

On the instrumental variable estimation with many weak and invalid instruments

Yiqi Lin^{a,b}, Frank Windmeijer^b, Xinyuan Song^a, Qingliang Fan^{c†}

^a*Department of Statistics, The Chinese University of Hong Kong, Hong Kong.*

^b*Department of Statistics, University of Oxford, Oxford, U.K.*

^c*Department of Economics, The Chinese University of Hong Kong, Hong Kong.*

Abstract.

We discuss the fundamental issue of identification in linear instrumental variable (IV) models with unknown IV validity. We revisit the popular majority and plurality rules and show that no identification condition can be “if and only if” in general. With the assumption of the “sparsest rule”, which is equivalent to the plurality rule but becomes operational in computation algorithms, we investigate and prove the advantages of non-convex penalized approaches over other IV estimators based on two-step selections, in terms of selection consistency and accommodation for individually weak IVs. Furthermore, we propose a surrogate sparsest penalty that aligns with the identification condition and provides oracle sparse structure simultaneously. Desirable theoretical properties are derived for the proposed estimator with weaker IV strength conditions compared to the previous literature. Finite sample properties are demonstrated using simulations and the selection and estimation method is applied to an empirical study concerning the effect of trade on economic growth.

Keywords: Invalid Instruments, Model Identification, Non-convex Penalty, Treatment Effect, Weak Instruments.

1. Introduction

Recently, estimation of causal effects with high-dimensional observational data has drawn much attention in many research fields such as economics, epidemiology and genomics. The instrumental variable (IV) method is widely used when the treatment variable of interest is endogenous. As shown in Figure 1, the ideal IV needs to be correlated with the endogenous treatment variable (C1), it should not have a direct effect on the outcome (C2) and should not be related to unobserved confounders that affect both outcome and treatment (C3).

Our research is motivated by the difficulty of finding IVs that satisfy all the above conditions. In applications, invalid IVs (violation of C2 or C3) ([Davey Smith and Ebrahim, 2003](#); [Kang et al., 2016](#); [Windmeijer et al., 2019](#)) and weak IVs (violation of C1) ([Bound et al., 1995](#); [Staiger and Stock, 1997](#)) are prevalent. A strand of literature studies the

†Correspondence: Qingliang Fan, Department of Economics, The Chinese University of Hong Kong, Shatin, N.T., Hong Kong. Email: michaelqfan@gmail.com

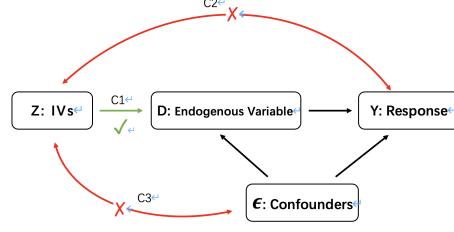


Figure 1. Relevance and Validity of IVs

“many weak IVs” problem (Stock et al., 2002; Chao and Swanson, 2005). With the increasing availability of large datasets, IV models are often high-dimensional (Belloni et al., 2012; Lin et al., 2015; Fan and Zhong, 2018), and have potentially weak IVs (Andrews et al., 2018), and invalid IVs (Guo et al., 2018; Windmeijer et al., 2021). Among those problems, we mainly focus on the invalid IV problem, while allowing for potential high-dimensionality and weak signals.

1.1. Related Works

Most related works fall into two main categories: robust estimation with invalid IVs and estimation which can select valid IVs without any prior knowledge of validity. The first strand of literature allows all IVs to be invalid. For example, Kolesár et al. (2015) restricted the direct effects of IVs on treatment and response being random effects and independent. In practice, this assumption might be difficult to justify. Lewbel (2012); Tchetgen et al. (2021); Guo and Bühlmann (2022) utilized conditional heteroskedasticity to achieve robustness with potentially all IVs invalid. However, their performances are not satisfactory once heteroskedasticity (identification condition) is not evident.

The second strand focused on unknown invalid IVs, while imposing certain identification conditions on the number of valid IVs. Kang et al. (2016) proposed a Lasso type estimator (sisVIVE). Windmeijer et al. (2019) pointed out the inconsistent variable selection of sisVIVE under a relative IV strength condition and proposed an adaptive Lasso estimator, which has asymptotic oracle properties under the assumption that more than half of the IVs are valid, also called the majority rule. Guo et al. (2018); Windmeijer et al. (2021) further developed two-step (the first one for relevance, the second one for validity) selection approaches, Two-Stage Hard Thresholding (TSHT) and Confidence Intervals IV (CIIV), respectively, under the plurality rule conditional on the set of relevant IVs. The plurality rule states that the valid IVs form the largest group. However, the approaches mentioned above are not robust to many weak IVs due to the restriction of the majority/plurality rule in strong IVs instead of all IVs. Our method closely follows this strand of literature. Instead of a two-step selection, we require the plurality rule for valid IVs amid a one-step selection procedure, thus considerably relaxing the requirement of valid IVs in theory and most practical scenarios.

The study of many (weak) IVs originated from empirical motivations but often assumed known validity. For example, Staiger and Stock (1997); Hansen et al. (2008); Newey and Windmeijer (2009); Hansen and Kozbur (2014) considered different estima-

tors in many (weak) valid IVs but fixed the number of known covariates. Kolesár et al. (2015); Kolesár (2018) allowed the number of covariates to grow with the sample size. We consider the weak IV issues that are prevalent in empirical studies.

1.2. The Main Results and Contributions

We propose a **Weak and Invalid IV robust Treatment effect (WIT)** estimator. The sparsest rule is sufficient for the identification and is operational in numerical optimization. The proposed procedure has a selection stage (regarding IV validity) and a post-selection estimation stage. The selection stage is a penalized IV-regression via minimax concave penalty (MCP, Zhang et al., 2010), a proper surrogate penalty aligned with the identification condition to achieve model selection consistency of valid IVs under much weaker technical conditions than existing methods (Guo et al., 2018; Windmeijer et al., 2021). In the estimation stage, we utilize the limited information maximum likelihood (LIML) estimator to handle the weak IVs (Staiger and Stock, 1997). An efficient computational algorithm for the optimal solution is provided. The computer codes for implementing the WIT estimator are available at <https://github.com/QoifoQ/WIT>.

The key contributions of this paper are summarized in the following.

- (a) We provide a self-contained framework to investigate the fundamental problem in model identification for linear IV models with unknown validity. Specifically, we study the identification condition from the general data generating process (DGP) framework. It addresses an earlier caveat on the if and only if (*iff*) condition statement of the plurality rule (Guo et al., 2018, Theorem 1). Furthermore, we develop a theorem on the impossibility result of the existence of an *iff* condition on the model identification in the linear IV model framework.
- (b) This study extends the IV estimation with unknown invalid IVs (namely, Kang et al., 2016; Guo et al., 2018; Windmeijer et al., 2019, 2021) to allow for many potentially weak IVs. We show that the sparsest rule, equivalent to the plurality rule on *the whole IV set*, could accommodate weak IVs in empirically relevant scenarios. Furthermore, we revisit the penalized approaches using the sparsest rule and propose a concept of proper surrogate sparsest penalty that targets identification conditions and provides sparse structure. We propose to deploy MCP as a surrogate sparsest penalty and ensure the targeted solution is the global minimizer. On the contrary, the existing methods (Kang et al., 2016; Windmeijer et al., 2019) do not fit the surrogate sparsest penalty and hence are mistargeting the model identification.
- (c) Our method is a one-step valid IV selection instead of the previous sequential two-step selections (Guo et al., 2018; Windmeijer et al., 2021). This allows us to utilize individually weak IVs instead of discarding them totally. We provide theoretical foundations to ensure the compatibility of weak IVs under a mild minimal signal condition. Formally, we establish the selection consistency of the valid IV set, the consistency, and asymptotic normality of the proposed treatment effect estimator under many potentially invalid and weak IVs, where both the number of valid and invalid IVs are increasing with the sample size n . We also provide the theoretical results for the case of a fixed and finite number of IVs.

The article is organized as follows. In Section 2, we describe the model with some invalid IVs and analyze identification conditions in a general way. In Section 3, we present the methodology and the novel WIT estimator. We establish the theorems to identify the valid IVs, estimation consistency, and asymptotic normality. Section 4 shows the finite sample performance of our proposed estimator using comprehensive numerical experiments. Section 5 applies our methods to an empirical international trade and growth study. Section 6 concludes. All the technical details and proofs are provided in the appendix.

2. Model and Identification Strategy

2.1. Potential Outcome Model with Some Invalid IVs

We consider the potential outcome model as in Small (2007); Kang et al. (2016). For $i = 1, 2, \dots, n$, let $Y_i^{(d,z)}$ be the potential outcome for object i having exposure d and instruments' values $\mathbf{z} \in \mathbb{R}^p$. The observed outcome for object i is denoted by the scalar Y_i , the treatment by the scalar D_i and the vector of p potential instruments by \mathbf{Z}_i . For two different sets of treatment values and IVs (d^A, \mathbf{z}^A) and (d^B, \mathbf{z}^B) , respectively, assume

$$Y_i^{(d^B, \mathbf{z}^B)} - Y_i^{(d^A, \mathbf{z}^A)} = (\mathbf{z}^B - \mathbf{z}^A)^\top \boldsymbol{\phi} + (d^B - d^A) \beta \text{ and } E(Y_i^{(0,0)} | \mathbf{Z}_i) = \mathbf{Z}_i^\top \boldsymbol{\theta}, \quad (1)$$

where $\boldsymbol{\phi}$ and $\boldsymbol{\theta}$ measure the direct effects of instruments on responses and the association between IVs and unmeasured confounders, respectively. Note that \mathbf{Z}_i could have non-linear transformations of the original variables (such as polynomials and B-splines), so a high-dimensional model is plausible (Appendix A1 contains a further discussion). A good instrument Z_j should not have a direct effect on the response and unmeasured confounders, i.e., $\phi_j = 0$ and $\theta_j = 0$, for $j = 1, 2, \dots, p$. For $i = 1, 2, \dots, n$, we have the random sample (Y_i, D_i, \mathbf{Z}_i) . Assuming a linear functional form between treatments D_i and instruments \mathbf{Z}_i , also called the first-stage specification, the above potential outcome model (1) results in the following observed data model,

$$\begin{aligned} Y_i &= D_i \beta + \mathbf{Z}_i^\top \boldsymbol{\alpha} + \epsilon_i, \\ D_i &= \mathbf{Z}_i^\top \boldsymbol{\gamma} + \eta_i. \end{aligned} \quad (2)$$

where $\epsilon_i = Y_i^{(0,0)} - E(Y_i^{(0,0)} | \mathbf{Z}_i)$ and $\boldsymbol{\alpha} = \boldsymbol{\phi} + \boldsymbol{\theta}$. Let $\boldsymbol{\alpha}^*, \beta^*$ and $\boldsymbol{\gamma}^*$ represent the true coefficients in (2).

Following Kang et al. (2016), we define the valid instruments as follows,

DEFINITION 1. For $j = 1, \dots, p$, the j -th instrument is valid if $\alpha_j^* = 0$.

Define the valid IV set $\mathcal{V}^* = \{j : \alpha_j^* = 0\}$ and invalid IV set $\mathcal{V}^{c*} = \{j : \alpha_j^* \neq 0\}$. Let $p_{\mathcal{V}^*} = |\mathcal{V}^*|$, $p_{\mathcal{V}^{c*}} = |\mathcal{V}^{c*}|$ and $p = p_{\mathcal{V}^*} + p_{\mathcal{V}^{c*}}$. Notably, $p_{\mathcal{V}^*} \geq 1$ refers to the existence of an excluded IV, thus satisfying the order condition (Wooldridge, 2010). Let the $n \times p$ matrix of observations on the instruments be denoted by \mathbf{Z} , and the n -vectors of outcomes and treatments by \mathbf{Y} and \mathbf{D} , respectively. We consider the cases of many and weak IVs in (2) and make the following model assumptions:

Assumption 1 (Many valid and invalid IVs): $p < n$, $p_{\mathcal{V}^{c*}}/n \rightarrow v_{p_{\mathcal{V}^{c*}}} + o(n^{-1/2})$ and $p_{\mathcal{V}^*}/n \rightarrow v_{p_{\mathcal{V}^*}} + o(n^{-1/2})$ for some non-negative constants $v_{p_{\mathcal{V}^{c*}}}$ and $v_{p_{\mathcal{V}^*}}$ such that $0 \leq v_{p_{\mathcal{V}^*}} + v_{p_{\mathcal{V}^{c*}}} < 1$.

Assumption 2: Assume \mathbf{Z} is standardized. It then has full column rank and $\|\mathbf{Z}_j\|_2^2 \leq n$ for $j = 1, 2, \dots, p$.

Assumption 3: Let $\mathbf{u}_i = (\epsilon_i, \eta_i)^\top$. $\mathbf{u}_i \mid \mathbf{Z}_i$ are i.i.d. and follow a multivariate normal distribution with mean zero and positive definite covariance matrix $\Sigma = \begin{pmatrix} \sigma_\epsilon^2 & \sigma_{\epsilon,\eta} \\ \sigma_{\epsilon,\eta} & \sigma_\eta^2 \end{pmatrix}$. The elements of Σ are finite and $\sigma_{\epsilon,\eta} \neq 0$.

Assumption 4 (Strength of valid IVs): The concentration parameter μ_n grows at the same rate as n , i.e., $\mu_n := \gamma_{\mathbf{Z}_{\mathcal{V}^*}}^*{}^\top \mathbf{Z}_{\mathcal{V}^*}^\top M_{\mathbf{Z}_{\mathcal{V}^{c*}}} \mathbf{Z}_{\mathcal{V}^*} \gamma_{\mathbf{Z}_{\mathcal{V}^*}}^* / \sigma_\eta^2 \rightarrow \mu_0 n$, for some $\mu_0 > 0$.

Assumption 1 is identical to the assumption of many instruments in [Kolesár et al. \(2015\)](#); [Kolesár \(2018\)](#). It relaxes the conventional many IVs assumptions ([Bekker, 1994](#); [Chao and Swanson, 2005](#)) that only allow the dimension of valid IVs $p_{\mathcal{V}^*}$ to grow with n . Also, it has not been considered in the literature on selecting valid IVs ([Kang et al., 2016](#); [Guo et al., 2018](#); [Windmeijer et al., 2021, 2019](#)). Assumption 2 is standard for data preprocessing and scaling \mathbf{Z}_j . Assumption 3 follows [Guo et al. \(2018\)](#); [Windmeijer et al. \(2021\)](#) to impose the homoskedasticity assumption and endogeneity issue of treatment D_i . The normality is not required for the selection stage, only for the estimation stage (asymptotics of embedded LIML estimator). Notably, Assumption 4 requires a strong-identified moment condition in terms of concentration parameter in the literature ([Bekker, 1994](#); [Newey and Windmeijer, 2009](#)). In the fixed p case, it indicates the presence of a constant coefficient γ_j , $\exists j$, and the rest of IVs could be weak ([Staiger and Stock, 1997](#)). Specifically, we model weakly correlated IVs as $\gamma = Cn^{-\tau}$, $0 < \tau \leq 1/2$, which is the “local to zero” case ([Staiger and Stock, 1997](#)). Essentially this is a mixture of constant γ -type and asymptotically diminishing γ -type instruments for fixed p . For $p \rightarrow \infty$ in the same order as n , we allow all the IVs to be weak in the “local to zero” case with specified rates. This IV strength assumption can be further weakened along the lines of [Hansen et al. \(2008\)](#) to have weak identification asymptotics. In this paper, we focus on the individually weak (diminishing to zero as in [Staiger and Stock, 1997](#)) signals model in high-dimensionality. Notice our model allows for much weaker individually weak IVs regardless of their validity (as long as the concentration parameter satisfies Assumption 4), unlike that of [Guo et al. \(2018\)](#). Nevertheless, the constant μ_0 can be a small number to accommodate empirically relevant finite samples with many individually weak IVs.

2.2. Identifiability of Model (2)

The following moment conditions can be derived from Model (2):

$$E\left(\mathbf{Z}^\top (\mathbf{D} - \mathbf{Z}\gamma^*)\right) = \mathbf{0}, \quad E\left(\mathbf{Z}^\top (\mathbf{Y} - \mathbf{D}\beta^* - \mathbf{Z}\alpha^*)\right) = \mathbf{0} \quad \Rightarrow \quad \mathbf{\Gamma}^* = \alpha^* + \beta^*\gamma^*, \quad (3)$$

where $\mathbf{\Gamma}^* = E(\mathbf{Z}^\top \mathbf{Z})^{-1} E(\mathbf{Z}^\top \mathbf{Y})$ and $\gamma^* = E(\mathbf{Z}^\top \mathbf{Z})^{-1} E(\mathbf{Z}^\top \mathbf{D})$, both are identified by the reduced form models. Without the exact knowledge about which IVs are valid, [Kang et al. \(2016\)](#) considered the identification of (α^*, β^*) via the unique mapping of

$$\beta_j^* = \mathbf{\Gamma}_j^* / \gamma_j^* = \beta^* + \alpha_j^* / \gamma_j^*. \quad (4)$$

Notice that the moment conditions (3) consist of p equations, but $(\boldsymbol{\alpha}^*, \beta^*) \in \mathbb{R}^{p+1}$ need to be estimated and is hence under-identified without further restrictions. Kang et al. (2016) proposed a sufficient condition, called majority rule (first proposed by Han, 2008), such that $p_{\mathcal{V}^*} \geq \lceil p/2 \rceil$, to identify the model parameters without any prior knowledge of the validity of individual IV. However, the majority rule could be restrictive in practice. Guo et al. (2018) further relaxed it to the plurality rule as follows:

$$\text{Plurality Rule: } \left| \mathcal{V}^* = \{j : \alpha_j^*/\gamma_j^* = 0\} \right| > \max_{c \neq 0} \left| \{j : \alpha_j^*/\gamma_j^* = c\} \right|, \quad (5)$$

which was stated as an “if and only if” condition of identification of $(\boldsymbol{\alpha}^*, \beta^*)$. We re-examine the identifiability and show that the “only if” part is true only when all IVs are valid. In general, no *iff* condition exists to identify model (3).

To illustrate the identification problem, we consider the model DGP. Given first-stage information: $\{\mathbf{D}, \mathbf{Z}, \boldsymbol{\gamma}^*\}$, without loss of generality, we denote the DGP with some $\{\beta^*, \boldsymbol{\alpha}^*, \boldsymbol{\epsilon}\}$ in (2) as DGP \mathcal{P}_0 that generates \mathbf{Y} . Given this \mathcal{P}_0 , for $j \in \mathcal{V}^{c*}$, we have $\mathbf{Z}_j \alpha_j^* = \frac{\alpha_j^*}{\gamma_j^*} (\mathbf{D} - \sum_{l \neq j} \mathbf{Z}_l \gamma_l^* - \boldsymbol{\eta})$. Denote $\mathcal{I}_c = \{j \in \mathcal{V}^{c*} : c = \alpha_j^*/\gamma_j^*\}$, where $c \neq 0$ and c could have up to $p_{\mathcal{V}^{c*}}$ different values. For compatibility, we denote $\mathcal{I}_0 = \mathcal{V}^*$. Thus, we can reformulate $\mathbf{Y} = \mathbf{D}\beta^* + \mathbf{Z}\boldsymbol{\alpha}^* + \boldsymbol{\epsilon}$ in (2) to:

$$\mathbf{Y} = \mathbf{D}\tilde{\beta}^c + \mathbf{Z}\tilde{\boldsymbol{\alpha}}^c + \tilde{\boldsymbol{\epsilon}}^c, \quad (6)$$

where $\{\tilde{\beta}^c, \tilde{\boldsymbol{\alpha}}^c, \tilde{\boldsymbol{\epsilon}}^c\} = \{\beta^* + c, \boldsymbol{\alpha}^* - c\boldsymbol{\gamma}^*, \boldsymbol{\epsilon} - c\boldsymbol{\eta}\}$, for some $j \in \mathcal{V}^{c*}$. Evidently, for different $c \neq 0$, it forms different DGPs $\mathcal{P}_c = \{\tilde{\beta}^c, \tilde{\boldsymbol{\alpha}}^c, \tilde{\boldsymbol{\epsilon}}^c\}$ that can generate the same \mathbf{Y} (given $\boldsymbol{\epsilon}$), which also satisfies the moment condition (3) as \mathcal{P}_0 since $E(\mathbf{Z}^\top \tilde{\boldsymbol{\epsilon}}^c) = \mathbf{0}$. Building on the argument of Guo et al. (2018), Theorem 1, the additional number of potential DGPs satisfying the moment condition (3) is the number of distinguished $c \neq 0$ for $j \in \mathcal{V}^{c*}$. We formally state this result in the following theorem.

THEOREM 1. *Suppose Assumptions 1-3 hold, given \mathcal{P}_0 and $\{\mathbf{D}, \mathbf{Z}, \boldsymbol{\gamma}^*, \boldsymbol{\eta}\}$, it can only produce additional $G = |\{c \neq 0 : \alpha_j^*/\gamma_j^* = c, j \in \mathcal{V}^{c*}\}|$ groups of different \mathcal{P}_c such that $\mathcal{V}^* \cup \{\cup_{c \neq 0} \mathcal{I}_c\} = \{1, 2, \dots, p\}$, $\mathcal{V}^* \cap \mathcal{I}_c = \emptyset$ for any $c \neq 0$ and $\mathcal{I}_c \cap \mathcal{I}_{\tilde{c}} = \emptyset$ for $c \neq \tilde{c}$, and $E(\mathbf{Z}^\top \tilde{\boldsymbol{\epsilon}}^c) = \mathbf{0}$. The sparsity structure regarding $\boldsymbol{\alpha}$ is non-overlapping for different solutions.*

Theorem 1 shows there is a collection of model DGPs

$$\mathcal{Q} = \{\mathcal{P} = \{\beta, \boldsymbol{\alpha}, \boldsymbol{\epsilon}\} : \boldsymbol{\alpha} \text{ is sparse, } E(\mathbf{Z}^\top \boldsymbol{\epsilon}) = \mathbf{0}\} \quad (7)$$

corresponding to the same observation \mathbf{Y} conditional on first-stage information. Given some \mathcal{P}_0 , there are additional $1 \leq G \leq p_{\mathcal{V}^{c*}}$ equivalent DGPs. All members in \mathcal{Q} are related through the transformation procedure (6) and $1 < |\mathcal{Q}| = G + 1 \leq p$. Notably, the non-overlapping sparse structure among all possible DGPs leads to the sparsest model regarding $\boldsymbol{\alpha}^*$ being equivalent to plurality rule $|\mathcal{V}^*| > \max_{c \neq 0} |\mathcal{I}_c|$ in the whole set of IVs. In the following, we discuss why it is not a necessary condition for identification.

DEFINITION 2. *Let \mathcal{H} be a collection of mappings $h : \mathcal{Q} \rightarrow \mathcal{P} \in \mathcal{Q}$ such that h maps a collection of DGPs \mathcal{Q} to one specific DGP $\mathcal{P} \in \mathcal{Q}$. Moreover, two different mappings h_i and h_j are image equivalent, denoted as $h_i \cong h_j$, in the sense of sharing an identical image. Also, the size of \mathcal{H} , i.e. $|\mathcal{H}|$, is defined as the number of distinct images.*

Through the above definition, clearly, $|\mathcal{H}| \leq G+1$, and any mapping $h \in \mathcal{H}$ such that $h : \mathcal{Q} \rightarrow \mathcal{P}_0$ can be treated as a sufficient condition to identify the model (2) according to \mathcal{P}_0 . The following theorem shows that, in general, no *iff* condition exists for identifying (α^*, β^*) in model (2).

THEOREM 2. *Under the same conditions as Theorem 1, $\exists i \in \{0, \dots, G\}$, let $\mathcal{F} = \{f : \mathcal{P} \in \mathcal{Q} \rightarrow \mathbb{R}; f(\mathcal{P}_i) < f(\mathcal{P}_j), \forall j \neq i\}$ and $\mathcal{G} = \{g = \operatorname{argmin}_{\mathcal{P} \in \mathcal{Q}} f(\mathcal{P}); f \in \mathcal{F}\}$, then we obtain:*

(a) $\mathcal{G} \subseteq \mathcal{H}$.

(b) *There does not exist a necessary condition for identifying (α^*, β^*) unless $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$ and $|\mathcal{H}| = 1$.*

REMARK 1. *In Theorem 2, $f \in \mathcal{F}$ maps $\mathcal{P} \in \mathcal{Q}$ on \mathbb{R} is only for convenience; in fact, it could map to any ordinal object. Moreover, $f \in \mathcal{F}$ means there exists a unique minimum value of the mapping, for all $\mathcal{P} \in \mathcal{Q}$. The elements in \mathcal{G} ($g(\mathcal{Q})$) typically correspond to some “most/least (.)” criteria, and they can serve as the sufficient conditions for identification. For instance, the plurality rule is equivalent to the most sparse α : $g(\mathcal{Q}) = \min \|\alpha\|_0$. One could consider other criteria such as the least standard deviation of the regression: $g(\mathcal{Q}) = \min \operatorname{Var}(\epsilon)$; the minimum l_2 norm of α : $g(\mathcal{Q}) = \min \|\alpha\|_2$, etc. Notably, the majority rule (Kang et al., 2016; Windmeijer et al., 2019) does not belong to \mathcal{G} in general since it can only hold for some specific \mathcal{Q} instead.*

COROLLARY 1. *Suppose the Assumptions 1-4 are satisfied. If all IVs are valid, then $|\mathcal{H}| = 1$. If there are invalid IVs, then $|\mathcal{H}| = G+1 > 1$.*

Theorem 2 describes the properties of the constituents of \mathcal{H} , while Corollary 1 states that none of the mappings in \mathcal{H} can be treated as a necessary condition unless $|\mathcal{H}| = 1$ and $h \in \mathcal{H}$ maps to \mathcal{P}_0 . The key implication of Theorem 2 (b) and Corollary 1 is that, unless all IVs are valid, the plurality rule (corresponding to $\min \|\alpha\|_0$) is not the only criteria for identification; the information of $\operatorname{Var}(\tilde{\epsilon}^c)$, e.g., can also be used as one of many criteria for identification.

In light of Theorem 2, is there a guidance for researchers to choose a proper identification condition $h \in \mathcal{H}$? Generally, seeking a proper/optimal identification condition requires a clear-defined loss function to measure how good the identification is and to perform an optimization in infinite functional space \mathcal{H} . That is infeasible to quantify. Even if we choose some common identification conditions, e.g., the ones in Remark 1, none of them can be uniformly optimal in all \mathcal{Q} . However, akin to the arguments in Guo et al. (2018); Windmeijer et al. (2021), it is often reasonable to impose the sparsest α (plurality rule) in practice because it aligns with most of the research designs to have valid IVs forming the largest group (in ideal situations all IVs would be valid). However, in practice, valid IVs tend to be weak, which is exactly the motivation of our study. As we show that individually weak IVs would not affect the model identification as long as we have some strong and valid IVs (specifically, Assumption 4 is satisfied), the sparsest rule assumption is practical and operational. Besides, sparsest α is only related to the number of valid IVs, while other criteria, $\min \|\alpha\|_2$, $\min \|\epsilon\|_2^2$, might involve more information that does not have a clear meaning in classical IV theory. Basically, in the remaining content, our target is to find all possible DGPs in \mathcal{Q} and identify the one

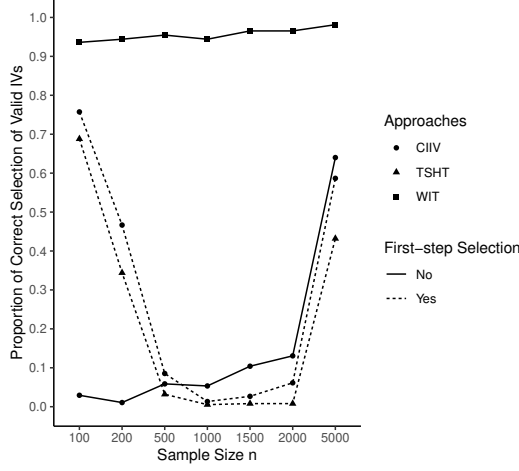


Figure 2. The proportion of correct selection of (subset) valid IVs based on 500 replications on each sample size.

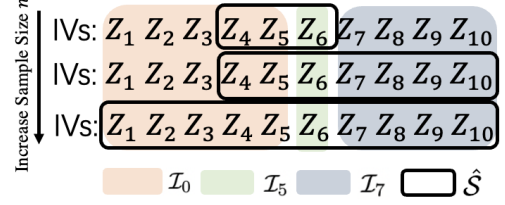


Figure 3. Illustration of Plurality rule based on first-step selection.

based on the sparsest rule. Other criteria for a true DGP may need some attention, say the reasonable value of $\text{Var}(\epsilon)$.

