

Confidence Intervals for Causal Effects with Invalid Instruments using Two-Stage Hard Thresholding

Zijian Guo

University of Pennsylvania, Philadelphia, USA

Hyunseung Kang

Stanford University, Stanford, USA

T. Tony Cai

University of Pennsylvania, Philadelphia, USA

Dylan S. Small

University of Pennsylvania, Philadelphia, USA

Summary.

The instrumental variable (IV) method is commonly used to estimate the causal effect of a treatment on an outcome by using IVs that satisfy the assumptions of association with treatment, no direct effect on the outcome and ignorability. A major challenge in IV analysis is to find said IVs, but typically one is unsure of whether all of the putative IVs are in fact valid (i.e. satisfy the assumptions). We propose a general inference procedure that provides honest inference in the presence of invalid IVs, even after controlling for a large number of covariates. The key step of our method is a novel selection procedure, which we call Two-Stage Hard Thresholding (TSHT), where we use hard thresholding to select the set of non-redundant instruments in the first stage and subsequently use hard thresholding to select the valid instruments in the second stage using the thresholding from the first stage. TSHT allows us to not only select invalid IVs, but also provides honest confidence intervals of the treatment effect at \sqrt{n} rate. We establish asymptotic properties of our procedure and demonstrate that our procedure performs well in simulation studies compared to traditional IV methods, especially when the instruments are invalid.

Keywords: Causal effect; confidence interval; exclusion restriction; high-dimensional covariates; invalid instruments; treatment effect.

1. Introduction

1.1. *Motivation: Invalid Instruments Even After Controlling for (Potentially Many) Confounders*

Instrumental variables (IV) analysis is a popular method to deduce causal effects in the presence of unmeasured confounding. An IV analysis requires variables called instruments that (A1) are related to the exposure (A2) have no direct pathway to the outcome and (A3) are not related to unmeasured variables that affect the exposure and the outcome (see Section 2.2 for details). Variables that satisfy these

assumptions are referred to as valid instruments. A major challenge in IV analysis is to find valid instruments

In practice, it is often the case that potential candidate instruments become more plausible as valid instruments after controlling for some, possibly high dimensional, covariates (see Hernán and Robins [2006] and Baiocchi et al. [2014] for discussion on control for covariates for an instrument to be valid). For example, a long-standing interest in economics is the causal effect of education on earnings and often, IV analysis is used to deduce the effect [Angrist and Krueger, 1991, Card, 1993, 1999]. A popular instrument in this analysis is a person’s proximity to a college when growing up [Card, 1999, 1993]. However, proximity to a college may be related to a person’s socioeconomic status, characteristics of a person’s high school and other covariates that may affect a person’s earnings. Thus, these covariates need to be controlled for in order for proximity to college to be a valid IV and with the growing trend toward collecting large data sets with many variables, this approach of finding instrumental variables that are valid after conditioning on covariates has increasing promise [Hernán and Robins, 2006, Swanson and Hernán, 2013, Baiocchi et al., 2014, Varian, 2014, Imbens, 2014].

Yet, despite the promise that large data sets may bring in terms of finding valid instruments by conditioning on potentially many covariates, some IVs may still turn out to be invalid and subsequent analysis assuming that all the IVs are valid after conditioning can be misleading [Murray, 2006]. For example, suppose for studying the causal effect of education on earnings, we used proximity as an IV and to make sure the IV satisfies (A3), we control for confounders like high school test scores of the student, high school size, individual’s genetic makeup, family education, and family’s socioeconomic status. But, if living close to college had other benefits beyond getting more education, say by being exposed to many programs available to high school students for job preparation and employers who come to the area to discuss employment opportunities for college students, then the IV, proximity to college, can directly affect individual’s earning potential and violate (A2) [Card, 1999]. This problem is also widely prevalent in other applications of instrumental variables, most notably in Mendelian randomization [Davey Smith and Ebrahim, 2003, 2004] where the instruments are genetic in nature and some instruments are likely to be invalid due to having pleiotropic effects [Lawlor et al., 2008, Burgess et al., 2015].

This paper tackles the problem of constructing confidence intervals for causal effects that are robust to invalid instruments even after controlling for possibly high dimensional covariates.

1.2. *Prior Work*

In non-IV settings with many high dimensional covariates, Zhang and Zhang [2014], Javanmard and Montanari [2014], van de Geer et al. [2014], Belloni et al. [2014] and Cai and Guo [2016a] provide honest confidence intervals for a treatment effect. In IV settings with high dimensional covariates (or IVs), Gautier and Tsybakov [2011], Belloni et al. [2012], Fan and Liao [2014] and Chernozhukov et al. [2015] provide

honest confidence intervals for a treatment effect, under the assumption that all the IVs are valid after controlling for said covariates. In invalid IV settings, Kolesár et al. [2015] and Bowden et al. [2015] provide inferential methods for treatment effects. However, the method requires that the effects of the IVs on the treatment be orthogonal to their direct effects on the outcome, a stringent assumption. Bowden et al. [2016], Burgess et al. [2016], Kang et al. [2016b] and Windmeijer et al. [2016] also work on the invalid IV setting, but without making the stringent orthogonality assumption. Unfortunately, all these papers focuses on the low dimensional setting and some only work in the case where the IVs are completely uncorrelated/orthogonal to each other unless modifications are made [Bowden et al., 2015, Burgess et al., 2016]. Furthermore, all the previous work only provides a consistent estimator of the treatment effect without any theoretical guarantees on inference; in fact, one of the simplest consistent estimators in this setting, the median estimator [Bowden et al., 2016, Burgess et al., 2016, Windmeijer et al., 2016] has been shown to be consistent, but not \sqrt{n} consistent [Windmeijer et al., 2016].

There are two major challenges in obtaining valid confidence intervals in our problem: (i) potentially high-dimensional covariates and (ii) the invalid IVs. The problem related with high-dimensional covariates can be dealt with by applying recent debiasing methods [Zhang and Zhang, 2014, Javanmard and Montanari, 2014, van de Geer et al., 2014, Cai and Guo, 2016a]. However, the general idea behind debiasing does not inherently resolve the invalid IV problem as even a single IV that is improperly assumed as valid while it is truly invalid can make these debiased estimates useless. To put it another way, debiasing as a method is only meant to asymptotically remove the bias of regression coefficients from ℓ_1 shrinkage estimators and to conduct proper inference on these de-biased coefficients. This methodological goal is different than in the invalid IV problem where the goal is to properly estimate a set of valid IVs, as even a single error of declaring an IV that is invalid as valid can lead to dishonest inference. In fact, the methodological challenge is not only to correctly select IVs, but also once selected, to do robust inference using the selected IVs.

1.3. Our Contributions

Although there are existing methods for estimating the treatment effect in the presence of possibly invalid IVs, there is a paucity of procedures for selecting the set of valid instruments and forming confidence intervals for the treatment effects with theoretical coverage guarantees. In this paper, we propose a novel two-stage hard thresholding (TSHT) procedure to estimate the set of valid instruments and form confidence intervals with theoretical coverage guarantee. As the name suggests, a key component of TSHT is the two sequential steps of hard-thresholding procedures common in high dimensional inference [Donoho and Johnstone, 1994, Donoho, 1995] to simultaneously allow for invalid IVs and endogeneity of the treatment. Specifically, in the first thresholding stage, we select non-redundant IVs (see Definition 2 for details) and in the second thresholding stage, we use the thresholded estimates from the first thresholding step as pilot estimates to guide the selection of the set

of valid instruments; see Section 3.3 for details. Using our two-stage variant of thresholding properly accounts for the selection of IVs and leads to $1/\sqrt{n}$ rate confidence intervals with desired coverage in both low and high dimensional settings where invalid IVs are present and without knowing a priori which of these IVs are invalid. Also, for the low dimensional covariate setting, our procedure is the first to have theoretical guarantees that it performs as well asymptotically as the oracle procedure that knows which instruments are valid. For the high dimensional covariate setting, our procedure is the first available procedure for forming confidence intervals with desired coverage when there may be invalid IVs.

The outline of the paper is as follows. After describing the model setup in Section 2, we formulate our TSHT procedure in Section 3. In Section 4, we develop the theoretical properties of our procedure. In Section 5, we investigate the performance of our procedure in a large simulation study and find that our confidence interval performs very similarly with respect to the oracle even if some of the underlying theoretical assumptions made in Section 4 are violated (see Sections 4 and 5 for details). In Section 6, we present an empirical study where we revisit the question of the causal effect of years of schooling on income using data from the Wisconsin Longitudinal Study. We provide conclusions and discussion in Section 7.

2. Model

2.1. Notation

To define causal effects, the potential outcome approach Neyman [1923], Rubin [1974] for instruments laid out in Holland [1988] is used. For each individual $i \in \{1, \dots, n\}$, let $Y_i^{(d, \mathbf{z})} \in \mathbb{R}$ be the potential outcome if the individual were to have exposure $d \in \mathbb{R}$ and instruments $\mathbf{z} \in \mathbb{R}^{p_z}$. Let $D_i^{(\mathbf{z})} \in \mathbb{R}$ be the potential exposure if the individual had instruments $\mathbf{z} \in \mathbb{R}^{p_z}$. For each individual, only one possible realization of $Y_i^{(d, \mathbf{z})}$ and $D_i^{(\mathbf{z})}$ is observed, denoted as Y_i and D_i , respectively, based on his observed instrument values $\mathbf{Z}_i \in \mathbb{R}^{p_z}$ and exposure D_i . We also denote pre-instrument covariates for each individual i as $\mathbf{X}_i \in \mathbb{R}^{p_x}$. In total, n sets of outcome, exposure, and instruments, denoted as $(Y_i, D_i, \mathbf{Z}_i, \mathbf{X}_i)$, are observed in an i.i.d. fashion.

We denote $\mathbf{Y} = (Y_1, \dots, Y_n)$ to be an n -dimensional vector of observed outcomes, $\mathbf{D} = (D_1, \dots, D_n)$ to be an n -dimensional vector of observed exposures/treatment, \mathbf{Z} to be a n by p_z matrix of instruments where row i consists of \mathbf{Z}_i , and \mathbf{X} to be an n by p_x matrix of covariates where row i consists of \mathbf{X}_i . Let \mathbf{W} be an n by $p = p_z + p_x$ matrix where \mathbf{W} is a result of concatenating the matrices \mathbf{Z} and \mathbf{X} and $\boldsymbol{\Sigma}^* = \mathbf{E}(\mathbf{W}_i \mathbf{W}_i^\top)$ is positive definite. For any vector $\mathbf{v} \in \mathbb{R}^p$, let \mathbf{v}_j denote the j th element of \mathbf{v} . Let $\|\mathbf{v}\|_1$, $\|\mathbf{v}\|_2$, and $\|\mathbf{v}\|_\infty$ denote the usual 1, 2 and ∞ -norms, respectively. Let $\|\mathbf{v}\|_0$ denote the number of non-zero elements in \mathbf{v} and $\text{supp}(\mathbf{v}) \subseteq \{1, \dots, p\}$, is defined as $\{j : \mathbf{v}_j \neq 0\}$.

For any n by p matrix $\mathbf{M} \in \mathbb{R}^{n \times p}$, we denote the (i, j) element of matrix \mathbf{M} as M_{ij} , the i th row as \mathbf{M}_i , and the j th column as \mathbf{M}_j . Let \mathbf{M}^\top be the transpose of \mathbf{M} and $\|\mathbf{M}\|_\infty$ represent the element-wise matrix sup norm of matrix \mathbf{M} . For a

sequence of random variables X_n , we use $X_n \xrightarrow{p} X$ and $X_n \xrightarrow{d} X$ to represent that X_n converges to X in probability and in distribution, respectively. Finally, for any two sequences a_n and b_n , we will write $a_n \gg b_n$ if $\limsup \frac{b_n}{a_n} = 0$ and write $a_n \ll b_n$ if $b_n \gg a_n$. Also, for a set J , $|J|$ denotes its cardinality.

2.2. Model and Instrumental Variables Assumptions

We consider the Additive Linear, Constant Effects (ALICE) model of [Holland, 1988] and extend it to allow for multiple valid and possibly invalid instruments as in Small [2007] and Kang et al. [2016b]. For two possible values of the exposure d', d and instruments \mathbf{z}', \mathbf{z} , we assume the following potential outcomes model

$$Y_i^{(d', \mathbf{z}')} - Y_i^{(d, \mathbf{z})} = (\mathbf{z}' - \mathbf{z})^\top \boldsymbol{\kappa}^* + (d' - d)\beta^*, \quad \mathbf{E}(Y_i^{(0,0)} \mid \mathbf{Z}_{i.}, \mathbf{X}_{i.}) = \mathbf{Z}_{i.}^\top \boldsymbol{\eta}^* + \mathbf{X}_{i.}^\top \boldsymbol{\phi}^* \quad (1)$$

where $\boldsymbol{\kappa}^*, \beta^*, \boldsymbol{\eta}^*$, and $\boldsymbol{\phi}^*$ are unknown parameters. The parameter β^* represents the causal parameter of interest, the causal effect (divided by $d' - d$) of changing the exposure from d' to d on the outcome. The parameter $\boldsymbol{\phi}^*$ represents the impact of covariates on the baseline potential outcome $Y_i^{(0,0)}$. The parameter $\boldsymbol{\kappa}^*$ represents violation of (A2), the direct effect of the instruments on the outcome. If (A2) holds, then $\boldsymbol{\kappa}^* = 0$. The parameter $\boldsymbol{\eta}^*$ represents violation of (A3), the presence of unmeasured confounding between the instrument and the outcome. If (A3) holds, then $\boldsymbol{\eta}^* = 0$.

Let $\boldsymbol{\pi}^* = \boldsymbol{\kappa}^* + \boldsymbol{\eta}^*$ and $\epsilon_{i1} = Y_i^{(0,0)} - \mathbf{E}(Y_i^{(0,0)} \mid \mathbf{Z}_{i.}, \mathbf{X}_{i.})$. When we combine equation (1) along with the definition of ϵ_{i1} , the observed data model becomes

$$Y_i = \mathbf{Z}_{i.}^\top \boldsymbol{\pi}^* + D_i \beta^* + \mathbf{X}_{i.}^\top \boldsymbol{\phi}^* + \epsilon_{i1}, \quad \mathbf{E}(\epsilon_{i1} \mid \mathbf{Z}_{i.}, \mathbf{X}_{i.}) = 0 \quad (2)$$

and we denote $\sigma^2 = \text{Var}(\epsilon_{i1} \mid \mathbf{Z}_{i.}, \mathbf{X}_{i.})$. The observed model is also known as the under-identified single-equation linear model in econometrics (page 83 of Wooldridge [2010]). This model is not a usual regression model because D_i might be correlated with ϵ_{i1} . In particular, the parameter β^* measures the causal effect of changing D on Y rather than an association. Also, the parameter $\boldsymbol{\pi}^*$ in model (2) combines both the violation of (A2), represented by $\boldsymbol{\kappa}^*$, and the violation of (A3), represented by $\boldsymbol{\eta}^*$. If both (A2) and (A3) are satisfied for IVs, typically referred to as valid IVs [Murray, 2006], then $\boldsymbol{\kappa}^* = \boldsymbol{\eta}^* = 0$ and $\boldsymbol{\pi}^* = 0$. Hence, $\boldsymbol{\pi}^*$ captures invalid IVs, i.e. the violations of (A2) and (A3). We formalize this notion with the following definition.

DEFINITION 1. Suppose we have p_z candidate instruments along with the models (1)–(2). We say that instrument $j = 1, \dots, p_z$ is valid, i.e. satisfies (A2) and (A3), if $\boldsymbol{\pi}_j^* = 0$.

We also assume a linear association/observational model between the endogenous variable D_i , the instruments $\mathbf{Z}_{i.}$, and the covariates $\mathbf{X}_{i.}$.

$$D_i = \mathbf{Z}_{i.}^\top \boldsymbol{\gamma}^* + \mathbf{X}_{i.}^\top \boldsymbol{\psi}^* + \epsilon_{i2}, \quad \mathbf{E}(\epsilon_{i2} \mid \mathbf{Z}_{i.}, \mathbf{X}_{i.}) = 0 \quad (3)$$

Each element γ_j^* is the partial correlation between the j th instrument and D . The parameter ψ^* represents the association between the covariates and D_i . Also, unlike the models (1)-(2), we do not need a causal model between D_i , \mathbf{Z}_i , and \mathbf{X}_i ; only the association model (3) is sufficient for our method. Finally, for notation, we let $s_{z2} = \|\boldsymbol{\pi}^*\|_0$, $s_{x2} = \|\boldsymbol{\phi}^*\|_0$, $s_{z1} = \|\boldsymbol{\gamma}^*\|_0$, $s_{x1} = \|\boldsymbol{\psi}^*\|_0$ and $s = \max\{s_{z2}, s_{x2}, s_{z1}, s_{x1}\}$. Finally, in both models, the instruments and the covariates are exogenous to the error terms; see Wooldridge [2010] for textbook discussion on exogeneity.

Based on model (3), we can define a set of instruments that satisfy (A1), or sometimes referred to as non-redundant instruments in the econometrics literature [Cheng and Liao, 2015].

DEFINITION 2. *Suppose we have p_z candidate instruments along with the model (3). We say that instrument $j = 1, \dots, p_z$ satisfies (A1), or is a non-redundant IV, if $\gamma_j^* \neq 0$ and denote \mathcal{S}^* to be the set of these instruments.*

Typically, satisfying (A1) has been defined in a global sense where (A1) is satisfied if $\boldsymbol{\gamma}^* \neq 0$ [Wooldridge, 2010]. However, this global definition can be misleading in the presence of multiple candidate instruments. For example, it is possible that $\gamma_1^* \gg 0$ while $\gamma_j^* = 0$ for all $j \neq 1$ so that only the first instrument has an effect on the exposure while the rest do not. Using the global definition would imply that all the p_z instruments satisfy (A1) while Definition 2 makes it explicit and, perhaps less ambiguous, that it is only the first instrument $j = 1$ that satisfies (A1). Nevertheless, both the traditional global definition and Definition 2 are equivalent if $\gamma_j^* \neq 0$ for all j , that is where we only include relevant instruments, which is typically the case in practice and is the scenario studied by Kang et al. [2016b]. Finally, when there is only one candidate instrument so that $p_z = 1$, both definitions are equivalent to the definition presented in Holland [1988] and both become a special case of the definition presented in Angrist et al. [1996] under an additive, linear, constant effects model. In short, Definition 2 agrees with most definitions of satisfying (A1) in the literature.

Combining Definitions 1 and 2, we can formally define the usual three core conditions, i.e. (A1)-(A3), that define instruments.

DEFINITION 3. *Suppose we have p_z candidate instruments along with the models (1)-(3). We say that the Z_j , $j = 1, \dots, p_z$, is an instrument if (A1) – (A3) are satisfied, i.e. if $\pi_j^* = 0$ and $\gamma_j^* \neq 0$. Let \mathcal{V}^* be the set of instruments.*

When there is only one instrument, $p_z = 1$, Definition 3 of an instrument is identical to the definition of an instrument in Holland [1988]. Specifically, Definition 2 satisfies assumption (A1) that the instrument is related to the exposure. Also, assumption (A2), the exclusion restriction, which means $Y_i^{(d, \mathbf{z})} = Y_i^{(d, \mathbf{z}')}$ for all $d, \mathbf{z}, \mathbf{z}'$, is equivalent to $\boldsymbol{\kappa}^* = 0$ and assumption (A3), no unmeasured confounding, which means $Y_i^{(d, \mathbf{z})}$ and $D_i^{(\mathbf{z})}$ are independent of Z_i for all d and \mathbf{z} , is equivalent to $\boldsymbol{\eta}^* = 0$, implying $\boldsymbol{\pi}^* = \boldsymbol{\kappa}^* + \boldsymbol{\eta}^* = \mathbf{0}$. Definition 3 is also a special case of the definition of an instrument in Angrist et al. [1996] where here we assume the model is additive, linear, and has a constant treatment effect β^* . Hence, when multiple

instruments, $p_z > 1$, are present, our models (1)–(3) and Definition 3 can be viewed as a generalization of the definition of instruments in Holland [1988].

Note that the models presented above are commonly used in applications of IVs in econometrics [Wooldridge, 2010] and applications of IVs in genetic epidemiology and Mendelian randomization [Didelez and Sheehan, 2007]. However, we generalize these widely used models in two important ways: (i) the model in (2) allows for possibly invalid instruments and (ii) we allow the number of covariates p_x (and even the number of instruments p_z) to be larger than the sample size n .

3. Confidence Interval Estimation via Two-Stage Hard Thresholding

3.1. General Approach

The construction of our confidence interval can be broken down into two parts. The first part, detailed in Section 3.2, is estimating ITT effects based on the models (2) and (3). As we will see, the first part primarily deals with the problem posed by potentially high dimensional covariates, specifically the bias that comes from penalized estimators for high dimensional regression. The second part, which is elaborated in Section 3.3, tackles the heart of the problem in this paper, the presence of invalid IVs even after conditioning on high dimensional controls. Here, we take a novel two-stage hard thresholding approach to correctly select the valid IVs. Specifically, in the first step, we estimate the set of IVs that satisfy (A1) and in the second step, we use these IVs as initial guides to find IVs that satisfy (A2) and (A3) using the estimated set in the first step. Combining the two parts gives our confidence interval estimation procedure and is summarized in Procedure 1.

Procedure 1 Two-Stage Hard Thresholding (TSHT) for Confidence Interval for β^* under Invalid IVs with High-Dimensional Covariates

Input: Outcome \mathbf{Y} , treatment \mathbf{D} , instrument \mathbf{Z} , covariates \mathbf{X} , significance level α

STEP 1: Estimate ITT effects (i.e. $\tilde{\Gamma}$, $\tilde{\gamma}$) via debiased scaled Lasso in (6)-(9)

STEP 2: Select valid IVs (i.e. \tilde{V}) via two-stage hard thresholding

STEP 2a: Estimate IVs satisfying (A1) (i.e. $\tilde{\mathcal{S}}$) via hard thresholding $\tilde{\gamma}$ in (10)

STEP 2b: For each IV satisfying (A1) (i.e. $j \in \tilde{\mathcal{S}}$), estimate IVs satisfying (A2) and (A3) via hard-thresholding $\tilde{\pi}^{[j]}$ in (11)-(12)

STEP 3: Combine STEP 1 and STEP 2 via (13)-(15) to obtain confidence interval

Output: $1 - \alpha$ Confidence interval for β^*

Procedure 1 provides a general recipe to construct confidence intervals in the presence of invalid IVs and high dimensional covariates. The key step in the procedure is STEP 2, where we utilize two-stage hard thresholding, to deal with the problem posed by invalid IVs; as such, we call our procedure TSHT procedure, akin to the acronym for two-stage least squares (TSLS) procedure in IV, arguably the most popular IV estimator in the literature. We also note that the Procedure 1 as stated can handle (i) low dimensional covariates, (ii) high dimensional covariates, and (iii) settings with all IVs having no direct effect and no unmeasured confound-

ing, which is unrealistic in practice, and can still obtain valid confidence intervals. However, depending on particular data sets one may have, the procedure can be modified for simplicity and, in some cases, efficiency; Sections 3.5 and 3.6 discusses these cases in detail.

3.2. Estimating ITT Effects

The first part of the confidence interval procedure involves estimation of ITT effects. Specifically, given the observed models (2) and (3), we can write reduced-forms models where both models are only functions of \mathbf{Z}_i and \mathbf{X}_i .

$$Y_i = \mathbf{Z}_i^\top \boldsymbol{\Gamma}^* + \mathbf{X}_i^\top \boldsymbol{\Psi}^* + e_{i1} \quad (4)$$

$$D_i = \mathbf{Z}_i^\top \boldsymbol{\gamma}^* + \mathbf{X}_i^\top \boldsymbol{\psi}^* + \epsilon_{i2} \quad (5)$$

Here, $\boldsymbol{\Gamma}^* = \beta^* \boldsymbol{\gamma}^* + \boldsymbol{\pi}^*$ and $\boldsymbol{\Psi}^* = \phi^* + \beta^* \boldsymbol{\psi}^*$ are the parameters of the reduced-form model with $\boldsymbol{\Gamma}^*$ representing the ITT effect of the instruments on the outcome and $\boldsymbol{\gamma}^*$ representing the ITT effect of the instruments on the exposure. The term $e_{i1} = \beta^* \epsilon_{i2} + \epsilon_{i1}$ is the reduced-form error term in (4). The errors have the property that $\mathbf{E}(e_{i1} | \mathbf{Z}_i, \mathbf{X}_i) = 0$ and $\mathbf{E}(\epsilon_{i2} | \mathbf{Z}_i, \mathbf{X}_i) = 0$ with the variances $\Theta_{11}^* = \text{Var}(e_{i1} | \mathbf{Z}_i, \mathbf{X}_i)$, $\Theta_{22}^* = \text{Var}(\epsilon_{i2} | \mathbf{Z}_i, \mathbf{X}_i)$, and $\Theta_{12}^* = \text{Cov}(e_{i1}, \epsilon_{i2} | \mathbf{Z}_i, \mathbf{X}_i)$. Thus, each equation in the reduced-form model is a usual (high dimensional) regression model with (high dimensional) covariates \mathbf{Z}_i and \mathbf{X}_i and outcomes Y_i and D_i , respectively.

There are many methods in the literature to estimate the parameters of high dimensional regression models like the reduced-form models in (4) and (5). One approach is the scaled Lasso estimator proposed by Sun and Zhang [2012],

$$\begin{aligned} & \{\hat{\boldsymbol{\Gamma}}, \hat{\boldsymbol{\Psi}}, \hat{\Theta}_{11}\} \\ &= \underset{\boldsymbol{\Gamma} \in \mathbb{R}^{p_z}, \boldsymbol{\Psi} \in \mathbb{R}^{p_x}, \Theta_{11} \in \mathbb{R}^+}{\text{argmin}} \frac{\|\mathbf{Y} - \mathbf{Z}\boldsymbol{\Gamma} - \mathbf{X}\boldsymbol{\Psi}\|_2^2}{2n\sqrt{\Theta_{11}}} + \frac{\sqrt{\Theta_{11}}}{2} + \frac{\lambda_0}{\sqrt{n}} \left(\sum_{j=1}^{p_z} \|\mathbf{Z}_{\cdot j}\|_2 |\Gamma_j| + \sum_{j=1}^{p_x} \|\mathbf{X}_{\cdot j}\|_2 |\Psi_j| \right) \end{aligned} \quad (6)$$

for the reduced model in (4) and

$$\begin{aligned} & \{\hat{\boldsymbol{\gamma}}, \hat{\boldsymbol{\psi}}, \hat{\Theta}_{22}\} \\ &= \underset{\boldsymbol{\Gamma} \in \mathbb{R}^{p_z}, \boldsymbol{\Psi} \in \mathbb{R}^{p_x}, \Theta_{22} \in \mathbb{R}^+}{\text{argmin}} \frac{\|\mathbf{D} - \mathbf{Z}\boldsymbol{\gamma} - \mathbf{X}\boldsymbol{\psi}\|_2^2}{2n\sqrt{\Theta_{22}}} + \frac{\sqrt{\Theta_{22}}}{2} + \frac{\lambda_0}{\sqrt{n}} \left(\sum_{j=1}^{p_z} \|\mathbf{Z}_{\cdot j}\|_2 |\gamma_j| + \sum_{j=1}^{p_x} \|\mathbf{X}_{\cdot j}\|_2 |\psi_j| \right) \end{aligned} \quad (7)$$

for the reduced model in (5). The term λ_0 in both estimation problems (6) and (7) represents the penalty term in the scaled Lasso estimator and we choose $\lambda_0 = \sqrt{a_0 \log p/n}$ for some constant $a_0 > 2$; in practice, we find that setting $a_0 = 2$ or 2.05 works well. Also, we can estimate Θ_{12}^* from the estimation problems (6) and (7) by $\hat{\Theta}_{12} = 1/n \left(\mathbf{Y} - \mathbf{Z}\hat{\boldsymbol{\Gamma}} - \mathbf{X}\hat{\boldsymbol{\Psi}} \right)^\top (\mathbf{D} - \mathbf{Z}\hat{\boldsymbol{\gamma}} - \mathbf{X}\hat{\boldsymbol{\psi}})$.

Unfortunately, most penalized estimators for high dimensional regression problems are biased and the scaled Lasso estimators are no exception. In our case, using the estimates, say $\hat{\boldsymbol{\Gamma}}$ and $\hat{\boldsymbol{\gamma}}$, are biased for the parameters that they estimate $\boldsymbol{\Gamma}^*$ and

γ^* . Thankfully, recent works by Zhang and Zhang [2014], Javanmard and Montanari [2014], van de Geer et al. [2014] and Cai and Guo [2016a] allow us to debias these biased estimates. Specifically, let \mathbf{W} be the concatenated matrix of the instruments \mathbf{Z} and the covariates \mathbf{X} . Suppose we solve p_z optimization problems where the solution to each p_z optimization problem, denoted as $\hat{\mathbf{u}}^{[j]} \in \mathbb{R}^p$, $j = 1, \dots, p_z$, is

$$\hat{\mathbf{u}}^{[j]} = \underset{\mathbf{u} \in \mathbb{R}^p}{\operatorname{argmin}} \frac{1}{n} \|\mathbf{W}\mathbf{u}\|_2^2 \quad \text{s.t.} \quad \left\| \frac{1}{n} \mathbf{W}^\top \mathbf{W}\mathbf{u} - \mathbf{I}_{\cdot j} \right\|_\infty \leq \lambda_n, \quad (8)$$

with $\mathbf{I}_{\cdot j}$ denoting the j -th column of the identity matrix \mathbf{I} . The tuning parameter λ_n is chosen to be $12M_1^2 \sqrt{\log p/n}$ with M_1 defined as the largest eigenvalue of Σ^* . Let $\hat{\mathbf{U}}$ denote the concatenation of the p_z solutions to the optimization problem, i.e. $\hat{\mathbf{U}} = (\hat{\mathbf{u}}^{[1]}, \dots, \hat{\mathbf{u}}^{[p_z]})^\top$. Then, the debiased estimates of $\hat{\Gamma}$ and $\hat{\gamma}$, denoted as $\tilde{\Gamma}$ and $\tilde{\gamma}$, are

$$\tilde{\Gamma} = \hat{\Gamma} + \frac{1}{n} \hat{\mathbf{U}} \mathbf{W}^\top (\mathbf{Y} - \mathbf{Z}\hat{\Gamma} - \mathbf{X}\hat{\Psi}), \quad \tilde{\gamma} = \hat{\gamma} + \frac{1}{n} \hat{\mathbf{U}} \mathbf{W}^\top (\mathbf{D} - \mathbf{Z}\hat{\gamma} - \mathbf{X}\hat{\psi}). \quad (9)$$

In short, we used scaled Lasso along with de-biasing methods on the reduced-form models to obtain de-biased estimates $\tilde{\Gamma}$ and $\tilde{\gamma}$ of the intent-to-treat effects of the instruments on the outcome and the exposure, respectively.

3.3. Two-Stage Hard Thresholding

The second part of the confidence interval procedure deals with the problem posed by invalid IVs. Specifically, we need to select valid IVs among p_z candidate IVs that satisfy all (A1)-(A3) assumptions, that is the set \mathcal{V}^* in Definition 3. As discussed before, we do this by first, finding IVs that satisfy (A1), that is the set \mathcal{S}^* in Definition 2 consisting of j s where $\gamma_j^* \neq 0$, by thresholding the de-biased estimate $\tilde{\gamma}$

$$\tilde{\mathcal{S}} = \left\{ j : |\tilde{\gamma}_j| \geq \frac{\sqrt{\hat{\Theta}_{22}} \|\mathbf{W}\hat{\mathbf{u}}^{[j]}\|_2}{\sqrt{n}} \sqrt{\frac{a_0 \log p_z}{n}} \right\}. \quad (10)$$

where $\tilde{\mathcal{S}}$ denotes an estimate of \mathcal{S}^* . The threshold is based on the noise level of $\tilde{\gamma}_j$ in (9) (represented by $\sqrt{\hat{\Theta}_{22}} \|\mathbf{W}\hat{\mathbf{u}}^{[j]}\|_2/n$), adjusted by dimensionality of the instrument size (represented by $\sqrt{a_0 \log p_z}$).

The second thresholding step involves selecting IVs that satisfy (A2) and (A3). Specifically, by Definition 1, the set of instruments that satisfy (A2) and (A3) are those j s where $\pi_j^* = 0$. Consequently, to estimate π^* , we take each instrument j in $\tilde{\mathcal{S}}$ that satisfy (A1) and we define $\hat{\beta}^{[j]}$ to be a “pilot” estimate of β^* by using this IV and dividing the reduced-form estimates, i.e. $\hat{\beta}^{[j]} = \tilde{\Gamma}_j / \tilde{\gamma}_j$, and $\hat{\pi}^{[j]}$ to be the estimate of π^* using this j th instrument’s estimate of β^* , i.e. $\hat{\pi}^{[j]} = \tilde{\Gamma} - \hat{\beta}^{[j]} \tilde{\gamma}$; we also construct corresponding pilot estimates of σ^2 , i.e. $\hat{\sigma}^{2[j]} = \hat{\Theta}_{11} + (\hat{\beta}^{[j]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[j]} \hat{\Theta}_{12}$. Then, for each $\hat{\pi}^{[j]}$ in $j \in \tilde{\mathcal{S}}$, we threshold each element of $\hat{\pi}^{[j]}$ to create the thresholded

estimate $\tilde{\pi}^{[j]}$,

$$\tilde{\pi}_k^{[j]} = \hat{\pi}_k^{[j]} \mathbf{1} \left(k \in \tilde{\mathcal{S}} \cap |\hat{\pi}_k^{[j]}| \geq a_0 \sqrt{\widehat{\sigma}^{2[j]}} \frac{\|\mathbf{W}(\hat{\mathbf{u}}^{[k]} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_j} \hat{\mathbf{u}}^{[j]})\|_2}{\sqrt{n}} \sqrt{\frac{\log p_z}{n}} \right) \quad (11)$$

for all $1 \leq k \leq p_z$. Each thresholded estimate $\tilde{\pi}^{[j]}$ is obtained by looking at the elements of the un-thresholded estimate, $\hat{\pi}^{[j]}$, and examining whether each element of it exceeds the noise threshold, denoted by the term $\sqrt{\widehat{\sigma}^{2[j]}} \|\mathbf{W}(\hat{\mathbf{u}}^{[k]} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_j} \hat{\mathbf{u}}^{[j]})\|_2/n$, adjusting for the multiplicity of the selection procedure by the term $a_0 \sqrt{\log p_z}$. Among the $|\tilde{\mathcal{S}}|$ candidate estimates of π^* based on each instrument in $\tilde{\mathcal{S}}$, i.e. $\tilde{\pi}^{[j]}$, and we choose $\tilde{\pi}^{[j]}$ with most valid instruments, or equivalently choose $j^* \in \tilde{\mathcal{S}}$ where $j^* = \operatorname{argmin} \|\tilde{\pi}^{[j]}\|_0$; if there is a non-unique solution, we choose $\tilde{\pi}^{[j]}$ with the smallest ℓ_1 norm, the closest convex norm of ℓ_0 .

Intuitively, the second-stage thresholding selects the invalid IVs and valid IVs as follows. Among the $|\tilde{\mathcal{S}}|$ pilot estimates $\tilde{\pi}^{[j]}$, the best estimate of π^* is the one that uses a valid IV from the set $\tilde{\mathcal{S}}$. In particular, if the j th pilot estimate is actually based on a valid IV, then all the invalid IVs will be included in the support of the thresholded estimate $\tilde{\pi}^{[j]}$ because their π^* will be away from zero and all the valid instruments will be excluded from the support because their π^* are zero. On the other hand, if the j th pilot estimate is based on an invalid IV, the pilot estimate $\tilde{\pi}^{[j]}$ will be biased in the sense that the valid IVs will no longer have $\tilde{\pi}^{[j]}$ that will be thresholded to zero and most of the elements of $\tilde{\pi}^{[j]}$ will be away from zero. Consequently, many IVs will be declared invalid based on $\tilde{\pi}^{[j]}$ and when we minimize with respect to the number of non-zero elements of the vector, i.e. $\min \|\tilde{\pi}^{[j]}\|_0$ among all pilot estimates, we should be able to select the best estimate of π^* . We remark that the latter ℓ_0 minimization is reminiscent of Theorem 1 in Kang et al. [2016b] where a necessary and sufficient condition for identification of β^* under invalid instruments is by looking at the largest subset of valid instruments that converge on a unique (i.e. identified) value; the search for the largest subset of valid IVs is essentially a minimization of ℓ_0 norm, which counts the number of invalid IVs, and hence, there is some sense that our procedure is both sufficient and necessary way to estimate β^* .

