

# Inference for high-dimensional instrumental variables regression

David Gold

Department of Statistics  
University of Washington

October 2, 2017

# Overview

- Paper: “Inference for high-dimensional nested regression” with Johannes Lederer (Statistics & Biostatistics) and Jing Tao (Economics). Available at <https://arxiv.org/abs/1708.05499>
- Contributions: We develop methods to estimate and conduct statistical inference for low-dimensional components of a high-dimensional regression vector under endogeneity of the regressors.

# Endogeneity

- Consider the linear model

$$Y = X\beta + U,$$

where  $Y$  is the response,  $X$  is a regressor whose effect on  $Y$  is of interest, and  $U$  is a noise term.

- In many empirical settings (especially in econometrics),  $X$  is **endogenous**:

$$E[XU] \neq 0.$$

- Endogeneity prevents us from discerning co-variation between  $X$  and  $Y$  due to a structural relationship from co-variation that is due to the correlation of random quantities that affect each or both variables.

# Endogeneity: Sources

Three main sources of endogeneity are:

- **confounding**
- **errors in measurements** of predictor variables, and
- the **interplay of mutually influential processes** that also exhibit random variation.

# Endogeneity: Example 1 — Supply & demand

Supply and demand curves

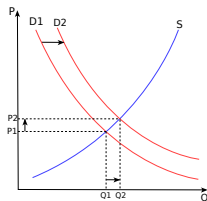
both live in the *quantity* ( $Q$ )  $\times$  *price* ( $P$ ) space.

Consider the linear supply and

demand curve equations with random disturbances:

$$Q = \theta_D P + u_D \quad (\text{Demand})$$

$$Q = \theta_S P + u_S \quad (\text{Supply}).$$



Some algebra shows that  $P$  has non-trivial covariance with  $u_D$  due to the mutual influence of  $Q$  and  $P$  on one another.

As a result, OLS estimates of  $\theta_D$  (and of  $\theta_S$ ) are biased.

## Endogeneity: Example 2 — Returns to schooling

- Suppose we are interested in the question of how educational attainment influences adult wages.
- We adopt the model

$$\underbrace{Y}_{\log(\text{wage})} = \underbrace{X}_{\text{education}} \beta + \underbrace{W}_{\substack{\text{academic} \\ \text{aptitude}}} \gamma + U,$$

but do not observe  $W$ . Rather, we can only fit the model

$$Y = X\beta + \tilde{U}, \quad \tilde{U} := W\gamma + U.$$

- If  $E[XW] \neq 0$ , then  $X$  is endogenous relative to  $\tilde{U}$ .
- See Angrist and Krueger (1991).

# Instrumental Variables

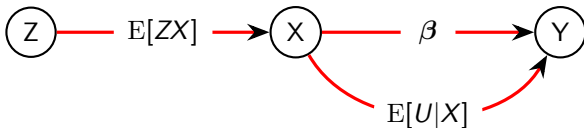
- Consider again the linear model  $Y = X\beta + U$
- Suppose that  $Z$  is a third variable that satisfies

$$E[ZU] = 0, \quad E[ZX] \neq 0.$$

- Then, we can write

$$\begin{aligned} E[ZY] &= E[ZX]\beta + E[ZU] \\ \Rightarrow \quad \beta &= E[ZY]/E[ZX] \end{aligned} .$$

- Such a  $Z$  is called an **instrumental variable** (IV).



# Model

Our model of interest is the *linear instrumental variables model* (Amemiya (1985)):

$$y_i = \mathbf{x}_i^\top \boldsymbol{\beta}_0 + u_i,$$

$$x_{ij} = \mathbf{z}_i^\top \boldsymbol{\alpha}_0^j + v_{ij} =: \underbrace{\mathbf{d}_i}_{=\mathbf{z}_i^\top \boldsymbol{\alpha}^j} + v_{ij},$$

where:

- the vector  $\mathbf{x}_i \in \mathbb{R}^{p_x}$  consists of the *endogenous variables*  $x_{i1}, \dots, x_{ip_x}$ ;
- the vector  $\boldsymbol{\beta}_0 \in \mathbb{R}^{p_x}$  is the parameter of interest;
- the vector  $\mathbf{z}_i \in \mathbb{R}^{p_z}$  consists of the *instrumental variables*  $z_{i1}, \dots, z_{ip_z}$ , which satisfy  $E[\mathbf{z}_i] = \mathbf{0}$ ;
- and the vectors  $\boldsymbol{\alpha}_0^j \in \mathbb{R}^{p_z}$  are regression parameters up to which the respective conditional means  $d_{ij} := E[x_{ij} | \mathbf{z}_i] = \mathbf{z}_i^\top \boldsymbol{\alpha}^j$



# Model

We assume that the noise elements  $u_i, \mathbf{v}_i$  satisfy

$$(u_i, \mathbf{v}_i) | \mathbf{z}_i \sim \mathcal{N}_{1+p_x}(\mathbf{0}, \Sigma_{u\mathbf{v}}), \quad \Sigma_{u\mathbf{v}} := \begin{pmatrix} \sigma_u^2 & \sigma_{u\mathbf{v}}^\top \\ \sigma_{u\mathbf{v}} & \Sigma_{\mathbf{v}} \end{pmatrix},$$

where

- $\sigma_{u\mathbf{v}} := (\sigma_{u\mathbf{v}^1}, \dots, \sigma_{u\mathbf{v}^{p_x}})^\top$  consists of the noise covariances  $\sigma_{u\mathbf{v}^j} := \text{cov}(\mathbf{u}, \mathbf{v}^j)$
- $\Sigma_{\mathbf{v}}$  is an unstructured covariance matrix with diagonal entries  $\sigma_{\mathbf{v}^j}^2 := \text{var}(\mathbf{v}^j)$  for  $j \in [p_x]$ .

