

# On the instrumental variable estimation with potentially many (weak) and some invalid instruments

Yiqi Lin<sup>a,b</sup>, Frank Windmeijer<sup>b</sup>, Xinyuan Song<sup>a</sup>, Qingliang Fan<sup>c</sup>

Presented by Yiqi LIN

Department of Statistics, The Chinese University of Hong Kong<sup>a</sup>

Department of Statistics, University of Oxford<sup>b</sup>

Department of Economics, The Chinese University of Hong Kong<sup>c</sup>

June 10, 2022

# Requirement of IVs

Good IV should satisfy the following conditions, illustrated as follows.

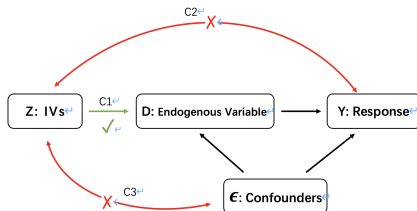


Figure: Illustration of Validity and Relevance.

## IVs Requirements

- 1 **Relevance Condition C1** : related to exposure (may strong or weak).
- 2 **Exogenous Condition C2**: not related to unmeasured variables that affect the exposure and the outcome.
- 3 **Exclusion Restriction C3**: have no direct pathway to the outcome.

# Model

Potential outcome model as in Kang et al. (2016); Small (2007). For  $i = 1, 2, \dots, n$ , we have the random sample  $(Y_i, D_i, \mathbf{Z}_i)$ . Let  $Y_i^{(D_i, \mathbf{Z}_i)}$  to be the potential outcome for the object  $i$  having exposure  $D_i$  and instrumental variables  $\mathbf{Z}_i \in \mathbb{R}^p$ .

## Potential Outcome Model

Given two different sets of treatment variables  $D_i^A, D_i^B$  and corresponding instruments  $\mathbf{Z}_i^A, \mathbf{Z}_i^B$ , assume

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

- 1  $D \in \mathbb{R}, \beta^* \in \mathbb{R}$  represents the constant causal parameter of interest.
- 2  $\mathbf{Z} \in \mathbb{R}^p, \phi \in \mathbb{R}^L$  represents the violation of **Exclusion Restriction**.
- 3  $\theta \in \mathbb{R}^p$  represents the violation of **Exogenous Condition**.

# Model

Note that  $\mathbf{Z}_i$  could have non-linear transformations of the original variables so a high-dimensional model is plausible. A good instrument  $\mathbf{Z}_j$  should not have a direct effect on the response and unmeasured confounders, i.e.,  $\phi_j = 0$  and  $\theta_j = 0$ .

## Model

Assuming the linear functional form between treatment effects  $D_i$  and instruments  $\mathbf{Z}_i$ , the above potential outcome model (1) can be rewritten as follows,

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

where  $\epsilon_i = Y_i^{(0,0)} - E(Y_i^{(0,0)} | \mathbf{Z}_i)$ ,  $\boldsymbol{\alpha} = \boldsymbol{\phi} + \boldsymbol{\theta}$ .

## Definition

- Relevant IV (satisfies C1): if  $\gamma_j^* \neq 0, j = 1, 2, \dots, p$ .
- Valid IV (satisfies C2 and C3): if  $\alpha_j^* = 0, j = 1, 2, \dots, p$ .

# Assumptions

Define the valid IV set  $\mathcal{V}^* = \{j : \alpha_j^* = 0\}$  and invalid IV set  $\mathcal{V}^{c*} = \{j : \alpha_j^* \neq 0\}$ . Let  $L = |\mathcal{V}^*|$ ,  $K = |\mathcal{V}^{c*}|$  and  $p = K + L$ . Notably,  $L \geq 1$  refers to the existence of excluded IV, namely the order condition (Wooldridge, 2010). We consider many (weak) IVs cases and make the following model assumptions:

## Assumptions

**Assumption 1** (Many valid and invalid IVs):  $p < n$ ,  $K/n \rightarrow v_K + o(n^{-1/2})$  and  $L/n \rightarrow v_L + o(n^{-1/2})$  for some  $v_K$  and  $v_L$  such that  $v_L + v_K < 1$ .

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

**Assumption 3:**  $\boldsymbol{\gamma}^* = (E(\mathbf{Z}_i \mathbf{Z}_i^\top))^{-1} E(\mathbf{Z}_i \mathbf{D}_i)$ ,  $\gamma_j^* \neq 0$  for  $j = 1, 2, \dots, p$  and given  $n$ .

