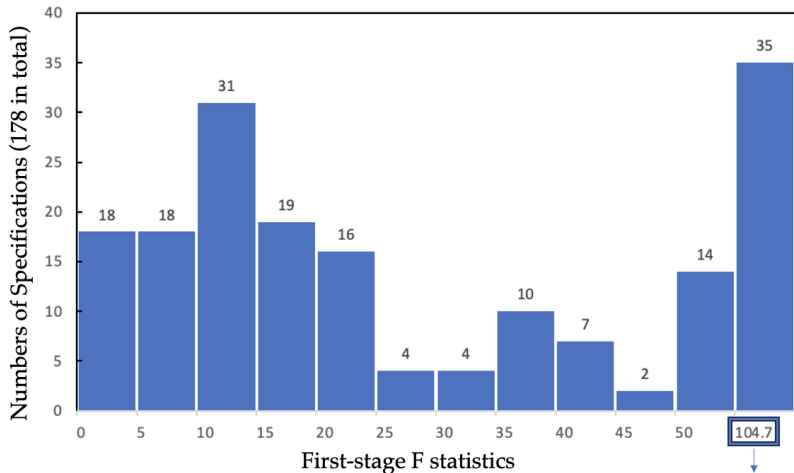


# Identification-robust inference for the LATE with high-dimensional covariates

Yukun Ma  
Vanderbilt University

January 17, 2024

## Motivation: Weak IV in Empirical Practice



*new threshold by Lee et al. (2022)  
under the conventional critical value 1.96*

Figure: *American Economic Review* 2018-2022

► [heterscadesticity](#)

## Motivation: Availability of Large Datasets

- Big datasets are becoming increasingly available nowadays,
  - ▶ high data volume, 15.5 ZB in 2015, 97 ZB in 2022 (representing a 525% increase); (1 ZB =  $10^{12}$  GB)
  - ▶ [high-dimensional controls](#), allowing for more flexible functional forms, e.g., polynomial terms and interaction effects.

# Abstract

- New inference procedure for local average treatment effect (**LATE**) when
  - ▶ identification may be **weak** (e.g. few compliers),
  - ▶ model incorporate **high-dimensional covariates** (e.g. many controls).
- The proposed test statistic has **uniformly correct asymptotic size**,
  - ▶ inversion of the proposed test statistic for inference on LATE.
- Revisit 2 empirical studies:
  - ▶ Hornung (2015, JEEA) and Ambrus et al. (2020, AER),
  - ↪ in both cases, the proposed method is more efficient, yielding narrower confidence regions -whereas competitors often report larger confidence intervals.

# Table of Contents

- 1 Introduction
- 2 Overview
- 3 Theory
- 4 Simulation
- 5 Applications
- 6 Takeaways

# Table of Contents

- 1 Introduction
- 2 Overview
- 3 Theory
- 4 Simulation
- 5 Applications
- 6 Takeaways

## LATE

- LATE: The effect of a treatment on compliers who adhere to the treatment assigned to their sample group.
- Assume we have  $N$  observations
  - ▶  $Y_i$  : outcome of interest for unit  $i$ .
  - ▶  $D_i \in \{0, 1\}$  : receipt of treatment.
  - ▶  $Z_i \in \{0, 1\}$  : offer of treatment.

- Imbens and Angrist (1994) propose

$$\text{LATE} = \frac{\mathbf{E}_P[Y|Z = 1] - \mathbf{E}_P[Y|Z = 0]}{\mathbf{E}_P[D|Z = 1] - \mathbf{E}_P[D|Z = 0]} = \frac{ITT}{FS}.$$

- Incorporate covariates into LATE estimation, e.g. Abadie (2003)
  - ▶  $\mathbf{X}_i$  :  $p$ -dimensional covariates.
- Weak identification in LATE:  $FS \rightarrow 0$ .