Finally, we note that it is crucial to construct pilot estimates of π^* from the IVs in the first thresholding step, that is $\tilde{\mathcal{S}}$, as each of these IVs represent strong IVs and have non-zero effects on the exposure; using IVs that are not in $\tilde{\mathcal{S}}$ may lead to poor estimates of the direct effect on the outcome since a redundant instrument j , whose true γ_j^* is zero, can lead to a large, unstable value of $|\tilde{\gamma}_k/\tilde{\gamma}_j|$ and the threshold value in (11), which will make it difficult to distinguish truly invalid IVs from noise.

3.4. Confidence Interval Estimation

After the two thresholding steps, we estimate the set of valid instruments $\tilde{\mathcal{V}} \subseteq \{1, \dots, p_z\}$ as those elements of $\tilde{\boldsymbol{\pi}}^{[j^*]}$ that are zero,

$$\tilde{\mathcal{V}} = \tilde{\mathcal{S}} \setminus \text{supp}(\tilde{\boldsymbol{\pi}}^{[j^*]}) \quad (12)$$

Then, using the estimated $\tilde{\mathcal{V}}$, we obtain our estimate of β^*

$$\hat{\beta} = \frac{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\Gamma}_j}{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2}, \quad (13)$$

along with an estimate of its standard error

$$\hat{V} = \frac{\left\| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \frac{1}{\sqrt{n}} \mathbf{W} \hat{\mathbf{u}}^{[j]} \right\|_2^2}{\left(\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 \right)^2} \hat{\sigma}^2 \quad \text{and} \quad \hat{\sigma}^2 = \hat{\Theta}_{11} + \hat{\beta}^2 \hat{\Theta}_{22} - 2\hat{\beta} \hat{\Theta}_{12}, \quad (14)$$

and the usual form for the confidence interval for β^* ,

$$\left(\hat{\beta} - z_{1-\alpha/2} \sqrt{\hat{V}/n}, \quad \hat{\beta} + z_{1-\alpha/2} \sqrt{\hat{V}/n} \right), \quad (15)$$

where $z_{1-\alpha/2}$ is the $1 - \alpha/2$ quantile of the standard normal distribution.

In the equations for $\hat{\beta}$ and the standard error \hat{V} , we see some familiar expressions from the traditional IV literature. First, $\hat{\beta}$ has the “correct” form in that if, by chance, we correctly estimated the set of valid instruments \mathcal{V}^* and our debiased estimates of the reduced-form parameters, $\tilde{\boldsymbol{\Gamma}}$ and $\tilde{\boldsymbol{\gamma}}$, are perfect estimates of the reduced-form parameters, our estimate of β^* in (13) would become $\hat{\beta} = \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\Gamma}_j / \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 = \sum_{j \in \mathcal{V}^*} \gamma_j^{*2} \beta^* / \sum_{j \in \mathcal{V}^*} \gamma_j^{*2} = \beta^*$. Hence, our estimator in (13) would identify β^* . Clearly, we would never have a perfect estimate of the set \mathcal{V}^* or of the reduced-form parameters in finite sample and Section 4 describes the properties of our estimate $\hat{\beta}$ under these uncertainties. Second, in the standard error formula in (14), the $\hat{\Theta}_{11} + \hat{\beta}^2 \hat{\Theta}_{22} - 2\hat{\beta} \hat{\Theta}_{12}$ is of the similar form to the usual IV estimator of σ^2 , variance of the error term in our original model (2). But, our standard error estimator is scaled by terms that depend on the estimated set of valid instruments $\tilde{\mathcal{V}}$.

3.5. Special Case of Procedure 1: Valid IVs After Controlling for High Dimensional Covariates

To better understand the components of our inference procedure 1, it is instructive to go through some specific cases of estimating β^* that is common in the literature as these special cases can greatly simplify the procedure and remove unnecessary components. The first case is when the instruments are assumed to be valid (i.e. no direct effect and no unmeasured confounding) after conditioning on high dimensional

covariates. This setup was considered in Gautier and Tsybakov [2011], Belloni et al. [2012], Fan and Liao [2014], Chernozhukov et al. [2015]. Under this case, our procedure doesn't have to go through STEP 2b, the estimation of π^* , illustrated in Section 3.3. Instead, we can simply replace STEP 2b with $\tilde{\mathcal{V}} = \tilde{\mathcal{S}}$ and the resulting estimator for β^* is

$$\hat{\beta}_H = \frac{\sum_{j \in \tilde{\mathcal{S}}} \tilde{\gamma}_j \tilde{\Gamma}_j}{\sum_{j \in \tilde{\mathcal{S}}} \tilde{\gamma}_j^2}. \quad (16)$$

The corresponding confidence interval for β^* would be

$$\left(\hat{\beta}_H - z_{1-\alpha/2} \sqrt{\hat{V}_H/n}, \quad \hat{\beta}_H + z_{1-\alpha/2} \sqrt{\hat{V}_H/n} \right). \quad (17)$$

Here, \hat{V}_H is \hat{V} in (14) except we replace $\tilde{\mathcal{V}} = \tilde{\mathcal{S}}$ and $\hat{\beta}$ with $\hat{\beta}_H$.

3.6. Special Case of Procedure 1: Invalid Instruments After Controlling for Low Dimensional Covariates

The second special case worth examining is the problem of invalid instruments after controlling for low dimensional covariates. While the plausibility of candidate IVs adherence to assumptions (A1)-(A3), especially (A3), is higher with many covariates, the low dimensional setting has recently received much attention and is discussed in Bowden et al. [2015, 2016], Burgess et al. [2016], Kang et al. [2016a], and Windmeijer et al. [2016]. As we will see below, our procedure simplifies greatly and with a minor modification, our estimator, unlike the estimators proposed in said prior literature, achieves optimal performance.

Specifically, under the low-dimensional scenario, there is no need to use the debiased scaled lasso in STEP 1 of Procedure 1. Instead, we can replace STEP 1 with the simple ordinary least square (OLS) estimates of the reduced-forms, $(\tilde{\Gamma}, \tilde{\Psi})^\top = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{Y}$ and $(\tilde{\gamma}, \tilde{\psi})^\top = (\mathbf{W}^\top \mathbf{W})^{-1} \mathbf{W}^\top \mathbf{D}$, and of the covariance terms $\hat{\Theta}_{11} = \|\mathbf{Y} - \mathbf{Z}\hat{\Gamma} - \mathbf{X}\hat{\Psi}\|_2^2/n$, $\hat{\Theta}_{22} = \|\mathbf{D} - \mathbf{Z}\hat{\gamma} - \mathbf{X}\hat{\psi}\|_2^2/n$, and $\hat{\Theta}_{12} = (\mathbf{Y} - \mathbf{Z}\hat{\Gamma} - \mathbf{X}\hat{\Psi})^\top (\mathbf{D} - \mathbf{Z}\hat{\gamma} - \mathbf{X}\hat{\psi})/n$. As a result of using OLS in STEP 1, we need to replace $\hat{\mathbf{u}}^{[j]}$ from (8) with $\hat{\mathbf{u}}^{[j]} = (\hat{\Sigma})_j^{-1}$, $\hat{\Sigma} = \mathbf{W}^\top \mathbf{W}/n$, and replace the log terms in our thresholds in (10) and (11) from $\sqrt{\log p_z}$ to $\sqrt{\log n}$.

We can then proceed to use the estimator defined in (13). Alternatively, we can use a modified version of $\hat{\beta}$, denoted as $\hat{\beta}_E$, using the weighting matrix $\mathbf{A} = \hat{\Sigma}_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}} - \hat{\Sigma}_{\tilde{\mathcal{V}}, \tilde{\mathcal{V}}^c} \hat{\Sigma}_{\tilde{\mathcal{V}}^c, \tilde{\mathcal{V}}^c}^{-1} \hat{\Sigma}_{\tilde{\mathcal{V}}^c, \tilde{\mathcal{V}}}$

$$\hat{\beta}_E = \frac{\tilde{\gamma}_{\tilde{\mathcal{V}}}^\top \mathbf{A} \tilde{\Gamma}_{\tilde{\mathcal{V}}}}{\tilde{\gamma}_{\tilde{\mathcal{V}}}^\top \mathbf{A} \tilde{\gamma}_{\tilde{\mathcal{V}}}} \quad (18)$$

along with the estimated standard error $\hat{V}_E = \hat{\sigma}^2 / \tilde{\gamma}_{\tilde{\mathcal{V}}}^\top \mathbf{A} \tilde{\gamma}_{\tilde{\mathcal{V}}}$ where $\hat{\sigma}^2 = \hat{\Theta}_{11} + \hat{\beta}_E^2 \hat{\Theta}_{22} - 2\hat{\beta}_E \hat{\Theta}_{12}$ and confidence interval

$$\left(\hat{\beta}_E - z_{1-\alpha/2} \sqrt{\hat{V}_E/n}, \quad \hat{\beta}_E + z_{1-\alpha/2} \sqrt{\hat{V}_E/n} \right). \quad (19)$$

Note that $\hat{\beta}_E$ in (18) is reduced to $\hat{\beta}$ in (13) by setting $A = I$. As we will see in Section 4.1, our estimator $\hat{\beta}_E$, compared to other estimators in prior work, achieves optimal performance in the sense that our performance is asymptotically identical to the TSLS estimator for β^* that knows which IVs are valid a priori, i.e. the set \mathcal{V}^* .

4. Theoretical Results

In this section, we investigate the properties of the confidence interval proposed in Procedure 1. We first consider in Section 4.1 the coverage property in the case of invalid IVs with low dimensional covariates where p_x and p_z are fixed. In Section 4.2, we establish the coverage property for the general case, invalid IVs even after controlling for many covariates.

4.1. Invalid IVs After Controlling for Low Dimensional Covariates

We state the following mild assumption commonly used in the invalid IV literature

(IN1) (50% Rule) The number of valid IVs is more than half of the number of non-redundant IVs, that is $|\mathcal{V}^*| > \frac{1}{2}|\mathcal{S}^*|$.

We denote the assumption as “IN” since the assumption is specific to the case of invalid IVs. Assumption (IN1) is the generalization of the 50% rule in Kang et al. [2016b] and Han [2008] in the presence of possibly redundant IVs. In a nutshell, (IN1) states that if the number of invalid instruments is not too large, then we can detect the invalid IVs from valid IVs, without knowing a priori which IVs are valid or invalid; see Kang et al. [2016b] for a detailed discussion of this assumption and how this type of proportion-based assumption is a necessary component for identification of model parameters under invalid instruments.

Under the 50% assumption alone, Theorem 1 states that we can show that our procedure produces confidence intervals with desired coverage and optimal length in low dimensional settings where p_x and p_z are fixed.

THEOREM 1. *Suppose that the assumption (IN1) holds. Then the following property holds for the estimator $\hat{\beta}_E$ defined in (18),*

$$\sqrt{n} \left(\hat{\beta}_E - \beta^* \right) \xrightarrow{d} N \left(0, \frac{\sigma^2}{\gamma_{\mathcal{V}^*}^{*\top} \left(\Sigma_{\mathcal{V}^* \mathcal{V}^*}^* - \Sigma_{\mathcal{V}^* (\mathcal{V}^*)^c}^* \Sigma_{(\mathcal{V}^*)^c (\mathcal{V}^*)^c}^{*-1} \Sigma_{(\mathcal{V}^*)^c \mathcal{V}^*}^* \right) \gamma_{\mathcal{V}^*}^*} \right). \quad (20)$$

Consequently, the confidence interval given in (19) has asymptotically coverage probability $1 - \alpha$, i.e.,

$$\mathbf{P} \left\{ \beta \in \left(\hat{\beta}_E - z_{1-\alpha/2} \sqrt{\hat{V}_E/n}, \quad \hat{\beta}_E + z_{1-\alpha/2} \sqrt{\hat{V}_E/n} \right) \right\} \rightarrow 1 - \alpha. \quad (21)$$

We note that the proposed estimator $\widehat{\beta}_E$ has the same asymptotic variance as the oracle TSLS estimator with the prior knowledge of \mathcal{V}^* , which is shown to be efficient under the homoskedastic variance assumption (Theorem 5.2 in Wooldridge [2010]); consequently, our confidence interval asymptotically performs like the oracle TSLS confidence interval and is of optimal length. But, unlike TSLS, we achieve this oracle performance without prior knowledge of \mathcal{V}^* . We also note the estimators proposed in prior work [Bowden et al., 2015, 2016, Burgess et al., 2016, Kang et al., 2016a, Windmeijer et al., 2016] do not achieve oracle performance and TSLS-like efficiency.

4.2. Invalid IVs After Controlling for High Dimensional Covariates

We now consider the coverage property for the general case, invalid IVs even after controlling for many confounders. We first introduce the usual regularity assumptions used in high-dimensional statistical inference [Bickel et al., 2009, Bühlmann and van de Geer, 2011, Cai and Guo, 2016a].

- (R1) (Coherence): The matrix Σ^* satisfies $1/M_1 \leq \lambda_{\min}(\Sigma^*) \leq \lambda_{\max}(\Sigma^*) \leq M_1$ for some constant $M_1 > 1$ and has bounded sub-Gaussian norm.
- (R2) (Normality): The error terms in (4) and (5) follow a bivariate normal distribution.
- (R3) (Global IV Strength): The IVs are globally strong with $\|\gamma_{\mathcal{V}^*}^*\|_2 = \sqrt{\sum_{j \in \mathcal{V}^*} \gamma_j^2} \geq \delta \gg s_{z1} \log p / \sqrt{n}$, where \mathcal{V}^* is the set of valid IVs defined in Definition 3.

Assumption (R1) places a condition on the spectrum of the design matrix \mathbf{W} and the tail distribution of $\mathbf{W}_{i\cdot}$, which is related to the restricted eigenvalue condition in Bickel et al. [2009]. For simplicity, we also assume that the sub-Gaussian norm of $\mathbf{W}_{i\cdot}$ is upper bounded by M_1 , that is, $\sup_{\mathbf{v} \in S^{p-1}} \sup_{q \geq 1} (\mathbf{E} |\mathbf{v}^\top \mathbf{W}_{i\cdot}|^q / q)^{\frac{1}{q}} \leq M_1$ where S^{p-1} is the unit sphere in \mathbb{R}^p ; see Vershynin [2012] for details on sub-Gaussian random variables and bounds. Assumption (R2) states that the errors (e_{i1}, e_{i2}) are bivariate normal. Here, we make the normality assumption out of simplicity, similar to the work on inference in weak IV literature where error terms are typically assumed to be normal (e.g. Section 2 of Moreira [2003] and Section 2.2.1 Andrews et al. [2007]). Finally, Assumption (R3) states that the global strength of instruments, measured by the ℓ_2 norm of γ^* among valid IVs \mathcal{V}^* , is bounded away from zero. This type of global strength assumption is commonly made in the IV literature under the guise as a concentration parameter, which is a measure of strength of the instrument (see Section 5 for details) and is the weighted ℓ_2 norm of $\gamma_{\mathcal{V}^*}^*$, and is often referred to as “traditional/strong” asymptotics [Stock et al., 2002, Wooldridge, 2010]. Recent works by Belloni et al. [2012] and Chernozhukov et al. [2015], which considered the setting where all IVs were valid after conditioning on high dimensional covariates, also make this type of assumption, specifically condition SM in Belloni et al. [2012] and condition RF in the supplementary materials of Chernozhukov et al. [2015]. Essentially, both these works require $\|\gamma^*\|_2$ to be

bounded away from zero by a constant and are actually stronger than our (R3). In practice, (R3) is satisfied so long as there is at least one IV that has a constant non-zero effect on the treatment, or a non-zero effect that doesn't diminish with sample size. However, if the IVs are arbitrary weak in the sense of Staiger and Stock [1997], then (R3), let alone the said assumptions in high dimensional valid IV literature [Belloni et al., 2012], do not hold, and we leave this as a future topic of research to deal with arbitrary weak IVs in invalid IV settings.

Section A in the supplementary materials shows that if the IVs are valid after conditioning on many covariates, then Assumptions (R1)-(R3) are sufficient for the confidence interval proposed in (17) to have correct coverage. However, when IVs are invalid after conditioning on said controls, we need to make two additional assumptions that are not in the usual high dimensional inference or instrumental variables literature and may be of theoretical interest in future work.

(IN2) (Individual IV Strength) For IVs in \mathcal{S}^* , $\min_{j \in \mathcal{S}^*} |\gamma_j^*| \geq \delta_{\min} \gg \sqrt{\log p/n}$.

(IN3) (Strong violation) Among IVs in the set $\mathcal{S}^* \setminus \mathcal{V}^*$, we have

$$\min_{j \in \mathcal{S}^* \setminus \mathcal{V}^*} \left| \frac{\pi_j^*}{\gamma_j^*} \right| \geq \frac{12(1 + |\beta^*|)}{\delta_{\min}} \sqrt{\frac{M_1 \log p_z}{\lambda_{\min}(\Theta^*)n}}. \quad (22)$$

Assumption (IN2) requires individual IV strength to be bounded away from zero so that all IVs in selected $\tilde{\mathcal{S}}$ are strong. This assumption is needed primarily for cleaner technical exposition and our simulation studies in Section 5 and Section C in the supplementary materials demonstrate that (IN2) is largely unnecessary for our confidence interval to guarantee coverage. In the literature, (IN2) is similar to the “beta-min” condition assumption in high dimensional linear regression without IVs, with the exception that this condition is not imposed on our inferential quantity of interest, β^* . Also, (IN2) is different from Assumption (R3), where (R3) only requires the global IV strength to be bounded away from zero. Next, Assumption (IN3) requires the ratios π_j^*/γ_j^* for invalid IVs to be large and this assumption is needed to correctly identify IVs that violate (A2) and (A3). Specifically, for any IV with $|\pi_j^*/\gamma_j^*|$ being non-zero but small, it's difficult to distinguish such a weakly invalid IV from valid IVs where $\pi_j^*/\gamma_j^* = 0$. If a weakly invalid IV is mistakenly declared as valid, the bias from this mistake is of the order $\sqrt{\log p_z/n}$, which has consequences, not for consistency of the point estimate, but for a \sqrt{n} confidence interval; see Theorem 2 and Section 7 for more discussions.

With (R1)-(R3) and (IN1)-(IN3), our general Procedure 1 produces a consistent and asymptotic normal estimate of β^* even if IVs are invalid after conditioning on high dimensional controls.

THEOREM 2. *Suppose the assumptions (R1) – (R3) and (IN1) – (IN2) hold. As $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$, with probability larger than $1 - c(p^{-c} + \exp(-cn))$,*

$$|\hat{\beta} - \beta^*| \leq C \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}}, \quad (23)$$

where $c, C > 0$ are constants independent of n and p . In addition, if (IN3) holds, we have

$$\sqrt{n}(\hat{\beta} - \beta^*) = T^{\beta^*} + \Delta^{\beta^*} \quad (24)$$

where $T^{\beta^*} \mid \mathbf{W} \sim N(0, V)$, $V = \sigma^2 / \left(\sum_{j \in \mathcal{V}^*} (\gamma_j^*)^2 \right)^2 \left\| \sum_{j \in \mathcal{V}^*} \gamma_j^* \mathbf{W} \hat{\mathbf{u}}^{[j]} / \sqrt{n} \right\|_2^2$, and $\Delta^{\beta^*} / \sqrt{V} \xrightarrow{p} 0$ as $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$. Consequently, the confidence interval given in (15) has asymptotically coverage probability $1 - \alpha$, i.e.,

$$\mathbf{P} \left\{ \beta^* \in \left(\hat{\beta} - z_{1-\alpha/2} \sqrt{\hat{V}/n}, \hat{\beta} + z_{1-\alpha/2} \sqrt{\hat{V}/n} \right) \right\} \rightarrow 1 - \alpha. \quad (25)$$

In Theorem 2, the consistency of our estimator in (23) is established without (IN3) because the bias term $\sqrt{\log p_z/n}$ discussed above is still going to zero. Section C.2 of the supplementary materials numerically demonstrates this by showing that the estimator $\hat{\beta}$ is still consistent even when (IN3) is violated. However, for \sqrt{n} asymptotic normality, Theorem 2 requires Assumption (IN3) to eliminate said bias so that our confidence interval (15) has correct coverage even with invalid IVs and high dimensional controls.

5. Simulation

5.1. Setup

In addition to the theoretical analysis of our method in Section 4, we also conduct a simulation study to investigate (i) the performance of our method and other comparators and (ii) sensitivity of our method to violations of the regularity assumptions mentioned above, most notably (IN2) and (IN3). The data generating process for the simulation follows the models (2) and (3) in Section 2.2 with $p_z = 100$ instruments and $p_x = 150$ covariates where \mathbf{W}_i is a multivariate normal with mean zero and covariance $\Sigma_{ij}^* = 0.5^{|i-j|}$ for $1 \leq i, j \leq 250$. The parameters for the models are: $\beta^* = 1$, $\phi^* = (0.6, 0.7, 0.8, \dots, 1.5, 0, 0, \dots, 0) \in \mathbb{R}^{150}$ so that $s_{x1} = 10$, $\psi^* = (1.1, 1.2, 1.3, \dots, 2.0, 0, 0, \dots, 0) \in \mathbb{R}^{150}$ so that $s_{x2} = 10$, and variance-covariance of the error terms are $\text{Var}(\epsilon_{i1}) = \text{Var}(\epsilon_{i2}) = 1.5$, and $\text{Cov}(\epsilon_{i1}, \epsilon_{i2}) = 0.75$. Instruments that satisfy Assumption (A1) are $\mathcal{S}^* = \{1, \dots, 7\}$ and instruments that satisfy all three IV assumptions (A1)-(A3) are $\mathcal{V}^* = \{1, 2, 3, 4, 5\}$; thus instruments 6 and 7 only satisfy (A1), but do not satisfy (A2) and (A3). We fix these values throughout the entire simulation study.

The parameters we vary in the simulation study are: the sample size n , the strength of IVs via γ^* , and violations of (A2) and (A3) via π^* . For sample size, we let $n = (100, 200, 300, 1000, 3000)$. For IV strength, we set $\gamma_{\mathcal{V}^*}^* = K(1, 1, 1, 1, \rho_1)$ and $\gamma_{\mathcal{S}^* \setminus \mathcal{V}^*}^* = K(1, 1)$ and $\gamma_{(\mathcal{S}^*)^c} = \mathbf{0}$, where we vary K (to be discussed later) and $\rho_1 = (0, 0.1, 0.2)$ across simulations. The value K controls the global strength of instruments, with higher $|K|$ indicating strong instruments in a global sense. The value ρ_1 controls the relative individual strength of instruments, specifically between the first four instruments in \mathcal{V}^* and the fifth instrument. For example, $\rho_1 = 0.2$ implies that the fifth IV's individual strength is only 20% of the other four

valid instruments, i.e IVs 1 to 4. Also, varying ρ_1 would simulate the adherence of regularity assumption (IN2).

To specify K across simulations, we introduce a quantity we call the oracle concentration parameter (OCP) denoted as $C(\gamma^*, \mathcal{V}^*, n)$

$$C(\gamma^*, \mathcal{V}^*, n) = n \frac{\gamma_{\mathcal{V}^*}^{*\top} \left(\Sigma_{\mathcal{V}^* \mathcal{V}^*}^* - \Sigma_{\mathcal{V}^* (\mathcal{V}^*)^c}^* \Sigma_{(\mathcal{V}^*)^c (\mathcal{V}^*)^c}^{*-1} \Sigma_{(\mathcal{V}^*)^c \mathcal{V}^*}^* \right) \gamma_{\mathcal{V}^*}^*}{|\mathcal{V}^*| \Theta_{22}^*}, \quad (26)$$

where Σ_{IJ}^* denotes the submatrix containing Σ_{ij}^* for $i \in I$ and $j \in J$ and $\gamma_{\mathcal{V}^*}^*$ denotes the subvector containing γ_j^* for $j \in \mathcal{V}^*$. In Section C.1 of the Supplementary Materials, we discuss the OCP and its relation to the usual concentration parameter in the IV literature [Stock and Wright, 2000]. In short, we define the OCP because the usual concentration parameter can be misleading when there are unknown redundant and invalid instruments and the OCP serves as a proxy for the usual concentration parameter.

Having defined the OCP, we can specify K as a function of n and $C(\gamma^*, \mathcal{V}^*, n)$. Specifically, if n is set at a baseline of 100 and the simulation parameters $\mathcal{V}^*, \rho_1, \Sigma^*$ and Θ_{22}^* are specified as above, we can find K for a particular value of the expected oracle concentration parameter $C(\gamma^*, \mathcal{V}^*, 100)$. Thus, by varying $C(\gamma^*, \mathcal{V}^*, 100) = (50, 100, 150, 200, 250, 500, 1000)$, we vary K .

Finally, we vary π^* , which controls the validity of the IVs by defining $\pi_j^* = \rho_2 \gamma_j^*$ for $j = 6, 7$ and $\pi_j^* = 0$ for all other j so that ρ_2 controls the magnitude of the violation of IV assumptions (A2) and (A3) from the 6th and 7th instruments. In the ideal case, we would have $\rho_2 = 0$ so that $\mathcal{S}^* = \mathcal{V}^* = \{1, 2, 3, 4, 5, 6, 7\}$. But, $\rho_2 \neq 0$ implies that the last two instruments do not satisfy (A2) and (A3). As such, we vary π^* by varying $\rho_2 = (0, 1, 2)$. Also, varying ρ_2 would simulate the adherence of regularity assumption (IN3).

In summary, we vary n , the strength of IVs via γ^* , and violations of (A2) and (A3) via π^* in our simulation study, with ρ_1 and ρ_2 simulating the adherence to the new regularity assumptions in the paper, (IN2), and (IN3), respectively. For the setting $n \leq p$, we compare our procedure to $\hat{\beta}_H$, which assumes IVs are valid. For the setting $n \geq p$, we add two additional comparators, the two-stage least squares (TSLS) and OLS. TSLS is the most popular IV method where one regresses \mathbf{D} on \mathbf{Z} and \mathbf{X} , and uses the predicted value of \mathbf{D} in the regression of \mathbf{Y} on \mathbf{X} and \mathbf{D} . Note that the way we implement TSLS mimics most practitioners' use of TSLS by simply assuming all the instruments \mathbf{Z} are valid. OLS is defined as where one regresses \mathbf{Y} on \mathbf{D} and \mathbf{X} . OLS will be biased because of confounding on \mathbf{D} . Finally, for both low and high dimensional settings, we have the oracle TSLS where an oracle provides us with the true set of valid IVs, which will not occur in practice. Our simulations are repeated 500 times.

5.2. Results

We present the most representative results from our simulation study; all the simulations are in Sections C.2, C.3 and C.4 of the supplementary materials. First,

Figure 1 considers the high dimensional setting with $n = 200$ and three comparators, our procedure $\hat{\beta}$ that is robust to invalid IVs, our procedure $\hat{\beta}_H$ that assumes all valid IVs, and the oracle TSLS. Columns “Weak” and “Strong” in the figure represent cases where $\rho_1 = 0.2$ and $\rho_1 = 0$, respectively. Columns “Valid” and “Invalid” represent cases where $\rho_2 = 0$ and $\rho_2 = 2$, respectively. The row “MAE” in the figure represents the median absolute error of the estimators, which measures the performance of the point estimators. The row “Coverage” represents the coverage performance of the confidence intervals. Finally, the row “Length” represents the average length of confidence intervals across simulations.

Both estimators $\hat{\beta}$ and $\hat{\beta}_H$ perform well in terms of estimation accuracy, coverage and length of confidence intervals and have similar performance to the benchmark, $\hat{\beta}_{\text{oracle}}$, when all the instruments are valid (i.e. first and second columns of Figure 1). For example, in the MAE and length plots, the solid lines, which represent our estimator, the dashed lines, which represent our estimator assuming all valid IVs after conditioning on covariates, and the dotted lines, which represent the oracle, overlap with each other. However, if the instruments are invalid (i.e. the third and fourth columns of Figure 1), $\hat{\beta}_H$ is not consistent and loses coverage, which makes sense since $\hat{\beta}_H$ assumes all the IVs are valid after conditioning. However, our proposed estimator $\hat{\beta}$ allows for possibly invalid instruments and performs as well as the oracle in terms of estimation accuracy and coverage. The average length of our robust confidence interval is only slightly larger than that of the oracle.

Figure 2 represents the same setting as Figure 1 except we now consider a larger sample size $n = 1000$. Even though n is larger than p , we still consider this to be in the many controls/high dimensional setting because the ratio of p to n is away from zero at $1/4$. As expected, the estimators $\hat{\beta}$ and $\hat{\beta}_H$ along with the traditional TSLS estimator perform similarly to the oracle benchmark in terms of estimation accuracy, coverage and the length of confidence intervals when all the instruments are actually valid. For example, in the MAE plot of Figure 2, the solid, dashed, green and dotted lines, representing $\hat{\beta}$, $\hat{\beta}_H$, TSLS and the oracle, respectively, overlap with each other. Note that OLS cannot deal with confounding and hence, produces a biased estimate. However, when the instruments are invalid, the traditional TSLS estimator and $\hat{\beta}_H$ are biased and fail to have the correct coverage. In contrast, the proposed estimator $\hat{\beta}$ performs as well as the oracle estimator in terms of estimation accuracy and coverage, with the length of the proposed estimator being slightly longer than that for the oracle.

Finally, Figure 3 represents the setting where invalid instruments are present after conditioning on low dimensional covariates where $p_z = 9$ and $p_x = 10$ so that no coefficients for ϕ^* and ψ^* are zero and the sample size is $n = 1000$. If we use the estimator $\hat{\beta}_E$ defined in (18) and the confidence interval (19), the proposed procedure performs almost the same as the oracle in terms of accuracy, coverage property and length, which supports the theory established in Theorem 1. Note that the performance of our procedure under the low dimensional setting with invalid IVs does not rely on assumptions (R1)-(R3) and, more importantly, (IN2)-(IN3).

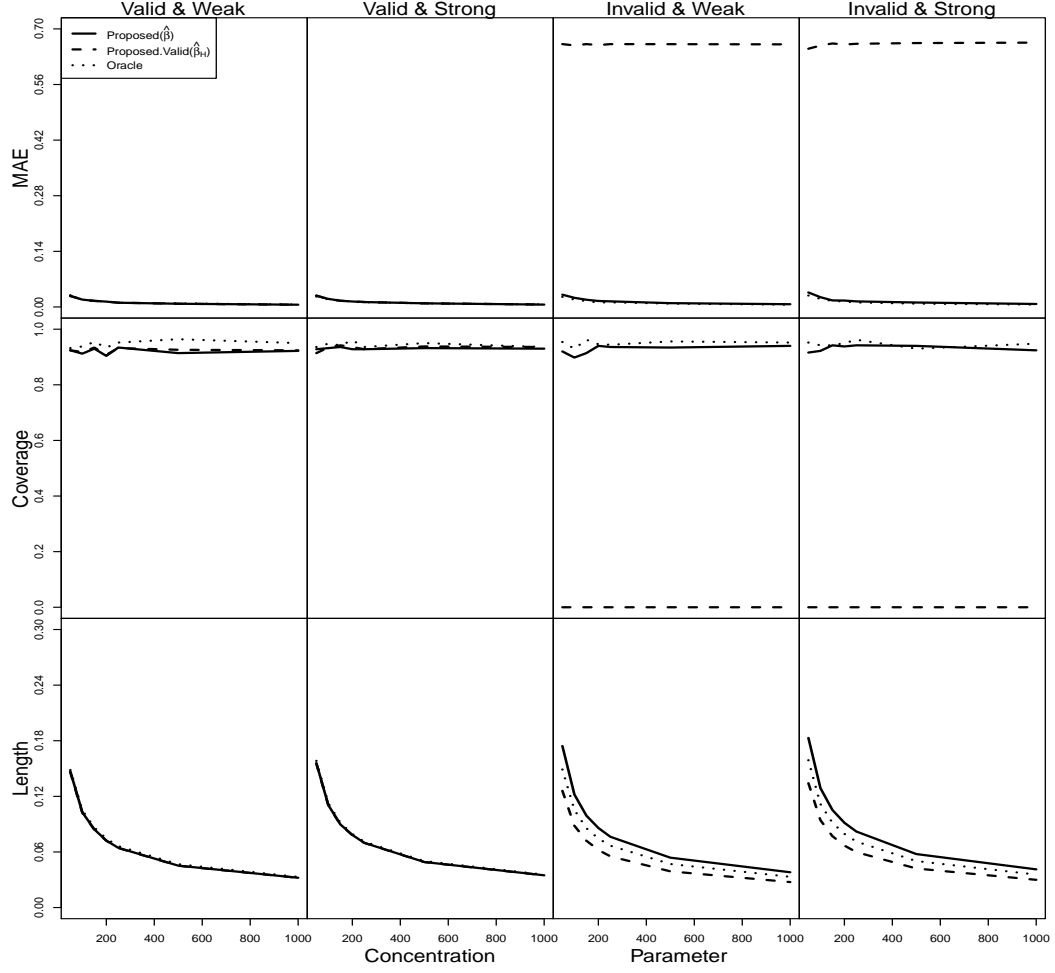


Fig. 1: Comparison of different methods when $p_z = 100$, $p_x = 150$ and $n = 200$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of the confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Valid & Weak represents the case $\rho_1 = 0.2$ and $\rho_2 = 0$. The column labeled with Valid & Strong represents the case $\rho_1 = 0$ and $\rho_2 = 0$. The column labeled with Invalid & Weak represents the case $\rho_1 = 0.2$ and $\rho_2 = 2$. Finally, the column labeled with Invalid & Strong represents the case $\rho_1 = 0$ and $\rho_2 = 2$.