**Assumption 4:** Let  $\mathbf{u}_i = (\epsilon_i, \eta_i)^\top$  and  $\mathbf{u}_i \mid \mathbf{Z}$  follows i.i.d. Normal distribution with mean zeros and positive definite covariance matrix  $\boldsymbol{\Sigma} = \begin{pmatrix} \sigma_\epsilon^2 & \sigma_{\epsilon, \eta} \\ \sigma_{\epsilon, \eta} & \sigma_\eta^2 \end{pmatrix}$ . The elements of  $\boldsymbol{\Sigma}$  are finite and  $\sigma_{\epsilon, \eta} \neq 0$  imposes the endogeneity of treatment  $\mathbf{D}_i$ .

**Assumption 5** (Strength of valid IVs): Concentration parameter  $\mu_n$  grows at the same rate as  $n$ , i.e.,  $\mu_n = \boldsymbol{\gamma}_{\mathbf{Z}_{\mathcal{V}^*}}^{*\top} \mathbf{Z}_{\mathcal{V}^*}^\top \mathbf{M}_{\mathbf{Z}_{\mathcal{V}^{c*}}} \mathbf{Z}_{\mathcal{V}^*} \boldsymbol{\gamma}_{\mathbf{Z}_{\mathcal{V}^*}}^* / \sigma_\eta^2 \rightarrow \mu_0 n$ , for some  $\mu_0 > 0$

# Identifiability of Model

- Moment condition available in Assumption 4:

$$\begin{aligned}E(\mathbf{Z}^T \varepsilon) &= E[\mathbf{Z}^T (\mathbf{Y} - \mathbf{Z}\alpha^* - \mathbf{D}\beta^*)] = \mathbf{0} \\E(\mathbf{Z}^T \mathbf{Y}) &= E(\mathbf{Z}^T \mathbf{Z}) \alpha^* + E(\mathbf{Z}^T \mathbf{D}) \beta^* \\ \Rightarrow \mathbf{\Gamma}_{p \times 1}^* &= \mathbf{\alpha}_{p \times 1}^* + \mathbf{\gamma}_{p \times 1}^* \beta_{1 \times 1}^*,\end{aligned} \tag{3}$$

where  $\mathbf{\Gamma}^* = E(\mathbf{Z}^T \mathbf{Z})^{-1} E(\mathbf{Z}^T \mathbf{Y})$  and  $\mathbf{\gamma}^* = E(\mathbf{Z}^T \mathbf{Z})^{-1} E(\mathbf{Z}^T \mathbf{D})$

- Both  $\mathbf{\Gamma}^*$  and  $\mathbf{\gamma}^*$  can be identified based on observed data.
- Define Expectation of individual IV Estimator:  $\beta_j^* \triangleq \frac{\Gamma_j^*}{\gamma_j^*} = \beta^* + \frac{\alpha_j^*}{\gamma_j^*}$ .

## Plurality Rule (Sufficient and Necessary Conditions?)

The parameters  $\beta^*$  and  $\alpha^*$  are identified if and only if the following hold:

$$|\mathcal{V}^* = \{j : \alpha_j^* / \gamma_j^* = 0\}| > \max_{c \neq 0} |\{j : \alpha_j^* / \gamma_j^* = c\}| \tag{4}$$

We re-examine the identifiability and show that the “only if” part is true only when all IVs are valid. In general, there is no iff condition in term of model

# View of Data Generating Process (DGP)

Given first stage information:  $\{\mathbf{D}, \mathbf{Z}, \boldsymbol{\gamma}^*\}$ , without loss of generality, we denote DGP with some  $\{\beta^*, \boldsymbol{\alpha}^*, \epsilon\}$  in (2) as DGP  $\mathcal{P}_0$  that generates  $\mathbf{Y}$ .

## Transformation of DGP

Given this  $\mathcal{P}_0$ , for  $j \in \mathcal{V}^{c*}$ , we have  $\mathbf{Z}_j \boldsymbol{\alpha}_j^* = \frac{\alpha_j^*}{\gamma_j^*} \left( \mathbf{D} - \sum_{l \neq j} \mathbf{Z}_l \boldsymbol{\gamma}_l^* - \boldsymbol{\eta} \right)$ . Denote  $\mathcal{I}_c = \{j \in \mathcal{V}^{c*} : \alpha_j^* / \gamma_j^* = c\}$ , where  $c \neq 0$  and  $c$  could have up to  $K$  different values. For compatibility, we denote  $\mathcal{I}_0 = \mathcal{V}^*$ . Thus, we are able to reformulate

$$\mathbf{Y} = \mathbf{D}\beta^* + \mathbf{Z}\boldsymbol{\alpha}^* + \epsilon \Rightarrow \mathbf{Y} = \mathbf{D}\tilde{\beta}^c + \mathbf{Z}\tilde{\boldsymbol{\alpha}}^c + \tilde{\epsilon}^c,$$

where  $\{\tilde{\beta}^c, \tilde{\boldsymbol{\alpha}}^c, \tilde{\epsilon}^c\} = \{\beta^* + c, \boldsymbol{\alpha}^* - c\boldsymbol{\gamma}^*, \epsilon - c\boldsymbol{\eta}\}$  and  $c = \alpha_j^* / \gamma_j^*$ , for some  $j \in \mathcal{V}^{c*}$ .

