$\leq 0.001$
	(10, 15)	0.961	0.067	$\leq 0.001$	0.927	0.064	$\leq 0.001$

Within a configuration study, the random quantities  $\mathbf{z}_i$ ,  $u_i$ , and  $\mathbf{v}_i$  are re-drawn for each trial; the quantities  $\beta$ ,  $\{\alpha^j\}_{j=1}^{p_x}$ ,  $\Sigma_z$ ,  $\sigma_{uv}$ ,  $n$ ,  $p_x$ ,  $p_z$  are held fixed.

For the estimation step of each trial, we compute the first- and second-stage Lasso as defined in Section 4.1 using the `glmnet` package (Friedman et al., 2010). Tuning parameters  $r$  for the Lasso estimators are selected by 10-fold cross-validation over a grid  $\{r_\ell\}_{\ell=1}^L$ , where  $L = 100$ ,  $r_L = .01r_1$ , and  $r_1$  is the least quantity for which the respective Lasso estimator is identically 0. The tuning parameter  $r_\beta$  is chosen by similar cross-validation procedure.

The rows  $\hat{\theta}_j$  of  $\hat{\Theta}$  are obtained as solutions to the respective Problem 3.2 with tuning parameter  $\mu_j$  chosen as follows. For each  $j \in [p_x]$ , we set  $\mu_j := \kappa \times \inf_{\theta \in \mathbb{R}^{p_x}} \|\hat{\Sigma}_d \theta - \mathbf{e}_j\|_\infty$ , where  $\mathbf{e}_j$  denotes the  $j$ th canonical basis vector in  $p_x$  dimensions and  $\kappa > 1$  is chosen at our discretion. Note that, under this choice of  $\mu_j$ , the respective Problem 3.2 is guaranteed feasible. The factor  $\kappa$  is chosen to balance the performance of  $\hat{\Theta}$  as a surrogate inverse for  $\hat{\Sigma}_d$ , for which a smaller  $\kappa$  is desirable, with the size of the objective function  $\|\theta\|_1$ , for which a larger  $\kappa$  is desirable. The following results were obtained under  $\kappa = 1.2$ . To obtain the infimum, we cast  $\min_{\theta \in \mathbb{R}^{p_x}} \|\hat{\Sigma}_d \theta - \mathbf{e}_j\|_\infty$  as a linear programming problem, which we solve using MOSEK optimization software (MOSEK AsP, 2017).

In a given trial  $t$ ,  $t = 1, \dots, T$ , we set  $\tau = 0.05$  and compute the respective  $100(1 - \tau)\%$  confidence interval

$$\hat{\mathcal{I}}_{\tau,t,j} = [\tilde{\beta}_{t,j} - z_\tau \widehat{SE}(\tilde{\beta}_{t,j}), \tilde{\beta}_{t,j} + z_\tau \widehat{SE}(\tilde{\beta}_{t,j})],$$

for each component  $\tilde{\beta}_{t,j}$  of  $\tilde{\beta}_t$ , where  $z_\tau = \Phi^{-1}(1 - \tau/2)$  and

$$\widehat{SE}(\tilde{\beta}_j)^2 = \mathbb{E}_n[(y_i - \hat{\beta}X)^2 \langle \hat{\theta}_j, \hat{\mathbf{d}}_i \rangle^2].$$

For each configuration of  $n$ ,  $p_x$ ,  $p_z$ ,  $s_\beta$ ,  $s_A$ ,  $\Sigma_z$ , we generate  $T = 100$  trials and calculate the average coverage  $\widehat{cv\hat{g}}$  for the 95% confidence intervals  $\hat{\mathcal{I}}_{\tau,j}$  about components of  $\tilde{\beta}$  and the average interval length  $\widehat{len}$  given by

$$\widehat{cv\hat{g}} = \frac{1}{p_x} \sum_{j=1}^{p_x} \frac{1}{T} \sum_{t=1}^T \mathbf{1}\{\beta_{t,j} \in \hat{\mathcal{I}}_{\tau,t,j}\}, \quad \widehat{len} = \frac{1}{p_x} \sum_{j=1}^{p_x} \frac{1}{T} \sum_{t=1}^T \text{len}(\hat{\mathcal{I}}_{\tau,t,j}).$$

For each configuration, we also provide the average mean squared error of the second-stage Lasso estimator given by  $\widehat{MSE} = \frac{1}{T} \sum_{t=1}^T \frac{1}{p_x} \sum_{j=1}^{p_x} (\tilde{\beta}_{t,j} - \beta_{t,j})^2$ . We present the results for the study described above in Table 1 in Section 5.2.

## 5.2. Specifications and results

We conduct simulations according to the design described in Section 5.1 for all configurations belonging to

$$\underbrace{\begin{pmatrix} (100, 125, 150) \\ (200, 250, 275) \\ (300, 400, 500) \\ (500, 600, 700) \end{pmatrix}}_{(n, p_x, p_z)} \times \underbrace{\begin{pmatrix} (3, 5) \\ (5, 10) \\ (10, 15) \end{pmatrix}}_{(s_\beta, s_A)} \times \underbrace{\begin{pmatrix} \Sigma_z^{CS} \\ \Sigma_z^{TZ} \end{pmatrix}}_{\Sigma_z}.$$

The results, which are presented in Table 1, show that our estimator achieves close to nominal coverage under a variety of configurations. We also see that arguably the greatest determinant of coverage is the relative magnitude of  $p_x$  and  $p_z$  to the size of the active set  $s_\beta$ . As the latter grows, coverage diverges from the nominal level. This phenomenon is expected and has been found in ordinary linear regression models as well van de Geer et al. (2014) and Javanmard and Montanari (2014), since the bounds for the estimation error of the Lasso are proportional to the size of the active set. Nevertheless, the

performance improves significantly when we increase the sample size. Finally, we observe that the covariance structure of the instrumental variables  $\mathbf{z}_i$  has a strong influence on coverage: the Toeplitz structure features greater correlation among the instrumental variables in general, and this is reflected in coverage that tends to be farther from the nominal level than in the case of the circulant-symmetric covariance structure. These results suggest that our proposed method of inference for the low-dimensional components of a high-dimensional regression vector is relevant to practical scenarios that may exhibit non-trivial degrees of correlation between the noise components and nontrivial correlation among the instrumental variables.

## 6. Conclusion

In this paper, we propose inference methods for the components of a high-dimensional instrumental variables regression parameter despite possible endogeneity of each regressor. We allow both the number of instruments and the number of regressors to be greater than the sample size. We construct asymptotically valid confidence intervals for the components of the second-stage regression coefficients. Though our estimator is not a nonlinear generalized method of moments (GMM) estimator (Hansen, 1982), we expect that our results can be extended to that more general setting.

Our Sections 2 and 3 comprise a general pipeline for estimation and inference, while the remainder is then an exemplification with a Lasso approach. Therefore, it would be interesting to use our pipeline with other regularized estimators as well (one could also use different estimators for the two different stages). Candidates include, for example, SCAD (Fan and Li, 2001) and MCP (Zhang, 2010). We finally refer to our software on [github.com/LedererLab/HDIV](https://github.com/LedererLab/HDIV).

## Appendix A. Materials required for Section 3

### A.1. Materials required for Section 3.1

**Proof of Lemma 3.1.** Note that

$$\begin{aligned}\tilde{\beta} &= \hat{\beta} - \hat{\Theta} \mathbb{E}_n[\hat{\mathbf{h}}(\hat{\beta})] \\ &= \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{y} - \mathbf{X} \hat{\beta})/n \\ &= \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X}[\beta - \hat{\beta}] + \mathbf{u})/n \\ &= \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top (\hat{\mathbf{D}}[\beta - \hat{\beta}] + [\mathbf{X} - \hat{\mathbf{D}}][\beta - \hat{\beta}] + \mathbf{u})/n \\ &= \beta + \hat{\Theta} \hat{\mathbf{D}}^\top \mathbf{u}/n + \underbrace{\hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})(\beta - \hat{\beta})/n}_{f_3/\sqrt{n}} + \underbrace{(\hat{\Theta} \hat{\Sigma}_d - \mathbf{I})(\beta - \hat{\beta})}_{f_4/\sqrt{n}}.\end{aligned}$$

Now decompose the second term on the right-hand side above as follows

$$\begin{aligned}\hat{\Theta} \hat{\mathbf{D}}^\top \mathbf{u}/n &= \hat{\Theta} \mathbf{D}^\top \mathbf{u}/n + \hat{\Theta} (\hat{\mathbf{D}} - \mathbf{D})^\top \mathbf{u}/n \\ &= \underbrace{\hat{\Theta} \mathbf{D}^\top \mathbf{u}/n}_{f_1/\sqrt{n}} + \underbrace{(\hat{\Theta} - \Theta) \mathbf{D}^\top \mathbf{u}/n}_{f_2/\sqrt{n}} + \underbrace{\hat{\Theta} (\hat{\mathbf{D}} - \mathbf{D})^\top \mathbf{u}/n}_{f_2/\sqrt{n}}\end{aligned}$$

to complete the proof.  $\square$

### A.2. Materials required for Section 3.2

We present the properties of the estimators  $\hat{\theta}_j$  required for the results of Sections 3.3 and 4.3. Lemma A.1, which gives an  $\ell_\infty$  bound for the estimation error  $\hat{\theta}_j - \theta_j$ , is comparable to Cai et al. (2011, Theorem 4); the proofs are similar but depend on different conditions on the covariance estimator  $\hat{\Sigma}_d$ . The proof of Lemma A.2 follows part of that of Cai et al. (2011, Theorem 6).

**Lemma A.1.** Suppose that: (i) the quantity  $\|\Theta\|_{L_1}$  is bounded above by a constant  $m_\Theta < \infty$ ; and (ii)  $\hat{\Theta}$  is an estimate of  $\Theta = \Sigma_d^{-1} = \text{cov}(\mathbf{d}_i)^{-1}$  with rows  $\hat{\theta}_j$  obtained as solutions to Problem 3.2. Then, on the set  $\mathcal{T}_\Theta(\mu)$  as defined in (8),

$$\|\hat{\theta}_j - \theta_j\|_\infty \leq 2m_\Theta \mu$$

for each  $j \in [p_x]$ .

**Proof of Lemma A.1.** First, observe that the conditions of the present lemma entail that

$$\|\hat{\Theta} \hat{\Sigma}_d - \mathbf{I}\|_\infty \leq \mu, \quad \|\hat{\Theta} \hat{\Sigma}_d - \mathbf{I}\|_\infty \leq \mu.$$

Now, on the set  $\mathcal{T}_\Theta(\mu)$ , each row  $\theta_j$  is feasible for the respective Specific [Problem 3.2](#). It then follows from the optimality of  $\hat{\theta}_j$  that  $\|\hat{\theta}_j\|_1 \leq \|\theta_j\|_1$  for each  $j \in [p_x]$  and hence that  $\max_{j \in [p_x]} \|\hat{\theta}_j\|_1 \leq \|\Theta\|_{L_1}$ . Next, observe that

$$\begin{aligned} \Theta - \hat{\Theta} &= \Theta(I - \Sigma_d \hat{\Theta}) = \Theta(I + (\hat{\Sigma}_d - \Sigma_d) \hat{\Theta} - \hat{\Sigma}_d \hat{\Theta}) \\ &= \Theta(I - \hat{\Sigma}_d \hat{\Theta}) - \Theta(\Sigma_d - \hat{\Sigma}_d) \hat{\Theta} \\ &= \underbrace{\Theta(I - \hat{\Sigma}_d \hat{\Theta})}_{I_1} - \underbrace{\Theta(\Sigma_d - \hat{\Sigma}_d) \hat{\Theta}}_{I_2}. \end{aligned}$$

From Hölder,

$$\|I_2\|_\infty \leq \|I - \Theta \hat{\Sigma}_d\|_\infty \|\hat{\Theta}\|_{L_1} \leq m_\Theta \mu.$$

The matrix  $\ell_\infty$ - and ( $L_1$ - norm if the argument is symmetric) are invariant under transposition of their arguments. Use this fact and Hölder to obtain

$$\begin{aligned} \|I_1\|_\infty &= \|(I - \hat{\Sigma}_d \hat{\Theta})^\top \Theta^\top\|_\infty \\ &\leq \|(I - \hat{\Sigma}_d \hat{\Theta})^\top\|_\infty \|\Theta^\top\|_{L_1} \\ &\leq \|I - \hat{\Sigma}_d \hat{\Theta}\|_\infty \|\Theta\|_{L_1} \leq m_\Theta \mu, \end{aligned}$$

where the final line follows from the fact that both  $I - \hat{\Sigma}_d \hat{\Theta}$  and  $\Theta$  are symmetric. Thus

$$\|\Theta - \hat{\Theta}\|_\infty \leq \|I_1\|_\infty + \|I_2\|_\infty \leq 2m_\Theta \mu,$$

as required.  $\square$

**Lemma A.2.** Suppose in addition to the conditions of [Lemma A.1](#) that  $\Theta$  belongs to the uniformity class  $\mathcal{U}(m_\Theta, q, s_\Theta)$ . Then,

$$\|\hat{\theta}_j - \theta_j\|_1 \leq 2c_q(2m_\Theta \mu)^{1-q} s_\Theta$$

for each  $j \in [p_x]$ , where  $c_q := 1 + 2^{1-q} + 3^{1-q}$ .