2.3. The Sparsest (α) Rule

The sparsest rule is conceptually equivalent to the plurality rule on the whole IV set, considering the non-overlapping sparse solutions given by Theorem 1. To relax the majority rule, Guo et al. (2018) proposed to use the plurality rule based on the relevant IV set:

$$\left| \mathcal{V}_{\mathcal{S}^*}^* = \{j \in \mathcal{S}^* : \alpha_j^* / \gamma_j^* = 0\} \right| > \max_{c \neq 0} \left| \{j \in \mathcal{S}^* : \alpha_j^* / \gamma_j^* = c\} \right|, \quad (8)$$

where \mathcal{S}^* is the strong IVs estimated by $\hat{\mathcal{S}}$ via first-step hard thresholding, i.e., $\hat{\mathcal{S}} = \{j : \hat{\gamma}_j > \sqrt{\hat{\text{Var}}(\hat{\gamma}_j) \cdot \sqrt{2.01 \log \max\{p, n\}}}\}$. Thus, TSHT and CIIV explicitly leverage on $\hat{\mathcal{S}}$ -based plurality rule to estimate $\mathcal{V}_{\hat{\mathcal{S}}}^*$ and β^* .

Unlike earlier literature on invalid IVs, our paper utilizes the information on weak IVs. For one, the weak IV can be used to estimate β^* . When we do not have strong IVs, the weak IV robust methods such as LIML are useful (Andrews et al., 2018). Second, weak IVs can be used in the identification of the valid IVs set as we show in Theorem 3. When weak IVs are present, the plurality rule applied after first-step selection may be unstable in estimating \mathcal{V}^* , as illustrated in the following example.

EXAMPLE 1 (WEAK AND INVALID IVs). Let $\gamma^* = (\mathbf{0.04}_3, \mathbf{0.5}_2, 0.2, \mathbf{0.1}_4)^\top$ and $\alpha^* = (\mathbf{0}_5, 1, \mathbf{0.7}_4)^\top$. There are therefore three groups: $\mathcal{I}_0 = \mathcal{V}^* = \{1, 2, 3, 4, 5\}$, $\mathcal{I}_5 = \{6\}$, $\mathcal{I}_7 = \{7, 8, 9, 10\}$ and plurality rule $|\mathcal{I}_0| > \max_{c=5,7} |\mathcal{I}_c|$ holds in the whole IVs set. This setup satisfies the individually weak IVs in fixed p (Discuss later in Corollary 2 and Appendix A3). Fig. 2 shows that the selection of valid IVs by CIIV and TSHT

breaks down in finite samples, e.g., for $n \in [500, 2000]$. This is because the solution of the plurality rule after first-stage selection in the finite sample (which may not hold in practice even though it is in theory) is quite sensitive to IV strength and sample sizes. On the other hand, weak IVs also deteriorate the performance of CIIV without first-step selection. Notably, the proposed WIT estimator significantly outperforms others. Fig. 3 demonstrates the relevant set \mathcal{S}^* selected by plurality rule-based TSHT and CIIV. It clearly shows $\hat{\mathcal{S}}$ is unstable and changes with sample size, even though the plurality rule holds in the whole IV set.

The mixture of weak and invalid IVs is ubiquitous in practice, especially in the many IVs case. For the sake of using all instruments' information for estimating β^* and identification of \mathcal{V}^* , we allow some individual IV strength to be local to zero (Chao and Swanson, 2005), say $\gamma_j^* \rightarrow 0$, or a small fixed constant that cannot pass the first-stage threshold (Guo et al., 2018) unless with a very large sample size. However, we can see that in (4), plurality rule-based methods that rely on first-stage selection are problematic, since $\mathcal{I}_0 = \{j : \alpha_j^*/\gamma_j^* = 0\}$ is ill-defined asymptotically due to the problem of “0/0” if γ_j^* is local to zero.

To the end of using weak IVs information and improving finite sample performance, it motivates us to turn to the sparsest rule that is also operational in computation algorithms. Back to the multiple DGPs \mathcal{Q} , recall $\mathcal{P}_c = \{\beta^c, \tilde{\alpha}^c, \tilde{\epsilon}^c\} = \{\beta^* + c, \alpha^* - c\gamma^*, \epsilon - c\eta\}$, where $\tilde{\alpha}_{\mathcal{I}_c}^c = \mathbf{0}$. For other elements in $\tilde{\alpha}^c$ (corresponding to a different DGP in \mathcal{Q}) and $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$, we obtain

$$|\tilde{\alpha}_j^c| = |\alpha_j^* - c\gamma_j^*| = |\alpha_j^*/\gamma_j^* - c| \cdot |\gamma_j^*| = |\tilde{c} - c| \cdot |\gamma_j^*|. \quad (9)$$

The above $|\tilde{\alpha}_j^c|$ needs to be distinguished from 0 on the ground of the non-overlapping structure stated in Theorem 1. To facilitate the discovery of all solutions in \mathcal{Q} , we assume:

Assumption 5: $|\tilde{\alpha}_j^c| > \kappa^c(n)$ for $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$ and $|\alpha_{\mathcal{V}_{c^*}}^*|_{\min} > \kappa(n)$, where $\kappa(n)$ and $\kappa^c(n)$ are a generally vanishing rate specified by some estimator under consideration to separate zero and non-zero terms.

The above condition is known as the “beta-min” condition (van de Geer and Bühlmann, 2009; Loh and Wainwright, 2015). Notably, as shown in (9), $|\tilde{\alpha}_j^c| = |\tilde{c} - c| \cdot |\gamma_j^*| > \kappa^c(n)$ depends on the product of $|\tilde{c} - c|$ and $|\gamma_j^*|$. As discussed in Guo et al. (2018), $|\tilde{c} - c|$ cannot be too small to separate different solutions in \mathcal{Q} , but a larger gap $|\tilde{c} - c|$ is helpful to mitigate the problem of small or local to zero $|\gamma_j^*|$ in favor of our model; see Appendix A2 for a more detailed discussion of Assumption 5.

Hence, the identification condition known as the sparsest rule is formally defined as **Assumption 6:** (The Sparsest Rule): $\alpha^* = \operatorname{argmin}_{\mathcal{P}=\{\beta, \alpha, \epsilon\} \in \mathcal{Q}} \|\alpha\|_0$.

EXAMPLE 1 (CONTINUED). Following the procedure (6), we are able to reformulate two additional solutions of (3) given the DGP of Example 1, $\alpha^* = (\mathbf{0}_5, 1, \mathbf{0.7}_4)^\top$: $\tilde{\alpha}^5 = (-\mathbf{0.2}_3, -\mathbf{2.5}_2, 0, \mathbf{0.2}_4)^\top$ and $\tilde{\alpha}^7 = (-\mathbf{0.28}_3, -\mathbf{3.5}_2, -0.4, \mathbf{0}_4)^\top$. Thus, the sparsest rule $\operatorname{argmin}_{\alpha \in \{\alpha^*, \tilde{\alpha}^5, \tilde{\alpha}^7\}} \|\alpha\|_0$ picks α^* up, and Assumption 5 is easy to satisfy since fixed minimum absolute values except 0 are 0.7, 0.2, 0.28 in $\alpha^*, \tilde{\alpha}^5, \tilde{\alpha}^7$, respectively. This example shows the first-stage signal should not interfere with the valid IV selection in the structural form equation in (2), as long as the first stage has sufficient information

(concentration parameter requirement in Assumption 4). Therefore, the most sparse rule using the whole IVs set is desirable. It is also shown to be stable in numerical studies. The detailed performance of the proposed method under this example refers to Case 1(II) in Section 4.1.

In the following subsection we revisit the penalized approaches by Kang et al. (2016) and Windmeijer et al. (2019) and discuss a class of existing estimators in terms of penalization, identification, and computation. We also discuss the general penalization approach aligning model identification with its objective function.

2.4. Penalization Approaches with Embedded Surrogate Sparsest Rule

A Lasso penalization approach was first used in unknown IV validity context by Kang et al. (2016). We extend this to a general formulation and discuss the properties of different classes of penalties.

Consider a general penalized estimator based on moment conditions (3),

$$(\hat{\alpha}^{\text{pen}}, \hat{\beta}^{\text{pen}}) = \underset{\alpha, \beta}{\operatorname{argmin}} \underbrace{\frac{1}{2n} \|P_Z(\mathbf{Y} - \mathbf{Z}\alpha - \mathbf{D}\beta)\|_2^2}_{(I)} + \underbrace{p_\lambda^{\text{pen}}(\alpha)}_{(II)}. \quad (10)$$

where $p_\lambda^{\text{pen}}(\alpha) = \sum_{j=1}^p p_\lambda^{\text{pen}}(\alpha_j)$ and $p_\lambda^{\text{pen}}(\cdot)$ is a general penalty function with tuning parameter $\lambda > 0$ and $p_\lambda^{\text{pen}'}(\cdot)$ is its derivative that satisfy: $\lim_{x \rightarrow 0^+} p_\lambda^{\text{pen}'}(x) = \lambda$, $p_\lambda^{\text{pen}}(0) = 0$, $p_\lambda^{\text{pen}}(x) = p_\lambda^{\text{pen}}(-x)$, $(x - y)(p_\lambda^{\text{pen}}(x) - p_\lambda^{\text{pen}}(y)) \geq 0$, and $p_\lambda^{\text{pen}'}(\cdot)$ is continuous on $(0, \infty)$.

In the RHS of (10), (I) and (II) correspond to two requirements for the collection of valid DGPs in \mathcal{Q} defined in (7). (I) is a scaled finite sample version of $E((\mathbf{Z}^\top \epsilon)^\top (\mathbf{Z}^\top \mathbf{Z})^{-1} (\mathbf{Z}^\top \epsilon))$, which is a $(\mathbf{Z}^\top \mathbf{Z})^{-1}$ weighted quadratic term of condition $E(\mathbf{Z}^\top \epsilon) = \mathbf{0}$, and (II) is imposed to ensure sparsity structure in $\hat{\alpha}$.

Further, regarding (I), one can reformulate (10) with respect to $\hat{\alpha}^{\text{pen}}$ as

$$\hat{\alpha}^{\text{pen}} = \underset{\alpha}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{Y} - \tilde{\mathbf{Z}}\alpha\|_2^2 + p_\lambda^{\text{pen}}(\alpha), \quad (11)$$

where $\tilde{\mathbf{Z}} = M_{\hat{\mathbf{D}}} \mathbf{Z}$ and $\hat{\mathbf{D}} = P_Z \mathbf{D} = \mathbf{Z} \hat{\gamma}$, where $\hat{\gamma}$ is the least squares estimator of γ , see Kang et al. (2016). The design matrix $\tilde{\mathbf{Z}}$ is rank-deficient with rank $p - 1$ since $\tilde{\mathbf{Z}} \hat{\gamma} = \mathbf{0}$. However, we show that it does not affect the α support recovery using a proper penalty function. On the other hand, $\tilde{\mathbf{Z}}$ is a function of η, γ^* and \mathbf{Z} , hence is correlated with ϵ . This inherited endogeneity initially stems from $\hat{\mathbf{D}}$, in which $E(\hat{\mathbf{D}}^\top \epsilon) = \sigma_{\epsilon, \eta}^2 p/n$ does not vanish in the many IVs model (Assumption 1). The following lemma implies that the level of endogeneity of each $\tilde{\mathbf{Z}}_j$ is limited.

LEMMA 1. Suppose Assumptions 1-4 hold and denote average gram matrix $\mathbf{Q}_n = \mathbf{Z}^\top \mathbf{Z}/n$. The endogeneity level of j -th transformed IV $\tilde{\mathbf{Z}}_j$ follows

$$\tilde{\mathbf{Z}}_j^\top \epsilon/n = \underbrace{\sigma_{\epsilon, \eta}^2 p/n}_{E(\hat{\mathbf{D}}^\top \epsilon/n)} \cdot \underbrace{\frac{\mathbf{Q}_{nj}^\top \gamma^*}{\gamma^{*\top} \mathbf{Q}_n \gamma^* + \sigma_\eta^2 p/n}}_{\text{dilution weight}} + O_p(n^{-1/2}). \quad (12)$$

REMARK 2. Under Assumption 1, $p/n \rightarrow v_{p_{V^*}} + v_{p_{V^{c*}}} < 1$ does not vanish as $n \rightarrow \infty$. This dilution weight is related to \mathbf{Q}_n and first-stage signal γ^* . In general the dilution weight is $o(1)$ and hence negligible except for the existence of dominated γ_j^* . However, in the fixed p case, since $p/n \rightarrow 0$, the endogeneity of $\tilde{\mathbf{Z}}$ disappears asymptotically.

Concerning (II) in (10), Theorem 1 shows that model (2) can be identified by different strategies with non-overlapping results. On the ground of the sparsest rule assumption, the role of penalty on α , i.e. $p_\lambda^{\text{pen}}(\alpha)$, should not only impose a sparsity structure but also serve as an objective function corresponding to the identification condition we choose. For example, the penalty $\lambda \|\alpha\|_0$ matches the sparsest rule.

To see the roles of a proper penalty function clearly, we rewrite (10) into an equivalent constrained objective function with the optimal penalty $\|\alpha\|_0$ regarding the sparsest rule:

$$(\hat{\alpha}^{\text{opt}}, \hat{\beta}^{\text{opt}}) = \underset{\alpha, \beta}{\operatorname{argmin}} \|\alpha\|_0 \quad \text{s.t.} \quad \|P_{\mathbf{Z}}(\mathbf{Y} - \mathbf{D}\beta - \mathbf{Z}\alpha)\|_2^2 < \delta, \quad (13)$$

where δ is the tolerance level which we specify in Section 4. The constraint above narrows the feasible solutions into \mathcal{Q} because it aligns with the Sargan test (Sargan, 1958) statistics $\|P_{\mathbf{Z}}(\mathbf{Y} - \mathbf{D}\beta - \mathbf{Z}\alpha)\|_2^2 / \|(\mathbf{Y} - \mathbf{D}\beta - \mathbf{Z}\alpha) / \sqrt{n}\|_2^2 = O_p(1)$ under null hypothesis $E(\mathbf{Z}^\top \epsilon) = \mathbf{0}$ as required in \mathcal{Q} ; otherwise, the constraint becomes $O_p(n)$ that cannot be bounded by δ . Thus, a properly chosen δ in (13) leads to an equivalent optimization problem

$$(\hat{\alpha}^{\text{opt}}, \hat{\beta}^{\text{opt}}) = \underset{\mathcal{P}=\{\beta, \alpha, \epsilon\} \in \mathcal{Q}}{\operatorname{argmin}} \|\alpha\|_0$$

that matches the identification condition. Therefore, the primary optimization object in (13) should also serve as an identification condition: the sparsest rule.

Due to computational NP-hardness for $\|\alpha\|_0$ in (13), a surrogate penalty function is needed. Kang et al. (2016) proposed to replace the optimal l_0 -norm with Lasso in (10), denoting their estimator sisVIVE as $(\hat{\alpha}^{\text{sis}}, \hat{\beta}^{\text{sis}})$. And \mathcal{V}^* is estimated as $\hat{\mathcal{V}}^{\text{sis}} = \{j : \hat{\alpha}_j^{\text{sis}} = 0\}$. However, the surrogate ℓ_1 penalty brings the following issues.

- (a) Failure in consistent variable selection under some deterministic conditions, namely the sign-aware invalid IV strength (SAIS) condition (Windmeijer et al., 2019, Proposition 2):

$$\left| \hat{\gamma}_{\mathcal{V}^{c*}}^\top \operatorname{sgn}(\alpha_{\mathcal{V}^{c*}}^*) \right| > \|\hat{\gamma}_{\mathcal{V}^*}\|_1 \quad (14)$$

The SAIS could well hold in practice, under which sisVIVE cannot achieve \mathcal{P}_0 .

- (b) Unclear dependency of regularization condition of $\tilde{\mathbf{Z}}$: Kang et al. (2016), Theorem 2, proposed a non-asymptotic error bound $|\hat{\beta}^{\text{sis}} - \beta^*|$ for sisVIVE. Under some regularity of restricted isometry property (RIP) constants of \mathbf{Z} and $P_{\hat{\mathbf{D}}}\mathbf{Z}$,

$$\|\hat{\beta}^{\text{sis}} - \beta^*\|_2 \leq \frac{|\hat{\mathbf{D}}^\top \epsilon|}{\|\hat{\mathbf{D}}\|_2^2} + \frac{1}{\|\hat{\mathbf{D}}\|_2} \left(\frac{(4/3\sqrt{5})\lambda \sqrt{p_{\mathcal{V}^*} \delta_{2p_{\mathcal{V}^*}}^+ (P_{\hat{\mathbf{D}}}\mathbf{Z})}}{2\delta_{2p_{\mathcal{V}^*}}^-(\mathbf{Z}) - \delta_{2p_{\mathcal{V}^*}}^+(\mathbf{Z}) - 2\delta_{2p_{\mathcal{V}^*}}^+(P_{\hat{\mathbf{D}}}\mathbf{Z})} \right),$$

where $\delta_k^{+/-}(\mathbf{H})$ refers to the upper and lower RIP constant of matrix \mathbf{H} . The dependence of RIP constant of $P_{\hat{\mathbf{D}}}\mathbf{Z}$ is not clear due to the randomness nature

of $\widehat{\mathbf{D}}$. Moreover, it is unclear what the impact of a potential failure of RIP is on selecting valid IVs.

- (c) The objective function deviates from the original sparsest rule: Multiple non-overlapping sparse solutions in (11) differentiate it from the standard Lasso problem, whereas the unique sparse solution satisfying (11) should share the same optimization target of l_1 and l_0 penalty. As shown in Theorem 2 and Remark 1, $g_1(\mathcal{P}) = \|\boldsymbol{\alpha}\|_0$ and $g_2(\mathcal{P}) = \|\boldsymbol{\alpha}\|_1$ correspond to incompatible identification conditions unless satisfying an additional strong requirement

$$\boldsymbol{\alpha}^* = \operatorname{argmin}_{\mathcal{P}=\{\beta, \boldsymbol{\alpha}, \boldsymbol{\epsilon}\} \in \mathcal{Q}} g_j(\mathcal{P}), \forall j = 1, 2 \iff \|\boldsymbol{\alpha}^* - c\boldsymbol{\gamma}^*\|_1 > \|\boldsymbol{\alpha}^*\|_1, \forall c \neq 0, \quad (15)$$

which further impedes sisVIVE in estimating $\beta^* \in \mathcal{P}_0$.

Among the above problems, (a) and (b) are directly linked to the Lasso problem within the framework of invalid IVs, while (c) reveals the root of the problem beyond Lasso: a proper surrogate penalty in (11) should align with the identification condition.

Windmeijer et al. (2019) proposed to use Adaptive Lasso (Zou, 2006) with a properly constructed initial estimator through median estimator to overcome the SAIS problem in (a). It also addresses (c) simultaneously. However, it requires the more stringent majority rule (see Remark 1), and all IVs are strong in the fixed p case. Furthermore, it suffers from the same sensitivity issue on weak IVs as TSHT and CIIV.

The following proposition explains what should be a proper surrogate sparsest penalty.

PROPOSITION 1. (The proper Surrogate Sparsest penalty) *Suppose Assumptions 1-6 are satisfied. If $p_\lambda^{\text{pen}}(\boldsymbol{\alpha})$ is the surrogate sparsest rule in the sense that it gives sparse solutions and*

$$\boldsymbol{\alpha}^* = \operatorname{argmin}_{\mathcal{P}=\{\beta, \boldsymbol{\alpha}, \boldsymbol{\epsilon}\} \in \mathcal{Q}} \|\boldsymbol{\alpha}\|_0 = \operatorname{argmin}_{\mathcal{P}=\{\beta, \boldsymbol{\alpha}, \boldsymbol{\epsilon}\} \in \mathcal{Q}} p_\lambda^{\text{pen}}(\boldsymbol{\alpha}), \quad (16)$$

then $p_\lambda^{\text{pen}}(\cdot)$ must be concave and $p_\lambda^{\text{pen}'}(t) = O(\lambda\kappa(n))$ for any $t > \kappa(n)$.

Such a requirement of surrogate sparsest penalty coincides with the folded-concave penalized method (Fan and Li, 2001; Zhang et al., 2010). Take MCP, the method we deploy, for example. In standard sparse linear regression, MCP requires $\lambda = \lambda(n) = O(\sqrt{\log p/n})$ and $p_\lambda^{\text{MCP}'}(t) = 0$ when $t > C\lambda$, for some constant C , which satisfies Proposition 1. We specify that such property holds in invalid IVs cases in the next section and demonstrate it also circumvents the flaws of sisVIVE shown in (a) and (b).

REMARK 3. *Proposition 1 shows the proper surrogate penalty must be concave and thus excludes Adaptive Lasso. Notice Adaptive Lasso penalty with adaptive weight constructed by a consistent estimate of $\boldsymbol{\alpha}^*$ could still satisfy (16) by adding more restrictive conditions, such as the majority rule, to identify $\boldsymbol{\alpha}^*$. The motivation of the surrogate sparsest penalty deploying concave penalties differentiates from debiasing purpose. Some other debias-oriented techniques, like the debiased Lasso (Javanmard and Montanari, 2018), cannot fit the identification condition and hence would deflect the objective function (13).*

In a nutshell, the proper surrogate sparsest penalty for (10) is to align the targeted identification condition and the global solution for $\mathcal{P}_0 \in \mathcal{Q}$.

3. WIT Estimator

3.1. Estimation Procedure

We adopt the penalized regression framework (10) and deploy a concave penalty in (11), the MCP in particular, which is a nearly unbiased estimator. Numerically, MCP penalty is shown to be the best choice in terms of the convexity of the penalized loss. Besides, it has the consistent variable selection property without adding incoherence conditions on the design matrix (Loh and Wainwright, 2017; Feng and Zhang, 2019), which suits the two-stage type of estimation problem better than Lasso. Formally, the selection stage is

$$\hat{\alpha}^{\text{MCP}} = \underset{\alpha}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{Y} - \tilde{\mathbf{Z}}\alpha\|_2^2 + p_{\lambda}^{\text{MCP}}(\alpha), \quad (17)$$

where $\tilde{\mathbf{Z}} = M_{\hat{\mathbf{D}}}\mathbf{Z}$, and $\hat{\mathbf{D}} = P_{\mathbf{Z}}\mathbf{D} = \mathbf{Z}\hat{\gamma}$ are the same as in (11). $p_{\lambda}^{\text{MCP}}(\alpha) = \sum_{j=1}^p p_{\lambda}^{\text{MCP}}(\alpha_j) = \sum_{j=1}^p \int_0^{|\alpha_j|} (\lambda - t/\rho)_+ dt$ is the MCP penalty and $\rho > 1$ is the tuning parameter, which also controls convexity level $1/\rho$, and its corresponding derivative is $p_{\lambda}^{\text{MCP}'}(t) = (\lambda - |t|/\rho)_+$. Unlike Lasso, MCP has no penalty once $|\alpha_j| > \lambda\rho$. It has nearly unbiased coefficients estimation, and achieves exact support recovery regardless of SAIS condition (see more intuitions in Appendix A4). Therefore, a consistent estimation of the valid IVs set, i.e. $\hat{\mathcal{V}} = \{j : \hat{\alpha}_j^{\text{MCP}} = 0\}$ and $\Pr(\hat{\mathcal{V}} = \mathcal{V}^*) \xrightarrow{p} 1$, is expected to hold under weaker conditions. Next, we show WIT combines the advantages of penalized TSLS and LIML estimators in different stages.

The LIML estimator is consistent not only in classic many (weak) IVs (Bekker, 1994; Hansen et al., 2008), but also in many IVs and many included covariates (Kolesár et al., 2015; Kolesár, 2018). However, simultaneous estimations in $\hat{\kappa}_{\text{liml}}$ and $\hat{\mathcal{V}}^c$ in (19) are difficult to analyze. In the selection stage (10), we use the moment-based objective function. If we do not consider the penalty term (II), the moment-based part (I) of (10) coincides with TSLS. Furthermore, the bias in TSLS has a limited effect on consistent variable selections (see Theorem 3). In the estimation step, however, due to LIML's superior finite sample performance and the issue of TSLS in the presence of many (or weak) IVs even when \mathcal{V}^* is known (Sawa, 1969; Chao and Swanson, 2005), we embed the LIML estimator to estimate β^* on the basis of estimated valid IVs set via (17). The performance of oracle-TSLS shown in simulations verifies this choice.

Consequently, we proposed the **Weak** and some **Invalid** instruments robust **Treatment** effect (WIT) estimator in the estimation stage,

$$\left(\hat{\beta}^{\text{WIT}}, \hat{\alpha}_{\mathbf{Z}_{\hat{\mathcal{V}}^c}}^{\text{WIT}}\right)^{\top} = \left(\left[\mathbf{D}, \mathbf{Z}_{\hat{\mathcal{V}}^c}\right]^{\top} (\mathbf{I} - \hat{\kappa}_{\text{liml}} M_{\mathbf{Z}}) \left[\mathbf{D}, \mathbf{Z}_{\hat{\mathcal{V}}^c}\right]\right)^{-1} \left(\left[\mathbf{D}, \mathbf{Z}_{\hat{\mathcal{V}}^c}\right]^{\top} (\mathbf{I} - \hat{\kappa}_{\text{liml}} M_{\mathbf{Z}}) \mathbf{Y}\right), \quad (18)$$

$$\hat{\kappa}_{\text{liml}} = \min_{\beta} \left\{ G(\beta) = \left((\mathbf{Y} - \mathbf{D}\beta)^{\top} M_{\mathbf{Z}} (\mathbf{Y} - \mathbf{D}\beta)\right)^{-1} \left((\mathbf{Y} - \mathbf{D}\beta)^{\top} M_{\mathbf{Z}_{\hat{\mathcal{V}}^c}} (\mathbf{Y} - \mathbf{D}\beta)\right) \right\}, \quad (19)$$

Note, (18) belongs to the general k -class estimators (Nagar, 1959), whose properties vary upon the choice of $\hat{\kappa}$, i.e. $\hat{\kappa} = 0$ refers to OLS and $\hat{\kappa} = 1$ reduces to the TSLS estimator.

(19) has a closed-form solution: $\hat{\kappa}_{\text{liml}} = \lambda_{\min} \left(\{[\mathbf{Y}, \mathbf{D}]^{\top} M_{\mathbf{Z}} [\mathbf{Y}, \mathbf{D}]\}^{-1} \{[\mathbf{Y}, \mathbf{D}]^{\top} M_{\mathbf{Z}_{\hat{\mathcal{V}}^c}} [\mathbf{Y}, \mathbf{D}]\} \right)$,

where $\lambda_{\min}(\cdot)$ means the smallest eigenvalue. Focusing on $\hat{\beta}^{\text{WIT}}$ as the primary interest, we reformulate (18) and (19) based on the residuals of $\mathbf{Y}, \mathbf{D}, \mathbf{Z}_{\hat{\mathcal{V}}}$ on $\mathbf{Z}_{\hat{\mathcal{V}}^c}$. Denote

$\mathbf{Y}_\perp = M_{\mathbf{Z}_{\hat{\mathcal{V}}^c}} \mathbf{Y}$, $\mathbf{D}_\perp = M_{\mathbf{Z}_{\hat{\mathcal{V}}^c}} \mathbf{D}$ and $\mathbf{Z}_{\hat{\mathcal{V}}_\perp} = M_{\mathbf{Z}_{\hat{\mathcal{V}}^c}} \mathbf{Z}_{\hat{\mathcal{V}}}$ and notice $M_{\mathbf{Z}_{\hat{\mathcal{V}}^c}} M_{\mathbf{Z}_{\hat{\mathcal{V}}_\perp}} = M_{\mathbf{Z}}$, thus it is equivalent to derive asymptotic results on the following model (20).

$$\hat{\beta}^{\text{WIT}} = \left(\mathbf{D}_\perp^\top (I - \hat{\kappa}_{\text{liml}} M_{\mathbf{Z}_{\hat{\mathcal{V}}_\perp}}) \mathbf{D}_\perp \right)^{-1} \left(\mathbf{D}_\perp^\top (I - \hat{\kappa}_{\text{liml}} M_{\mathbf{Z}_{\hat{\mathcal{V}}_\perp}}) \mathbf{Y}_\perp \right), \quad (20)$$

$$\hat{\kappa}_{\text{liml}} = \lambda_{\min} \left(\{[\mathbf{Y}_\perp, \mathbf{D}_\perp]^\top M_{\mathbf{Z}_{\hat{\mathcal{V}}_\perp}} [\mathbf{Y}_\perp, \mathbf{D}_\perp]\}^{-1} \{[\mathbf{Y}_\perp, \mathbf{D}_\perp]^\top [\mathbf{Y}_\perp, \mathbf{D}_\perp]\} \right). \quad (21)$$

3.2. Asymptotic Behavior of WIT Estimator

Throughout this section, we aim to recover the one specific element in \mathcal{Q} , denoted as $(\beta^*, \alpha^*, \epsilon)$ temporally. Though a slight abuse of notation, we use $\hat{\alpha}$ to denote a local solution of (17) with MCP.