Evidently, for different  $c \neq 0$ , it forms different DGPs  $\mathcal{P}_c = \{\tilde{\beta}^c, \tilde{\boldsymbol{\alpha}}^c, \tilde{\epsilon}^c\}$  that can generate the same  $\mathbf{Y}$  (given  $\epsilon$ ) which also satisfies the moment condition (2) as  $\mathcal{P}_0$  since  $E(\mathbf{Z}^\top \tilde{\epsilon}^c) = \mathbf{0}$ .

## Example

$$\text{(Structural equation)} \mathbf{Y} = \mathbf{D}\beta^* + \mathbf{Z}_1\alpha_1^* + \mathbf{Z}_2\alpha_2^* + \epsilon$$

$$\text{(First Stage equation)} \mathbf{D} = \mathbf{Z}_1\gamma_1^* + \mathbf{Z}_2\gamma_2^* + \mathbf{Z}_3\gamma_3^* + \eta$$

Then we rearrange first stage equation:

$$\mathbf{Z}_1\gamma_1^* = \mathbf{D} - \mathbf{Z}_2\gamma_2^* - \mathbf{Z}_3\gamma_3^* - \eta$$

$$\Rightarrow \mathbf{Z}_1\alpha_1^* = \mathbf{Z}_1\gamma_1^*\left(\frac{\alpha_1^*}{\gamma_1^*}\right) = \mathbf{D}\frac{\alpha_1^*}{\gamma_1^*} - \mathbf{Z}_2\gamma_2^*\frac{\alpha_1^*}{\gamma_1^*} - \mathbf{Z}_3\gamma_3^*\frac{\alpha_1^*}{\gamma_1^*} - \eta\frac{\alpha_1^*}{\gamma_1^*}$$

$$\Rightarrow \mathbf{Y} = \mathbf{D}\beta^* + \left(\mathbf{D}\frac{\alpha_1^*}{\gamma_1^*} - \mathbf{Z}_2\gamma_2^*\frac{\alpha_1^*}{\gamma_1^*} - \mathbf{Z}_3\gamma_3^*\frac{\alpha_1^*}{\gamma_1^*} - \eta\frac{\alpha_1^*}{\gamma_1^*}\right) + \mathbf{Z}_2\alpha_2^* + \epsilon$$

$$\Rightarrow \mathbf{Y} = \mathbf{D}\left(\beta^* + \frac{\alpha_1^*}{\gamma_1^*}\right) + \mathbf{Z}_2\left(\alpha_2^* - \frac{\alpha_1^*}{\gamma_1^*}\right) + \mathbf{Z}_3\left(-\frac{\alpha_1^*}{\gamma_1^*}\right) + \left(\epsilon - \frac{\alpha_1^*}{\gamma_1^*}\eta\right)$$

Then It forms a new DGP:  $\tilde{\beta} = \beta^* + \frac{\alpha_1^*}{\gamma_1^*}$ ,  $\tilde{\alpha} = (0, \alpha_2^* - \frac{\alpha_1^*}{\gamma_1^*}, -\frac{\alpha_1^*}{\gamma_1^*})$  and  $\tilde{\epsilon} = \epsilon - \frac{\alpha_1^*}{\gamma_1^*}\eta$ .

For comparison:  $\alpha^* = (\alpha_1^*, \alpha_2^*, 0)$ .



# Identifiability of Model (Cont.)

## Theorem 1

Suppose Assumption 1-4 holds, given  $\mathcal{P}_0$  and  $\{\mathbf{D}, \mathbf{Z}, \gamma^*, \eta\}$ , 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  $\mathcal{V}^* \cup \{\cup_{c \neq 0} \mathcal{I}_c\} = \{1, 2, \dots, p\}$ ,  $\mathcal{V}^* \cap \{\cup_{c \neq 0} \mathcal{I}_c\} = \emptyset$  and  $E(\mathbf{Z}^\top \tilde{\epsilon}^c) = \mathbf{0}$ . The sparsity structure regarding  $\alpha$  is non-overlapping for different solutions.

Theorem 1 tells us there is a collection of model DGPs

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

corresponding to the same observation  $\mathbf{Y}$  conditional on first stage information. Let  $\mathcal{H}$  be a collection of mappings  $h : \mathcal{Q} \rightarrow \mathcal{P} \in \mathcal{Q}$

## Theorem 2

Under same conditions in Theorem 1, let  $\mathcal{F} = \{f : \mathcal{P} \in \mathcal{Q} \rightarrow \mathbb{R}; f(\mathcal{P}_i) < f(\mathcal{P}_j) \forall j \neq i \text{ and } \exists i \in \{0, \dots, G\}\}$  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 never exist a necessary condition of identifying  $(\alpha^*, \beta^*)$  unless  $\exists h \in \mathcal{H} : \mathcal{Q} \rightarrow \mathcal{P}_0$  and  $|\mathcal{H}| = 1$ .