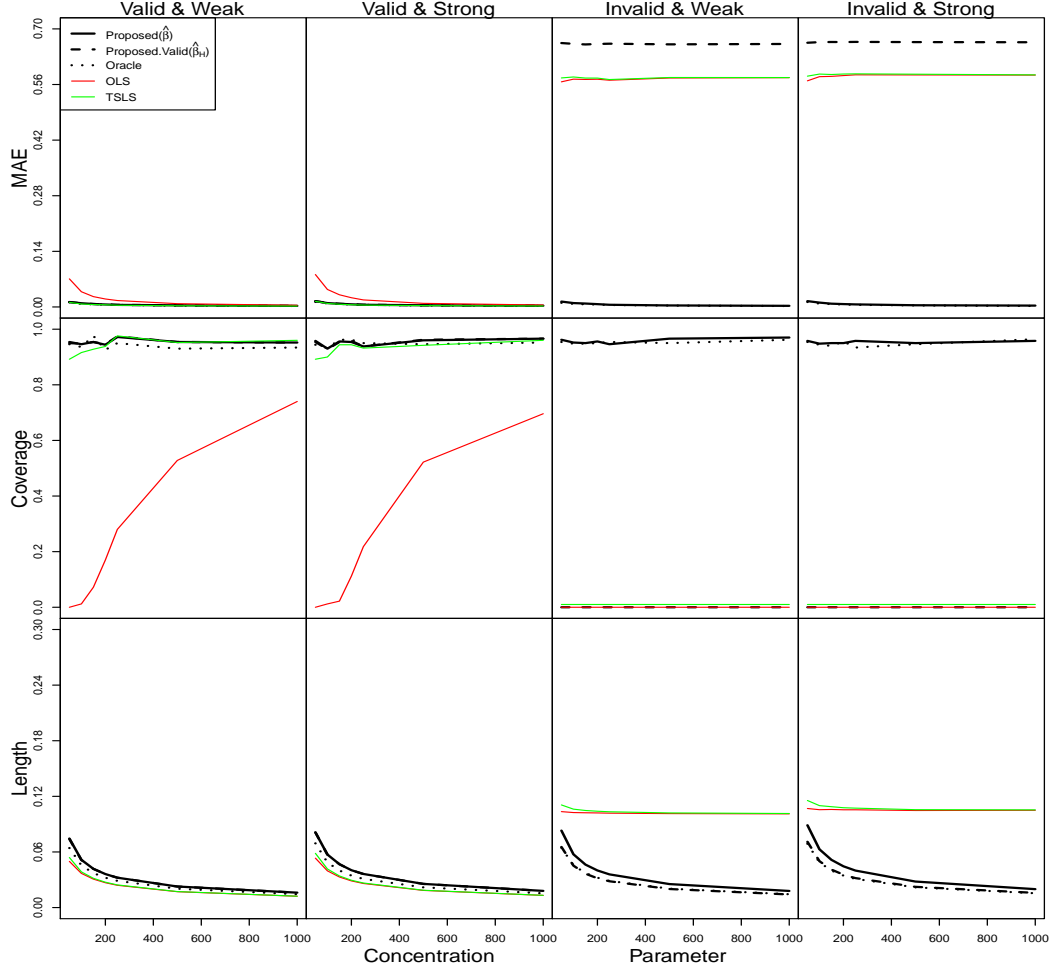


Fig. 2: Comparison of different methods when $p_z = 100$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Valid & Weak represents the case $\rho_1 = 0.2$ and $\rho_2 = 0$. The column labeled with Valid & Strong represents the case $\rho_1 = 0$ and $\rho_2 = 0$. The column labeled with Invalid & Weak represents the case $\rho_1 = 0.2$ and $\rho_2 = 2$. Finally, the column labeled with Invalid & Strong represents the case $\rho_1 = 0$ and $\rho_2 = 2$.

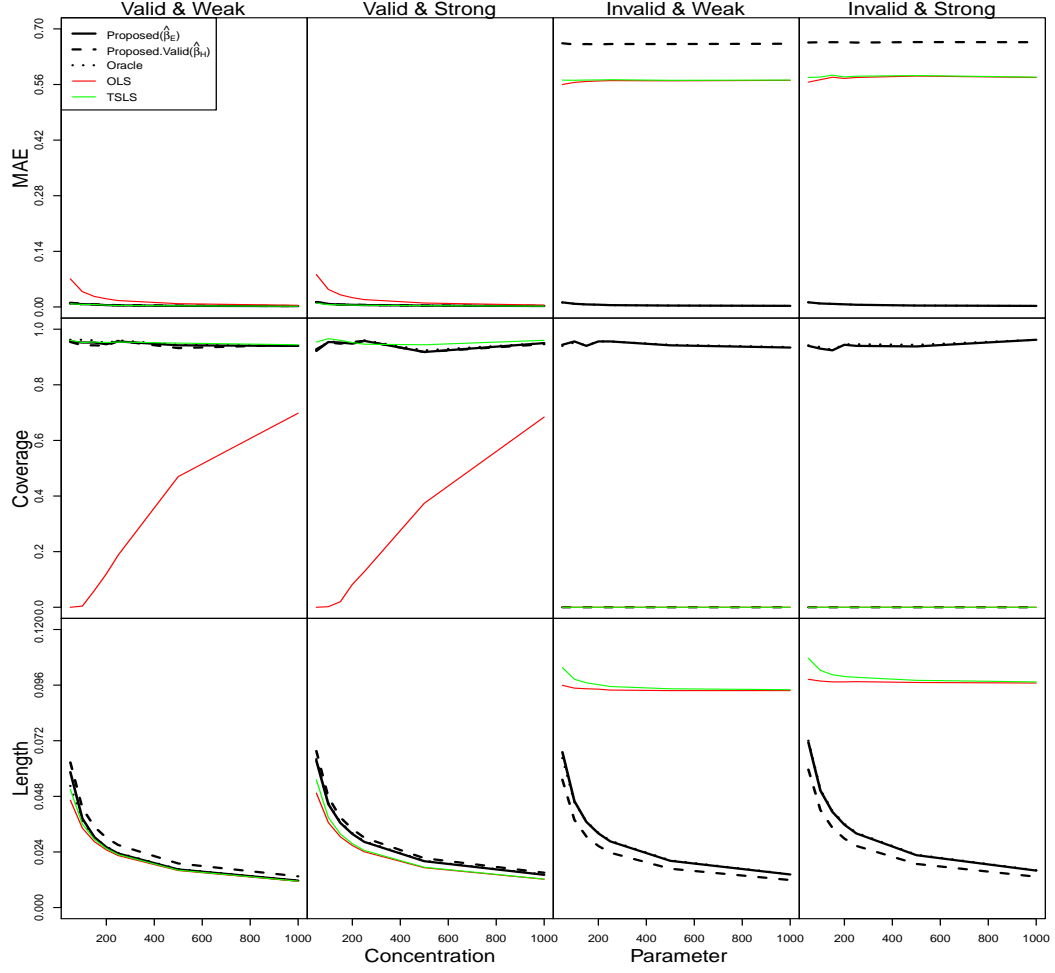


Fig. 3: Comparison of different methods when $p_z = 9$, $p_x = 10$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Valid & Weak represents the case $\rho_1 = 0.2$ and $\rho_2 = 0$. The column labeled with Valid & Strong represents the case $\rho_1 = 0$ and $\rho_2 = 0$. The column labeled with Invalid & Weak represents the case $\rho_1 = 0.2$ and $\rho_2 = 2$. Finally, the column labeled with Invalid & Strong represents the case $\rho_1 = 0$ and $\rho_2 = 2$.

5.3. Violation of (IN2) and (IN3)

This section summarizes the simulation study mostly in Section C in the supplementary materials that assess the violation of regularity assumptions (IN2) and (IN3), which were used to prove properties of our procedure in high dimensions. First, in Figures 1 and 2 of the main text, our proposed estimator $\hat{\beta}$ performs similarly for individually weak and strong IVs (i.e. for different ρ_1 s). Even with $\rho_1 = 0.1$ (shown in supplementary materials), our procedure still retains the desired level of coverage, reiterating our comment in Section 4.2 that the proposed estimator is not sensitive to the individual instrument strength assumption (IN2) required in the theoretical analysis. Also, when $\rho_2 = 2$ so that the invalid IVs violate the assumptions (A2) and (A3) in a strong sense and subsequently, (IN3) is satisfied, our method is similar to the oracle; this is in contrast to the estimator $\hat{\beta}_H$ and TSLS which are biased and lose coverage when any invalid instruments are present, even if the instruments strongly violate (A2) and (A3).

However, as discussed in Section 4.2, in high dimensions, if the instruments violate (A2) and (A3) weakly and consequently, violate (IN3), our procedure tends to suffer. In particular, in Section C.2 of the supplementary materials, when the invalid IVs weakly violated instrument with $\rho_2 = 1$, our method does not achieve 95 percent coverage even though the point estimate is still consistent, which was expected from Theorem 2. Fortunately, the under-coverage is at most 5 percent and is still much higher than TSLS or other methods assuming valid IVs after conditioning; see figures A2, A5, A8 and A11 of the supplementary materials for details. Also, in Section C.4 of the supplementary materials, when the number of candidate instruments are small and we have high dimensional covariates, which we think is the most typical setting in practice, the simulation results match closely with the simulation results presented in Figures 1 and 2. Specifically, in Figure A17 of the supplementary materials, the proposed confidence intervals have 95% coverage even though (IN3) is violated. Overall, the simulation study demonstrates that our procedure of selecting valid instruments can still improve the coverage property significantly compared to naive and popular methods, even if assumption (IN3) is violated.

6. Application: Causal Effect of Years of Education on Annual Earnings

To demonstrate our procedure 1 in real settings, we analyze the causal effect of years of education on yearly earnings, which has been studied extensively in economics using IV methods [Angrist and Krueger, 1991, Card, 1993, 1999]. The data comes from the Wisconsin Longitudinal Study (WLS), a longitudinal study that has kept track of American high school graduates from Wisconsin since 1957, and we examine the relationship between graduates' earnings and education from the 1974 survey Hauser [2005], roughly 20 years after they graduated from high school. Our analysis includes $N = 3772$ individuals, 1784 males and 1988 females. For our outcome, we use imputed log total yearly earnings prepared by WLS (see WLS documentation and Hauser [2005] for details) and for the treatment, we use the total years of education, all from the 1974 survey. The median total earnings is \$9,200 with a

25% quartile of \$1,000 and a 75% quartile of \$15,320 in 1974 dollars. The mean years of total education is 13.7 years with a standard deviation of 2.3 years.

We incorporate many covariates, including sex, graduate’s hometown population, educational attainment of graduates’ parents, graduates’ family income, relative income in graduates’ hometown, graduates’ high school denomination, high school class size, all measured in 1957 when the participants were high school seniors. We also include 81 genetic covariates, specifically single nucleotide polymorphisms (SNPs), that were part of WLS to further control for potential variations between graduates; see Section B in Supplementary Materials for details on the non-genetic and genetic covariates. In summary, our data analysis includes 7 non-genetic covariates and 81 genetic covariates.

We used five instruments in our analysis, all derived from past studies of education on earnings [Card, 1993, Blundell et al., 2005, Gary-Bobo et al., 2006]. They are (i) total number of sisters, (ii) total number of brothers, (iii) individual’s birth order in the family, all from Gary-Bobo et al. [2006], (iv) proximity to college from Card [1993], and (v) teacher’s interest in individual’s college education from Blundell et al. [2005], all measured in 1957. Although all these IVs have been suggested to be valid with varying explanations as to why they satisfy (A2) and (A3) after controlling for the aforementioned covariates, in practice, we are always uncertain due to the lack of complete socioeconomic knowledge about the effect of these IVs. Our method should provide some protection against this uncertainty compared to traditional methods where they simply assume that all five IVs are valid. Also, the first-stage F-test produces an F-statistic of 90.3 with a p-value less than 10^{-16} , which indicates very strong set of instruments. For more details on the instruments, see Section B of the Supplementary Materials.

Table 1 summarizes the results of our data analysis. OLS refers to running a regression of the treatment and the covariates on the outcome and looking at the slope coefficient of the treatment variable. TSLS refers to running two-stage least squares as described in Section 5 under the operating assumption that all the five instruments are valid; this is the usual and most popular analysis in the IV literature. Finally, we run the Procedure 1.

| Method | Point Estimate | 95% Confidence Interval |
|--------|----------------|-------------------------|
| OLS | 0.097 | (0.051, 0.143) |
| TSLS | 0.169 | (0.029, 0.301) |
| TSHT | 0.062 | (0.046, 0.077) |

Table 1: Estimates of the Effect of Years of Education on Log Earnings. OLS is ordinary least squares, TSLS is two-stage least squares, and TSHT is Procedure 1.

The OLS estimate suggests a positive association between education and earnings, with statistically significant result at $\alpha = 0.05$ level. This agrees with previous literature which suggests a statistically significant positive association between years of education and log earnings Card [1999]. However, OLS does not completely control for confounding even after controlling for covariates. TSLS provides an al-

ternative method of controlling for confounding by using instruments so long as all the instruments satisfy the three core assumptions and the inclusion of covariates helps make these assumptions more plausible. Unfortunately, we notice that the TSLS estimate in Table 1 is inconsistent with previous studies’ estimates among individuals from the U.S. between 1950s to 1970s, which range from 0.06 to 0.13 (see Table 4 in Card [1999]). Our method, which addresses the concern for invalid instruments with TSLS, provides an estimate of 0.062, which is more consistent with previous studies’ estimates of the effect of years of education on earnings.

The data analysis suggests that our method can be a useful tool in IV analysis when there is concern for invalid instruments, even after attempting to mitigate this problem via covariates. Our method provides much more accurate estimates of the returns on education than TSLS, which naively assumes all the instruments are valid.

7. Conclusion and Discussion

We present a method to estimate the effect of the treatment on the outcome using instrumental variables where we do not make the assumption that all the instruments are valid. Our approach is based on the novel TSHT procedure, which is shown to succeed in selecting valid IVs in the presence of possibly invalid IVs. Our approach provides robust confidence intervals in the presence of invalid IVs even after controlling for many covariates. In simulation and in real data settings, our approach provides a more robust analysis than the traditional IV approaches, most notably TSLS, by providing some protection against possibly invalid instruments.

As discussed in Section 4.2, our theoretical analysis for the case of invalid IVs even after controlling for high-dimensional covariates require Assumptions (IN2) and (IN3). While (IN2) is not crucial in practice as our simulation study demonstrates and is made for a cleaner technical exposition, we believe (IN3) is most likely necessary for invalid IV problems and this is echoed in the model selection literature by Leeb and Pötscher [2005] who pointed out that “in general no model selector can be uniformly consistent for the most parsimonious true model” and hence the post-model-selection inference is generally non-uniform. Consequently, the set of competing models has to be “well separated” such that we can consistently select a correct model and Assumption (IN3) serves as this “well separated” condition in our invalid IV problem. While some recent work in high dimensional inference [Zhang and Zhang, 2014, Javanmard and Montanari, 2014, van de Geer et al., 2014, Chernozhukov et al., 2015, Cai and Guo, 2016a] do not make this “well separated” assumption, as we stressed before, our invalid IV problem is of different nature than the prior work because a single invalid IV declared as valid can ruin inference while said prior works assume covariates are exogenous and moments are known perfectly.

Finally, in practice, we believe that violation of (IN3) in high dimensions will not drastically harm inference and our CI will still have coverage around $1 - \alpha$, which is much better than TSLS and prior work assuming valid IVs after conditioning on many covariates, which have no coverage. In particular, our empirical investigations generally show that the under-coverage is no more than 5% and we think this is

partly due to the fact that (i) our procedure will still pick up the strongly invalid IVs and (ii) if the instruments are weakly invalid, the bias from them via π^* will be relatively small. It is certainly possible that advanced methods can weaken (IN3) and we leave this as a direction for further research.

Acknowledgments

The research of Hyunseung Kang was supported in part by NSF Grant DMS-1502437. The research of T. Tony Cai was supported in part by NSF Grants DMS-1208982 and DMS-1403708, and NIH Grant R01 CA127334. The research of Dylan S. Small was supported in part by NSF Grant SES-1260782.

References

- Donald W. K. Andrews, Marcelo J. Moreira, and James H. Stock. Performance of conditional wald tests in $\{IV\}$ regression with weak instruments. *Journal of Econometrics*, 139(1):116–132, 2007.
- Joshua D Angrist and Alan B Krueger. Does compulsory school attendance affect schooling and earnings? *Quarterly Journal of Economics*, 106(4):979–1014, 1991.
- Joshua D. Angrist, Guido W. Imbens, and Donald B. Rubin. Identification of causal effects using instrumental variables. *Journal of the American Statistical Association*, 91(434):444–455, 1996.
- Michael Baiocchi, Jing Cheng, and Dylan S. Small. Instrumental variable methods for causal inference. *Statistics in Medicine*, 33(13):2297–2340, 2014.
- A. Belloni, D. Chen, V. Chernozhukov, and C. Hansen. Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica*, 80(6):2369–2429, 2012.
- Alexandre Belloni, Victor Chernozhukov, and Christian Hansen. Inference on treatment effects after selection among high-dimensional controls. *The Review of Economic Studies*, 81(2):608–650, 2014.
- Peter J Bickel, Ya’acov Ritov, and Alexandre B Tsybakov. Simultaneous analysis of lasso and dantzig selector. *The Annals of Statistics*, 37(4):1705–1732, 2009.
- Richard Blundell, Lorraine Dearden, and Barbara Sianesi. Evaluating the effect of education on earnings: models, methods and results from the national child development survey. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 168(3):473–512, 2005.
- J. Bowden, G. Davey Smith, and S. Burgess. Mendelian randomization with invalid instruments: effect estimation and bias detection through egger regression. *International Journal of Epidemiology*, 44(2):512–525, 2015.

- Jack Bowden, George Davey Smith, Philip C. Haycock, and Stephen Burgess. Consistent estimation in mendelian randomization with some invalid instruments using a weighted median estimator. *Genetic Epidemiology*, 40(4):304–314, 2016.
- Peter Bühlmann and Sara van de Geer. *Statistics for high-dimensional data: methods, theory and applications*. Springer, 2011.
- Stephen Burgess, Nicholas J. Timpson, Shah Ebrahim, and George Davey Smith. Mendelian randomization: where are we now and where are we going? *International Journal of Epidemiology*, 44(2):379–388, 2015.
- Stephen Burgess, Jack Bowden, Frank Dudbridge, and Simon G. Thompson. Robust instrumental variable methods using multiple candidate instruments with application to mendelian randomization. *arXiv*, 2016.
- T Tony Cai and Zijian Guo. Confidence intervals for high-dimensional linear regression: Minimax rates and adaptivity. *The Annals of statistics*, To appear, 2016a.
- T Tony Cai and Zijian Guo. Supplement to “confidence intervals for high-dimensional linear regression: Minimax rates and adaptivity”. *The Annals of statistics*, To appear, 2016b.
- David Card. Using geographic variation in college proximity to estimate the return to schooling. Working Paper 4483, National Bureau of Economic Research, October 1993.
- David Card. Chapter 30 - the causal effect of education on earnings. In Orley C. Ashenfelter and David Card, editors, *Handbook of Labor Economics*, volume 3, Part A, pages 1801 – 1863. Elsevier, 1999.
- Xu Cheng and Zhipeng Liao. Select the valid and relevant moments: An information-based lasso for gmm with many moments. *Journal of Econometrics*, 186(2):443–464, 2015.
- Victor Chernozhukov, Christian Hansen, and Martin Spindler. Post-selection and post-regularization inference in linear models with many controls and instruments. 2015.
- George Davey Smith and Shah Ebrahim. Mendelian randomization: can genetic epidemiology contribute to understanding environmental determinants of disease? *International Journal of Epidemiology*, 32(1):1–22, 2003.
- George Davey Smith and Shah Ebrahim. Mendelian randomization: prospects, potentials, and limitations. *International Journal of Epidemiology*, 33(1):30–42, 2004.
- Vanessa Didelez and Nuala Sheehan. Mendelian randomization as an instrumental variable approach to causal inference. *Statistical methods in medical research*, 16(4):309–330, 2007.

- David L Donoho. De-noising by soft-thresholding. *IEEE transactions on information theory*, 41(3):613–627, 1995.
- David L Donoho and Jain M Johnstone. Ideal spatial adaptation by wavelet shrinkage. *Biometrika*, 81(3):425–455, 1994.
- Jianqing Fan and Yuan Liao. Endogeneity in high dimensions. *Annals of statistics*, 42(3):872, 2014.
- Robert Gary-Bobo, Nathalie Picard, and Ana Prieto. Birth order and sibship sex composition as instruments in the study of education and earnings. CEPR Discussion Papers 5514, C.E.P.R. Discussion Papers, 2006.
- Eric Gautier and Alexandre B. Tsybakov. High-dimensional instrumental variables regression and confidence sets. *arXiv preprint arXiv:1105.2454*, 2011.
- Chirok Han. Detecting invalid instruments using l1-gmm. *Economics Letters*, 101(3):285–287, 2008.
- Robert M. Hauser. Survey response in the long run: The wisconsin longitudinal study. *Field Methods*, 17(1):3–29, 2005.
- Miguel A. Hernán and James M. Robins. Instruments for causal inference: An epidemiologist’s dream? *Epidemiology*, 17(4):360–372, 2006.
- Paul W. Holland. Causal inference, path analysis, and recursive structural equations models. *Sociological Methodology*, 18(1):449–484, 1988.
- Guido W. Imbens. Instrumental variables: An econometrician’s perspective. *Statistical Science*, 29(3):323–358, 2014.
- Adel Javanmard and Andrea Montanari. Confidence intervals and hypothesis testing for high-dimensional regression. *The Journal of Machine Learning Research*, 15(1):2869–2909, 2014.
- Hyunseung Kang, T. Tony Cai, and D. S. Small. A simple and robust confidence interval for causal effects with possibly invalid instruments. *arXiv preprint arXiv:1504.03718*, 2016a.
- Hyunseung Kang, Anru Zhang, T Tony Cai, and Dylan S Small. Instrumental variables estimation with some invalid instruments and its application to mendelian randomization. *Journal of the American Statistical Association*, 111:132–144, 2016b.
- Michal Kolesár, Raj Chetty, John N. Friedman, Edward L. Glaeser, and Guido W. Imbens. Identification and inference with many invalid instruments. *Journal of Business & Economic Statistics*, 33(4):474–484, 2015.
- Debbie A. Lawlor, Roger M. Harbord, Jonathan A. C. Sterne, Nic Timpson, and George Davey Smith. Mendelian randomization: Using genes as instruments for making causal inferences in epidemiology. *Statistics in Medicine*, 27(8):1133–1163, 2008.

- Hannes Leeb and Benedikt M Pötscher. Model selection and inference: Facts and fiction. *Econometric Theory*, 21(01):21–59, 2005.
- Marcelo J. Moreira. A conditional likelihood ratio test for structural models. *Econometrica*, 71(4):1027–1048, 2003.
- Michael P. Murray. Avoiding invalid instruments and coping with weak instruments. *The Journal of Economic Perspectives*, 20(4):111–132, 2006.
- Jerzy Neyman. On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, 5(4):465–472, 1923.
- Zhao Ren, Tingni Sun, Cun-Hui Zhang, and Harrison H Zhou. Asymptotic normality and optimalities in estimation of large gaussian graphical model. *arXiv preprint arXiv:1309.6024*, 2013.
- Nicholas S. Roetker, James A. Yonker, Chee Lee, Vicky Chang, Jacob J. Basson, Carol L. Roan, Taissa S. Hauser, Robert M. Hauser, and Craig S. Atwood. Multi-gene interactions and the prediction of depression in the wisconsin longitudinal study. *British Medical Journal Open*, 2(4):e000944, 2012.
- Donald B. Rubin. Estimating causal effects of treatments in randomized and non-randomized studies. *Journal of Educational Psychology*, 66(5):688, 1974.
- Dylan S. Small. Sensitivity analysis for instrumental variables regression with over-identifying restrictions. *Journal of the American Statistical Association*, 102(479):1049–1058, 2007.
- D Staiger and JH Stock. Instrumental variables regression with weak instruments. *Econometrica*, 65(3):557–586, 1997.
- James H Stock and Jonathan H Wright. Gmm with weak identification. *Econometrica*, pages 1055–1096, 2000.
- James H Stock and Motohiro Yogo. Asymptotic distributions of instrumental variables statistics with many instruments. *Identification and inference for econometric models: essays in honor of Thomas Rothenberg*, 2005.
- James H. Stock, Jonathan H. Wright, and Motohiro Yogo. A survey of weak instruments and weak identification in generalized method of moments. *Journal of Business & Economic Statistics*, 20(4), 2002.
- Tingni Sun and Cun-Hui Zhang. Scaled sparse linear regression. *Biometrika*, 101(2):269–284, 2012.
- Sonja A. Swanson and Miguel A. Hernán. Commentary: How to report instrumental variables analyses (suggestions welcome). *Epidemiology*, 24:370–374, 2013.
- Nicholas J. Timpson, Debbie A. Lawlor, Roger M. Harbord, Tom R Gaunt, Ian NM Day, Lyle J. Palmer, Andrew T. Hattersley, Shah Ebrahim, Gordon Lowe, Ann Rumley, and George Davey Smith. C-reactive protein and its role in metabolic

- syndrome: Mendelian randomisation study. *The Lancet*, 366(9501):1954–1959, 2005.
- Sara van de Geer, Peter Bühlmann, Yaacov Ritov, and Ruben Dezeure. On asymptotically optimal confidence regions and tests for high-dimensional models. *The Annals of Statistics*, 42(3):1166–1202, 2014.
- Hal R Varian. Big data: New tricks for econometrics. *The Journal of Economic Perspectives*, 28(2):3–27, 2014.
- Roman Vershynin. Introduction to the non-asymptotic analysis of random matrices. In Yonina Eldar and Gitta Kutyniok, editors, *Compressed Sensing: Theory and Applications*, pages 210–268. Cambridge University Press, 2012.
- Frank Windmeijer, Helmut Farbmacher, Neil Davies, and George Davey Smith. On the use of the lasso for instrumental variables estimation with some invalid instruments. Technical report, Department of Economics, University of Bristol, UK, 2016.
- Jeffrey M. Wooldridge. *Econometric Analysis of Cross Section and Panel Data*. MIT press, 2nd ed. edition, 2010.
- Cun-Hui Zhang and Stephanie S Zhang. Confidence intervals for low dimensional parameters in high dimensional linear models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):217–242, 2014.
- Shuheng Zhou. Restricted eigenvalue conditions on subgaussian random matrices. *arXiv preprint arXiv:0912.4045*, 2009.

A. Supplementary Materials: Theory for Valid IVs After Controlling for High Dimensional Covariates

In this section, we state the theoretical results for valid IVs after controlling for High Dimensional Covariates. Under the assumptions (R1)-(R3), Theorem 3 shows if the instruments are valid after conditioning on many covariates, then the estimator $\hat{\beta}_H$ in our procedure is consistent and asymptotically normal.

THEOREM 3. *Suppose we have valid IVs, that is $\pi^* = 0$ in (2), and the assumptions (R1) – (R3) hold. The following property holds for the estimator $\hat{\beta}_H$,*

$$\sqrt{n} \left(\hat{\beta}_H - \beta^* \right) = T^{\beta^*} + \Delta^{\beta^*}, \quad (27)$$

where $T^{\beta^*} \mid \mathbf{W} \sim N(0, V_H)$, $V_H = \sigma^2 / \|\gamma^*\|_2^4 \left\| \sum_{j \in S^*} \gamma_j^* \mathbf{W} \hat{\mathbf{u}}^{[j]} / \sqrt{n} \right\|_2^2$ and $\Delta^{\beta^*} / \sqrt{V_H} \xrightarrow{p} 0$ as $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$.

Theorem 3 states that if the IVs satisfy the exclusion restriction and no unmeasured confounding after conditioning on many covariates, $\hat{\beta}_H$ defined in (16) is a consistent and the dominating part of the scaled difference $\sqrt{n}(\hat{\beta}_H - \beta)$ is normal. Based on the asymptotic normality established in (27), the following theorem justifies the coverage property of the confidence interval proposed in (17) under the assumption that all instruments have no direct effect and are unconfounded after conditioning on many covariates.

THEOREM 4. *Suppose we have valid IVs, that is $\pi^* = 0$ in (2) and the assumptions (R1) – (R3) hold. Assuming $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$, the confidence interval given in (17) has asymptotically coverage probability $1 - \alpha$, i.e.,*

$$\mathbf{P} \left\{ \beta^* \in \left(\hat{\beta}_H - z_{1-\alpha/2} \sqrt{\hat{V}_H/n}, \hat{\beta}_H + z_{1-\alpha/2} \sqrt{\hat{V}_H/n} \right) \right\} \rightarrow 1 - \alpha, \quad (28)$$

Theorem 4 is similar to a result given in Chernozhukov et al. [2015], who studied IV estimators in high dimensional regime where all the instruments are valid after conditioning. However, there are some notable differences between our results and those in Chernozhukov et al. [2015] in terms of sparsity and instrument-covariate modeling assumptions that are required to achieve $1 - \alpha$ coverage. A simulation study is carried out in Section 5 to compare our procedure to that of the oracle.

B. Supplementary Materials: Details on Data Analysis

B.1. Genetic Data

This is the list of genetic covariates, i.e. 81 SNPs, used in the data analysis: rs1018381, rs1042838, rs11178997, rs1137100, rs11564774, rs11902591, rs12152850, rs12313279, rs12602084, rs12664989, rs12913832, rs1421085, rs1424954, rs1435252, rs144848, rs1501299, rs1535_chr11, rs17070145, rs174575, rs17529477, rs17561, rs17571, rs1799913, rs1799945, rs1799966, rs1799998, rs1800046, rs1800497, rs1800795, rs1800955,

| Covariate | Statistics |
|-------------------------------------|--|
| Sex | 1784 males (47.3%), 1988 females (52.7%) |
| Hometown Population | |
| Under 1,000 | 585 (15.5%) |
| 1,000 - 2,499 | 584 (15.5%) |
| 2,500 - 9,999 | 771 (20.4 %) |
| 10,000 - 24,999 | 375 (9.9%) |
| 25,000 - 49,999 | 602 (16.0 %) |
| 50,000 - 99,999 | 226 (7.1%) |
| 100,000 - 150,000 | 112 (3.0%) |
| Over 150,000 | 477 (12.6%) |
| Father's Total Years of Education | 10.4 years (SD: 3.2) |
| Mother's Total Years of Education | 10.8 years (SD: 3.0) |
| Family Income Relative to Community | |
| Considerably below average | 15 (0.4%) |
| Somewhat below average | 246 (6.5 %) |
| Average | 2670 (70.8%) |
| Somewhat above average | 783 (20.8%) |
| Considerably above average | 58 (1.5%) |
| School Type | 479 catholic/private (12.7%) , 3293 public (87.3%) |
| Graduating Class Size | 166.8 students (SD: 131.5) |

Table 2: Covariates used in the data analysis. All the covariates were measured in 1957. SD stands for standard deviation.

rs1805420, rs1937_chr10, rs2059693, rs2061174, rs2071219, rs2237436, rs2241766, rs2242592, rs2254298, rs2306604, rs25533, rs2760118, rs2779562, rs2963238, rs35118453, rs363050, rs3749386, rs3788862, rs3797297, rs3802657, rs3853248, rs3990403, rs4073366, rs4245147, rs429358, rs4502885, rs45537037, rs4680, rs4986852, rs5031016, rs592389, rs6152_chrX, rs6166, rs6169, rs6265, rs6277, rs6312, rs6314 rs6318 rs669, rs707555, rs7412, rs760761, rs7761133, rs7997012, rs8039957, rs8076005, rs8191992, rs821616, rs878567, rs908867.

Similar to another analysis with the WLS genetic data, we remove individuals with more than 10% missing genotype data Roetker et al. [2012]. We also remove SNPs with no variation in our data set or SNPs with more than 20% missing. The other missing values were imputed using the most frequent genotype for that SNP. Following a Mendelian randomization study by Timpson et al. [2005], we code the SNPs using an additive model.

B.2. Covariates in the Data Analysis

We include the following nine covariates for our data analysis: sex, graduate's hometown population, educational attainment of graduates' parents, graduates' family income, relative income in graduates' hometown, graduates' high school denomination, high school class size, all measured in 1957 when the participants were high school seniors. The details of these covariates are in Table 2

| Instrument | Statistics |
|---|----------------------------|
| Total number of sisters | Median: 1 (25Q: 1, 75Q: 2) |
| Total number of brothers | Median: 1 (25Q: 1, 75Q: 2) |
| Birth order | Median: 1 (25Q: 1, 75Q: 3) |
| Distance from High School (HS) to College | |
| HS more than 15 miles from any college | 1,487 (39.4%) |
| HS 15 miles or less from extension center | 168 (4.4 %) |
| HS 15 miles or less from state college | 196 (5.2%) |
| HS less than 15 miles from private college | 159 (4.2 %) |
| HS in city with extension center | 446 (1.2 %) |
| HS in city with stage college | 283 (7.5 %) |
| HS in city with private college | 205 (5.4 %) |
| HS 15 miles or less from state university | 239 (6.3 %) |
| HS in city with state university (Milwaukee or Madison) | 589 (15.6%) |
| Teacher's Interest in Individual's College Education | |
| Discouraged to attend college | 47 (1.2%) |
| Had no effect or no response | 1825 (48.4%) |
| Encouraged to attend college | 1900 (50.4%) |

Table 3: Instruments used in the data analysis. All the instruments were measured in 1957. 25Q and 75Q stand for 25% and 75% quartile, respectively. There were 119 individuals who had no response to teacher's interest in individual's college education.

B.3. Instruments in the Data Analysis

They are five instruments in the analysis: (i) total number of sisters, (ii) total number of brothers, (iii) individual's birth order in the family, (iv) proximity to college, and (v) teacher's interest in individuals' college education, all measured in 1957.

A small proportion of individuals had missing teacher's interest in individual's college education. To accommodate this, we created a binary covariate to indicate missingness and imputed missing observations as the average value for that covariate. In total, there were 119 individuals with missing values for teacher's interest in education out of 3772 individuals.

C. Supplementary Materials: Details of the Simulation Study

C.1. Oracle Concentration Parameter

In the IV literature, the concentration parameter is a measure of instruments' global strength [Stock and Wright, 2000] where a high value of the concentration parameter denotes strong instruments; the concentration parameter is also the expected value of the partial F-test where we regress \mathbf{D} on the instruments \mathbf{Z} and covariates \mathbf{X} and conduct the usual F-test on the coefficients of \mathbf{Z} . We can also define the population version of the concentration parameter by taking the expected value of the concentration parameter over many samples [Stock and Yogo, 2005]. Unfortunately, the traditional notions of concentration parameter cannot deal with invalid

instruments because the traditional concentration parameter assumes, a priori, that we know exactly which instruments are valid and invalid. The oracle concentration parameter resolves this problem by making a simple modification to the traditional concentration parameter where an oracle provides us with the true set of valid instruments \mathcal{V}^* and we examine the strength, specifically the elements of γ^* in \mathcal{V}^* only among the valid instruments. In fact, much like the usual concentration parameter, the oracle concentration parameter grows with sample size n . For example, when sample size increases from 100 to 200, the concentration parameter will double, that is, $C(\gamma^*, \mathcal{V}^*, 200) = 2C(\gamma^*, \mathcal{V}^*, 100)$.

We call it the oracle expected concentration parameter because $C(\gamma^*, \mathcal{V}^*, n)$ is usually not computable in the real application when there is no prior information about which instruments are valid or invalid, unless, of course, an oracle provides such information to the data analyst.

C.2. Extended simulations for high-dimensional instruments and covariates

In this section, we present the extended simulation analysis by varying more parameters from the main text. In particular, we keep the data generating model and the parameter settings identical, except we vary ρ_1 , ρ_2 , and n more broadly than what is presented in the main text. These results are in Figures A1 to A15.

Specifically, across the figures, the column titled “WeakIV1” represents the case $\rho_1 = 0.1$, the column titled “WeakIV2” represents the case $\rho_1 = 0.2$ and the column titled with “Strong” represents the case $\rho_1 = 0$. The row indexed with “MAE” represents the Median Absolute Error of the proposed estimators, which measures the performance of the point estimator. The row indexed with “Coverage” represents the coverage performance of the confidence intervals proposed. The row indexed with “Length” represents the length performance of corresponding confidence intervals. Finally, different figures vary with sample size n with $n = (100, 200, 300, 1000, 3000)$ and violation of instrumental variable, where $\rho_2 = 0$ represents the case of all valid instruments, $\rho_2 = 1$ represents the case of weakly violated instruments and $\rho_2 = 2$ represents the case of strongly violated instruments.

The extended simulation results presented here match closely with the main simulation results presented in the main paper. Our proposed procedure performs as well as the oracle in terms of accuracy and coverage property. In addition, when there are invalid instruments, the estimator $\hat{\beta}_H$ and TSLS are biased and lose the coverage while our procedure incorporating invalid instruments performs as well as the oracle.

The case that our method fails to do as well as advertised is in the coverage of the confidence interval when $\rho_2 = 1$ where the violation of instrumental variable is not very strong. Here, our method does not achieve 95 percent coverage even though it is still consistent. The under-coverage is up to 5 percent. We believe this is due to the failure of selecting valid instruments. Specifically, based on the theoretical results in Theorem 2 of the paper, the constructed confidence intervals might not achieve the guaranteed coverage if Assumption (IN3) in the main paper is not satisfied, that is, the violation of instruments is weak. We leave it as a direction

of future research to construct honest confidence intervals in the case of weakly violated instruments.