**Proof of Lemma A.2.** See the proof of line (14) of [Cai et al. \(2011, Theorem 6\)](#).  $\square$

### A.3. Materials required for Section 3.3

**Proof of Theorem 3.4.** The proof of the first claim consists of two steps. First, we show that the quantity

$$Z_{j,n} := \frac{1}{\sqrt{n}} \sum_{i=1}^n \langle \theta_j, \mathbf{d}_i \rangle u_i / \omega_j.$$

and  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j$  share the same weak limit. Second, we show that  $Z_{j,n} \rightsquigarrow \mathcal{N}(0, 1)$ . To establish the first step, we claim that

$$\lim_{n \rightarrow \infty} \mathbb{P}\{\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j \leq t\} \leq \lim_{n \rightarrow \infty} \mathbb{P}\{Z_{j,n} \leq t\}, \quad (\text{A.1})$$

for all  $t \in \mathbb{R}$ . An analogous lower bound follows by a matching argument, which shows that  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j$  and  $Z_{j,n}$  share the same weak limit. To show the claim above, let  $t \in \mathbb{R}$  be given, fix a controlled  $\epsilon > 0$ , and note that, by the decomposition of [Lemma 3.1](#), we have

$$\begin{aligned} \mathbb{P}\{\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j \leq t\} &\leq \mathbb{P}\left\{Z_{j,n} + \sum_{\ell=1}^4 f_{\ell,j}/\omega_j \leq t + 4\epsilon\right\} \\ &\leq \mathbb{P}\{Z_{j,n} \leq t + \epsilon\} + \sum_{\ell=1}^4 \mathbb{P}\{f_{\ell,j}/\omega_j > \epsilon\}. \end{aligned}$$

By specification of  $\sigma_u$  and  $\Theta_{jj}$  in Conditions (5) and (6) and Condition (3) of the present theorem, it follows that  $\omega_j$  is bounded strictly away from 0 uniformly in  $n$ . The assumptions of the present theorem then entail that  $\mathbb{P}\{f_{\ell,j}/\omega_j > \epsilon\} = o(1)$  for all  $\epsilon > 0$  and each  $\ell$ . Letting  $\epsilon$  tend to zero shows the claim of (A.1). It follows from the analogous lower bound that

$$\lim_{n \rightarrow \infty} \mathbb{P}\{\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j \leq t\} = \lim_{n \rightarrow \infty} \mathbb{P}\{Z_{j,n} \leq t\}$$

for all  $t \in \mathbb{R}$ , thus completing the first step.

Next, we show that, under each of Conditions (5) and (6) in the statement of the present theorem,  $Z_{j,n} \rightsquigarrow Z_j \sim \mathcal{N}(0, 1)$ . To this end, we define the quantity

$$w_j^2 := \boldsymbol{\theta}_j^\top \bar{\Sigma}_d \boldsymbol{\theta}_j = \frac{1}{n} \sum_{i=1}^n \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2,$$

where we recall that  $\bar{\Sigma}_d = \mathbf{D}^\top \mathbf{D}/n$ . We claim first that

$$\tilde{Z}_{j,n} := \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle u_i}{w_j \sigma_u} \rightsquigarrow Z_j \sim \mathcal{N}(0, 1)$$

under each of Conditions (5) and (6) and second that  $\sigma_u w_j / \omega_j \rightarrow_p 1$ . Note that  $Z_{j,n} = \frac{w_j \sigma_u}{\omega_j} \tilde{Z}_{j,n}$ , hence the desired limit follows from an application Slutsky's Lemma.

To show the first claim under Condition (5), note that, by specification of  $w_j$ , we have under Assumption 2.4 that

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle u_i}{w_j \sigma_u} \mid \mathbf{Z} \sim \mathcal{N}(0, 1).$$

Thus

$$\lim_{n \rightarrow \infty} \mathbb{P}(\tilde{Z}_{j,n} \leq t) = \lim_{n \rightarrow \infty} \mathbb{E}[\mathbb{P}(\tilde{Z}_{j,n} \leq t \mid \mathbf{Z})] = \lim_{n \rightarrow \infty} \mathbb{E}[\Phi(t) \mid \mathbf{Z}] = \Phi(t)$$

for all  $t \in \mathbb{R}$ , where  $\Phi$  denotes the c.d.f. of a standard Normal random variable. This shows the desired weak limit under Condition (5).

We use the Lindeberg–Feller Central Limit Theorem to show the limit under Condition (6). To begin, write

$$\tilde{Z}_{j,n} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \xi_i, \quad \xi_i := \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle u_i / (w_j \sigma_u).$$

Note that

$$\mathbb{E}[\xi_i] = \mathbb{E}[\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle / (w_j \sigma_u) \mathbb{E}[u_i \mid \mathbf{z}_i]] = 0$$

and that

$$\sigma_n^2 := \sum_{i=1}^n \mathbb{E}[\xi_i^2] = \mathbb{E} \left[ \frac{1}{w_j^2} \sum_{i=1}^n \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 \mathbb{E}[(u_i / \sigma_u)^2 \mid \mathbf{z}_i] \right] = n.$$

To demonstrate Lindeberg's condition, let  $\delta > 0$  be arbitrary and write

$$\begin{aligned} \sigma_n^{-2} \sum_{i=1}^n \mathbb{E}[\xi_i^2 \mathbb{1}_{\{|\xi_i| > \delta \sigma_n\}}] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[\mathbb{E}[\xi_i^2 \mathbb{1}_{\{|\xi_i| > \delta \sqrt{n}\}} \mid \mathbf{z}_i]] \\ &= \mathbb{E} \left[ \frac{1}{n} \frac{1}{w_j^2} \sum_{i=1}^n \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 \mathbb{E}[(u_i / \sigma_u)^2 \mathbb{1}_{\{|\xi_i| > \delta \sqrt{n}\}} \mid \mathbf{z}_i] \right], \end{aligned}$$

where we have substituted the definition of  $\xi_i$  and extracted the factor  $\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 / w_j^2$  from the conditional expectation to obtain the second line. Now substitute the definition of  $w_j$  below and note

$$\frac{1}{w_j^2} \sum_{i=1}^n \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 = \frac{n}{\sum_{i=1}^n \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2} \times \sum_{i=1}^n \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 = n.$$

Substitute the above result into the second line of two displays prior and use the law of iterated expectation to obtain

$$\sigma_n^{-2} \sum_{i=1}^n \mathbb{E}[\xi_i^2 \mathbb{1}_{\{|\xi_i| > \delta \sigma_n\}}] = \mathbb{E}[(u_1 / \sigma_u)^2 \mathbb{1}_{\{|\xi_1| > \delta \sqrt{n}\}}],$$

where we cite condition (6) that the  $\mathbf{z}_i$  and  $u_i$  are i.i.d. with  $\mathbb{E}[u_i^2 \mid \mathbf{z}_i] = \sigma_u^2$  to reduce to the case of  $u_1$ . For brevity, we write  $\tilde{u}_1 := u_1 / \sigma_u$ . Introduce the set  $\mathcal{T} := \{|\langle \boldsymbol{\theta}_j, \mathbf{d}_1 \rangle| \leq n^\zeta\}$  and note, since  $|\xi_1| \leq w_j^{-1} |\langle \boldsymbol{\theta}_j, \mathbf{d}_1 \rangle| |\tilde{u}_1|$ , that

$$\{|\xi_1| > \delta \sqrt{n}\} \cap \mathcal{T} \subseteq \{|\tilde{u}_1| > \delta w_j n^{1/2-\zeta}\},$$

and hence that  $\mathbb{1}_{\{|\xi_1| > \delta \sqrt{n}\}} \cap \mathcal{T} \leq \mathbb{1}_{\{|\tilde{u}_1| > \delta w_j n^{1/2-\zeta}\}}$ . Combine this inequality and the law of total probability with the result of two displays previous to obtain

$$\sigma_n^{-2} \sum_{i=1}^n \mathbb{E}[\xi_i^2 \mathbb{1}_{\{|\xi_i| > \delta \sigma_n\}}] = \mathbb{E}[\tilde{u}_1^2 (\mathbb{1}_{\{|\xi_1| > \delta \sqrt{n}\}} \cap \mathcal{T} + \mathbb{1}_{\{|\xi_1| > \delta \sqrt{n}\}} \cap \mathcal{T}^c)]$$



$$= \underbrace{E[\tilde{u}_1^2 \mathbb{1}_{\{|\tilde{u}_1| > \delta w_j n^{1/2-\zeta}\}}]}_{I_1} + \underbrace{E[\tilde{u}_1^2 \mathbb{1}_{\mathcal{T}^c}]}_{I_2},$$

where the substitution of indicators in the final line above is permitted since  $\tilde{u}_1^2 \geq 0$ . To show Lindeberg's condition, it suffices to show that  $I_1$  and  $I_2$  are each  $o(1)$ . To treat  $I_1$ , consider the event  $\{w_j \leq \vartheta^{1/2}/\sqrt{2}\}$  and write

$$\begin{aligned} \mathbb{1}_{\{|\tilde{u}_1| > \delta w_j n^{1/2-\zeta}\}} &= \mathbb{1}_{\{|\tilde{u}_1| > \delta w_j n^{1/2-\zeta}\} \cap \{w_j \leq \vartheta^{1/2}/\sqrt{2}\}} \\ &\quad + \mathbb{1}_{\{|\tilde{u}_1| > \delta w_j n^{1/2-\zeta}\} \cap \{w_j > \vartheta^{1/2}/\sqrt{2}\}} \\ &\leq \mathbb{1}_{\{w_j \leq \vartheta^{1/2}/\sqrt{2}\}} + \mathbb{1}_{\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\}} \end{aligned}$$

so that

$$I_1 \leq \underbrace{E[\tilde{u}_1^2 \mathbb{1}_{\{w_j \leq \vartheta^{1/2}/\sqrt{2}\}}]}_{I_{1a}} + \underbrace{E[\tilde{u}_1^2 \mathbb{1}_{\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\}}]}_{I_{1b}}.$$

Observe that

$$\begin{aligned} I_{1a} &= E[\mathbb{1}_{\{w_j \leq \vartheta^{1/2}/\sqrt{2}\}} E[\tilde{u}_1^2 \mid \mathbf{z}_i]] \\ &\lesssim P\{w_j \leq \vartheta^{1/2}/\sqrt{2}\} \\ &= P\{w_j^2 - \Theta_{jj} \leq \vartheta/2 - \Theta_{jj}\} \\ &\leq P\{w_j^2 - \Theta_{jj} \leq -\vartheta/2\} \\ &\leq P\{|w_j^2 - \Theta_{jj}| \geq \vartheta/2\}, \end{aligned}$$

where the inference to the second line of the display above follows from specification of  $E[\tilde{u}_1^2 \mid \mathbf{z}_i]$  in condition (6) of the present Theorem. Now note that

$$\begin{aligned} |w_j^2 - \Theta_{jj}| &= |\boldsymbol{\theta}_j^\top \bar{\boldsymbol{\Sigma}}_{\mathbf{d}} \boldsymbol{\theta}_j - \Theta_{jj}| \\ &= |\boldsymbol{\theta}_j^\top (\bar{\boldsymbol{\Sigma}}_{\mathbf{d}} - \boldsymbol{\Sigma}_{\mathbf{d}}) \boldsymbol{\theta}_j| \leq \|\boldsymbol{\theta}_j\|_1^2 \|\bar{\boldsymbol{\Sigma}}_{\mathbf{d}} - \boldsymbol{\Sigma}_{\mathbf{d}}\|_\infty. \end{aligned}$$

Thus

$$\begin{aligned} I_{1a} &\leq P\{|w_j^2 - \Theta_{jj}| \geq \vartheta/2\} \\ &\leq P\{\|\bar{\boldsymbol{\Sigma}}_{\mathbf{d}} - \boldsymbol{\Sigma}_{\mathbf{d}}\|_\infty \geq \vartheta/(2\|\boldsymbol{\theta}_j\|_1^2)\} = o(1) \end{aligned} \tag{A.2}$$

by Conditions (1) and (4) of the present theorem.

We now show that

$$I_{1b} = E[\tilde{u}_1^2 \mathbb{1}_{\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\}}] = o(1).$$

Consider the event  $\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\}$ , which is the index of the indicator function above. Exponentiate both sides of the inequality in that set by  $\nu$  to obtain that

$$\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\} = \{|\tilde{u}_1|^\nu > (\delta \vartheta^{1/2})^\nu n^{\nu(1/2-\zeta)}\}.$$

Conclude of the respective indicator functions that

$$\mathbb{1}_{\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\}} = \mathbb{1}_{\{|\tilde{u}_1|^\nu > (\delta \vartheta^{1/2})^\nu n^{\nu(1/2-\zeta)}\}}.$$

Note by direct manipulation that, given the event that indexes the right-hand indicator above, it holds that

$$(\delta \vartheta^{1/2})^{-\nu} n^{-\nu(1/2-\zeta)} |\tilde{u}_1|^\nu > 1.$$

Thus, on the support of the right-hand side two displays prior, it holds that

$$\mathbb{1}_{\{(\delta \vartheta^{1/2})^{-\nu} n^{-\nu(1/2-\zeta)} |\tilde{u}_1|^\nu > 1\}} \leq (\delta \vartheta^{1/2})^{-\nu} n^{-\nu(1/2-\zeta)} |\tilde{u}_1|^\nu.$$

Hence, by the monotonicity and linearity of expectations,

$$\begin{aligned} I_{1b} &= E[\tilde{u}_1^2 \mathbb{1}_{\{|\tilde{u}_1| > \delta \vartheta^{1/2} n^{1/2-\zeta}\}}] \\ &\leq E[\tilde{u}_1^2 (\delta \vartheta^{1/2})^{-\nu} n^{-\nu(1/2-\zeta)} |\tilde{u}_1|^\nu] = (\delta \vartheta^{1/2})^{-\nu} n^{-\nu(1/2-\zeta)} E[|\tilde{u}_1|^{2+\nu}]. \end{aligned}$$

Condition (6) of the present theorem stipulates that  $E[|\tilde{u}_1|^{2+\nu}]$  is of constant order, hence,

$$\lim_{n \rightarrow \infty} I_{1b} \lesssim \lim_{n \rightarrow \infty} (\delta \vartheta^{1/2})^{-\nu} n^{-\nu(1/2-\zeta)} = 0$$

by the specification that  $\nu > 0$  and  $\zeta < 1/2$  in Condition (6) of the present theorem.