# The Sparsest Rule

The sparsest rule is conceptually equivalent to plurality rule, since the non-overlapping sparse solutions given by Theorem 1.

- 1 For the sake of stable estimation, Guo et al. (2018) proposed to use plurality rule based on relevant IV set:

$$|\mathcal{V}_{S^*}^* = \{j \in \mathcal{S}^* : \alpha_j^* / \gamma_j^* = 0\}| > \max_{c \neq 0} |\{j \in \mathcal{S}^* : \alpha_j^* / \gamma_j^* = c\}|.$$

- 2  $\mathcal{S}^*$  is the relevant IVs estimated via first-step hard thresholding (Guo et al., 2018) and the individual margin to select relevant IVs is  $\sqrt{\widehat{\text{Var}}(\hat{\gamma}_j)} \cdot \sqrt{2.01 \log p \vee n} \asymp \sigma_\eta \sqrt{2.01 \log p \vee n/n}$ .

Thus, TSHT proposed by Guo et al. (2018) and CIIV proposed by Windmeijer et al. (2021) all explicitly leverage  $\mathcal{S}^*$ -based plurality rule to estimate  $\mathcal{V}_{S^*}^*$  and  $\beta^*$ .

# Drawback of plurality based-method

However, weak IVs sometimes is determinant in identification while limited in estimation.

## Example

Consider a toy example with mixed weak IV. Let  $\gamma^* = (\mathbf{0.04}_3, \mathbf{0.5}_2, 0.2, \mathbf{0.1}_4)^\top$  and  $\alpha^* = (\mathbf{0}_5, 1, \mathbf{0.7}_4)^\top$ . Obviously it forms 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 whole IVs set.

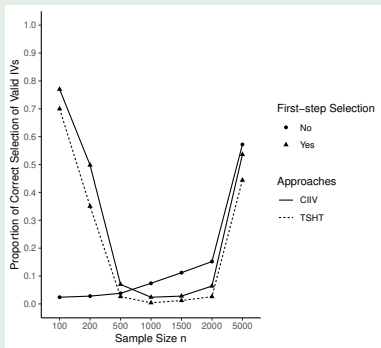


Figure: Proportion of correct selection of (subset) valid IVs.

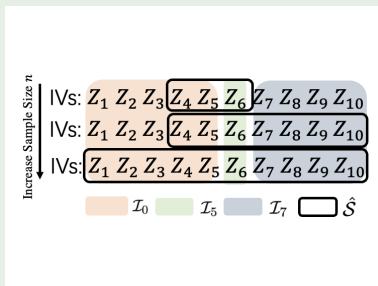


Figure: Plurality based on first stage selection.

# The Sparsest Rule (cont.)

- Mixture of weak and invalid IVs is ubiquitous in practice, especially in many IVs.
- For 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.

Recall  $\mathcal{P}_c = \{\tilde{\beta}^c, \tilde{\alpha}^c, \tilde{\epsilon}^c\} = \{\beta^* + c, \alpha^* - c\gamma^*, \epsilon - c\eta\}$ , where  $\tilde{\alpha}_{\mathcal{I}_c}^c = \mathbf{0}$ . For other elements in  $\tilde{\alpha}^c$  (corresponds to a different DGP in  $\mathcal{Q}$ ) and  $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$ , we obtain,

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

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

## Assumptions

Assumption 6.  $|\tilde{\alpha}_j^c| > \kappa(n)$  for  $j \in \{j : \alpha_j^*/\gamma_j^* = \tilde{c} \neq c\}$  and  $|\alpha_{\mathcal{V}^{c*}}^*|_{\min} > \kappa(n)$ , where  $\kappa(n)$  is a generally vanishing rate that specified by particular estimator to separate 0 and other non-zero terms.

Assumption 7. (The Sparsest Rule):  $\alpha^* = \operatorname{argmin}_{\mathcal{P}=\{\beta, \alpha, \epsilon\} \in \mathcal{Q}} \|\alpha\|_0$

# The Sparsest Rule (Cont.)

For penalized regression, this condition is known as "beta-min" condition. Notably  $|\tilde{\alpha}_j^c| = |\tilde{c} - c| \cdot |\gamma_j^*| > \kappa(n)$  depends on product of  $|\tilde{c} - c|$  and  $|\gamma_j^*|$ .