C.3. *Low-dimensional instruments and covariates*

The set up is similar to Section 5.1 except for the number of instruments and covariates. The data generating process for the simulation follows the models (2) and (3) in Section 2.2 with $p_z = 9$ instruments and $p_x = 10$ covariates where \mathbf{W}_i is a multivariate normal with mean zero and covariance $\Sigma_{ij}^* = 0.5^{|i-j|}$ for $1 \leq i, j \leq 19$. The parameters for the models are: $\beta^* = 1$, $\phi^* = (0.6, 0.7, 0.8, \dots, 1.5) \in \mathbb{R}^{10}$ so that $s_{x1} = 10$, $\psi^* = (1.1, 1.2, 1.3, \dots, 2.0) \in \mathbb{R}^{10}$ so that $s_{x2} = 10$. Instruments that satisfy Assumption (A1) are $\mathcal{S}^* = \{1, \dots, 7\}$ and instruments that satisfy all three IV assumptions (A1)-(A3) are $\mathcal{V}^* = \{1, 2, 3, 4, 5\}$; thus instruments 6 and 7 only satisfy (A1), but do not satisfy (A2) and (A3). The sample size is 1000.

The extended simulation results presented here match closely with the main simulation results presented in the main paper and the simulation results presented in Section C.2. Our proposed procedure performs as well as the oracle in terms of accuracy and coverage property. In addition, when there are invalid instruments, the estimator $\hat{\beta}_H$ and TSLS are biased and lose the coverage while our procedure incorporating invalid instruments performs as well as the oracle.

If we use the estimator $\hat{\beta}_E$ defined in (18) and the confidence interval (19), the proposed procedure performs almost the same as the oracle in terms of accuracy, coverage property and length, which supports the theory established in Theorem 1. Specifically, in Figure A14, the proposed confidence intervals have 95% coverage even though (IN3) is violated.

C.4. *Low-dimensional instruments and high-dimensional covariates*

The set up is similar to Section 5.1 except for the number of instruments. The data generating process for the simulation follows the models (2) and (3) in Section 2.2 with $p_z = 9$ instruments and $p_x = 150$ covariates where \mathbf{W}_i is a multivariate normal with mean zero and covariance $\Sigma_{ij}^* = 0.5^{|i-j|}$ for $1 \leq i, j \leq 159$. The parameters for the models are: $\beta^* = 1$, $\phi^* = (0.6, 0.7, 0.8, \dots, 1.5, 0, 0, \dots, 0) \in \mathbb{R}^{150}$ so that $s_{x1} = 10$, $\psi^* = (1.1, 1.2, 1.3, \dots, 2.0, 0, 0, \dots, 0) \in \mathbb{R}^{150}$ so that $s_{x2} = 10$. Instruments that satisfy Assumption (A1) are $\mathcal{S}^* = \{1, \dots, 7\}$ and instruments that satisfy all three IV assumptions (A1)-(A3) are $\mathcal{V}^* = \{1, 2, 3, 4, 5\}$; thus instruments 6 and 7 only satisfy (A1), but do not satisfy (A2) and (A3). The sample size is 1000.

The extended simulation results presented here match closely with the main simulation results presented in the main paper and the simulation results presented in Section C.2. Our proposed procedure performs as well as the oracle in terms of accuracy and coverage property. In addition, when there are invalid instruments, the estimator $\hat{\beta}_H$ and TSLS are biased and lose the coverage while our procedure incorporating invalid instruments performs as well as the oracle. Specifically, in Figure A17, the proposed confidence intervals have 95% coverage even though (IN3) is violated.

D. Supplementary Materials: Proofs of Theorems

In this section, we provide detailed proofs of Theorem 1 3, 4 and 2. Proof of extra lemmas are presented in next section. Before presenting the proof, we will introduce the notations used throughout the proof.

D.1. Notations

For any vector $\mathbf{v} \in \mathbb{R}^p$, let \mathbf{v}_j denote the j th element of \mathbf{v} . Let $\|\mathbf{v}\|_1$, $\|\mathbf{v}\|_2$, and $\|\mathbf{v}\|_\infty$ be the usual 1, 2 and ∞ -norms, respectively. Let $\|\mathbf{v}\|_0$ denote the 0-norm, i.e. the number of non-zero elements in \mathbf{v} . The support of \mathbf{v} , denoted as $\text{supp}(\mathbf{v}) \subseteq \{1, \dots, p\}$, is defined as the set containing the non-zero elements of the vector \mathbf{v} , i.e. $j \in \text{supp}(\mathbf{v})$ if and only if $\mathbf{v}_j \neq 0$. Also, for a vector $\mathbf{v} \in \mathbb{R}^p$ and set $J \subseteq \{1, \dots, p\}$, we denote $\mathbf{v}_J \in \mathbb{R}^p$ to be the vector where all the elements except whose indices are in J are zero. For a set J , $|J|$ denotes its cardinality.

For any n by p matrix $\mathbf{M} \in \mathbb{R}^{n \times p}$, we denote the (i, j) element of matrix \mathbf{M} as \mathbf{M}_{ij} , the i th row as $\mathbf{M}_{i\cdot}$, and the j th column as $\mathbf{M}_{\cdot j}$. Let \mathbf{M}^\top be the transpose of \mathbf{M} . Finally, $\|\mathbf{M}\|_\infty$ represents the element-wise matrix sup norm of matrix \mathbf{M} .

For a sequence of random variables X_n , we use $X_n \xrightarrow{p} X$ and $X_n \xrightarrow{d} X$ to represent that X_n converges to X in probability and in distribution, respectively. For any two sequences a_n and b_n , we will write $a_n \gg b_n$ if $\limsup \frac{b_n}{a_n} = 0$ and write $a_n \ll b_n$ if $b_n \gg a_n$. We use c and C to denote generic positive constants that may vary from place to place.

Throughout the whole proof section, we will use $\beta, \gamma, \mathbf{\Gamma}, \psi, \Psi, \pi, \mathbf{\Theta}_{11}, \mathbf{\Theta}_{22}, \mathbf{\Theta}_{12}, \Sigma, T^\beta, \Delta^\beta$ to stand for $\beta^*, \gamma^*, \mathbf{\Gamma}^*, \psi^*, \Psi^*, \pi^*, \mathbf{\Theta}_{11}^*, \mathbf{\Theta}_{22}^*, \mathbf{\Theta}_{12}^*, \Sigma^*, T^{\beta^*}, \Delta^{\beta^*}$ respectively and define

$$\hat{\mathbf{v}}^{[j]} = \mathbf{W}^\top \hat{\mathbf{u}}^{[j]} \quad \text{for } 1 \leq j \leq p_z.$$

We also introduce the notation $\mathbf{\Omega} = \Sigma^{-1}$, $\sigma_1 = \sqrt{\mathbf{\Theta}_{11}}$, $\sigma_2 = \sqrt{\mathbf{\Theta}_{22}}$ and $\mathbf{\Pi}_i = (e_{i1}, e_{i2})$. Let $M_2 = \max\{1/\lambda_{\min}(\mathbf{\Theta}), \lambda_{\max}(\mathbf{\Theta})\}$ and hence $1/M_2 \leq \lambda_{\min}(\mathbf{\Theta}) \leq \lambda_{\max}(\mathbf{\Theta}) \leq M_2$. We normalize the columns of \mathbf{W} as $\mathbf{H}_{\cdot j} = \sqrt{n} \mathbf{W}_{\cdot j} / \|\mathbf{W}_{\cdot j}\|_2$ for $j \in [p]$. Let $\text{Diag} = \text{diag}(\|\mathbf{W}_{\cdot j}\|_2 / \sqrt{n})_{1 \leq j \leq p}$ denote the $p \times p$ diagonal matrix with (j, j) entry to be $\|\mathbf{W}_{\cdot j}\|_2 / \sqrt{n}$. We set $\lambda_0 = \sqrt{2.05 \log p / n} = (1 + \gamma_0) \sqrt{2\delta_0 \log p / n}$, where $\delta_0 = \sqrt{1.025} > 1$ and $\gamma_0 = (1.025)^{\frac{1}{4}} - 1 > 0$. Take $\epsilon_0 = 2.01/\gamma_0 + 1$, $\nu_0 = 0.01$, $\tau_0 = 0.01$, $C_1 = 2.25$, $c_0 = 1/6$ and $C_0 = 3$. We also assume that $\log p / n \rightarrow 0$ and $\delta_0 \log p > 2$. Rather than use the constants directly in the following discussion, we use $\delta_0, \pi_0, \epsilon_0, \nu_0, C_1, C_0$ and c_0 to represent the above fixed constants in the following discussion. We review the following definition of restricted eigenvalue introduced in Bickel et al. [2009],

$$\kappa(X, k, \alpha_0) = \min_{\substack{J_0 \subseteq \{1, \dots, p\}, \\ |J_0| \leq k}} \min_{\substack{\delta \neq 0, \\ \|\delta_{J_0^c}\|_1 \leq \alpha_0 \|\delta_{J_0}\|_1}} \frac{\|X\delta\|_2}{\sqrt{n} \|\delta_{J_0}\|_2}. \quad (29)$$

Define the oracle estimator of σ_1 and σ_2 as

$$\sigma_1^{ora} = \frac{1}{\sqrt{n}} \|Y - \mathbf{Z}\mathbf{\Gamma} - \mathbf{X}\Psi\|_2 \text{ and } \sigma_2^{ora} = \frac{1}{\sqrt{n}} \|D - \mathbf{Z}\gamma - \mathbf{X}\psi\|_2,$$

and

$$\tau = \sqrt{1 + \epsilon_0} \frac{2\sqrt{s}\lambda_0}{\kappa(\mathbf{H}, 4s, 1 + 2\epsilon_0)}. \quad (30)$$

D.2. Proof of Theorem 1

Define

$$A(\mathcal{V}) = \widehat{\Sigma}_{\mathcal{V}, \mathcal{V}} - \widehat{\Sigma}_{\mathcal{V}, \mathcal{V}^c} \widehat{\Sigma}_{\mathcal{V}^c, \mathcal{V}^c}^{-1} \widehat{\Sigma}_{\mathcal{V}^c, \mathcal{V}} \quad \text{and} \quad A^*(\mathcal{V}) = \Sigma_{\mathcal{V}, \mathcal{V}} - \Sigma_{\mathcal{V}, \mathcal{V}^c} \Sigma_{\mathcal{V}^c, \mathcal{V}^c}^{-1} \Sigma_{\mathcal{V}^c, \mathcal{V}}.$$

We introduce the following lemmas to facilitate the proof.

LEMMA 1. *Under the assumptions of Theorem 1, we have*

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\tilde{\mathcal{V}} = \mathcal{V}^* \right) = 1 \quad (31)$$

LEMMA 2. *Under the assumptions of Theorem 1, we have*

$$\sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\Gamma}_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} - \beta \right) \xrightarrow{d} N \left(0, \frac{\Theta_{11} + \beta^2 \Theta_{22} - 2\beta \Theta_{12}}{\gamma_{\mathcal{V}^*}^\top A^*(\mathcal{V}^*) \gamma_{\mathcal{V}^*}} \right). \quad (32)$$

The estimator defined in (18) can be expressed as $\widehat{\beta}_E = \frac{\tilde{\gamma}_{\tilde{\mathcal{V}}}^\top A(\tilde{\mathcal{V}}) \tilde{\Gamma}_{\tilde{\mathcal{V}}}}{\tilde{\gamma}_{\tilde{\mathcal{V}}}^\top A(\tilde{\mathcal{V}}) \tilde{\gamma}_{\tilde{\mathcal{V}}}}$, and hence the difference $\sqrt{n} (\widehat{\beta}_E - \beta)$ can be expressed as

$$\sqrt{n} (\widehat{\beta}_E - \beta) = \sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\Gamma}_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} - \beta \right) \mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}^*} + \sum_{\mathcal{V} \neq \mathcal{V}^*} \sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}}^\top A(\mathcal{V}) \tilde{\Gamma}_{\mathcal{V}}}{\tilde{\gamma}_{\mathcal{V}}^\top A(\mathcal{V}) \tilde{\gamma}_{\mathcal{V}}} - \beta \right) \mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}} \quad (33)$$

By Lemma 1, we have $\mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}^*} \xrightarrow{p} 1$ and $\mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}} \xrightarrow{p} 0$ if $\mathcal{V} \neq \mathcal{V}^*$. Combined with Lemma 2 and Slutsky's theorem, we establish

$$\sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\Gamma}_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} - \beta \right) \mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}^*} \xrightarrow{d} N \left(0, \frac{\Theta_{11} + \beta^2 \Theta_{22} - 2\beta \Theta_{12}}{\gamma_{\mathcal{V}^*}^\top A^*(\mathcal{V}^*) \gamma_{\mathcal{V}^*}} \right).$$

Note that for any $\epsilon_0 > 0$,

$$\mathbf{P} \left(\left| \sqrt{n} (\widehat{\beta}_E - \beta) - \sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\Gamma}_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} - \beta \right) \mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}^*} \right| \geq \epsilon_0 \right) \leq \mathbf{P} (\tilde{\mathcal{V}} \neq \mathcal{V}^*) \quad (34)$$

and it follows from Lemma 1 that

$$\sqrt{n} (\widehat{\beta}_E - \beta) - \sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\Gamma}_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} - \beta \right) \mathbf{1}_{\tilde{\mathcal{V}} = \mathcal{V}^*} \xrightarrow{p} 0. \quad (35)$$

By Lemma 3.7 in Wooldridge [2010], we establish (20).

D.3. Preliminary lemmas for high dimension case

We first define the following events for the random design \mathbf{W} (the normalized \mathbf{H}) and the error $\mathbf{\Pi}$,

$$\begin{aligned}
G_1 &= \left\{ \frac{2}{5} \frac{1}{\sqrt{M_1}} < \frac{\|\mathbf{W}_{\cdot j}\|_2}{\sqrt{n}} < \frac{7}{5} \sqrt{M_1} \text{ for } 1 \leq j \leq p \right\}, \\
G_2 &= \left\{ \left| \frac{(\sigma_i^{ora})^2}{\sigma_i^2} - 1 \right| \leq 2\sqrt{\frac{\log p}{n}} + 2\frac{\log p}{n} \text{ for } i = 1, 2 \right\}, \\
G_3 &= \left\{ \left| \frac{\gamma^\top \hat{\Sigma} \gamma}{\gamma^\top \Sigma \gamma} - 1 \right| \leq 12\sqrt{\frac{\log p}{n}} \text{ and } \left| \frac{\Omega_{j\cdot}^\top \hat{\Sigma} \Omega_{j\cdot}}{\Omega_{jj}} - 1 \right| \leq 12\sqrt{\frac{\log p}{n}}, 1 \leq j \leq p_z \right\}, \\
G_4 &= \left\{ \kappa(\mathbf{H}, 4s, 1 + 2\epsilon_0) \geq \frac{1}{2\sqrt{M_1}} \right\}, \\
G_5 &= \left\{ \frac{\|\mathbf{H}^\top \mathbf{\Pi}_i\|_\infty}{n} \leq \sigma_i \sqrt{\frac{2\delta_0 \log p}{n}} \text{ for } i = 1, 2 \right\}, \\
S_1 &= \left\{ \frac{\|\mathbf{H}^\top \mathbf{\Pi}_i\|_\infty}{n} \leq \sigma_i^{ora} \lambda_0 \frac{\epsilon_0 - 1}{\epsilon_0 + 1} (1 - \tau) \text{ for } i = 1, 2 \right\}, \\
S_2 &= \{(1 - \nu_0) \hat{\sigma}_i \leq \sigma_i \leq (1 + \nu_0) \hat{\sigma}_i \text{ for } i = 1, 2\},
\end{aligned} \tag{36}$$

and

$$\begin{aligned}
A_1 &= \left\{ \|e_j^\top \Omega \hat{\Sigma} - e_j^\top\|_\infty \leq \lambda_n, j = 1, 2, \dots, p_z \right\}, \text{ where } \lambda_n = 2eC_0 M_1^2 \sqrt{\frac{\log p}{n}}, \\
A_2 &= \left\{ |\tilde{\gamma}_j - \gamma_j| \leq \frac{\|\hat{\mathbf{v}}^{[j]}\|_2 \sigma_2}{\sqrt{n}} \sqrt{2.05 \log p_z} \text{ for } 1 \leq j \leq p_z \right\}, \\
A_3 &= \left\{ \max_{1 \leq j \leq p_z} \left\| \frac{1}{n} (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_i \right\|_\infty \leq \left(1 + 12\sqrt{\frac{\log p}{n}} \right) M_1 \sqrt{\frac{2.05 \log p_z}{n}} \sigma_i, \text{ for } i = 1, 2 \right\}, \\
A_4 &= \left\{ \frac{2}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{\cdot 2} \leq \frac{2\sqrt{\log p}}{\sqrt{n}} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2 \sqrt{\Theta_{22}} \right\}, \\
A_5 &= \left\{ \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top (\mathbf{\Pi}_{\cdot 1} + \beta \mathbf{\Pi}_{\cdot 2}) \leq \frac{\sqrt{\log p}}{\sqrt{n}} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2 \sqrt{\Theta_{11} + \beta^2 \Theta_{22} + 2\beta \Theta_{12}} \right\},
\end{aligned} \tag{37}$$

where $\hat{\Sigma} = \frac{1}{n} \mathbf{W}^\top \mathbf{W}$ and $\hat{\mathbf{v}}^{[j]} = \mathbf{W}^\top \hat{\mathbf{u}}^{[j]}$. Define

$$G = \cap_{i=1}^5 G_i \quad \text{and} \quad S = \cap_{i=1}^2 S_i \quad \text{and} \quad A = \cap_{i=1}^5 A_i.$$

We introduce the following lemmas to control the probability of events G , S and A . The detailed proofs of the following lemmas are presented in Section E.4 and E.5.

LEMMA 3. If $s \leq cn/\log p$, then

$$\mathbf{P}(G) \geq 1 - \frac{6}{p} - 2p^{1-C_1} - \frac{1}{2\sqrt{\pi\delta_0\log p}}p^{1-\delta_0} - 2\exp\left(-\frac{c'n}{M_1^3}\right), \quad (38)$$

and

$$\mathbf{P}(G \cap S) \geq \mathbf{P}(G) - 2\exp\left(-\left(\frac{g_0 + 1 - \sqrt{2g_0 + 1}}{2}\right)n\right) - c''\frac{1}{\sqrt{\log p}}p^{1-\delta_0}, \quad (39)$$

where $g_0 = \nu_0/(2 + 3\nu_0)$ and c, c', c_* and c'' are universal positive constants, not depending on n and p . We also have

$$\mathbf{P}(A_1) \geq 1 - 2p_z p^{1-c_0 C_0^2}, \quad \text{and} \quad \mathbf{P}(A_4 \cap A_5) \geq 1 - p^{-c}, \quad (40)$$

$$\min\{\mathbf{P}(A_2), \mathbf{P}(A_3)\} \geq \mathbf{P}((A_1 \cap G_1 \cap G_3)) - \frac{1}{2\sqrt{\pi\log p_z}}p_z^{-0.02}. \quad (41)$$

LEMMA 4. On the event $A_1 \cap G_1 \cap G_3$, we have

$$\frac{(1 - \lambda_n)^2}{2M_1} \leq \frac{\|\widehat{\mathbf{v}}^{[j]}\|_2^2}{n} \leq \left(1 + 12\sqrt{\frac{\log p}{n}}\right)M_1, \quad \text{for } 1 \leq j \leq p_z. \quad (42)$$

If $s_{z1}\sqrt{\log p/n} \rightarrow 0$, on the event G_3 , we have

$$\frac{1}{n} \left\| \sum_{j \in S^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 \geq \frac{M_1 \|\gamma\|_2^2 (1 - s_{z1}\lambda_n)^2}{1 - 12\sqrt{\frac{\log p}{n}}} \quad \text{and} \quad \frac{1}{n} \left\| \sum_{j \in \mathcal{V}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 \geq \frac{M_1 \|\gamma_{\mathcal{V}^*}\|_2^2 (1 - s_{z1}\lambda_n)^2}{1 - 12\sqrt{\frac{\log p}{n}}}. \quad (43)$$

Furthermore, we have

$$\frac{M_1 (1 - s_{z1}\lambda_n)^2}{\|\gamma\|_2^2 \left(1 - 12\sqrt{\frac{\log p}{n}}\right)} \frac{1}{M_2} \leq V_H \leq \frac{4s_{z1}M_1^2 M_2 (1 + \beta^2)}{\|\gamma\|_2^2}, \quad (44)$$

and

$$\frac{M_1 (1 - s_{z1}\lambda_n)^2}{\|\gamma_{\mathcal{V}^*}\|_2^2 \left(1 - 12\sqrt{\frac{\log p}{n}}\right)} \frac{1}{M_2} \leq V \leq \frac{4s_{z1}M_1^2 M_2 (1 + \beta^2)}{\|\gamma_{\mathcal{V}^*}\|_2^2}. \quad (45)$$

D.4. Proof of Theorem 3

The proof of Theorem 3 is based on Lemma 5 and the following expression for the estimator $\widehat{\beta}_H$, $\widehat{\beta}_H = \widehat{\gamma}^\top \widehat{\Gamma} / \|\widehat{\gamma}\|_2^2$, where $\|\widehat{\gamma}\|_2^2 = \sum_{j \in \widehat{\mathcal{S}}} \widetilde{\gamma}_j^2$ and $\widehat{\gamma}^\top \widehat{\Gamma} = \sum_{j \in \widehat{\mathcal{S}}} \widetilde{\gamma}_j \widetilde{\Gamma}_j$.

LEMMA 5. Suppose that $\sqrt{s_{z1}}s\log p/\sqrt{n} \rightarrow 0$, $\boldsymbol{\pi}^* = 0$ and the assumptions (R1)–(R3) hold. Then we have the following decompositions,

$$\sqrt{n} \left(\widehat{\|\boldsymbol{\gamma}\|_2^2} - \|\boldsymbol{\gamma}\|_2^2 \right) = \frac{2}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top \boldsymbol{\Pi}_{\cdot 2} + R^\gamma, \quad (46)$$

and

$$\sqrt{n} \left(\widehat{\boldsymbol{\gamma}^\top \boldsymbol{\Gamma}} - \boldsymbol{\gamma}^\top \boldsymbol{\Gamma} \right) = \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top (\boldsymbol{\Pi}_{\cdot 1} + \beta \boldsymbol{\Pi}_{\cdot 2}) + R^{\text{inter}}, \quad (47)$$

where

$$\begin{aligned} \frac{2}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top \boldsymbol{\Pi}_{\cdot 2} &\sim N \left(0, \frac{4}{n} \left\| \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 \boldsymbol{\Theta}_{22} \right), \\ \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top (\boldsymbol{\Pi}_{\cdot 1} + \beta \boldsymbol{\Pi}_{\cdot 2}) &\sim N \left(0, \frac{1}{n} \left\| \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 (\boldsymbol{\Theta}_{11} + \beta^2 \boldsymbol{\Theta}_{22} + 2\beta \boldsymbol{\Theta}_{12}) \right), \end{aligned} \quad (48)$$

and on the event $A \cap S \cap G$, we have

$$\max \{ |R^\gamma|, |R^{\text{inter}}| \} \leq C(|\beta| + 1) \|\boldsymbol{\gamma}\|_2 \sqrt{s_{z1}} s \frac{\log p}{\sqrt{n}} + C s_{z1} \frac{\log p_z}{\sqrt{n}}. \quad (50)$$

Then on the event $A \cap S \cap G$, we have

$$\max \left\{ \left| \widehat{\|\boldsymbol{\gamma}\|_2^2} - \|\boldsymbol{\gamma}\|_2^2 \right|, \left| \widehat{\boldsymbol{\gamma}^\top \boldsymbol{\Gamma}} - \boldsymbol{\gamma}^\top \boldsymbol{\Gamma} \right| \right\} \leq C \|\boldsymbol{\gamma}\|_2 s_{z1} \sqrt{\frac{\log p}{n}} + C s_{z1} \frac{\log p_z}{n} \leq C \|\boldsymbol{\gamma}\|_2 s_{z1} \sqrt{\frac{\log p}{n}}. \quad (51)$$

In the following, we will prove (27) in the main paper. Note that

$$\widetilde{\beta} - \beta = -\frac{\beta}{\|\boldsymbol{\gamma}\|_2^2} \left(\widehat{\|\boldsymbol{\gamma}\|_2^2} - \|\boldsymbol{\gamma}\|_2^2 \right) + \frac{1}{\|\boldsymbol{\gamma}\|_2^2} \left(\widehat{\boldsymbol{\gamma}^\top \boldsymbol{\Gamma}} - \boldsymbol{\gamma}^\top \boldsymbol{\Gamma} \right) + \frac{\|\boldsymbol{\gamma}\|_2^2 - \widehat{\|\boldsymbol{\gamma}\|_2^2}}{\|\boldsymbol{\gamma}\|_2^2} \left(\frac{\widehat{\boldsymbol{\gamma}^\top \boldsymbol{\Gamma}}}{\widehat{\|\boldsymbol{\gamma}\|_2^2}} - \frac{\boldsymbol{\gamma}^\top \boldsymbol{\Gamma}}{\|\boldsymbol{\gamma}\|_2^2} \right). \quad (52)$$

By Lemma 5, we have the following decomposition,

$$\sqrt{n} \left(\widetilde{\beta} - \beta \right) = T^\beta + \Delta^\beta, \quad (53)$$

where

$$\begin{aligned} T^\beta &= -\frac{\beta}{\|\boldsymbol{\gamma}\|_2^2} \frac{2}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top \boldsymbol{\Pi}_{\cdot 2} + \frac{1}{\|\boldsymbol{\gamma}\|_2^2} \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top (\boldsymbol{\Pi}_{\cdot 1} + \beta \boldsymbol{\Pi}_{\cdot 2}) \\ &= \frac{1}{\|\boldsymbol{\gamma}\|_2^2} \frac{1}{\sqrt{n}} \sum_{j \in \mathcal{S}^*} \boldsymbol{\gamma}_j (\widehat{\mathbf{v}}^{[j]})^\top (\boldsymbol{\Pi}_{\cdot 1} - \beta \boldsymbol{\Pi}_{\cdot 2}), \end{aligned}$$

and $\Delta^\beta = \text{Res}_1 + \text{Res}_2$ with

$$\text{Res}_1 = \frac{1}{\|\gamma\|_2^2} (-\beta R^\gamma + R^{\text{inter}}) \text{ and } \text{Res}_2 = \sqrt{n} \frac{\|\gamma\|_2^2 - \widehat{\|\gamma\|_2^2}}{\|\gamma\|_2^2} \left(\frac{\widehat{\gamma^\top \Gamma}}{\widehat{\|\gamma\|_2^2}} - \frac{\gamma^\top \Gamma}{\|\gamma\|_2^2} \right).$$

By the distribution of $\mathbf{\Pi}$, we establish that

$$T^\beta \mid \mathbf{W} \sim N \left(0, \frac{1}{n\|\gamma\|_2^4} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 (\boldsymbol{\Theta}_{11} + \beta^2 \boldsymbol{\Theta}_{22} - 2\beta \boldsymbol{\Theta}_{12}) \right). \quad (54)$$

By Lemma 5, on the event $G \cap S \cap A$, we have

$$\frac{1}{\sqrt{V_H}} |\text{Res}_1| \leq C \frac{1}{\|\gamma\|_2} (|\beta| |R^\gamma| + |R^{\text{inter}}|) \leq C (|\beta| + 1) \sqrt{s_{z1}} s \frac{\log p}{\sqrt{n}} + C \frac{1}{\|\gamma\|_2} \frac{s_{z1} \log p}{\sqrt{n}}. \quad (55)$$

Note that on the event $G \cap S \cap A$,

$$\frac{1}{\sqrt{V_H}} \text{Res}_2 \leq C \sqrt{n} \frac{\|\gamma\|_2^2 - \widehat{\|\gamma\|_2^2}}{\|\gamma\|_2} \times \frac{(\widehat{\gamma^\top \Gamma} - \gamma^\top \Gamma) + \beta (\|\gamma\|_2^2 - \widehat{\|\gamma\|_2^2})}{\|\gamma\|_2^2 + (\widehat{\|\gamma\|_2^2} - \|\gamma\|_2^2)} \leq C \frac{s_{z1}^3 (\log p)^{\frac{3}{2}}}{n}, \quad (56)$$

where the last inequality follows from (51). Combined with (55), by $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$, we can establish that on the event $G \cap S \cap A$,

$$|\Delta^\beta / \sqrt{V_H}| \leq C \sqrt{s_{z1}} s \frac{\log p}{\sqrt{n}} + C \frac{1}{\|\gamma\|_2} \frac{s_{z1} \log p}{\sqrt{n}}. \quad (57)$$

Since $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$, we establish $\Delta^\beta / \sqrt{V_H} \xrightarrow{p} 0$. Combined with (54), we establish (27).

D.5. Proof of Theorem 4

We first introduce the following lemma to establish the coverage property.

LEMMA 6. *Suppose that $\boldsymbol{\pi}^* = 0$ and the assumptions (R1) – (R3) hold. As $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$, then we have*

$$\frac{\widehat{V}_H}{V_H} \xrightarrow{p} 1. \quad (58)$$

By (54), we have $\frac{T^\beta}{\sqrt{V_H}} \sim N(0, 1)$. Combined with (57) and Lemma 6, we have

$$\sqrt{n} \frac{\widehat{\beta}_H - \beta}{\sqrt{\widehat{V}_H}} = \frac{T^\beta + \Delta^\beta}{\sqrt{V_H}} \times \frac{\sqrt{V_H}}{\sqrt{\widehat{V}_H}} \xrightarrow{d} N(0, 1). \quad (59)$$

and hence the coverage property (28) follows.

D.6. Proof of Theorem 2

The proof of the theorem follows from the following lemma, which characterizes the behavior of the selection process (10) and (12) in the main paper.

LEMMA 7. *Suppose that $\sqrt{s_{z1}}s \log p / \sqrt{n} \rightarrow 0$ and the assumptions (R1) – (R3) and (IN1) – (IN2) are satisfied. With probability larger than $1 - c(p^{-c} + \exp(-cn))$, we have*

$$\tilde{\mathcal{V}} \subset \left\{ i \in \mathcal{S}^* : \left| \frac{\pi_j}{\gamma_j} \right| \leq 2C_* \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}} \right\} \quad \text{and} \quad |\tilde{\mathcal{V}}| > \frac{1}{2} |\mathcal{S}^*|. \quad (60)$$

Under the extra assumption (IN3) in the main paper, with probability larger than $1 - c(p^{-c} + \exp(-cn))$

$$\tilde{\mathcal{V}} = \mathcal{V}^*. \quad (61)$$

We have the following decomposition,

$$\begin{aligned} \hat{\beta} - \beta &= \frac{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\Gamma}_j}{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2} - \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \Gamma_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} + \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \Gamma_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} - \beta \\ &= \left(\frac{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\Gamma}_j}{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2} - \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \Gamma_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} \right) + \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \pi_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2}, \end{aligned} \quad (62)$$

where the first term is taken as the variance term and the second term is taken as the bias term. In the following, we are going to analyze the bias and the variance term separately. For the bias term, we have

$$\left| \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \pi_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} \right| = \left| \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \frac{\pi_j}{\gamma_j}}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} \right| \leq \max_{j \in \tilde{\mathcal{V}}} \left| \frac{\pi_j}{\gamma_j} \right| \leq 2C_* \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}}, \quad (63)$$

where the last inequality follows from (60). The following lemma controls the variance term.

LEMMA 8. *Suppose that $\sqrt{s_{z1}}s \log p / \sqrt{n} \rightarrow 0$ and the assumptions (R1) – (R3) and (IN1) – (IN3) are satisfied. On the event $A \cap S \cap G$, we have*

$$\left| \frac{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\Gamma}_j}{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2} - \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \Gamma_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} \right| \leq C \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}}. \quad (64)$$

Combining (63) and (64), we can establish (23) in the main paper. Under the stronger assumption (22) in the main paper, we can establish (61) and the decomposition (62) holds as

$$\hat{\beta} - \beta = \left(\frac{\sum_{j \in \mathcal{V}^*} \tilde{\gamma}_j \tilde{\Gamma}_j}{\sum_{j \in \mathcal{V}^*} \tilde{\gamma}_j^2} - \frac{\sum_{j \in \mathcal{V}^*} \gamma_j \Gamma_j}{\sum_{j \in \mathcal{V}^*} \gamma_j^2} \right).$$

Based on the this expression, (24) in the main paper will follow the same argument with (27), which is presented in Section D.4. We introduce the following lemma to establish the coverage property.

LEMMA 9. Suppose the assumptions (R1) – (R5) and (IN1) – (IN3) are satisfied. As $\sqrt{s_{z1}}s \log p / \sqrt{n} \rightarrow 0$, we have

$$\frac{\widehat{V}}{\overline{V}} \xrightarrow{p} 1. \quad (65)$$

Similarly to the proof of Theorem 4 in Section D.5, we establish

$$\sqrt{n} \frac{\widehat{\beta} - \beta}{\sqrt{\widehat{V}}} = \frac{T^\beta + \Delta^\beta}{\sqrt{\overline{V}}} \times \frac{\sqrt{\overline{V}}}{\sqrt{\widehat{V}}} \xrightarrow{d} N(0, 1), \quad (66)$$

and hence establish the coverage property (25) in the main paper.

E. Supplementary Materials: Proof of Extra Lemmas

In this section, we prove extra lemmas used in the proof of main theorems.