All local solutions of (17) we consider are characterized by the Karush–Kuhn–Tucker (KKT) or first-order condition, i.e.

$$\tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\alpha})/n = \frac{\partial}{\partial \mathbf{t}} \sum_{j=1}^p p'_\lambda(t_j) \Big|_{\mathbf{t}=\hat{\alpha}}. \quad (22)$$

Explicitly, to the end of finding valid IVs via comparing with true signal α^* , we rewrite (22) as

$$\begin{cases} \left(\lambda - \frac{1}{\rho} |\hat{\alpha}_j| \right)_+ \leq \text{sgn}(\hat{\alpha}_j) \tilde{\mathbf{Z}}_j^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\alpha})/n \leq \lambda, & j \in \hat{\mathcal{V}}^c \\ \left| \tilde{\mathbf{Z}}_j^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\alpha})/n \right| \leq \lambda, & j \in \hat{\mathcal{V}} \end{cases} \quad (23)$$

where the inequalities in the first line stem from the convexity of MCP penalty and in the last line originate in the sub-derivative of the MCP penalty at the origin.

As discussed in Section 2, $\tilde{\mathbf{Z}}$ is a function of $(\mathbf{Z}, \gamma^*, \boldsymbol{\eta})$ and thus endogenous with ϵ . The fact that $\tilde{\mathbf{Z}}$ inherits the randomness of $\boldsymbol{\eta}$ distinguishes itself from the general assumptions put on the design matrix, and obscures the feasibility of conditions required to achieve exact support recovery in the literature of penalized least squares estimator (PLSE) (Feng and Zhang, 2019; Loh and Wainwright, 2017; Zhang and Zhang, 2012). The sisVIVE method imposed the RIP condition directly on $\tilde{\mathbf{Z}}$ to establish an error bound. However, the restricted eigenvalue (RE) condition (Bickel et al., 2009) is the weakest condition (van de Geer and Bühlmann, 2009) available to guarantee rate minimax performance in prediction and coefficient estimation, as well as to establish variable selection consistency for Lasso penalty. Feng and Zhang (2019) further adopted the RE condition for non-convex penalty analysis. We then state the conditions on the design matrix $\tilde{\mathbf{Z}}$ of (17) in the following. Define restricted cone $\mathcal{C}(\mathcal{V}^*; \xi) = \{\mathbf{u} : \|\mathbf{u}_{\mathcal{V}^*}\|_1 \leq \xi \|\mathbf{u}_{\mathcal{V}^{c*}}\|_1\}$ for some $\xi > 0$ that estimation error $\hat{\alpha} - \alpha^*$ belongs to. The restricted eigenvalue $K_{\mathcal{C}}$ for $\tilde{\mathbf{Z}}$ is defined as $K_{\mathcal{C}} = K_{\mathcal{C}}(\mathcal{V}^*, \xi) := \inf_{\mathbf{u}} \{\|\tilde{\mathbf{Z}} \mathbf{u}\|_2 / (\|\mathbf{u}\|_2 n^{1/2}) : \mathbf{u} \in \mathcal{C}(\mathcal{V}^*; \xi)\}$ and the RE condition refers to the condition that $K_{\mathcal{C}}$ for $\tilde{\mathbf{Z}}$ should be bounded away from zero.

LEMMA 2. (RE condition of $\tilde{\mathbf{Z}}$) Under Assumptions 1-3, for any given $\gamma^* \neq \mathbf{0}$, there exists a constant $\xi \in (0, \|\hat{\gamma}_{\mathcal{V}^*}\|_1 / \|\hat{\gamma}_{\mathcal{V}^{c*}}\|_1)$ and further, a restricted cone $\mathcal{C}(\mathcal{V}^*; \xi)$ defined by chosen ξ such that $K_{\mathcal{C}}^2 > 0$ holds strictly.

Lemma 2 elaborates that the RE condition on $\tilde{\mathbf{Z}}$ holds without any additional assumptions on $\tilde{\mathbf{Z}}$, unlike the extra RIP condition for sisVIVE. Moreover, this restricted cone is invariant of scaling, thus indicating the accommodation of many weak IVs. These two features suggest the theoretical advantages of penalized methods (11) over existing methods.

Next, we discuss the selection of valid IVs by comparing the local solution of (23) with the oracle (moment-based) counterpart. Define $\hat{\boldsymbol{\alpha}}_{\mathcal{V}^*}^{\text{or}} = \mathbf{0}$ and

$$\hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}} = (\tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}})^{-1} \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \mathbf{Y} \quad \text{or} \quad \hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}} = (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} [\mathbf{Z}_{\mathcal{V}^{c*}}^\top (\mathbf{Y} - \hat{\mathbf{D}} \hat{\beta}_{\text{or}}^{\text{TSL}})], \quad (24)$$

where $\hat{\beta}_{\text{or}}^{\text{TSL}} = [\mathbf{D}^\top (\mathbf{P}_{\mathbf{Z}} - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D}]^{-1} [\mathbf{D}^\top (\mathbf{P}_{\mathbf{Z}} - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{Y}]$ and the two versions of $\hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}}$ are equivalent. Notice this $\hat{\beta}_{\text{or}}^{\text{TSL}}$ is not for the final treatment effect estimation, but to illustrate the selection stage consistency only. To this end, we show the supremum norm of $\mathbf{R}^{\text{or}} = \tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\alpha}}^{\text{or}})/n$ is bounded by an inflated order of $O(\sqrt{\log p_{\mathcal{V}^*}/n})$. Denote $\tilde{\mathbf{D}} = (\mathbf{P}_{\mathbf{Z}} - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D}$ and $\tilde{\mathbf{Q}}_n = \mathbf{Z}_{\mathcal{V}^*}^\top (\mathbf{P}_{\mathbf{Z}} - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{Z}_{\mathcal{V}^*}/n$, we derive the following lemma.

LEMMA 3. *Suppose Assumptions 1-4 hold and let*

$$\zeta \asymp \frac{p_{\mathcal{V}^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n \boldsymbol{\gamma}_{\mathcal{V}^*}^*\|_\infty}{\boldsymbol{\gamma}_{\mathcal{V}^*}^{*\top} \tilde{\mathbf{Q}}_n \boldsymbol{\gamma}_{\mathcal{V}^*}^*} + \sqrt{\frac{\log p_{\mathcal{V}^*}}{n}}. \quad (25)$$

Then, the supremum norms of residual \mathbf{R}^{or} are bounded by ζ , i.e.,

$$\|\mathbf{R}^{\text{or}}\|_\infty \leq \left\| \frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\boldsymbol{\epsilon}}}{n} \right\|_\infty + \left\| \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n}}{\frac{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}}{n}} \right\|_\infty \leq \zeta \quad (26)$$

holds with probability approaching 1.

Based on Lemmas 2 and 3, we consider the set $\mathcal{B}(\lambda, \rho) = \{\hat{\boldsymbol{\alpha}} \text{ in (23)} : \lambda \geq \zeta, \rho > K_{\mathcal{C}}^{-2}(\mathcal{V}^*, \xi) \vee 1\}$, in which ζ is defined in (25) and ξ is guaranteed by Lemma 2, as a collection of all local solutions $\hat{\boldsymbol{\alpha}}$ computed in (23) through a broad class of MCP under a certain penalty level λ and convexity $1/\rho$. Given that the computed local solutions in practice are through a discrete path with some starting point (see Section 3.3), we further consider the computable solution set $\mathcal{B}_0(\lambda, \rho)$, introduced by Feng and Zhang (2019), i.e.,

$$\mathcal{B}_0(\lambda, \rho) = \{\hat{\boldsymbol{\alpha}} : \hat{\boldsymbol{\alpha}} \text{ and starting point } \hat{\boldsymbol{\alpha}}^{(0)} \text{ are connected in } \mathcal{B}(\lambda, \rho)\}. \quad (27)$$

The connection between $\mathcal{B}_0(\lambda, \rho)$ and $\mathcal{B}(\lambda, \rho)$ is that $\exists \hat{\boldsymbol{\alpha}}^{(l)} \in \mathcal{B}(\lambda, \rho)$ with penalty level $\lambda^{(l)}$ increasing with the index $l = 1, 2, \dots$, such that $\hat{\boldsymbol{\alpha}}^{(0)} - \boldsymbol{\alpha}^* \in \mathcal{C}(\mathcal{V}^*, \xi)$, $\hat{\boldsymbol{\alpha}} = \hat{\boldsymbol{\alpha}}^{(l)}$ for large enough l and $\|\hat{\boldsymbol{\alpha}}^{(l)} - \hat{\boldsymbol{\alpha}}^{(l-1)}\|_1 < a_0 \lambda^{(l)}$, where a_0 is specified in Lemma A3 of Appendix B9. Thus, $\mathcal{B}_0(\lambda, \rho)$ is a collection of approximations of $\boldsymbol{\alpha}$ in all DGPs.

Denote $\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSL}}) = \frac{\mathbf{D}^\top (\mathbf{P}_{\mathbf{Z}} - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \boldsymbol{\epsilon}}{\mathbf{D}^\top (\mathbf{P}_{\mathbf{Z}} - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D}} = \frac{\mathbf{D}_\perp^\top \mathbf{P}_{\mathbf{Z}_\perp} \boldsymbol{\epsilon}_\perp}{\mathbf{D}_\perp^\top \mathbf{P}_{\mathbf{Z}_\perp} \mathbf{D}_\perp}$ and $\tilde{\boldsymbol{\gamma}}_{\mathcal{V}^{c*}}^* = \boldsymbol{\gamma}_{\mathcal{V}^{c*}}^* + (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} \mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^*} \boldsymbol{\gamma}_{\mathcal{V}^*}^*$.

Also let $\tilde{\mathbf{Q}}_n^c$ and $\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSL}})$ be defined as \mathcal{P}_c version of $\tilde{\mathbf{Q}}_n$ and $\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSL}})$. Then, we provide the asymptotic result of selection consistency of WIT.

THEOREM 3. (Selection Consistency) *Specify $\kappa(n)$ and $\kappa^c(n)$ in Assumption 5 as*

$$\kappa(n) \asymp \underbrace{\sqrt{\frac{\log p_{\mathcal{V}^*}}{n}}}_{T_1} + \underbrace{\frac{p_{\mathcal{V}^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*\|_\infty}{\gamma_{\mathcal{V}^*}^{*\top} \tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*}}_{T_2} + \underbrace{|\text{Bias}(\hat{\beta}_{or}^{TSLS})| \|\tilde{\gamma}_{\mathcal{V}^*}^*\|_\infty}_{T_3}, \quad (28)$$

$$\kappa^c(n) \asymp (1+c) \left\{ \sqrt{\frac{\log |\mathcal{I}_c|}{n}} + \frac{|\mathcal{I}_c|}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n^c \gamma_{\mathcal{I}_c}^*\|_\infty}{\gamma_{\mathcal{I}_c}^{*\top} \tilde{\mathbf{Q}}_n^c \gamma_{\mathcal{I}_c}^*} \right\} + |\text{Bias}(\hat{\beta}_{or}^{c TSLS})| \|\tilde{\gamma}_{\mathcal{I}_c}^*\|_\infty, \quad (29)$$

where $T_1 \rightarrow 0$ as $n \rightarrow \infty$. Suppose Assumptions 1-6 hold and consider computable local solutions specified in (27), we have

$$\hat{\alpha}^{MCP} = \underset{\hat{\alpha} \in \mathcal{B}_0(\lambda, \rho)}{\text{argmin}} \|\hat{\alpha}\|_0, \Pr(\hat{\mathcal{V}} = \mathcal{V}^*, \hat{\alpha}^{MCP} = \hat{\alpha}^{or}) \xrightarrow{p} 1. \quad (30)$$

In Theorem 3, T_1 is similar to a standard rate $\sqrt{\log p/n}$ in penalized linear regression, while T_2 and T_3 are the additional terms that only happen to many IVs context and vanish fast in the finite strong IVs case (see Corollary 2). This result is new to the literature. We further provide justification of Theorem 3 in commonly considered cases (fixed p or diverging p with many individually weak IVs) in Appendix A3.

PROPOSITION 2. *Under the same assumptions of Theorem 3, if there does not exist dominant scaled γ_j^* , i.e. $\|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_\infty / \|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_1 = o(\|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_1 / (p_{\mathcal{V}^*} \|\tilde{\mathbf{Q}}_n^{1/2}\|_\infty))$, then $T_2 \rightarrow 0$.*

Proposition 2 shows that T_2 is limited in the general case where dominant scaled γ_j^* does not exist. For example, if we assume $\mathbf{Q}_n = \mathbf{I}$ and $\gamma_{\mathcal{V}^*}^* = C \mathbf{1}_{p_{\mathcal{V}^*}}$, where C is a constant or diminishing to zero, then $\|\gamma_{\mathcal{V}^*}^*\|_\infty / \|\gamma_{\mathcal{V}^*}^*\|_1 = o(\|\gamma_{\mathcal{V}^*}^*\|_1 / p_{\mathcal{V}^*}) = o(C p_{\mathcal{V}^*} / p_{\mathcal{V}^*}) = o(C)$ holds and it leads $T_2 \rightarrow 0$.

PROPOSITION 3. (Approximation of Bias($\hat{\beta}_{or}^{TSLS}$)) *Let $s = \max(\mu_n, p_{\mathcal{V}^*})$, under the Assumptions 1-4, we obtain*

$$E \left[\text{Bias}(\hat{\beta}_{or}^{TSLS}) \right] = \frac{\sigma_{\epsilon\eta}}{\sigma_\eta^2} \left(\frac{p_{\mathcal{V}^*}}{(\mu_n + p_{\mathcal{V}^*})} - \frac{2\mu_n^2}{(\mu_n + p_{\mathcal{V}^*})^3} \right) + o(s^{-1}). \quad (31)$$

REMARK 4. *The rate of concentration parameter μ_n will affect T_3 through $|\text{Bias}(\hat{\beta}_{or}^{TSLS})|$ under many IVs setting. Suppose Assumption 4 holds, that $\mu_n \xrightarrow{p} \mu_0 n$, the leading term in (31) is $\frac{\sigma_{\epsilon\eta}}{\sigma_\eta^2} \frac{\nu_{p_{\mathcal{V}^*}}}{\mu_0 + \nu_{p_{\mathcal{V}^*}}} \ll \frac{\sigma_{\epsilon\eta}}{\sigma_\eta^2}$ for moderate μ_0 since $0 < \nu_{p_{\mathcal{V}^*}} < 1$ while μ_0 could be larger than 1. While under many weak IVs setting (Chao and Swanson, 2005; Hansen et al., 2008; Newey and Windmeijer, 2009), $\mu_n/n \xrightarrow{p} 0$ and the leading term in (31) becomes $\sigma_{\epsilon\eta}/\sigma_\eta^2$. Thus, many weak IVs setting imposes some difficulty (a higher T_3) for selecting valid IVs in Theorem 3. We further discuss the eligibility of many weak IVs in Remark 5.*

The following theorem describes the asymptotic behavior of the WIT estimator for many valid and invalid IVs cases by combining Theorem 1 and invariant likelihood arguments in [Kolesár \(2018\)](#). We further denote two statistics

$$\mathbf{S} = \frac{1}{n-p}(\mathbf{Y}, \mathbf{D})^\top M_{\mathbf{Z}}(\mathbf{Y}, \mathbf{D}), \quad \mathbf{T} = \frac{1}{n}(\mathbf{Y}_\perp, \mathbf{D}_\perp)^\top M_{\mathbf{Z}_\perp}(\mathbf{Y}_\perp, \mathbf{D}_\perp) \quad (32)$$

as the estimates of the covariance matrix of reduced-form error $\boldsymbol{\Omega} = \text{Cov}(\boldsymbol{\epsilon} + \beta^* \boldsymbol{\eta}, \boldsymbol{\eta})$ and a variant of concentration parameter, respectively. Also, let $m_{\max} = \lambda_{\max}(\mathbf{S}^{-1} \mathbf{T})$, $\hat{\mu}_n = \max(m_{\max} - p_{\mathcal{V}^*}/n, 0)$ and $\hat{\boldsymbol{\Omega}} = \frac{n-p}{n-p_{\mathcal{V}^*}/n} \mathbf{S} + \frac{n}{n-p_{\mathcal{V}^*}/n} (\mathbf{T} - \frac{\hat{\mu}_n}{\hat{\mathbf{a}}^\top \mathbf{S}^{-1} \hat{\mathbf{a}}} \hat{\mathbf{a}} \hat{\mathbf{a}}^\top)$, where $\hat{\mathbf{a}} = (\hat{\beta}^{\text{WIT}}, 1)$ and $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$.

THEOREM 4. *Under the same conditions as in Theorem 3, we obtain:*

- (a) (Consistency): $\hat{\beta}^{\text{WIT}} \xrightarrow{p} \beta^*$ with $\hat{\kappa}_{\text{liml}} = \frac{1-v_{p_{\mathcal{V}^*}}}{1-v_{p_{\mathcal{V}^*}}-v_{p_{\mathcal{V}^*}}} + o_p(1)$.
- (b) (Asymptotic normality): $\sqrt{n}(\hat{\beta}^{\text{WIT}} - \beta^*) \xrightarrow{d} \mathcal{N}\left(0, \mu_0^{-2} [\sigma_\epsilon^2 \mu_0 + \frac{v_{p_{\mathcal{V}^*}}(1-v_{p_{\mathcal{V}^*}})}{1-v_{p_{\mathcal{V}^*}}-v_{p_{\mathcal{V}^*}}} |\boldsymbol{\Sigma}|]\right)$.
- (c) (Consistent variance estimator):

$$\begin{aligned} \widehat{\text{Var}}(\hat{\beta}^{\text{WIT}}) &= \frac{\hat{\mathbf{b}}^\top \hat{\boldsymbol{\Omega}} \hat{\mathbf{b}}(\hat{\mu}_n + p_{\mathcal{V}^*}/n)}{-\hat{\mu}_n} \left(\hat{Q}_S \hat{\boldsymbol{\Omega}}_{22} - \mathbf{T}_{22} + \frac{\hat{c}}{1 - \hat{c} \hat{\mathbf{a}}^\top \hat{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{a}}} \right)^{-1} \\ &\xrightarrow{p} \mu_0^{-2} \left[\sigma_\epsilon^2 \mu_0 + \frac{v_{p_{\mathcal{V}^*}}(1-v_{p_{\mathcal{V}^*}})}{1-v_{p_{\mathcal{V}^*}}-v_{p_{\mathcal{V}^*}}} |\boldsymbol{\Sigma}| \right], \end{aligned}$$

where $\hat{\mathbf{b}} = (1, -\hat{\beta}^{\text{WIT}})$ and $\hat{Q}_S = \frac{\hat{\mathbf{b}}^\top \mathbf{T} \hat{\mathbf{b}}}{\hat{\mathbf{b}}^\top \hat{\boldsymbol{\Omega}} \hat{\mathbf{b}}}$.

Notably, when the number of invalid IVs $p_{\mathcal{V}^*}$ is a constant, the variance estimator above is reduced to the one that [Bekker \(1994\)](#) derived for typical many IVs case. [Hansen et al. \(2008\)](#) showed that it is still valid for many weak IVs asymptotics.

REMARK 5. (Many Weak IVs Asymptotics) *Regarding the many weak IVs asymptotic sequence considered in [Chao and Swanson \(2005\)](#); [Hansen et al. \(2008\)](#); [Newey and Windmeijer \(2009\)](#), we need to modify the number of invalid IVs $p_{\mathcal{V}^*}$ as fixed. Thus, Assumption 4 can be relaxed to $\mu_n/n \rightarrow 0$ but $\mu_n/\sqrt{n} \rightarrow \infty$ and Theorem 4 holds according to [Hansen et al. \(2008\)](#), Theorems 1 and 3.*

In the following, we show WIT is a more powerful tool than the existing methods requiring the majority rule ([Kang et al., 2016](#); [Windmeijer et al., 2019](#)) even under the finite IVs with a mixture of strong and weak settings. WIT achieves the same asymptotic results as [Windmeijer et al. \(2019\)](#) under more relaxed conditions. Specifically, under finite IVs, Assumption 1 can be reduced to Assumption 1' as follows.

Assumption 1' (Finite Number of IVs): $p_{\mathcal{V}^*} \geq 1$ and $p_{\mathcal{V}^*} \geq 1$ are fixed constants, and $p_{\mathcal{V}^*} + p_{\mathcal{V}^*} = p < n$.

In the finite IVs case, the T_2 and T_3 terms in Theorem 3 for selecting valid IVs go to zero fast and the required vanishing rate reduces to $\kappa(n) \asymp n^{-1/2}$. We present the asymptotic properties for the WIT estimator in the following. Consider the following IV signals mixture case. Let $\gamma_j^* = Cn^{-\tau_j}$ for $j = 1, 2, \dots, p$. Thus, let $\tau_{\mathcal{V}^*} = \arg\max_{\tau_j} \{\tau_j : j \in \mathcal{V}^*\}$, $\tau_{\mathcal{I}_c} = \arg\max_{\tau_j} \{\tau_j : j \in \mathcal{I}_c\}$ and $\tau_{\mathcal{I}_{\tilde{c}}} = \arg\max_{\tau_j} \{\tau_j : j \in \mathcal{I}_{\tilde{c}}\}$, where $c \neq \tilde{c}$.

COROLLARY 2. (Finite p with Mixture of Strong and Weak IVs) Suppose Assumptions 1', 2-4 and 6 hold. Additionally, we assume each γ_j^* is at least a weak IV such that $\gamma_j^* = O(n^{-\tau})$ and $0 \leq \tau \leq 1/2$ for $j = 1, 2, \dots, p$. For any fixed $\min_{j \in \mathcal{V}^*} \alpha_j^* > 0$, if $\tau_{\mathcal{V}^*} + 2\tau_{\mathcal{I}_c} < 1$, $\tau_{\mathcal{V}^*} + 2\tau_{\mathcal{I}_e} < 1$ and $\tau_{\mathcal{I}_c} + \tau_{\mathcal{I}_e} < 2/3$, then we have

(a) (Selection consistency): $\hat{\alpha}^{MCP} = \underset{\hat{\alpha} \in \mathcal{B}_0(\lambda, \rho)}{\operatorname{argmin}} \|\hat{\alpha}\|_0, \Pr(\hat{\mathcal{V}} = \mathcal{V}^*, \hat{\alpha}^{MCP} = \hat{\alpha}^{or}) \xrightarrow{p} 1.$

(b) (Consistency & Equivalence of WIT and TSLS): $\hat{\beta}^{WIT} \xrightarrow{p} \beta^*$ with $\hat{\kappa}_{liml} = 1 + o_p(1).$

(c) (Asymptotic normality): $\sqrt{n}(\hat{\beta}^{WIT} - \beta^*) \xrightarrow{d} \mathcal{N}(0, \mu_0^{-1} \sigma_\epsilon^2).$

(d) (Consistent variance estimator):

$$\widehat{\operatorname{Var}}(\hat{\beta}^{WIT}) = \frac{\hat{\mathbf{b}}^\top \hat{\Omega} \hat{\mathbf{b}} (\hat{\mu}_n + p\mathcal{V}^*/n)}{-\hat{\mu}_n} \left(\hat{Q}_S \hat{\Omega}_{22} - \mathbf{T}_{22} + \frac{\hat{c}}{1 - \hat{c}} \frac{\hat{Q}_S}{\hat{\mathbf{a}}^\top \hat{\Omega}^{-1} \hat{\mathbf{a}}} \right)^{-1} \xrightarrow{p} \mu_0^{-1} \sigma_\epsilon^2,$$

where $\hat{\mathbf{b}} = (1, -\hat{\beta}^{WIT})$ and $\hat{Q}_S = \frac{\hat{\mathbf{b}}^\top \mathbf{T} \hat{\mathbf{b}}}{\hat{\mathbf{b}}^\top \hat{\Omega} \hat{\mathbf{b}}}.$

3.3. Computational Implementation of WIT Estimator

Through Proposition 1, we know that MCP belongs to the surrogate sparsest penalty under Assumptions 1-6, and it ensures the global solution in (17) matches the sparsest rule. However, each element in \mathcal{Q} is the local solution of (17) and we can only obtain one local solution from one initial value practically. A multiple starting points strategy is needed to achieve the global solution. Enumerating the whole $\alpha^* \in \mathbb{R}^p$ is impossible. Therefore, it is important to develop an efficient algorithm for MCP penalty.

In light of practical use, we adopt the iterative local adaptive majorize-minimization (I-LAMM) algorithm (Fan et al., 2018), which satisfies (27) as shown in Feng and Zhang (2019) Section 2.1, with different initial values to achieve the local solution of $\hat{\alpha}$ in (17). See more technical details and derivation in Appendix A5.

Motivated by the individual IV estimator (Windmeijer et al., 2021) such that $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j \xrightarrow{p} \beta^* + \alpha_j^* / \gamma_j^*$, where $\hat{\Gamma} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}$ and $\hat{\gamma} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{D}$, we can construct

$$\check{\alpha}(\check{\beta}) = \hat{\Gamma} - \check{\beta} \hat{\gamma} = \hat{\Gamma} - \beta^* \hat{\gamma} + (\beta^* - \check{\beta}) \hat{\gamma}, \quad (33)$$

with any initial value $\check{\beta}$. Thus $\|\check{\alpha}(\check{\beta}) - \alpha^*\|_1 \xrightarrow{p} |\beta^* - \check{\beta}| \cdot \|\hat{\gamma}\|_1$. It is valid to replace (β^*, α^*) by $(\check{\beta}^c, \check{\alpha}^c)$, we know that $\|\check{\alpha}(\check{\beta}) - \check{\alpha}^c\|_1$ is asymptotically controlled by $|\check{\beta}^c - \check{\beta}| \cdot \|\hat{\gamma}\|_1$. Thus, varying $\check{\beta} \in \mathbb{R}^1$ is equivalent to varying $\check{\alpha}(\check{\beta}) \in \mathbb{R}^p$ of a close DGP \mathcal{Q} .

Here we provide a simple heuristic procedure to efficiently choose a proper $\check{\beta}$. Let $\hat{\beta}_{[j]}$ be the sorted $\hat{\beta}_j$. A fuzzy MCP regression is conducted to clustering $\hat{\beta}_j$:

$$\bar{\beta} = \underset{\beta}{\operatorname{argmin}} \sum_{j=1}^p \|\beta_j - \hat{\beta}_{[j]}\|_2^2 + \sum_{j=2}^p p_{\bar{\lambda}}^{MCP} (|\beta_j - \beta_{j-1}|), \quad (34)$$

where $\bar{\beta} \in \mathbb{R}^p$, $\bar{\lambda}$ is a prespecified penalty level and $\bar{\beta}$ is initialized by $\bar{\beta}_j = \hat{\beta}_{[j]}$. Therefore, $\bar{\beta}$ consists of less than p distinct values. Thus we can choose the initial $\hat{\alpha}$ in (27) such that

$$\hat{\alpha}^{(0)}(\check{\beta}) = \check{\alpha}(\check{\beta}), \quad (35)$$

where $\check{\beta}$'s are chosen in $\bar{\beta}$ with the priority given to the largest (remaining) cluster and $\hat{\alpha}_j^{(0)}(\check{\beta}) = 0$ for j to be the unsorted index of $\hat{\beta}_j$ such that $\check{\beta} = \bar{\beta}_j$.

Following (27), which provides a theoretical guideline for the tuning parameter, we look for a data-driven tuning procedure that has good performance in practice. From numerical studies, results are not sensitive to the choice of ρ , which is fixed to 2 for most applications. However, λ is required to be tuned. Cross-validation is implemented in sisVIVE, but is known to be time-consuming and select too many IVs as invalid. Windmeijer et al. (2019, 2021) proposed to use the Sargan test under low dimensions to choose the tuning parameter consistently and obtain a good finite sample performance. However, the Sargan test is designed for fixed p and cannot handle many IVs. Therefore, we propose the modified Cragg-Donald (MCD, Kolesár, 2018) test-based tuning procedure, which extends the Sargan test to allow high-dimensional covariates and IVs.