- 1 As discussed in Guo et al. (2018),  $|\tilde{c} - c|$  can not be too small to separate different solutions, but the larger gap  $|\tilde{c} - c|$  are helpful to mitigate the small or local to zero  $|\gamma_j^*|$  in favor of our model

## Example (continued)

- 1 Follow the procedure, we are able to reformulate total three solutions of (3):  $\alpha^* = (0_5, 1, 0.7_4)^\top$ ,  $\tilde{\alpha}^5 = (-0.2_3, -2.5_2, 0, 0.2_4)^\top$  and  $\tilde{\alpha}^7 = (-0.28_3, -3.5_2, -0.4, 0_4)^\top$ .
- 2 Thus, the sparsest rule  $\operatorname{argmin}_{\alpha \in \{\alpha^*, \tilde{\alpha}^5, \tilde{\alpha}^7\}} \|\alpha\|_0$  picks  $\alpha^*$  up, and Assumption 6 is easy to satisfy since fixed minimum absolute value except 0 are 0.7, 0.2, 0.28 in  $\alpha^*$ ,  $\tilde{\alpha}^5$ ,  $\tilde{\alpha}^7$ , respectively.

# Penalization Approaches with Embedded Surrogate Sparsest Rule

Consider the general penalized TSLS estimator based on moment conditions:

$$\left( \hat{\alpha}^{\text{pen}}, \hat{\beta}^{\text{pen}} \right) = \underset{\alpha, \beta}{\operatorname{argmin}} \underbrace{\frac{1}{2n} \|P_Z(\mathbf{Y} - \mathbf{Z}\alpha - \mathbf{D}\beta)\|_2^2}_{(I)} + \underbrace{p_{\lambda}^{\text{pen}}(\alpha)}_{(II)}. \quad (5)$$

They have the two different functions:

- 1 Approximate  $Q$ : (I) is a scaled finite sample version of  $E\left([\mathbf{Z}^{\top}\epsilon]^{\top} [\mathbf{Z}^{\top}\mathbf{Z}]^{-1} [\mathbf{Z}^{\top}\epsilon]\right)$ , which is a  $[\mathbf{Z}^{\top}\mathbf{Z}]^{-1}$  weighed quadratic term of condition  $E(\mathbf{Z}^{\top}\epsilon) = \mathbf{0}$ , and (II) is imposed to ensure sparsity structure in  $\alpha$ .
- 2 Rewrite (4) into equivalent constrained objective function form with the optimal penalty  $\|\alpha\|_0$  with respect to the sparsest rule:

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

The constraint above narrows the feasible solutions into  $Q$  because Sargan test statistics  $\|P_Z(\mathbf{Y} - \mathbf{D}\beta - \mathbf{Z}\alpha)\|_2^2 / \|(\mathbf{Y} - \mathbf{D}\beta - \mathbf{Z}\alpha)/\sqrt{n}\|_2^2 = O_p(1)$  (Sargan, 1958) under null hypothesis  $E(\mathbf{Z}^{\top}\epsilon) = \mathbf{0}$  as required in  $Q$ , otherwise  $O_p(n)$  that cannot be bounded by  $\delta$ . Therefore, the primary object of optimization in (4) is identification condition: the sparsest rule.

# Penalized method in literature

Kang et al. (2016) propose to use Lasso in (4), call sisVIVE. It has three problems:

- 1 Failure in consistent variable selection under some deterministic conditions, namely the sign-aware invalid IV strength (SAIS) condition:

$$\left| \hat{\gamma}_{\nu^{c*}}^T \text{sgn}(\alpha_{\nu^{c*}}) \right| > \|\hat{\gamma}_{\nu^*}\|_1 \quad (\text{Windmeijer et al., 2019, Proposition 2})$$

The SAIS is a rather common situation in practice, under which sisVIVE cannot achieve  $\mathcal{P}_0$ .

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

$$\left\| \hat{\beta}^{\text{sis}} - \beta^* \right\|_2 \leq \frac{|\hat{\mathbf{D}}^T \epsilon|}{\|\hat{\mathbf{D}}\|_2^2} + \frac{1}{\|\hat{\mathbf{D}}\|_2} \left( \frac{(4/3\sqrt{5})\lambda\sqrt{L\delta_{2L}^+(\mathbf{P}_{\hat{D}}\mathbf{Z})}}{2\delta_{2L}^-(\mathbf{Z}) - \delta_{2L}^+(\mathbf{Z}) - 2\delta_{2L}^+(\mathbf{P}_{\hat{D}}\mathbf{Z})} \right)$$

where  $\delta_k^{+/-}(\mathbf{H})$  refers to upper and lower RIP constant of matrix  $\mathbf{H}$ .

- 3 Objective function deviates from the original sparsest rule: As shown in Theorem 2 and Remark 1,  $g_1(\mathcal{P}) = \|\alpha\|_0$  and  $g_2(\mathcal{P}) = \|\alpha\|_1$  correspond to incompatible identification conditions unless satisfying an additional strong requirement

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

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

# The Surrogate Sparsest Penalty

How to solve?