To show as much for  $I_2$ , observe that

$$\begin{aligned} \lim_{n \rightarrow \infty} I_2 &= \lim_{n \rightarrow \infty} E[\mathbb{1}\{|\langle \theta_j, \mathbf{d}_i \rangle| > n^\zeta\} E[\tilde{u}_1^2 | \mathbf{z}_i]] \\ &\lesssim \lim_{n \rightarrow \infty} P\{|\langle \theta_j, \mathbf{d}_i \rangle| > n^\zeta\} = 0 \end{aligned}$$

by Condition (a) of the present theorem. This concludes the demonstration of Lindeberg's condition. It follows that  $\tilde{Z}_{j,n} \rightsquigarrow \mathcal{N}(0, 1)$ . To show as much for  $Z_{j,n}$  and hence for  $\sqrt{n}(\tilde{\beta}_j - \beta_j)/\omega_j$ , it suffices to show that  $w_j \sigma_u / \omega_j \rightarrow_p 1$ . Note that  $\omega_j = \sigma_u \Theta_{jj}$  and hence that  $w_j \sigma_u / \omega_j = w_j / \Theta_{jj}$ . Since  $\Theta_{jj}$  is bounded strictly away from zero uniformly in  $n$ , we have  $|w_j / \Theta_{jj} - 1| = |w_j - \Theta_{jj}| / \Theta_{jj}$  and hence that it suffices to show that  $|w_j - \Theta_{jj}| = o_p(1)$ . But, as we established above, we have for arbitrary  $\epsilon > 0$

$$P\{|w_j - \Theta_{jj}| > \epsilon\} \leq P\{\|\bar{\Sigma}_d - \Sigma_d\|_\infty > \epsilon / m_\Theta^2\} = o(1)$$

by Condition (1) of the present theorem. Thus  $w_j \sigma_u / \omega_j \rightarrow_p 1$  hence  $Z_{j,n} \rightsquigarrow \mathcal{N}(0, 1)$  under each Condition (5) and (6) of the present theorem.

It remains to show that the limit holds when  $\omega_j$  is replaced by an estimator  $\hat{\omega}_j$  that satisfies  $|\hat{\omega}_j - \omega_j| = o_p(1)$ . Suppose that  $\hat{\omega}_j$  is such an estimator. We claim that  $\hat{\omega}_j / \omega_j - 1 = o_p(1)$ . To see as much, note first that, from the hypotheses of the present theorem,  $\omega_j = \sigma_u \Theta_{jj}$  is bounded strictly away from zero uniformly in  $n$ . It then follows that  $|\hat{\omega}_j / \omega_j - 1| = |\hat{\omega}_j - \omega_j| / \omega_j = o_p(1)$  by the specification of  $\hat{\omega}_j$ . By the continuous mapping theorem, it holds that  $\omega_j / \hat{\omega}_j \rightarrow_p 1$  and hence by Slutsky's Lemma that

$$\sqrt{n}(\tilde{\beta}_j - \beta_j) / \hat{\omega}_j = (\omega_j / \hat{\omega}_j) Z_{j,n} \rightsquigarrow \mathcal{N}(0, 1),$$

as claimed.  $\square$

**Proof of Lemma 3.5.** We first show that, under the conditions of the present Lemma,

$$\|\bar{\Sigma}_d - \Sigma_d\|_\infty = o_p(1),$$

thereby demonstrating Condition (1). Note that Lemma C.5 gives

$$P\{\|\bar{\Sigma}_d - \Sigma_d\|_\infty > a\sqrt{(\log(p_x \vee n))/n}\} \leq 2(p_x \vee n)^{2-a^2/(6e^2\kappa^2)}$$

for  $n$  sufficiently large, where  $a$  is a controlled constant and  $\kappa = m_A^2(2\tau_z^2 + \|\Sigma_z\|_\infty / \log 2)$ . Condition (1) then follows given that  $\sqrt{(\log p_x)/n} = o(1)$  and by choosing  $a = O(1)$  large enough so that  $a^2/(6e^2\kappa^2) > 2$ .

We now show that, under the conditions of the present Lemma,  $P\{|\langle \theta_j, \mathbf{d}_i \rangle| > n^\zeta\} = o(1)$  for  $0 < \zeta < 1/2$  and hence that Condition (a) is satisfied. First, write

$$P\{|\langle \theta_j, \mathbf{d}_i \rangle| > n^\zeta\} \leq P\{\|\theta_j\|_1 \|\mathbf{d}_i\|_\infty > n^\zeta\} \leq P\{m_\Theta \|\mathbf{d}_i\|_\infty > n^\zeta\}.$$

Next, (van der Vaart and Wellner, 1996, Chapters 2.1.3 and 2.2, Pages 90–91 and 95–97) and the proof of Lemma C.5 to infer that

$$\|\max_{j \in [p_x]} |d_{ij}| \|_{\psi_2} \leq C_2 \sqrt{\log p_x + 1} \max_{j \in [p_x]} \|d_{ij}\|_{\psi_2} \leq C_2 \tau_z m_A \sqrt{\log p_x + 1},$$

where  $C_2$  is an absolute constant. The exponential Markov bound then yields

$$P\{\|\mathbf{d}_i\|_\infty > t\} \leq \exp\left(1 - \frac{C_3 t^2}{C_2 \tau_z m_A \sqrt{\log p_x + 1}}\right)$$

for  $t \geq 0$ , where  $C_3$  is an absolute constant. Combine the above result with that of three displays previous, choose  $t = n^\zeta / m_\Theta$  and cite the growth conditions of the present lemma to conclude that

$$P\{|\langle \theta_j, \mathbf{d}_i \rangle| > n^\zeta\} \lesssim \exp\left(-\frac{n^{2\zeta}}{\sqrt{\log p_x}}\right) \xrightarrow{n \rightarrow \infty} 0,$$

as required for Condition (a).  $\square$

## Appendix B. Materials required for Section 4

### B.1. Materials required for Section 4.2

Our guarantees for estimating  $\mathbf{A}$  and  $\beta$  consist of two parts. The first is the *oracle inequality*, which bounds the  $\ell_1$  estimation error of a generic Lasso estimator conditional on the occurrence of a special set  $\mathcal{T}$ . The oracle inequality is a fixture of the  $\ell_1$  regularized estimation literature; see for instance (Bühlmann and van de Geer, 2011, Chapter 6). We present it for the sake of completeness. The oracle inequality itself is specific to neither the first- nor the second- stage

estimators of the present work. Indeed, we require the result to derive bounds for both estimators. As such, we present the theorem in terms of a generic model that shares notation with neither the first- nor second- stage models described in Section 2.1 except for the number of observations  $n$ .

**Theorem B.1** (Oracle Inequality). Consider the generic linear model

$$\mathbf{g} = \mathbf{W}\boldsymbol{\gamma} + \mathbf{h},$$

where  $\mathbf{g} \in \mathbb{R}^n$  is a vector of univariate responses,  $\mathbf{W} \in \mathbb{R}^{n \times p}$  is a design matrix with rows  $\mathbf{w}_i$ ,  $\boldsymbol{\gamma} \in \mathbb{R}^n$  is a noise vector with arbitrary distribution. Let  $\hat{\boldsymbol{\gamma}}$  denote the Lasso estimator given by

$$\hat{\boldsymbol{\gamma}} \in \operatorname{argmin}_{\mathbf{a} \in \mathbb{R}^p} \left\{ \|\mathbf{g} - \mathbf{W}\mathbf{a}\|_2^2 / (2n) + r \|\mathbf{a}\|_1 \right\}.$$

Let  $S_{\boldsymbol{\gamma}} := \operatorname{supp} \boldsymbol{\gamma}$ , and let  $s_{\boldsymbol{\gamma}} := |S_{\boldsymbol{\gamma}}|$ . Suppose that  $\mathbf{W}$  satisfies the compatibility Condition with respect to the index set  $S_{\boldsymbol{\gamma}}$  and compatibility constant  $\phi_{\boldsymbol{\gamma}} > 0$ . Then, on the set

$$\mathcal{T}(r) := \left\{ 4 \|\mathbf{W}^{\top} \mathbf{h} / n\|_{\infty} \leq r \right\}, \quad (\text{B.1})$$

the bound

$$\|\mathbf{W}(\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma})\|_2^2 / n + r \|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}\|_1 \leq 4s_{\boldsymbol{\gamma}} r^2 / \phi_{\boldsymbol{\gamma}}^2$$

holds.

**Proof of Theorem B.1.** The proof is algebra. See Theorem 6.1 of Bühlmann and van de Geer (2011).  $\square$

Theorem B.1 provides a deterministic guarantee for the  $\ell_1$  estimation error of a generic Lasso estimator  $\hat{\boldsymbol{\gamma}}$  on the set  $\mathcal{T}(r) = \{4 \|\mathbf{W}^{\top} \mathbf{h} / n\|_{\infty} \leq r\}$ . Consequently, it holds that

$$\mathbb{P}\{\|\hat{\boldsymbol{\gamma}} - \boldsymbol{\gamma}\|_1 > 4s_{\boldsymbol{\gamma}} r / \phi_{\boldsymbol{\gamma}}^2\} \leq \mathbb{P}\mathcal{T}(r)^c.$$

The quantity  $\|\mathbf{W}^{\top} \mathbf{h} / n\|_{\infty}$  is sometimes called the *empirical process term*; for instance, (Bühlmann and van de Geer, 2011, Chapter 6). Thus, upper bounds for  $\mathbb{P}\mathcal{T}(r)^c$  yield probabilistic guarantees for the  $\ell_1$  estimation error. We provide such bounds in Section C of the main Appendix.

## B.2. Materials required for Section 4.2.2

Recall that the model for  $\mathbf{x}^j$  is given by

$$\mathbf{x}^j = \mathbf{Z}\boldsymbol{\alpha}^j + \mathbf{v}^j,$$

where  $\mathbf{v}^j$  has nontrivial covariance with the noise  $\mathbf{u}$ . It suffices for our purposes to take a naïve approach to bounding the quantity  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} = \max_{j \in [p_{\mathbf{x}}]} \|\hat{\boldsymbol{\alpha}}^j - \boldsymbol{\alpha}^j\|_1$ . That is, we simultaneously bound the estimation error of each individual task. One could use a more complex approach such as Liu et al. (2015) to treat different patterns of joint sparsity amongst the first-stage regression vectors. We make the following assumption.

The following lemma provides finite-sample guarantees for  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$  under the choice of tuning parameters in Definition 4.4.

The following generic bound for  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$  can be combined with concentration results for specific distributions of the first-stage noise elements. We present it separately for the sake of modularity with respect to such assumptions.

**Lemma B.2** (Generic Bound for  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$ ). Suppose that Assumption 4.3 holds. For each  $j \in [p_{\mathbf{x}}]$ , let the sets  $\mathcal{T}_{\mathbf{v}^j}(\lambda)$  be as defined in Lemma C.3. It then holds on the set  $\bigcap_{j \in [p_{\mathbf{x}}]} \mathcal{T}_{\mathbf{v}^j}(r_j)$  that

$$\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \leq 4s_{\mathbf{A}} r_{\mathbf{A}} / \phi_{\mathbf{A}}^2 = 4 \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}},$$

where  $r_j$  is the tuning parameter for the respective first-stage Lasso problem,  $r_{\mathbf{A}}$  and  $s_{\mathbf{A}}$  are as defined in Section 2 and  $\phi_{\mathbf{A}}$  is as defined in Assumption 4.3.

**Proof of Lemma B.2.** For each  $j \in [p_{\mathbf{x}}]$ , Bühlmann and van de Geer (2011, Theorem 6.1, Lemma 6.2) entails that

$$\begin{aligned} \mathcal{T}_{\mathbf{v}^j}(r_j) &= \left\{ 4 \|\mathbf{Z}^{\top} \mathbf{v}^j / n\|_{\infty} \leq r_j \right\} \subseteq \left\{ \|\hat{\boldsymbol{\alpha}}^j - \boldsymbol{\alpha}^j\|_1 \leq 4s_{\boldsymbol{\alpha}^j} r_j / \phi_{\boldsymbol{\alpha}^j}^2 \right\} \\ &\subseteq \left\{ \|\hat{\boldsymbol{\alpha}}^j - \boldsymbol{\alpha}^j\|_1 \leq 4s_{\mathbf{A}} r_{\mathbf{A}} / \phi_{\mathbf{A}}^2 \right\}, \end{aligned}$$

where the latter containment follows by specification of  $s_{\mathbf{A}}$ ,  $r_{\mathbf{A}}$ , and  $\phi_{\mathbf{A}}$ . Take intersections over  $j \in [p_{\mathbf{x}}]$  on both sides of the above display to conclude.  $\square$

**Lemma B.3** (Bound for  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$ ). Suppose that [Assumptions 2.4](#) and [4.3](#) hold. Set  $r_j$  according to [Definition 4.4](#). Then,

$$P\{\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} > 4 \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} c_{\mathbf{V}} (\|\hat{\Sigma}_{\mathbf{Z}}\|_{\infty} (\log p_{\mathbf{Z}}/n)^{1/2}) \} \leq e p_{\mathbf{Z}}^{2 - C_0 \min_{j \in [p_{\mathbf{X}}]} \{c_j^2 / \|\mathbf{v}_j\|_{\psi_2}^2\}},$$

where  $c_{\mathbf{V}}$  is as specified in [Definition 4.4](#) and  $C_0$  is as defined in [Lemma C.1](#).

**Proof of Lemma B.3.** [Lemma B.2](#) entails that

$$P\{\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} > 4 s_{\mathbf{A}} r_{\mathbf{A}} / \phi_{\mathbf{A}}^2\} \leq P\left\{\bigcap_{j \in [p_{\mathbf{X}}]} \mathcal{T}_{\mathbf{w}}(r_j)\right\}^c.$$

Apply the estimate of [Lemma C.3](#) for the right-hand side to conclude.  $\square$

Note that, in order to obtain a rate of convergence, we must choose each  $c_j$  so that the quantity  $\min_{j \in [p_{\mathbf{X}}]} c_j^2 C_0 / \tau_{\mathbf{w}}^2 - 2$  is bounded strictly away from zero. Such a task may not be feasible in practice. Our empirical results, which we present in [Section 5](#), suggest that cross-validated choices of the tuning parameters for the first- and second-stage Lasso estimators suffice for good behavior of the resultant updated estimator  $\hat{\beta}$ . For the sake of our theory, we assume that such appropriate choices of  $c_j$  have been made. Given such an assumption, [Lemma B.3](#) entails that  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} = O_p(s_{\mathbf{A}} \sqrt{\log(p_{\mathbf{Z}}/n)})$ , essentially identical to the Lasso rate for single task regression problems.