# Model: Regularity conditions

- The first-stage regression parameters  $\alpha_{0,j}$  satisfy

$$\max_{j \in [p_x]} \|\alpha_{0,j}\|_1 \leq M_A$$

for a universal constant  $M_A$ ;

- The second-stage regression parameter  $\beta_0$  satisfies

$$\|\beta_0\|_1 \leq M_\beta$$

for a universal constant  $M_\beta$ ;

- The variances  $\sigma_u^2$  and  $\sigma_{\nu_j}^2$  are bounded strictly away from zero and infinity for  $j \in [p_x]$ .
- The minimal and maximal eigenvalues of  $\Sigma_z = E[\mathbf{z}_i \mathbf{z}_i^\top]$  are strictly bounded away from zero and infinity, respectively.

# Model: Matrix notation

In matrix notation, we write

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta}_0 + \mathbf{u}$$

and

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

where:

- the vectors  $\mathbf{y}, \mathbf{u} \in \mathbb{R}^n$  consist of the responses  $y_i$  and the noise components  $u_i$ , respectively;
- the matrix  $\mathbf{X} \in \mathbb{R}^{n \times p_x}$  has columns  $\mathbf{x}^j$  given by  $\mathbf{x}^j = (x_{1j}, \dots, x_{nj})^\top$ ;
- the matrix  $\mathbf{D} = \mathbb{E}[\mathbf{X}|\mathbf{Z}] \in \mathbb{R}^{n \times p_x}$  has columns  $\mathbf{d}^j$  given by  $\mathbf{d}^j = (d_{1j}, \dots, d_{nj})^\top$ ;
- the matrix  $\mathbf{Z} \in \mathbb{R}^{n \times p_z}$  has columns  $\mathbf{z}^k$  given by  $\mathbf{z}^k = (z_{1k}, \dots, z_{nk})^\top$ ;
- the matrix  $\mathbf{A} \in \mathbb{R}^{p_z \times p_x}$  has columns given by  $\boldsymbol{\alpha}_0^j$ .

# Methodology overview

- IV methods were first used by Wright (1928) to estimate supply and demand elasticities for flaxseed (though see Stock and Trebbi (2003)).
- Early work in IV methods considered the *two-stage least squares (2SLS)* estimator (Theil (1953), Basmann (1957), Sargan (1958))

$$\hat{\beta}^{2SLS} = (\hat{\mathbf{D}}^\top \hat{\mathbf{D}})^{-1} \hat{\mathbf{D}} \mathbf{y},$$

where  $\hat{\mathbf{D}}$  is the OLS prediction of  $\mathbf{D} = \mathbb{E}[\mathbf{X}|\mathbf{Z}]$ , and the *limited information maximum likelihood (LIML)* estimator (Anderson and Rubin (1950)); see also Anderson (2005).

- More recently, Belloni et al. (2011) and Belloni et al. (2012) study the use of the Lasso in the first-stage prediction of  $\mathbf{D}$  when the first-stage regression vectors  $\alpha_{0,j}$  are high-dimensional.

## Two-stage estimation

Similar to the 2SLS estimator in the low-dimensional linear IV model, we use a two-stage Lasso estimator to contend with both high-dimensionality and endogeneity:

- Definition (first-stage Lasso estimator):

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

- Definition (second-stage Lasso estimator):

$$\hat{\beta} \in \arg \min_{\beta \in \mathbb{R}^{p_x}} \left\{ \|\mathbf{y} - \hat{\mathbf{D}}\alpha\|_2^2 / (2n) + r_\beta \|\beta\|_1 \right\},$$

where  $\hat{\mathbf{D}} = \mathbf{Z}\hat{\mathbf{A}}$  and  $\hat{\mathbf{A}} = (\hat{\alpha}_1, \dots, \hat{\alpha}_{p_z})$ .

# Compatibility condition

- Definition (Compatibility condition): For a given index set  $S \in [p]$  define the set

$$\mathcal{C}(S) := \{ \boldsymbol{\delta} \in \mathbb{R}^p \setminus \mathbf{0} : \|\boldsymbol{\delta}\|_1 \leq C, \|\boldsymbol{\delta}_{S^c}\|_1 \leq 3\|\boldsymbol{\delta}_S\|_1 \}.$$

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

$$\phi^2 \leq \inf_{\boldsymbol{\delta} \in \mathcal{C}(S)} \frac{\sqrt{|S|} \|\mathbf{M}\boldsymbol{\delta}\|_2^2}{n \|\boldsymbol{\delta}\|_1^2}$$

holds. We call such a quantity  $\phi^2$  the *compatibility constant*.

- Related to the *restricted eigenvalue condition* of Bickel et al. (2009); see also (Bühlmann and van de Geer, 2011, Chapter 6) and van de Geer and Bühlmann (2009).

# Estimation error bounds: First-stage

- Assumption (First-stage compatibility conditions) For each active set  $S_j := \text{supp } \alpha^j$  of the first-stage model, there exists a constant  $\phi_j > 0$  such that  $\mathbf{Z}$  satisfies the compatibility condition relative to  $S_j$  and  $\phi_j$ .
- Lemma** (First-stage estimation error) Write

$$\delta_{\hat{\mathbf{A}}} := \max_{j \in [p_x]} \|\hat{\alpha}_j - \alpha_{0,j}\|_1.$$

Suppose that the above assumption holds. Set  $r_j = \sigma_{vj} c \sqrt{\hat{\sigma}_{\mathbf{Z}} \log p_{\mathbf{Z}} / n}$ , and set  $r_{\mathbf{A}} = \max_{j \in [p_x]} r_j$ . Then,

$$\mathbb{P}\{\delta_{\hat{\mathbf{A}}} > 4s_{\mathbf{A}}r_{\mathbf{A}}/\phi_{\mathbf{A}}^2\} \leq 2p_{\mathbf{Z}}^{1-c_{\text{ep}}},$$

where  $c_{\text{ep}} := c^2/32 - 1$ , where  $s_{\mathbf{A}} = \max_{j \in [p_x]} |S_j|$  and  $\phi_{\mathbf{A}} = \min_{j \in [p_x]} \phi_j$ .