- 1 and 2 corresponds to non-ignorable bias in Lasso and RIP conditions, respectively. It could be avoided by other penalties.
- 3 reveals the root of the problem: a proper surrogate penalty in (4) should align with identification condition.

One solution: Windmeijer et al. (2019) proposed to use Adaptive Lasso (Zou, 2006) with proper constructed initial estimator through median estimator. However, it requires the more stringent majority rule (than the sparsest rule, see Remark 1) and suffers from the same sensitivity issue on weak IVs as TSHT and CIIV.

## Proposition 1 (The proper Surrogate Sparsest penalty)

Suppose Assumptions 1-7 are satisfied. If  $p_{\lambda}^{\text{pen}}(\alpha)$  is surrogate sparsest rule in the sense of that it gives sparse solutions and

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

then  $p_{\lambda}^{\text{pen}}(\cdot)$  must be concave and  $p_{\lambda}^{\text{pen}}(t) = O(\lambda \kappa(n))$  for any  $t > \kappa(n)$ .



# WIT Estimator

We adopt the penalized method framework (10) and deploy a concave penalty in (11), the MCP in particular, which is nearly unbiased estimator.

## WIT Estimator

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

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

$$\hat{\kappa}_{\text{liml}} = \min_{\beta} \left\{ G(\beta) = \left( (\mathbf{Y} - \mathbf{D}\beta)^{\top} \mathbf{M}_{\mathbf{Z}} (\mathbf{Y} - \mathbf{D}\beta) \right)^{-1} \left( (\mathbf{Y} - \mathbf{D}\beta)^{\top} \mathbf{M}_{\mathbf{Z}_{\hat{\nu}^c}} (\mathbf{Y} - \mathbf{D}\beta) \right) \right\}$$

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  $\tilde{\mathbf{Z}}$  formally is defined as  $K_{\mathcal{C}} = K_{\mathcal{C}}(\mathcal{V}^*, \xi) := \inf_{\mathbf{u}} \left\{ \|\tilde{\mathbf{Z}}\mathbf{u}\|_2 / \left( \|\mathbf{u}\|_2 n^{1/2} \right) : \mathbf{u} \in \mathcal{C}(\mathcal{V}^*; \xi) \right\}$  and RE condition refers to that  $K_{\mathcal{C}}$  for  $\tilde{\mathbf{Z}}$  should be bounded away from zero.

## LEMMA 2. (RE condition of $\tilde{\mathbf{Z}}$ )

Under assumption A1 – 4, For any given  $\gamma^* \neq \mathbf{0}$ , there always exists a constant  $\xi \in (0, \|\hat{\gamma}_{\mathcal{V}^*}\|_1 / \|\hat{\gamma}_{\mathcal{V}^{c*}}\|_1)$  and further define the restricted cone  $\mathcal{C}(\mathcal{V}^*; \xi)$  such that  $K_{\mathcal{C}}^2 > 0$  holds strictly.

Lemma 2 elaborates RE condition of  $\tilde{\mathbf{Z}}$  holds without any additional assumptions on  $\tilde{\mathbf{Z}}$ , unlike sisVIVE. Moreover, this restricted cone is invariant of scaling and, thus, indicates accommodation of many weak IVs cases.

## Theorem 3 (Selection Consistency of Valid IVs)

Let  $\kappa(n)$  in Assumption 6 are specified as

$$\kappa(n) \asymp \underbrace{\sqrt{\frac{\log p}{n}}}_{T_1} + \underbrace{\frac{p}{n} \cdot \frac{\|Q_n \gamma^*\|_\infty}{\gamma^{*\top} Q_n \gamma^*}}_{T_2} + \underbrace{|\text{Bias}(\hat{\beta}_{or}^{TSLS})| \cdot \|\bar{\gamma}_{\mathcal{V}^{c*}}\|_\infty}_{T_3}, \quad (6)$$

where  $T_1 \rightarrow 0$  as  $n \rightarrow \infty$ . Suppose Assumptions 1-7 hold, we have

$$\hat{\alpha}^{\text{MCP}} = \underset{\hat{\alpha} \in \mathcal{B}_0(\lambda, \rho)}{\text{argmin}} \|\hat{\alpha}\|_0, \quad \Pr(\hat{\mathcal{V}} = \mathcal{V}^*, \hat{\alpha}^{\text{MCP}} = \hat{\alpha}^{\text{or}}) \xrightarrow{p} 1. \quad (7)$$