# LATE

- LATE: The effect of a treatment on compliers who adhere to the treatment assigned to their sample group.
- Assume we have  $N$  observations
  - ▶  $Y_i$  : outcome of interest for unit  $i$ .
  - ▶  $D_i \in \{0, 1\}$  : receipt of treatment.
  - ▶  $Z_i \in \{0, 1\}$  : offer of treatment.

- Imbens and Angrist (1994) propose

$$\theta := \text{LATE} = \frac{\mathbf{E}_P[Y|Z = 1] - \mathbf{E}_P[Y|Z = 0]}{\mathbf{E}_P[D|Z = 1] - \mathbf{E}_P[D|Z = 0]} = \frac{ITT}{FS} := \frac{\delta}{\pi}.$$

- Incorporate covariates into LATE estimation, e.g. Abadie (2003)
  - ▶  $\mathbf{X}_i$  :  $p$ -dimensional covariates.
- Weak identification in LATE:  $\pi \rightarrow 0$



## Weak identification

- Issue: When instruments  $Z$  are weakly correlated with endogenous regressors  $D$ , conventional methods for IV estimation and inference become unreliable.

$$\theta = \frac{\delta}{\pi},$$

when  $\hat{\pi}$  is close to zero,  $\hat{\theta}$  is **highly nonlinear** in  $\hat{\pi}$

- Trick: Fieller-type transformation & test inversion.  
Given  $H_0 : \theta = \theta_0$ , we have  $\delta - \theta_0\pi = 0$ . The Anderson-Rubin (AR) test,

$$AR(\theta) = (\delta - \theta\pi)' \Omega(\theta)^{-1} (\delta - \theta\pi),$$

follows a  $\chi^2$  distribution under  $H_0$ .

## Weak identification

- Issue: When instruments  $Z$  are weakly correlated with endogenous regressors  $D$ , conventional methods for IV estimation and inference become unreliable.

$$\theta = \frac{\delta}{\pi},$$

when  $\hat{\pi}$  is close to zero,  $\hat{\theta}$  is **highly nonlinear** in  $\hat{\pi}$

- Trick: Fieller-type transformation & test inversion.  
Given  $H_0 : \theta = \theta_0$ , we have  $\delta - \theta_0\pi = 0$ . The Anderson-Rubin (AR) test,

$$AR(\theta) = (\delta - \theta\pi)' \Omega(\theta)^{-1} (\delta - \theta\pi),$$

follows a  $\chi^2$  distribution under  $H_0$ .

## Relations to the Literature: Weak Identification

- Inference procedures depend on the observed process only through its value, and potentially derivative, at the point  $\theta_0$ : Anderson–Rubin statistic, Stock and Wright (2000), Kleibergen (2005).
- Methods depend on the full path of the observed process: Moreira (2003) and Andrews and Mikusheva (2016).

	Low-dimensional Model
Strong Identification	z-test
Weak Identification	Stock and Wright (2000), Kleibergen (2005), Andrews and Mikusheva (2016)

## Relations to the Literature: Weak Identification

- Inference procedures depend on the observed process only through its value, and potentially derivative, at the point  $\theta_0$ : Anderson–Rubin statistic, Stock and Wright (2000), Kleibergen (2005).
- Methods depend on the full path of the observed process: Moreira (2003) and Andrews and Mikusheva (2016).

	Low-dimensional Model
Strong Identification	z-test
Weak Identification	Stock and Wright (2000), Kleibergen (2005), Andrews and Mikusheva (2016)

↪ **Limitation:** None of them considers high-dimensional model (model with many covariates).

## Relations to the Literature: ML Methods

In high-dimensional models, ML methods are commonly employed for model selection:

- Chernozhukov et al. (2018) introduce the double/debiased machine learning (DML) method, a combination of the Neyman orthogonality condition and cross-fitting method.
- Belloni et al. (2017) present an efficient estimator and confidence bands for the LATE with nonparametric/high-dimensional components.