# Estimation error bounds: Second-stage

- First-stage estimation error bounds are straightforward union bounds.
- Second-stage bounds are trickier because we need to account for the predictions  $\widehat{\mathbf{D}}$ .
- To take advantage of standard Lasso oracle inequalities for the second-stage Lasso estimator, we must write

$$\mathbf{y} = \mathbf{X}\beta_0 + \mathbf{u} = \widehat{\mathbf{D}}\beta_0 + \tilde{\mathbf{u}},$$

where  $\tilde{\mathbf{u}} := \mathbf{u} + [(\mathbf{D} - \widehat{\mathbf{D}}) + \mathbf{V}]\beta_0$ .



# Estimation error bounds: Second-stage

- Using standard oracle inequalities, we find that, on the set  $\{4\|\hat{\mathbf{D}}^\top \tilde{\mathbf{u}}/n\|_\infty \leq r_\beta\}$ , we have

$$\|\hat{\beta} - \beta_0\|_1 \leq 4s_\beta r_\beta / \phi_\beta,$$

where  $s_\beta = |\text{supp } \beta_0|$ .

- Using Hölder's inequality many times, we also find that

$$\begin{aligned} \|\hat{\mathbf{D}}^\top \tilde{\mathbf{u}}/n\|_\infty &\leq \delta_{\hat{\mathbf{A}}} \hat{\sigma}_{\mathbf{Z}} (\delta_{\hat{\mathbf{A}}} \|\beta\|_1 + M_{\mathbf{A}}) \\ &\quad + (\delta_{\hat{\mathbf{A}}} + M_{\mathbf{A}}) (\hat{\sigma}_{\mathbf{Z}, \mathbf{v}} \|\beta\|_1 + \hat{\sigma}_{\mathbf{Z}, \mathbf{u}}), \end{aligned}$$

where  $\delta_{\hat{\mathbf{A}}} = \max_{j \in [p_x]} \|\hat{\alpha}_j - \alpha_{0,j}\|_1$  and  $\hat{\sigma}_{\mathbf{Z}, \cdot} := \|\mathbf{Z}^\top \cdot\|_\infty$ .

# Estimation error bounds: Second-stage

- Assumption (Second-stage compatibility condition): There exists a constant  $\phi_\beta$  such that  $\hat{\mathbf{D}}$  satisfies the compatibility condition with respect to  $S_\beta = \text{supp } \beta_0$  and  $\phi_\beta$ .
- Lemma** (Second-stage estimation error): Suppose that the two Assumptions above hold. For each  $j \in [p_x]$ , set  $r_j$  and  $r_A$  as in the previous Lemma; set

$$\lambda_{\mathbf{V}} = \sigma_{\mathbf{V}} c \sqrt{\hat{\sigma}_{\mathbf{Z}} \log p_{\mathbf{Z}} / n}, \quad \lambda_{\mathbf{U}} = \sigma_{\mathbf{U}} c \sqrt{\hat{\sigma}_{\mathbf{Z}} \log p_{\mathbf{Z}} / n},$$

$$r_\beta = 16\psi_{\mathbf{A}} r_{\mathbf{A}} \hat{\sigma}_{\mathbf{Z}} (4M_\beta \psi_{\mathbf{A}} r_{\mathbf{A}} + M_{\mathbf{A}}) + (4\psi_{\mathbf{A}} r_{\mathbf{A}} + M_{\mathbf{A}}) (M_\beta \lambda_{\mathbf{V}} + \lambda_{\mathbf{U}}).$$

Then,

$$\mathbb{P}\{\|\hat{\beta} - \beta\|_1 > 4s_\beta r_\beta / \phi_\beta^2\} \leq 2p_{\mathbf{Z}}^{1-c_{\text{ep}}} + 2p_{\mathbf{Z}}^{-c_{\text{ep}}}.$$

# Estimation error bounds: Second-stage compatibility condition

- van de Geer and Bühlmann (2009) and others discuss the tenability of the compatibility condition in the context of the ordinary linear model.
- Since  $\hat{\mathbf{D}}$  differs from an ordinary design matrix, we need to do extra work.

# Estimation error bounds: Second-stage compatibility condition

**Lemma** (Second-stage compatibility constant) Let  $S \subseteq [p_x]$  be an arbitrary index set with  $s = |S|$ , and let  $c_S > 0$  and  $C > 0$  be arbitrary. For a given matrix  $\mathbf{M} \in \mathbb{R}^{n \times p}$ , define the quantity

$$\phi_{\dagger}^2(\mathbf{M}, S, c_S, C) = \inf_{\substack{\delta \in \mathcal{C}^p(S, c_S): \\ \|\delta\|_1 \leq C}} \frac{s \|\mathbf{M}\delta\|_2^2}{n \|\delta_S\|_1^2}.$$

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

$$\begin{aligned} & \mathbb{P} \left\{ \phi_{\dagger}^2(\widehat{\mathbf{D}}, S, c_S, C) < (1 - \epsilon_1)(\Lambda_{\min}(\Sigma_{\mathbf{d}}) - (1 + c_S)^2 \epsilon_2) \right\} \\ & \leq \mathbb{P} \left\{ (2M_{\mathbf{A}} \delta_{\widehat{\mathbf{A}}} + \delta_{\widehat{\mathbf{A}}}^2) \hat{\sigma}_{\mathbf{Z}} C^2 > \epsilon_1 \right\} + \mathbb{P} \left\{ s \|\bar{\Sigma}_{\mathbf{d}} - \Sigma_{\mathbf{d}}\|_{\infty} > \epsilon_2 \right\}, \end{aligned}$$

where  $\bar{\Sigma}_{\mathbf{d}} = \mathbf{D}^{\top} \mathbf{D} / n$  and  $\Lambda_{\min}(\Sigma_{\mathbf{d}})$  denotes the minimal eigenvalue of  $\Sigma_{\mathbf{d}}$ .