E.1. Proof of Lemma 1

Define $\mathcal{I} = \{1, 2, \dots, p_z\}$. We first note the following expression for $\tilde{\gamma}_j$ and $\tilde{\Gamma}_j$ for $i \in \mathcal{I}$,

$$\sqrt{n}(\tilde{\gamma}_j - \gamma_j) = \left(\widehat{\Sigma}^{-1}\right)_{j,\cdot} \frac{1}{\sqrt{n}} \mathbf{W}^\top \mathbf{\Pi}_{\cdot 2} \quad \text{and} \quad \sqrt{n}(\tilde{\Gamma}_j - \Gamma_j) = \left(\widehat{\Sigma}^{-1}\right)_{j,\cdot} \frac{1}{\sqrt{n}} \mathbf{W}^\top \mathbf{\Pi}_{\cdot 1} \quad (67)$$

and the following limiting theorem (Theorem 3.1 in Wooldridge [2010]),

$$\tilde{\gamma} \xrightarrow{p} \gamma \quad \text{and} \quad \tilde{\Gamma} \xrightarrow{p} \Gamma, \quad (68)$$

$$\sqrt{n}(\tilde{\gamma} - \gamma) \xrightarrow{d} N\left(0, \Theta_{22}(\Sigma^{-1})_{\mathcal{I}, \mathcal{I}}\right) \quad \text{and} \quad \sqrt{n}(\tilde{\Gamma} - \Gamma) \xrightarrow{d} N\left(0, \Theta_{11}(\Sigma^{-1})_{\mathcal{I}, \mathcal{I}}\right). \quad (69)$$

Note that

$$\frac{\sqrt{\widehat{\Theta}_{22}} \|\mathbf{W}(\widehat{\Sigma}^{-1})_{\cdot j}\|_2}{\sqrt{n}} \xrightarrow{p} \sqrt{\Theta_{22}(\Sigma^{-1})_{jj}}. \quad (70)$$

We define the following events

$$\begin{aligned} \mathcal{B}_1 &= \{\tilde{\mathcal{S}} = \mathcal{S}^*\} \\ \mathcal{B}_2 &= \left\{ \max_{j \in \mathcal{V}^*} \|\tilde{\pi}^{[j]}\|_0 < \frac{|\mathcal{S}^*|}{2} < \min_{j \in \mathcal{S}^* \setminus \mathcal{V}^*} \|\tilde{\pi}^{[j]}\|_0 \right\} \\ \mathcal{B}_3 &= \left\{ \text{supp}(\tilde{\pi}_{\mathcal{S}^*}^{[j]}) = \text{supp}(\pi_{\mathcal{S}^*}) \quad \text{for } j \in \mathcal{V}^* \right\} \end{aligned} \quad (71)$$

On the event $\mathcal{B} = \mathcal{B}_1 \cap \mathcal{B}_2 \cap \mathcal{B}_3$, we have $\tilde{\mathcal{V}} = \mathcal{V}^*$ and it is sufficient to show that

$$\lim_{n \rightarrow \infty} \mathbf{P}(\mathcal{B}) = 0 \quad (72)$$

For $j \in \mathcal{S}$, we have

$$|\tilde{\gamma}_j| - \frac{\sqrt{\hat{\Theta}_{22}} \|\mathbf{W}(\hat{\Sigma}^{-1})_{\cdot j}\|_2}{\sqrt{n}} \sqrt{\frac{a_0 \log n}{n}} \xrightarrow{p} |\gamma_j| > 0, \quad (73)$$

where the convergence follows from (68) and (70). For $j \in \mathcal{S}^c$, we have

$$\sqrt{\frac{n}{a_0 \log n}} |\tilde{\gamma}_j| - \frac{\sqrt{\hat{\Theta}_{22}} \|\mathbf{W}(\hat{\Sigma}^{-1})_{\cdot j}\|_2}{\sqrt{n}} \xrightarrow{p} -\sqrt{\Theta_{22}(\Sigma^{-1})_{jj}} < 0, \quad (74)$$

where the convergence follows from (69) and (70). Combining (73) and (74), we establish that

$$\lim_{n \rightarrow \infty} \mathbf{P}(\mathcal{B}_1) = 1. \quad (75)$$

In the following, we control $\lim_{n \rightarrow \infty} \mathbf{P}(\mathcal{B})$. Without loss of generality, we assume $1 \in \tilde{\mathcal{S}}$ and focus on the case $i = 1$. In the following, we are going to analyze the performance of $\hat{\beta}^{[1]}$ and $\hat{\pi}_j^{[1]}$. Note that

$$\hat{\beta}^{[1]} - \frac{\Gamma_1}{\gamma_1} \xrightarrow{p} 0. \quad (76)$$

and hence we have

$$\begin{aligned} & \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \frac{\|\mathbf{W}((\hat{\Sigma}^{-1})_{k\cdot} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_1}(\hat{\Sigma}^{-1})_{1\cdot})\|_2}{\sqrt{n}} \\ & \xrightarrow{p} \sqrt{\Theta_{11} + \left(\frac{\Gamma_1}{\gamma_1}\right)^2 \Theta_{22} - 2\frac{\Gamma_1}{\gamma_1} \Theta_{12}} \sqrt{(\Sigma^{-1})_{kk} + \left(\frac{\gamma_k}{\gamma_1}\right)^2 (\Sigma^{-1})_{11} - 2\frac{\gamma_k}{\gamma_1} (\Sigma^{-1})_{k1}}. \end{aligned} \quad (77)$$

We also have the following expression

$$\begin{aligned} \hat{\pi}_k^{[1]} - \left(\Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k\right) &= \left(\tilde{\Gamma}_k - \frac{\tilde{\Gamma}_1}{\tilde{\gamma}_1} \tilde{\gamma}_k\right) - \left(\Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k\right) \\ &= \left(\tilde{\Gamma}_k - \Gamma_k\right) - \frac{\Gamma_1}{\gamma_1} (\tilde{\gamma}_k - \gamma_k) - \frac{\gamma_k}{\gamma_1^2} \left(\gamma_1 (\tilde{\Gamma}_1 - \Gamma_1) - \Gamma_1 (\tilde{\gamma}_1 - \gamma_1)\right) \\ &\quad + \left(\frac{\tilde{\Gamma}_1}{\tilde{\gamma}_1} - \frac{\Gamma_1}{\gamma_1}\right) \left(\frac{\gamma_k}{\gamma_1} (\tilde{\gamma}_1 - \gamma_1) - (\tilde{\gamma}_k - \gamma_k)\right) \end{aligned} \quad (78)$$

Note that

$$\begin{aligned} & \sqrt{n} \left(\left(\tilde{\Gamma}_k - \Gamma_k\right) - \frac{\Gamma_1}{\gamma_1} (\tilde{\gamma}_k - \gamma_k) - \frac{\gamma_k}{\gamma_1^2} \left(\gamma_1 (\tilde{\Gamma}_1 - \Gamma_1) - \Gamma_1 (\tilde{\gamma}_1 - \gamma_1)\right) \right) \\ &= \left(\left(\hat{\Sigma}^{-1}\right)_{\cdot k} - \frac{\gamma_k}{\gamma_1} \left(\hat{\Sigma}^{-1}\right)_{\cdot 1} \right) \frac{1}{\sqrt{n}} \mathbf{W}^\top \left(\boldsymbol{\Pi}_{\cdot 2} - \frac{\Gamma_1}{\gamma_1} \boldsymbol{\Pi}_{\cdot 1} \right) \\ & \xrightarrow{d} N \left(0, \Theta_{11} + \left(\frac{\Gamma_1}{\gamma_1}\right)^2 \Theta_{22} - 2\frac{\Gamma_1}{\gamma_1} \Theta_{12} (\Sigma^{-1})_{kk} + \left(\frac{\gamma_k}{\gamma_1}\right)^2 (\Sigma^{-1})_{11} - 2\frac{\gamma_k}{\gamma_1} (\Sigma^{-1})_{k1} \right), \end{aligned} \quad (79)$$

where the convergence follows from Theorem 3.1 in Wooldridge [2010]. By (68) and (69), we have

$$\left(\frac{\tilde{\Gamma}_1}{\tilde{\gamma}_1} - \frac{\Gamma_1}{\gamma_1} \right) \left(\frac{\gamma_k}{\gamma_1} (\tilde{\gamma}_1 - \gamma_1) - (\tilde{\gamma}_k - \gamma_k) \right) \xrightarrow{p} 0.$$

Combined with (76) and (79), we have

$$\frac{\sqrt{n}}{\sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \frac{\|\mathbf{W}((\hat{\Sigma}^{-1})_{k\cdot} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_1} (\hat{\Sigma}^{-1})_{1\cdot})\|_2}{\sqrt{n}}} \left(\hat{\pi}_k^{[1]} - \left(\Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k \right) \right) \xrightarrow{p} N(0, 1) \quad (80)$$

and hence

$$\hat{\pi}_k^{[1]} \xrightarrow{p} \Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k. \quad (81)$$

We divide the discussion into the following three cases,

- $1 \in \mathcal{S}^* \setminus \mathcal{V}^*$ and $k \in \mathcal{V}^*$;
- $1 \in \mathcal{V}^*$ and $k \in \mathcal{V}^*$;
- $1 \in \mathcal{V}^*$ and $k \in \mathcal{S}^* \setminus \mathcal{V}^*$.

$1 \in \mathcal{S}^* \setminus \mathcal{V}^*$ and $k \in \mathcal{V}^*$

In this case, $\Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k = \frac{\pi_1}{\gamma_1} \gamma_k \neq 0$. Hence, we have

$$\begin{aligned} & \left| \hat{\pi}_k^{[1]} \right| - 2.05 \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \frac{\|\mathbf{W}((\hat{\Sigma}^{-1})_{k\cdot} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_1} (\hat{\Sigma}^{-1})_{1\cdot})\|_2}{\sqrt{n}} \sqrt{\frac{\log n}{n}} \\ & \xrightarrow{p} \left| \frac{\pi_1}{\gamma_1} \gamma_k \right| > 0 \end{aligned} \quad (82)$$

and

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(k \in \text{supp} \left(\tilde{\pi}_{\mathcal{S}^*}^{[1]} \right) \right) = 1. \quad (83)$$

Hence, by the assumption (IN1),

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\|\tilde{\pi}_{\mathcal{S}^*}^{[1]}\|_0 > \frac{|\mathcal{V}^*|}{2} \right) = 1. \quad (84)$$

$1 \in \mathcal{V}^*$ and $k \in \mathcal{V}^*$

In this case, $\Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k = 0$. By (80), we have

$$\frac{\left| \hat{\pi}_k^{[1]} \right|}{2.05 \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \frac{\|\mathbf{W}((\hat{\Sigma}^{-1})_{k\cdot} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_1} (\hat{\Sigma}^{-1})_{1\cdot})\|_2}{\sqrt{n}} \sqrt{\frac{\log n}{n}}} \xrightarrow{p} 0 \quad (85)$$

and

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(k \notin \text{supp} \left(\tilde{\pi}_{\mathcal{S}^*}^{[1]} \right) \right) = 1. \quad (86)$$

Hence, by the assumption (IN1),

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\|\tilde{\pi}_{\mathcal{S}^*}^{[1]}\|_0 < \frac{|\mathcal{V}^*|}{2} \right) = 1. \quad (87)$$

$1 \in \mathcal{V}^*$ and $k \in \mathcal{S}^* \setminus \mathcal{V}^*$

In this case, $\Gamma_k - \frac{\Gamma_1}{\gamma_1} \gamma_k = \pi_k \neq 0$. Hence, we have

$$\left| \hat{\pi}_k^{[1]} \right| - 2.05 \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \frac{\|\mathbf{W}((\hat{\Sigma}^{-1})_{k\cdot} - \frac{\tilde{\gamma}_k}{\tilde{\gamma}_1} (\hat{\Sigma}^{-1})_{1\cdot})\|_2}{\sqrt{n}} \sqrt{\frac{\log n}{n}} \xrightarrow{p} |\pi_k| > 0 \quad (88)$$

and hence

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(k \in \text{supp} \left(\tilde{\pi}_{\mathcal{S}^*}^{[1]} \right) \right) = 1. \quad (89)$$

Since we can replace the index 1 with any index $j \in \tilde{\mathcal{V}}$, then (84) and (87) can be correspondingly replaced by

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\|\tilde{\pi}_{\mathcal{S}^*}^{[j]}\|_0 > \frac{|\mathcal{V}^*|}{2} \right) = 1 \quad \text{for } j \in \mathcal{S}^* \setminus \mathcal{V}^*; \quad (90)$$

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\|\tilde{\pi}_{\mathcal{S}^*}^{[j]}\|_0 < \frac{|\mathcal{V}^*|}{2} \right) = 1 \quad \text{for } j \in \mathcal{V}^*. \quad (91)$$

By (90) and (91), we can establish $\lim_{n \rightarrow \infty} \mathbf{P}(\mathcal{B}_1 \cap \mathcal{B}_2) = 1$. By (86) and (89), we have

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\text{supp} \left(\tilde{\pi}_{\mathcal{S}^*}^{[1]} \right) = \text{supp}(\pi_{\mathcal{S}^*}) \right) = 1. \quad (92)$$

Similarly, we can obtain that

$$\lim_{n \rightarrow \infty} \mathbf{P} \left(\text{supp} \left(\tilde{\pi}_{\mathcal{S}^*}^{[j]} \right) = \text{supp}(\pi_{\mathcal{S}^*}) \quad \text{for } j \in \mathcal{V}^* \right) = 1, \quad (93)$$

and hence $\lim_{n \rightarrow \infty} \mathbf{P}(\mathcal{B}_1 \cap \mathcal{B}_2 \cap \mathcal{B}_3) = 1$.

E.2. Proof of Lemma 2

Note that

$$\sqrt{n} \left(\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\Gamma}_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} - \beta \right) = \frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \left(\hat{\Sigma}^{-1} \right)_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} \frac{1}{\sqrt{n}} \mathbf{W}^\top (\mathbf{\Pi}_{\cdot 2} - \beta \mathbf{\Pi}_{\cdot 1}). \quad (94)$$

Since

$$\frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \left(\hat{\Sigma}^{-1} \right)_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} \xrightarrow{p} \frac{\gamma_{\mathcal{V}^*}^\top A^*(\mathcal{V}^*) \left(\Sigma^{-1} \right)_{\mathcal{V}^*}}{\gamma_{\mathcal{V}^*}^\top A^*(\mathcal{V}^*) \gamma_{\mathcal{V}^*}} \quad (95)$$

and

$$\frac{1}{\sqrt{n}} \mathbf{W}^\top (\mathbf{\Pi}_{\cdot 2} - \beta \mathbf{\Pi}_{\cdot 1}) \xrightarrow{d} N \left(0, \left(\Theta_{11} + \beta^2 \Theta_{22} - 2\beta \Theta_{12} \right) \Sigma \right), \quad (96)$$

we have

$$\begin{aligned} & \frac{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) (\hat{\Sigma}^{-1})_{\mathcal{V}^*}}{\tilde{\gamma}_{\mathcal{V}^*}^\top A(\mathcal{V}^*) \tilde{\gamma}_{\mathcal{V}^*}} \frac{1}{\sqrt{n}} \mathbf{W}^\top (\mathbf{\Pi}_{.2} - \beta \mathbf{\Pi}_{.1}) \\ & \xrightarrow{d} N \left(0, \frac{\gamma_{\mathcal{V}^*}^\top A^*(\mathcal{V}^*) (\Sigma^{-1})_{\mathcal{V}^* \mathcal{V}^*} A^*(\mathcal{V}^*) \gamma_{\mathcal{V}^*}}{(\gamma_{\mathcal{V}^*}^\top A^*(\mathcal{V}^*) \gamma_{\mathcal{V}^*})^2} (\Theta_{11} + \beta^2 \Theta_{22} - 2\beta \Theta_{12}) \right) \end{aligned} \quad (97)$$

Since $A^*(\mathcal{V}^*) (\Sigma^{-1})_{\mathcal{V}^* \mathcal{V}^*} A^*(\mathcal{V}^*) = A^*(\mathcal{V}^*)$, we establish (32).

E.3. Lemmas for scaled Lasso and de-biasing Lasso

We introduce the following lemmas for scaled Lasso and de-biasing Lasso used in the later proofs. Lemma 10 establishes the convergence rate of the scaled Lasso method, which is based on the analysis in Sun and Zhang [2012].

LEMMA 10. *On the event $G \cap S$, if $s \leq cn/\log p$, then*

$$\|\hat{\Gamma} - \Gamma\|_1 + \|\hat{\Psi} - \Psi\|_1 \leq Cs \sqrt{\frac{\log p}{n}} \sigma_1, \quad \|\hat{\gamma} - \gamma\|_1 + \|\hat{\psi} - \psi\|_1 \leq Cs \sqrt{\frac{\log p}{n}} \sigma_2, \quad (98)$$

$$\frac{1}{\sqrt{n}} \|\mathbf{Z}(\hat{\Gamma} - \Gamma) + \mathbf{X}(\hat{\Psi} - \Psi)\|_2 \leq C \sqrt{\frac{s \log p}{n}} \sigma_1, \quad (99)$$

and

$$\frac{1}{\sqrt{n}} \|\mathbf{Z}(\hat{\gamma} - \gamma) + \mathbf{X}(\hat{\psi} - \psi)\|_2 \leq C \sqrt{\frac{s \log p}{n}} \sigma_2. \quad (100)$$

The following lemma is the key result for the de-biasing Lasso estimator, established in Zhang and Zhang [2014], Javanmard and Montanari [2014], van de Geer et al. [2014].

LEMMA 11. *We have the following expressions for the proposed de-biased estimator,*

$$\tilde{\Gamma} - \Gamma = D^\Gamma + \Delta^\Gamma, \quad (101)$$

where

$$D_j^\Gamma = \frac{1}{n} (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.1} \quad \text{and} \quad \Delta_j^\Gamma = \left(\frac{1}{n} (\hat{\mathbf{u}}^{[j]})^\top \hat{\Sigma} - e_j^\top \right) \begin{pmatrix} \hat{\Gamma} - \Gamma \\ \hat{\Psi} - \Psi \end{pmatrix}, \quad i = 1, \dots, p_z. \quad (102)$$

We also have

$$\tilde{\gamma} - \gamma = D^\gamma + \Delta^\gamma, \quad (103)$$

where

$$D_j^\gamma = \frac{1}{n} (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2} \quad \text{and} \quad \Delta_j^\gamma = \left(\frac{1}{n} (\hat{\mathbf{u}}^{[j]})^\top \hat{\Sigma} - e_j^\top \right) \begin{pmatrix} \hat{\gamma} - \gamma \\ \hat{\psi} - \psi \end{pmatrix}, \quad i = 1, \dots, p_z. \quad (104)$$

On the event $S \cap G \cap A$, we have

$$\max \{ \|\Delta^\gamma\|_\infty, \|\Delta^\Gamma\|_\infty \} \leq Cs \frac{\log p}{n} \max \{ \sigma_1, \sigma_2 \}. \quad (105)$$

E.4. Proof of Lemma 3

The proof of Lemma 3 is a generalization of Lemma 4 in Cai and Guo [2016a]. In the following, we extend the Gaussian design in Cai and Guo [2016a] to sub-gaussian design considered in this paper. Since the error of the regression is still assumed to be Gaussian, it is sufficient to establish the probability bound of G_1, G_3, G_4 and A_1 for the sub-gaussian design matrix and control the events A_2 and A_3 . The probability bound of the event A_1 for the sub-gaussian design is established in Lemma 4 of Cai and Guo [2016b]. By Corollary 5.17 in Vershynin [2012] and the union bound, we have

$$\mathbf{P} \left(\max_{1 \leq j \leq p} \left| \frac{1}{n} (\|\mathbf{W}_{\cdot j}\|_2^2 - \mathbf{E}\|\mathbf{W}_{\cdot j}\|_2^2) \right| \geq \epsilon \right) \leq 2p \exp \left(-\frac{1}{6} \min \left\{ \frac{\epsilon^2}{K^2}, \frac{\epsilon}{K} \right\} n \right),$$

where $K = 4M_1$. Taking $\epsilon = 12M_1 \sqrt{\log p/n}$, we have

$$\mathbf{P} \left(\max_{1 \leq j \leq p} \left| \frac{1}{n} (\|\mathbf{W}_{\cdot j}\|_2^2 - \mathbf{E}\|\mathbf{W}_{\cdot j}\|_2^2) \right| \geq 12M_1 \sqrt{\frac{\log p}{n}} \right) \leq 2p^{-\frac{1}{2}} \quad \text{and} \quad \mathbf{P}(G_1) \geq 1 - 2p^{-\frac{1}{2}}. \quad (106)$$

Similarly, we have $\mathbf{P} \left(\left| \frac{1}{n} (\|\mathbf{W}u\|_2^2 - \mathbf{E}\|\mathbf{W}u\|_2^2) \right| \geq 12M_1 \|u\|^2 \sqrt{\log p/n} \right) \leq 2p^{-\frac{3}{2}}$, and

$$\mathbf{P} \left(\left| \frac{1}{n} \left(\frac{\|\mathbf{W}u\|_2^2}{\mathbf{E}\|\mathbf{W}u\|_2^2} - 1 \right) \right| \geq 12M_1 \frac{\|u\|^2}{\mathbf{E}\|\mathbf{W}u\|_2^2} \sqrt{\log p/n} \right) \leq 2p^{-\frac{3}{2}},$$

and hence

$$\mathbf{P}(G_3) \geq 1 - (p_z + 2)p^{-\frac{3}{2}}.$$

By Theorem 1.6 in Zhou [2009], if $n \geq 1/\theta^2 \times c' M_1^3 \max \{12(2 + \gamma_0)^2 M_1 s \log(5ep/4s), 9 \log p\}$, then with probability at least $1 - 2 \exp(-c\theta^2 n/M_1^3)$, for all δ such that there exists $|J_0| \leq 4s$ and $\|\delta_{J_0^c}\|_1 \leq \gamma_0 \|\delta_{J_0}\|_1$, we have $\|Z\delta\|_2/(\sqrt{n}\|\Sigma^{\frac{1}{2}}\delta\|_2) \geq 1 - \theta$. By taking $\theta = \frac{1}{2}$, if $n \geq 4c' M_1^3 \max \{12(2 + \gamma_0)^2 M_1 s \log(5ep/4s), 9 \log p\}$, then $\mathbf{P}(G_4) \geq 1 - 2 \exp(-cn/M_1^3)$. In the following, we control the events A_2 and A_3 ,

$$\begin{aligned} \mathbf{P}(A_2^c) &\leq \mathbf{P} \left(\max_{1 \leq i \leq q} \frac{|D_j^\gamma|}{\sqrt{\text{Var}(D_j^\gamma)}} \geq \sqrt{2.02 \log p_z} \right) + \mathbf{P} \left(\max_{1 \leq j \leq p_z} \frac{|\Delta_j^\gamma|}{\sqrt{\text{Var}(D_j^\gamma)}} \geq 0.01 \sqrt{\log p_z} \right) \\ &\leq \frac{1}{2\sqrt{\pi \log p_z}} p_z^{-0.02} + \mathbf{P}((S \cap G \cap A_1)^c), \end{aligned}$$

where the first inequality follows from (103) and the second inequality follows from (104) and (105). The control of $\mathbf{P}(A_4 \cap A_5)$ follows from (48) and (49). Note that

$$\begin{aligned} \mathbf{P}(A_3^c) &\leq \mathbf{P}((A_1 \cap G_1 \cap G_3)^c) + \mathbf{P}(A_3^c \cap A_1 \cap G_1 \cap G_3) \\ &\leq \mathbf{P}((A_1 \cap G_1 \cap G_3)^c) + 2\mathbf{P} \left(\max_{1 \leq j \leq p_z} \frac{1}{\|\hat{\mathbf{v}}^{[j]}\|_2 \sigma_1} \left| (\hat{\mathbf{v}}^{[j]})^\top \boldsymbol{\Pi}_{\cdot 1} \right| \geq \sqrt{2.05 \log p_z} \right) \\ &\leq \mathbf{P}((A_1 \cap G_1 \cap G_3)^c) + \frac{1}{\sqrt{\pi \log p_z}} p_z^{-0.05}, \end{aligned}$$

where the second inequality follows from (42) and the last inequality follows from the fact that $1/(\|\hat{\mathbf{v}}^{[j]}\|_2 \sigma_1) \times (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{\cdot 1}$ conditioning on \mathbf{W} is normally distributed.

E.5. Proof of Lemma 4

In the following, we only establish the results for $\hat{\mathbf{v}}^{[1]}$ and the same argument extends to $\hat{\mathbf{v}}^{[j]}$ where $1 \leq j \leq p_z$. Since $\lambda_n = 2eC_0M_1^2\sqrt{\log p/n}$ is chosen such that $\mathbf{\Omega}_1$ belongs to the feasible set, we have

$$\frac{\|\hat{\mathbf{v}}^{[1]}\|_2^2}{n} \leq \frac{\|\mathbf{W}\mathbf{\Omega}_1\|_2^2}{n}. \quad (107)$$

By Lemma 12 in Javanmard and Montanari [2014], we have

$$\frac{\|\hat{\mathbf{v}}^{[1]}\|_2^2}{n} \geq \frac{(1 - \lambda_n)^2}{\hat{\Sigma}_{11}}. \quad (108)$$

By the definition of G_1 and G_3 , we establish (42). Let $\mathcal{I} = \{1, 2, \dots, p_z\}$ and assume that $M \in \mathbb{R}^{p \times p_z}$ belongs to the feasible set $\|\hat{\Sigma}\mathbf{\Omega} - \mathbf{I}_{\mathcal{I}}\|_\infty \leq \lambda_n$, where $\mathbf{I}_{\mathcal{I}}$ denotes the sub-matrix of the identity matrix containing the column with index $i \in \mathcal{I}$, that is, $\|\hat{\Sigma}M - \mathbf{I}_{\mathcal{I}}\|_\infty \leq \lambda_n$, and hence

$$\|\hat{\Sigma}M\gamma - \gamma\|_\infty = \|(\hat{\Sigma}M - \mathbf{I}_{\mathcal{I}})\gamma\|_\infty \leq \|\hat{\Sigma}M - \mathbf{I}_{\mathcal{I}}\|_\infty \|\gamma\|_1 \leq \lambda_n \|\gamma\|_1. \quad (109)$$

Note that

$$\left| \gamma^\top \hat{\Sigma}M\gamma - \|\gamma\|_2^2 \right| = \left| \gamma^\top (\hat{\Sigma}M\gamma - \gamma) \right| \leq \|\gamma\|_1 \|\hat{\Sigma}M\gamma - \gamma\|_\infty \leq \lambda_n \|\gamma\|_1^2, \quad (110)$$

where the last inequality follows from (109). The inequality (110) informs that $M\gamma$ is in the feasible set

$$\left| \gamma^\top \hat{\Sigma}(M\gamma) - \|\gamma\|_2^2 \right| \leq \lambda_n \|\gamma\|_1^2. \quad (111)$$

We define μ^* as

$$\begin{aligned} \mu^* &= \arg \min_{\mu} \mu^\top \hat{\Sigma} \mu \\ \text{subject to } & \left| \gamma^\top \hat{\Sigma} \mu - \|\gamma\|_2^2 \right| \leq \lambda_n \|\gamma\|_1^2 \end{aligned} \quad (112)$$

By (111), we have the following inequality,

$$\frac{1}{n} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2^2 = \gamma^\top M^\top \hat{\Sigma} M \gamma \geq (\mu^*)^\top \hat{\Sigma} \mu^*. \quad (113)$$

In the following, we will show that $(\mu^*)^\top \hat{\Sigma} \mu^* = \langle \mu^*, \hat{\Sigma} \mu^* \rangle$ is further lower bounded. Since μ^* is feasible in the constrained set of (112), we have $\|\gamma\|_2^2 - \gamma^\top \hat{\Sigma} \mu^* - \lambda_n \|\gamma\|_1^2 \leq 0$, and hence for any positive constant $c > 0$, we have

$$\begin{aligned} \langle \mu^*, \hat{\Sigma} \mu^* \rangle &\geq \langle \mu^*, \hat{\Sigma} \mu^* \rangle + c \left(\|\gamma\|_2^2 - \gamma^\top \hat{\Sigma} \mu^* - \lambda_n \|\gamma\|_1^2 \right) \\ &\geq \min_{\mu} \left(\langle \mu, \hat{\Sigma} \mu \rangle + c \left(\|\gamma\|_2^2 - \gamma^\top \hat{\Sigma} \mu - \lambda_n \|\gamma\|_1^2 \right) \right) = -\frac{c^2}{4} \langle \gamma, \hat{\Sigma} \gamma \rangle + c \left(\|\gamma\|_2^2 - \lambda_n \|\gamma\|_1^2 \right). \end{aligned} \quad (114)$$

Note that $\|\gamma\|_1^2 \lambda_n \leq s_{z1} \lambda_n \|\gamma\|_2^2 = C s_{z1} \sqrt{\log p/n} \|\gamma\|_2^2 \ll \|\gamma\|_2^2$, where the last inequality holds when $s_{z1} \sqrt{\log p/n} \rightarrow 0$. By (114), we have

$$\begin{aligned} \langle \mu^*, \widehat{\Sigma} \mu^* \rangle &\geq \max_{c>0} -\frac{c^2}{4} \langle \gamma, \widehat{\Sigma} \gamma \rangle + c (\|\gamma\|_2^2 - \lambda_n \|\gamma\|_1^2) \\ &= \frac{(\|\gamma\|_2^2 - \lambda_n \|\gamma\|_1^2)^2}{\langle \gamma, \widehat{\Sigma} \gamma \rangle} \geq \frac{\|\gamma\|_2^4 (1 - s_{z1} \lambda_n)^2}{\langle \gamma, \widehat{\Sigma} \gamma \rangle}. \end{aligned} \quad (115)$$

On the event G_3 , we establish (43) for $\left\| \sum_{j \in \mathcal{S}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 / n$. The same argument holds for $\left\| \sum_{j \in \mathcal{V}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 / n$. Note that

$$\frac{1}{M_2} \leq \boldsymbol{\Theta}_{11} + \beta^2 \boldsymbol{\Theta}_{22} - 2\beta \boldsymbol{\Theta}_{12} = \begin{pmatrix} 1 & -\beta \end{pmatrix} \boldsymbol{\Theta} \begin{pmatrix} 1 \\ -\beta \end{pmatrix} \leq M_2 (1 + \beta^2). \quad (116)$$

Combined with (43), we establish the first inequality of (44). Note that

$$\frac{1}{n} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2^2 \leq \left(2M_2 \sum_{j \in \mathcal{V}^*} |\gamma_j| \right)^2 \leq s_{z1} \|\gamma\|_2^2. \quad (117)$$

Combined with (116), we establish the second inequality of (44). By the similar argument, we can establish (45).

E.6. Proof of Lemma 5

In the following proof, we will use the shorthand $\langle a, b \rangle_J = \sum_{j \in J} a_j b_j$. We have the following decompositions for $\|\widehat{\gamma}\|_2^2 - \|\gamma\|_2^2$ and $\widehat{\gamma}^\top \widehat{\Gamma} - \gamma^\top \Gamma$,

$$\begin{aligned} \|\widehat{\gamma}\|_2^2 - \|\gamma\|_2^2 &= 2\langle \gamma, D^\gamma \rangle_{\widehat{\mathcal{S}}} + 2\langle \gamma, \Delta^\gamma \rangle_{\widehat{\mathcal{S}}} + \langle D^\gamma, D^\gamma \rangle_{\widehat{\mathcal{S}}} + \langle \Delta^\gamma, \Delta^\gamma \rangle_{\widehat{\mathcal{S}}} + 2\langle D^\gamma, \Delta^\gamma \rangle_{\widehat{\mathcal{S}}} \\ &\quad - \left(\sum_{j \in \mathcal{S}^* \setminus \widehat{\mathcal{S}}} \gamma_j^2 - \sum_{j \in \widehat{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j^2 \right), \end{aligned} \quad (118)$$

and

$$\begin{aligned} \widehat{\gamma}^\top \widehat{\Gamma} - \gamma^\top \Gamma &= \langle \gamma, D^\Gamma \rangle_{\widehat{\mathcal{S}}} + \langle \Gamma, D^\gamma \rangle_{\widehat{\mathcal{S}}} + \langle \gamma, \Delta^\Gamma \rangle_{\widehat{\mathcal{S}}} + \langle \Gamma, \Delta^\gamma \rangle_{\widehat{\mathcal{S}}} + \langle D^\gamma, D^\Gamma \rangle_{\widehat{\mathcal{S}}} + \langle \Delta^\gamma, \Delta^\Gamma \rangle_{\widehat{\mathcal{S}}} \\ &\quad + \langle D^\gamma, \Delta^\Gamma \rangle_{\widehat{\mathcal{S}}} + \langle D^\Gamma, \Delta^\gamma \rangle_{\widehat{\mathcal{S}}} - \left(\sum_{j \in \mathcal{S}^* \setminus \widehat{\mathcal{S}}} \gamma_j \Gamma_j - \sum_{j \in \widehat{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j \Gamma_j \right). \end{aligned} \quad (119)$$

Recall that $\widehat{\mathbf{v}}^{[j]} = \mathbf{W}^\top \widehat{\mathbf{u}}^{[j]}$, then we have the following expression

$$\langle \gamma, D^\gamma \rangle_{\widehat{\mathcal{S}}} = \frac{1}{n} \sum_{j \in \widehat{\mathcal{S}}} \gamma_j (\widehat{\mathbf{v}}^{[j]})^\top \Pi_{\cdot 2}, \quad (120)$$

and

$$\langle \gamma, D^{\mathbf{r}} \rangle_{\tilde{\mathcal{S}}} + \langle \mathbf{\Gamma}, D^{\gamma} \rangle_{\tilde{\mathcal{S}}} = \frac{1}{n} \sum_{j \in \tilde{\mathcal{S}}} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}). \quad (121)$$

Note that $\tilde{\mathcal{S}}$ is correlated with the error $\mathbf{\Pi}_{.1}$ and $\mathbf{\Pi}_{.2}$. However, we can compare $\tilde{\mathcal{S}}$ with the true support \mathcal{S}^* ,

$$\begin{aligned} \langle \gamma, D^{\gamma} \rangle_{\tilde{\mathcal{S}}} - \langle \gamma, D^{\gamma} \rangle_{\mathcal{S}^*} &= \frac{1}{n} \sum_{j \in \tilde{\mathcal{S}}} \gamma_j (\hat{\mathbf{v}}^{[j]})^{\top} \mathbf{\Pi}_{.2} - \frac{1}{n} \sum_{j \in \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^{\top} \mathbf{\Pi}_{.2} \\ &= \frac{1}{n} \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^{\top} \mathbf{\Pi}_{.2} - \frac{1}{n} \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} \gamma_j (\hat{\mathbf{v}}^{[j]})^{\top} \mathbf{\Pi}_{.2}, \end{aligned} \quad (122)$$

and

$$\begin{aligned} &(\langle \gamma, D^{\mathbf{r}} \rangle_{\tilde{\mathcal{S}}} + \langle \mathbf{\Gamma}, D^{\gamma} \rangle_{\tilde{\mathcal{S}}}) - (\langle \gamma, D^{\mathbf{r}} \rangle_{\mathcal{S}^*} + \langle \mathbf{\Gamma}, D^{\gamma} \rangle_{\mathcal{S}^*}) \\ &= \frac{1}{n} \sum_{j \in \tilde{\mathcal{S}}} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}) - \frac{1}{n} \sum_{j \in \mathcal{S}^*} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}) \\ &= \frac{1}{n} \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}) - \frac{1}{n} \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}). \end{aligned} \quad (123)$$

Hence, the residual terms are

$$\begin{aligned} R^{\gamma} &= \sqrt{n} (2\langle \gamma, \Delta^{\gamma} \rangle_{\tilde{\mathcal{S}}} + \langle D^{\gamma}, D^{\gamma} \rangle_{\tilde{\mathcal{S}}} + \langle \Delta^{\gamma}, \Delta^{\gamma} \rangle_{\tilde{\mathcal{S}}} + 2\langle D^{\gamma}, \Delta^{\gamma} \rangle_{\tilde{\mathcal{S}}}) \\ &+ \sqrt{n} \left(\frac{1}{n} \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^{\top} \mathbf{\Pi}_{.2} - \frac{1}{n} \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} \gamma_j (\hat{\mathbf{v}}^{[j]})^{\top} \mathbf{\Pi}_{.2} \right) - \sqrt{n} \left(\sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} \gamma_j^2 - \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j^2 \right), \end{aligned} \quad (124)$$

and

$$\begin{aligned} R^{\text{inter}} &= \sqrt{n} (\langle \gamma, \Delta^{\mathbf{r}} \rangle_{\tilde{\mathcal{S}}} + \langle \mathbf{\Gamma}, \Delta^{\gamma} \rangle_{\tilde{\mathcal{S}}} + \langle D^{\gamma}, D^{\mathbf{r}} \rangle_{\tilde{\mathcal{S}}} + \langle \Delta^{\gamma}, \Delta^{\mathbf{r}} \rangle_{\tilde{\mathcal{S}}} + \langle D^{\gamma}, \Delta^{\mathbf{r}} \rangle_{\tilde{\mathcal{S}}} + \langle D^{\mathbf{r}}, \Delta^{\gamma} \rangle_{\tilde{\mathcal{S}}}) \\ &+ \sqrt{n} \left(\frac{1}{n} \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}) - \frac{1}{n} \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} (\hat{\mathbf{v}}^{[j]})^{\top} (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}) \right) \\ &- \sqrt{n} \left(\sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} \gamma_j \mathbf{\Gamma}_j - \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j \mathbf{\Gamma}_j \right). \end{aligned} \quad (125)$$

Define $\mathcal{S}_0^* = \left\{ j : |\gamma_j| > \sqrt{2.05 \log p_z} \sqrt{\text{Var}(D_j^{\gamma})} \right\}$ to be the set of strong signals, on the event A_2 , we have

$$\mathcal{S}_0^* \subset \tilde{\mathcal{S}} \subset \mathcal{S}^*, \quad \text{and} \quad |\tilde{\mathcal{S}}| \leq s_{z1}. \quad (126)$$

On the event A_3 , we have

$$\max \{ \|D^\Gamma\|_\infty, \|D^\gamma\|_\infty \} \leq \left(1 + 12\sqrt{\frac{\log p}{n}}\right) M_1 \sqrt{\frac{2.05 \log p_z}{n}} \max\{\sigma_1, \sigma_2\}. \quad (127)$$

On the event $S \cap G \cap A$,

$$\max \{ \|\Delta^\gamma\|_\infty, \|\Delta^\Gamma\|_\infty \} \leq C s \frac{\log p}{n} \max\{\sigma_1, \sigma_2\}. \quad (128)$$

Combing (126), (127) and (128), we have on the event $S \cap G \cap A$,

$$\max \{ \langle D^\gamma, D^\gamma \rangle_{\tilde{\mathcal{S}}}, \langle D^\Gamma, D^\Gamma \rangle_{\tilde{\mathcal{S}}} \} \leq C s_{z1} \frac{\log p_z}{n}, \quad (129)$$

$$\max \{ \langle \Delta^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}, \langle \Delta^\Gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{S}}} \} \leq C s_{z1} \left(s \frac{\log p}{n} \right)^2. \quad (130)$$

Note that

$$|\langle D^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}| \leq \sqrt{\langle D^\gamma, D^\gamma \rangle_{\tilde{\mathcal{S}}} \langle \Delta^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}} \leq \frac{1}{2} (\langle D^\gamma, D^\gamma \rangle_{\tilde{\mathcal{S}}} + \langle \Delta^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}).$$

Hence, we have

$$|\langle D^\gamma, D^\gamma \rangle_{\tilde{\mathcal{S}}} + \langle \Delta^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}} + 2\langle D^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}| \leq C s_{z1} \frac{\log p_z}{n} + C s_{z1} \left(s \frac{\log p}{n} \right)^2, \quad (131)$$

and

$$|\langle D^\gamma, D^\Gamma \rangle_{\tilde{\mathcal{S}}} + \langle \Delta^\gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{S}}} + \langle D^\gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{S}}} + \langle D^\Gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}| \leq C s_{z1} \frac{\log p_z}{n} + C s_{z1} \left(s \frac{\log p}{n} \right)^2. \quad (132)$$

We also have the following control

$$2|\langle \gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}| \leq \|\gamma\|_2 \sqrt{\langle \Delta^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}} \leq C \|\gamma\|_2 \sqrt{s_{z1} s} \frac{\log p}{n}, \quad (133)$$

and

$$|\langle \gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{S}}} + \langle \Gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{S}}}| \leq C (\|\gamma\|_2 + \|\Gamma\|_2) \sqrt{s_{z1} s} \frac{\log p}{n}. \quad (134)$$

On the event $S \cap G \cap A$, we have $\tilde{\mathcal{S}} \setminus \mathcal{S}^* = \emptyset$ and hence $\frac{1}{n} \sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2} = 0$, $\sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j \Gamma_j = 0$ and $\sum_{j \in \tilde{\mathcal{S}} \setminus \mathcal{S}^*} \gamma_j^2 = 0$; On the event $S \cap G \cap A$, we also have

$$\begin{aligned} & \left| \frac{1}{n} \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2} \right| \leq \frac{1}{n} s_{z1} \max_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} |\gamma_j| |(\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2}| \\ & \leq s_{z1} \sqrt{2.05 \log p_z} \sqrt{\text{Var}(D_j^\gamma)} \left(1 + 12\sqrt{\frac{\log p}{n}}\right) M_1 \sqrt{\frac{2.05 \log p_z}{n}} \sigma_2 \leq \frac{s_{z1} \log p_z}{n}. \end{aligned} \quad (135)$$

On the event $S \cap G \cap A$, we get

$$\left| \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} \gamma_j^2 \right| \leq s_{z1} \frac{\log p_z}{n}. \quad (136)$$

By (124), (131), (133), (135) and (136), we establish that on the event $S \cap G \cap A$,

$$|R^\gamma| \leq C s_{z1} \frac{\log p_z}{\sqrt{n}} + C \|\gamma\|_2 \sqrt{s_{z1} s} \frac{\log p}{\sqrt{n}}. \quad (137)$$

Similarly, we can establish that and

$$|R^{\text{inter}}| \leq C s_{z1} \frac{\log p_z}{\sqrt{n}} + C (\|\gamma\|_2 + \|\Gamma\|_2) \sqrt{s_{z1} s} \frac{\log p}{\sqrt{n}}. \quad (138)$$

We can establish (47) and (50) by taking $\Gamma_j = \beta \gamma_j$. Note that

$$\begin{aligned} \frac{2\sqrt{\log p}}{\sqrt{n}} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2 \sqrt{\Theta_{22}} &= 2\sqrt{\frac{\log p}{n}} \sqrt{\frac{\Theta_{22}}{\Theta_{11} + \beta^2 \Theta_{22} - 2\beta \Theta_{12}}} \sqrt{V_H} \|\gamma\|_2^2, \\ \frac{\sqrt{\log p}}{\sqrt{n}} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2 \sqrt{\Theta_{11} + \beta^2 \Theta_{22} + 2\beta \Theta_{12}} &= \sqrt{\frac{\log p}{n}} \sqrt{\frac{\Theta_{11} + \beta^2 \Theta_{22} + 2\beta \Theta_{12}}{\Theta_{11} + \beta^2 \Theta_{22} - 2\beta \Theta_{12}}} \sqrt{V_H} \|\gamma\|_2^2. \end{aligned} \quad (139)$$

By the definition of A_4 and A_5 in (37) and Lemma 4, we establish

$$\max \left\{ \left| \frac{2}{n} \sum_{j \in \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2} \right|, \left| \frac{1}{n} \sum_{j \in \mathcal{S}^*} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top (\mathbf{\Pi}_{.1} + \beta \mathbf{\Pi}_{.2}) \right| \right\} \leq C s_{z1} \sqrt{\frac{\log p}{n}} \|\gamma\|_2 \quad (140)$$

Combined with (50), we establish (51).