	Low-dimensional Model	High-dimensional Model
Strong Identification	z-test	ML methods
Weak Identification	Stock and Wright (2000), Kleibergen (2005), Andrews and Mikusheva(2016)	

## Relations to the Literature: ML Methods

In high-dimensional models, ML methods are commonly employed for model selection:

- Chernozhukov et al. (2018) introduce the double/debiased machine learning (DML) method, a combination of the Neyman orthogonality condition and cross-fitting method.
- Belloni et al. (2017) present an efficient estimator and confidence bands for the LATE with nonparametric/high-dimensional components.

	Low-dimensional Model	High-dimensional Model
Strong Identification	z-test	ML methods
Weak Identification	Stock and Wright (2000), Kleibergen (2005), Andrews and Mikusheva(2016)	

# Comparison of the Literature

## Weak Identification lit

- identification-robust test statistics
- use traditional methods to handle  $\mathbf{X}_i$

## Machine Learning lit

- normal z-test
- use ML to handle  $\mathbf{X}_i$

## Drawbacks

- overfitting
- multicollinearity
- cannot perform well under weakly identified scenarios

↪ My proposed method takes advantage from both literature:

- identification-robust test statistics
- use ML to handle the high-dimensional  $\mathbf{X}_i$

# Comparison of the Literature

## Weak Identification lit

- identification-robust test statistics
- use traditional methods to handle  $\mathbf{X}_i$

## Drawbacks

- overfitting
- multicollinearity

## Machine Learning lit

- normal z-test
- use ML to handle  $\mathbf{X}_i$

- cannot perform well under weakly identified scenarios

↪ My proposed method takes advantage from both literature:

- identification-robust test statistics
- use ML to handle the high-dimensional  $\mathbf{X}_i$



# Contributions

	Low-dimensional Model	High-dimensional Model
Strong Identification	z-test	ML methods
Weak Identification	Stock and Wright (2000), Kleibergen (2005), Andrews and Mikusheva(2016)	My proposed method

**Technical Contribution:** Instead of proposing a consistent estimator for  $\theta$ , I present an empirical process along with its uniformly consistent estimator,

↪ the proposed test statistic is shown to be uniformly size-correct.

# Table of Contents

- 1 Introduction
- 2 Overview**
- 3 Theory
- 4 Simulation
- 5 Applications
- 6 Takeaways

## Setup

- Model the random vector  $\mathbf{W} = (Y, D, Z, \mathbf{X}')'$  as follows,

$$\text{First stage} \quad D = m_0(Z, \mathbf{X}) + v, \quad \mathbf{E}_P[v|Z, \mathbf{X}] = 0$$

$$\text{Reduced form} \quad Y = g_0(Z, \mathbf{X}) + u, \quad \mathbf{E}_P[u|Z, \mathbf{X}] = 0$$

$$\text{Propensity score} \quad Z = p_0(\mathbf{X}) + e, \quad \mathbf{E}_P[e|\mathbf{X}] = 0$$

- $m_0, g_0, p_0$ : no need to impose any parametric assumptions by [Blandhol et al. \(2022\)](#).
- The LATE framework proposed by [Tan \(2006\)](#) is given by

$$\theta = \frac{\mathbf{E}_P[g(1, \mathbf{X}) - g(0, \mathbf{X}) + \frac{Z}{p(\mathbf{X})} (Y - g(1, \mathbf{X})) - \frac{1-Z}{1-p(\mathbf{X})} (Y - g(0, \mathbf{X}))]}{\mathbf{E}_P[m(1, \mathbf{X}) - m(0, \mathbf{X}) + \frac{Z}{p(\mathbf{X})} (D - m(1, \mathbf{X})) - \frac{1-Z}{1-p(\mathbf{X})} (D - m(0, \mathbf{X}))]}$$