The following bound is required for [Lemma 4.5](#)

**Lemma B.4** (Control of  $\|\hat{\mathbf{D}}^{\top} \tilde{\mathbf{u}}/n\|_{\infty}$ ). Let  $\tilde{\mathbf{u}} = \mathbf{u} + [(\mathbf{D} - \hat{\mathbf{D}}) + \mathbf{V}]\beta$ . Then,

$$\begin{aligned} \|\hat{\mathbf{D}}^{\top} \tilde{\mathbf{u}}/n\|_{\infty} &\leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\hat{\Sigma}_{\mathbf{Z}}\|_{\infty} (\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\beta\|_1 + \|\mathbf{A}\|_{L_1}) \\ &\quad + (\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\mathbf{A}\|_{L_1}) (\|\mathbf{Z}^{\top} \mathbf{V}/n\|_{\infty} \|\beta\|_1 + \|\mathbf{Z}^{\top} \mathbf{u}/n\|_{\infty}). \end{aligned}$$

**Proof of Lemma B.4.** Write  $\hat{\mathbf{D}}^{\top} = (\hat{\mathbf{D}} - \mathbf{D})^{\top} + \mathbf{D}^{\top}$  to find that

$$\begin{aligned} \|\hat{\mathbf{D}}^{\top} \tilde{\mathbf{u}}/n\|_{\infty} &= \|[(\hat{\mathbf{D}} - \mathbf{D})^{\top} + \mathbf{D}^{\top}][(\mathbf{D} - \hat{\mathbf{D}})\beta + (\mathbf{V}\beta + \mathbf{u})]/n\|_{\infty} \\ &\leq \|(\hat{\mathbf{D}} - \mathbf{D})^{\top}(\mathbf{D} - \hat{\mathbf{D}})\beta/n\|_{\infty} + \|(\hat{\mathbf{D}} - \mathbf{D})^{\top}(\mathbf{V}\beta + \mathbf{u})/n\|_{\infty} \\ &\quad + \|\mathbf{D}^{\top}(\mathbf{D} - \hat{\mathbf{D}})\beta/n\|_{\infty} + \|\mathbf{D}^{\top}(\mathbf{V}\beta + \mathbf{u})/n\|_{\infty} \\ &:= I_1 + I_2 + I_3 + I_4 \end{aligned} \tag{B.2}$$

We treat each quantity in the right-hand side above in turn.

For  $I_1$ , write

$$I_1 = \|(\hat{\mathbf{D}} - \mathbf{D})^{\top}(\mathbf{D} - \hat{\mathbf{D}})\beta/n\|_{\infty} \leq \|(\hat{\mathbf{D}} - \mathbf{D})^{\top}(\hat{\mathbf{D}} - \mathbf{D})/n\|_{\infty} \|\beta\|_1.$$

Recall that  $\hat{\mathbf{D}} - \mathbf{D} = \mathbf{Z}(\hat{\mathbf{A}} - \mathbf{A})$  and write

$$\begin{aligned} \|(\hat{\mathbf{D}} - \mathbf{D})^{\top}(\hat{\mathbf{D}} - \mathbf{D})/n\|_{\infty} &= \|(\hat{\mathbf{A}} - \mathbf{A})^{\top} \hat{\Sigma}_{\mathbf{Z}}(\hat{\mathbf{A}} - \mathbf{A})\|_{\infty} \\ &\leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2 \|\hat{\Sigma}_{\mathbf{Z}}\|_{\infty} = \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2 \|\hat{\Sigma}_{\mathbf{Z}}\|_{\infty}, \end{aligned}$$

where the second line follows from repeated application of Hölder's inequality. Combine the two previous displays to conclude that

$$I_1 \leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2 \|\hat{\Sigma}_{\mathbf{Z}}\|_{\infty} \|\beta\|_1. \tag{B.3}$$

For  $I_2$ , write

$$\begin{aligned} I_2 &= \|(\hat{\mathbf{D}} - \mathbf{D})^{\top}(\mathbf{V}\beta + \mathbf{u})/n\|_{\infty} \\ &= \|(\hat{\mathbf{A}} - \mathbf{A})^{\top} \mathbf{Z}^{\top}(\mathbf{V}\beta + \mathbf{u})/n\|_{\infty} \\ &\leq \underbrace{\|(\hat{\mathbf{A}} - \mathbf{A})^{\top} \mathbf{Z}^{\top} \mathbf{V}\beta/n\|_{\infty}}_{I_{2,a}} + \underbrace{\|(\hat{\mathbf{A}} - \mathbf{A})^{\top} \mathbf{Z}^{\top} \mathbf{u}/n\|_{\infty}}_{I_{2,b}}. \end{aligned}$$

Applications of Hölder's inequality yield

$$I_{2,a} \leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\mathbf{Z}^{\top} \mathbf{V}/n\|_{\infty} \|\beta\|_1$$

and

$$I_{2,b} \leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\mathbf{Z}^{\top} \mathbf{u}/n\|_{\infty} = \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\mathbf{Z}^{\top} \mathbf{u}/n\|_{\infty}.$$

Combine the previous three displays to conclude that

$$I_2 \leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} (\|\mathbf{Z}^{\top} \mathbf{V}/n\|_{\infty} \|\beta\|_1 + \|\mathbf{Z}^{\top} \mathbf{u}/n\|_{\infty}), \tag{B.4}$$

For the quantity  $I_3$ , write

$$I_3 = \|\mathbf{D}^\top(\widehat{\mathbf{D}} - \mathbf{D})\beta/n\|_\infty \leq \|\mathbf{D}^\top(\widehat{\mathbf{D}} - \mathbf{D})/n\|_\infty \|\beta\|_1$$

and observe that

$$\mathbf{D}^\top(\widehat{\mathbf{D}} - \mathbf{D})/n = \mathbf{A}^\top \mathbf{Z}^\top \mathbf{Z}(\widehat{\mathbf{A}} - \mathbf{A})/n = \mathbf{A}^\top \widehat{\Sigma}_z(\widehat{\mathbf{A}} - \mathbf{A}),$$

which yields

$$\begin{aligned} \|\mathbf{D}^\top(\widehat{\mathbf{D}} - \mathbf{D})/n\|_\infty &= \|\mathbf{A}^\top \widehat{\Sigma}_z(\widehat{\mathbf{A}} - \mathbf{A})\|_\infty \\ &\leq \|\mathbf{A}\|_{L_1} \|\widehat{\Sigma}_z\|_\infty \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} \end{aligned}$$

after repeated application of Hölder's inequality. Conclude from the previous three displays that

$$I_3 \leq \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\widehat{\Sigma}_z\|_\infty \|\mathbf{A}\|_{L_1}. \quad (\text{B.5})$$

For the quantity  $I_4$ , write

$$\begin{aligned} I_4 &= \|\mathbf{D}^\top(\mathbf{V}\beta + \mathbf{u})/n\|_\infty \\ &= \|\mathbf{A}^\top \mathbf{Z}^\top(\mathbf{V}\beta + \mathbf{u})/n\|_\infty \\ &\leq \|\mathbf{A}\|_{L_1} (\|\mathbf{Z}^\top \mathbf{V}\beta/n\|_\infty + \|\mathbf{Z}^\top \mathbf{u}/n\|_\infty) \\ &\leq (\|\mathbf{Z}^\top \mathbf{V}/n\|_\infty \|\beta\|_1 + \|\mathbf{Z}^\top \mathbf{u}/n\|_\infty) \|\mathbf{A}\|_{L_1}. \end{aligned} \quad (\text{B.6})$$

The original claim follows from line (B.2) and lines (B.3)–(B.6).  $\square$

The following generic bound for  $\|\hat{\beta} - \beta\|_1$  can be combined with concentration results for specific distributions of the first- and second-stage noise elements. We present it separately for the sake of modularity with respect to such assumptions.

**Lemma B.5** (Generic Bound for  $\|\hat{\beta} - \beta\|_1$ ). Suppose that [Assumption 4.3](#) holds. Let  $\lambda_v, \lambda_u > 0$  be arbitrary. Set  $r_\beta$  according to [Definition 4.4](#). Then, on the set  $\mathcal{T}_v(\mathbf{r}) \cap \mathcal{T}_v(\lambda_v) \cap \mathcal{T}_u(\lambda_u)$ , where  $\mathcal{T}_v$  and  $\mathcal{T}_u$  are defined as in [Lemmas C.2 and C.3](#), respectively, and  $\mathbf{r}$  is the tuple of first-stage tuning parameters, we have

$$\|\hat{\beta} - \beta\|_1 \leq 4 \frac{s_\beta}{\phi_\beta^2} r_\beta.$$

**Proof of Lemma B.5.** By [Theorem B.1](#), we have

$$\mathcal{T}_{\tilde{\mathbf{u}}} = \{4\|\widehat{\mathbf{D}}^\top \tilde{\mathbf{u}}/n\|_\infty \leq r\} \subseteq \{\|\hat{\beta} - \beta\|_1 \leq 4s_\beta r/\phi^2\}.$$

It therefore suffices to show that  $\mathcal{T}_v(\mathbf{r}) \cap \mathcal{T}_u(\lambda_u) \subseteq \mathcal{T}_{\tilde{\mathbf{u}}}(r_\beta)$  for the present choice of  $r_\beta$ . [Lemma B.4](#) gives the bound

$$\begin{aligned} \|\widehat{\mathbf{D}}^\top \tilde{\mathbf{u}}/n\|_\infty &\leq \underbrace{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\widehat{\Sigma}_z\|_\infty (m_\beta \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} + m_A)}_{I_1} \\ &\quad + \underbrace{(\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} + m_A)(\|\mathbf{Z}^\top \mathbf{V}/n\|_\infty \|\beta\|_1 + \|\mathbf{Z}^\top \mathbf{u}/n\|_\infty)}_{I_2}. \end{aligned}$$

Cite [Lemma B.2](#) to conclude that, on the set  $\mathcal{T}_v(\mathbf{r})$ ,

$$I_1 \leq 4 \frac{s_A}{\phi_A^2} r_A \|\widehat{\Sigma}_z\|_\infty (4m_\beta \frac{s_A}{\phi_A^2} r_A + m_A).$$

Note that, on the set  $\mathcal{T}_v(\mathbf{r}) \cap \mathcal{T}_v(\lambda_v) \cap \mathcal{T}_u(\lambda_u)$ ,

$$I_2 \leq \frac{1}{4} (4 \frac{s_A}{\phi_A^2} r_A + m_A) (m_\beta \lambda_v + \lambda_u)$$

by specification. Multiply the two previous displays by 4 and combine with the third-previous display to conclude that  $\mathcal{T}_v(\mathbf{r}) \cap \mathcal{T}_v(\lambda_v) \cap \mathcal{T}_u(\lambda_u) \subseteq \mathcal{T}_{\tilde{\mathbf{u}}}(r_\beta)$  for the present choice of  $r_\beta$ , as required.  $\square$

**Proof of Lemma 4.5.** Note first that

$$\mathcal{T}_v(\lambda_v) = \bigcap_{j=1}^{p_x} \{4\|\mathbf{Z}^\top \mathbf{v}^j/n\|_\infty > r_j\} \subseteq \mathcal{T}_v(\mathbf{r}).$$



Let  $\mathcal{T}_{JJ}$  be the set on which [Assumption 4.3](#) holds. [Lemma B.5](#) then entails that

$$\{\|\hat{\beta} - \beta\|_1 > 4 \frac{s\beta}{\phi_\beta^2} r_\beta\} \cap \mathcal{T}_{JJ} \subseteq (\mathcal{T}_V(\mathbf{r}) \cap \mathcal{T}_U(\lambda_U))^c = \mathcal{T}_V(\mathbf{r})^c \cup \mathcal{T}_U(\lambda_U)^c.$$

Thus,

$$\mathbb{P}\{\|\hat{\beta} - \beta\|_1 > 4 \frac{s\beta}{\phi_\beta^2} r_\beta\} \leq \mathbb{P} \mathcal{T}_V(\mathbf{r})^c + \mathbb{P} \mathcal{T}_U(\lambda_U)^c + t_n.$$

Now substitute the present choices of tuning parameters and cite the estimates of [Lemmas C.2](#) and [C.3](#).  $\square$

**Lemma B.6** (Second-Stage Compatibility Constant). *Let  $S \subseteq [p]$  be an arbitrary index set with  $s = |S|$ . For a given matrix  $\mathbf{M} \in \mathbb{R}^{n \times p}$ , define the quantity*

$$\phi^2(\mathbf{M}, S) = \inf_{\delta \in \mathcal{C}(S)} \frac{s \|\mathbf{M}\delta\|_2^2}{n \|\delta_S\|_1^2}.$$

Let  $\epsilon_1, \epsilon_2 > 0$  be arbitrary. Then,

$$\begin{aligned} \mathbb{P}\{\phi^2(\hat{\mathbf{D}}, S) < \Lambda_{\min}(\Sigma_d) - \epsilon_2 - \epsilon_1\} \\ \leq \mathbb{P}\{16s(2m_A \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2) \|\hat{\Sigma}_Z\|_\infty > \epsilon_1\} \\ + \mathbb{P}\{16s \|\bar{\Sigma}_d - \Sigma_d\|_\infty > \epsilon_2\}, \end{aligned}$$

where  $\bar{\Sigma}_d = \mathbf{D}^\top \mathbf{D}/n$ .