E.7. Proof of Lemma 8

The proof of Lemma 8 is similar to that of Lemma 5 and we will present it here. Similar to (118) and (119), we can obtain the following expressions,

$$\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 = \frac{1}{n} \sum_{j \in \tilde{\mathcal{V}}} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2} + \langle D^\gamma, D^\gamma \rangle_{\tilde{\mathcal{V}}} + \langle \Delta^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{V}}} + 2\langle D^\gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{V}}}, \quad (141)$$

and

$$\begin{aligned} \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\Gamma}_j - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j \Gamma_j &= \frac{1}{n} \sum_{j \in \tilde{\mathcal{V}}} (\hat{\mathbf{v}}^{[j]})^\top (\gamma_j \mathbf{\Pi}_{.1} + \Gamma_j \mathbf{\Pi}_{.2}) + \langle \gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{V}}} + \langle \Gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{V}}}, \\ &\quad + \langle D^\gamma, D^\Gamma \rangle_{\tilde{\mathcal{V}}} + \langle \Delta^\gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{V}}} + \langle D^\gamma, \Delta^\Gamma \rangle_{\tilde{\mathcal{V}}} + \langle D^\Gamma, \Delta^\gamma \rangle_{\tilde{\mathcal{V}}}, \end{aligned} \quad (142)$$

On the event $A \cap S \cap G$, we have

$$\frac{1}{n} \sum_{j \in \tilde{\mathcal{V}}} \gamma_j (\hat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{.2} \leq \sqrt{|\tilde{\mathcal{V}}|} \|\gamma\|_2 \max_{j \in \tilde{\mathcal{V}}} \frac{\|\hat{\mathbf{v}}^{[j]}\|_2}{n} \sqrt{2.05 \log p_z} \leq C \|\gamma\|_2 \sqrt{\frac{|\tilde{\mathcal{V}}| \log p_z}{n}}, \quad (143)$$

where the last inequality follows from (42). Combined with the fact $\mathbf{\Gamma}_j = \gamma_j(\beta + \pi_j/\gamma_j)$ and $|\pi_j/\gamma_j| \leq C\sqrt{\log p_z/n}$, we also establish that

$$\frac{1}{n} \sum_{j \in \tilde{\mathcal{V}}} (\hat{\mathbf{v}}^{[j]})^\top (\gamma_j \mathbf{\Pi}_{.1} + \mathbf{\Gamma}_j \mathbf{\Pi}_{.2}) \leq \sqrt{|\tilde{\mathcal{V}}|} \|\gamma\|_2 \max_{j \in \tilde{\mathcal{V}}} \frac{\|\hat{\mathbf{v}}^{[j]}\|_2}{n} \sqrt{2.05 \log p_z} \leq C \|\gamma\|_2 \sqrt{\frac{|\tilde{\mathcal{V}}| \log p_z}{n}}. \quad (144)$$

By (143), (144), (131) and (132), we obtain the following inequalities

$$\begin{aligned} & \max \left\{ \left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\mathbf{\Gamma}}_j - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j \mathbf{\Gamma}_j \right|, \left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \right| \right\} \\ & \leq C s_{z1} \frac{\log p_z}{n} + C s_{z1} \left(s \frac{\log p}{n} \right)^2 + C \|\gamma_{\tilde{\mathcal{V}}}\|_2 \sqrt{\frac{2|\tilde{\mathcal{V}}| \log p_z}{n}}. \end{aligned} \quad (145)$$

Since $|\tilde{\mathcal{V}}| > |\mathcal{V}^*|/2$ and $\min_{j \in \mathcal{S}^*} |\gamma_j| \geq \delta_{\min} \gg \sqrt{\log p_z/n}$, we have $\left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \right| \ll \sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2$. We have the following decomposition,

$$\begin{aligned} & \left| \frac{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\mathbf{\Gamma}}_j}{\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2} - \frac{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j \mathbf{\Gamma}_j}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} \right| \\ & \leq \frac{\left(\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \right) \left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\mathbf{\Gamma}}_j - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j \mathbf{\Gamma}_j \right| + \left| \sum_{j \in \tilde{\mathcal{V}}} \gamma_j \mathbf{\Gamma}_j \right| \left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \right|}{\left(\sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 \right) \left(\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \right)} \\ & \leq C \frac{\max \left\{ \left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \tilde{\mathbf{\Gamma}}_j - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j \mathbf{\Gamma}_j \right|, \left| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2 - \sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2 \right| \right\}}{\sum_{j \in \tilde{\mathcal{V}}} \gamma_j^2} \leq C \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}}, \end{aligned} \quad (146)$$

where the last inequality follows from (145) and the facts that $|\tilde{\mathcal{V}}| > |\mathcal{V}^*|/2$, $\min_{j \in \mathcal{S}^*} |\gamma_j| \geq \delta_{\min} \gg \sqrt{\log p_z/n}$ and $s \log p / \sqrt{n} \rightarrow 0$.

E.8. Proof of Lemma 7

Since $\min_{j \in \mathcal{S}^*} |\gamma_j| \geq \delta_{\min} \gg \sqrt{\log p/n}$, on the event $A \cap S \cap G$, we have

$$\tilde{\mathcal{S}} = \mathcal{S}^*.$$

Without loss of generality, we assume $1 \in \tilde{\mathcal{S}}$ and focus on the case $i = 1$. In the following, we are going to analyze the performance of $\hat{\beta}^{[1]}$ and $\tilde{\pi}_j^{[1]}$. In the following,

we first analyze $\widehat{\beta}^{[1]}$,

$$\sqrt{n} \left(\widehat{\beta}^{[1]} - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \right) = T^{\beta,1} + \Delta^{\beta,1}, \quad (147)$$

where

$$T^{\beta,1} = \frac{1}{\sqrt{n}\gamma_1} (\widehat{\mathbf{v}}^{[1]})^\top \left(\mathbf{\Pi}_{\cdot 1} - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \mathbf{\Pi}_{\cdot 2} \right) \quad \text{and} \quad \Delta^{\beta,1} = R_1 + R_2, \quad (148)$$

with

$$R_1 = \frac{\sqrt{n}}{\gamma_1} \left(\Delta_1^\Gamma - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \Delta_1^\gamma \right), \quad \text{and} \quad R_2 = \frac{-(D_1^\gamma + \Delta_1^\gamma)}{\gamma_1 + (D_1^\gamma + \Delta_1^\gamma)} (T^{\beta,1} + R_1). \quad (149)$$

To analyze $\widehat{\pi}^{[1]}$, we first analyze the following estimator,

$$\widehat{\pi}^{[1]} = \widetilde{\Gamma} - \widehat{\beta}^{[1]} \widetilde{\gamma}. \quad (150)$$

Note that

$$\begin{aligned} \widehat{\pi}_j^{[1]} - \pi_j &= -\frac{\pi_1}{\gamma_1} \gamma_j + \left(\widetilde{\Gamma}_j - \Gamma_j \right) - \left(\beta + \frac{\pi_1}{\gamma_1} \right) (\widetilde{\gamma}_j - \gamma_j) - \gamma_j \left(\widehat{\beta}^{[1]} - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \right) \\ &\quad - \left(\widehat{\beta}^{[1]} - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \right) (\widetilde{\gamma}_j - \gamma_j). \end{aligned} \quad (151)$$

By (101) and (102), we have

$$\sqrt{n} \left(\widetilde{\Gamma}_j - \Gamma_j \right) = \frac{1}{\sqrt{n}} (\widehat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{\cdot 1} + \sqrt{n} \Delta_j^\Gamma. \quad (152)$$

By (103) and (104), we have

$$\sqrt{n} (\widetilde{\gamma}_j - \gamma_j) = \frac{1}{\sqrt{n}} (\widehat{\mathbf{v}}^{[j]})^\top \mathbf{\Pi}_{\cdot 2} + \sqrt{n} \Delta_j^\gamma. \quad (153)$$

By plugging (148), (152) and (153) into (151), we have the following decomposition of $\widehat{\pi}_j^{[1]} - \pi_j$

$$\sqrt{n} \left(\widehat{\pi}_j^{[1]} - \pi_j \right) = -\sqrt{n} \frac{\pi_1}{\gamma_1} \gamma_j + T^{\pi_j} + \Delta^{\pi_j}, \quad (154)$$

where

$$T^{\pi_j} = \frac{1}{\sqrt{n}} \left((\widehat{\mathbf{v}}^{[j]})^\top - \frac{\gamma_j}{\gamma_1} (\widehat{\mathbf{v}}^{[1]})^\top \right) \left(\mathbf{\Pi}_{\cdot 1} - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \mathbf{\Pi}_{\cdot 2} \right),$$

and

$$\Delta^{\pi_j} = \sqrt{n} \left(\Delta_j^\Gamma - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \Delta_j^\gamma - \gamma_j \Delta^{\beta,1} \right) - \left(T^{\beta,1} + \Delta^{\beta,1} \right) (\widetilde{\gamma}_j - \gamma_j). \quad (155)$$

Define the events for $i \in \mathcal{S}^*$,

$$F^i = \left\{ \max_{j \in \mathcal{S}^*, j \neq i} |T^{\pi_j}| \leq 2.02 \sqrt{\log p_z} \sqrt{\boldsymbol{\Theta}_{11} + \left(\beta + \frac{\pi_i}{\gamma_i} \right)^2 \boldsymbol{\Theta}_{22} - 2 \left(\beta + \frac{\pi_i}{\gamma_i} \right) \boldsymbol{\Theta}_{12}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_i}{\gamma_i} \hat{\mathbf{v}}^{[i]} \right\|_2}{\sqrt{n}} \right\}$$

Then for $F = \cap_{i \in \mathcal{S}^*} F^i$, we have

$$\mathbf{P}(F) \geq 1 - C s_{z1}^2 p_z^{-2.04} \geq 1 - c p^{-c}. \quad (156)$$

The proof of Lemma 7 relies on the following lemmas. The following lemma provides upper bound and lower bound for the variance term and the proof of the following lemma can be found in Section E.11.

LEMMA 12. *On the event $A \cap S \cap G$, we have*

$$\begin{aligned} & \sqrt{\boldsymbol{\Theta}_{11} + \left(\beta + \left| \frac{\pi_1}{\gamma_1} \right| \right)^2 \boldsymbol{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \boldsymbol{\Theta}_{12}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_j}{\gamma_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}} \\ & \leq 1.1 \sqrt{M_1 M_2} \left(1 + \left| \frac{\gamma_j}{\gamma_1} \right| \right) \sqrt{1 + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2}, \end{aligned} \quad (157)$$

and

$$\begin{aligned} & \sqrt{\boldsymbol{\Theta}_{11} + \left(\beta + \left| \frac{\pi_1}{\gamma_1} \right| \right)^2 \boldsymbol{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \boldsymbol{\Theta}_{12}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_j}{\gamma_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}} \\ & \geq 0.45 \sqrt{\frac{M_1}{M_2}} \left(1 + \left| \frac{\gamma_j}{\gamma_1} \right| \right) \sqrt{1 + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2}. \end{aligned} \quad (158)$$

LEMMA 13. *On the event $A \cap S \cap G \cap F^1$, for large n , we have*

$$0.995 \leq \frac{\sqrt{\hat{\boldsymbol{\Theta}}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\boldsymbol{\Theta}}_{22} - 2 \hat{\beta}^{[1]} \hat{\boldsymbol{\Theta}}_{12}} \left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{\boldsymbol{\Theta}_{11} + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2 \boldsymbol{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \boldsymbol{\Theta}_{12}} \left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_j}{\gamma_1} \hat{\mathbf{v}}^{[1]} \right\|_2} \leq 1.005. \quad (159)$$

On the event $A \cap S \cap G \cap F^1$, we have

$$\max_{j \in \mathcal{S}^*} \frac{1}{\sqrt{n}} |T^{\pi_j}| \leq 2.02 \sqrt{\frac{\log p_z}{n}} \sqrt{\boldsymbol{\Theta}_{11} + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2 \boldsymbol{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \boldsymbol{\Theta}_{12}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_j}{\gamma_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}; \quad (160)$$

and

$$\max_{j \in \mathcal{S}^*} \frac{1}{\sqrt{n}} |\Delta^{\pi_j}| \leq \frac{1}{300} \sqrt{\frac{\log p_z}{n}} \sqrt{\boldsymbol{\Theta}_{11} + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2 \boldsymbol{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \boldsymbol{\Theta}_{12}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_j}{\gamma_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}. \quad (161)$$

The ratio $\frac{\pi_1}{\gamma_1}$ can be divided into the following three cases,

- (a) strongly invalid instrument case, $|\pi_1/\gamma_1| \geq C_*(1/\delta_{\min})\sqrt{\log p_z/n}$, where $C_* = 12(1 + |\beta|)\sqrt{M_1/M_2}$;
- (b) weakly invalid instrument case, $|\pi_1/\gamma_1| < C_*(1/\delta_{\min})\sqrt{\log p_z/n}$;
- (c) valid instrument case, $\pi_1/\gamma_1 = 0$.

We are going to show that our procedure (12) in the main paper will rule out the strong invalid instrument case and a stronger assumption (22) in the main paper will help use rule out the weakly invalid instrument case. In the following, we will analyze the three cases separately.

strongly invalid instrument case

In this case, we assume that $|\pi_1/\gamma_1| \geq C_*(1/\delta_{\min})\sqrt{\log p_z/n}$. For $j \in \mathcal{V}^*$, (154) can be re-expressed as

$$\sqrt{n} \left(\hat{\pi}_j^{[1]} - 0 \right) = -\sqrt{n} \frac{\pi_1}{\gamma_1} \gamma_j + T^{\pi_j} + \Delta^{\pi_j}. \quad (162)$$

We are going to show that on the event $A \cap S \cap G \cap F^1$,

$$\|\tilde{\pi}^{[1]}\|_0 > \frac{|\mathcal{S}^*|}{2}. \quad (163)$$

It is sufficient to show for $j \in \mathcal{V}^*$

$$\left| -\frac{\pi_1}{\gamma_1} \gamma_j + \frac{1}{\sqrt{n}} (T^{\pi_j} + \Delta^{\pi_j}) \right| \geq 2.05 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}, \quad (164)$$

which can be reduced to

$$\max_{j \in \mathcal{V}^*} \frac{1}{\sqrt{n}} |T^{\pi_j} + \Delta^{\pi_j}| \leq 2.05 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}, \quad (165)$$

and

$$\left| \frac{\pi_1}{\gamma_1} \gamma_j \right| \geq 4.1 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}. \quad (166)$$

By (159), (160) and (161), we establish (165). By (157) and (159), we have

$$\begin{aligned} & 4.1 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}} \\ & \leq 4.1 \times 1.005 \times 1.1 \sqrt{M_1 M_2} \left(1 + \left| \frac{\gamma_j}{\gamma_1} \right| \right) \sqrt{1 + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2} \sqrt{\frac{\log p_z}{n}} \\ & \leq |\gamma_j| 4.1 \times 1.005 \times 1.1 \sqrt{M_1 M_2} \left(\left| \frac{1}{\gamma_j} \right| + \left| \frac{1}{\gamma_1} \right| \right) \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right) \sqrt{\frac{\log p_z}{n}}. \end{aligned} \quad (167)$$

The last term can be further upper bounded by

$$\begin{aligned}
 & |\gamma_j| \frac{1}{\delta_{\min}} 8.2 \times 1.005 \times 1.1 \sqrt{M_1 M_2} \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right) \sqrt{\frac{\log p_z}{n}} \\
 & \leq |\gamma_j| \frac{1}{\delta_{\min}} 8.2 \times 1.005 \times 1.1 \sqrt{M_1 M_2} (1 + |\beta|) \sqrt{\frac{\log p_z}{n}} \\
 & \quad + |\gamma_j| \frac{1}{\delta_{\min}} 8.2 \times 1.005 \times 1.1 \sqrt{M_1 M_2} \left| \frac{\pi_1}{\gamma_1} \right| \sqrt{\frac{\log p_z}{n}} \\
 & \leq 0.99 \frac{C_*}{\delta_{\min}} |\gamma_j| \sqrt{\frac{\log p_z}{n}} + C \frac{\sqrt{\frac{\log p_z}{n}}}{\delta_{\min}} \left| \frac{\pi_1}{\gamma_1} \gamma_j \right|,
 \end{aligned} \tag{168}$$

where the first inequality follows from triangle inequality and the second inequality follows from the definition of C_* . Since $\delta_{\min} \gg \sqrt{\log p/n}$ and $|\pi_1/\gamma_1| \geq C_*(1/\delta_{\min})\sqrt{\log p_z/n}$, by (167) and (168), we conclude (166).

weakly invalid instrument case

In this case, we assume $0 < |\pi_1/\gamma_1| < C_*(1/\delta_{\min})\sqrt{\log p_z/n}$. We have the following expression of (154),

$$\hat{\pi}_j^{[1]} = \pi_j - \frac{\pi_1}{\gamma_1} \gamma_j + \frac{1}{\sqrt{n}} (T^{\pi_j} + \Delta^{\pi_j}). \tag{169}$$

We are going to show that on the event $A \cap S \cap G \cap F^1$,

$$\left\{ j \in \mathcal{S}^* : \left| \frac{\pi_j}{\gamma_j} \right| \geq 2C_* \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}} \right\} \subset \text{supp}(\tilde{\pi}^{[1]}). \tag{170}$$

It is sufficient to show the following inequality if $|\pi_j/\gamma_j| > 2C_*(1/\delta_{\min})\sqrt{\log p_z/n}$,

$$\left| \pi_j - \frac{\pi_1}{\gamma_1} \gamma_j + \frac{1}{\sqrt{n}} (T^{\pi_j} + \Delta^{\pi_j}) \right| \geq 2.05 \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}, \tag{171}$$

which can be reduced to

$$\max_{j \in \text{supp}(\pi) \cap \mathcal{S}^*} \frac{1}{\sqrt{n}} |T^{\pi_j} + \Delta^{\pi_j}| \leq 2.05 \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}, \tag{172}$$

$$\left| \pi_j - \frac{\pi_1}{\gamma_1} \gamma_j \right| \geq 4.1 \sqrt{\hat{\Theta}_{11} + (\hat{\beta}^{[1]})^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}. \tag{173}$$

By (159), (160) and (161), we establish (172). Since

$$\frac{1}{|\gamma_j|} \left| \pi_j - \frac{\pi_1}{\gamma_1} \gamma_j \right| \geq \left| \frac{\pi_j}{\gamma_j} \right| - \left| \frac{\pi_1}{\gamma_1} \right| \geq C_* \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}},$$

the assumption $0 < |\pi_1/\gamma_1| < C_*(1/\delta_{\min})\sqrt{\log p_z/n}$ and (168) lead to (173).

valid instrument case

In this case, the instrumental variable is valid with $\pi_1/\gamma_1 = 0$. For $j \in \mathcal{V}^*$, (154) can be re-expressed as

$$\sqrt{n} \left(\hat{\pi}_j^{[1]} - 0 \right) = T^{\pi_j} + \Delta^{\pi_j}. \quad (174)$$

By (159), (160) and (161), on the event $A \cap S \cap G \cap F^1$,

$$\max_{j \in \mathcal{V}^*} \frac{1}{\sqrt{n}} |T^{\pi_j} + \Delta^{\pi_j}| \leq 2.05 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}. \quad (175)$$

and hence

$$\text{supp} \left(\tilde{\pi}^{[1]} \right) \subset \text{supp}(\pi) \quad \text{and} \quad \|\tilde{\pi}^{[1]}\|_0 < \frac{|\mathcal{S}^*|}{2}. \quad (176)$$

For $|\pi_j/\gamma_j| \geq C_*(1/\delta_{\min})\sqrt{\log p_z/n}$, by (168), we obtain

$$|\pi_j| \geq 4.1 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}. \quad (177)$$

Combined with (175), we have

$$\left| \hat{\pi}_j^{[1]} \right| \geq 2.05 \sqrt{\hat{\Theta}_{11} + \left(\hat{\beta}^{[1]} \right)^2 \hat{\Theta}_{22} - 2\hat{\beta}^{[1]} \hat{\Theta}_{12}} \sqrt{\frac{\log p_z}{n}} \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\tilde{\gamma}_j}{\tilde{\gamma}_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}}. \quad (178)$$

Hence

$$\left\{ j \in \mathcal{S}^* : \left| \frac{\pi_j}{\gamma_j} \right| \geq C_* \frac{1}{\delta_{\min}} \sqrt{\frac{\log p_z}{n}} \right\} \subset \text{supp} \left(\tilde{\pi}^{[1]} \right). \quad (179)$$

By comparing (163) and (176), we rule out the strong invalid instrumental variable case and obtain $|\tilde{\mathcal{V}}| > |\mathcal{S}^*|/2$ in (60). Further by (170) and (179), we establish (60).

With a stronger assumption (22) in the main paper, the weak invalid instrument case is also ruled out and (179) leads to (61).

E.9. Proof of Lemma 6

The proof of this lemma follows from the following results. Under the regularity assumptions (R1) – (R3), as $\sqrt{s_{z1}}s \log p/\sqrt{n} \rightarrow 0$, we have

$$\max_{1 \leq i, j \leq 2} |\hat{\Theta}_{ij} - \Theta_{ij}| \xrightarrow{p} 0; \quad (180)$$

and

$$\frac{\|\widehat{\gamma}\|_2^2}{\|\gamma\|_2^2} \xrightarrow{p} 1 \quad \text{and} \quad \frac{\left\| \sum_{j \in \tilde{\mathcal{S}}} \tilde{\gamma}_j \hat{\mathbf{v}}^{[j]} \right\|_2}{\left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2} \xrightarrow{p} 1. \quad (181)$$

By (27) and (180), we establish that

$$\frac{\sqrt{\widehat{\Theta}_{11} + \widehat{\beta}^2 \widehat{\Theta}_{22} - 2\widehat{\beta} \widehat{\Theta}_{12}}}{\sqrt{\Theta_{11} + (\beta)^2 \Theta_{22} - 2\beta \Theta_{12}}} \xrightarrow{p} 1.$$

Combined with (181), we establish (58).

Proof of (180) A stronger version of this proposition has already been proved in Ren et al. [2013], where part of it was already established in Sun and Zhang [2012]. To be self-contained, we will provide the sketch of the proof in the following.

The difference between $\widehat{\Theta} - \Theta$ can be decomposed as,

$$\widehat{\Theta} - \Theta = \Theta^{\text{ora}} - \Theta + \widehat{\Theta} - \Theta^{\text{ora}}, \quad (182)$$

where $\Theta_{11}^{\text{ora}} = \frac{1}{n} \|Y - \mathbf{Z}\Gamma - \mathbf{X}\Psi\|_2^2$, $\Theta_{22}^{\text{ora}} = \frac{1}{n} \|D - \mathbf{Z}\gamma - \mathbf{X}\psi\|_2^2$ and $\Theta_{12}^{\text{ora}} = \frac{1}{n} (Y - \mathbf{Z}\Gamma - \mathbf{X}\Psi)^\top (D - \mathbf{Z}\gamma - \mathbf{X}\psi)$. In the following, we only provide the detailed analysis of $\widehat{\Theta}_{12} - \Theta_{12}^{\text{ora}}$. The other differences can be established in a similar way and the difference between $\Theta^{\text{ora}} - \Theta$ can be established by central limit theorem.

$$\widehat{\Theta}_{12} - \Theta_{12}^{\text{ora}} = \frac{1}{n} \begin{pmatrix} \widehat{\gamma} - \gamma \\ \widehat{\psi} - \psi \end{pmatrix}^\top \mathbf{W}^\top \mathbf{W} \begin{pmatrix} \widehat{\Gamma} - \Gamma \\ \widehat{\Psi} - \Psi \end{pmatrix} + \frac{1}{n} \mathbf{\Pi}_{\cdot 2}^\top \mathbf{W} \begin{pmatrix} \widehat{\Gamma} - \Gamma \\ \widehat{\Psi} - \Psi \end{pmatrix} + \frac{1}{n} \mathbf{\Pi}_{\cdot 1}^\top \mathbf{W} \begin{pmatrix} \widehat{\gamma} - \gamma \\ \widehat{\psi} - \psi \end{pmatrix}. \quad (183)$$

By (183), we have

$$\begin{aligned} \left| \widehat{\Theta}_{12} - \Theta_{12}^{\text{ora}} \right| &\leq \frac{1}{\sqrt{n}} \left\| \mathbf{W} \begin{pmatrix} \widehat{\gamma} - \gamma \\ \widehat{\psi} - \psi \end{pmatrix} \right\|_2 \frac{1}{\sqrt{n}} \left\| \mathbf{W} \begin{pmatrix} \widehat{\Gamma} - \Gamma \\ \widehat{\Psi} - \Psi \end{pmatrix} \right\|_2 + \frac{1}{n} \|\mathbf{\Pi}_{\cdot 2}^\top \mathbf{W}\|_\infty \left\| \begin{pmatrix} \widehat{\Gamma} - \Gamma \\ \widehat{\Psi} - \Psi \end{pmatrix} \right\|_1 \\ &\quad + \frac{1}{n} \|\mathbf{\Pi}_{\cdot 1}^\top \mathbf{W}\|_\infty \left\| \begin{pmatrix} \widehat{\gamma} - \gamma \\ \widehat{\psi} - \psi \end{pmatrix} \right\|_1. \end{aligned} \quad (184)$$

The following of the proof follows from Lemma 10 and definition of event G .

Proof of (181) For a given $0 < \epsilon_0 < 1$, we have

$$\mathbf{P} \left(\left| \frac{\widehat{\|\gamma\|_2^2}}{\|\gamma\|_2^2} - 1 \right| \geq \epsilon_0 \right) \leq \mathbf{P} \left(\left| \frac{\widehat{\|\gamma\|_2^2} - \|\gamma\|_2^2}{\|\gamma\|_2^2} \right| \geq \frac{\epsilon_0}{1 - \epsilon_0} \right).$$

By Lemma 5, on the event $A \cap S \cap G$, we have

$$\left| \frac{\widehat{\|\gamma\|_2^2} - \|\gamma\|_2^2}{\|\gamma\|_2^2} \right| \leq C \frac{1}{\|\gamma\|_2^2} \left(s_{z1} \frac{\log p_z}{n} + C \|\gamma\|_2 \sqrt{\frac{2s_{z1} \log p_z}{n}} \right).$$

Since $\|\gamma\|_2^2 \gg (s \log p / \sqrt{n})^2$, we obtain that

$$\left| \frac{\widehat{\|\gamma\|_2^2} - \|\gamma\|_2^2}{\|\gamma\|_2^2} \right| \leq \frac{\epsilon_0}{1 - \epsilon_0} \text{ and } \mathbf{P} \left(\left| \frac{\widehat{\|\gamma\|_2^2}}{\|\gamma\|_2^2} - 1 \right| \geq \epsilon_0 \right) \leq \mathbf{P}((A \cap S \cap G)^c).$$

Combined with Lemma 3, we establish the first convergence result of (181). On the event $S \cap G \cap A$, we have

$$\left| \frac{\left\| \sum_{j \in \tilde{\mathcal{S}}} \tilde{\gamma}_j \hat{\mathbf{v}}^{[j]} \right\|_2}{\left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2} - 1 \right| \leq \frac{\sum_{j \in \tilde{\mathcal{S}}} |\tilde{\gamma}_j - \gamma_j| \frac{\|\hat{\mathbf{v}}^{[j]}\|}{\sqrt{n}} + \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} |\gamma_j| \frac{\|\hat{\mathbf{v}}^{[j]}\|}{\sqrt{n}}}{\frac{1}{\sqrt{n}} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2}$$

By Lemma 4, we have

$$\begin{aligned} & \frac{\sum_{j \in \tilde{\mathcal{S}}} |\tilde{\gamma}_j - \gamma_j| \frac{\|\hat{\mathbf{v}}^{[j]}\|}{\sqrt{n}} + \sum_{j \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} |\gamma_j| \frac{\|\hat{\mathbf{v}}^{[j]}\|}{\sqrt{n}}}{\frac{1}{\sqrt{n}} \left\| \sum_{j \in \mathcal{S}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2} \\ & \leq \frac{\left(\sum_{j \in \tilde{\mathcal{S}}} |\tilde{\gamma}_j - \gamma_j| + \sum_{i \in \mathcal{S}^* \setminus \tilde{\mathcal{S}}} |\gamma_i| \right) \sqrt{\left(1 + 12\sqrt{\frac{\log p}{n}} \right) M_1}}{\sqrt{\frac{M_1 \|\gamma\|_2^2 (1 - s_{z1} \lambda_n)^2}{1 - 12\sqrt{\frac{\log p}{n}}}}} \leq C s_{z1} \sqrt{\frac{\log p}{n}} \leq \epsilon_0, \end{aligned}$$

and hence the second convergence result of (181) follows from the following inequality.

E.10. *Proof of Lemma 9*

Define $\|\gamma\|_2^2 = \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j^2$ and $\|\gamma_{\mathcal{V}^*}\|_2^2 = \sum_{j \in \mathcal{V}^*} \gamma_j^2$. The proof of this lemma is further based on the following results. Under the assumptions (R1)–(R5) and (IN1)–(IN3). As $\sqrt{s_{z1}} s \log p / \sqrt{n} \rightarrow 0$, we have

$$\frac{\widetilde{\|\gamma\|_2^2}}{\|\gamma_{\mathcal{V}^*}\|_2^2} \xrightarrow{p} 1 \quad \text{and} \quad \frac{\left\| \sum_{j \in \tilde{\mathcal{V}}} \tilde{\gamma}_j \hat{\mathbf{v}}^{[j]} \right\|_2}{\left\| \sum_{j \in \mathcal{V}^*} \gamma_j \hat{\mathbf{v}}^{[j]} \right\|_2} \xrightarrow{p} 1. \quad (185)$$

By (180) and (23) in the main paper, we establish that

$$\frac{\sqrt{\hat{\Theta}_{11} + \hat{\beta}^2 \hat{\Theta}_{22} - 2\hat{\beta} \hat{\Theta}_{12}}}{\sqrt{\Theta_{11} + (\beta)^2 \Theta_{22} - 2\beta \Theta_{12}}} \xrightarrow{p} 1.$$

Combined with (185), we establish (65).