# Inference for high-dimensional regression parameters

- Knight and Fu (2000) show that the limiting distribution of  $\sqrt{n}(\hat{\beta} - \beta_0)$  need be neither normal nor pivotal with respect to  $\beta_0$  in the low-dimensional case; see also Pötscher and Leeb (2009).
- There exist a number of recent proposals for conducting statistical inference for high-dimensional regression parameters: Zhang and Zhang (2014), Javanmard and Montanari (2014), van de Geer et al. (2014), Ning and Liu (2014), Neykov et al. (2015).

# One-step update: Newton-Raphson method

- Suppose that we have a system of  $p_x$  equations

$$\mathbf{f}(\boldsymbol{\beta}) = \mathbf{0}.$$

- The Newton-Raphson method is to take steps

$$\boldsymbol{\beta}^{k+1} = \boldsymbol{\beta}^k - \left[ \frac{\partial \mathbf{f}}{\partial \boldsymbol{\beta}} \Big|_{\boldsymbol{\beta}^k} \right]^{-1} \mathbf{f}(\boldsymbol{\beta}^k)$$

in the direction of the solution in  $\boldsymbol{\beta}$ .

# One-step update

- In the linear model with Gaussian errors, the  $p_x$ -dimensional vector of score functions  $\mathbf{h}(\mathbf{x}_i; \boldsymbol{\beta})$  satisfy

$$\mathbb{E}[\mathbf{h}(\mathbf{x}_i; \boldsymbol{\beta}_0)] = \mathbb{E}[-\mathbf{x}_i(y_i - \mathbf{x}_i^\top \boldsymbol{\beta}_0)] = \mathbf{0}.$$

- The **one-step update** to an initial estimator  $\hat{\boldsymbol{\beta}}$

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

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

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

is one Newton-Raphson step in the direction of the solution in  $\boldsymbol{\beta}$  to the empirical analogue

$$\mathbb{E}_n[\mathbf{h}(\mathbf{x}_i; \boldsymbol{\beta})] = -\mathbf{X}^\top (\mathbf{y} - \mathbf{X} \boldsymbol{\beta}) / n = \mathbf{0}.$$

# One-step update

- If  $\hat{\Sigma}_x$  is invertible, set  $\hat{\Theta} = \hat{\Sigma}_x^{-1}$ . Then,

$$\tilde{\beta} = \hat{\beta} + \hat{\Theta} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\beta}) / n = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.$$

- When  $p_x > n$ ,  $\hat{\Sigma}_x$  is not invertible, and one finds

$$\begin{aligned} \tilde{\beta} &= \hat{\beta} + \hat{\Theta} \mathbf{X}^\top (\mathbf{y} - \mathbf{X} \hat{\beta}) / n \\ &= \hat{\beta} + \hat{\Theta} \mathbf{X}^\top (\mathbf{X}(\beta_0 - \hat{\beta}) + \mathbf{u}) / n \\ &= \beta_0 + \hat{\Theta} \mathbf{X}^\top \mathbf{u} / n + (\hat{\Theta} \hat{\Sigma}_x - \mathbf{I})(\beta_0 - \hat{\beta}). \end{aligned}$$

Thus

$$\sqrt{n}(\tilde{\beta} - \beta_0) = \hat{\Theta} \mathbf{X}^\top \mathbf{u} / \sqrt{n} + \underbrace{\sqrt{n}(\hat{\Theta} \hat{\Sigma}_x - \mathbf{I})(\beta_0 - \hat{\beta})}_{\Delta}.$$



# One-step update

- The relationship

$$\begin{aligned}\sqrt{n}(\tilde{\beta}_j - \beta_{0,j}) &= \hat{\boldsymbol{\theta}}_j^\top \mathbf{X}^\top \mathbf{u} / \sqrt{n} + \Delta_j \\ &= \frac{1}{\sqrt{n}} \sum_i \langle \hat{\boldsymbol{\theta}}_j, \mathbf{x}_i \rangle u_i + \Delta_j\end{aligned}$$

suggests a method for conducting statistical inference for  $\beta_{0,j}$ .

- If  $\Delta = o_P(1)$ , and if  $\hat{\boldsymbol{\theta}}_j, \mathbf{x}_i \perp u_i$ , then

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

where

$$\omega_j^2 = \mathbb{E}[\langle \hat{\boldsymbol{\theta}}_j, \mathbf{x}_i \rangle^2 u_i^2].$$

- This strategy is studied in Javanmard and Montanari (2014); van de Geer et al. (2014), though these authors rather proceed directly from the KKT conditions for the Lasso.