**Proof of Lemma B.6.** Let  $S, s$  be as in the statement of [Lemma B.6](#), and let  $\delta \in \mathbb{R}^{px} \setminus \{\mathbf{0}\}$  satisfying  $\|\delta_{S^c}\|_1 \leq 3\|\delta_S\|_1$  be arbitrary. Write  $\hat{\mathbf{D}} = \mathbf{D} + (\hat{\mathbf{D}} - \mathbf{D})$ , so that

$$\begin{aligned} \|\hat{\mathbf{D}}\delta\|_2^2 &= \|[\mathbf{D} + (\hat{\mathbf{D}} - \mathbf{D})]\delta\|_2^2 \\ &= \langle [\mathbf{D} + (\hat{\mathbf{D}} - \mathbf{D})]\delta, [\mathbf{D} + (\hat{\mathbf{D}} - \mathbf{D})]\delta \rangle \\ &= \|\mathbf{D}\delta\|_2^2 + 2\langle \mathbf{D}^\top (\hat{\mathbf{D}} - \mathbf{D})\delta, \delta \rangle + \langle (\hat{\mathbf{D}} - \mathbf{D})^\top (\hat{\mathbf{D}} - \mathbf{D})\delta, \delta \rangle. \end{aligned}$$

Thus,

$$\begin{aligned} \|\hat{\mathbf{D}}\delta\|_2^2/n &\geq \|\mathbf{D}\delta\|_2^2/n - \underbrace{2|\langle \mathbf{D}^\top (\hat{\mathbf{D}} - \mathbf{D})\delta/n, \delta \rangle|}_{I_1} \\ &\quad - \underbrace{|\langle (\hat{\mathbf{D}} - \mathbf{D})^\top (\hat{\mathbf{D}} - \mathbf{D})\delta/n, \delta \rangle|}_{I_2}. \end{aligned} \tag{B.7}$$

We now obtain bounds for the quantities  $I_1, I_2$ . From repeated applications of Hölder's inequality, write

$$\begin{aligned} I_1 &\lesssim |\langle \mathbf{D}^\top (\hat{\mathbf{D}} - \mathbf{D})\delta/n, \delta \rangle| \leq \|\mathbf{D}^\top (\hat{\mathbf{D}} - \mathbf{D})/n\|_\infty \|\delta\|_1^2 \\ &= \|\mathbf{A}^\top \mathbf{Z}^\top \mathbf{Z} (\hat{\mathbf{A}} - \mathbf{A})/n\|_\infty \|\delta\|_1^2 \\ &\leq \|\mathbf{A}\|_{L_1} \|\mathbf{Z}^\top \mathbf{Z}/n\|_\infty \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\delta\|_1^2 \\ &\leq m_A \|\hat{\Sigma}_Z\|_\infty \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\delta\|_1^2 \end{aligned}$$

and

$$\begin{aligned} I_2 &= |\langle (\hat{\mathbf{D}} - \mathbf{D})^\top (\hat{\mathbf{D}} - \mathbf{D})\delta/n, \delta \rangle| \\ &\leq \|(\hat{\mathbf{D}} - \mathbf{D})^\top (\hat{\mathbf{D}} - \mathbf{D})/n\|_\infty \|\delta\|_1^2 \\ &= \|(\hat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} (\hat{\mathbf{A}} - \mathbf{A})/n\|_\infty \|\delta\|_1^2 \\ &\leq \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\mathbf{Z}^\top \mathbf{Z}/n\|_\infty \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\delta\|_1^2 \\ &\leq \|\hat{\Sigma}_Z\|_\infty \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2 \|\delta\|_1^2. \end{aligned}$$

Combine the previous two displays with [\(B.7\)](#) to find that

$$\|\hat{\mathbf{D}}\delta\|_2^2/n \geq \|\mathbf{D}\delta\|_2^2/n - (2m_A \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2) \|\hat{\Sigma}_Z\|_\infty \|\delta\|_1^2.$$

By assumption, we have  $\|\delta\|_1 \leq 4\|\delta_S\|_1$ . Substitute this expression in the right-hand side above and multiply through by  $s/\|\delta_S\|_1^2$  to obtain

$$\frac{s \|\hat{\mathbf{D}}\delta\|_2^2}{n \|\delta_S\|_1^2} \geq \frac{s \|\mathbf{D}\delta\|_2^2}{n \|\delta_S\|_1^2} - 16s(2m_A \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2) \|\hat{\Sigma}_Z\|_\infty.$$

Thus, on the set  $\{16s(2m_A\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2)\|\widehat{\Sigma}_z\|_\infty \leq \epsilon_1\}$ , we have

$$\begin{aligned} \frac{s\|\widehat{\mathbf{D}}\delta\|_2^2}{n\|\delta_S\|_1^2} &\geq \frac{s\|\mathbf{D}\delta\|_2^2}{n\|\delta_S\|_1^2} - \epsilon_1 \\ &= \left( \frac{s\delta^\top \Sigma_d \delta}{\|\delta_S\|_1^2} - \frac{s\delta^\top (\bar{\Sigma}_d - \Sigma_d) \delta}{\|\delta_S\|_1^2} \right) - \epsilon_1 \\ &\geq \frac{s\delta^\top \Sigma_d \delta}{\|\delta_S\|_1^2} - \frac{s\|\bar{\Sigma}_d - \Sigma_d\|_\infty \|\delta\|_1^2}{\|\delta_S\|_1^2} - \epsilon_1 \\ &\geq \frac{s\delta^\top \Sigma_d \delta}{\|\delta_S\|_1^2} - 16s\|\bar{\Sigma}_d - \Sigma_d\|_\infty - \epsilon_1. \end{aligned}$$

From Cauchy-Schwartz we have  $\|\delta_S\|_1 \leq \sqrt{s}\|\delta_S\|_2$  and hence that  $\|\delta_S\|_1^2 \leq s\|\delta\|_2^2$ . Substitute this bound into the first term on the right-hand side above to obtain

$$\begin{aligned} \frac{s\|\widehat{\mathbf{D}}\delta\|_2^2}{n\|\delta_S\|_1^2} &\geq \frac{\delta^\top \Sigma_d \delta}{\|\delta\|_2^2} - 16s\|\bar{\Sigma}_d - \Sigma_d\|_\infty - \epsilon_1 \\ &\geq \Lambda_{\min}(\Sigma_d) - 16s\|\bar{\Sigma}_d - \Sigma_d\|_\infty - \epsilon_1, \end{aligned}$$

where  $\Lambda_{\min}(\Sigma_d)$  denotes the minimal eigenvalue of  $\Sigma_d$ . The right-hand side above does not depend on  $\delta$ , so we may take the infimum of the left-hand side above over  $\delta \in \mathcal{C}(S)$  to write, for any  $\epsilon_1, \epsilon_2 > 0$  as in the statement of the present Lemma,

$$\begin{aligned} \mathbb{P}\{\phi^2(\widehat{\mathbf{D}}, S) < \Lambda_{\min}(\Sigma_d) - \epsilon_2 - \epsilon_1\} \\ \leq \mathbb{P}\{16s(2m_A\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2)\|\widehat{\Sigma}_z\|_\infty > \epsilon_1\} \\ + \mathbb{P}\{16s\|\bar{\Sigma}_d - \Sigma_d\|_\infty > \epsilon_2\}, \end{aligned}$$

as claimed.  $\square$

**Lemma B.6** may be combined with results for the maximum first-stage estimation error  $\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$  and the maximum entry-wise difference  $\|\bar{\Sigma}_d - \Sigma_d\|_\infty$  to obtain specific bounds for  $\phi^2(\widehat{\mathbf{D}}, S_\beta)$  under different error and design matrix regimes, such as in **Lemma 4.6**.

**Proof of Lemma 4.6.** Set  $\epsilon_1$  in the statement of **Lemma B.6** as

$$\epsilon_1 = 128s_\beta(m_A \frac{s_A}{\phi_A^2} r_A + 2 \frac{s_A^2}{\phi_A^4} r_A^2) \|\widehat{\Sigma}_z\|_\infty,$$

so that

$$\mathbb{P}\{16s_\beta(2m_A\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2)\|\widehat{\Sigma}_z\|_\infty > \epsilon_1\} \leq \mathbb{P}\{\|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} > 4 \frac{s_A}{\phi_A^2} r_A\}.$$

The present choices of tuning parameters, along with the estimates of **Lemma B.3** and **Lemma B.6**, entail that

$$\begin{aligned} \mathbb{P}\left\{\phi^2(\widehat{\mathbf{D}}, S_\beta) < \Lambda_{\min}(\Sigma_d) - \epsilon_2 - 128s_\beta(m_A \frac{s_A}{\phi_A^2} r_A + 2 \frac{s_A^2}{\phi_A^4} r_A^2) \|\widehat{\Sigma}_z\|_\infty\right\} \\ \leq ep_z^{2-C_0 \min_{j \in [p_X]} \{c_j^2 / \|\mathbf{v}^j\|_{\psi_2}^2\}} + \mathbb{P}\{16s\|\bar{\Sigma}_d - \Sigma_d\|_\infty > \epsilon_2\}. \end{aligned}$$

Cite the growth assumptions of the present lemma to observe that, for  $n$  sufficiently large, we have  $m_A \frac{s_A}{\phi_A^2} r_A + 2(\frac{s_A}{\phi_A^2})^2 r_A^2 \leq 3m_A \frac{s_A}{\phi_A^2} r_A$ , from which it follows that, for such  $n$ ,

$$\begin{aligned} \mathbb{P}\{\phi^2(\widehat{\mathbf{D}}, S_\beta) < \Lambda_{\min}(\Sigma_d) - \epsilon_2 - 384m_A s_\beta \frac{s_A}{\phi_A^2} c_V \|\widehat{\Sigma}_z\|_\infty^{3/2} \sqrt{(\log p_Z)/n}\} \\ \leq ep_z^{2-C_0 \min_{j \in [p_X]} \{c_j^2 / \|\mathbf{v}^j\|_{\psi_2}^2\}} + \mathbb{P}\{16s_\beta\|\bar{\Sigma}_d - \Sigma_d\|_\infty > \epsilon_2\}, \end{aligned}$$

Cite a slight modification of Lemma C.5 to conclude that

$$\begin{aligned} \mathbb{P}\{\phi^2(\widehat{\mathbf{D}}, S_\beta) < \Lambda_{\min}(\Sigma_d) - (a + 384m_A s_\beta \frac{s_A}{\phi_A^2} c_V \|\widehat{\Sigma}_z\|_\infty^{3/2}) \sqrt{(\log p_Z)/n}\} \\ \leq ep_z^{2-C_0 \min_{j \in [p_X]} \{c_j^2 / \|\mathbf{v}^j\|_{\psi_2}^2\}} + 2p_X^{2-a^2/(6e^2\kappa^2)}, \end{aligned}$$

where  $a > 0$  is a controlled quantity, as claimed.  $\square$

### B.3. Materials required for Section 4.3

Lemmas B.7, B.9, B.11, and B.13 provide finite-sample bounds for the quantities  $\|\mathbf{f}_\ell\|_\infty$  for  $\ell \in [4]$  that are generic over various noise regimes. We present them separately for the sake of modularity with respect to such assumptions. Lemmas B.8, B.10, B.12, and B.14 in turn provide specific rates for the  $\|\mathbf{f}_\ell\|_\infty$  under the sub-Gaussian noise regime of Assumption 2.4.

**Lemma B.7** (Control of  $\mathbf{f}_1$ ). Suppose that Assumption 4.3 and Conditions (2), (3) of Assumption 4.7 hold and that  $\widehat{\Theta}$  is chosen according to Assumption 4.8. Then, on the set  $\mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}}) \cap \mathcal{T}_{\Theta}(\mu)$ , where  $\lambda_{\mathbf{u}} > 0$  is arbitrary, the remainder term

$$\mathbf{f}_1 = (\widehat{\Theta} - \Theta)^\top \mathbf{D}^\top \mathbf{u} / \sqrt{n}$$

satisfies

$$\|\mathbf{f}_1\|_\infty \leq 2^{-q} \sqrt{n} m_{\mathbf{A}} c_q(m_{\Theta} \mu)^{1-q} s_{\Theta} \lambda_{\mathbf{u}},$$

where  $c_q$  is as in Lemma A.2.

**Proof of Lemma B.7.** Lemma A.2 entails that, on the set  $\mathcal{T}_{\Theta}(\mu)$ ,

$$\max_{j \in [p_{\mathbf{x}}]} \|\widehat{\theta}_j - \theta_j\|_1 \leq 2c_q(2m_{\Theta} \mu)^{1-q} s_{\Theta}.$$

We therefore find that

$$\begin{aligned} \|(\widehat{\Theta} - \Theta)^\top \mathbf{D}^\top \mathbf{u} / \sqrt{n}\|_\infty &\leq \sqrt{n} \|\widehat{\Theta} - \Theta\|_{L_1} \|\mathbf{D}^\top \mathbf{u} / n\|_\infty \\ &\leq \sqrt{n} \|\widehat{\Theta} - \Theta\|_{L_1} \|\mathbf{A}\|_{L_1} \|\mathbf{Z}^\top \mathbf{u} / n\|_\infty \\ &\leq \sqrt{n} m_{\mathbf{A}} \|\widehat{\Theta} - \Theta\|_{L_1} \|\mathbf{Z}^\top \mathbf{u} / n\|_\infty \\ &\leq 2\sqrt{n} m_{\mathbf{A}} c_q(2m_{\Theta} \mu)^{1-q} s_{\Theta} \|\mathbf{Z}^\top \mathbf{u} / n\|_\infty. \end{aligned}$$

On the set  $\mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}})$  we have  $\|\mathbf{Z}^\top \mathbf{u} / n\|_\infty \leq \lambda_{\mathbf{u}} / 4$ . From this bound and the previous display we conclude that, on the set  $\mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}}) \cap \mathcal{T}_{\Theta}(\mu)$ ,

$$\|(\Theta - \widehat{\Theta})^\top \mathbf{D}^\top \mathbf{u} / \sqrt{n}\|_\infty \leq 2^{-q} \sqrt{n} m_{\mathbf{A}} c_q(m_{\Theta} \mu)^{1-q} s_{\Theta} \lambda_{\mathbf{u}},$$

as claimed.  $\square$

**Lemma B.8** (Control of  $\mathbf{f}_1$ , sub-Gaussian Noise). Suppose that (i) Assumption 4.3 and Conditions (2), (3) of Assumption 4.7 hold and (ii) Assumption 2.4 holds. Choose  $\widehat{\Theta}$  according to Assumption 4.8. Then,

$$\mathbb{P} \left\{ \|\mathbf{f}_1\|_\infty > 2^{-q} m_{\mathbf{A}} c_q c_{\mathbf{u}}(m_{\Theta} \mu)^{1-q} s_{\Theta} \sqrt{m_{\mathbf{Z}} \log p_{\mathbf{Z}}} \right\} \leq e p_{\mathbf{Z}}^{1-c_{\mathbf{u}}^2 C_0 / \tau_{\mathbf{u}}^2} + \mathbb{P} \mathcal{T}_{\Theta}(\mu)^c,$$

where  $c_{\mathbf{u}}$  is as in Definition 4.4 and  $C_0$  is as defined in Lemma C.1. If Conditions (1), (4), (5), and (6)a of Assumption 4.7 also hold, then  $\|\mathbf{f}_1\|_\infty = o_{\mathbb{P}}(1)$ .