Proof of (185) For a given $0 < \epsilon_0 < 1$, we have

$$\mathbf{P} \left(\left| \frac{\widetilde{\|\gamma\|_2^2}}{\|\gamma_{\mathcal{V}^*}\|_2^2} - 1 \right| \geq \epsilon_0 \right) \leq \mathbf{P} \left(\left| \frac{\widetilde{\|\gamma\|_2^2} - \|\gamma_{\mathcal{V}^*}\|_2^2}{\|\gamma_{\mathcal{V}^*}\|_2^2} \right| \geq \frac{\epsilon_0}{1 - \epsilon_0} \right).$$

By Lemma 7, on the event $A \cap S \cap G \cap F$, we have $\tilde{\mathcal{V}} = \mathcal{V}^*$ and (145) leads to

$$\left| \frac{\widetilde{\|\gamma\|_2^2} - \|\gamma_{\mathcal{V}^*}\|_2^2}{\|\gamma_{\mathcal{V}^*}\|_2^2} \right| \leq C \frac{1}{\|\gamma_{\mathcal{V}^*}\|_2^2} \left(s_{z1} \frac{\log p_z}{n} + C s_{z1} \left(s \frac{\log p}{n} \right)^2 + C \|\gamma_{\tilde{\mathcal{V}}}\|_2 \sqrt{\frac{2|\tilde{\mathcal{V}}| \log p_z}{n}} \right). \quad (186)$$

Since $\|\gamma_{\mathcal{V}^*}\|_2^2 \gg (s \log p / \sqrt{n})^2$, we obtain that

$$\left| \frac{\|\widehat{\gamma}\|_2^2 - \|\gamma_{\mathcal{V}^*}\|_2^2}{\|\gamma_{\mathcal{V}^*}\|_2^2} \right| \leq \frac{\epsilon_0}{1 - \epsilon_0} \text{ and } \mathbf{P} \left(\left| \frac{\|\widehat{\gamma}\|_2^2}{\|\gamma_{\mathcal{V}^*}\|_2^2} - 1 \right| \geq \epsilon_0 \right) \leq \mathbf{P}((A \cap S \cap G \cap F)^c).$$

Combined with Lemma 3 and (156), we establish the first converge result of (185).

On the event $S \cap G \cap A \cap F$, we have

$$\left| \frac{\left\| \sum_{j \in \mathcal{V}} \tilde{\gamma}_j \widehat{\mathbf{v}}^{[j]} \right\|_2}{\left\| \sum_{j \in \mathcal{V}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2} - 1 \right| \leq \frac{\sum_{j \in \mathcal{V}^*} |\tilde{\gamma}_j - \gamma_j| \frac{\|\widehat{\mathbf{v}}^{[j]}\|}{\sqrt{n}}}{\frac{1}{\sqrt{n}} \left\| \sum_{j \in \mathcal{V}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2}.$$

By Lemma 4, we have

$$\frac{\sum_{j \in \mathcal{V}^*} |\tilde{\gamma}_j - \gamma_j| \frac{\|\widehat{\mathbf{v}}^{[j]}\|}{\sqrt{n}}}{\frac{1}{\sqrt{n}} \left\| \sum_{j \in \mathcal{V}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2} \leq \frac{\left(\sum_{j \in \mathcal{V}^*} |\tilde{\gamma}_j - \gamma_j| \right) \sqrt{\left(1 + 12 \sqrt{\frac{\log p}{n}} \right) M_1}}{\sqrt{\frac{M_1 \|\gamma\|_2^2 (1 - s_{z1} \lambda_n)^2}{1 - 12 \sqrt{\frac{\log p}{n}}}}} \leq C s_{z1} \sqrt{\frac{\log p}{n}}. \quad (187)$$

Hence the second converge result of (185) follows from

$$\mathbf{P} \left(\left| \frac{\left\| \sum_{j \in \mathcal{V}} \tilde{\gamma}_j \widehat{\mathbf{v}}^{[j]} \right\|_2}{\left\| \sum_{j \in \mathcal{V}^*} \gamma_j \widehat{\mathbf{v}}^{[j]} \right\|_2} - 1 \right| \geq \epsilon_0 \right) \leq \mathbf{P}((S \cap G \cap A \cap F)^c).$$

E.11. Proof of Lemmas 10, 11, 12 and 13

Proof of Lemma 10 We only establish the first half of (98) and (99). The proof of the second half of (98) and (100) will be similar. The proof has been established in Sun and Zhang [2012] for fixed designs under certain assumptions for the design. In the following, we will check that the assumptions in Corollary 1 in Sun and Zhang [2012] are satisfied with high probability for the subgaussian random designs considered in this paper and then apply equation (23) in Sun and Zhang [2012]. By the definition of τ^* in Sun and Zhang [2012], we have $\tau^* \leq \tau$ where τ is defined in (30). Hence, on the event S_1 , equation (23) in Sun and Zhang [2012] holds. By the relationship between ℓ_1 cone invertibility factor and the restricted eigenvalue established in Lemma 13 of Cai and Guo [2016b], we obtain that on the event $S \cap G$,

$$\|\widehat{\mathbf{\Gamma}} - \mathbf{\Gamma}\|_1 + \|\widehat{\mathbf{\Psi}} - \mathbf{\Psi}\|_1 \leq C \frac{s \lambda_0 \sigma_1}{\kappa^2(\mathbf{H}, 4s, 1 + 2\epsilon_0)}. \quad (188)$$

Similar to the proof of Lemma 13 in Cai and Guo [2016b], we establish

$$\kappa^2(\mathbf{H}, 4s, 1 + 2\epsilon_0) \geq \frac{n}{\max \|\mathbf{W}_{\cdot j}\|_2^2} \kappa^2 \left(\mathbf{W}, 4s, (1 + 2\epsilon_0) \left(\frac{\max \|\mathbf{W}_{\cdot j}\|_2}{\min \|\mathbf{W}_{\cdot j}\|_2} \right) \right). \quad (189)$$

Hence, on the event $G \cap S$, we establish the first half of (98). Since

$$\frac{1}{n} \|\mathbf{Z}(\hat{\mathbf{\Gamma}} - \mathbf{\Gamma}) + \mathbf{X}(\hat{\mathbf{\Psi}} - \mathbf{\Psi})\|_2^2 \leq \left\| \frac{1}{n} \mathbf{W}^\top \mathbf{W} \begin{pmatrix} \hat{\mathbf{\Gamma}} - \mathbf{\Gamma} \\ \hat{\mathbf{\Psi}} - \mathbf{\Psi} \end{pmatrix} \right\|_\infty \left(\|\hat{\mathbf{\Gamma}} - \mathbf{\Gamma}\|_1 + \|\hat{\mathbf{\Psi}} - \mathbf{\Psi}\|_1 \right),$$

we establish (99).

Proof of Lemma 11 The decompositions (101) and (103) are established by the definitions of $\tilde{\mathbf{\Gamma}}$ and $\tilde{\gamma}$. The error bound (105) follows from the following inequality

$$|\Delta_j^\gamma| \leq \left\| \left(\frac{1}{n} (\hat{\mathbf{u}}^{[j]})^\top \hat{\mathbf{\Sigma}} - e_j^\top \right) \right\|_\infty \left\| \begin{pmatrix} \hat{\gamma} - \gamma \\ \hat{\psi} - \psi \end{pmatrix} \right\|_1.$$

Proof of Lemma 12

This lemma can be established by a similar argument with Lemma 4. On the event $A \cap S \cap G$, we have

$$\frac{1}{\sqrt{M_2}} \sqrt{1 + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2} \leq \sqrt{\mathbf{\Theta}_{11} + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2 \mathbf{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \mathbf{\Theta}_{12}} \leq \sqrt{M_2} \sqrt{1 + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2}, \quad (190)$$

and

$$\frac{\sqrt{M_1} \sqrt{1 + \left(\frac{\gamma_j}{\gamma_1} \right)^2} |1 - 2\lambda_n|}{\sqrt{1 - 12\sqrt{\frac{\log p}{n}}}} \leq \frac{\left\| \hat{\mathbf{v}}^{[j]} - \frac{\gamma_j}{\gamma_1} \hat{\mathbf{v}}^{[1]} \right\|_2}{\sqrt{n}} \leq \sqrt{M_1} \left(1 + \left| \frac{\gamma_j}{\gamma_1} \right| \right) \sqrt{1 + 12\sqrt{\frac{\log p}{n}}}. \quad (191)$$

Hence, (157) and (158) follow from the above inequalities (190) and (191).

Proof of Lemma 13

(159) follows from the standard convergence analysis and (160) follows from high probability statement of Gaussian random variable. It remains to establish (161). We will analyze the expression (155) term by term. Note that on the event $A \cap S \cap G$,

$$\begin{aligned} |T^{\beta,1}| &\leq \frac{\sqrt{\log p_z}}{\sqrt{n} |\gamma_1|} \|\hat{\mathbf{v}}^{[1]}\|_2 \sqrt{\mathbf{\Theta}_{11} + \left(\beta + \frac{\pi_1}{\gamma_1} \right)^2 \mathbf{\Theta}_{22} - 2 \left(\beta + \frac{\pi_1}{\gamma_1} \right) \mathbf{\Theta}_{12}} \\ &\leq C \frac{1}{|\gamma_1|} \frac{\|\hat{\mathbf{v}}^{[1]}\|_2}{\sqrt{n}} \sqrt{\log p_z} \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right); \\ |R_1| &\leq C \frac{1}{|\gamma_1|} \left(\left| \beta + \frac{\pi_1}{\gamma_1} \right| + 1 \right) s \frac{\log p}{\sqrt{n}}; \\ |R_2| &\leq C \frac{1}{|\gamma_1|} \left(\sqrt{\frac{\log p}{n}} + s \frac{\log p}{n} \right) \left(|T^{\beta,1}| + |R_1| \right) \\ &\leq C \frac{1}{|\gamma_1|} \sqrt{\frac{\log p}{n}} \frac{1}{|\gamma_1|} \frac{\|\hat{\mathbf{v}}^{[1]}\|_2}{\sqrt{n}} \sqrt{\log p_z} \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right). \end{aligned} \quad (192)$$

Hence

$$\begin{aligned}
& \frac{1}{\sqrt{n}} \left| \gamma_j \Delta^{\beta,1} + \left(T^{\beta,1} + \Delta^{\beta,1} \right) (\tilde{\gamma}_j - \gamma_j) \right| \leq C \frac{|\gamma_j|}{\sqrt{n}} (|R_1| + |R_2|) + \sqrt{\frac{\log p_z}{n}} |T^{\beta,1}| \\
& \leq C \frac{|\gamma_j|}{|\gamma_1|} \left(\left| \beta + \frac{\pi_1}{\gamma_1} \right| + 1 \right) s \frac{\log p}{n} + C \frac{|\gamma_j|}{|\gamma_1|} \frac{\sqrt{\log p \log p_z}}{n} \frac{1}{|\gamma_1|} \frac{\|\hat{\mathbf{v}}^{[1]}\|_2}{\sqrt{n}} \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right) \\
& \quad + C \frac{1}{|\gamma_1|} \frac{\|\hat{\mathbf{v}}^{[1]}\|_2}{\sqrt{n}} \frac{\sqrt{\log p \log p_z}}{n} \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right) \\
& \leq C \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right) \left(\frac{|\gamma_j|}{|\gamma_1|} s \frac{\log p}{n} + \left(1 + \frac{|\gamma_j|}{|\gamma_1|} \right) \frac{1}{|\gamma_1|} \frac{\sqrt{\log p \log p_z}}{n} \right). \tag{193}
\end{aligned}$$

Since

$$\left| \Delta_j^\Gamma - \left(\beta + \frac{\pi_1}{\gamma_1} \right) \Delta_j^\gamma \right| \leq C \left(\left| \beta + \frac{\pi_1}{\gamma_1} \right| + 1 \right) s \frac{\log p}{n},$$

we have

$$\max_{j \in \mathcal{S}^*} \frac{1}{\sqrt{n}} |\Delta^{\pi_j}| \leq \left(1 + \left| \beta + \frac{\pi_1}{\gamma_1} \right| \right) \left(1 + \frac{|\gamma_j|}{|\gamma_1|} \right) \left(s \frac{\log p}{n} + \frac{1}{|\gamma_1|} \frac{\sqrt{\log p \log p_z}}{n} \right) \tag{194}$$

By the assumption $\min_{j \in \mathcal{S}^*} |\gamma_j| \gg \sqrt{\log p/n}$ and (158), we establish (161).

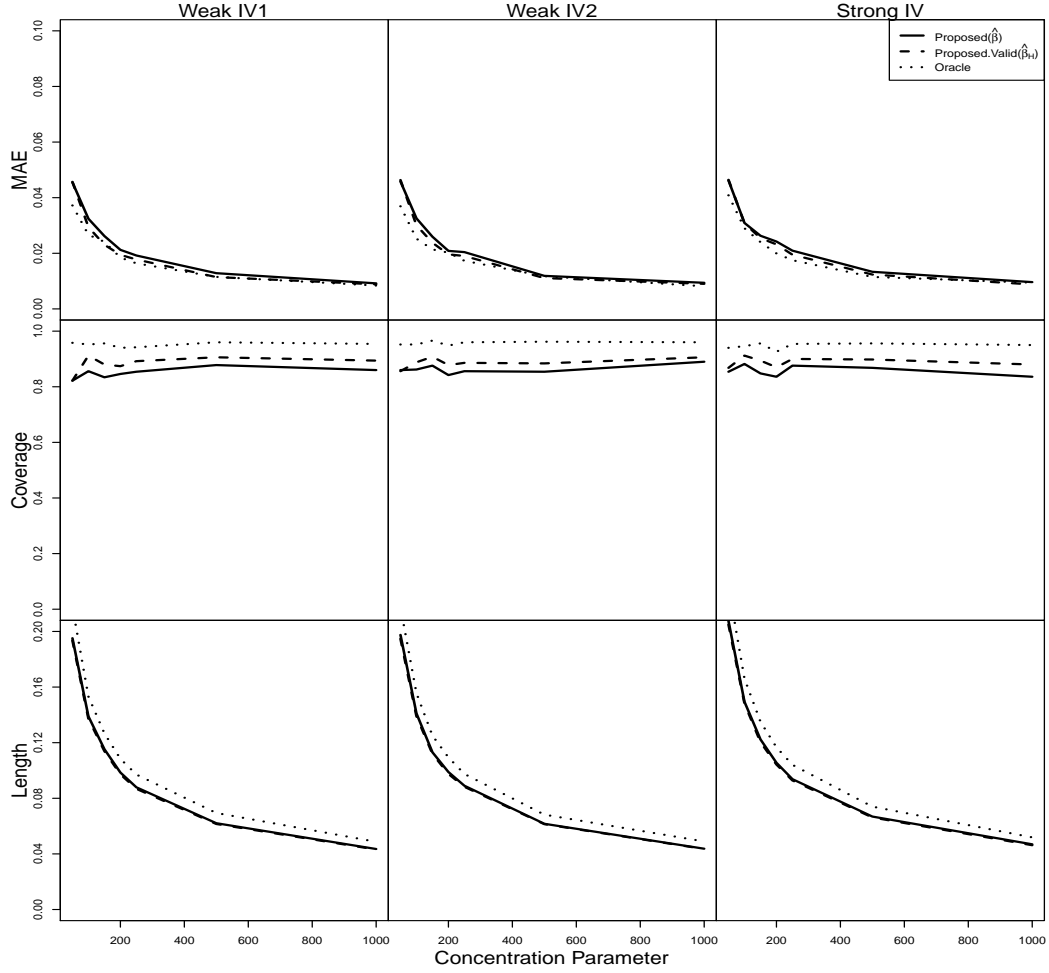


Fig. A1: Comparison of different methods in the case $\rho_2 = 0$, $p_z = 100$, $p_x = 150$ and $n = 100$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

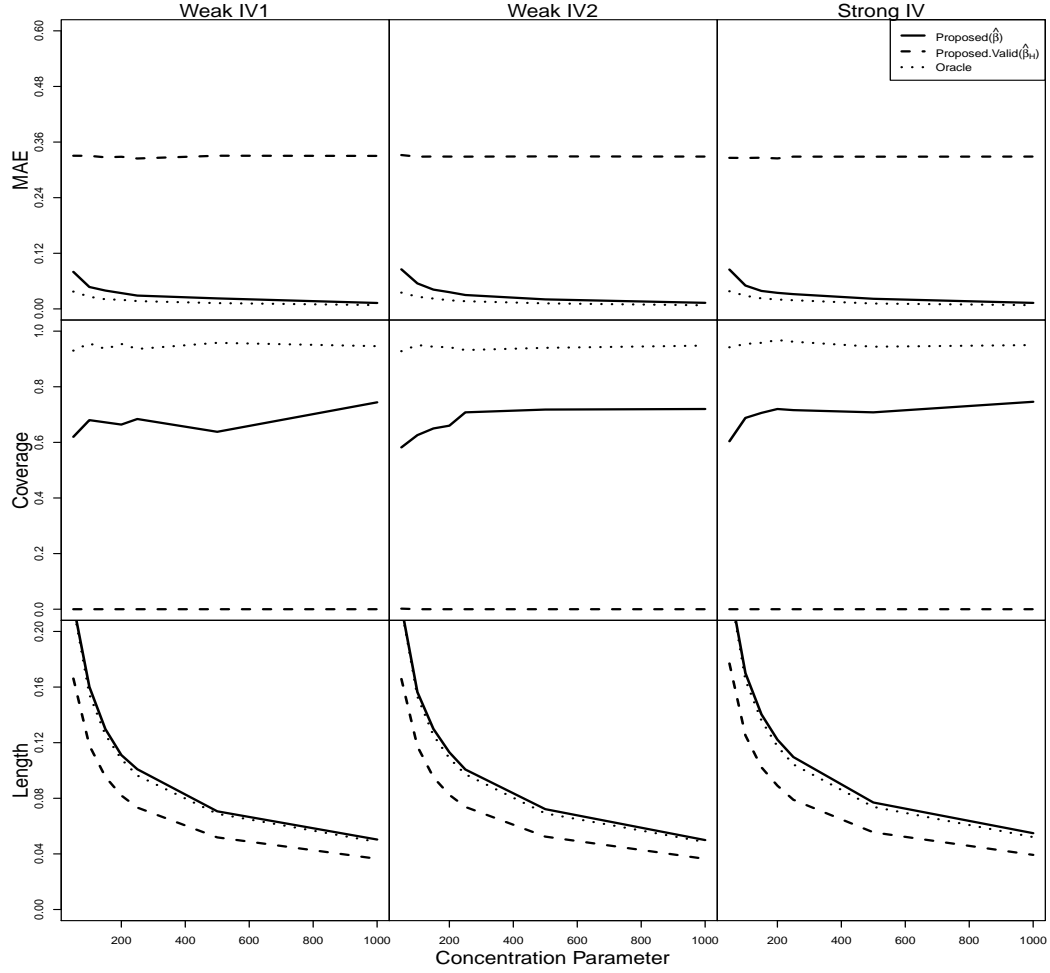


Fig. A2: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 100$, $p_x = 150$ and $n = 100$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

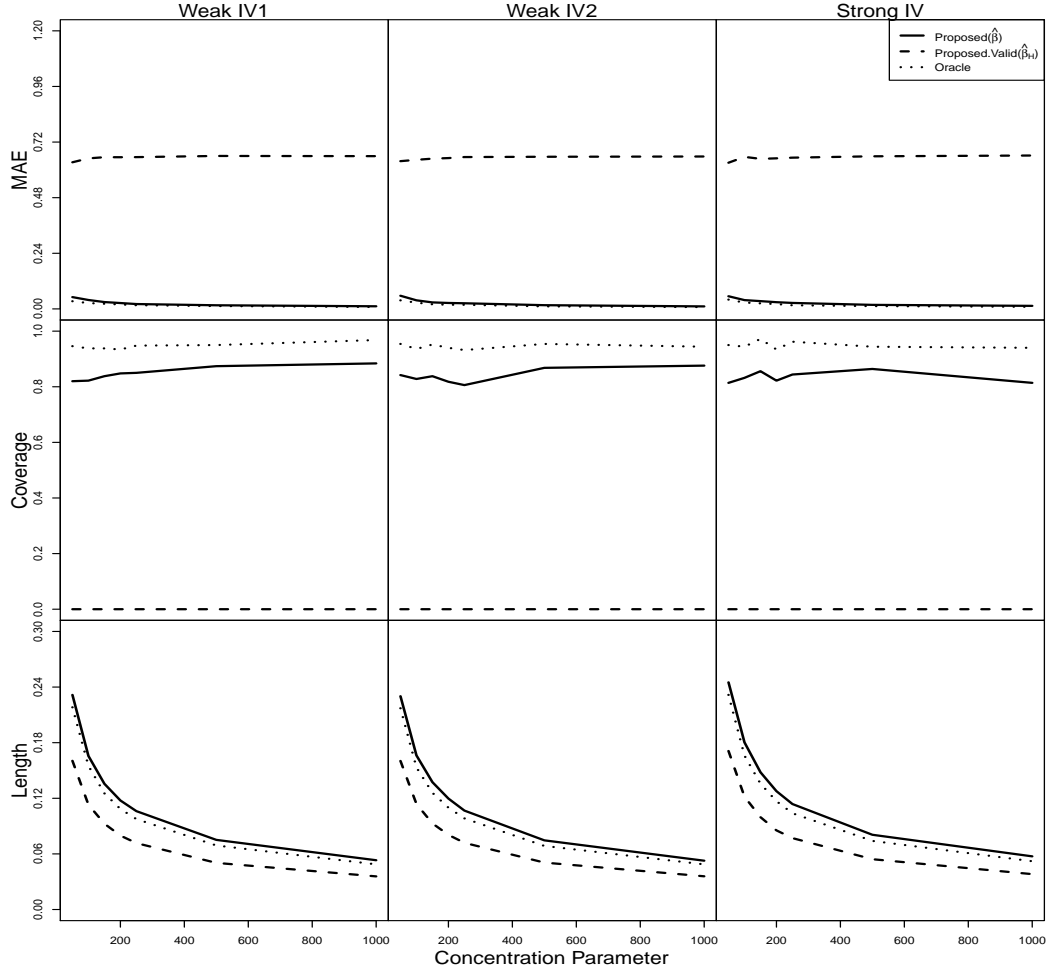


Fig. A3: Comparison of different methods in the case $\rho_2 = 2$, $p_z = 100$, $p_x = 150$ and $n = 100$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

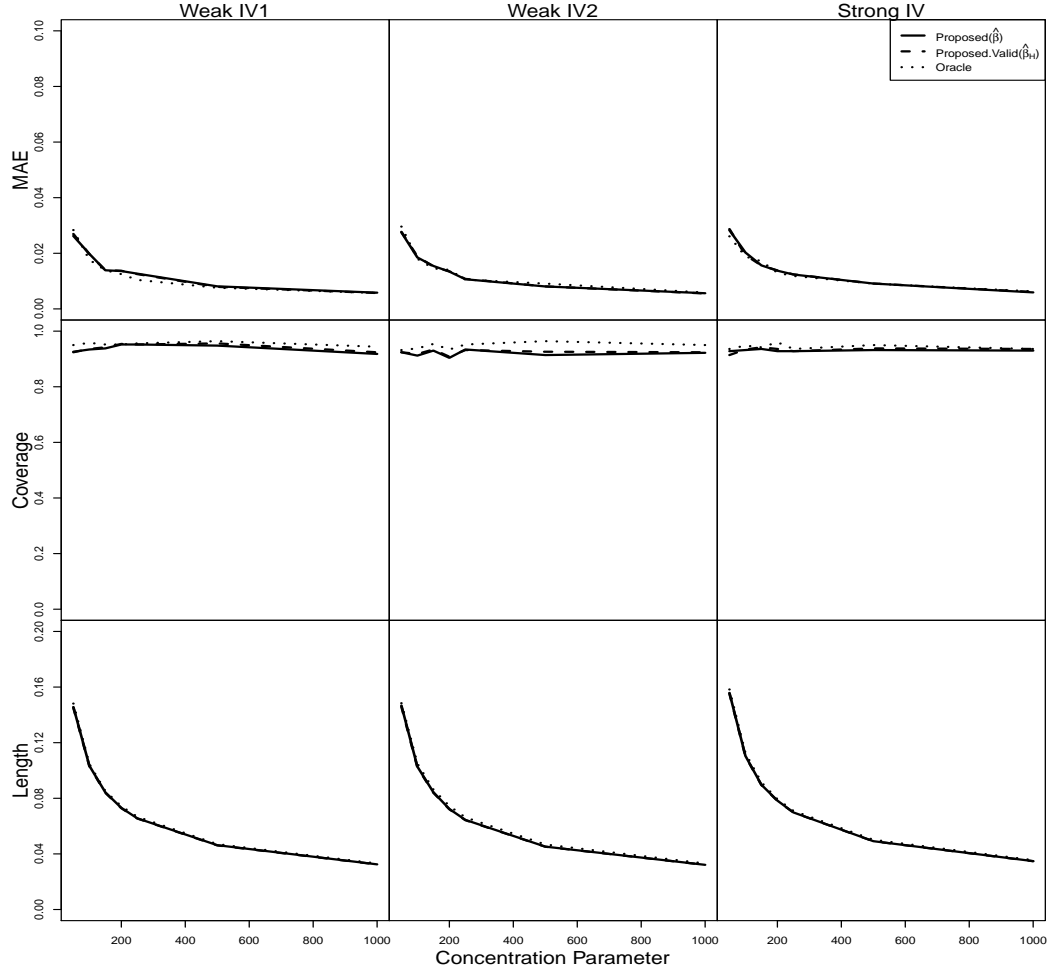


Fig. A4: Comparison of different methods in the case $\rho_2 = 0$, $p_z = 100$, $p_x = 150$ and $n = 200$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

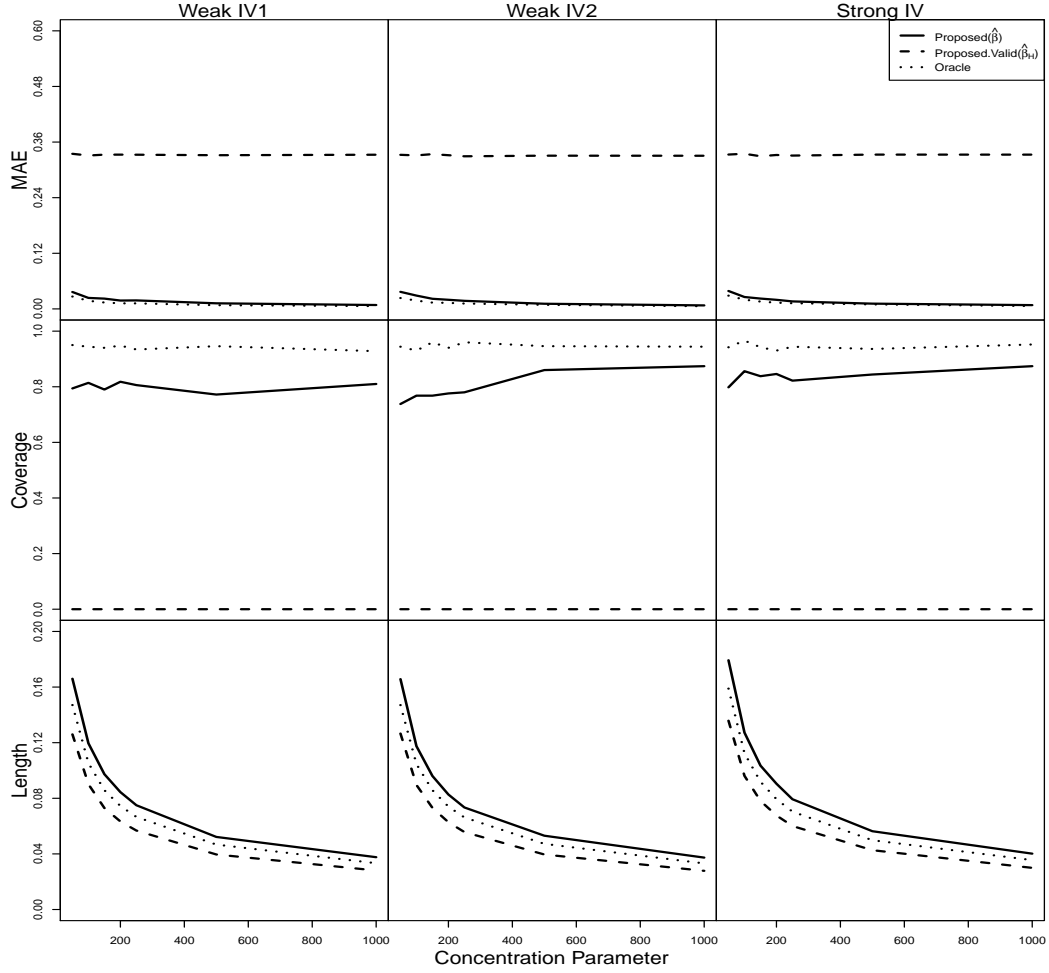


Fig. A5: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 100$, $p_x = 150$ and $n = 200$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

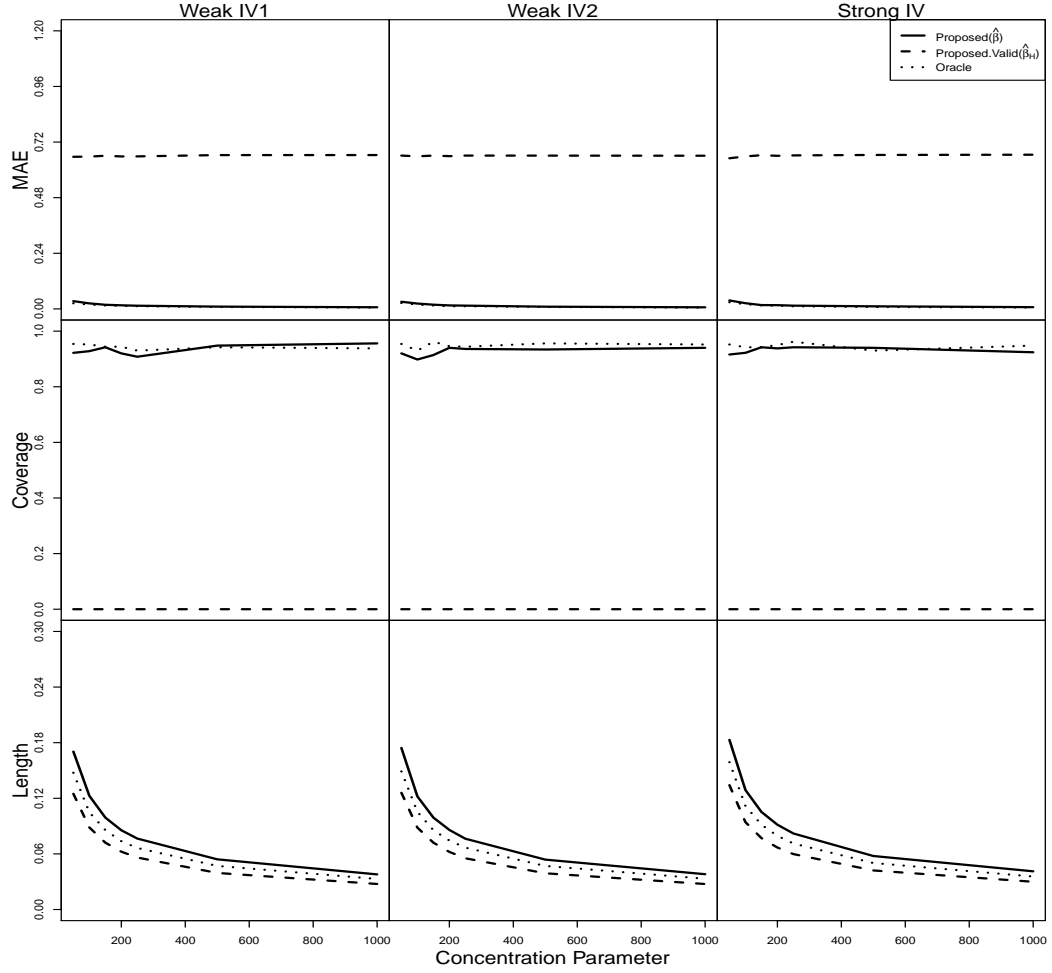


Fig. A6: Comparison of different methods in the case $\rho_2 = 2$, $p_z = 100$, $p_x = 150$ and $n = 200$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

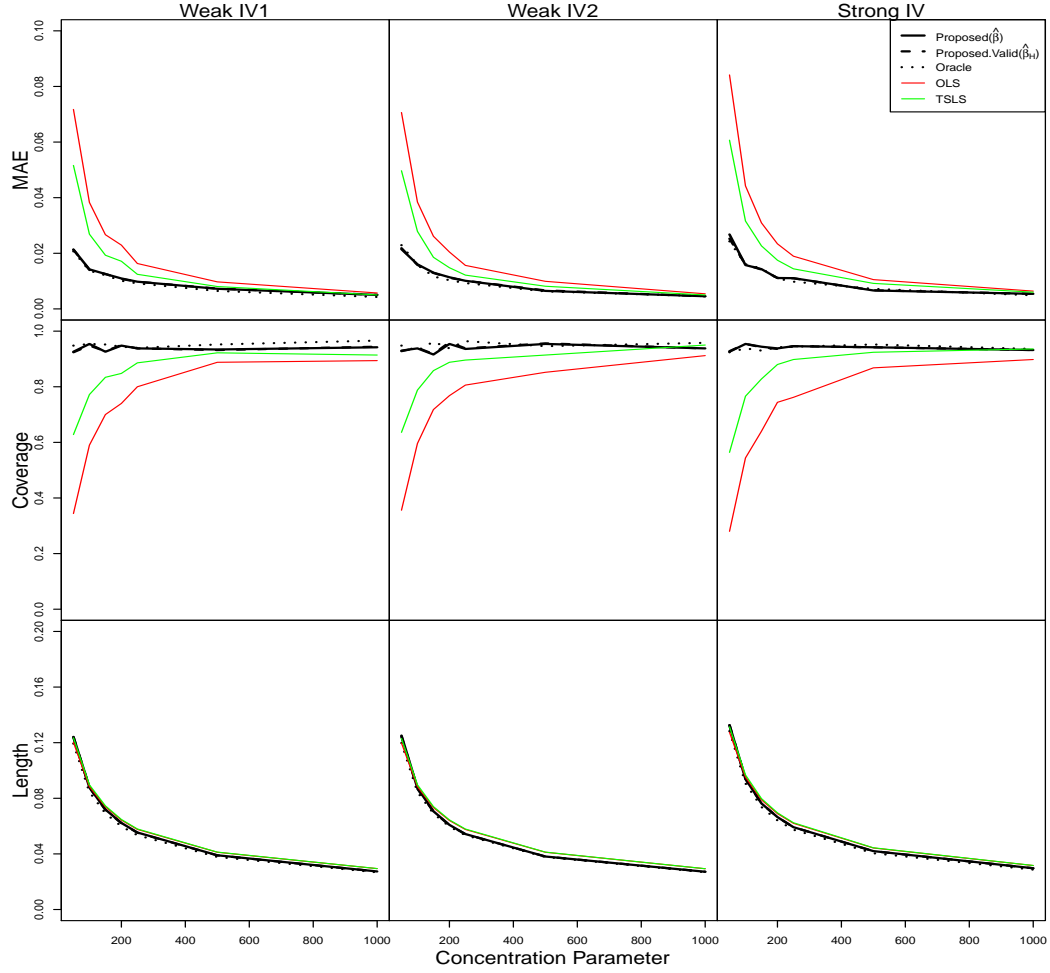


Fig. A7: Comparison of different methods in the case $\rho_2 = 0$, $p_z = 100$, $p_x = 150$ and $n = 300$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