## One-step update: With endogeneity

- If the errors are not Gaussian, then  $\mathbf{h}(\mathbf{x}_i; \boldsymbol{\beta}) = -\mathbf{x}_i(\mathbf{y}_i - \mathbf{x}_i^\top \boldsymbol{\beta})$  is not a score per se but may still satisfy  $E\mathbf{h}(\mathbf{x}_i; \boldsymbol{\beta}_0) = \mathbf{0}$  as part of an **orthogonality condition**

$$E[\mathbf{x}_i u_i] = \mathbf{0}.$$

- In our case, the  $\mathbf{x}_i$  are endogenous, so the above doesn't hold. Instead,

$$E[\mathbf{z}_i | u_i] = \mathbf{0} \Rightarrow E[\mathbf{z}_i u_i] = \mathbf{0},$$

and hence

$$E[\mathbf{d}_i u_i] = \mathbf{0}.$$

- Can we devise a one-step estimator that takes such information into account?
- (Note that we must substitute a prediction for  $\mathbf{d}_i$ , since the latter is unavailable...)

# One-step update: With endogeneity

- The empirical analogue of  $\mathbb{E}[\hat{\mathbf{d}}_i u_i] = \mathbf{0}$  as a function of  $\beta$  is

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

- The one-step update to an initial estimator  $\hat{\beta}$  in the direction to the solution of the above is

$$\begin{aligned} \tilde{\beta} &= \hat{\beta} - \hat{\Theta} \mathbb{E}_n[\tilde{\mathbf{h}}(\beta)] \\ &= \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top(\mathbf{y} - \mathbf{X}\hat{\beta})/n \\ &= \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top(\mathbf{X}[\beta_0 - \hat{\beta}] + \mathbf{u})/n \\ &= \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top(\hat{\mathbf{D}}[\beta_0 - \hat{\beta}] + [\mathbf{X} - \hat{\mathbf{D}}][\beta_0 - \hat{\beta}] + \mathbf{u})/n \\ &= \beta_0 + \underbrace{\hat{\Theta} \hat{\mathbf{D}}^\top \mathbf{u}/n}_{\Delta_3/\sqrt{n}} + \underbrace{\hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{X} - \hat{\mathbf{D}})(\beta_0 - \hat{\beta})/n}_{\Delta_4/\sqrt{n}} + (\hat{\Theta} \hat{\Sigma}_d - \mathbf{I})(\beta_0 - \hat{\beta}) . \end{aligned}$$

# One-step update: With endogeneity

- What about  $\hat{\Theta}$ ? Note that  $\frac{\partial \mathbb{E}_n[\tilde{h}(\beta)]}{\partial \beta} = \hat{D}^\top X/n$ , but

$$\Delta_4 = (\hat{\Theta} \hat{\Sigma}_d - I)(\beta_0 - \hat{\beta}).$$

Taking  $\hat{\Theta}$  as an estimator of  $\Theta := E[\hat{\Sigma}_d]^{-1} =: \Sigma_d^{-1}$  works.

- Now write

$$\begin{aligned} \sqrt{n}(\tilde{\beta} - \beta_0) &= \hat{\Theta} \hat{D}^\top u / \sqrt{n} + \Delta_3 + \Delta_4 \\ &= \Theta D^\top u / \sqrt{n} + \underbrace{(\hat{\Theta} - \Theta) D^\top u / \sqrt{n}}_{\Delta_1} + \underbrace{\hat{\Theta} (\hat{D} - D)^\top u / \sqrt{n}}_{\Delta_2} \\ &\quad + \Delta_3 + \Delta_4. \end{aligned}$$

# One-step update: Recap

- Goal: Inference for components  $\beta_{0,j}$  under high-dimensionality of  $\beta_0$  and endogeneity of  $\mathbf{x}_i$ .
- Method: Adapt recent literature on one-step update for Lasso to present case to define

$$\tilde{\beta} = \hat{\beta} - \hat{\Theta} \mathbb{E}_n[\tilde{\mathbf{h}}(\hat{\beta})] = \hat{\beta} + \hat{\Theta} \hat{\mathbf{D}}^\top (\mathbf{y} - \mathbf{X} \hat{\beta}) / n,$$

which also relates to work by Amemiya (1974, 1977); Hansen (1982); Chamberlain (1987); Newey (1990) on *generalized method of moments* (GMM) IV estimators.

- Upshot: We obtained the relationship

$$\begin{aligned} \sqrt{n}(\tilde{\beta}_j - \beta_{0,j}) &= \boldsymbol{\theta}_j \mathbf{D}^\top \mathbf{u} / \sqrt{n} + \sum_{\ell=1}^4 \Delta_{\ell,j} \\ &= \frac{1}{\sqrt{n}} \sum_i \langle \boldsymbol{\theta}_i, \mathbf{d}_i \rangle u_i + \sum_{\ell=1}^4 \Delta_{\ell,j}. \end{aligned}$$

# Inverse estimation

- van de Geer et al. (2014) construct  $\hat{\Theta}$  from the Nodewise Lasso estimates (Meinshausen and Bhlmann (2006))

$$\hat{\gamma}_j := \arg \min_{\gamma \in \mathbb{R}^{p_x-1}} \|\mathbf{x}_j - \mathbf{X}_{-j}\gamma\|_2^2/n + 2\lambda_j \|\gamma\|_1.$$

- Javanmard and Montanari (2014) construct  $\hat{\Theta}$  with rows  $\hat{\theta}_j$  as solutions to

$$\underset{\theta \in \mathbb{R}^{p_x}}{\text{minimize}} \quad Q(\theta) \quad \text{subject to } \|\hat{\Sigma}_x - \theta\|_\infty \leq \mu$$

with  $Q(\theta) = \theta^\top \hat{\Sigma}_x \theta$ .

- We use essentially the CLIME estimator of Cai et al. (2011), which solves the above program with  $Q(\theta) = \|\theta\|_1$ .