Specifically, consider a local solution $\hat{\alpha}$ in (27). Denote $p_{\hat{\gamma}} = |\{j : \hat{\alpha}_j = 0\}|$ and $p_{\hat{\gamma}^c} = |\{j : \hat{\alpha}_j \neq 0\}|$. Let m_{\min} be the minimum eigenvalue of $\mathbf{S}^{-1}\mathbf{T}$, where \mathbf{S} and \mathbf{T} are defined as a function of $\hat{\alpha}$ in (32). Then, the MCD test is given by nm_{\min} . According to Kolesár (2018) Proposition 4, the MCD test with asymptotic size ϱ_n would reject the null of $\alpha_{\hat{\gamma}} = \mathbf{0}$, when

$$nm_{\min} > \chi_{p_{\hat{\gamma}^c}-1}^2 \left\{ \Phi \left(\sqrt{\frac{n - p_{\hat{\gamma}^c}}{n - p_{\hat{\gamma}^c} - p_{\hat{\gamma}}}} \Phi^{-1}(\varrho_n) \right) \right\}, \quad (36)$$

where $\Phi(\cdot)$ denotes the CDF of the standard normal distribution and $\chi_{p_{\hat{\gamma}^c}-1}^2(\varrho_n)$ is $1 - \varrho_n$ quantile of $\chi_{p_{\hat{\gamma}^c}-1}^2$ distribution. This property holds regardless of whether $p_{\mathcal{V}^{c*}}$ is fixed or grows with n . For the sake of the model selection consistency, the size of the MCD test needs to be $o(1)$. Following Belloni et al. (2012); Windmeijer et al. (2019, 2021), we adopt a scaled rate $\varrho_n = 0.5/\log n$ that works well in simulation studies. Thus, with $\hat{\alpha}^{(0)}$ in (35) and a sequence of $\lambda^{\text{seq}} = \mathbf{C}\sqrt{\log p/n}$ where $\mathbf{C} = 0.1 \times (1, \dots, 20)^\top$, we propose to use the MCD test to select the proper $\lambda \in \lambda^{\text{seq}}$ that is not rejected by (36) with size ϱ_n and largest $p_{\hat{\gamma}}$.

To sum up, we provide Algorithm 1 to demonstrate the detailed implementation.

4. Numerical Simulations

In this section, we conduct numerical studies to evaluate the finite sample performance of the proposed WIT estimator. In the design of the simulation experiments, we consider scenarios corresponding to different empirically relevant problems.

We consider the same model in Section 2,

$$\mathbf{Y} = \mathbf{D}\beta^* + \mathbf{Z}\alpha^* + \epsilon, \quad \mathbf{D} = \mathbf{Z}\gamma^* + \eta.$$

Throughout all settings, we fix true treatment effect $\beta^* = 1$. \mathbf{Z} is the $n \times p$ potential IV matrix and $\mathbf{Z}_i \stackrel{i.i.d.}{\sim} \mathcal{N}(\mathbf{0}, \Sigma^{\mathbf{Z}})$, where $\Sigma_{jj}^{\mathbf{Z}} = 0.3$ and $\Sigma_{jk}^{\mathbf{Z}} = 0.3|j - k|^{0.8}$ for $i = 1, \dots, n$ and $k, j = 1, \dots, p$. Denote $\epsilon = (\epsilon_1, \epsilon_2, \dots, \epsilon_n)^\top$ and $\eta = (\eta_1, \eta_2, \dots, \eta_n)^\top$ and generate $(\epsilon_i, \eta_i)^\top \stackrel{i.i.d.}{\sim} \mathcal{N}\left(\mathbf{0}, \begin{pmatrix} \sigma_\epsilon^2 & \sigma_{\epsilon, \eta} \\ \sigma_{\epsilon, \eta} & \sigma_\eta^2 \end{pmatrix}\right)$. We let $\sigma_\epsilon^2 = 1$ and $\text{corr}(\epsilon_i, \eta_i) = 0.6$ in all settings but vary σ_η^2 to get different concentration parameters concerning strong or weak IVs cases.

Algorithm 1 WIT estimator with MCD test tuning strategy

Input: $Y, Z, D, \lambda^{\text{seq}}, \varrho_n$, and J

- 1: Calculate $\bar{\beta}$ in (34), initialize $\hat{\alpha}^{\text{MCP}} = \mathbf{1}$ and $\mathbf{I} = \mathbf{1}$
- 2: **for** $\hat{\alpha}^{(0)} = \mathbf{0}, \hat{\alpha}^{(0)}(\hat{\beta}_j)$ (35) **do** $\triangleright \hat{\beta}_j \in \beta$ in priority of the largest cluster, for $j = 1, \dots, J$
- 3: **for** λ in λ^{seq} **do**
- 4: Calculate $\hat{\alpha}$ through I-LAMM in Algorithm 2
- 5: **if** $\hat{\alpha}$ is not rejected by MCD test (36) with size ϱ_n **then**
- 6: $\mathbf{I}[l] = 0$ for $l \in \{l : \hat{\alpha}_l = 0\}$
- 7: **if** $|\{j : \hat{\alpha}_j = 0\}| > |\{j : \hat{\alpha}_j^{\text{MCP}} = 0\}|$ **then**
- 8: $\hat{\alpha}^{\text{MCP}} = \hat{\alpha}$
- 9: **end if**
- 10: **end if**
- 11: **end for**
- 12: **if** $\|\mathbf{I}\|_1 \leq |\{j : \{j : \hat{\alpha}_j^{\text{MCP}} = 0\}|$ **then** \triangleright Impossible for the existence of more sparse $\hat{\alpha}$
- 13: Break
- 14: **end if**
- 15: **end for**

Output: $\hat{\alpha}^{\text{MCP}}$

We compare the WIT estimator with other popular estimators in the literature. Specifically, sisVIVE is computed by R package `sisVIVE`; Post-ALasso (Windmeijer et al., 2019), TSHT and CIIV are implemented using codes on Github (Guo et al., 2018; Windmeijer et al., 2021). TSLS, LIML, oracle-TSLS and oracle-LIML (the truly valid set \mathcal{V}^* is known a priori) are also included. Regarding our proposed WIT estimator, the MCD tuning strategy is implemented to determine λ , and we fix $\rho = 2$. In the I-LAMM algorithm, we take $\delta_c = 10^{-3}$ and $\delta_t = 10^{-5}$ as the tolerance levels. We report results based on 500 simulations.

We measure the performance of all estimators in terms of median absolute deviation (MAD), standard deviation (STD), and coverage probability (CP) based on 95% confidence interval. Moreover, we provide measurements on the estimation of α^* and IV model selection. Specifically, We measure the performance of invalid IVs selection by false positive rate (FPR) and false negative rate (FNR). To be concrete, denote the number of incorrect selections of valid and invalid IVs as FP and FN, respectively, and the number of correct selections of valid and invalid as TP and TN, respectively. Thus, $\text{FPR} = \text{FP}/(\text{FP} + \text{TN})$ and $\text{FNR} = \text{FN}/(\text{FN} + \text{TP})$.

4.1. Case 1: Low dimension

We first consider the low dimension scenario:

Case 1(I): $\gamma^* = (\mathbf{0.5}_4, \mathbf{0.6}_6)^\top$ and $\alpha^* = (\mathbf{0}_5, \mathbf{0.4}_3, \mathbf{0.8}_2)^\top$.

Case 1(II): $\gamma^* = (\mathbf{0.04}_3, \mathbf{0.5}_2, 0.2, \mathbf{0.1}_4)^\top$ and $\alpha^* = (\mathbf{0}_5, 1, \mathbf{0.7}_4)^\top$.

In the above two cases, we maintain $\sigma_\eta^2 = 1$ and vary the sample size $n = 200$ to 500. Case 1(I) refers to all strong IVs case, but SAIS (14) condition holds. Case 1(II) refers to Example 1 with mixed strong and weak IVs.

Table 1 presents the detailed results. In Case 1(I), high FPR or FNP indicates that

Table 1. Simulation results in low dimension

| Case | Approaches | $n = 200$ | | | | $n = 500$ | | | |
|-------|-------------|-----------|-------|-------|-------|-----------|-------|-------|-------|
| | | MAD | CP | FPR | FNR | MAD | CP | FPR | FNR |
| 1(I) | TSLs | 0.532 | 0 | - | - | 0.530 | 0 | - | - |
| | LIML | 0.952 | 0.004 | - | - | 0.978 | 0 | - | - |
| | oracle-TSLs | 0.038 | 0.938 | - | - | 0.023 | 0.958 | - | - |
| | oracle-LIML | 0.038 | 0.938 | - | - | 0.023 | 0.958 | - | - |
| | TSHT | 0.060 | 0.828 | 0.158 | 0.018 | 0.023 | 0.950 | 0 | 0.002 |
| | CIIV | 0.045 | 0.810 | 0.072 | 0.009 | 0.023 | 0.946 | 0.004 | 0.003 |
| | sisVIVE | 0.539 | - | 0.428 | 0.957 | 0.589 | - | 0.479 | 1 |
| | Post-Alasso | 0.532 | 0 | 1 | 0 | 0.530 | 0 | 0.996 | 0 |
| | WIT | 0.046 | 0.818 | 0.068 | 0.065 | 0.024 | 0.948 | 0.004 | 0.020 |
| | TSLs | 1.098 | 0 | - | - | 1.111 | 0 | - | - |
| 1(II) | LIML | 7.437 | 0.292 | - | - | 7.798 | 0.030 | - | - |
| | oracle-TSLs | 0.072 | 0.938 | - | - | 0.046 | 0.948 | - | - |
| | oracle-LIML | 0.072 | 0.950 | - | - | 0.046 | 0.958 | - | - |
| | TSHT | 0.110 | 0.914 | 0.122 | 0.585 | 0.742 | 0.598 | 0.423 | 0.712 |
| | CIIV | 0.099 | 0.724 | 0.088 | 0.642 | 4.321 | 0.334 | 0.360 | 0.824 |
| | sisVIVE | 0.259 | - | 0.018 | 0.234 | 0.154 | - | 0 | 0.168 |
| | Post-Alasso | 1.831 | 0.016 | 0.429 | 0.258 | 3.533 | 0 | 0.560 | 0.387 |
| | WIT | 0.079 | 0.914 | 0.016 | 0.034 | 0.049 | 0.920 | 0.016 | 0.016 |

sisVIVE and Post-Alasso mistarget α^* because of lack of majority rule and SIAS holds. Their performances do not improve much with sample size n . Due to finite strong IVs, WIT performs similarly to TSHT, CIIV, and oracle-LIML. In low dimension settings, the oracle-TSLs is very close to oracle-LIML. Case 1(II) shows how weak IVs break the strong IVs-based plurality rule. As shown in Example 1, TSHT and CIIV worsen in terms of all measures though the sample size increase from $n = 200$ to 500. Post-Alasso also fails due to the failure of majority rule. As the analysis in Example 1, such a mixture of weak IVs does not impede penalized methods. The WIT estimator outperforms even when $n = 200$ and approaches oracle-LIML when n goes to 500. Interestingly, the comparably low FPR and MAD imply sisVIVE correctly target true α^* since the additional requirement of matching objective function (15) happens to hold in this example. However, its FNR and MAD are worse than WIT due to the conservative cross-validation tuning strategy and the non-ignorable bias of Lasso, respectively.

Further, for closer comparison, we present a replication of the simulation design considered in Windmeijer et al. (2021) and its weak IVs variant:

Case 1(III) : $\gamma^* = (\mathbf{0.4}_{21})^\top$ and $\alpha^* = (\mathbf{0}_9, \mathbf{0.4}_6, \mathbf{0.2}_6)^\top$.

Case 1(IV) : $\gamma^* = (\mathbf{0.15}_{21})^\top$ and $\alpha^* = (\mathbf{0}_9, \mathbf{0.4}_6, \mathbf{0.2}_6)^\top$.

We now vary sample size $n = 500$ to 1000 and fix $\sigma_\eta^2 = 1$ to strictly follow their design. Between them, Case 1(III) corresponds to the exact setting, while Case 1(IV) scales down the magnitude of γ^* to introduce small coefficients in the first stage.

Table 2 shows the results. In Case 1(III), CIIV outperforms TSHT because CIIV can utilize available information better (Windmeijer et al., 2021, Section 7). The WIT estimator performs similar to CIIV and approaches oracle-LIML. sisVIVE and Post-Alasso fail again due to a lack of majority rule. In Case 1(IV), scaling down the first-stage

Table 2. Simulation results in low dimension: A replication of experiment (Windmeijer et al., 2021)

| Case | Approaches | $n = 500$ | | | | $n = 1000$ | | | |
|--------|-------------|-----------|-------|-------|-------|------------|-------|-------|-------|
| | | MAD | CP | FPR | FNR | MAD | CP | FPR | FNR |
| 1(III) | TSLS | 0.436 | 0 | - | - | 0.435 | 0 | - | - |
| | LIML | 0.729 | 0 | - | - | 0.739 | 0 | - | - |
| | oracle-TSLS | 0.021 | 0.936 | - | - | 0.014 | 0.942 | - | - |
| | oracle-LIML | 0.021 | 0.932 | - | - | 0.014 | 0.944 | - | - |
| | TSHT | 0.142 | 0.404 | 0.398 | 0.150 | 0.016 | 0.924 | 0.023 | 0.004 |
| | CIIV | 0.037 | 0.710 | 0.125 | 0.032 | 0.017 | 0.894 | 0.031 | 0.002 |
| | sisVIVE | 0.445 | - | 0.463 | 0.972 | 0.465 | - | 0.482 | 0.999 |
| | Post-Alasso | 0.436 | 0 | 1 | 0 | 0.435 | 0 | 0.999 | 0 |
| | WIT | 0.036 | 0.708 | 0.121 | 0.099 | 0.016 | 0.910 | 0.020 | 0.027 |
| | TSLS | 1.124 | 0 | - | - | 1.144 | 0 | - | - |
| 1(IV) | LIML | 1.952 | 0 | - | - | 1.976 | 0 | - | - |
| | oracle-TSLS | 0.060 | 0.936 | - | - | 0.044 | 0.942 | - | - |
| | oracle-LIML | 0.056 | 0.948 | - | - | 0.042 | 0.962 | - | - |
| | TSHT | 0.532 | 0.058 | 0.342 | 0.457 | 0.155 | 0.660 | 0.310 | 0.208 |
| | CIIV | 1.213 | 0.224 | 0.337 | 0.670 | 0.100 | 0.574 | 0.300 | 0.426 |
| | sisVIVE | 1.101 | - | 0.392 | 0.936 | 1.175 | - | 0.428 | 0.996 |
| | Post-Alasso | 1.112 | 0 | 0.945 | 0.010 | 1.029 | 0 | 0.652 | 0.205 |
| | WIT | 0.102 | 0.634 | 0.198 | 0.220 | 0.047 | 0.898 | 0.051 | 0.064 |

coefficients causes some problems for CIIV and TSHT, since the first-stage selection thresholding $\sigma_\eta \sqrt{2.01 \log p/n} = 0.111 < 0.15$, which might break the plurality rule numerically. TSHT and CIIV perform poorly when $n = 500$ and improve when $n = 1000$ when the issue of violating the plurality rule is mitigated. Among penalized methods, sisVIVE and Post-Alasso mistarget and perform like TSLS because an additional requirement for sisVIVE (15) and majority rule fail simultaneously. Distinguished from them, the WIT estimator outperforms with acceptable MAD when $n = 500$. The FPR and FNR improve when the sample size increases. Fig. 4 presents all replications in Case 1(IV) when $n = 1000$. It shows that WIT is nearly oracle besides the mild number of incorrect selections. By contrast, CIIV and TSHT fail in selections more frequently.

4.2. Case 2: High dimension (many IVs)

To assess performance in many IVs, we consider the following examples:

Case 2(I) (increasing $p = 0.5n$): $\gamma^* = (\mathbf{1.5}/\sqrt{n})_p^\top$ and $\alpha^* = (\mathbf{0}_{0.6p}, \mathbf{0.5}_{0.4p})^\top$.

Case 2(II) (increasing $p = 0.6n$): $\gamma^* = (\mathbf{1.5}/\sqrt{n})_p^\top$ and $\alpha^* = (\mathbf{0}_{0.4p}, -\mathbf{0.5}_{0.2p}, \mathbf{1}_{0.3p}, -\mathbf{1}_{0.1p})^\top$.

The numbers of valid and invalid IVs are growing with the sample size. To verify the theoretical result, we maintain the ratio of concentration parameter to sample size n at a low constant level, i.e. $\mu_n/n = 0.5$, by adjusting σ_η^2 . We vary sample size n from 500 to 1000, and let the first-stage coefficients go to 0. Due to the computational burden in CIIV for many IVs case, we omit its results.

Table 3 provides the detailed estimation results. Case 2(I) satisfies the majority rule, and only two groups are present. All weak IVs narrow available choices in TSHT. When $n = 1000$, a low FPR but high FNR indicates that TSHT only selected a limited

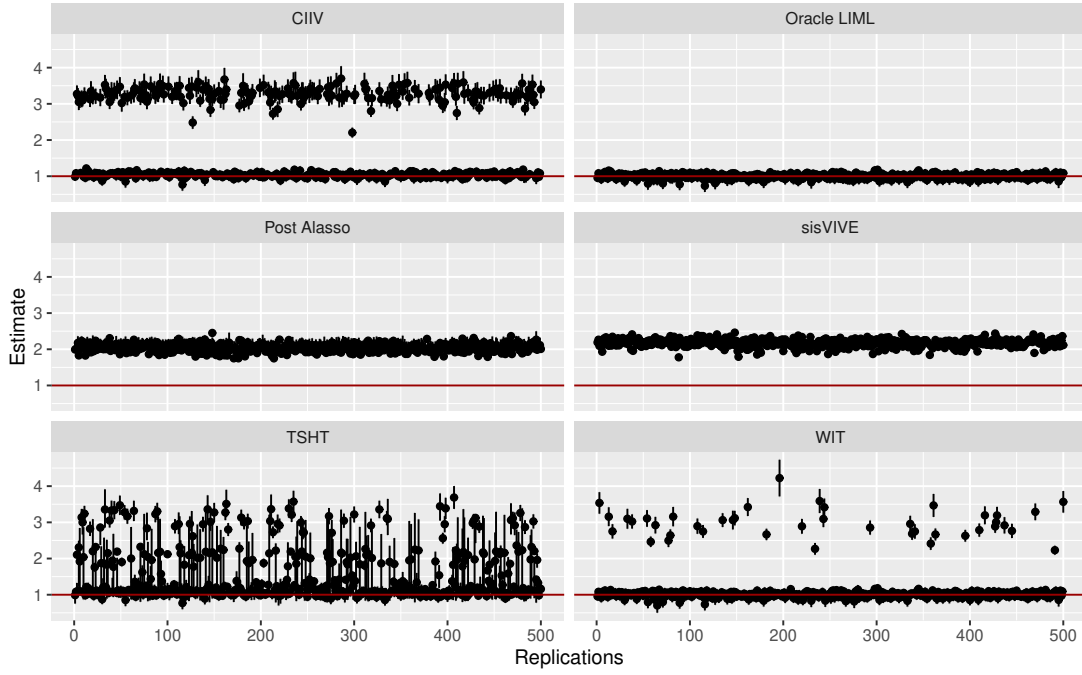


Figure 4. Scatter plot of estimations of β^* with confidence intervals of Case 1(IV) for $n = 1000$. The red solid lines show the true value of $\beta^* = 1$.

Table 3. Simulation Results in High dimension (many IVs)

| Case | Approaches | $n = 500$ | | | | $n = 1000$ | | | |
|-------|-------------|-----------|-------|-------|-------|------------|-------|-------|-------|
| | | MAD | CP | FPR | FNR | MAD | CP | FPR | FNR |
| 2(I) | TSLS | 2.170 | 0 | - | - | 3.015 | 0 | - | - |
| | LIML | 8.266 | 0 | - | - | 11.195 | 0 | - | - |
| | oracle-TSLS | 0.177 | 0 | - | - | 0.176 | 0 | - | - |
| | oracle-LIML | 0.044 | 0.950 | - | - | 0.027 | 0.956 | - | - |
| | TSHT | 0.972 | 0.108 | 0.221 | 0.698 | 0.376 | 0.112 | 0.002 | 0.649 |
| | sisVIVE | 0.552 | - | 0 | 0.214 | 0.551 | - | 0 | 0.123 |
| | Post-Alasso | 1.801 | 0 | 0.540 | 0 | 0.403 | 0 | 0.083 | 0 |
| | WIT | 0.047 | 0.936 | 0.003 | 0.013 | 0.029 | 0.946 | 0 | 0.001 |
| 2(II) | TSLS | 1.300 | 0.096 | - | - | 1.733 | 0.010 | - | - |
| | LIML | 102.843 | 0.586 | - | - | 132.231 | 0.212 | - | - |
| | oracle-TSLS | 0.191 | 0.008 | - | - | 0.196 | 0 | - | - |
| | oracle-LIML | 0.054 | 0.958 | - | - | 0.035 | 0.936 | - | - |
| | TSHT | 1.592 | 0.254 | 0.293 | 0.937 | 2.577 | 0.174 | 0.296 | 0.955 |
| | sisVIVE | 0.280 | - | 0 | 0.201 | 0.284 | - | 0 | 0.118 |
| | Post-Alasso | 0.195 | 0 | 0.030 | 0 | 0.198 | 0 | 0.002 | 0 |
| | WIT | 0.085 | 0.760 | 0.005 | 0.010 | 0.044 | 0.904 | 0 | 0.010 |

number of valid IVs. Both low FPR and FNR implies sisVIVE targets the correct solutions. Meanwhile, a high MAD shows that the TSLS based method is biased in many IVs. The majority rule holds to ensure Post-Alasso is consistent. However, when $n = 500$, Post-Alasso is severely biased with relatively high FPR. It may result from the sensitivity problem of the initial estimator of Post-Alasso in weak IVs. Among these approaches, the WIT estimator behaves highly similar to oracle-LIML and achieves the best performance in every measure, even when $n = 500$. On the contrary, oracle-TSLS has a much larger bias than LIML in many IVs cases, and the coverage rate is poor.

Case 2(II) has more invalid IVs with the sparsest rule. TSHT breaks down in this case since the strong IVs-based plurality rule is unlikely to hold. sisVIVE correctly identifies solutions by chance in this example. However, its estimate suffers from the bias of Lasso and TSLS-based bias in many IVs. Without majority rule, Post-Alasso performs comparably to the WIT estimator in terms of FPR and FNR. But it is only a coincidence that the initial estimator in Post-Alasso is consistent in this example since $E(\text{median}(\beta_j^*)) = \text{median}(\{\mathbf{1}_{0.4p}, -\mathbf{9.540}_{0.2p}, \mathbf{22.081}_{0.3p}, -\mathbf{20.081}_{0.1p}\}) = 1$. Compared with sisVIVE, Post-Alasso removes the bias of Lasso. Compared with Post-Alasso, WIT circumvents the majority rule and bias from TSLS. Thus, WIT outperforms in many IVs cases and approaches oracle-LIML when the sample size increases.

5. Application to Trade and Economic Growth

In this section, we revisit a classic empirical study on trade and growth ([Frankel and Romer, 1999](#), FR99 henceforth). This problem remains a frontier and intensely debated issue in the field of international economics and also has strong guidance for policy-making. We investigate the causal effect of trade on income using more comprehensive and updated data, taking into account that trade is an endogenous variable (it correlates with unobserved common factors driving both trade and growth), and some instruments might be invalid.

The structural equation considered in FR99 is,

$$\log(Y_i) = \alpha + \beta T_i + \mathbf{S}_{i\cdot}^\top \boldsymbol{\psi}_1 + \epsilon_i,$$

where for each country i , Y_i is the GDP per worker, T_i is the share of international trade to GDP, $\mathbf{S}_{i\cdot}$ denotes the size of the country, such as area and population, and ϵ_i is the error term.

FR99 proposed to construct an IV (called a proxy for trade) based on the gravity theory of trade ([Anderson, 1979](#)). The logic of IV validity in aggregate variables is that the geographical variables, such as common border and distance between countries, indirectly affect growth through the channel of convenience for trade.

Following the same logic, [Fan and Zhong \(2018\)](#) extended the IV set to include more geographic and meteorological variables. The first stage equation is

$$T_i = \mathbf{Z}_{i\cdot}^\top \boldsymbol{\gamma} + \mathbf{S}_{i\cdot}^\top \boldsymbol{\psi}_2 + \nu_i, \quad (37)$$

where $\mathbf{Z}_{i\cdot}$ is a vector of instruments that we elaborate in Section 5.1 below. In this study, we expand the candidate IV set even further. On the one hand, more information contained in newly introduced IVs could increase the accuracy of estimating the causal

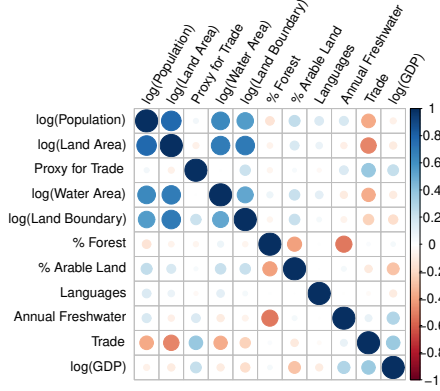


Figure 5. Correlation of all variables

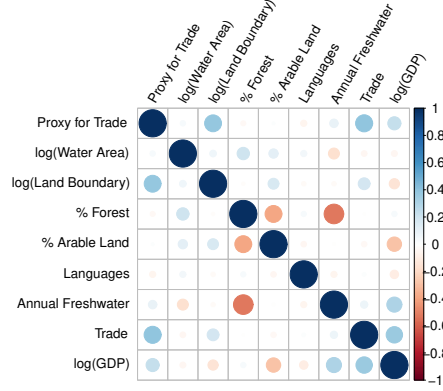


Figure 6. Correlation of transformed variables

effect of trade. On the other hand, some of the newly introduced IVs might be invalid. Also, with a large IV set with potentially invalid IVs, it is desirable to have a robust estimate of the treatment effect when the signal in the first stage is weak. This problem can be addressed by the proposed WIT estimator as discussed in previous sections.

5.1. Data Description and Empirical Model

We collect cross-sectional data from 158 countries in the year 2017. Table A1 (in Appendix A6) presents the summary statistics of the main variables.

We first standardize all the variables, then we formulate the structural equation as:

$$\log(Y_i) = T_i\beta + \mathbf{Z}_i^\top \boldsymbol{\alpha} + \mathbf{S}_i^\top \boldsymbol{\psi}_1 + \epsilon_i, \quad \text{for } i = 1, 2, \dots, 158, \quad (38)$$

where \mathbf{S}_i consisting of $\log(\text{Population})$ and $\log(\text{Land Area})$ for country i serve as control variables following FR99, \mathbf{Z}_i contain the potential IVs including all from Fan and Zhong (2018) and more geo-economic variables from the World Bank database, and $\boldsymbol{\alpha}$ indicates whether the IV is invalid. In the first stage, we consider the linear reduced form equation (37). Fig. 5 shows the plot of correlations between the variables.

Following Fan and Wu (2022), we partial out the effect from control variables \mathbf{S} in (38), which does not affect the estimation of $\boldsymbol{\alpha}$ and β . Denote $\mathbf{M}_\mathbf{S}$ as the projection matrix on the orthogonal space with respect to the column space of \mathbf{S} , we transform $\{\mathbf{Y}, \mathbf{T}, \mathbf{Z}\}$ to $\{\mathbf{M}_\mathbf{S}\mathbf{Y}, \mathbf{M}_\mathbf{S}\mathbf{T}, \mathbf{M}_\mathbf{S}\mathbf{Z}\}$ and denote the transformed observations as $\{\ddot{\mathbf{Y}}, \ddot{\mathbf{T}}, \ddot{\mathbf{Z}}\}$. Equivalently for the error terms, $\ddot{\epsilon}$ and $\ddot{\nu}$. Hence, we work on the transformed model:

$$\ddot{Y}_i = \ddot{T}_i\beta + \ddot{\mathbf{Z}}_i^\top \boldsymbol{\alpha} + \ddot{\epsilon}_i, \quad \ddot{T}_i = \ddot{\mathbf{Z}}_i^\top \boldsymbol{\gamma} + \ddot{\nu}_i.$$

The correlation matrix of the transformed variables is plotted in Fig 6.

Table 4. Empirical Results of Various Estimators

| | $\hat{\beta} \left(\widehat{\text{Var}}^{1/2}(\hat{\beta}) \right)$ | 95% CI | Valid IVs $\hat{\mathcal{V}}$ | Relevant IVs $\hat{\mathcal{S}}$ | Sargan Test |
|-------------|--|-----------------|-------------------------------|----------------------------------|-------------|
| OLS | 0.413(0.084) | (0.246, 0.581) | - | - | - |
| FR99 | 0.673(0.220) | (0.228, 1.117) | - | - | 0.999 |
| LIML | 2.969(1.503) | (0.023, 5.916) | - | - | 0.001 |
| TSHT | 0.861(0.245) | (0.380, 1.342) | {1} | {1} | 0.999 |
| CIIV* | 2.635(1.974) | (-1.233, 6.504) | {2,4,5,6,7} | - | 0.385 |
| sisVIVE | 0.819(-) | - | {1,2,4} | - | 0.418 |
| Post-Alasso | 0.964(0.251) | (0.471, 1.457) | {1,2,4,5,6} | - | 0.086 |
| WIT | 0.974(0.323) | (0.340, 1.609) | {1,2,4,6} | - | 0.275 |

Note: CIIV* stands for CIIV method without first-stage IVs selection because it reports that “Less than two IVs are individually relevant, treat all IVs as strong”. Sargan test p -value is shown in the last column. The selection of relevant IVs $\hat{\mathcal{S}}$ is only implemented in TSHT and CIIV.

5.2. Empirical Results

We explore the causal effect of trade using the proposed WIT estimator and also provide the estimation results from other popular estimators, including sisVIVE, TSHT, CIIV, OLS, LIML, and FR99 (TSLS using one IV \hat{T}) for comparison.