**Proof of Lemma B.8.** Lemma B.7 entails that

$$\mathbb{P} \left\{ \|\mathbf{f}_1\|_\infty > 2^{-q} \sqrt{n} m_{\mathbf{A}} c_q(m_{\Theta} \mu)^{1-q} s_{\Theta} \lambda_{\mathbf{u}} \right\} \leq \mathbb{P} \mathcal{T}_{\mathbf{u}}^c + \mathbb{P} \mathcal{T}_{\Theta}(\mu)^c.$$

Substitute  $\lambda_{\mathbf{u}}$  chosen according to Definition 4.4 into the display above and cite the estimate of Lemma C.2 to deduce the original claim.  $\square$

**Lemma B.9** (Control of  $\mathbf{f}_2$ ). Suppose that Assumption 4.3 and Condition (3) of Assumption 4.7 hold. Choose  $\widehat{\Theta}$  according to Assumption 4.8, set  $r_{\mathbf{A}}$  according to Definition 4.4, and let  $\lambda_{\mathbf{u}} > 0$  be arbitrary. Then, on the set  $\mathcal{T}_{\mathbf{V}}(\mathbf{r}) \cap \mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}}) \cap \mathcal{T}_{\Theta}(\mu)$ , the remainder term

$$\mathbf{f}_2 = \widehat{\Theta}^\top (\widehat{\mathbf{D}} - \mathbf{D})^\top \mathbf{u} / \sqrt{n}$$

satisfies

$$\|\mathbf{f}_2\|_\infty \leq \sqrt{n} m_{\Theta} \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} \lambda_{\mathbf{u}}$$

for  $n$  sufficiently large.

**Proof of Lemma B.9.** Observe that

$$\widehat{\Theta}^\top (\widehat{\mathbf{D}} - \mathbf{D})^\top \mathbf{u} / \sqrt{n} = \sqrt{n} \widehat{\Theta}^\top (\widehat{\mathbf{A}} - \mathbf{A})^\top (\mathbf{Z}^\top \mathbf{u} / n).$$

On the set  $\mathcal{T}_\Theta(\mu)$ , each row  $\theta$  is feasible for [Problem 3.2](#). Then,  $\|\hat{\theta}_j\|_1 \leq \|\theta_j\|_1$  for each  $j \in [p_x]$  by specification. [Lemmas B.2](#) and C.2 then entail that, on the set  $\mathcal{T}_A(\mathbf{r}) \cap \mathcal{T}_u(\lambda_u) \cap \mathcal{T}_\Theta(\mu)$ ,

$$\begin{aligned} \|\hat{\Theta}^\top(\hat{\mathbf{D}} - \mathbf{D})^\top \mathbf{u} / \sqrt{n}\|_\infty &\leq \max_{j,k \in [p_x]} \sqrt{n} \|\hat{\theta}_j\|_1 \|\hat{\alpha}^k - \alpha^k\|_1 \|\mathbf{Z}^\top \mathbf{u} / n\|_\infty \\ &\leq \sqrt{nm_\Theta} \frac{S_A}{\phi_A^2} r_A \lambda_u, \end{aligned}$$

as claimed.  $\square$

**Lemma B.10** (Control of  $\mathbf{f}_2$ , sub-Gaussian Noise). Suppose that (i) [Assumption 4.3](#) and Condition (3) of [Assumption 4.7](#) hold and (ii) [Assumption 2.4](#) holds. Choose  $\hat{\Theta}$ ,  $\mathbf{r} = (r_1, \dots, r_{p_x})$ , and  $r_A$  according to [Assumption 4.8](#). Then,

$$\begin{aligned} P\{\|\mathbf{f}_2\|_\infty > m_\Theta c_V c_u m_Z \frac{S_A}{\phi_A^2} \log p_Z / \sqrt{n}\} \\ \leq e p_Z^{2-c_V^2 c_0 / \tau_V^2} + e p_Z^{1-c_u^2 c_0 / \tau_u^2} + P\{\|\hat{\Sigma}_Z\|_\infty > m_Z\} + P\mathcal{T}_\Theta(\mu)^c, \end{aligned}$$

where  $c_V$  is as in [Definition 4.4](#),  $c_u$  is as in [Definition 4.4](#), and  $C_0$  is as defined in [Lemma C.1](#). Consequently, if Conditions (1), (4), (5), and (6)b of [Assumption 4.7](#) also hold, then  $\|\mathbf{f}_2\|_\infty = o_p(1)$ .

**Proof of Lemma B.10.** [Lemma B.9](#) entails that

$$P\left\{\|\mathbf{f}_2\|_\infty > \frac{1}{4} \sqrt{nm_\Theta} \frac{S_A}{\phi_A^2} r_A \lambda_u\right\} \leq P\mathcal{T}_A(\mathbf{r})^c + P\mathcal{T}_u(\lambda_u)^c + P\mathcal{T}_\Theta(\mu)^c.$$

Substitute the present choices of  $r_A$  and  $\lambda_u$  into the previous display and cite the estimates of [Lemmas C.2](#) and C.3 to deduce the original claim.  $\square$

**Lemma B.11** (Control of  $\mathbf{f}_3$ ). Suppose that [Assumptions 4.3](#) and [4.3](#) and Conditions (2) and (3) of [Assumption 4.7](#) hold. Choose  $\hat{\Theta}$  according to [Assumption 4.8](#); let  $\mathbf{r} = (r_1, \dots, r_{p_x}) > \mathbf{0}$ ,  $\lambda_u > 0$ , and  $\lambda_V > 0$  be arbitrary. Set

$$r_\beta = 16 \frac{S_A}{\phi_A^2} r_A \|\hat{\Sigma}_Z\|_\infty (4m_\beta \frac{S_A}{\phi_A^2} r_A + m_A) + (4 \frac{S_A}{\phi_A^2} r_A + m_A) (m_\beta \lambda_V + \lambda_u).$$

Then, on the set

$$\mathcal{T}_V(\mathbf{r}) \cap \mathcal{T}_u(\lambda_u) \cap \mathcal{T}_\Theta(\mu),$$

the remainder term

$$\mathbf{f}_3 = \hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})(\beta - \hat{\beta}) / \sqrt{n}$$

satisfies

$$\|\mathbf{f}_3\|_\infty \leq 8m_\Theta m_A \sqrt{n} (4 \frac{S_A}{\phi_A^2} r_A + \lambda_u) \frac{S_\beta}{\phi_\beta^2} r_\beta.$$

**Proof of Lemma B.11.** We first observe that

$$\|\hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})(\beta - \hat{\beta}) / \sqrt{n}\|_\infty \leq \sqrt{n} \|\hat{\Theta}\|_{L_1} \|\hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}}) / n\|_\infty \|\hat{\beta} - \beta\|_1.$$

Now,

$$\begin{aligned} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}}) / n &= \hat{\mathbf{D}}^\top (\mathbf{D} + \mathbf{u} - \hat{\mathbf{D}}) / n \\ &= \hat{\mathbf{D}}^\top (\mathbf{D} - \hat{\mathbf{D}}) / n + \hat{\mathbf{D}}^\top \mathbf{u} / n \\ &= \hat{\mathbf{A}}^\top (\mathbf{Z}^\top \mathbf{Z} / n) (\mathbf{A} - \hat{\mathbf{A}}) + \hat{\mathbf{D}}^\top \mathbf{u} / n. \end{aligned}$$

For the first term on the right-hand side above, write

$$\begin{aligned} \|\hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{Z}^\top \mathbf{Z} / n) (\mathbf{A} - \hat{\mathbf{A}})\|_\infty &\leq \|\hat{\mathbf{A}}^\top (\mathbf{Z}^\top \mathbf{Z} / n) (\mathbf{A} - \hat{\mathbf{A}})\|_\infty \\ &\quad + \|(\hat{\mathbf{A}} - \mathbf{A})^\top (\mathbf{Z}^\top \mathbf{Z} / n) (\mathbf{A} - \hat{\mathbf{A}})\|_\infty \\ &\leq \|\hat{\Sigma}_Z\|_\infty [m_A \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2]. \end{aligned}$$

On the set  $\mathcal{T}_V(\mathbf{r})$ , the right-hand side above is less than or equal to  $2m_A \cdot \|\hat{\Sigma}_Z\|_\infty \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1}$  for  $n$  sufficiently large by [Lemma B.2](#) and the hypotheses of the present lemma. For the second term on the right-hand side of two displays previous, write

$$\|\hat{\mathbf{D}}^\top \mathbf{u} / n\|_\infty = \|\hat{\mathbf{A}}^\top \mathbf{Z}^\top \mathbf{u} / n\|_\infty$$

$$\begin{aligned} &= \|(\mathbf{A} + [\hat{\mathbf{A}} - \mathbf{A}])^\top (\mathbf{Z}^\top \mathbf{u}/n)\|_\infty \\ &\leq 2m_{\mathbf{A}} \|\mathbf{Z}^\top \mathbf{u}/n\|_\infty, \end{aligned}$$

where the final line holds on the set  $\mathcal{T}_{\mathbf{V}}(\mathbf{r})$  for such  $n$  by Lemma B.2. Thus, on the set  $\mathcal{T}_{\mathbf{V}}(\mathbf{r}) \cap \mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}})$ , we have

$$\begin{aligned} \|\hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})/n\|_\infty &\leq 2m_{\mathbf{A}} (\|\hat{\Sigma}_{\mathbf{z}}\|_\infty \|\hat{\mathbf{A}} - \mathbf{A}\|_{L_1} + \|\mathbf{Z}^\top \mathbf{u}/n\|_\infty) \\ &\leq 2m_{\mathbf{A}} (4\|\hat{\Sigma}_{\mathbf{z}}\|_\infty \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + \lambda_{\mathbf{u}}), \end{aligned}$$

where the latter substitutions are justified by Lemma B.2 and the definition of  $\mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}})$ .

On the set  $\mathcal{T}_{\Theta}(\mu)$ , each row  $\theta$  is feasible for Problem 3.2. Then,  $\|\hat{\theta}_j\|_1 \leq \|\theta_j\|_1$  for each  $j \in [p_{\mathbf{x}}]$  by specification. Finally, Lemma B.5 entails that, for the present choice of  $r_{\beta}$ ,

$$\|\hat{\beta} - \beta\|_1 \leq 4 \frac{s_{\beta}}{\phi_{\beta}^2} r_{\beta}$$

on the set  $\mathcal{T}_{\mathbf{V}}(\lambda_{\mathbf{V}}) \cap \mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}})$ .

Combining the foregoing results, we see that, on the set

$$\mathcal{T}_{\mathbf{A}}(\mathbf{r}) \cap \mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}}) \cap \mathcal{T}_{\mathbf{V}}(\lambda_{\mathbf{V}}) \cap \mathcal{T}_{\Theta}(\mu),$$

it holds that

$$\|\hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})(\beta - \hat{\beta})/\sqrt{n}\|_\infty \leq 8m_{\Theta} m_{\mathbf{A}} \sqrt{n} (4\|\hat{\Sigma}_{\mathbf{z}}\|_\infty \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + \lambda_{\mathbf{u}}) \frac{s_{\beta}}{\phi_{\beta}^2} r_{\beta}.$$

Now note that, under the present choices of  $\mathbf{r}$  and  $\lambda_{\mathbf{V}}$ , the set  $\mathcal{T}_{\mathbf{V}}(\lambda_{\mathbf{V}})$  is contained in  $\mathcal{T}_{\mathbf{V}}(\mathbf{r})$ .  $\square$

**Lemma B.12** (Control of  $\mathbf{f}_3$ , sub-Gaussian Noise). Suppose that Conditions (1), (2), and (3) of Assumption 4.7 hold and (ii) Assumption 2.4 holds. Choose  $\hat{\Theta}$ ,  $\hat{\mathbf{A}}$  and  $\hat{\beta}$  according to Assumption 4.8. Then,

$$\begin{aligned} \mathbb{P}\{\|\mathbf{f}_3\|_\infty > 8m_{\Theta} m_{\mathbf{A}} \sqrt{n} (4m_{\mathbf{Z}} \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + \lambda_{\mathbf{u}}) \frac{s_{\beta}}{\phi_{\beta}^2} r_{\beta}\} \\ \leq ep_{\mathbf{z}}^{2-C_0 \min_{j \in [p_{\mathbf{x}}]} \{c_j^2 / \|\mathbf{v}^j\|_{\psi_2}^2\}} + ep_{\mathbf{z}}^{1-c_{\mathbf{u}}^2 C_0 / \tau_{\mathbf{u}}^2} + \mathbb{P} \mathcal{T}_{\Theta}(\mu)^c, \end{aligned}$$

where  $c_{\mathbf{V}}$  is as in Definition 4.4,  $c_{\mathbf{u}}$  is as in Definition 4.4, and  $C_0$  is as defined in Lemma C.1. Consequently, if Conditions (4), (5), and (b) of Assumption 4.7 also hold, then  $\|\mathbf{f}_3\|_1 = o_{\mathbb{P}}(1)$ .