# Inverse estimation

- The CLIME estimator has established rates for  $\ell_1$  estimation error, whereas the Javanmard and Montanari (2014) estimator does not.
- In order to obtain the  $\ell_1$  rates for the CLIME, we must impose additional assumptions on  $\hat{\Theta}$ .

# Inverse estimation

- Definition (Uniformity class of Cai et al. (2011)): The *uniformity class* of population inverses  $\hat{\Theta}$  is defined to be

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

- Theorem** ( $\ell_1$  estimation error of  $\hat{\theta}_j$ ): Suppose that: (i)  $\hat{\Theta} \in \mathcal{U}(M_{\Theta}, q, s_{\Theta})$ ; (ii)  $\hat{\theta}_j$  is the CLIME estimator of  $\theta_j$ ; (iii)  $\theta_j$  is feasible for the CLIME program. Then,

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

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



## Remainder terms

- Under model regularity and rate conditions, the remainder terms  $\Delta_\ell$  each satisfy  $\Delta_\ell = o_P(1)$ .
- For sub-Gaussian  $\mathbf{Z}$ , choosing  $\mu$  on the order of  $\sqrt{\log(p_x)/n}$  achieves balance between feasibility of  $\theta_j$  for inverse program and the needs of the rates for the remainder terms.

# Asymptotic normality

- Under model regularity and rate conditions, we then have

$$\sqrt{n}(\tilde{\beta}_j - \beta_{0,j}) = n^{-1/2} \sum_i \langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle u_i + o_P(1).$$

- Now define

$$\begin{aligned} \omega_j^2 &= \mathbb{E}[\langle \boldsymbol{\theta}_j, \mathbf{d}_i \rangle^2 u_i^2] \\ &= \boldsymbol{\theta}_j^\top \mathbb{E}[\mathbf{d}_i \mathbf{d}_i^\top] \boldsymbol{\theta}_j \mathbb{E}[u_i^2] \\ &= \sigma_u^2 \boldsymbol{\Theta}_{jj}. \end{aligned}$$

# Asymptotic normality

**Theorem** (Asymptotic normality of  $\sqrt{n}(\tilde{\beta}_j - \beta_{0,j})/\omega_j$ ): Let  $\tilde{\beta}$  denote the one-step second-stage Lasso estimator of  $\beta_0$ , and let  $\omega_j$  be as above. Then, under model regularity and rate conditions, it holds that

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

Moreover, the result continues to hold if  $\omega_j$  is replaced with

$$\hat{\omega}_j := \hat{\sigma}_u \hat{\Theta}_{jj},$$

where  $\hat{\sigma}_u$  is an estimator of the second-stage noise level that satisfies  $\hat{\sigma}_u/\sigma_u - 1 = o_P(1)$ .

# Design: Data-generating mechanism

- Choose  $\alpha_{0,j}, \beta_0$  to have  $s_{\alpha^j}, s_{\beta}$  non-zero entries equal to 1.
- Draw observations  $(y_i, \mathbf{x}_i, \mathbf{z}_i)$  according to

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

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

- For all configurations, we set

$$\Sigma_{uv} = \begin{pmatrix} \sigma_u & \sigma_{uv} \mathbf{1}^\top \\ \sigma_{uv} \mathbf{1} & \sigma_v \mathbf{I} \end{pmatrix},$$

where  $\sigma_u, \sigma_v = 1$  are held fixed but  $\sigma_{uv}$  varies over configurations.

# Design: Choices of $\Sigma_z$

- The first choice of  $\Sigma_z$  is of a Toeplitz (TZ) structure given by

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

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

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

# Design: Metrics

- We generate  $N = 100$  trials for each configuration of

$$n, p_x, p_z, s_\beta, s_A, b, a, \Sigma_z, \sigma_{uv}.$$

- In each trial, for each component  $\tilde{\beta}_j$  of  $\tilde{\beta}$  we compute the respective  $(1 - \alpha)\%$  confidence interval

$$\hat{\mathcal{I}}_{\alpha,j} = [\tilde{\beta}_j - z_{0.05} \widehat{\text{SE}}(\tilde{\beta}_j), \tilde{\beta}_j + z_{0.05} \widehat{\text{SE}}(\tilde{\beta}_j)],$$

where  $z_{0.05} = \Phi^{-1}(1 - 0.05/2)$  and  $\widehat{\text{SE}}(\tilde{\beta}_j)^2 = \mathbb{E}_n[(y_i - \hat{\beta}\mathbf{X})^2 \langle \hat{\theta}_j, \hat{\mathbf{d}}_i \rangle^2]$

- For each configuration, we calculate the average coverage  $\widehat{\text{cvg}}$  for the 95% confidence intervals  $\hat{\mathcal{I}}_{\alpha,j}$  about components of  $\tilde{\beta}$  and the average interval length  $\widehat{\text{len}}$  given by

$$\widehat{\text{cvg}} = \frac{1}{p_x} \sum_{j=1}^{p_x} \frac{1}{N} \sum_{i=1}^N \mathbb{1}\{\beta_j \in \hat{\mathcal{I}}_{\alpha,j}\}, \quad \widehat{\text{len}} = \frac{1}{p_x} \sum_{j=1}^{p_x} \frac{1}{N} \sum_{i=1}^N \text{len}(\hat{\mathcal{I}}_{\alpha,j}).$$