Table 4 provides the detailed results of estimation and inference. The p -value of the Hausman test for endogeneity is 1.81e-3 using the proxy for trade as IV. The OLS estimate is likely biased due to the endogeneity of trade. The FR99 result using the proxy $Z_1 = \hat{T}$ as an instrument gives a smaller treatment effect estimate compared to WIT. LIML using all potential IVs (without distinguishing the invalid ones) likely overestimates the treatment effect. The 0.001 p -value of the Sargan test rejects the null that all potential IVs are valid.

For the penalized IV regression approaches: the sisVIVE, Post-Alasso and WIT estimators jointly select 3 valid IVs: Z_1 , Z_2 and Z_4 . Z_6 : Languages is selected by WIT as valid but not by sisVIVE. Regarding the p -value of Sargan test for WIT and sisVIVE, 0.275, and 0.418, it supports that Z_6 is valid. Compared with Post-Alasso, Z_5 : % Arable Land is chosen by Post-Alasso while not by WIT estimator. In view of the marginal correlation in Fig. 6, Z_5 is nearly uncorrelated to trade but significantly correlated with $\log(\text{GDP})$. Concerning the Sargan test, a p -value of $0.086 < 0.1$ for Post-Alasso indicates Z_5 is not very credible to be valid. From an economic perspective, more arable land generates higher crop yields and maintains a higher agriculture sector labor force, which directly affects GDP. Thus it has a channel to affect GDP directly. Put together, we conclude Z_5 should be an invalid IV.

Regarding individual IV estimator-based approaches, their selection results might not be accurate compared to the WIT estimator. These methods impose a stronger condition of relevant IVs than WIT, which prohibits their selection in weaker IVs, narrows the available choices, and cannot handle weak IVs problems. Specifically, the first-stage threshold of TSHT only selects proxy for trade, a consensus valid instrument. While CIIV*, which treats all IVs as strong, apparently chooses the incorrect set since it excludes proxy of trade. It shows the sensitivity problem of individual IV estimator-based approaches when weak IVs are present.

Four out of seven IVs are estimated as valid by WIT, suggesting that the majority rule holds. This supports that the selection result of Post-Alasso is close to WIT. Even

though $\hat{\beta} = 0.964$ in Post-Alasso is almost identical to WIT estimator, Post-Alasso utilizes TSLS as second-stage estimator instead of LIML in WIT estimator. Nevertheless, a closer check finds that the first-stage F-value is $3.682 < 5$, indicating the weak IVs problem in Post-Alasso estimated valid IVs. Using LIML as the second stage instead provides a much different estimate: $1.441(0.421)$ that is far away from the result of the WIT estimator.

6. Conclusion

We extended the study of IV models with unknown invalid IVs to allow for many weak IVs. We provided a complete framework to investigate the identification issue of such models and showed the impossibility of the existence of an *iff* identification condition. Sticking to the sparsest rule, we proposed the surrogate sparsest penalty that fits the identification condition. We proposed a novel WIT estimator that addresses the issues that can lead to poor performance of sisVIVE and Post-Alasso, and can outperform the plurality rule-based TSHT and CIIV. Simulations and real data analysis support the theoretical findings and the advantages of the proposed method over existing approaches.

Acknowledgements

References

- Anderson, J. E. (1979) A theoretical foundation for the gravity equation. *Am. Econ. Rev.*, **69**, 106–116.
- Andrews, I., Stock, J. and Sun, L. (2018) Weak instruments in IV regression: Theory and practice. *Annu. Rev. Econ.*
- Bekker, P. A. (1994) Alternative approximations to the distributions of instrumental variable estimators. *Econometrica*, 657–681.
- Belloni, A., Chen, D., Chernozhukov, V. and Hansen, C. (2012) Sparse models and methods for optimal instruments with an application to eminent domain. *Econometrica*, **80**, 2369–2429.
- Bickel, P. J., Ritov, Y. and Tsybakov, A. B. (2009) Simultaneous analysis of lasso and dantzig selector. *The Ann. Statist.*, **37**, 1705–1732.
- Bound, J., Jaeger, D. A. and Baker, R. M. (1995) Problems with instrumental variables estimation when the correlation between the instruments and the endogenous explanatory variable is weak. *J. Am. Statist. Ass.*, **90**, 443–450.
- Bun, M. J. and Windmeijer, F. (2011) A comparison of bias approximations for the two-stage least squares (2sls) estimator. *Econ. Lett.*, **113**, 76–79.
- Chao, J. C. and Swanson, N. R. (2005) Consistent estimation with a large number of weak instruments. *Econometrica*, **73**, 1673–1692.

- Daubechies, I., Defrise, M. and De Mol, C. (2004) An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Communications on Pure and Applied Mathematics*, **57**, 1413–1457.
- Davey Smith, G. and Ebrahim, S. (2003) ‘Mendelian randomization’: can genetic epidemiology contribute to understanding environmental determinants of disease? *International journal of epidemiology*, **32**, 1–22.
- Fan, J. and Li, R. (2001) Variable selection via nonconcave penalized likelihood and its oracle properties. *J. Am. Statist. Ass.*, **96**, 1348–1360.
- Fan, J., Liu, H., Sun, Q. and Zhang, T. (2018) I-lamm for sparse learning: Simultaneous control of algorithmic complexity and statistical error. *Ann. Statist.*, **46**, 814.
- Fan, Q. and Wu, Y. (2022) Endogenous treatment effect estimation with a large and mixed set of instruments and control variables. *Rev. Econ. Statist.*, forthcoming.
- Fan, Q. and Zhong, W. (2018) Nonparametric additive instrumental variable estimator: A group shrinkage estimation perspective. *J. Bus. Econ. Statist.*, **36**, 388–399.
- Feng, L. and Zhang, C.-H. (2019) Sorted concave penalized regression. *The Ann. Statist.*, **47**, 3069–3098.
- Frankel, J. A. and Romer, D. H. (1999) Does trade cause growth? *Am. Econ. Rev.*, **89**, 379–399.
- van de Geer, S. A. and Bühlmann, P. (2009) On the conditions used to prove oracle results for the lasso. *Electron. J. Statist.*, **3**, 1360–1392.
- Guo, Z. and Bühlmann, P. (2022) Two stage curvature identification with machine learning: Causal inference with possibly invalid instrumental variables. *arXiv preprint arxiv.2203.12808*.
- Guo, Z., Kang, H., Tony Cai, T. and Small, D. S. (2018) Confidence intervals for causal effects with invalid instruments by using two-stage hard thresholding with voting. *J. R. Statist. Soc. B*, **80**, 793–815.
- Han, C. (2008) Detecting invalid instruments using l_1 -GMM. *Economics Letters*, **3**, 285–287.
- Hansen, C., Hausman, J. and Newey, W. (2008) Estimation with many instrumental variables. *J. Bus. Econ. Statist.*, **26**, 398–422.
- Hansen, C. and Kozbur, D. (2014) Instrumental variables estimation with many weak instruments using regularized JIVE. *J. Econometrics*, **182**, 290–308.
- Javanmard, A. and Montanari, A. (2018) Debiasing the lasso: Optimal sample size for gaussian designs. *The Ann. Statist.*, **46**, 2593–2622.
- Kang, H., Zhang, A., Cai, T. T. and Small, D. S. (2016) Instrumental variables estimation with some invalid instruments and its application to mendelian randomization. *J. Am. Statist. Ass.*, **111**, 132–144.

- Kolesár, M. (2018) Minimum distance approach to inference with many instruments. *J. Econometrics*, **204**, 86–100.
- Kolesár, M., Chetty, R., Friedman, J., Glaeser, E. and Imbens, G. W. (2015) Identification and inference with many invalid instruments. *J. Bus. Econ. Statist.*, **33**, 474–484.
- Lewbel, A. (2012) Using heteroscedasticity to identify and estimate mismeasured and endogenous regressor models. *J. Bus. Econ. Statist.*, **30**, 67–80.
- Lin, W., Feng, R. and Li, H. (2015) Regularization methods for high-dimensional instrumental variables regression with an application to genetical genomics. *J. Am. Statist. Ass.*, **110**, 270–288.
- Loh, P.-L. and Wainwright, M. J. (2015) Regularized m-estimators with nonconvexity: Statistical and algorithmic theory for local optima. *J. Mach. Learn. Res.*, **16**, 559–616.
- (2017) Support recovery without incoherence: A case for nonconvex regularization. *The Ann. Statist.*, **45**, 2455–2482.
- Nagar, A. L. (1959) The bias and moment matrix of the general k-class estimators of the parameters in simultaneous equations. *Econometrica*, 575–595.
- Newey, W. K. and Windmeijer, F. (2009) Generalized method of moments with many weak moment conditions. *Econometrica*, **77**, 687–719.
- Sargan, J. D. (1958) The estimation of economic relationships using instrumental variables. *Econometrica*, 393–415.
- Sawa, T. (1969) The exact sampling distribution of ordinary least squares and two-stage least squares estimators. *J. Am. Statist. Ass.*, **64**, 923–937.
- Small, D. S. (2007) Sensitivity analysis for instrumental variables regression with overidentifying restrictions. *J. Am. Statist. Ass.*, **102**, 1049–1058.
- Staiger, D. and Stock, J. H. (1997) Instrumental variables regression with weak instruments. *Econometrica*, 557–586.
- Stock, J., Yogo, M. and Wright, J. (2002) A survey of weak instruments and weak identification in generalized method of moments. *J. Bus. Econ. Statist.*, **20**, 518–529.
- Tchetgen, E. T., Sun, B. and Walter, S. (2021) The GENIUS approach to robust mendelian randomization inference. *Statistical Science*, **36**, 443–464.
- Windmeijer, F., Farbmacher, H., Davies, N. and Davey Smith, G. (2019) On the use of the lasso for instrumental variables estimation with some invalid instruments. *J. Am. Statist. Ass.*, **114**, 1339–1350.
- Windmeijer, F., Liang, X., Hartwig, F. P. and Bowden, J. (2021) The confidence interval method for selecting valid instrumental variables. *J. R. Statist. Soc. B*, **83**, 752–776.

- Wooldridge, J. M. (2010) *Econometric analysis of cross section and panel data*. MIT press.
- Zhang, C.-H. and Zhang, T. (2012) A general theory of concave regularization for high-dimensional sparse estimation problems. *Statistical Science*, **27**, 576–593.
- Zhang, C.-H. et al. (2010) Nearly unbiased variable selection under minimax concave penalty. *The Ann. Statist.*, **38**, 894–942.
- Zhao, P. and Yu, B. (2006) On model selection consistency of lasso. *J. Mach. Learn. Res.*, **7**, 2541–2563.
- Zou, H. (2006) The adaptive lasso and its oracle properties. *J. Am. Statist. Ass.*, **101**, 1418–1429.
- Zou, H. and Li, R. (2008) One-step sparse estimates in nonconcave penalized likelihood models. *Ann. Statist.*, **36**, 1509.

Web-based supporting materials for “On the instrumental variable estimation with many weak and invalid instruments”

Yiqi Lin^{a,b}, Frank Windmeijer^b, Xinyuan Song^a, Qingliang Fan^c

^aDepartment of Statistics, The Chinese University of Hong Kong, Hong Kong.

^bDepartment of Statistics, University of Oxford, Oxford, U.K.

^cDepartment of Economics, The Chinese University of Hong Kong, Hong Kong.

This supplementary material mainly includes the following parts: [Appendix A](#) provides additional details of the main paper. [Appendix B](#) contains all technical proofs. Throughout the online supplementary material, we allow constant C to be a generic positive constant that may differ in different cases.

Appendix A. Additional Demonstrations

Appendix A1. Discussion of \mathbf{Z}_i and its non-linear transformation

In Section 2.1, it is mentioned that \mathbf{Z}_i can be non-linear transformations of original variables such as polynomials and B-splines to form a high-dimensional model and provide flexible model fitting. In general, Eq. (2) can be reformulated as

$$Y_i = D_i\beta + w_1(\mathbf{Z}_i) + \epsilon_i \approx D_i\beta + \phi(\mathbf{Z}_i)^\top \boldsymbol{\alpha} + \epsilon_i, \quad (\text{A1})$$

$$D_i = w_2(\mathbf{Z}_i) + \eta_i \approx \phi(\mathbf{Z}_i)^\top \boldsymbol{\gamma} + \eta_i, \quad (\text{A2})$$

where $w_1(\mathbf{Z}_i) \neq 0$ and $w_2(\mathbf{Z}_i) \neq 0$ are different unknown functions that can be approximated by the same basis expansion through $\phi(\cdot)$. With the same argument in (2), the $\boldsymbol{\alpha}$ in (A1) should consist of zero and non-zero terms to ensure the existence of an identifiable sub-model.

A subtle yet important point here is that we should use the same basis $\phi(\cdot)$ to approximate different w_1 and w_2 . Otherwise, using different ϕ_1 and ϕ_2 in (A1) and (A2), respectively, requires strong prior information to justify why the components in $\{\phi_2(\mathbf{Z}_i)\} \setminus \{\phi_1(\mathbf{Z}_i)\}$ can be treated as valid IVs. In the following we show that this could be a stringent assumption for model identification. For simplicity of illustrating this point, we consider $\{\phi_1(\mathbf{Z}_i)\} \subseteq \{\phi_2(\mathbf{Z}_i)\}$ firstly. It can be extended to general $\phi_1 \neq \phi_2$ since coefficients of $\{\phi_1(\mathbf{Z}_i)\} \setminus \{\phi_2(\mathbf{Z}_i)\}$ can be estimated directly by their moment conditions. Nonetheless, in the simplified case, it is easy to see that

$$\begin{aligned} Y_i &= D_i\beta^* + \phi_1(\mathbf{Z}_i)^\top \boldsymbol{\alpha}^* + \epsilon_i \\ &= D_i\beta^* + \phi_1(\mathbf{Z}_i)^\top \boldsymbol{\alpha}^* + (\{\phi_2(\mathbf{Z}_i)\} \setminus \{\phi_1(\mathbf{Z}_i)\})^\top \mathbf{0} + \epsilon_i \\ &= D_i\beta^* + \phi_2(\mathbf{Z}_i)^\top (\boldsymbol{\alpha}^{*\top}, \mathbf{0}^\top)^\top + \epsilon_i, \\ D_i &= \phi_2(\mathbf{Z}_i)^\top \boldsymbol{\gamma}^* + \eta_i, \end{aligned}$$

where we denote the $\boldsymbol{\alpha}^{\text{new}*} = (\boldsymbol{\alpha}^{*\top}, \mathbf{0}^\top)^\top$. Therefore, by Theorem 1, we immediately know that all other possible $\tilde{\boldsymbol{\alpha}}^{\text{new}}$ in the remaining valid DGPs must have non-zero coefficients in $\{\phi_2(\mathbf{Z}_i)\} \setminus \{\phi_1(\mathbf{Z}_i)\}$ but are impossible to be selected because zero coefficients for $\{\phi_2(\mathbf{Z}_i)\} \setminus \{\phi_1(\mathbf{Z}_i)\}$ are set by default in (A1). The identification condition essentially assumes the fixed functional form of ϕ_1 in (A1).

Appendix A2. Further Discussions on Assumption 5

Recall that Assumption 5: $|\tilde{\alpha}_j^c| = |\alpha_j^* - c\gamma_j^*| = |\tilde{c} - c| \cdot |\gamma_j^*| > \kappa(n)$ and $|\alpha_{\mathcal{V}^{c*}}^*|_{\min} > \kappa(n)$, is a “beta-min” condition for penalized method. $|\tilde{c} - c| \neq 0$ is required for individual IV estimator-based approaches: TSHT and CIIV. We consider the finite p case where TSHT and CIIV primarily work on and thus, $\kappa(n)$ is specified as $\sqrt{\log p \mathcal{V}^*/n} \asymp n^{-1/2}$.

For individual IV estimator-based approaches, Guo et al. (2018); Windmeijer et al. (2021) rely on $\hat{\beta}_j = \hat{\Gamma}_j / \hat{\gamma}_j \xrightarrow{p} \beta^* + \alpha_j^* / \gamma_j^*$, where $\hat{\Gamma} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{Y}$ and $\hat{\Gamma} = (\mathbf{Z}^\top \mathbf{Z})^{-1} \mathbf{Z}^\top \mathbf{D}$. For grouping of different values of α_j^* / γ_j^* , the assumption $|\tilde{c} - c| \neq 0$ is natural. However, the penalized methods directly target on $\alpha \in \mathcal{Q}$ through (11) instead of α_j^* / γ_j^* . Hence it leads to a different requirement: $|\tilde{c} - c| \cdot |\gamma_j^*| \neq 0$.

In Guo et al. (2018), first-stage thresholding to construct relevant IVs set $\hat{\mathcal{S}}$ rules out the weak IVs because the small value in $\hat{\gamma}_j$ appearing in the denominator of $\hat{\beta}_j$ can cause trouble to the grouping of $\hat{\beta}_j$ subsequently. However, the weak IVs can be useful in detecting valid IVs by having a large $\tilde{c} = \alpha_j^* / \gamma_j^*$ for a fixed α_j^* . E.g., let $c = \alpha_i^* / \gamma_i^*$, $\tilde{c} = \alpha_j^* / \gamma_j^*$ and $\gamma_i = C_1 / \sqrt{n}$, $\gamma_j = C_2 / \sqrt{n}$. Even though $|\tilde{c} - c| = \sqrt{n} |\alpha_i^* / C_1 - \alpha_j^* / C_2| \neq 0$ grows with n and is easy to distinguish, TSHT or CIIV cannot utilize that because they discard variables with small γ_i and γ_j . Whereas, $|\tilde{c} - c| \cdot |\gamma_j^*| = |C_2| \cdot |\alpha_i^* / C_1 - \alpha_j^* / C_2| \neq 0$ are able to be utilized for penalized methods.

Take simulation study Case 2(1) for another empirically-relevant example: $\gamma^* = (\mathbf{1.5} / \sqrt{n})_p^\top$ and $\alpha^* = (\mathbf{0}_{0.6p}, \mathbf{0.5}_{0.4p})^\top$, where $p = 0.5n$ is a typical “many (individually) weak IVs” case with a low concentration parameter. We can see from the simulation results that WIT can distinguish the valid IVs under this scenario. The accommodation of WIT for weak IVs is also numerically demonstrated by simulation Case 1(II), 1(IV), and real data application in trade and economic growth.

Appendix A3. Justification of some common cases in fixed or many IVs.

This section provides some specific cases where we directly verify that Assumption 5 is not stringent and allows the existence of weak IVs.

First we consider fixed p and a mixture of strong and weak IVs. The main result is shown in Corollary 2. Recall $\gamma_j^* = Cn^{-\tau_j}$ for $j = 1, 2, \dots, p$. Thus, let $\tau_{\mathcal{V}^*} = \arg\max_{\tau_j} \{\tau_j : j \in \mathcal{V}^*\}$, $\tau_{\mathcal{I}_c} = \arg\max_{\tau_j} \{\tau_j : j \in \mathcal{I}_c\}$ and $\tau_{\mathcal{I}_{\tilde{c}}} = \arg\max_{\tau_j} \{\tau_j : j \in \mathcal{I}_{\tilde{c}}\}$, where $c \neq \tilde{c}$. And we assume $0 \leq \tau_j \leq 1/2$, i.e., each potential IV is at least weak in terms of Staiger and Stock (1997). Hence, for any fixed $\alpha_j > 0$, $j \in \mathcal{V}^{c*}$, if $\tau_{\mathcal{V}^*} + 2\tau_{\mathcal{I}_c} < 1$, $\tau_{\mathcal{V}^*} + 2\tau_{\mathcal{I}_{\tilde{c}}} < 1$ and $\tau_{\mathcal{I}_c} + \tau_{\mathcal{I}_{\tilde{c}}} < 2/3$, then we can obtain the selection consistency and estimation consistency and asymptotic normality.

To be specific, it allows such mixture cases:

- (a) if $\tau_{\mathcal{V}^*} = 0$, i.e. valid IVs are strong, then invalid IVs can be semi-weak in the sense of $\tau_{\mathcal{I}_c} + \tau_{\mathcal{I}_{\tilde{c}}} < 2/3$.
- (b) all potentials IVs are semi-weak: $\tau_j = 1/3 - \delta$, where $\delta \rightarrow 0^+$.
- (c) if $\tau_{\mathcal{V}^*} = 1/2$, i.e., a mixture of strong and weak valid IVs (and at least one strong IV), then the invalid IVs can be semi-weak in the sense of $\tau_j < 1/4$ for $j \in \mathcal{V}^{c*}$.

Among above three cases, (a) extends the working scenario of Kang et al. (2016) that also avoids SAIS condition, (b) refers to simulation Case 1(4) and (c) corresponds to Example 1 and simulation Case 1(2).

Second, we turn to many (individually) weak IVs case. WLOG, we assume $\gamma_j = Cn^{-\tau_{V^*}}$ for $j \in \mathcal{V}^*$ and $\gamma_j = Cn^{-\tau_{V^{c*}}}$ for $j \in \mathcal{V}^{c*}$. Then according to Theorem 3, we known

$$\kappa(n) \asymp \underbrace{\sqrt{\frac{\log p_{V^*}}{n}}}_{T_1} + \underbrace{\frac{p_{V^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n \gamma_{V^*}^*\|_\infty}{\gamma_{V^*}^{*\top} \tilde{\mathbf{Q}}_n \gamma_{V^*}^*}}_{T_2} + \underbrace{|\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}})| \|\tilde{\gamma}_{V^{c*}}^*\|_\infty}_{T_3}.$$

Combined with Propositions 2 and 3, we have $T_2 \asymp o(n^{-\tau_{V^{c*}}})$ and $T_3 \approx O(\frac{\sigma_{\epsilon\eta}}{\sigma_\eta^2} \frac{\nu_{p_{V^*}}}{\mu_0 + \nu_{p_{V^*}}} \cdot n^{-\tau_{V^{c*}}})$. Thus, $\kappa(n) \asymp \sqrt{\frac{\log p_{V^*}}{n}} + \frac{\sigma_{\epsilon\eta}}{\sigma_\eta^2} \frac{\nu_{p_{V^*}}}{\mu_0 + \nu_{p_{V^*}}} n^{-\tau_{V^{c*}}} \rightarrow 0$ if $\tau_{V^{c*}} > 0$. Therefore, it suffices to have fixed $\alpha_j > 0$, $j \in \mathcal{V}^{c*}$ to satisfy

$$|\alpha_{V^{c*}}^*|_{\min} > C\kappa(n),$$

which is the “beta-min” condition.

Then we check whether $|\tilde{\alpha}_j^c| > \kappa^c(n)$ for $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$ holds or not. Notice $c = \alpha_j^* \cdot Cn^{\tau_{V^*}} \asymp n^{\tau_{V^*}}$,

$$\begin{aligned} \kappa^c(n) &\asymp (1+c) \left\{ \sqrt{\frac{\log |\mathcal{I}_c|}{n}} + \frac{|\mathcal{I}_c|}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n^c \gamma_{\mathcal{I}_c}^*\|_\infty}{\gamma_{\mathcal{I}_c}^{*\top} \tilde{\mathbf{Q}}_n^c \gamma_{\mathcal{I}_c}^*} \right\} + |\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}})| \|\tilde{\gamma}_{\mathcal{I}_c^c}^*\|_\infty \\ &\asymp \sqrt{\log |\mathcal{I}_c|} n^{\tau_{V^{c*}} - 1/2} + n^{3\tau_{V^{c*}} - \min(\tau_{V^{c*}}, \tau_{V^*}) - 1} / |\mathcal{I}_c|. \end{aligned}$$

For $|\tilde{\alpha}_j^c| = |\alpha_j^* - c\gamma_j^*|$ and $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$, we consider all possible cases:

- (a) $\tilde{c} = 0$, i.e. $j \in \{\mathcal{I}_c^c : \alpha_j^* = 0\}$: $|\tilde{\alpha}_j^c| = |c| \cdot |\gamma_j^*|$, and $\gamma_j^* = Cn^{-\tau_{V^*}}$. Hence $|\tilde{\alpha}_j^c| \asymp n^{\tau_{V^{c*}} - \tau_{V^*}}$.
- (b) $\tilde{c} \neq 0$, i.e. $j \in \{\mathcal{I}_c^c : \alpha_j^*/\gamma_j^* = \tilde{c}\}$: $|\tilde{\alpha}_j^c| = |c - \tilde{c}| \cdot |\gamma_j^*|$ and $\gamma_j^* = n^{-\tau_{V^{c*}}}$. Hence, $|\tilde{\alpha}_j^c| \asymp C$

Thus, $|\tilde{\alpha}_j^c| > \kappa^c(n)$ is equivalent to $0 \leq \tau_{V^*} \leq 1/2 - \delta$, $\delta \rightarrow 0^+$ and $0 \leq \tau_{V^{c*}} \leq 1/2$.

Thus, in many individually weak IVs settings, valid and invalid IVs can be simultaneously and individually weak. It is a much more relaxed condition than the fixed p case. This many individually weak IVs case has been numerically demonstrated in simulation study Case 2.

Appendix A4. SAIS condition and intuition of why MCP can circumvent it?

Correct selection of valid IVs is a more subtle and important issue in IV content. Recall Windmeijer et al. (2019) indicated that failure of consistent variable selection of sisVIVE is assured if the SAIS condition holds. The SAIS condition was first proposed in Windmeijer et al. (2019) derived from Irrepresentable Condition (IRC) directly. IRC is known as (almost) necessary and sufficient condition for variable selection consistency of Lasso (Zhao and Yu, 2006) for $n^{-1} \tilde{\mathbf{Z}}' \tilde{\mathbf{Z}}$, i.e.,

$$\max_{j \in \mathcal{V}^*} \left\| \left(\tilde{\mathbf{Z}}_{V^{c*}}^\top \tilde{\mathbf{Z}}_{V^{c*}} \right)^{-1} \tilde{\mathbf{Z}}_{V^{c*}}^\top \tilde{\mathbf{Z}}_j \right\|_1 \leq \xi < 1, \text{ for some } \xi \in [0, 1). \quad (\text{A3})$$

In the standard Lasso regression problem, the IRC only relates to the design matrix and holds for many standard designs (see corollaries in [Zhao and Yu, 2006](#)). However, in the context of IV model, the IRC on $\tilde{\mathbf{Z}}$ (instead of \mathbf{Z}) involves the first-stage signal estimate $\hat{\gamma}$, which further complicates the verifiability of the SAIS condition. Typically, two-stage IV regression modeling exacerbates the difficulty of detecting valid IVs through penalized methods than support recovery problem in a simple linear model. Moreover, among the penalty functions, Lasso penalty aggravates the problem if the first-stage coefficients related SAIS condition hold.

The MCP penalty inherits a much weaker condition for oracle property than IRC that Lasso required ([Zhao and Yu, 2006](#)). Theorem 6 proposed in ([Zhang and Zhang, 2012](#)) has generalized the IRC in Lasso to a concave penalty in a linear regression problem. We briefly state the key result by defining two quantities:

$$\begin{aligned}\theta_{\text{select}} &= \inf \left\{ \theta : \left\| \left(\frac{\tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}}{n} \right)^{-1} p'_\lambda(\boldsymbol{\varphi}_{\mathcal{V}^{c*}} + \hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}}) \right\|_\infty \leq \theta \lambda, \forall \|\boldsymbol{\varphi}_{\mathcal{V}^{c*}}\|_\infty \leq \theta \lambda \right\}, \\ \kappa_{\text{select}} &= \sup \left\{ \|\tilde{\mathbf{Z}}_{\mathcal{V}^*}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}} (\tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}})^{-1} p'_\lambda(\boldsymbol{\varphi}_{\mathcal{V}^{c*}} + \hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}})\|_\infty / \lambda : \|\boldsymbol{\varphi}_{\mathcal{V}^{c*}}\|_\infty \leq \theta_{\text{select}} \lambda \right\},\end{aligned}$$

where $\boldsymbol{\varphi}_{\mathcal{V}^{c*}}$ is a $|\mathcal{V}^{c*}|$ -vector, and let the $\hat{\boldsymbol{\alpha}}^{\text{or}}$ to be the oracle estimate, i.e., $\hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}} = (\tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}})^{-1} \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \mathbf{Y}$ and $\hat{\boldsymbol{\alpha}}_{\mathcal{V}^*}^{\text{or}} = \mathbf{0}$. The extended IRC for concave penalty required $\kappa_{\text{select}} < 1$ as the most crucial one to achieve selection consistency. When replacing MCP penalty $p_\lambda^{\text{MCP}}(\boldsymbol{\alpha})$ with Lasso penalty $\lambda \|\boldsymbol{\alpha}\|_1$ whose coordinate sub-derivative lies in $[-\lambda, \lambda]$, extended condition will be reduced to $\kappa_{\text{select}}(\boldsymbol{\alpha}) = \|\tilde{\mathbf{Z}}_{\mathcal{V}^*}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}} (\tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}})^{-1}\|_\infty < 1$ as identical to IRC of Lasso ([A3](#)). By the feature of nearly unbiasedness characteristic of MCP, $p'_\lambda(t) = 0 \ \forall |t| > \lambda\rho$. Once the mild condition $\min_{j \in \mathcal{V}^{c*}} |\hat{\alpha}_j^{\text{or}}| > \lambda\rho$ holds, it implies $\theta_{\text{select}} = 0$ with $\boldsymbol{\varphi}_{\mathcal{V}^{c*}} = \mathbf{0}$ and $\kappa_{\text{select}} = 0$ sequentially. Thus the extended IRC $\kappa_{\text{select}} < 1$ holds automatically for MCP but does not for Lasso. Consequently, this property is desirable that MCP achieves exact support recovery regardless of the constraint of SAIS condition.