**Proof of Lemma B.12.** Lemma B.11 entails that

$$\begin{aligned} \mathbb{P}\{\|\mathbf{f}_3\|_\infty > 8m_{\Theta} m_{\mathbf{A}} \sqrt{n} (4m_{\mathbf{Z}} \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + \lambda_{\mathbf{u}}) \frac{s_{\beta}}{\phi_{\beta}^2} r_{\beta}\} \\ \leq \mathbb{P} \mathcal{T}_{\mathbf{V}}(\lambda_{\mathbf{V}})^c + \mathbb{P} \mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}})^c + \mathbb{P} \mathcal{T}_{\Theta}(\mu)^c. \end{aligned}$$

Substitute the present choices of tuning parameters into the display above and cite the estimates of Lemmas C.2 and C.3, to deduce the first claim. Expand the present choices of tuning parameters to find

$$\begin{aligned} 8m_{\Theta} m_{\mathbf{A}} \sqrt{n} (4m_{\mathbf{Z}} \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + \lambda_{\mathbf{u}}) \frac{s_{\beta}}{\phi_{\beta}^2} r_{\beta} \\ \lesssim s_{\mathbf{A}}^3 s_{\beta} (\log p_{\mathbf{z}})^{3/2} / n + s_{\mathbf{A}}^2 s_{\beta} (\log p_{\mathbf{z}}) / \sqrt{n} \\ + s_{\mathbf{A}}^2 s_{\beta} (\log p_{\mathbf{z}})^{3/2} / n + s_{\mathbf{A}} s_{\beta} \log p_{\mathbf{z}} / \sqrt{n}, \end{aligned}$$

from which the latter claim follows.  $\square$

**Lemma B.13** (Control of  $\mathbf{f}_4$ ). Suppose that Assumption 4.3 and Condition (3) of Assumption 4.7 hold. Choose  $\hat{\Theta}$  according to Assumption 4.8; let  $\lambda_{\mathbf{V}}, \lambda_{\mathbf{u}} > 0$  be arbitrary; set

$$r_{\beta} = 16 \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} \|\hat{\Sigma}_{\mathbf{z}}\|_\infty (4m_{\beta} \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + m_{\mathbf{A}}) + (4 \frac{s_{\mathbf{A}}}{\phi_{\mathbf{A}}^2} r_{\mathbf{A}} + m_{\mathbf{A}}) (m_{\beta} \lambda_{\mathbf{V}} + \lambda_{\mathbf{u}}).$$

Then, on the set  $\mathcal{T}_{\mathbf{V}}(\lambda_{\mathbf{V}}) \cap \mathcal{T}_{\mathbf{u}}(\lambda_{\mathbf{u}}) \cap \mathcal{T}_{\Theta}(\mu)$ , the remainder term

$$\mathbf{f}_4 = \sqrt{n} (\hat{\Theta} \hat{\Sigma}_{\mathbf{d}} - \mathbf{I})(\beta - \hat{\beta})$$

satisfies

$$\|\mathbf{f}_4\|_\infty \leq 4\sqrt{n} \mu \frac{s_{\beta}}{\phi_{\beta}^2} r_{\beta}.$$



**Proof of Lemma B.13.** Note first that

$$\begin{aligned}\|\sqrt{n}(\widehat{\Theta}\widehat{\Sigma}_d - \mathbf{I})(\boldsymbol{\beta} - \hat{\boldsymbol{\beta}})\|_\infty &\leq \sqrt{n}\|\widehat{\Theta}\widehat{\Sigma}_d - \mathbf{I}\|_\infty\|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1 \\ &\leq \sqrt{n}\mu\|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1,\end{aligned}$$

where the latter inequality follows from the specification of  $\widehat{\Theta}$  and the fact that [Problem 3.2](#) is feasible given  $a$ . By [Lemma B.5](#), on the set  $\mathcal{T}_V(\lambda_V) \cap \mathcal{T}_U(\lambda_U) \cap \mathcal{T}_\Theta(\mu)$ ,

$$\|\boldsymbol{\beta} - \hat{\boldsymbol{\beta}}\|_1 \leq 4\frac{s_\beta}{\phi_\beta^2}r_\beta.$$

Combine the two previous displays to deduce the original claim.  $\square$

**Lemma B.14** (Control of  $\mathbf{f}_4$ , Gaussian Noise). Suppose that (i) [Assumption 4.3](#) and Condition (3) of [Assumption 4.7](#) hold and (ii) [Assumption 2.4](#) holds. Choose  $\widehat{\Theta}$ ,  $\widehat{\mathbf{A}}$  and  $\hat{\boldsymbol{\beta}}$  according to [Assumption 4.8](#). Then,

$$\mathbb{P}\{\|\mathbf{f}_4\|_\infty > 4\sqrt{n}\mu\frac{s_\beta}{\phi_\beta^2}r_\beta\} \leq ep_Z^{2-c_V^2c_0/\tau_V^2} + ep_Z^{1-c_U^2c_0/\tau_U^2} + \mathbb{P}\mathcal{T}_\Theta(\mu)^c,$$

where  $c_V$  is as in [Definition 4.4](#),  $c_U$  is as in [Definition 4.4](#), and  $C_0$  is as defined in [Lemma C.1](#). Consequently, if Conditions (4), (5), and (c) of [Assumption 4.7](#) also hold, then  $\|\mathbf{f}_4\|_1 = o_p(1)$ .

**Proof of Lemma B.14.** [Lemma B.13](#) entails that

$$\mathbb{P}\{\|\mathbf{f}_4\|_\infty > 4\sqrt{n}\mu\frac{s_\beta}{\phi_\beta^2}r_\beta\} \leq \mathbb{P}\mathcal{T}_V(\lambda_V)^c + \mathbb{P}\mathcal{T}_U(\lambda_U)^c + \mathbb{P}\mathcal{T}_\Theta(\mu).$$

Substitute the present choices of tuning parameters into the display above and cite the estimates of [Lemmas C.2](#) and [C.3](#) to deduce the first claim. Expand the present choices of tuning parameters to find

$$\sqrt{n}\mu\frac{s_\beta}{\phi_\beta^2}r_\beta \lesssim \mu s_\beta(s_A^2 \log p_Z / \sqrt{n} + s_A \sqrt{\log p_Z}),$$

from which the second claim follows.  $\square$

**Proof of Lemma 4.9.** The result follows from [Lemmas B.8](#), [B.10](#), [B.12](#), and [B.14](#).  $\square$

The following two Lemmas are required for [Lemma 4.10](#).

**Lemma B.15.** For  $\epsilon < m_A$ , it holds that

$$\begin{aligned}\mathbb{P}\{\|\Theta\widehat{\Sigma}_d - \mathbf{I}\|_\infty > t + 3m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \epsilon\} \\ \leq \mathbb{P}\{\|\Theta\widehat{\Sigma}_d - \mathbf{I}\|_\infty > t\} + \mathbb{P}\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} \leq \epsilon\},\end{aligned}$$

where  $\widehat{\Sigma}_d = \mathbb{E}_n[\mathbf{d}_i \mathbf{d}_i^\top]$ .

**Proof of Lemma B.15.** Note that

$$\begin{aligned}\widehat{\Sigma}_d &= \widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}/n = \mathbf{D}^\top \mathbf{D}/n + (\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{A}/n + \widehat{\mathbf{A}}^\top \mathbf{Z}^\top \mathbf{Z} (\widehat{\mathbf{A}} - \mathbf{A})/n \\ &= \widehat{\Sigma}_d + (\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{A}/n + \mathbf{A}^\top \mathbf{Z}^\top \mathbf{Z} (\widehat{\mathbf{A}} - \mathbf{A})/n \\ &\quad + (\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} (\widehat{\mathbf{A}} - \mathbf{A})/n,\end{aligned}$$

so that

$$\begin{aligned}\|\Theta\widehat{\Sigma}_d - \mathbf{I}\|_\infty &\leq \|\Theta(\mathbf{D}^\top \mathbf{D}/n) - \mathbf{I}\|_\infty + \|\Theta(\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{A}/n\|_\infty \\ &\quad + \|\Theta \mathbf{A}^\top \mathbf{Z}^\top \mathbf{Z} (\widehat{\mathbf{A}} - \mathbf{A})/n\|_\infty \\ &\quad + \|\Theta(\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} (\widehat{\mathbf{A}} - \mathbf{A})/n\|_\infty \\ &:= \|\Theta(\mathbf{D}^\top \mathbf{D}/n) - \mathbf{I}\|_\infty + \mathbf{I}_1 + \mathbf{I}_2 + \mathbf{I}_3.\end{aligned}$$

Note that

$$\begin{aligned}\mathbf{I}_1 &= \|\Theta(\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{A}/n\|_\infty \leq \|\Theta\|_{L_1} \|(\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z} \mathbf{A}/n\|_\infty \\ &\leq m_\Theta \|\widehat{\Sigma}_z\|_\infty \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} \|\mathbf{A}\|_{L_1} \\ &= m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1}.\end{aligned}$$

The same bound holds for  $\mathbf{I}_2$  by symmetry of the  $\ell_\infty$  norm under transposition of its argument.

For the term  $I_3$ , similar reasoning yields

$$\begin{aligned} I_3 &= \|\Theta(\widehat{\mathbf{A}} - \mathbf{A})^\top \mathbf{Z}^\top \mathbf{Z}(\widehat{\mathbf{A}} - \mathbf{A})/n\|_\infty \\ &\leq m_\Theta \|\widehat{\Sigma}_z\|_\infty \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2 = m_\Theta \|\widehat{\Sigma}_z\|_\infty \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1}^2. \end{aligned}$$

If  $\epsilon < m_A$ , then, on the set  $\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} \leq \epsilon\}$ , it holds that  $I_3 \leq m_\Theta m_A \cdot \|\widehat{\Sigma}_z\|_\infty \|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1}$ . Conclude that

$$P\{I_1 + I_2 + I_3 > 3m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \epsilon\} \leq P\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} > \epsilon\}$$

and hence that

$$\begin{aligned} P\{\|\Theta \widehat{\Sigma}_d - \mathbf{I}\|_\infty > t + 3m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \epsilon\} \\ \leq P\{\|\Theta \bar{\Sigma}_d - \mathbf{I}\|_\infty > t\} + P\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} > \epsilon\}, \end{aligned}$$

as claimed.  $\square$

**Lemma B.16.** Suppose that the  $\mathbf{z}_i$  satisfy [Assumption 2.2](#). Set

$$\mu = m_\Theta a \sqrt{\log(p_x)/n} + 3m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \epsilon$$

where  $a > 0$  and  $\epsilon > 0$  are controlled quantities. Then

$$P\mathcal{T}_\Theta(\mu) \leq 2p_x^{2-a^2/(6e^2\kappa^2)} + P\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} > \epsilon\},$$

where  $\kappa = m_A^2(2\tau_z^2 + \|\Sigma_z\|_\infty/\log 2)$ .

**Proof of Lemma B.16.** [Lemma B.15](#) entails that

$$\begin{aligned} P\{\|\Theta \widehat{\Sigma}_d - \mathbf{I}\|_\infty > t + 3m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \epsilon\} \\ \leq \underbrace{P\{\|\Theta \bar{\Sigma}_d - \mathbf{I}\|_\infty > t\}}_{I_1(t)} + P\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} > \epsilon\}. \end{aligned}$$

For  $t > 0$ . Now write

$$\Theta \bar{\Sigma}_d - \mathbf{I} = \Theta \Sigma_d + \Theta(\bar{\Sigma}_d - \Sigma_d) - \mathbf{I} = \Theta(\bar{\Sigma}_d - \Sigma_d)$$

to infer that

$$\|\Theta \bar{\Sigma}_d - \mathbf{I}\|_\infty = \|\Theta(\bar{\Sigma}_d - \Sigma_d)\|_\infty \leq \|\Theta\|_{L_1} \|\bar{\Sigma}_d - \Sigma_d\|_\infty \leq m_\Theta \|\bar{\Sigma}_d - \Sigma_d\|_\infty$$

and hence that

$$I_1(t) \leq P\{m_\Theta \|\bar{\Sigma}_d - \Sigma_d\|_\infty > t\}$$

Choose  $t = m_\Theta a \sqrt{\log(p_x)/n}$  for a controlled quantity  $a > 0$  and cite a slight modification of [Lemma C.5](#) to find

$$\begin{aligned} I_1(m_\Theta a \sqrt{\log(p_x)/n}) &\leq P\{\|\bar{\Sigma}_d - \Sigma_d\|_\infty > a \sqrt{\log(p_x)/n}\} \\ &\leq 2p_x^{2-a^2/(6e^2\kappa^2)}, \end{aligned}$$

where  $\kappa = m_A^2(2\tau_z^2 + \|\Sigma_z\|_\infty/\log 2)$ . Substitute the above bound into the first display of the present proof to conclude that

$$\begin{aligned} P\{\|\Theta \widehat{\Sigma}_d - \mathbf{I}\|_\infty > m_\Theta a \sqrt{\log(p_x)/n} + 3m_\Theta m_A \|\widehat{\Sigma}_z\|_\infty \epsilon\} \\ \leq 2p_x^{2-a^2/(6e^2\kappa^2)} + P\{\|\widehat{\mathbf{A}} - \mathbf{A}\|_{L_1} > \epsilon\}, \end{aligned}$$

as claimed.  $\square$

**Proof of Lemma 4.10.** Set  $\epsilon = 4 \frac{s_A}{\phi_A^2} c_V (\|\widehat{\Sigma}_z\|_\infty (\log p_z)/n)^{1/2}$ , plug this choice into the result of [Lemma B.16](#) and cite the estimate of [Lemma B.3](#) to find

$$\begin{aligned} P\{\|\Theta \widehat{\Sigma}_d - \mathbf{I}\|_\infty > \frac{m_\Theta}{\sqrt{n}} (a \sqrt{\log p_x} + 12m_A \frac{s_A}{\phi_A^2} c_V \|\widehat{\Sigma}_z\|_\infty^{3/2} \sqrt{\log p_z})\} \\ \leq 2p_x^{2-a^2/(6e^2\kappa^2)} + ep_z^{2-C_0 \min_{j \in [p_x]} \{c_j^2/\|\psi_j\|_{\psi_2}^2\}} \end{aligned}$$

for  $n$  sufficiently large, as claimed.  $\square$

**Proof of Lemma 4.11.** (i) We first show that  $\hat{\sigma}_u - \sigma_u = o_p(1)$ . To begin, write

$$\begin{aligned}\hat{\sigma}_u^2 &= \|\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}\|_2^2/n = \|\mathbf{u} + \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\|_2^2/n \\ &= \|\mathbf{u}\|_2^2/n + 2 \underbrace{\langle \mathbf{u}, \mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}) \rangle}_I_1 /n + \underbrace{\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\|_2^2}_I_2 /n.\end{aligned}$$

It follows that

$$\begin{aligned}(\|\mathbf{u}\|_2^2/n - \sigma_u^2) - 2|I_1| + I_2 &\leq \hat{\sigma}_u^2 - \sigma_u^2 \\ &\leq |\hat{\sigma}_u^2 - \sigma_u^2| \leq \|\mathbf{u}\|_2^2/n - \sigma_u^2 + 2|I_1| + I_2.\end{aligned}$$