# Configurations

We conduct simulations according the design described above for all configurations belonging to

$$\underbrace{\begin{pmatrix} (100, 250, 500) \\ (400, 800, 1000) \\ (1000, 1500, 1750) \end{pmatrix}}_{(n, p_x, p_z)} \times \underbrace{\begin{pmatrix} (3, 5) \\ (10, 15) \end{pmatrix}}_{(s_\beta, s_A)} \times \underbrace{\begin{pmatrix} .1 \\ 1 \end{pmatrix}}_b \times \underbrace{\begin{pmatrix} .1 \\ 1 \end{pmatrix}}_a \times \underbrace{\begin{pmatrix} \Sigma_z^{CS} \\ \Sigma_z^{TZ} \end{pmatrix}}_{\Sigma_z} \times \underbrace{\begin{pmatrix} .9 \\ .5 \\ .1 \end{pmatrix}}_{\sigma_{uv}},$$

with the exception of the parameter combination

$$(n, p_x, p_z, s_\beta, s_A) = (1000, 1500, 1750, 20, 25),$$

which we include in place of

$$(n, p_x, p_z, s_\beta, s_A) = (1000, 1500, 1750, 10, 15)$$

due to the magnitude of  $n$ ,  $p_x$ , and  $p_z$  in that particular combination.

# Results: Simulation results for circulant-symmetric $\Sigma_z$

$(s_\beta, s_A)$	$\sigma_{uv}$	$\widehat{\text{cvg}}(\widehat{\text{len}})$		
		$\beta$	$\beta_{s_\beta}$	$\beta_{\bar{s}_\beta}$
$(n, p_x, p_z) = (100, 250, 500)$				
(3, 5)	0.9	0.944 (0.59)	0.843 (0.60)	0.946 (0.59)
	0.5	0.944 (0.43)	0.830 (0.43)	0.946 (0.43)
	0.1	0.949 (0.21)	0.843 (0.21)	0.950 (0.21)
(10, 15)	0.9	0.909 (0.32)	0.569 (0.33)	0.923 (0.32)
	0.5	0.919 (0.27)	0.641 (0.27)	0.931 (0.27)
	0.1	0.924 (0.21)	0.727 (0.22)	0.932 (0.21)
$(n, p_x, p_z) = (400, 800, 1000)$				
(3, 5)	0.9	0.952 (0.48)	0.923 (0.48)	0.952 (0.48)
	0.5	0.950 (0.36)	0.950 (0.36)	0.950 (0.36)
	0.1	0.949 (0.16)	0.947 (0.16)	0.949 (0.16)
(10, 15)	0.9	0.949 (0.25)	0.867 (0.25)	0.950 (0.25)
	0.5	0.950 (0.19)	0.891 (0.19)	0.951 (0.19)
	0.1	0.948 (0.09)	0.880 (0.09)	0.949 (0.09)
$(n, p_x, p_z) = (1000, 1500, 1750)$				
(3, 5)	0.9	0.951 (0.39)	0.950 (0.40)	0.951 (0.39)
	0.5	0.950 (0.29)	0.963 (0.29)	0.950 (0.29)
	0.1	0.950 (0.13)	0.950 (0.13)	0.950 (0.13)
(20, 25)	0.9	0.948 (0.16)	0.917 (0.15)	0.949 (0.16)
	0.5	0.950 (0.12)	0.929 (0.12)	0.951 (0.12)
	0.1	0.958 (0.06)	0.929 (0.06)	0.959 (0.06)



# Results: Simulation results for Toeplitz $\Sigma_z$

$(s_\beta, s_A)$	$\sigma_{uv}$	$\widehat{\text{cvg}}(\widehat{\text{len}})$		
		$\beta$	$\beta_{s_\beta}$	$\beta_{\bar{s}_\beta}$
$(n, p_x, p_z) = (100, 250, 500)$				
(3, 5)	0.9	0.946 (0.58)	0.803 (0.59)	0.947 (0.58)
	0.5	0.948 (0.44)	0.780 (0.44)	0.950 (0.44)
	0.1	0.950 (0.21)	0.873 (0.22)	0.951 (0.21)
(10, 15)	0.9	0.913 (0.31)	0.585 (0.31)	0.927 (0.31)
	0.5	0.903 (0.26)	0.521 (0.26)	0.919 (0.26)
	0.1	0.901 (0.19)	0.659 (0.20)	0.911 (0.19)
$(n, p_x, p_z) = (400, 800, 1000)$				
(3, 5)	0.9	0.952 (0.48)	0.897 (0.50)	0.953 (0.48)
	0.5	0.952 (0.36)	0.927 (0.34)	0.952 (0.36)
	0.1	0.952 (0.16)	0.930 (0.16)	0.952 (0.16)
(10, 15)	0.9	0.955 (0.25)	0.861 (0.25)	0.956 (0.25)
	0.5	0.954 (0.18)	0.837 (0.18)	0.956 (0.18)
	0.1	0.956 (0.09)	0.875 (0.09)	0.957 (0.09)
$(n, p_x, p_z) = (1000, 1500, 1750)$				
(3, 5)	0.9	0.951 (0.37)	0.923 (0.35)	0.951 (0.37)
	0.5	0.951 (0.27)	0.940 (0.27)	0.951 (0.27)
	0.1	0.953 (0.12)	0.947 (0.12)	0.953 (0.12)
(20, 25)	0.9	0.952 (0.14)	0.838 (0.14)	0.954 (0.14)
	0.5	0.946 (0.10)	0.810 (0.10)	0.947 (0.10)
	0.1	0.890 (0.05)	0.678 (0.05)	0.892 (0.05)

## Limitations & next steps

- Tuning parameters: In practice, cross-validation seems to work fine, but we must develop theory for data-adaptive tuning parameter selection. We will investigate use of results in Cai et al. (2016), for instance, as well as ideas from Chichignoud et al. (2016); Bien et al. (2017).
- Real world applications: We have some ideas for conducting inference for large numbers of interaction effects in the Angrist and Krueger (1991) model that hitherto have been unexplored.
- Noise models: We are currently developing theory to handle non-Gaussian and heteroscedastic noise.
- Asymptotic efficiency: We are currently investigating the use of Jankova and van de Geer (2016) to demonstrate asymptotic efficiency of our method.
- Application to nonlinear regression: We will investigate the use of our method to approximate nonlinear regression functions, for instance under the “approximate sparsity” regime of Belloni et al. (2011, 2012).