Appendix A5. I-LAMM Algorithm for MCP penalty

Theoretically, the proposed WIT estimator enjoys better performance under weaker regulation conditions for low- and high-dimension cases with weak IVs. [Fan et al. \(2018\)](#) proposed the iterative local adaptive majorize-minimization (I-LAMM) algorithm for non-convex regularized model. I-LAMM first contracts the initial values in the neighborhood of the optimum solutions to serve as a better sparse coarse initial $\hat{\boldsymbol{\alpha}}^{(1)}$, then tightens it to the solution under precision tolerance. The computation is in polynomial time.

To be concrete, the I-LAMM algorithm combines the adaptive Local Linear Approximation (LLA) method and proximal gradient method or iterative shrinkage-thresholding (ISTA) algorithm ([Daubechies et al., 2004](#)). We adopt the I-LAMM to solve a sequence of optimization problems through LLA,

$$\hat{\boldsymbol{\alpha}}^{(t)} = \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \frac{1}{2n} \|\mathbf{Y} - \tilde{\mathbf{Z}}\boldsymbol{\alpha}\|_2^2 + \sum_{j=1}^p \left[p'_\lambda(|\hat{\alpha}_j^{(t-1)}|) |\alpha_j| \right]. \quad (\text{A4})$$

For $t = 1$ refers to contraction stage to obtain the better initial estimators which fall in the contraction region. And $t = 2, 3, \dots$, is the tightening stage, until it converges. It is worth noting Zou and Li (2008) used one-step LLA iteration to achieve the oracle estimator based on OLS estimate as initial. Nevertheless, in our case, the OLS estimate is not achievable under the perfect multicollinearity of design matrix $\tilde{\mathbf{Z}}$. Windmeijer et al. (2019) used the adaptive Lasso with embedded information in the median estimator as the one-step LLA to achieve the oracle property. Compared with the one-step method, iterations of $\boldsymbol{\lambda}^{(t-1)}$ (defined as below) in I-LAMM circumvent the stringent condition (namely, the majority rule) to obtain the root n initial estimator as the proper adaptive weight for ALasso.

For each iteration (A4), the ISTA method is implemented to achieve the closed-form updating formula, which is the reduced model in Fan et al. (2018)'s derivation. For a given iteration t , let $k = 0, 1, 2, \dots$ denotes the iteration in the proximal gradient updating, thus,

$$\begin{aligned}\hat{\boldsymbol{\alpha}}^{(t,k)} &= \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \left\{ \frac{1}{2n} \left\| \mathbf{Y} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\alpha}}^{(t,k-1)} \right\|_2^2 - \frac{1}{n} \left[\tilde{\mathbf{Z}} (\boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}^{(t,k-1)}) \right]^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\alpha}}^{(t,k-1)}) \right. \\ &\quad \left. + \frac{\phi}{2} \left\| \boldsymbol{\alpha} - \hat{\boldsymbol{\alpha}}^{(t,k-1)} \right\|_2^2 + \sum_{j=1}^p \left[p'_\lambda(|\hat{\alpha}_j^{(t-1)}|) |\alpha_j| \right] \right\} \\ &= \underset{\boldsymbol{\alpha}}{\operatorname{argmin}} \left\{ \frac{\phi}{2} \left\| \boldsymbol{\alpha} - \left(\hat{\boldsymbol{\alpha}}^{(t,k-1)} + \frac{1}{\phi n} \tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\alpha}}^{(t,k-1)}) \right) \right\|_2^2 + \sum_{j=1}^p \left[p'_\lambda(|\hat{\alpha}_j^{(t-1)}|) |\alpha_j| \right] \right\} \\ &= S \left(\hat{\boldsymbol{\alpha}}^{(t,k-1)} + \frac{1}{\phi n} \tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\alpha}}^{(t,k-1)}), \frac{1}{\phi} \boldsymbol{\lambda}^{(t-1)} \right),\end{aligned}\tag{A5}$$

where ϕ should be no smaller than the largest eigenvalue of $\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}$, or simply put, to ensure the majorization and $S(\mathbf{x}, \mathbf{a})$ denotes the component-wise soft-threshoding operator, i.e., $S(\mathbf{x}, \mathbf{a})_j = \operatorname{sgn}(x_j)(|x_j| - a_j)_+$, with $\boldsymbol{\lambda}^{(t-1)} = \left(p'_n(|\hat{\alpha}_1^{(t-1)}|), \dots, p'_n(|\hat{\alpha}_p^{(t-1)}|) \right)^\top$.

We adopt the first order optimality condition as a stopping criterion in the sub-problem. Let

$$\omega_{\boldsymbol{\lambda}^{(t-1)}}(\boldsymbol{\alpha}) = \min_{\boldsymbol{\xi} \in \partial|\boldsymbol{\alpha}|} \left\{ \left\| -\tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \boldsymbol{\alpha}) + \boldsymbol{\lambda}^{(t-1)} \odot \boldsymbol{\xi} \right\|_\infty \right\}$$

as a natural measure of suboptimality of $\boldsymbol{\alpha}$, where \odot is the Hadamard product. Once $\omega_{\hat{\boldsymbol{\lambda}}^{(t-1)}}(\boldsymbol{\alpha}^{(t,k)}) \leq \delta$, where δ is a pre-determined tolerance level and is assigned δ_c and δ_t for contraction and tightening stages, respectively, we stop the inner iteration and take $\hat{\boldsymbol{\alpha}}^{(t)} = \hat{\boldsymbol{\alpha}}^{(t,k)}$ as the δ -optimal solution in the sub-problem. Notably, it is an early stopped variant of ISTA method in each sub-problem to obtain $\hat{\boldsymbol{\alpha}}^{(t)}$ from $\hat{\boldsymbol{\alpha}}^{(t-1)}$. The following Algorithm 2 demonstrates the details of I-LAMM algorithm for WIT.

Appendix A6. Additional Information on Real Data Analysis

The following table A1 provide the detailed summary statistics of variables used in our empirical analysis.

Algorithm 2 I-LAMM algorithm for $\hat{\alpha}$ with MCP penalty

Input: $Y, \tilde{Z}, \hat{\alpha}^{(0)}, \lambda, \phi = \lambda_{\max}(\tilde{Z}^\top \tilde{Z}), \delta_c = 10^{-3}, \delta_t = 10^{-5}$

- 1: **for** $t = 1, 2, \dots$ **do**
- 2: $\lambda^{(t-1)} = (p'_\lambda(|\hat{\alpha}_1^{(t-1)}|), \dots, p'_\lambda(|\hat{\alpha}_p^{(t-1)}|))^\top$ $\triangleright p'_n$ is the derivative of MCP penalty
- 3: $\hat{\alpha}^{(t,0)} = \hat{\alpha}^{(t-1)}$
- 4: **for** $k = 1, 2, \dots$ **do**
- 5: $\hat{\alpha}^{(t,k)} = S(\hat{\alpha}^{(t,k-1)} + \frac{1}{n\phi} \tilde{Z}^\top (Y - \tilde{Z} \hat{\alpha}^{(t,k-1)}), \frac{1}{\phi} \lambda^{(t-1)})$ \triangleright LAMM updating
- 6: $\omega_{\lambda^{(t-1)}}(\hat{\alpha}^{(t,k)}) = \min_{\xi \in \partial|\hat{\alpha}^{(t,k)}|} \left\{ \left\| -\tilde{Z}^\top (Y - \tilde{Z} \hat{\alpha}^{(t,k)}) + \lambda^{(t-1)} \odot \xi \right\|_\infty \right\}$
- 7: **if** $\omega_{\lambda^{(t-1)}}(\hat{\alpha}^{(t,k)}) \leq \mathbf{1}(t=1)\delta_c + \mathbf{1}(t \neq 1)\delta_t$ **then**
- 8: $\hat{\alpha}^{(t)} = \hat{\alpha}^{(t,k)}$
- 9: **break**
- 10: **end if**
- 11: **end for**
- 12: **if** $\|\hat{\alpha}^{(t)} - \hat{\alpha}^{(t-1)}\|_\infty \leq \delta_t$ **then**
- 13: $\hat{\alpha} = \hat{\alpha}^{(t)}$
- 14: **break**
- 15: **end if**
- 16: **end for**

Output: $\hat{\alpha}$

Appendix B. Proofs

In this section we provide the proofs of the theoretical results in the main text. To proceed, some lemmas from previous literature are needed. We restate some of those results for the convenience of the readers.

Appendix B1. Ancillary Lemmas

LEMMA A1. (Lemmas 1 and 2 in [Bekker, 1994](#) and Lemma A.1 in [Kolesár et al., 2015](#)). Consider the quadratic form $Q = (M + U)^\top C(M + U)$, where $M \in \mathbb{R}^{n \times S}, C \in \mathbb{R}^{n \times n}$ are non-stochastic, C is symmetric and idempotent with rank J_n which may depend on n , and $U = (u_1, \dots, u_n)^\top$, with $u_i \sim_{i.i.d.} [0, \Omega]$. Let $a \in \mathbb{R}^S$ be a non-stochastic vector. Then,

(a) If u_i has finite fourth moment:

$$\begin{aligned}
\mathbb{E}[Q \mid C] &= M^\top C M + J_N \Omega \\
\text{var}(Qa \mid C) &= a^\top \Omega a M^\top C M + a^\top M^\top C M a \Omega + \Omega a a^\top M^\top C M + M C M a a^\top \Omega + J_N \left(a^\top \Omega a \Omega + \Omega a a^\top \Omega \right) \\
&\quad + d_C^\top d_C \left[\mathbb{E} \left((a^\top u)^2 u u^\top - a^\top \Omega a a^\top \Omega - a^\top \Omega a \Omega \right) + 2 d_C^\top C M a \mathbb{E} \left[(a^\top u) u u^\top \right] \right. \\
&\quad \left. + M^\top C d_C \mathbb{E} \left[(a^\top u)^2 u^\top \right] + \mathbb{E} \left[(a^\top u)^2 u \right] d_C^\top C M \right]
\end{aligned}$$

where $d_C = \text{diag}(C)$. If the distribution of u_i is normal, the last two lines of the variance component equal zero.

Table A1. Summary statistics of main variables

| | Notation | Type | Mean | Std | Median | Min | Max |
|-----------------------------|------------|---------------------|--------|--------|--------|--------|--------|
| log(GDP) | log(Y) | Response | 10.177 | 1.0102 | 10.416 | 7.463 | 12.026 |
| Trade | T | Endogenous Variable | 0.866 | 0.520 | 0.758 | 0.198 | 4.128 |
| log(Population) | S_1 | Control Variable | 1.382 | 1.803 | 1.480 | -3.037 | 6.674 |
| log(Land Area) | S_2 | Control Variable | 11.726 | 2.260 | 12.015 | 5.680 | 16.611 |
| \hat{T} (proxy for trade) | Z_1 | IV | 0.093 | 0.052 | 0.079 | 0.015 | 0.297 |
| log(Water Area) | Z_2 | IV | 6.756 | 3.654 | 7.768 | 0 | 13.700 |
| log(Land Boundaries) | Z_3 | IV | 6.507 | 2.920 | 7.549 | 0 | 10.005 |
| % Forest | Z_4 | IV | 29.89 | 22.380 | 30.62 | 0 | 98.26 |
| % Arable Land | Z_5 | IV | 40.947 | 21.549 | 42.062 | 0.558 | 82.560 |
| Languages | Z_6 | IV | 1.873 | 2.129 | 1 | 1 | 16 |
| Annual Freshwater | Z_7 | IV | 2.190 | 2.129 | 2.155 | -2.968 | 8.767 |

Source: FR99, the World Bank, and CIA world Factbook.

(b) Suppose that the distribution of u_i is normal, and as $n \rightarrow \infty$,

$$M^\top CM/N \rightarrow Q_{CM}, J_n/n \rightarrow \alpha_r$$

where the elements c_{is} of C may depend on N . Then,

$$\sqrt{n}(Qa/n - \mathbb{E}Qa/n) \xrightarrow{d} \mathcal{N}(0, V),$$

where $V = a^\top \Omega a Q_{CM} + a^\top Q_{CM} a \Omega + \Omega a a^\top Q_{CM} + Q_{CM} a a^\top \Omega + \alpha_r (a^\top \Omega a \Omega + \Omega a a^\top \Omega)$.

Appendix B2. Proof of Theorem 1

PROOF. Firstly, we prove that procedure (6) can generate the G different groups of \mathcal{P}_c satisfying the requirements. With direct calculation, $\mathcal{P}_c = \{\tilde{\beta}^c, \tilde{\alpha}^c, \tilde{\epsilon}^c\} = \{\beta^* + c, \alpha^* - c\gamma^*, \epsilon - c\eta\}$ and $E(\tilde{\epsilon}^c) = E(\epsilon) - cE(\eta) = \mathbf{0}$ for $c = \alpha_j^*/\gamma_j^*, j \in \mathcal{I}_c$. Therefore, $\tilde{\alpha}_{\mathcal{I}_c}^c = \alpha_{\mathcal{I}_c}^* - c\gamma_{\mathcal{I}_c}^* = \alpha_{\mathcal{I}_c}^* - \alpha_{\mathcal{I}_c}^* = \mathbf{0}$ and $\tilde{\alpha}_j^c = \alpha_j^* - c\gamma_j^* \neq 0$ for $j \notin \mathcal{I}_c$. Going through all possible $c = \{\alpha_j^*/\gamma_j^* : j \notin \mathcal{V}^*\}$, we conclude that procedure (6) has generated G groups additional DGP $_c$. The exhaustive and mutually exclusive property of \mathcal{I}_c and \mathcal{V}^* is also guaranteed by its construction.

Second, we show the proof by contradiction that there is no other possible DGP with the same sparse structure and zero mean structural error.

Assume there is an additional DGP $\{\tilde{\beta}, \tilde{\alpha}, \tilde{\epsilon}\}$ differentiating with \mathcal{P}_c and \mathcal{P}_0 but still has $E(\tilde{\epsilon}) = \mathbf{0}$ and zero(s) component in $\tilde{\alpha}$, i.e., $\tilde{\mathcal{I}} = \{j : \tilde{\alpha}_j = 0\} \neq \emptyset$. By the property of exhaustive and mutually exclusiveness of \mathcal{V}^* and $\{\mathcal{I}_c\}_{c \neq 0}$, WLOG, we assume $\mathcal{I}_g \cap \tilde{\mathcal{I}} \neq \emptyset$ for some g and DGP $_g = \{\tilde{\beta}^g, \tilde{\alpha}^g, \tilde{\epsilon}^g\}$. Since $E(\tilde{\epsilon}) = E(\tilde{\epsilon}^g) = \mathbf{0}$, it suffices to show the contradiction in terms of moment condition (3). Hence, $\{\tilde{\beta}^g, \tilde{\alpha}^g\}$ and $\{\tilde{\beta}, \tilde{\alpha}\}$ are solutions of

$$\Gamma^* = \alpha + \beta\gamma^*.$$

For $j \in \mathcal{I}_g \cap \tilde{\mathcal{I}}$, $\Gamma_j^* = \tilde{\beta}^g\gamma_j^* = \tilde{\beta}\gamma_j^*$ derives $\tilde{\beta}^g = \tilde{\beta}$. In turn, $\forall j \notin \mathcal{I}_g \cap \tilde{\mathcal{I}}$,

$$\Gamma_j^* - \tilde{\beta}\gamma_j^* = \tilde{\alpha}_j = \check{\alpha}_j.$$

Thus, $\{\tilde{\beta}^g, \tilde{\alpha}^g\}$ and $\{\check{\beta}, \check{\alpha}\}$ are equivalent. So as $\tilde{\epsilon}^g$ and $\check{\epsilon}$ since

$$\tilde{\epsilon}^c = Y - D\tilde{\beta}^g - Z\tilde{\alpha}^g = Y - D\check{\beta} - Z\check{\alpha} = \check{\epsilon}.$$

Hence DGP $\{\check{\beta}, \check{\alpha}, \check{\epsilon}\}$ is equivalent to DGP_g and forms a contradiction. It concludes that procedure (6) can produce all possible DGPs. \square

Appendix B3. Proof of Theorem 2

PROOF. Recall the construction that $\exists i \in \{0, \dots, G\}$, let $\mathcal{F} = \{f : \mathcal{P} \in \mathcal{Q} \rightarrow \mathbb{R}; f(\mathcal{P}_i) < f(\mathcal{P}_j), \forall j \neq i\}$ and $\mathcal{G} = \{g = \arg\min_{\mathcal{P} \in \mathcal{Q}} f(\mathcal{P}); f \in \mathcal{F}\}$. To show any element in \mathcal{G} is sufficient to identify one specific DGP in \mathcal{Q} , we assume \mathcal{P}_m and \mathcal{P}_n are both elements in \mathcal{Q} but with different minimum values under the same $f \in \mathcal{F}$ and its corresponding $g \in \mathcal{G}$. According to Theorem 1, β must be different in \mathcal{P}_m and \mathcal{P}_n if $\mathcal{P}_m \neq \mathcal{P}_n$. Hence,

- (1) $\mathcal{P}_m = \mathcal{P}_n$, it leads to $f(\mathcal{P}_m) = f(\mathcal{P}_n) < f(\mathcal{P}_j)$ for any $j \neq m, n$ and the rules \mathcal{P}_m and \mathcal{P}_n share the same β .
- (2) $\mathcal{P}_m \neq \mathcal{P}_n$, apply the minimum value assumption of f on $\mathcal{P}_m \neq \mathcal{P}_n$ respectively, we obtain

$$f(\mathcal{P}_m) < f(\mathcal{P}_i), \quad \forall i \neq m \quad \text{and} \quad f(\mathcal{P}_n) < f(\mathcal{P}_j), \quad \forall j \neq n. \quad (\text{B6})$$

Taking $i = m, j = n$, it forms a contradiction that

$$f(\mathcal{P}_m) < f(\mathcal{P}_n) < f(\mathcal{P}_m) \quad (\text{B7})$$

Therefore, \mathcal{P}_m and \mathcal{P}_n must be equivalent and identifiable under f .

Due to the arbitrary choice of $f \in \mathcal{F}$, it concludes part (a) $\mathcal{G} \subseteq \mathcal{H}$.

Now we move to part (b): it states that if $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$ and such h is the necessary condition for identifying β^* , then it must have $|\mathcal{H}| = 1$ and $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$. Again, we prove this by contradiction. Consider two cases, (1) $\forall h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_k$, where $\mathcal{P}_k \neq \mathcal{P}_0$ for some k , and (2) $1 < |\mathcal{H}| \leq G + 1$ with $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$:

- (1) $\forall h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_k$, where $\mathcal{P}_k \neq \mathcal{P}_0$ for some k . It directly forms a contradiction of $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$.
- (2) $1 < |\mathcal{H}| \leq G + 1$ with $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$. There exist at least two distinct mappings $h_m(\mathcal{Q}) \cong \mathcal{P}_m$, $h_n(\mathcal{Q}) \cong \mathcal{P}_n$ and $\mathcal{P}_m \neq \mathcal{P}_n$. By the assumption that there is a necessary condition $h_i \in \mathcal{H}$ of identifying β^* . Thus, h_i must be image equivalent to h_m or h_n . WLOG, we let $h_0 = h_m$ and $\mathcal{P}_m = \mathcal{P}_0 \neq \mathcal{P}_n$. Given the necessity of h_m , the contrapositive arguments must hold. That is, supposing $h_m(\mathcal{Q}) \neq \mathcal{P}_0$, it leads to there is no other mapping that maps \mathcal{Q} to \mathcal{P}_0 . However, the distinct image equivalent mapping h_n forms a contradiction that it is possible to pick $h_n(\mathcal{Q}) = \mathcal{P}_0$ since it only requires h_m to have a different image with h_n .

Together with above arguments, we conclude (b) that there is no necessary condition for the identification of β^* unless $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$ and $|\mathcal{H}| = 1$. \square

Appendix B4. Proof of Corollary 1

PROOF. To show Corollary 1, we only need to consider some specific constructions of $g \in \mathcal{G} \subseteq \mathcal{H}$. Let $f_i = \mathbf{1}(\alpha_i = 0)$ for $i = 1, 2, \dots, p$ to be the valid IV indicators. According to Theorem 1, the non-empty set \mathcal{I}_c and \mathcal{V}^* are exhaustive and have mutual exclusive “zeros structure” in $\tilde{\alpha}^c$ and α^* correspondingly. Thus, $\{f_i : i = 1, 2, \dots, p\}$ must enumerate all solutions. Hence,

$$G + 1 = |\{f_i : i = 1, 2, \dots, p\}| \leq |\mathcal{G}| \leq |\mathcal{H}| \leq G + 1 \quad (\text{B8})$$

holds. It leads to $|\mathcal{H}| = G + 1$. According to Theorem 2, there is not a necessary and sufficient condition to identify β^* , unless $G = 0$, i.e., all potential IVs are valid. \square

Appendix B5. Proof of Proposition 1

PROOF.

$$\begin{aligned} \alpha^* &= \underset{\mathcal{P}=\{\beta, \alpha, \epsilon\} \in \mathcal{Q}}{\operatorname{argmin}} p_{\lambda}^{\text{pen}}(\alpha), \\ \iff \sum_{j=1}^p p_{\lambda}^{\text{pen}}(\alpha_j^*) &< \sum_{j=1}^p p_{\lambda}^{\text{pen}}(\tilde{\alpha}_j^c) \\ \iff \sum_{j \in \mathcal{V}^{c*}} p_{\lambda}^{\text{pen}}(\alpha_j^*) &< \sum_{j \in \mathcal{I}_c^c} p_{\lambda}^{\text{pen}}(\tilde{\alpha}_j^c), \end{aligned} \quad (\text{B9})$$

where $\mathcal{I}_c = \{j : \alpha_j^*/\gamma_j^* = c, c \neq 0\}$ and \mathcal{I}_c^c is the complement of \mathcal{I}_c . By the Assumption 6: $\alpha^* = \underset{\mathcal{P}=\{\beta, \alpha, \epsilon\} \in \mathcal{Q}}{\operatorname{argmin}} \|\alpha\|_0$, we have $|\mathcal{V}^{c*}| < |\mathcal{I}_c^c|$.

Define $\pi(\mathcal{I}_c^c, \mathcal{V}^{c*})$ as $|\mathcal{V}^{c*}|$ -combination of \mathcal{I}_c^c . Now we rewrite (B9) as

$$\sum_{j \in \mathcal{V}^{c*}} p_{\lambda}^{\text{pen}}(\alpha_j^*) - \sum_{k \in \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} p_{\lambda}^{\text{pen}}(\tilde{\alpha}_k^c) < \sum_{l \in \mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} p_{\lambda}^{\text{pen}}(\tilde{\alpha}_l^c), \quad (\text{B10})$$

and it should hold for any $|\mathcal{V}^{c*}|$ -combination and \mathcal{P}_0 . Note that the RHS in (B10) are non-negative by definition of the penalty function and $|\mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})| > 0$.

For arbitrary $\pi(\mathcal{I}_c^c, \mathcal{V}^{c*})$ and \mathcal{P}_0 , we consider the worst case where

$$|\tilde{\alpha}_l^c| < |\tilde{\alpha}_k^c| \text{ and } |\tilde{\alpha}_k^c| < |\alpha_j^*|,$$

for $\forall k \in \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}), l \in \mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})$ and for each pair $(j, k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))$. Using Taylor expansion and mean value theorem, we obtain

$$\begin{aligned} &\left\{ \sum_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} |\alpha_j^*| - |\tilde{\alpha}_k^c| \right\} \min_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} [p_{\lambda}^{\text{pen}'}(\bar{\xi}_{(j,k)})] \\ &< \sum_{j \in \mathcal{V}^{c*}} p_{\lambda}^{\text{pen}}(\alpha_j^*) - \sum_{k \in \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} p_{\lambda}^{\text{pen}}(\tilde{\alpha}_k^c) \\ &< \sum_{l \in \mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} p_{\lambda}^{\text{pen}}(\tilde{\alpha}_l^c) \\ &< \|\tilde{\alpha}_{\mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})}^c\|_1 \max_{l \in \mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} [p_{\lambda}^{\text{pen}'}(\bar{\xi}_l)], \end{aligned}$$

where $\bar{\xi}_{(j,k)} \in (|\tilde{\alpha}_k^c|, |\alpha_j^*|)$ and $\bar{\xi}_l \in (0, |\tilde{\alpha}_l^c|)$. Thus, it leads to

$$\frac{\min_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} [p_\lambda^{\text{pen}'}(\bar{\xi}_{(j,k)})]}{\max_{l \in \mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} [p_\lambda^{\text{pen}'}(\bar{\xi}_l)]} < \frac{\|\tilde{\alpha}_{\mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})}^c\|_1}{\left\{ \sum_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} |\alpha_j^*| - |\tilde{\alpha}_k^c| \right\}}. \quad (\text{B11})$$

Because $p_\lambda^{\text{pen}'}$ is free of \mathcal{P}_0 , the RHS in (B11) could be much smaller than 1 in some extreme case leads to $p_\lambda^{\text{pen}'}$ being a bounded monotone decreasing function. That is to say, p_λ^{pen} should be a concave penalty.

Consider another case where we keep α^* fixed, but vary γ^* to make $\tilde{\alpha}^c$ to have the same order with $\kappa(n)$ defined in Assumption 5. Again by (B11),

$$\begin{aligned} \min_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} [p_\lambda^{\text{pen}'}(\bar{\xi}_{(j,k)})] &< \frac{\|\tilde{\alpha}_{\mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})}^c\|_1}{\left\{ \sum_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} |\alpha_j^*| - |\tilde{\alpha}_k^c| \right\}} \cdot \max_{l \in \mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})} [p_\lambda^{\text{pen}'}(\bar{\xi}_l)] \\ &< \frac{\lambda \|\tilde{\alpha}_{\mathcal{I}_c^c / \pi(\mathcal{I}_c^c, \mathcal{V}^{c*})}^c\|_1}{\left\{ \sum_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} |\alpha_j^*| - |\tilde{\alpha}_k^c| \right\}} \\ &\asymp \frac{\lambda \kappa(n) (|\mathcal{I}_c^c| - |\mathcal{V}^{c*}|)}{\left\{ \sum_{(j,k) \in (\mathcal{V}^{c*}, \pi(\mathcal{I}_c^c, \mathcal{V}^{c*}))} |\alpha_j^*| - |\tilde{\alpha}_k^c| \right\}} \asymp C \lambda \kappa(n). \end{aligned}$$

Thus, it concludes that $p_\lambda^{\text{pen}'}(t) = O(\lambda \kappa(n))$ for any $t > \kappa(n)$. \square

Appendix B6. Proof of Lemma 1

PROOF. $\tilde{\mathbf{Z}}_j = M_{\hat{\mathbf{D}}} \mathbf{Z}_j$, \mathbf{Z}_j is the j -th instrument. The target is to analyze that the rate of $\tilde{\mathbf{Z}}_j^\top \boldsymbol{\epsilon} / n$ is not inflated. To proceed, $|\tilde{\mathbf{Z}}_j^\top \boldsymbol{\epsilon} / n| \leq |\mathbf{Z}_j^\top \boldsymbol{\epsilon} / n| + |\mathbf{Z}_j^\top P_{\hat{\mathbf{D}}} \boldsymbol{\epsilon} / n|$. The first term is known to be $O_p(n^{-1/2})$. For the second term, we have,

$$|\mathbf{Z}_j^\top P_{\hat{\mathbf{D}}} \boldsymbol{\epsilon} / n| = \underbrace{\frac{|\mathbf{Z}_j^\top \hat{\mathbf{D}}|}{n}}_{(I)} \cdot \underbrace{\frac{|\boldsymbol{\epsilon}^\top \hat{\mathbf{D}}|}{n}}_{(II)} \cdot \left(\underbrace{\frac{|\hat{\mathbf{D}}^\top \hat{\mathbf{D}}|}{n}}_{(III)} \right)^{-1}. \quad (\text{B12})$$

Notably, $\hat{\mathbf{D}} = \mathbf{Z} \boldsymbol{\gamma}^* + P_{\mathbf{Z}} \boldsymbol{\eta}$ is a projected endogenous variable. It consists of the random and non-random parts. Then we analyze the size of the above three terms sequentially. For (I), we have

$$\frac{\mathbf{Z}_j^\top \hat{\mathbf{D}}}{n} = \frac{\mathbf{Z}_j^\top (\mathbf{Z} \boldsymbol{\gamma}^* + \boldsymbol{\eta})}{n} = (\mathbf{Q}_{nj})^\top \boldsymbol{\gamma}^* + \frac{\mathbf{Z}_j^\top \boldsymbol{\epsilon}}{n} = \mathbf{Q}_{nj}^\top \boldsymbol{\gamma}^* + O_p(n^{-1/2}). \quad (\text{B13})$$

With regard to (II),

$$\frac{\boldsymbol{\epsilon}^\top \hat{\mathbf{D}}}{n} = \frac{\boldsymbol{\epsilon}^\top P_{\mathbf{Z}} (\mathbf{Z} \boldsymbol{\gamma}^* + \boldsymbol{\eta})}{n} = \frac{\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*}{n} + \frac{\boldsymbol{\epsilon}^\top P_{\mathbf{Z}} \boldsymbol{\eta}}{n}. \quad (\text{B14})$$

Within this decomposition, we first have $E(\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n) = E(E(\boldsymbol{\epsilon}^\top | \mathbf{Z}) \mathbf{Z} \boldsymbol{\gamma}^*) = 0$ and

$$\text{var}(\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n) = E(\text{var}(\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n | \mathbf{Z})) = E(\sigma_\epsilon^2 \boldsymbol{\gamma}^{*\top} \mathbf{Z}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n^2) = O(\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n).$$

Thus we obtain $\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n = O_P(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n})$. Regarding the second term, we have $E(\boldsymbol{\epsilon}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n) = E(\text{tr}(P_{\mathbf{Z}} \boldsymbol{\eta} \boldsymbol{\epsilon}^\top)/n) = \text{tr}(E(P_{\mathbf{Z}} E(\boldsymbol{\eta} \boldsymbol{\epsilon}^\top | \mathbf{Z}))/n) = \sigma_{\epsilon, \eta}^2 p/n$ and

$$\begin{aligned} \text{var}(\boldsymbol{\epsilon}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n) &= E(\text{var}(\boldsymbol{\epsilon}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n | \mathbf{Z})) = E(2\sigma_\epsilon^2 \sigma_\eta^2 [\text{tr}(P_{\mathbf{Z}})/n^2]) + E([n^{-2} \sum_i (P_{\mathbf{Z}})_{ii}^2] ([\sigma_{\epsilon, \eta}^2]^2 - 2\sigma_\epsilon^2 \sigma_\eta^2)) \\ &\leq O(p/n^2) + O(p/n^2) = O(p/n^2), \end{aligned}$$

where the inequality holds since $\text{tr}(P_{\mathbf{Z}}) = p$, $(P_{\mathbf{Z}})_{ii} \in (0, 1)$ and $\sum_i (P_{\mathbf{Z}})_{ii}^2 \leq \sum_i (P_{\mathbf{Z}})_{ii} = p$. Thus, we conclude the size of (II) is $O_P(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2}) + \sigma_{\epsilon, \eta}^2 p/n$.