We claim that  $I_1$  and  $I_2$  are each  $o_p(1)$ . It then follows that  $\hat{\sigma}_u^2 - \sigma_u^2 = o_p(1)$ , since  $\|\mathbf{u}\|_2^2/n - \sigma_u^2 = o_p(1)$  by assumption. To show the claim, note first that

$$I_1^2 \leq (\|\mathbf{u}\|_2^2/n)(\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\|_2^2/n) = (\|\mathbf{u}\|_2^2/n)I_2.$$

From the assumption that  $\|\mathbf{u}\|_2^2/n - \sigma_u^2 = o_p(1)$  we infer that  $\|\mathbf{u}\|_2^2/n = O_p(1)$ . Thus, it suffices to show that  $I_2 = o_p(1)$ . To this end, note that

$$\|\mathbf{X}(\hat{\boldsymbol{\beta}} - \boldsymbol{\beta})\|_2 \leq \|\mathbf{X}\|_{L_2} \|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1$$

hence  $I_2 \leq \|\mathbf{X}\|_{L_2}^2/n(\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1^2)$ . Now cite Lemma C.6 to obtain that

$$\|\mathbf{X}\|_{L_2}^2/n = O_p(\max_{j \in [p_{\mathbf{x}}]} E[\chi_{ij}^2] + a\sqrt{(\log p_{\mathbf{x}})/n})$$

for a suitably chosen  $a$  of constant order. It then follows from the growth conditions of the present lemma and the estimate of Lemma 4.5 for  $\|\hat{\boldsymbol{\beta}} - \boldsymbol{\beta}\|_1$  that  $I_2^2 = O_p(1)o_p(1) = o_p(1)$ , as required. Thus,  $\hat{\sigma}_u^2 - \sigma_u^2 = o_p(1)$ . To show as much for  $\hat{\sigma}_u - \sigma_u$ , it suffices to show that  $P\{|\hat{\sigma}_u - \sigma_u| > \epsilon\} \rightarrow 0$  for each  $\epsilon > 0$ . We will show that  $P\{\hat{\sigma}_u - \sigma_u > \epsilon\} \rightarrow 0$ ; the matching limit follows from an analogous argument. Fix  $\epsilon > 0$  and note that

$$\begin{aligned}P\{\hat{\sigma}_u - \sigma_u > \epsilon\} &= P\{\hat{\sigma}_u > \sigma_u + \epsilon\} \\ &= P\{\hat{\sigma}_u^2 > \sigma_u^2 + 2\sigma_u\epsilon + \epsilon^2\} \\ &= P\{\hat{\sigma}_u^2 - \sigma_u^2 > 2\sigma_u\epsilon + \epsilon^2\} \rightarrow 0\end{aligned}$$

as  $n \rightarrow \infty$  since, by Assumption 2.4,  $\sigma_u$  is bounded strictly away from zero uniformly in  $n$ . The previous display and the matching limit for  $P\{\sigma_u - \hat{\sigma}_u > \epsilon\}$  entail that  $\hat{\sigma}_u - \sigma_u = o_p(1)$ , as claimed.

(ii) We now show that  $\hat{\omega}_j^2 = \hat{\sigma}_u^2 \hat{\Theta}_{jj}$  satisfies  $\hat{\omega}_j - \omega_j = o_p(1)$ . We first show as much for  $\hat{\omega}_j^2 - \omega_j^2$ ; the original claim then follows from reasoning analogous to that above. To this end, note that, since the noise components  $u_i$  are homoscedastic, we have

$$\omega_j^2 = E[\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 u_i^2] = E[\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 E[u_i^2 | \mathbf{z}_i]] = \sigma_u^2 \Theta_{jj}.$$

Now write

$$\begin{aligned}\hat{\omega}_j^2 - \omega_j^2 &= \hat{\sigma}_u^2 \hat{\Theta}_{jj} - \sigma_u^2 \Theta_{jj} \\ &= \sigma_u^2 (\hat{\Theta}_{jj} - \Theta_{jj}) + (\hat{\sigma}_u^2 - \sigma_u^2) \hat{\Theta}_{jj} + (\sigma_u^2 - \hat{\sigma}_u^2) \Theta_{jj}.\end{aligned}$$

Next, recall that, on the set  $\mathcal{T}_{\Theta}(\mu)$ , it holds due to Lemma A.1 that  $\|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}_j\|_{\infty} \leq 2m_{\Theta}\mu$ . Letting  $\epsilon > 0$  be arbitrary, it then follows from the previous display that

$$P\{|\hat{\omega}_j^2 - \omega_j^2| > 2m_{\Theta}\mu(\sigma_u^2 + \epsilon) + \epsilon\Theta_{jj}\} \leq P\mathcal{T}_{\Theta}(\mu)^c + P\{|\hat{\sigma}_u^2 - \sigma_u^2| > \epsilon\}.$$

Now note that, due to the first claim of the present lemma and condition (4) of Assumption 4.7, we have

$$\begin{aligned}\lim_{n \rightarrow \infty} P\{|\hat{\omega}_j - \omega_j| > 2m_{\Theta}\mu(\sigma_u + \epsilon) + \epsilon\Theta_{jj}\} \\ \leq \lim_{n \rightarrow \infty} P\mathcal{T}_{\Theta}(\mu)^c + P\{|\hat{\sigma}_u - \sigma_u| > \epsilon\} = 0.\end{aligned}$$

Finally, cite Proposition C.4 and the present assumption that  $\mu = o(1)$  to find that

$$\lim_{n \rightarrow \infty} m_{\Theta}\mu(\sigma_u + \epsilon) + \epsilon\Theta_{jj} = 0.$$

Conclude that  $\omega_j^2 - \hat{\omega}_j^2 = o_p(1)$  and hence that  $\hat{\omega}_j - \omega_j = o_p(1)$ , as claimed.  $\square$

## Appendix C. Supplementary data

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.jeconom.2019.09.009>.

## References

- Amemiya, T., 1974. The nonlinear two-stage least-squares estimator. *Journal of Econometrics* 2, 105–110.
- Amemiya, T., 1977. The maximum likelihood and the nonlinear three-stage least squares estimator in the general nonlinear simultaneous equation model. *Econometrica* 45, 955–968.
- Angrist, J., Pischke, J.-S., 2009. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press.
- Belloni, A., Chen, D., Chernozhukov, V., Hansen, C., 2012. Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica* 80, 2369–2429.
- Belloni, A., Chernozhukov, V., 2013. Least squares after model selection in high-dimensional sparse models. *Bernoulli* 19, 521–547.
- Belloni, A., Chernozhukov, V., Hansen, C., 2011. Inference for high-dimensional sparse econometric models. In: *Advances in Economics and Econometrics: Tenth World Congress Volume 3*, Econometrics.
- Belloni, A., Chernozhukov, V., Hansen, C., Newey, W., 2018. Simultaneous confidence intervals for high-dimensional linear models with many endogenous variables. *arXiv preprint arXiv:1712.08102*.
- Bickel, P., Ritov, Y., Tsybakov, A., 2009. Simultaneous analysis of Lasso and Dantzig selector. *Annals of Statistics* 37, 1705–1732.
- Bühlmann, P., van de Geer, S., 2011. *Statistics for High-Dimensional Data*. Springer.
- Cai, T., Liu, W., Luo, X., 2011. A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation. *Journal of the American Statistical Association* 106, 594–607.
- Caner, M., Fan, Q., 2015. Hybrid generalized empirical likelihood estimators: Instrument selection with adaptive lasso. *Journal of Econometrics* 187, 256–274.
- Caner, M., Kock, A.B., 2018a. Asymptotically honest confidence regions for high dimensional parameters by the desparsified conservative Lasso. *Journal of Econometrics* 203, 143–168.
- Caner, M., Kock, A.B., 2018b. High dimensional linear GMM. *arXiv preprint arXiv:1811.08779*.
- Chamberlain, G., 1987. Asymptotic efficiency in estimation with conditional moment restrictions. *Journal of Econometrics* 34, 305–334.
- Cheng, X., Liao, Z., 2015. Select the valid and relevant moments: An information-based lasso for gmm with many moments. *Journal of Econometrics* 186, 443–464.
- Chételat, D., Lederer, J., Salmon, J., 2017. Optimal two-step prediction in regression. *Electronic Journal of Statistics* 11, 2519–2546.
- Chichignoud, M., Lederer, J., Wainwright, M., 2016. A practical scheme and fast algorithm to tune the Lasso with optimality guarantees. *Journal of Machine Learning Research* 17, 1–17.
- Dalalyan, A., Hebiri, M., Lederer, J., 2017. On the prediction performance of the Lasso. *Bernoulli* 23, 552–581.
- Fan, J., Li, R., 2001. Variable selection via nonconcave penalized likelihood and its oracle properties. *Journal of the American Statistical Association* 96, 1348–1360.
- Fan, J., Liao, Y., 2014. Endogeneity in high dimensions. *Annals of Statistics* 42, 872–917.
- Fan, Q., Zhong, W., 2018. Nonparametric additive instrumental variable estimator: a group shrinkage estimation perspective. *Journal of Business and Economic Statistics* 36 (3), 388–399.
- Friedman, J., Hastie, T., Tibshirani, R., 2010. Regularization paths for generalized linear models via coordinate descent. *Journal of Statistical Software* 33, 1–22.
- Gautier, E., Tsybakov, A., 2014. High-dimensional instrumental variables regression and confidence sets. *ArXiv:1105.2454v4*.
- Gautier, E., Tsybakov, A., 2018. High-dimensional instrumental variables regression and confidence sets. *ArXiv e-prints*.
- van de Geer, S., Bühlmann, P., 2009. On the conditions used to prove oracle results for the Lasso. *Electronic Journal of Statistics* 3, 1360–1392.
- van de Geer, S., Bühlmann, P., Ritov, Y., Dezeure, R., 2014. On asymptotically optimal confidence regions and tests for high-dimensional models. *Annals of Statistics* 42, 1166–1202.
- van de Geer, S., Lederer, J., 2013. The Lasso, correlated design, and improved oracle inequalities. In: *From Probability to Statistics and Back: High-Dimensional Models and Processes – A Festschrift in Honor of Jon A. Wellner*. In: *Collections*, vol. 9, Institute of Mathematical Statistics, Beachwood, Ohio, USA, pp. 303–316.
- van de Geer, S., Muro, A., 2014. On higher order isotropy conditions and lower bounds for sparse quadratic forms. *Electronic Journal of Statistics* 8, 3031–3061.
- Giraud, C., 2014. Introduction to high-dimensional statistics. In: *Monographs on Statistics and Applied Probability (Series)*, vol. 139, CRC Press, Taylor & Francis Group.
- Hansen, L., 1982. Large sample properties of generalized method of moments estimators. *Econometrica* 50, 1029–1054.
- Hansen, C., Kozbur, D., 2014. Instrumental variables estimation with many weak instruments using regularized jive. *Journal of Econometrics* 182, 290–308.
- Hastie, T., Tibshirani, R., Wainwright, M., 2015. *Statistical learning with sparsity: The Lasso and generalizations*. In: *Monographs on Statistics and Applied Probability (Series)*, vol. 143, Boca Raton: CRC Press, Taylor & Francis Group.
- Hebiri, M., Lederer, J., 2013. How correlations influence lasso prediction. *IEEE Trans. Inform. Theory* 59, 1846–1854.
- Imbens, G., Donald, S., Newey, W., 2003. Empirical likelihood estimation and consistent tests with conditional moment restrictions. *Journal of Econometrics* 117, 55–93.
- Jankova, J., van de Geer, S., 2018. Semiparametric efficiency bounds for high-dimensional models. *Annals of Statistics* 46 (5), 2336–2359.
- Javanmard, A., Montanari, A., 2014. Confidence intervals and hypothesis testing for high-dimensional regression. *Journal of Machine Learning Research* 15, 2869–2909.
- Knight, K., Fu, W., 2000. Asymptotics for lasso-type estimators. *Annals of Statistics* 28, 1356–1378.
- Lederer, J., Yu, L., Gaynanova, I., 2019. Oracle inequalities for high-dimensional prediction. *Bernoulli* 25, 1225–1255.
- Lin, W., Feng, R., Li, H., 2015. Regularization methods for high-dimensional instrumental variables regression with an application to genetical genomics. *Journal of the American Statistical Association* 110, 270–288.
- Liu, H., Wang, L., Zhao, T., 2015. Calibrated multivariate regression with application to neural semantic basis discovery. *J. Mach. Learn. Res.* 16, 1579–1606.
- MOSEK ASP, 2017. MOSEK Rmosek Package 8.1.0.34.
- Newey, W., 1990. Efficient instrumental variables estimation of nonlinear models. *Econometrica* 58, 809–837.
- Neykov, M., Ning, Y., Liu, J., Liu, H., 2015. A unified theory of confidence regions and testing for high dimensional estimating equations. *arXiv:1510.08986*.
- Ning, Y., Liu, H., 2017. A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *Annals of Statistics* 45, 158–195.
- Pötscher, B., Leeb, H., 2009. On the distribution of penalized maximum likelihood estimators: The LASSO, SCAD, and thresholding. *Journal of Multivariate Analysis* 100, 2065–2082.
- Rudelson, M., Zhou, S., 2012. Reconstruction from anisotropic random measurements. In: *JMLR Workshop Conf. Proc.*, Vol. 23, pp. 10.1–10.28.
- van der Vaart, A., Wellner, J., 1996. *Weak Convergence and Empirical Processes: With Applications to Statistics*. Springer-Verlag New York.

- Vershynin, R., 2012. Introduction to the non-asymptotic analysis of random matrices. In: Eldar, Y.C., Kutyniok, G. (Eds.), *Compressed Sensing: Theory and Applications*, Vol. 5. Cambridge University Press, pp. 210–268.
- Zhang, C.-H., 2010. Nearly unbiased variable selection under minimax concave penalty. *Annals of Statistics* 38, 894–942.
- Zhang, C.-H., Zhang, S., 2014. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 76, 217–242.
- Zhu, Y., 2018. Sparse linear models and l1-regularized 2sls with high-dimensional endogenous regressors and instruments. *Journal of Econometrics* 202 (2), 196–213.
- Zhuang, R., Lederer, J., 2018. Maximum regularized likelihood estimators: A general prediction theory and applications. *Statistics* 7, e186, sta4.186.