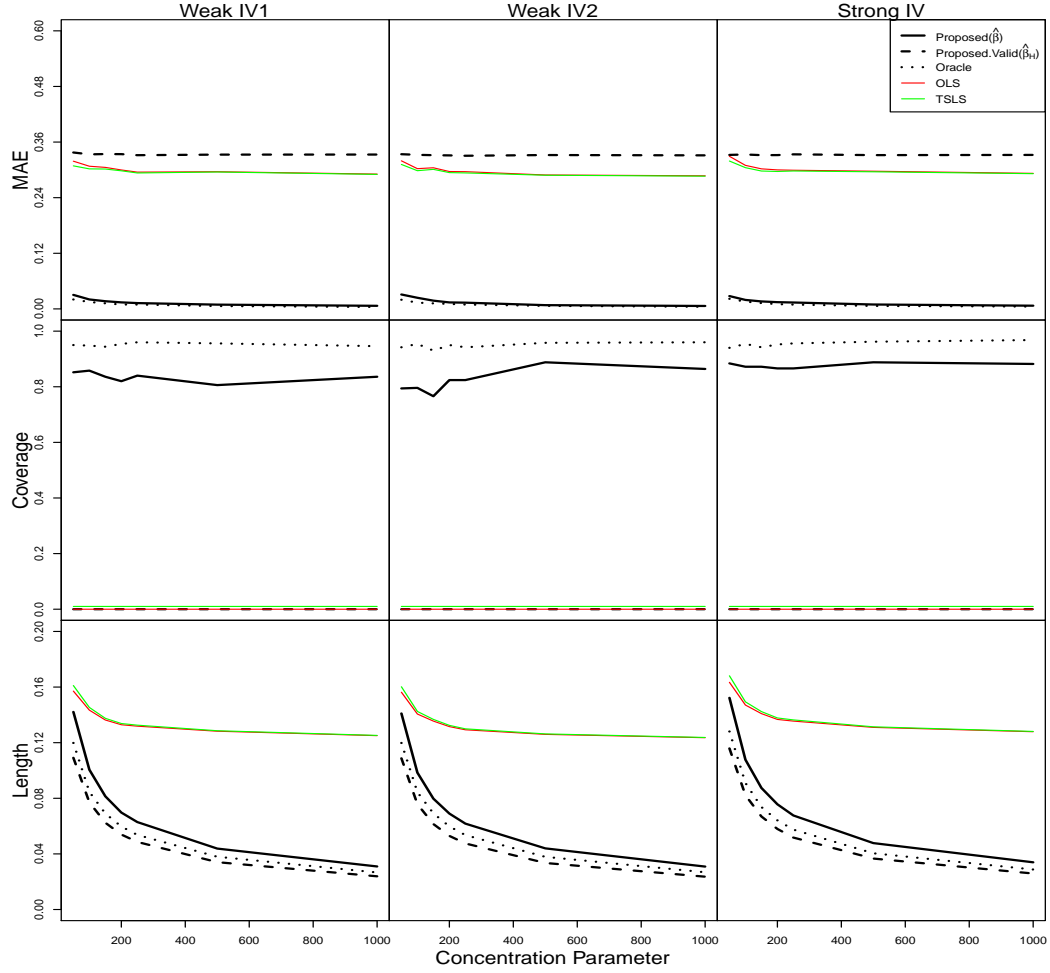


Fig. A8: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 100$, $p_x = 150$ and $n = 300$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

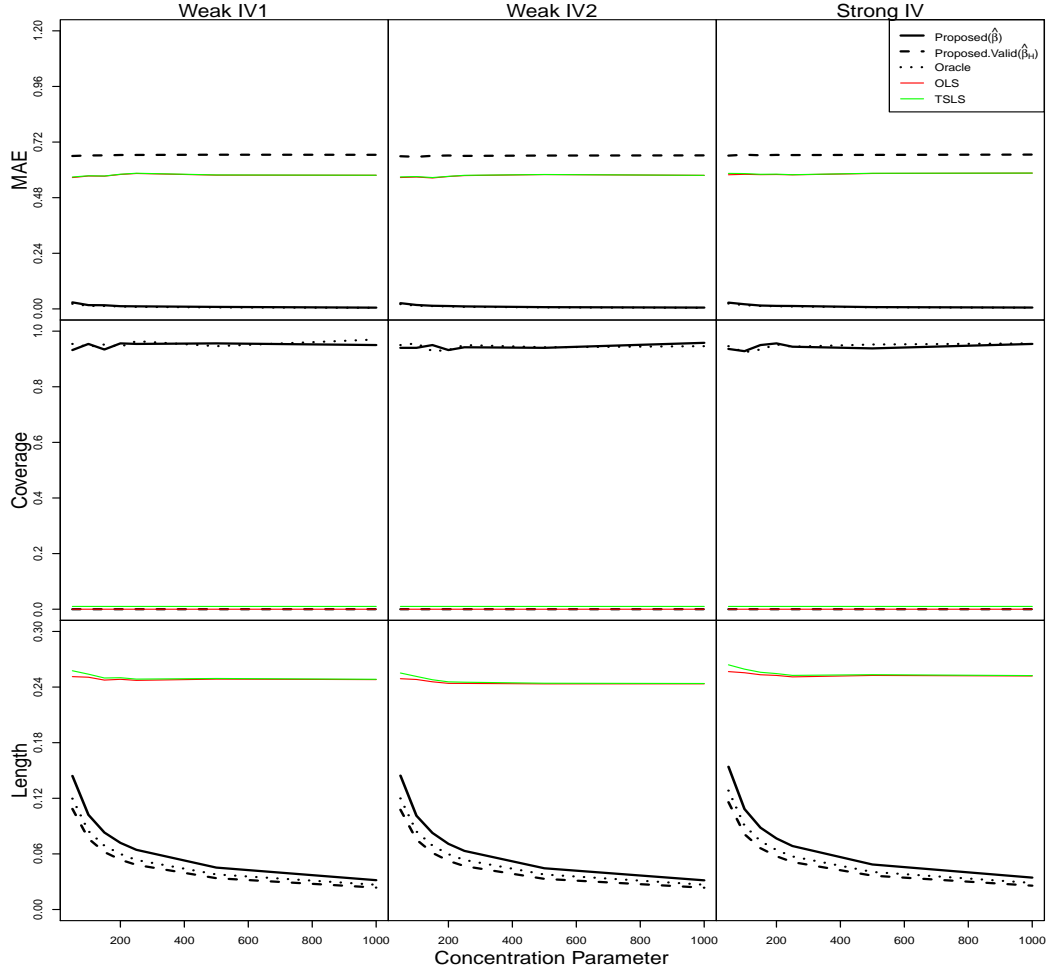


Fig. A9: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 100$, $p_x = 150$ and $n = 300$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

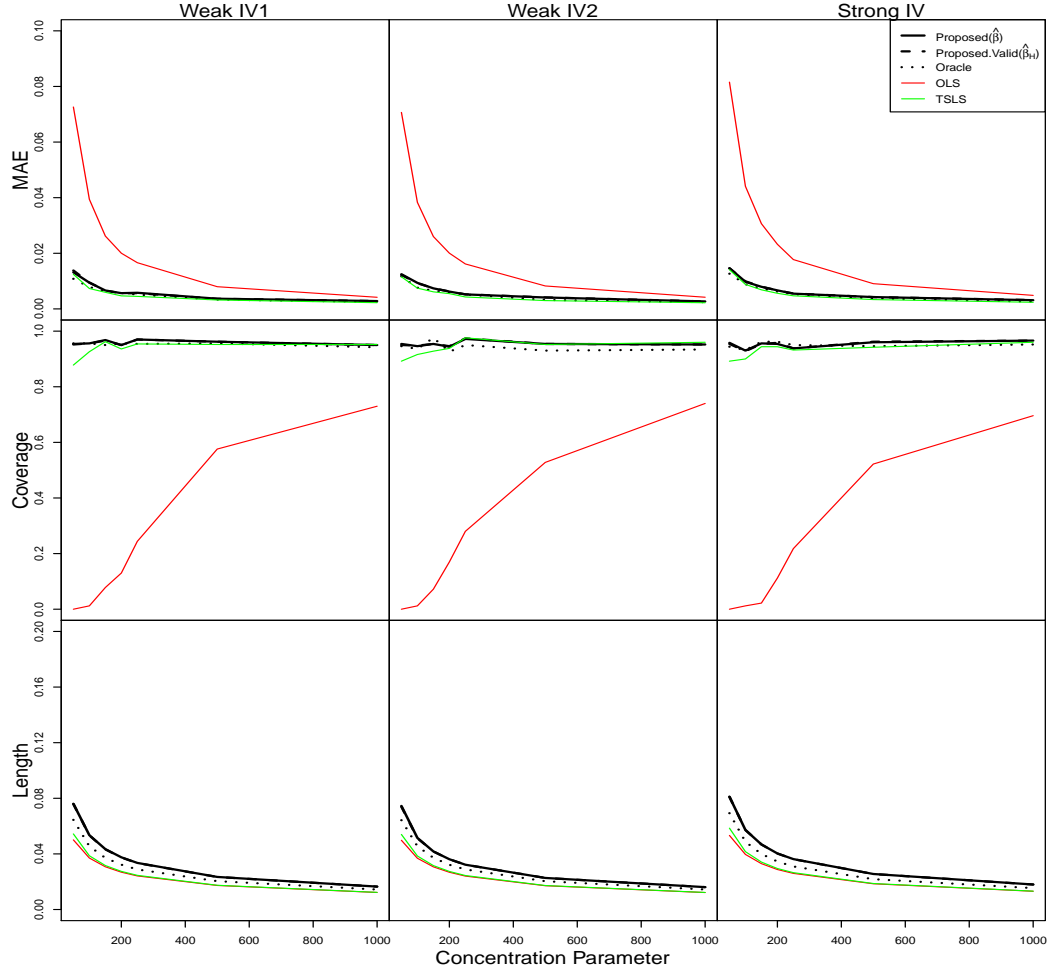


Fig. A10: Comparison of different methods in the case $\rho_2 = 0$, $p_z = 100$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

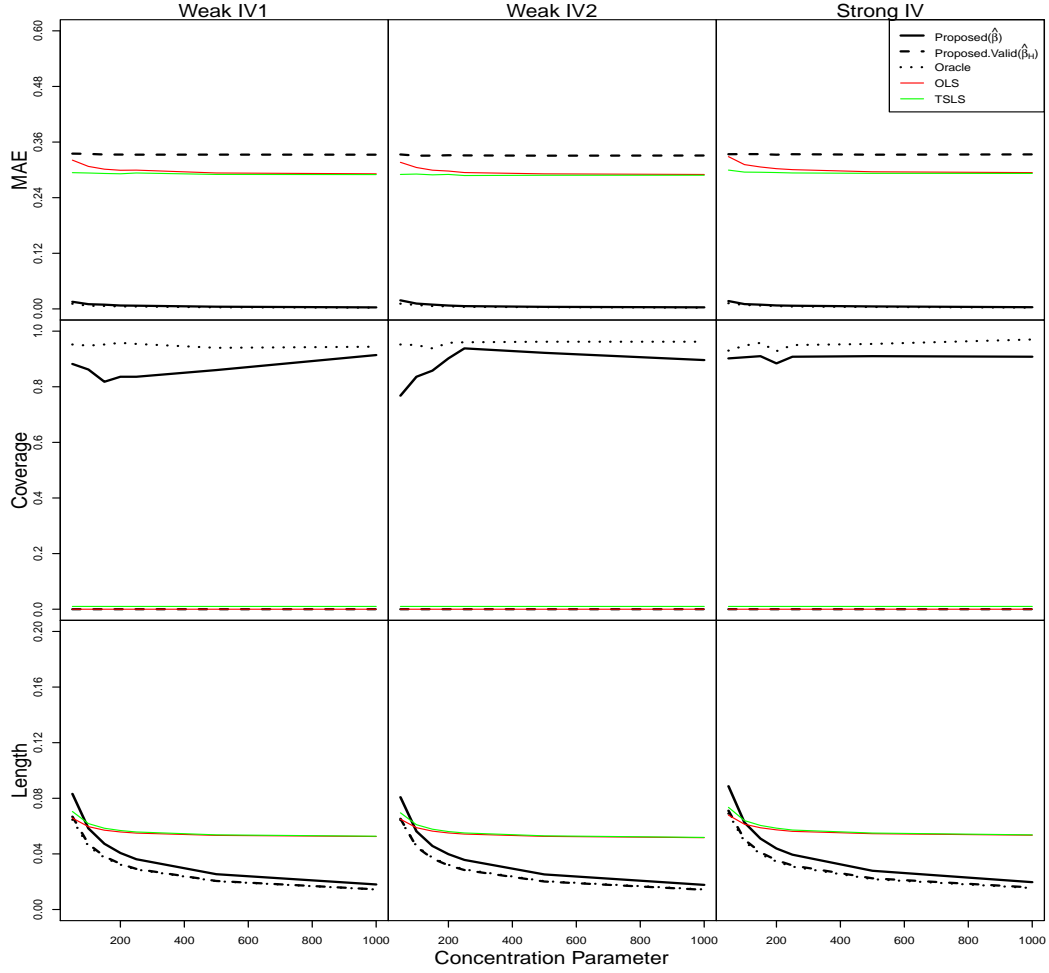


Fig. A11: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 100$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

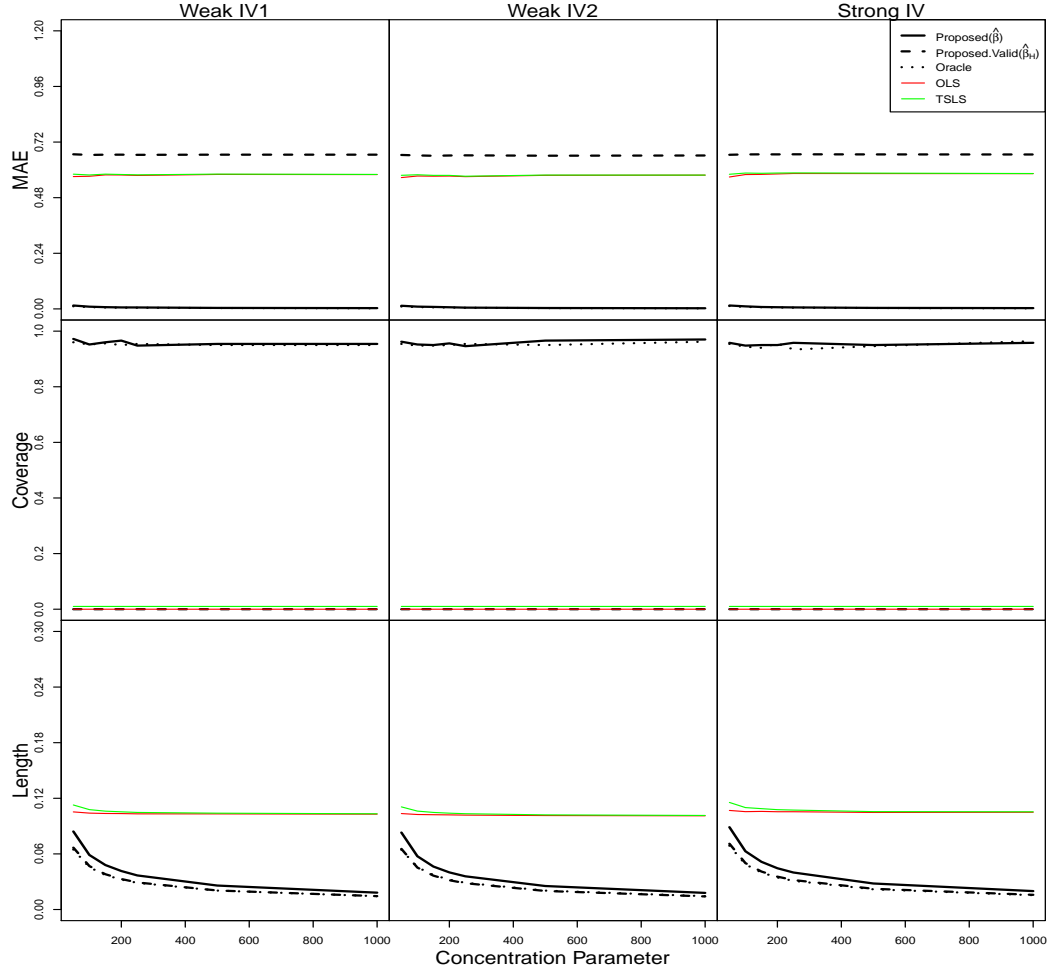


Fig. A12: Comparison of different methods in the case $\rho_2 = 2$, $p_z = 100$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

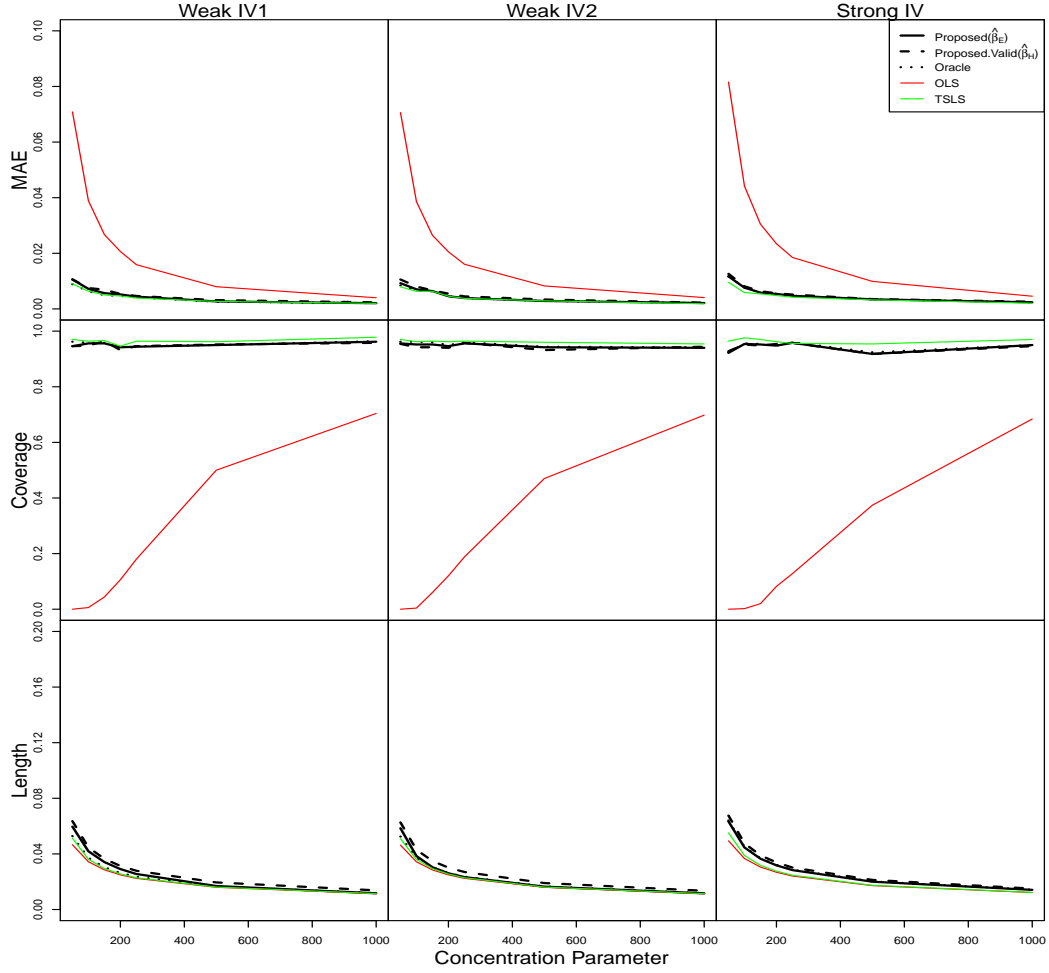


Fig. A13: Comparison of different methods in the case $\rho_2 = 0$, $p_z = 9$, $p_x = 10$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

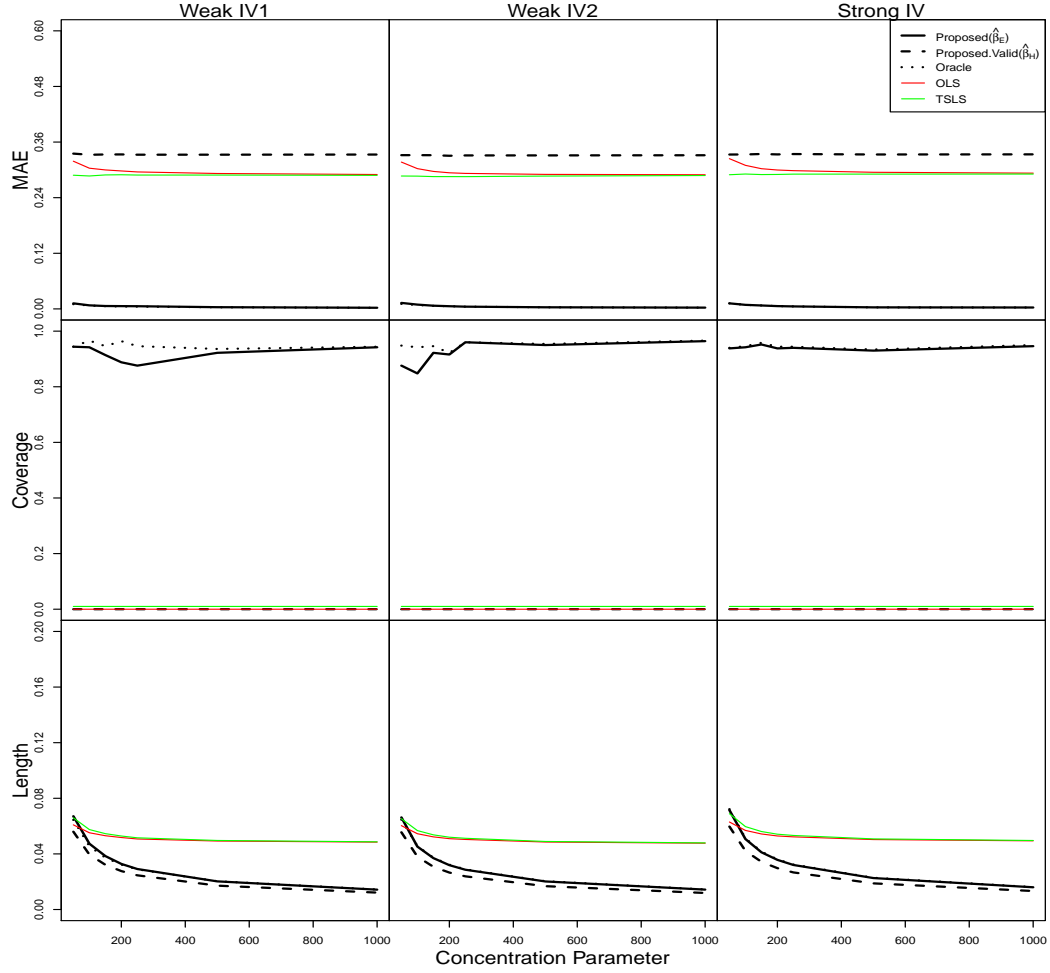


Fig. A14: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 9$, $p_x = 10$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

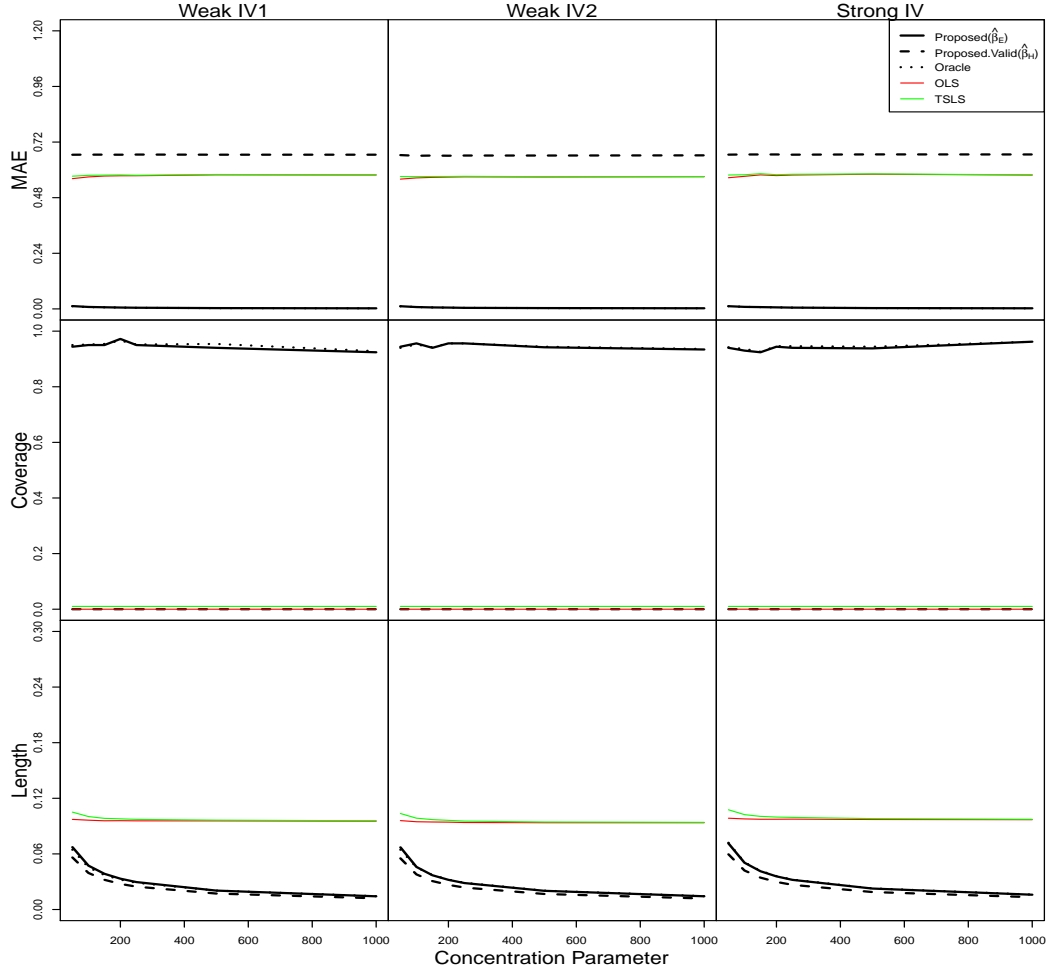


Fig. A15: Comparison of different methods in the case $\rho_2 = 2$, $p_z = 9$, $p_x = 10$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

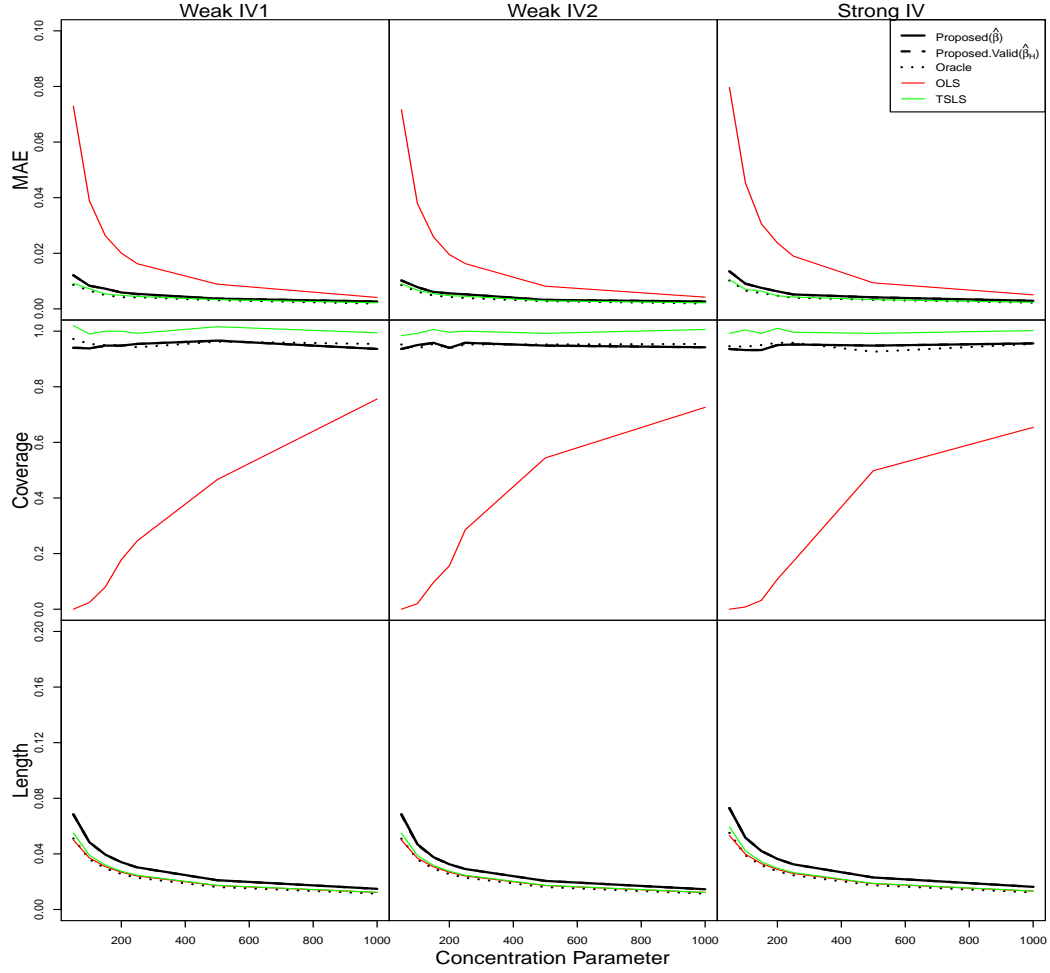


Fig. A16: Comparison of different methods in the case $\rho_2 = 0$, $p_z = 9$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

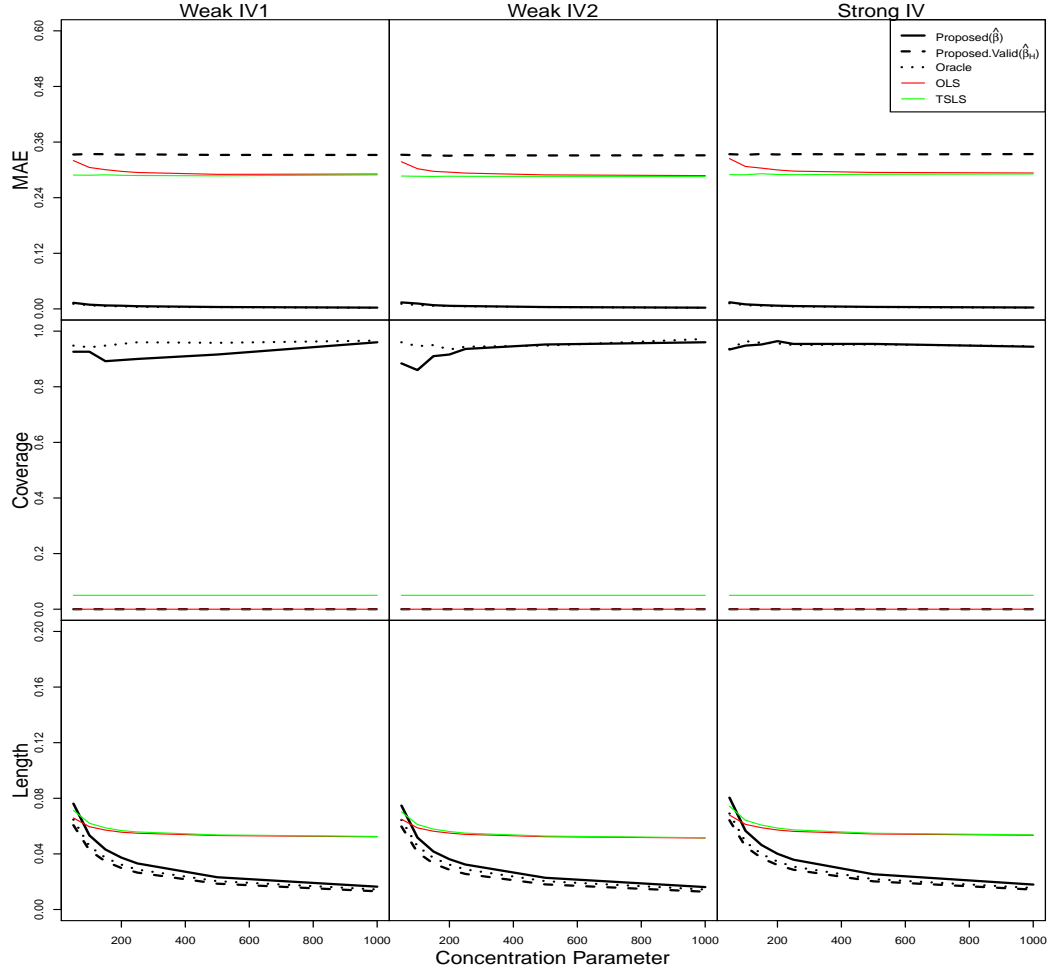


Fig. A17: Comparison of different methods in the case $\rho_2 = 1$, $p_z = 9$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.

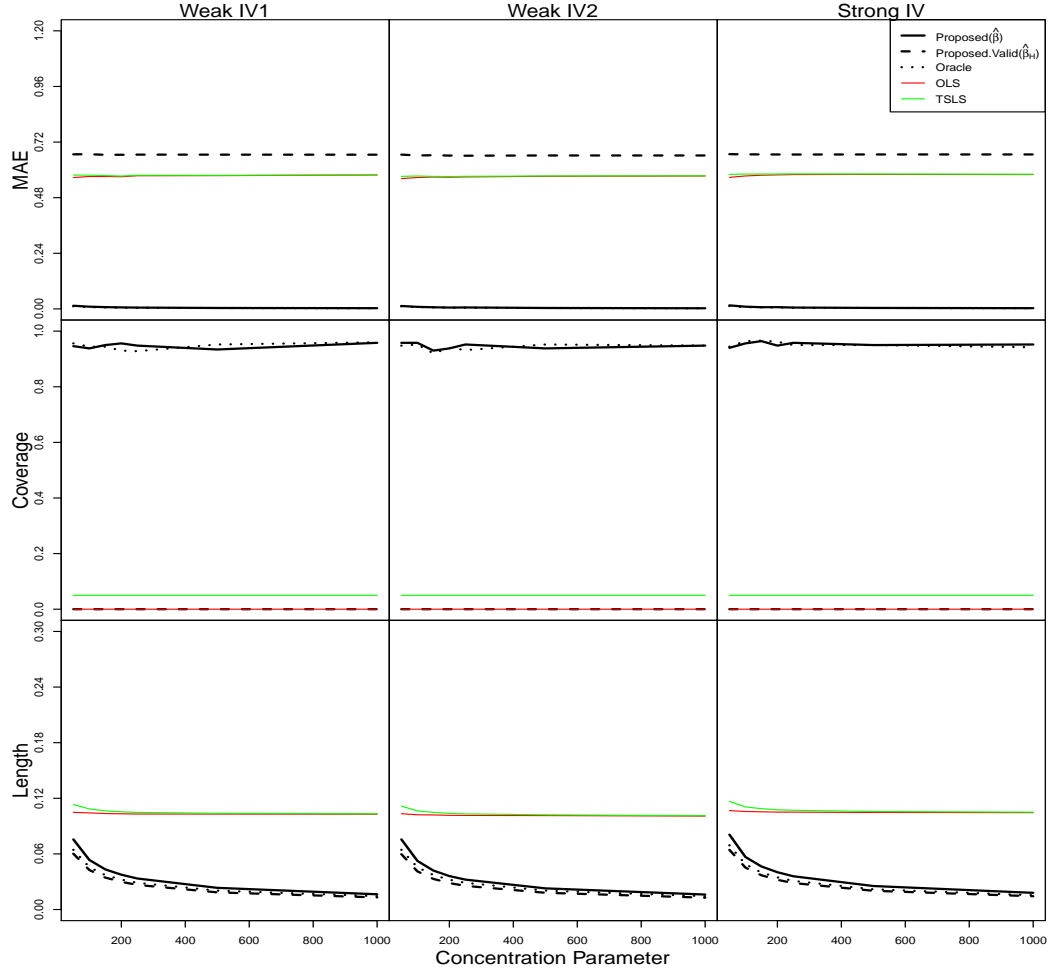


Fig. A18: Comparison of different methods in the case $\rho_2 = 2$, $p_z = 9$, $p_x = 150$ and $n = 1000$. The x -axis represents the concentration parameter. On the y -axis, MAE represents Median Absolute Error of the estimators, Coverage represents coverage of confidence intervals and Length represents the average length of confidence intervals. Proposed is our method allowing for invalid IVs and is represented by the solid line. Proposed.valid is our method that assumes all the IVs are valid and is represented by the dashed line. Oracle is the method that knows exactly which instruments are valid and is represented by the dotted line. The column labeled with Weak IV1 represents the case $\rho_1 = 0.1$, the column labeled with Weak IV2 represents the case $\rho_1 = 0.2$ and the column labeled with Strong IV represents the case $\rho_1 = 0$.