# Acknowledgments

- Johannes Lederer & Jing Tao
- Joseph Salmon & Mohamed Hebiri
- Donghui Mai
- UW Royalty Research Fund Grant
- Amazon Cloud Credits for Research Grant

## References I

- Amemiya, T. (1974). The nonlinear two-stage least-squares estimator. *Journal of Econometrics*, 2:105–110.
- Amemiya, T. (1977). The maximum likelihood and the nonlinear three-stage least squares estimator in the general nonlinear simultaneous equation model. *Econometrica*, 45(4):955–968.
- Amemiya, T. (1985). *Advanced Econometrics*. Harvard University Press.
- Anderson, T. W. (2005). Origins of the limited information maximum likelihood estimator and two-stage least squares estimators. *Journal of Econometrics*, 127:1–16.
- Anderson, T. W. and Rubin, H. (1950). The asymptotic properties of estimates of the parameters of a single equation in a complete system of stochastic equations. *The Annals of Mathematical Statistics*, 21(4):570–582.
- Angrist, J. D. and Krueger, A. B. (1991). Does compulsory school attendance affect schooling and earnings? *The Quarterly Journal of Economics*.

## References II

- Basmann, R. L. (1957). A generalized classical method of linear estimation of coefficients in a structural equation. *Econometrica*, 25(1):77–83.
- 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(6):2369–2429.
- Belloni, A., Chernozhukov, V., and Hansen, C. (2011). Inference for high-dimensional sparse econometric models. *Advances in Economics and Econometrics*.
- Bickel, P. J., Ritov, Y., and Tsybakov, A. B. (2009). Simultaneous analysis of lasso and dantzig selector. *Ann. Statist.*, 37(4):1705–1732.
- Bien, J., Gaynanova, I., Lederer, J., and Müller, C. L. (2017). Non-convex Global Minimization and False Discovery Rate Control for the TREX. *J. Comput. Graph. Statist.*
- Bühlmann, P. and van de Geer, S. (2011). *Statistics for High-Dimensional Data*. Springer.

## References III

- Cai, T., Liu, W., and Luo, X. (2011). A constrained  $\ell_1$  minimization approach to sparse precision matrix estimation. *J. Amer. Statist. Assoc.*, 106(464):594–607.
- Cai, T. T., Liu, W., and Zhou, H. H. (2016). Estimating sparse precision matrix: Optimal rates of convergence and adaptive estimation. *Ann. Statist.*, 44(2):455–488.
- Chamberlain, G. (1987). Asymptotic efficiency in estimation with conditional moment restrictions. *Journal of Econometrics*, 34(3):305–334.
- Chichignoud, M., Lederer, J., and Wainwright, M. J. (2016). A practical scheme and fast algorithm to tune the lasso with optimality guarantees. *J. Mach. Learn. Res.*, 17:1–17.
- Hansen, L. P. (1982). Large sample properties of generalized method of moments estimators. *Econometrica*, 50(4):1029–1054.
- Jankova, J. and van de Geer, S. (2016). Semi-parametric efficiency bounds for high-dimensional models. *arXiv:1601.00815*.

## References IV

- Javanmard, A. and Montanari, A. (2014). Confidence intervals and hypothesis testing for high-dimensional regression. *J. Mach. Learn. Res.*, 15:2869–2909.
- Knight, K. and Fu, W. (2000). Asymptotics for lasso-type estimators. *Ann. Statist.*, 28(5):1356–1378.
- Meinshausen, N. and Bhlmann, P. (2006). High-dimensional graphs and variable selection with the lasso. *Ann. Statist.*, 34(3):1436–1462.
- Newey, W. K. (1990). Efficient instrumental variables estimation of nonlinear models. *Econometrica*, 58(4):809–837.
- Neykov, M., Ning, Y., Liu, J. S., and Liu., H. (2015). A unified theory of confidence regions and testing for high dimensional estimating equations. *arXiv:1510.08986*.
- Ning, Y. and Liu, H. (2014). A general theory of hypothesis tests and confidence regions for sparse high dimensional models. *arXiv:1412.8765*.

## References V

- Pötscher, B. M. and Leeb, H. (2009). On the distribution of penalized maximum likelihood estimators: The LASSO, SCAD, and thresholding. *J. Multivar. Anal.*, 100(9):2065–2082.
- Sargan, J. D. (1958). The estimation of economic relationships using instrumental variables. *Econometrica*.
- Stock, J. H. and Trebbi, F. (2003). Who invented instrumental variable regression? *Journal of Economic Perspectives*, 17(3):177–194.
- Theil, H. (1953). Repeated least-squares applied to complete equation systems. *Centraal Planbureau Memorandum*.
- van de Geer, S. and Bühlmann, P. (2009). On the conditions used to prove oracle results for the lasso. *Electron. J. Stat.*, 3:1360–1392.
- van de Geer, S., Bühlmann, P., Ritov, Y., and Dezeure, R. (2014). On asymptotically optimal confidence regions and tests for high-dimensional models. *Ann. Statist.*, 42(3):1166–1202.
- Wright, P. G. (1928). *The tariff on animal and vegetable oils*. New York, The Macmillan Company.



## References VI

Zhang, C.-H. and Zhang, S. S. (2014). Confidence intervals for low dimensional parameters in high dimensional linear models. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, 76(1):217–242.