## Remark of $\kappa(n)$

### Proposition 1 (Magnitude of $T_2$ )

The magnitude of  $T_2$  is generally limited if there is not dominated  $\gamma_j$  in the sense of  $\|\gamma^*\|_\infty / \|\gamma^*\|_1 = o(\|\gamma\|_1/p)$ . Consider  $\mathbf{Q}_n = \mathbf{I}$  for simplicity,

$$T_2 = \frac{p}{n} \cdot \frac{\|\mathbf{Q}_n \gamma^*\|_\infty}{\gamma^{*\top} \mathbf{Q}_n \gamma^*} = \frac{p}{n} \cdot \frac{\|\gamma^*\|_\infty}{\|\gamma^*\|_2^2} \leq \frac{p}{n} \cdot \frac{p \|\gamma^*\|_\infty}{\|\gamma^*\|_1^2} \rightarrow 0.$$

### Proposition 2 (Approximation of Bias( $\hat{\beta}_{or}^{TSLs}$ ))

Let  $s = \max(\mu_n, L)$ , under the Assumptions 1-5, we obtain

$$E \left[ \text{Bias}(\hat{\beta}_{or}^{TSLs}) \right] = \frac{\sigma_{\epsilon\eta}}{\sigma_\eta^2} \left( \frac{L}{(\mu_n + L)} - \frac{2\mu_n^2}{(\mu_n + L)^3} \right) + o(s^{-1}). \quad (8)$$

## Theorem 4 (Consistency and Asymptotic Normality)

Under same condition in Theorem 3, together with Assumption A5 and conditional on  $\mathbf{Z}$ , we obtain:

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

$$\widehat{\text{Var}}(\hat{\beta}^{\text{WIT}}) = \frac{\hat{\mathbf{b}}^\top \hat{\boldsymbol{\Omega}} \hat{\mathbf{b}} (\hat{\mu}_n + L/n)}{-\hat{\mu}_n} \left( \hat{Q}_S \hat{\boldsymbol{\Omega}}_{22} - \boldsymbol{\tau}_{22} + \frac{\hat{c}}{1 - \hat{c}} \frac{\hat{Q}_S}{\hat{\mathbf{a}}^\top \hat{\boldsymbol{\Omega}}^{-1} \hat{\mathbf{a}}} \right)^{-1} \\ \xrightarrow{P} \mu_0^{-2} \left[ \sigma_\epsilon^2 \mu_0 + \frac{v_K(1-v_L)}{1-v_K-v_L} |\boldsymbol{\Sigma}| \right],$$

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

# Simulation

We consider the following four examples:

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

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

**Table:** Simulation results in low dimension

Case	Approaches	$n = 200$				$n = 500$			
		MAD	CP	FPR	FNR	MAD	CP	FPR	FNR
1(I)	TSLS	0.532	0.000	-	-	0.530	0	-	-
	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
1(II)	TSLS	1.072	0	-	-	1.100	0	-	-
	oracle-LIML	0.069	0.948	-	-	0.042	0.956	-	-
	TSHT	0.107	0.900	0.114	0.585	0.737	0.608	0.410	0.701
	CIIV	0.097	0.760	0.082	0.642	4.206	0.362	0.350	0.8084
	sisVIVE	0.247	-	0.007	0.216	0.155	-	0	0.176
	Post-Alasso	1.833	0	0.423	0.251	3.552	0	0.580	0.385
	WIT	0.077	0.856	0.055	0.056	0.045	0.918	0.015	0.030

# Simulation (Cont)

Further, we present a replication of simulation design in literature and its variant:

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

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

**Table:** Simulation results in low dimension: A replication of experiment

Case	Approaches	$n = 500$				$n = 1000$			
		MAD	CP	FPR	FNR	MAD	CP	FPR	FNR
1(III)	TSLs	0.436	0	-	-	0.435	0	-	-
	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
1(IV)	TSLs	1.124	0	-	-	1.144	0	-	-
	oracle-LIML	0.056	0.948	-	-	0.042	0.948	-	-
	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	0.526	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	0.048	0.844	0.068	0.090