Now we turn to (III), with similar argument,

$$\begin{aligned} \frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n} &= \frac{\mathbf{D}^\top P_{\mathbf{Z}} \mathbf{D}}{n} = \frac{\boldsymbol{\gamma}^{*\top} \mathbf{Z}^\top \mathbf{Z} \boldsymbol{\gamma}^*}{n} + 2 \frac{\boldsymbol{\eta}^\top \mathbf{Z} \boldsymbol{\gamma}^*}{n} + \frac{\boldsymbol{\eta}^\top P_{\mathbf{Z}} \boldsymbol{\eta}}{n} \\ &= \boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + o_p(1) + O_P(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2}) + \sigma_\eta^2 p/n \\ &= \boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + \sigma_\eta^2 p/n + O_P(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2}). \end{aligned}$$

Therefore, together with above three terms, we are able to derive the size of $|\mathbf{Z}_j^\top P_{\widehat{\mathbf{D}}} \boldsymbol{\epsilon}/n|$,

$$\begin{aligned} \mathbf{Z}_j^\top P_{\widehat{\mathbf{D}}} \boldsymbol{\epsilon}/n &= \frac{[\mathbf{Q}_{nj}^\top \boldsymbol{\gamma}^* + O_p(n^{-1/2})] \cdot [\sigma_{\epsilon, \eta}^2 p/n + O_P(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2})]}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + \sigma_\eta^2 p/n + O_P(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2})} \\ &= \sigma_{\epsilon, \eta}^2 p/n \cdot \frac{\mathbf{Q}_{nj}^\top \boldsymbol{\gamma}^*}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + \sigma_\eta^2 p/n} + O_p(n^{-1/2}) \end{aligned} \quad (\text{B15})$$

□

Appendix B7. Proof of Lemma 2

PROOF. We have the transformed $\tilde{\mathbf{Z}} = M_{\widehat{\mathbf{D}}} \mathbf{Z}$, where $\widehat{\mathbf{D}} = P_{\mathbf{Z}} \mathbf{D} = \mathbf{Z} \hat{\boldsymbol{\gamma}}$. Thus,

$$\begin{aligned} \tilde{\mathbf{Z}} &= \mathbf{Z} - \mathbf{Z} \hat{\boldsymbol{\gamma}} \left(\hat{\boldsymbol{\gamma}}^\top \mathbf{Z}^\top \mathbf{Z} \hat{\boldsymbol{\gamma}} \right)^{-1} \hat{\boldsymbol{\gamma}}^\top \mathbf{Z}^\top \mathbf{Z}, \\ \mathbf{C}_n &= \frac{\tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}}{n} = \left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n} \right) - \left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n} \right) \hat{\boldsymbol{\gamma}} \left(\hat{\boldsymbol{\gamma}}^\top \left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n} \right) \hat{\boldsymbol{\gamma}} \right)^{-1} \hat{\boldsymbol{\gamma}}^\top \left(\frac{\mathbf{Z}^\top \mathbf{Z}}{n} \right). \end{aligned} \quad (\text{B16})$$

Denote $\mathbf{Q}_n = \mathbf{Z}^\top \mathbf{Z}/n$, consider the square of restricted eigenvalue of $\tilde{\mathbf{Z}}$, $K_{\mathcal{C}}^2$, we have

$$\begin{aligned}
K_{\mathcal{C}}^2(\mathcal{V}^*, \xi) &= \inf_{\mathbf{u}} \{ \|\mathbf{u}^\top (n^{-1} \tilde{\mathbf{Z}}^\top \tilde{\mathbf{Z}}) \mathbf{u} / \|\mathbf{u}\|_2^2; \mathbf{u} \in \mathcal{C}(\mathcal{V}^*; \xi) \} \\
&= \inf_{\mathbf{u} \in \mathcal{C}(\mathcal{V}^*; \xi)} \{ \|\mathbf{u}^\top \mathbf{C}_n \mathbf{u} / \|\mathbf{u}\|_2^2 \} \\
&= \inf_{\mathbf{u} \in \mathcal{C}(\mathcal{V}^*; \xi)} \frac{\mathbf{u}^\top (\mathbf{Q}_n - \mathbf{Q}_n \hat{\gamma} (\hat{\gamma}^\top \mathbf{Q}_n \hat{\gamma})^{-1} \hat{\gamma}^\top \mathbf{Q}_n) \mathbf{u}}{\|\mathbf{u}\|_2^2} \\
&= \inf_{\mathbf{u} \in \mathcal{C}(\mathcal{V}^*; \xi)} \frac{\mathbf{u}^\top \mathbf{Q}_n \mathbf{u} \hat{\gamma}^\top \mathbf{Q}_n \hat{\gamma}^\top - (\hat{\gamma}^\top \mathbf{Q}_n \mathbf{u})^2}{\|\mathbf{u}\|_2^2 \hat{\gamma}^\top \mathbf{Q}_n \hat{\gamma}}
\end{aligned} \tag{B17}$$

Notice the denominator $\mathbf{u}^\top \mathbf{Q}_n \mathbf{u} \hat{\gamma}^\top \mathbf{Q}_n \hat{\gamma}^\top - (\hat{\gamma}^\top \mathbf{Q}_n \mathbf{u})^2 \geq 0$ by Cauchy–Schwarz inequality and the equality holds if and only if $\mathbf{u} = k \hat{\gamma}$ for any $k \neq 0$. Furthermore, the exact difference arises from this equation can be determined by the Lagrange’s identity. One can always take $\xi \in (0, \|\hat{\gamma}_{\mathcal{V}^*}\|_1 / \|\hat{\gamma}_{\mathcal{V}^{c*}}\|_1)$ such that the cone $\mathcal{C}(\mathcal{V}^*; \xi) = \{\mathbf{u} : \|\mathbf{u}_{\mathcal{V}^*}\|_1 \leq \xi \|\mathbf{u}_{\mathcal{V}^{c*}}\|_1\}$ excludes the membership of $k \hat{\gamma}$ because the cone $\mathcal{C}(\mathcal{V}^*; \xi)$ is invariant of scale. Therefore, the denominator $\mathbf{u}^\top \mathbf{Q}_n \mathbf{u} \hat{\gamma}^\top \mathbf{Q}_n \hat{\gamma}^\top - (\hat{\gamma}^\top \mathbf{Q}_n \mathbf{u})^2 > 0$ and $K_{\mathcal{C}}^2(\mathcal{V}^*, \xi) > 0$ holds strictly. \square

Appendix B8. Proof of Lemma 3

PROOF. Let $\mathbf{R}^* = \tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \alpha^*)/n$ and $\tilde{\mathbf{Z}} = M_{\widehat{\mathbf{D}}} \mathbf{Z}, \mathbf{D} = \mathbf{Z} \gamma^* + \boldsymbol{\eta}$, we have

$$\begin{aligned}
\mathbf{R}^* &= \tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \alpha^*)/n = \mathbf{Z}^\top M_{\widehat{\mathbf{D}}} (\mathbf{D} \beta^* + \boldsymbol{\epsilon})/n = \mathbf{Z}^\top \left(\mathbf{I} - \frac{\widehat{\mathbf{D}} \widehat{\mathbf{D}}^\top}{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}} \right) (\mathbf{D} \beta^* + \boldsymbol{\epsilon})/n \\
&= \beta^* \frac{\mathbf{Z}^\top \mathbf{D}}{n} + \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n} - \frac{\frac{\mathbf{Z}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top (\mathbf{D} \beta^* + \boldsymbol{\epsilon})}{n}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= \beta^* \frac{\mathbf{Z}^\top \mathbf{D}}{n} + \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n} - \beta^* \frac{\mathbf{Z}^\top \widehat{\mathbf{D}}}{n} - \frac{\frac{\mathbf{Z}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n} - \frac{\frac{\mathbf{Z}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}}.
\end{aligned} \tag{B18}$$

By direct algebra, using $\widehat{\mathbf{D}} = \mathbf{Z} \gamma^* + P_{\mathbf{Z}} \boldsymbol{\eta}$ and $\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}/n = \gamma^{*\top} \mathbf{Q}_n \gamma^* + 2 \boldsymbol{\eta}^\top \mathbf{Z} \gamma^*/n + \boldsymbol{\eta}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n$, we obtain

$$\mathbf{R}^* = \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n} - \frac{\frac{\mathbf{Z}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} = \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n} - \frac{(\mathbf{Q}_n \gamma^* + \mathbf{Z}^\top \boldsymbol{\eta}/n) (\boldsymbol{\epsilon}^\top \mathbf{Z} \gamma^*/n + \boldsymbol{\epsilon}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n)}{\gamma^{*\top} \mathbf{Q}_n \gamma^* + 2 \boldsymbol{\eta}^\top \mathbf{Z} \gamma^*/n + \boldsymbol{\eta}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n},$$

in which the expression is free of β^* . Thus, by triangle inequality, we attain the bound

$$\|\mathbf{R}\|_\infty \leq \left\| \frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n} \right\|_\infty + \left\| \frac{(\mathbf{Q}_n \gamma^* + \mathbf{Z}^\top \boldsymbol{\eta}/n) (\boldsymbol{\epsilon}^\top \mathbf{Z} \gamma^*/n + \boldsymbol{\epsilon}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n)}{\gamma^{*\top} \mathbf{Q}_n \gamma^* + 2 \boldsymbol{\eta}^\top \mathbf{Z} \gamma^*/n + \boldsymbol{\eta}^\top P_{\mathbf{Z}} \boldsymbol{\eta}/n} \right\|_\infty. \tag{B19}$$

Under standard argument through concentration inequality,

$$Pr\left(\left\|\frac{\mathbf{Z}^\top \boldsymbol{\epsilon}}{n}\right\|_\infty \geq t\right) = Pr\left(\max_{1 \leq j \leq p} |\mathbf{Z}_j^\top \boldsymbol{\epsilon}| \geq nt\right) \leq \sum_{1 \leq j \leq p} Pr(|\mathbf{Z}_j^\top \boldsymbol{\epsilon}| \geq nt) \leq 2p \exp\left(-\frac{nt^2}{2\sigma_\epsilon^2}\right).$$

Let $t = \sigma_\epsilon \sqrt{\frac{2}{n} \log(2p)}$, we obtain $\|\mathbf{Z}^\top \boldsymbol{\epsilon}/n\|_\infty = O_p\left(\sigma_\epsilon \sqrt{\frac{2}{n} \log(2p)}\right)$.

For the second term,

$$\left\|\frac{(\mathbf{Q}_n \boldsymbol{\gamma}^* + \mathbf{Z}^\top \boldsymbol{\eta}/n)(\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n + \boldsymbol{\epsilon}^\top \mathbf{P}_Z \boldsymbol{\eta}/n)}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + 2\boldsymbol{\eta}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n + \boldsymbol{\eta}^\top \mathbf{P}_Z \boldsymbol{\eta}/n}\right\|_\infty = \frac{\left|(\boldsymbol{\epsilon}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n + \boldsymbol{\epsilon}^\top \mathbf{P}_Z \boldsymbol{\eta}/n)\left\|\mathbf{Q}_n \boldsymbol{\gamma}^* + \mathbf{Z}^\top \boldsymbol{\eta}/n\right\|_\infty\right|}{\left|\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + 2\boldsymbol{\eta}^\top \mathbf{Z} \boldsymbol{\gamma}^*/n + \boldsymbol{\eta}^\top \mathbf{P}_Z \boldsymbol{\eta}/n\right|}$$

Similarly, we attain $\|\mathbf{Q}_n \boldsymbol{\gamma}^* + \mathbf{Z}^\top \boldsymbol{\eta}/n\|_\infty \leq \|\mathbf{Q}_n \boldsymbol{\gamma}^*\|_\infty + O_p\left(\sigma_\eta \sqrt{\frac{2}{n} \log(2p)}\right)$.

Also using the ancillary lemmas A1 and similar analyses in proof of Lemma 1, we know the second term in (B19) is upper-bounded by:

$$\begin{aligned} & \frac{\left(\|\mathbf{Q}_n \boldsymbol{\gamma}^*\|_\infty + O_p\left(\sigma_\eta \sqrt{\frac{2}{n} \log(2p)}\right)\right)\left(\sigma_{\epsilon, \eta} p/n + O_p\left(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2}\right)\right)}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + \sigma_\eta^2 p/n + O_p\left(\sqrt{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*/n} \vee \sqrt{p/n^2}\right)} \\ &= \sigma_{\epsilon, \eta} p/n \frac{\|\mathbf{Q}_n \boldsymbol{\gamma}^*\|_\infty}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + \sigma_\eta^2 p/n} + O_p\left(\sigma_\eta \sqrt{\frac{2}{n} \log(2p)}\right). \end{aligned} \quad (\text{B20})$$

Hence, together with above results, we obtain the final upper bound of $\|\mathbf{R}^*\|_\infty$:

$$\begin{aligned} \|\mathbf{R}^*\|_\infty &\leq \sigma_{\epsilon, \eta} p/n \frac{\|\mathbf{Q}_n \boldsymbol{\gamma}^*\|_\infty}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^* + \sigma_\eta^2 p/n} + O_p\left((\sigma_\eta + \sigma_\epsilon) \sqrt{\frac{2}{n} \log(2p)}\right) \\ &= O_p\left(\frac{p}{n} \cdot \frac{\|\mathbf{Q}_n \boldsymbol{\gamma}^*\|_\infty}{\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*} + \sqrt{\frac{\log p}{n}}\right), \end{aligned} \quad (\text{B21})$$

where the second line holds provided $\boldsymbol{\gamma}^{*\top} \mathbf{Q}_n \boldsymbol{\gamma}^*$ dominates $\sigma_{\epsilon, \eta} p/n$.

Similarly, $\mathbf{R}^{\text{or}} = \tilde{\mathbf{Z}}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \hat{\boldsymbol{\alpha}}^{\text{or}})/n = \tilde{\mathbf{Z}}^\top [\mathbf{Y} - \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}} (\tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}})^{-1} \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}}^\top \mathbf{Y}]/n$. Therefore, for the sake of controlling the supremum norm of \mathbf{R}^{or} , we only need to consider the valid IV \mathcal{V}^* since $\mathbf{R}_j^{\text{or}} = 0$ for $j \in \mathcal{V}^{c*}$. Recall $\hat{\beta}_{\text{or}}^{\text{TSLs}} = [\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D}]^{-1} [\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{Y}]$, we have

$$\begin{aligned} \hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}} &= (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} [\mathbf{Z}_{\mathcal{V}^{c*}}^\top (\mathbf{Y} - \hat{\mathbf{D}} \hat{\beta}_{\text{or}}^{\text{TSLs}})] \\ &= (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} \mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Y} - (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} \mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{D} [\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D}]^{-1} [\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{Y}] \\ &= \boldsymbol{\alpha}_{\mathcal{V}^{c*}}^* + (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} \mathbf{Z}_{\mathcal{V}^{c*}}^\top \left[\mathbf{I} - \mathbf{D} [\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D}]^{-1} \mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \right] \boldsymbol{\epsilon}. \end{aligned} \quad (\text{B22})$$

Thus, for any $j \in \mathcal{V}^*$, we obtain

$$\begin{aligned} \mathbf{R}_{\mathcal{V}^*}^{\text{or}} &= \tilde{\mathbf{Z}}_{\mathcal{V}^*}^\top \left\{ \mathbf{Y} - \tilde{\mathbf{Z}} (\hat{\boldsymbol{\alpha}}^{\text{or}} - \boldsymbol{\alpha}^* + \boldsymbol{\alpha}^*) \right\} / n \\ &= \tilde{\mathbf{Z}}_{\mathcal{V}^*}^\top (\mathbf{Y} - \tilde{\mathbf{Z}} \boldsymbol{\alpha}^*) / n + \tilde{\mathbf{Z}}_{\mathcal{V}^*}^\top \tilde{\mathbf{Z}}_{\mathcal{V}^{c*}} (\boldsymbol{\alpha}_{\mathcal{V}^{c*}}^* - \hat{\boldsymbol{\alpha}}_{\mathcal{V}^{c*}}^{\text{or}}) / n \\ &= \mathbf{R}_{\mathcal{V}^*}^* + \mathbf{Z}_{\mathcal{V}^*}^\top \mathbf{M}_{\hat{\mathbf{D}}} \mathbf{P}_{\mathbf{Z}_{\mathcal{V}^{c*}}} \left(\mathbf{D} \cdot \text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) - \boldsymbol{\epsilon} \right) / n, \end{aligned} \quad (\text{B23})$$

where $\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) = \frac{\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{Z_{V^{c*}}}) \boldsymbol{\epsilon}}{\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{Z_{V^{c*}}}) \mathbf{D}}$.

To explore the second term, we denote $\bar{\mathbf{D}} = \mathbf{P}_{Z_{V^{c*}}} \mathbf{D}$, $\bar{\boldsymbol{\epsilon}} = \mathbf{P}_{Z_{V^{c*}}} \boldsymbol{\epsilon}$. Due to the blockwise formula for the projection matrix, we have

$$\mathbf{P}_Z - \mathbf{P}_{Z_{V^{c*}}} = \mathbf{P}_{M_{Z_{V^{c*}}} \mathbf{Z}_{V^*}},$$

which is also a projection matrix. Thus, we denote $\tilde{\mathbf{D}} = \mathbf{P}_{M_{Z_{V^{c*}}} \mathbf{Z}_{V^*}} \mathbf{D}$, $\tilde{\boldsymbol{\epsilon}} = \mathbf{P}_{M_{Z_{V^{c*}}} \mathbf{Z}_{V^*}} \boldsymbol{\epsilon}$ and $\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) = \frac{\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{Z_{V^{c*}}}) \boldsymbol{\epsilon}}{\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{Z_{V^{c*}}}) \mathbf{D}} = \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}}$. Thus we have

$$\hat{\mathbf{D}} - \bar{\mathbf{D}} = (\mathbf{P}_Z - \mathbf{P}_{Z_{V^{c*}}}) \mathbf{D} = \tilde{\mathbf{D}}. \quad (\text{B24})$$

Therefore, the second term inside the RHS of (B23) can be reformulated as:

$$\begin{aligned} & \mathbf{Z}_{V^*}^\top M_{\hat{\mathbf{D}}} \mathbf{P}_{Z_{V^{c*}}} (\mathbf{D} \cdot \text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) - \boldsymbol{\epsilon}) / n \\ &= \mathbf{Z}_{V^*}^\top \left(\mathbf{I} - \frac{\hat{\mathbf{D}} \hat{\mathbf{D}}^\top}{\hat{\mathbf{D}}^\top \hat{\mathbf{D}}} \right) (\bar{\mathbf{D}} \text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) - \bar{\boldsymbol{\epsilon}}) / n \\ &= \frac{(I) + (II) + (III) + (IV)}{\hat{\mathbf{D}}^\top \hat{\mathbf{D}} / n}, \end{aligned} \quad (\text{B25})$$

where

$$\begin{aligned} (I) &= \text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) \cdot \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}^\top}{n} \cdot \frac{\hat{\mathbf{D}}^\top \bar{\mathbf{D}}}{n} = \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \cdot \frac{\hat{\mathbf{D}}^\top \hat{\mathbf{D}}}{n} \cdot \frac{\mathbf{Z}_{V^*}^\top \bar{\mathbf{D}}}{n}, \\ (II) &= - \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}^\top}{n} \cdot \frac{\hat{\mathbf{D}} \bar{\boldsymbol{\epsilon}}}{n} = - \frac{\hat{\mathbf{D}}^\top \hat{\mathbf{D}}}{n} \cdot \frac{\mathbf{Z}_{V^*}^\top \bar{\boldsymbol{\epsilon}}}{n}, \\ (III) &= - \text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) \cdot \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}}{n} \cdot \frac{\hat{\mathbf{D}}^\top \bar{\mathbf{D}}}{n} = - \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \cdot \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}}{n} \cdot \frac{\hat{\mathbf{D}}^\top \bar{\mathbf{D}}}{n}, \\ (IV) &= \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}}{n} \cdot \frac{\hat{\mathbf{D}}^\top \bar{\boldsymbol{\epsilon}}}{n}. \end{aligned}$$

Now, using identity (B24), we replace $\frac{\hat{\mathbf{D}}^\top \bar{\mathbf{D}}}{n}$ in (III) with $\frac{\hat{\mathbf{D}}^\top (\hat{\mathbf{D}} - \tilde{\mathbf{D}})}{n}$. Therefore we obtain

$$\begin{aligned} (I) + (III) &= \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \cdot \frac{\hat{\mathbf{D}}^\top \hat{\mathbf{D}}}{n} \cdot \frac{\mathbf{Z}_{V^*}^\top (\bar{\mathbf{D}} - \hat{\mathbf{D}})}{n} + \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \cdot \frac{\hat{\mathbf{D}}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}}{n} \\ &= \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \cdot \frac{\hat{\mathbf{D}}^\top \hat{\mathbf{D}}}{n} \cdot \frac{-\mathbf{Z}_{V^*}^\top \tilde{\mathbf{D}}}{n} + \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \cdot \frac{\hat{\mathbf{D}}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\mathbf{Z}_{V^*}^\top \hat{\mathbf{D}}}{n}. \end{aligned} \quad (\text{B26})$$

Hence, (B25) becomes

$$\begin{aligned}
& \frac{(I) + (II) + (III) + (IV)}{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}/n} \\
&= \frac{-\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n} \left(\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} + \frac{\mathbf{Z}_{\mathcal{V}^*}^\top \bar{\boldsymbol{\epsilon}}}{n} \right) + \frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \left(\frac{\widehat{\mathbf{D}}^\top \bar{\boldsymbol{\epsilon}}}{n} + \frac{\widehat{\mathbf{D}}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \tilde{\boldsymbol{\epsilon}}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} \right)}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= -\frac{\mathbf{Z}_{\mathcal{V}^*}^\top}{n} (P_{\tilde{\mathbf{D}}} \tilde{\boldsymbol{\epsilon}} + \bar{\boldsymbol{\epsilon}}) + \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top}{n} (P_{\tilde{\mathbf{D}}} \tilde{\boldsymbol{\epsilon}} + \bar{\boldsymbol{\epsilon}})}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}},
\end{aligned} \tag{B27}$$

where $P_{\tilde{\mathbf{D}}} = \tilde{\mathbf{D}} \tilde{\mathbf{D}}^\top / \tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}$ is a projection matrix of $\tilde{\mathbf{D}}$. Notice

$$\begin{aligned}
P_{\tilde{\mathbf{D}}} \tilde{\boldsymbol{\epsilon}} + \bar{\boldsymbol{\epsilon}} &= \{P_{\tilde{\mathbf{D}}} (P_{\mathbf{Z}} - P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) + P_{\mathbf{Z}_{\mathcal{V}^{c*}}}\} \boldsymbol{\epsilon} = (P_{\tilde{\mathbf{D}}} + P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \boldsymbol{\epsilon}, \\
\widehat{\mathbf{D}}^\top \tilde{\mathbf{D}} &= \mathbf{D}^\top P_{\mathbf{Z}} (P_{\mathbf{Z}} - P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D} = \mathbf{D}^\top (P_{\mathbf{Z}} - P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{D} = \tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}.
\end{aligned} \tag{B28}$$

Thus, we are able to further simplify (B27) as

$$\begin{aligned}
& -\frac{\mathbf{Z}_{\mathcal{V}^*}^\top}{n} (P_{\tilde{\mathbf{D}}} \tilde{\boldsymbol{\epsilon}} + \bar{\boldsymbol{\epsilon}}) + \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top}{n} (P_{\tilde{\mathbf{D}}} \tilde{\boldsymbol{\epsilon}} + \bar{\boldsymbol{\epsilon}})}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= \frac{-\mathbf{Z}_{\mathcal{V}^*}^\top (P_{\tilde{\mathbf{D}}} + P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \boldsymbol{\epsilon}}{n} + \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top}{n} (P_{\tilde{\mathbf{D}}} + P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \boldsymbol{\epsilon}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= \frac{-\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n} + \frac{-\mathbf{Z}_{\mathcal{V}^*}^\top P_{\mathbf{Z}_{\mathcal{V}^{c*}}} \boldsymbol{\epsilon}}{n} + \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot \left\{ \frac{\widehat{\mathbf{D}}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n} + \frac{\widehat{\mathbf{D}}^\top P_{\mathbf{Z}_{\mathcal{V}^{c*}}} \boldsymbol{\epsilon}}{n} \right\}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= \frac{-\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n} + \frac{-\mathbf{Z}_{\mathcal{V}^*}^\top P_{\mathbf{Z}_{\mathcal{V}^{c*}}} \boldsymbol{\epsilon}}{n} + \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot (\tilde{\mathbf{D}}^\top + \widehat{\mathbf{D}}^\top P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \boldsymbol{\epsilon}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}} \\
&= \frac{-\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n} + \frac{-\mathbf{Z}_{\mathcal{V}^*}^\top P_{\mathbf{Z}_{\mathcal{V}^{c*}}} \boldsymbol{\epsilon}}{n} + \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \widehat{\mathbf{D}}}{n} \cdot \frac{\widehat{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n}}{\frac{\widehat{\mathbf{D}}^\top \widehat{\mathbf{D}}}{n}}
\end{aligned} \tag{B29}$$

Combining (B29), (B18) and (B23), we obtain

$$R_{\mathcal{V}^*}^{\text{or}} = \frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\boldsymbol{\epsilon}}}{n} - \frac{\frac{\mathbf{Z}_{\mathcal{V}^*}^\top \tilde{\mathbf{D}}}{n} \cdot \frac{\tilde{\mathbf{D}}^\top \boldsymbol{\epsilon}}{n}}{\frac{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}}{n}}, \tag{B30}$$

which has a similar structure to (B18). Consequently, with an analogous argument, we derive

$$\|\mathbf{R}^{\text{or}}\|_\infty = \|\mathbf{R}_{\mathcal{V}^*}^{\text{or}}\|_\infty = O_p \left(\frac{p_{\mathcal{V}^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*\|_\infty}{\gamma_{\mathcal{V}^*}^{*\top} \tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*} + \sqrt{\frac{\log p_{\mathcal{V}^*}}{n}} \right), \tag{B31}$$

where $\tilde{\mathbf{Q}}_n = \mathbf{Z}_{\mathcal{V}^*}^\top (P_{\mathbf{Z}} - P_{\mathbf{Z}_{\mathcal{V}^{c*}}}) \mathbf{Z}_{\mathcal{V}^*} / n$. \square

Appendix B9. Proof of Theorem 3

The proof of Theorem 3 is similar to the technique of Feng and Zhang (2019). However, we pay more careful attention to the issue of endogeneity of $\tilde{\mathbf{Z}}$ and ξ , which is discussed in Lemma 2.