# Simulation (Cont.)

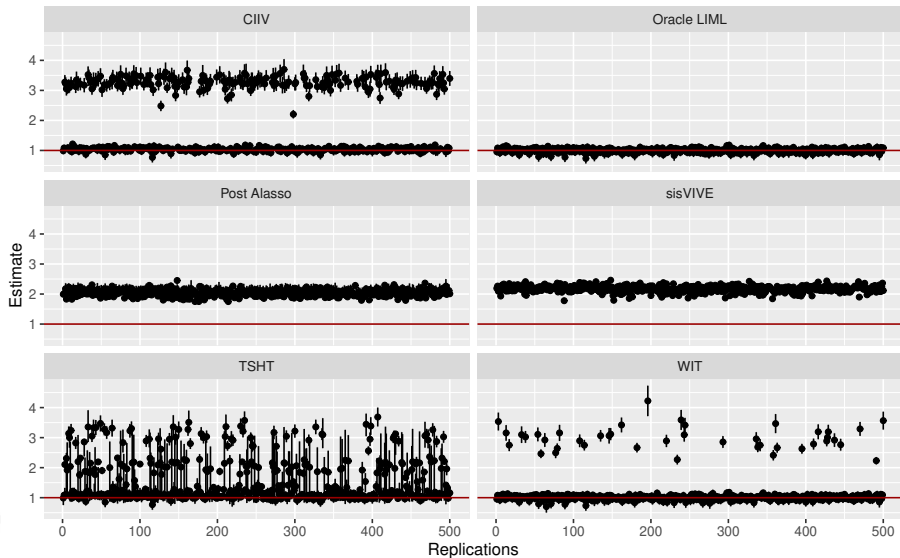


Figure: Scatter plot of estimations of  $\beta^*$  with confidence intervals of Case 1 (IV)

# Real Application

- We revisit the classic empirical study in trade and growth (Frankel and Romer, 1999, FR99 henceforth).
- We investigate the causal effect of trade on income using a 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,

$$\log(Y_i) = \alpha + \beta T_i + \psi S_i + \epsilon_i \quad (9)$$

where for each country  $i$ ,  $Y_i$  is the GDP per worker,  $T_i$  is the share of international trade to GDP,  $S_i$  is the size of the country, such as area, population, and  $\epsilon_i$  is the error term.

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



# Trade on GDP

Following the same logic, Fan and Zhong (2018) extended the IV set to include more geographic and meteorological variables. The reduced form equation is

$$T_i = \gamma^T \mathbf{Z}_i + \nu_i, \quad (10)$$

where  $\mathbf{Z}_i$  is a vector of instruments.

**Table:** 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.

# Trade on GDP

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

$$\log(Y_i) = T_i\beta + \mathbf{Z}_i^\top \alpha + \mathbf{S}_i^\top \psi + \epsilon_i \quad \text{for } i = 1, 2, \dots, 158, \quad (11)$$

Partial out of Control:  $\ddot{Y}_i = \ddot{T}_i\beta + \ddot{\mathbf{Z}}_i^\top \alpha + \ddot{\epsilon}_i, \quad \ddot{T}_i = \ddot{\mathbf{Z}}_i^\top \phi + \ddot{\nu}_i.$

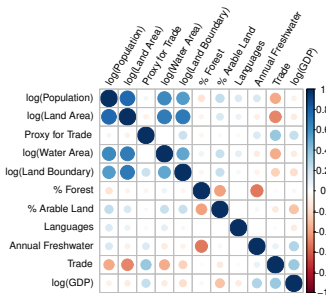


Figure: Correlation of all variables

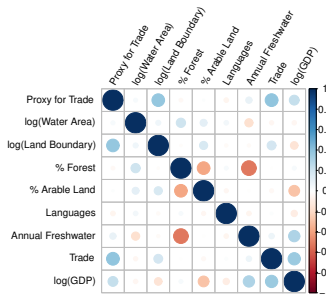


Figure: Correlation of transformed variables

**Table:** 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 means  $p$ -value of Sargan test and selection of relevant IVs  $\hat{\mathcal{S}}$  is only be implemented in TSHT and CIIV.

# Observations

## Observation:

- $p$ -value of the Hausman test for endogeneity is 0.000181 using the proxy for trade as IV.
- LIML using all potential IVs (without distinguishing the invalid ones) likely overestimates the treatment effect. The 0.001  $p$ -value of Sargan test strongly reject the null of all potential IVs are valid.
- $Z_5$  should be a invalid IV:
  - 1 In view of marginal correlation in Fig. 6,  $Z_5$  is nearly uncorelated to trade but significantly correlated with  $\log(\text{GDP})$ .
  - 2 Concerning the Sargen Test,  $p$ -value of  $0.086 < 0.1$  in Post-Lasso indicates  $Z_5$  is not very credible to be valid.
  - 3 In the economic perspective, more arable land generates higher crop yields and maintains a higher agriculture sector labor force, which directly affects GDP.
- Strong IVs based method fails.

- Anderson, J. E. (1979). A theoretical foundation for the gravity equation. *Am. Econ. Rev.*, 69(1):106–116.
- Fan, Q. and Zhong, W. (2018). Nonparametric additive instrumental variable estimator: A group shrinkage estimation perspective. *J. Bus. Econ. Statist.*, 36(3):388–399.
- 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(513):132–144.
- Small, D. S. (2007). Sensitivity analysis for instrumental variables regression with overidentifying restrictions. *J. Am. Statist. Ass.*, 102(479):1049–1058.
- 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(527):1339–1350.
- Wooldridge, J. M. (2010). *Econometric analysis of cross section and panel data*. MIT press.