To proceed, we first prove a useful inequality, known as basic inequality in sparse regression literature. Denote α^0 as α^* or $\hat{\alpha}^{\text{or}}$, sharing the same support on \mathcal{V}^* , and \mathbf{R}^0 is \mathbf{R}^* or \mathbf{R}^{or} defined in Lemma 2 upon the choice of α^0 .

LEMMA A2. Suppose $\hat{\alpha}$ is a solution of (23) and denote $\Delta = \hat{\alpha} - \alpha^0$ and let

$$\omega(\alpha) = \left[\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\alpha} - \tilde{\mathbf{Z}}^{\top}(\mathbf{Y} - \tilde{\mathbf{Z}}\alpha)/n \right] / \lambda \quad (\text{B32})$$

to measure the scaled violation of first order condition in (22), then

$$\Delta^{\top} \mathbf{C}_n \Delta \leq -\lambda \Delta^{\top} \omega(\alpha^0) + 1/\rho \|\Delta\|_2^2. \quad (\text{B33})$$

Further, choose a proper sub-derivative at the origin: $\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\alpha_j^0} = \lambda \text{sgn}(\Delta_j), \forall j \in \mathcal{V}^*$,

$$\Delta^{\top} \mathbf{C}_n \Delta + (\lambda - \|\mathbf{R}_{\mathcal{V}^*}\|_{\infty}) \|\Delta_{\mathcal{V}^*}\|_1 \leq -\lambda \Delta_{\mathcal{V}^*}^{\top} \omega_{\mathcal{V}^*}(\alpha^0) + 1/\rho \|\Delta\|_2^2. \quad (\text{B34})$$

PROOF. Because $\omega(\hat{\alpha}) = \mathbf{0}$ in (22), we have $\tilde{\mathbf{Z}}^{\top} \mathbf{Y}/n = \frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\hat{\alpha}} + \tilde{\mathbf{Z}}^{\top} \tilde{\mathbf{Z}} \hat{\alpha}/n$. Recall $\mathbf{C}_n = \tilde{\mathbf{Z}}^{\top} \tilde{\mathbf{Z}}/n$. Thus, we replace $\tilde{\mathbf{Z}}^{\top} \mathbf{Y}/n$ in $\omega(\alpha^*)$ and obtain

$$\mathbf{C}_n \Delta = -\lambda \omega(\alpha^0) + \frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\alpha^0} - \frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\hat{\alpha}}.$$

Multiply Δ^{\top} on both sides, we have

$$\begin{aligned} \Delta^{\top} \mathbf{C}_n \Delta &= -\lambda \Delta^{\top} \omega(\alpha^0) + \Delta^{\top} \left(\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\alpha^0} - \frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\hat{\alpha}} \right) \\ &\leq -\lambda \Delta^{\top} \omega(\alpha^0) + 1/\rho \|\Delta\|_2^2, \end{aligned} \quad (\text{B35})$$

where the second line follows the convexity level of MCP up to $1/\rho$ and concludes (B33). Moreover, we further examine the terms in $\omega(\alpha^0)$ with respect to \mathcal{V}^* . We obtain,

$$\begin{aligned} \omega_{\mathcal{V}^*}(\alpha^0) &= \left[\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\alpha_{\mathcal{V}^*}^0} - \tilde{\mathbf{Z}}_{\mathcal{V}^*}^{\top}(\mathbf{Y} - \tilde{\mathbf{Z}}\alpha^0)/n \right] / \lambda \\ &= \left[\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\mathbf{0}_{\mathcal{V}^*}} - \mathbf{R}_{\mathcal{V}^*}^0 \right] / \lambda. \end{aligned} \quad (\text{B36})$$

Thus, we rewrite (B33) as

$$\begin{aligned} -\lambda \Delta_{\mathcal{V}^*}^{\top} \omega_{\mathcal{V}^*}(\alpha^0) + 1/\rho \|\Delta\|_2^2 &\geq \Delta^{\top} \mathbf{C}_n \Delta + \lambda \Delta_{\mathcal{V}^*}^{\top} \omega_{\mathcal{V}^*}(\alpha^0) \\ &= \Delta^{\top} \mathbf{C}_n \Delta + \lambda \Delta_{\mathcal{V}^*}^{\top} \left[\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\mathbf{0}_{\mathcal{V}^*}} - \mathbf{R}_{\mathcal{V}^*}^0 \right] / \lambda \\ &\geq \Delta^{\top} \mathbf{C}_n \Delta + \Delta_{\mathcal{V}^*}^{\top} \left[\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\mathbf{0}_{\mathcal{V}^*}} \right] - \|\Delta_{\mathcal{V}^*}\|_1 \|\mathbf{R}_{\mathcal{V}^*}^0\|_{\infty} \\ &= \Delta^{\top} \mathbf{C}_n \Delta + (\lambda - \|\mathbf{R}_{\mathcal{V}^*}^0\|_{\infty}) \|\Delta_{\mathcal{V}^*}\|_1, \end{aligned} \quad (\text{B37})$$

where the last equality holds for a proper sub-derivative at the origin, i.e., for $j \in \mathcal{V}^*$, $\frac{\partial}{\partial \mathbf{t}} p_{\lambda}^{\text{MCP}}(\mathbf{t})|_{\mathbf{t}=\alpha_j^0} = \lambda \text{sgn}(\Delta_j)$. \square

Under the event $\Omega = \{\omega_{\mathcal{V}^*}(\hat{\alpha}^{or}) = \mathbf{0}\}$, we now prove the estimation error Δ belongs to the cone $\mathcal{C}(\mathcal{V}^*; \xi)$, where ξ is chosen in Lemma 2. Recall $\mathcal{B}(\lambda, \rho) = \{\hat{\alpha} \text{ in (23)} : \lambda \geq \zeta, \rho > K_{\mathcal{C}}^{-2}(\mathcal{V}^*, \xi) \vee 1\}$ as a collection of $\hat{\alpha}$ computed in (23) through a broad class of MCP, in which ζ is define in (25).

LEMMA A3. *Under the event $\Omega = \{\omega_{\mathcal{V}^*}(\hat{\alpha}^{or}) = \mathbf{0}\}$, consider $\hat{\alpha}_{\lambda_1}, \hat{\alpha}_{\lambda_2} \in \mathcal{B}(\lambda, \rho)$ with different penalty levels λ_1 and λ_2 , and denote their estimation errors as $\Delta_1 = \hat{\alpha}_{\lambda_1} - \hat{\alpha}^{or}$ and $\Delta_2 = \hat{\alpha}_{\lambda_2} - \hat{\alpha}^{or}$, respectively. Define $a_1 = 1 - \|\mathbf{R}^{or}_{\mathcal{V}^*}\|_{\infty}/\lambda$, $a_2 = a_1\xi\rho/[2(\xi+1)]$, $a_3 = a_1\xi/(\xi+1+a_1)$, and $a_0 = a_2 \wedge \{a_2a_3/(1 \vee \xi)\}$. Then, once $\Delta_1 \in \mathcal{C}(\mathcal{V}^*; \xi)$ and $\|\Delta_1 - \Delta_2\|_1 \leq a_0\lambda$ hold, we conclude $\Delta_2 \in \mathcal{C}(\mathcal{V}^*; \xi)$.*

PROOF. Let $a_1 = 1 - \|\mathbf{R}^{or}_{\mathcal{V}^*}\|_{\infty}/\lambda > 0$. Applying Lemma A2 to Δ_2 and under the event Ω , we have

$$\begin{aligned} \Delta_2^{\top} C_n \Delta_2 + (\lambda - \|\mathbf{R}^{or}_{\mathcal{V}^*}\|_{\infty}) \|\Delta_{2\mathcal{V}^*}\|_1 &\leq 1/\rho \|\Delta_2\|_2^2 \\ \iff \Delta_2^{\top} C_n \Delta_2 + a_1 \lambda \|\Delta_{2\mathcal{V}^*}\|_1 &\leq 1/\rho \|\Delta_2\|_2^2. \end{aligned} \quad (\text{B38})$$

Consider the first case $\|\Delta_1\|_1 \vee \|\Delta_1 - \Delta_2\|_1 \leq a_2\lambda$, where $a_2 = a_1\xi\rho/[2(\xi+1)]$. We have

$$\|\Delta_2\|_2^2 \leq \|\Delta_2\|_{\infty} \cdot \|\Delta_2\|_1 \leq (\|\Delta_2 - \Delta_1\|_1 + \|\Delta_1\|_1) \cdot \|\Delta_2\|_1 \leq 2a_2\lambda \|\Delta_2\|_1.$$

The above inequalities yield

$$a_1 \lambda \|\Delta_{2\mathcal{V}^*}\|_1 \leq 1/\rho \|\Delta_2\|_2^2 \leq \{a_1\xi/(\xi+1)\} (\|\Delta_{2\mathcal{V}^*}\|_1 + \|\Delta_{2\mathcal{V}^*}\|_1), \quad (\text{B39})$$

which is equivalent to $\Delta_2 \in \mathcal{C}(\mathcal{V}^*, \xi)$ by algebra in the first case.

Consider the second case that $\|\Delta_1\|_1 \geq a_2\lambda$ and $\|\Delta_1 - \Delta_2\|_1 \leq \lambda a_2 a_3 / (1 \vee \xi)$, where $a_3 = a_1\xi/(\xi+1+a_1)$. Similarly, applying Lemma A2 to Δ_1 and using $\Delta_1 \in \mathcal{C}(\mathcal{V}^*; \xi)$, we obtain

$$a_1 \lambda \|\Delta_{1\mathcal{V}^*}\|_1 \leq 0 \quad \Rightarrow \quad \Delta_{1\mathcal{V}^*} = \mathbf{0}.$$

The triangle inequalities

$$\|\Delta_{1\mathcal{V}^*} - \Delta_{2\mathcal{V}^*}\|_1 + \|\Delta_{2\mathcal{V}^*}\|_1 \geq \|\Delta_{1\mathcal{V}^*}\|_1, \quad \|\Delta_{2\mathcal{V}^*} - \Delta_{1\mathcal{V}^*}\|_1 + \|\Delta_{1\mathcal{V}^*}\|_1 \geq \|\Delta_{2\mathcal{V}^*}\|_1$$

give rise to

$$\begin{aligned} &\|\Delta_{2\mathcal{V}^*}\|_1 - \xi \|\Delta_{2\mathcal{V}^*}\|_1 \\ &\leq \|\Delta_{1\mathcal{V}^*}\|_1 - \xi \|\Delta_{1\mathcal{V}^*}\|_1 + (1 \vee \xi) \|\Delta_2 - \Delta_1\|_1 \\ &\leq \|\Delta_{1\mathcal{V}^*}\|_1 - \xi \|\Delta_{1\mathcal{V}^*}\|_1 + a_3 \|\Delta_1\|_1 \\ &= (a_3 - \xi) \|\Delta_{1\mathcal{V}^*}\|_1 = \frac{-\xi(\xi+1)}{\xi+1+a_1} \|\Delta_{1\mathcal{V}^*}\|_1 \leq 0, \end{aligned} \quad (\text{B40})$$

where the second inequality follows the assumptions of $\|\Delta_1\|_1$ and $\|\Delta_1 - \Delta_2\|_1$.

Thus, together with above two cases, it concludes the $\Delta_2 \in \mathcal{C}(\mathcal{V}^*, \xi)$ when $\|\Delta_1 - \Delta_2\|_1 < a_0\lambda$, where $a_0 = a_2 \wedge \{a_2a_3/(1 \vee \xi)\}$ \square

Based on the above two lemmas, we are able to derive the theoretical results stated in Theorem 3.

PROOF (OF THEOREM 3). Consider the local solution $\hat{\alpha}$ in $\mathcal{B}_0(\lambda, \rho)$ and denote $\hat{\mathcal{V}} = \{j : \hat{\alpha}_j = 0\}$ and event $\Phi = \{\hat{\mathcal{V}} = \mathcal{V}^*\}$ of most interest. Thus,

$$\Pr(\Phi) = \Pr(\Phi, \Omega) + \Pr(\Phi, \Omega^c) \geq \Pr(\Phi|\Omega) \Pr(\Omega). \quad (\text{B41})$$

Firstly, conditional on the event Ω , we denote $\Delta = \hat{\alpha} - \hat{\alpha}^{\text{or}}$ and immediately have $\Delta \in \mathcal{C}(\mathcal{V}^*; \xi)$ by Lemma A3. Applying Lemma A2 to Δ , we have

$$\Delta^\top C_n \Delta + (\lambda - \|\mathbf{R}_{\mathcal{V}^*}^{\text{or}}\|_\infty) \|\Delta_{\mathcal{V}^*}\|_1 \leq -\lambda \Delta_{\mathcal{V}^{c*}}^\top \omega_{\mathcal{V}^{c*}}(\hat{\alpha}^{\text{or}}) + 1/\rho \|\Delta\|_2^2. \quad (\text{B42})$$

By Cauchy-Schwarz inequality,

$$-\lambda \Delta_{\mathcal{V}^{c*}}^\top \omega_{\mathcal{V}^{c*}}(\alpha^{\text{or}}) = \lambda \Delta_{\mathcal{V}^{c*}}^\top [-\omega_{\mathcal{V}^{c*}}(\alpha^{\text{or}})] \leq \lambda \|\Delta_{\mathcal{V}^{c*}}\|_2 \|\omega_{\mathcal{V}^{c*}}(\alpha^{\text{or}})\|_2 = 0$$

follows the definition of Ω . Rearranging (B42) yields,

$$\begin{aligned} 0 &\geq \Delta^\top Q_n \Delta - 1/\rho \|\Delta\|_2^2 + (\lambda - \|\mathbf{R}_{\mathcal{V}^*}^{\text{or}}\|_\infty) \|\Delta_{\mathcal{V}^*}\|_1 \\ &\geq (K_{\mathcal{C}}^2(\mathcal{V}^*, \xi) - 1/\rho) \|\Delta\|_2^2 + (\lambda - \|\mathbf{R}_{\mathcal{V}^*}^{\text{or}}\|_\infty) \|\Delta_{\mathcal{V}^*}\|_1 \geq 0, \end{aligned} \quad (\text{B43})$$

where the second line follows membership of cone $\mathcal{C}(\mathcal{V}^*; \xi)$ of Δ and the RE condition of \tilde{Z} in Lemma 2. Inequality (B43) forces $\|\Delta_{\mathcal{V}^*}\|_1 = 0$ and $\|\Delta\|_2^2 = 0$, i.e., $\mathcal{V}^* = \hat{\mathcal{V}}$, in probability because $\|\mathbf{R}_{\mathcal{V}^*}^{\text{or}}\|_\infty < \lambda$ holds with probability approaching 1 in Lemma 3.

Therefore, it remains to investigate event Ω .

$$\omega_{\mathcal{V}^{c*}}(\alpha^{\text{or}}) = \left[\frac{\partial}{\partial t} p_\lambda^{\text{MCP}}(t) \Big|_{t=\alpha_{\mathcal{V}^{c*}}^{\text{or}}} - \tilde{Z}_{\mathcal{V}^{c*}}^\top (Y - \tilde{Z} \hat{\alpha}^{\text{or}}) \right] = \frac{\partial}{\partial t} p_\lambda^{\text{MCP}}(t) \Big|_{t=\hat{\alpha}_{\mathcal{V}^{c*}}^{\text{or}}} \quad (\text{B44})$$

follows definition of $\hat{\alpha}^{\text{or}}$. By the definition of Ω and characteristics of MCP,

$$\Omega = \{\omega_{\mathcal{V}^{c*}}(\hat{\alpha}^{\text{or}}) = \frac{\partial}{\partial t} p_\lambda^{\text{MCP}}(t) \Big|_{t=\hat{\alpha}_{\mathcal{V}^{c*}}^{\text{or}}} = 0\} = \{|\hat{\alpha}_{\mathcal{V}^{c*}}^{\text{or}}|_{\min} > \lambda/\rho\}.$$

Rearranging (B22), which yields

$$\left\{ |\alpha_{\mathcal{V}^{c*}}^*|_{\min} > \lambda/\rho + \|(Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top \epsilon\|_\infty + \left\| (Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top D \frac{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) \epsilon}{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) D} \right\|_\infty \right\} \subseteq \Omega. \quad (\text{B45})$$

Thus, it suffices to examine

$$\left\| (Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top D \frac{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) \epsilon}{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) D} \right\|_\infty = \left| \frac{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) \epsilon}{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) D} \right| \cdot \|(Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top D\|_\infty.$$

The first term in the RHS measures the estimation error of TSLS estimator, i.e.,

$$\text{Bias}(\hat{\beta}_{or}^{TSLS}) = \frac{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) \epsilon}{D^\top (P_Z - P_{Z_{\mathcal{V}^{c*}}}) D} = \frac{\tilde{D}^\top \tilde{\epsilon}}{\tilde{D}^\top \tilde{D}},$$

is the unvanishing term under many (weak) IVs setting. While for the second term,

$$\begin{aligned} &\|(Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top D\|_\infty \\ &= \|\gamma_{\mathcal{V}^{c*}}^* + (Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^*} \gamma_{\mathcal{V}^*}^*\|_\infty + \|(Z_{\mathcal{V}^{c*}}^\top Z_{\mathcal{V}^{c*}})^{-1} Z_{\mathcal{V}^{c*}}^\top \eta\|_\infty \\ &= \|\tilde{\gamma}_{\mathcal{V}^{c*}}^*\|_\infty + O_p\left(\sigma_\eta \sqrt{\frac{2 \log(2p_{\mathcal{V}^{c*}})}{n}}\right), \end{aligned}$$

where $\bar{\gamma}_{\mathcal{V}^{c*}}^* = \gamma_{\mathcal{V}^{c*}}^* + (\mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^{c*}})^{-1} \mathbf{Z}_{\mathcal{V}^{c*}}^\top \mathbf{Z}_{\mathcal{V}^*} \gamma_{\mathcal{V}^*}^*$.

Thereby, (B45) reduces to

$$\left\{ |\alpha_{\mathcal{V}^{c*}}^*|_{\min} > \lambda/\rho + O_p\left(\sigma_\epsilon \sqrt{\frac{2\log(2p_{\mathcal{V}^{c*}})}{n}}\right) + |\text{Bias}(\hat{\beta}_{or}^{TSLS})| \cdot \left[\|\bar{\gamma}_{\mathcal{V}^{c*}}^*\|_\infty + O_p\left(\sigma_\eta \sqrt{\frac{2\log(2p_{\mathcal{V}^{c*}})}{n}}\right) \right] \right\} \subseteq \Omega.$$

Combining with $\lambda > \zeta$, we now specify (B45) as

$$\left\{ |\alpha_{\mathcal{V}^{c*}}^*|_{\min} > C\kappa(n) \right\} \subseteq \Omega,$$

where $\kappa(n) = \sqrt{\frac{\log p_{\mathcal{V}^*}}{n}} + \frac{p_{\mathcal{V}^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*\|_\infty}{\gamma_{\mathcal{V}^*}^{*\top} \tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*} + |\text{Bias}(\hat{\beta}_{or}^{TSLS})| \cdot \|\bar{\gamma}_{\mathcal{V}^{c*}}^*\|_\infty$. Thus, under the condition that event Ω holds in the finite sample or in probability, we achieve consistency of the selection of valid IVs.

We then turn to \mathcal{P}_c . Recall that $\tilde{\epsilon}^c = \epsilon - c\eta$, so $\sqrt{\text{Var}(\tilde{\epsilon}^c)} = \sqrt{\sigma_\epsilon^2 + c^2\sigma_\eta^2 - 2c\sigma_{\epsilon,\eta}} \asymp O(1+c)$ and $\text{Cov}(\tilde{\epsilon}^c, \eta) = \sigma_\epsilon^2 - c\sigma_\eta^2 \asymp O(1+c)$. After some direct derivation, we obtain

$$\kappa^c(n) \asymp (1+c) \sqrt{\frac{\log |\mathcal{I}_c|}{n}} + (1+c) \frac{|\mathcal{I}_c|}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n^c \gamma_{\mathcal{I}_c}^*\|_\infty}{\gamma_{\mathcal{I}_c}^{*\top} \tilde{\mathbf{Q}}_n^c \gamma_{\mathcal{I}_c}^*} + |\text{Bias}(\hat{\beta}_{or}^{TSLS})| \cdot \|\bar{\gamma}_{\mathcal{I}_c}^*\|_\infty, \quad (\text{B46})$$

where $\tilde{\mathbf{Q}}_n^c$ and $\text{Bias}(\hat{\beta}_{or}^{TSLS})$ are defined as \mathcal{P}_c version of $\tilde{\mathbf{Q}}_n$ and $\text{Bias}(\hat{\beta}_{or}^{TSLS})$. With similar argument, we also achieve the consistency of selection of valid counterparts in \mathcal{P}_c if $|\alpha_j^c| > \kappa^c(n)$ for $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$ holds.

□

Appendix B10. Proof of Proposition 2

PROOF.

$$T_2 = \frac{p_{\mathcal{V}^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*\|_\infty}{\gamma_{\mathcal{V}^*}^{*\top} \tilde{\mathbf{Q}}_n \gamma_{\mathcal{V}^*}^*} \leq \frac{p_{\mathcal{V}^*}}{n} \cdot \frac{\|\tilde{\mathbf{Q}}_n^{1/2}\|_\infty \|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_\infty}{\|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_2^2} \leq \frac{p_{\mathcal{V}^*}}{n} \cdot \frac{p_{\mathcal{V}^*} \|\tilde{\mathbf{Q}}_n^{1/2}\|_\infty \|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_\infty}{\|\tilde{\mathbf{Q}}_n^{1/2} \gamma_{\mathcal{V}^*}^*\|_1^2} \rightarrow 0$$

□

Appendix B11. Proof of Proposition 3

This proof is extended from Bun and Windmeijer (2011)'s higher order approximation arguments.

PROOF. Recall

$$\text{Bias}(\hat{\beta}_{or}^{TSLS}) = \frac{\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{Z_{\mathcal{V}^{c*}}}) \epsilon}{\mathbf{D}^\top (\mathbf{P}_Z - \mathbf{P}_{Z_{\mathcal{V}^{c*}}}) \mathbf{D}} = \frac{\tilde{\mathbf{D}}^\top \tilde{\epsilon}}{\tilde{\mathbf{D}}^\top \tilde{\mathbf{D}}} := \frac{c}{d}, \quad (\text{B47})$$

we have $\bar{c} := E(c) = \sigma_{\epsilon, \eta}^2(p - p_{\mathcal{V}^*}) = \sigma_{\epsilon, \eta}^2 p_{\mathcal{V}^*}$ and $\bar{d} := E(d) = \sigma_{\eta}^2(\mu_n + L)$. That is free of the number of invalid IVs $p_{\mathcal{V}^*}$. Let $s = \max(\mu_n, p_{\mathcal{V}^*})$, we have

$$\begin{aligned} \text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) &= \frac{\bar{c}}{\bar{d}} + \frac{c - \bar{c}}{\bar{d}} - \frac{\bar{c}(d - \bar{d})}{\bar{d}^2} - \frac{(c - \bar{c})(d - \bar{d})}{\bar{d}^2} + \frac{\bar{c}(d - \bar{d})^2}{\bar{d}^3} + O_p\left(s^{-\frac{3}{2}}\right) \\ E[\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}})] &= \frac{\sigma_{\epsilon, \eta}}{\sigma_{\eta}^2} \left(\frac{p_{\mathcal{V}^*}}{(\mu_n + p_{\mathcal{V}^*})} - \frac{2\mu_n^2}{(\mu_n + p_{\mathcal{V}^*})^3} \right) + o(s^{-1}) \end{aligned} \quad (\text{B48})$$

follows from Section 3 in [Bun and Windmeijer \(2011\)](#). \square

Appendix B12. Proof of Theorem 4

PROOF. Under the conditions of Theorem 3, we have $\Pr(\hat{\mathcal{V}} = \mathcal{V}^*) \xrightarrow{p} 1$. Thus, $\hat{\beta}_{\text{WIT}} \xrightarrow{p} \hat{\beta}_{\text{or}}^{\text{liml}}$, where $\hat{\beta}_{\text{or}}^{\text{liml}}$ stands for LIML estimator with known \mathcal{V}^* a priori. Thus, (a) follows Corollary 1(iv) in [\(Kolesár et al., 2015\)](#) with $\min \text{eig}(\Sigma^{-1}\Lambda) = 0$ in their context. (b) and (c) follow [Kolesár \(2018\)](#), Proposition 1. \square

Appendix B13. Proof of Corollary 2

PROOF. Notice $p_{\mathcal{V}^*}/n < p/n \rightarrow 0$ and $\mu_n/n \xrightarrow{p} \mu_0$, therefore the threshold $T_2 \rightarrow 0$ in Theorem 3. Likewise $T_3 \rightarrow 0$ follows $\text{Bias}(\hat{\beta}_{\text{or}}^{\text{TSLs}}) \xrightarrow{p} \frac{\sigma_{\epsilon, \eta}}{\sigma_{\eta}^2} \left(\frac{p_{\mathcal{V}^*}}{(\mu_n + p_{\mathcal{V}^*})} - \frac{2\mu_n^2}{(\mu_n + p_{\mathcal{V}^*})^3} \right) = o(1)$. Thus, $\kappa(n) \asymp n^{-1/2}$ in (28) diminishes to 0. Thus any fixed $\min_{j \in \mathcal{V}^*} \alpha_j^*$ would pass $\kappa(n)$ asymptotically. Let $c = \alpha_l^*/\gamma_l^* = C_1/n^{\tau_1}$ for $l \in \mathcal{I}_c$. Then, for $\kappa^c(n)$, we have: $|E \text{Bias}(\hat{\beta}_{\text{or}}^c)| \approx \frac{\text{Cov}(\tilde{\epsilon}^c, \boldsymbol{\eta})}{\text{Var} \boldsymbol{\eta}} \cdot C_1^{-2} n^{2\tau_1-1} \asymp n^{3\tau_1-1}$ according to Proposition 3 and $\|\tilde{\gamma}_{\mathcal{I}_c}^*\|_{\infty} \leq C$ due to Assumption 4. Then, for the first two terms in $\kappa^c(n)$:

$$(1+c)\sqrt{\frac{\log |\mathcal{I}_c|}{n}} + (1+c)\frac{|\mathcal{I}_c|}{n} \cdot \frac{\|\tilde{\tilde{\mathbf{Q}}}_n^c \boldsymbol{\gamma}_{\mathcal{I}_c}^*\|_{\infty}}{\boldsymbol{\gamma}_{\mathcal{I}_c}^{*\top} \tilde{\tilde{\mathbf{Q}}}_n^c \boldsymbol{\gamma}_{\mathcal{I}_c}^*} \asymp n^{\tau_1-1/2} + n^{\tau_1-1}.$$

Hence, we conclude $\kappa^c(n) \asymp n^{\max(\tau_1-1/2, 3\tau_1-1)}$. For $|\tilde{\alpha}_j^c| = |\alpha_j^* - c\gamma_j^*|$ and $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$. We consider all possible cases:

- (a) $\tilde{c} = 0$, i.e. $j \in \{\mathcal{I}_c^c : \alpha_j^* = 0\}$: $|\tilde{\alpha}_j^c| = |c| \cdot |\gamma_j^*|$, and $\gamma_j^* = n^{-\tau_2}$. Hence $|\tilde{\alpha}_j^c| \asymp n^{\tau_1-\tau_2}$.
- (b) $\tilde{c} \neq 0$, i.e. $j \in \{\mathcal{I}_c^c : \alpha_j^*/\gamma_j^* = \tilde{c}\}$: $|\tilde{\alpha}_j^c| = |c - \tilde{c}| \cdot |\gamma_j^*|$ and $\gamma_j^* = n^{-\tau_3}$. Hence, $|\tilde{\alpha}_j^c| \asymp |C_1 - C_3 n^{\tau_1-\tau_3}|$

Thus, $|\tilde{\alpha}_j^c| > \kappa^c(n)$ is equivalent to $2\tau_1 + \tau_2 < 1$ and $2\tau_1 + \tau_3 < 1$. By the symmetry of τ_1 and τ_3 of invalid IVs, we have $2\tau_3 + \tau_1 < 1$. Hence $\tau_1 + \tau_3 < 2/3$. Therefore, Assumption 5 holds automatically. Then it follows Theorem 4. \square