## Setup

- Model the random vector  $\mathbf{W} = (Y, D, Z, \mathbf{X}')'$  as follows,

$$\text{First stage} \quad D = m_0(Z, \mathbf{X}) + v, \quad \mathbf{E}_P[v|Z, \mathbf{X}] = 0$$

$$\text{Reduced form} \quad Y = g_0(Z, \mathbf{X}) + u, \quad \mathbf{E}_P[u|Z, \mathbf{X}] = 0$$

$$\text{Propensity score} \quad Z = p_0(\mathbf{X}) + e, \quad \mathbf{E}_P[e|\mathbf{X}] = 0$$

- The LATE framework proposed by [Tan \(2006\)](#) is given by

$$\theta = \frac{\mathbf{E}_P[g(1, \mathbf{X}) - g(0, \mathbf{X}) + \frac{Z}{p(\mathbf{X})}(Y - g(1, \mathbf{X})) - \frac{1-Z}{1-p(\mathbf{X})}(Y - g(0, \mathbf{X}))]}{\mathbf{E}_P[m(1, \mathbf{X}) - m(0, \mathbf{X}) + \frac{Z}{p(\mathbf{X})}(D - m(1, \mathbf{X})) - \frac{1-Z}{1-p(\mathbf{X})}(D - m(0, \mathbf{X}))]} = \frac{\mathbf{E}_P[\mathbf{a}]}{\mathbf{E}_P[\mathbf{b}]}$$

$$\triangleright \mathbf{a} := g(1, \mathbf{X}) - g(0, \mathbf{X}) + \frac{Z}{p(\mathbf{X})}(Y - g(1, \mathbf{X})) - \frac{1-Z}{1-p(\mathbf{X})}(Y - g(0, \mathbf{X}))$$

$$\triangleright \mathbf{b} := m(1, \mathbf{X}) - m(0, \mathbf{X}) + \frac{Z}{p(\mathbf{X})}(D - m(1, \mathbf{X})) - \frac{1-Z}{1-p(\mathbf{X})}(D - m(0, \mathbf{X}))$$

# Score Function

- Consider a function

$$\psi(\mathbf{W}; \theta, g, m, p) = \overbrace{g(1, X) - g(0, X) + \frac{Z(Y - g(1, X)) - (1 - Z)(Y - g(0, X))}{p(X) - (1 - p(X))}}^a - \theta \times \underbrace{\left( m(1, X) - m(0, X) + \frac{Z(D - m(1, X)) - (1 - Z)(D - m(0, X))}{p(X) - (1 - p(X))} \right)}_b,$$

with

- ▶ target parameter  $\theta \in \Theta \subset \mathbb{R}$  is the LATE.
  - ▶ nuisance parameter  $\eta = (g, m, p) \in T$  for a convex<sup>1</sup> set  $T$ .
  - ▶  $\psi$  is a score function.
- Two-stage procedure:
    - Stage 1 estimating nuisance parameter  $\eta$ ,
    - Stage 2 making inference for the target parameter  $\theta$ .

---

<sup>1</sup>To ensure that  $\psi(\mathbf{W}; \theta_0, \eta_0 + r(\eta - \eta_0))$  is well defined for all  $r \in [0, 1)$ .

# Score Function

- Consider a function

$$\psi(\mathbf{W}; \theta, g, m, p) = \overbrace{g(1, X) - g(0, X) + \frac{Z(Y - g(1, X))}{p(X)} - \frac{(1 - Z)(Y - g(0, X))}{1 - p(X)}}^a - \theta \times \underbrace{\left( m(1, X) - m(0, X) + \frac{Z(D - m(1, X))}{p(X)} - \frac{(1 - Z)(D - m(0, X))}{1 - p(X)} \right)}_b,$$

with

- ▶ target parameter  $\theta \in \Theta \subset \mathbb{R}$  is the LATE.
  - ▶ nuisance parameter  $\eta = (g, m, p) \in T$  for a convex<sup>1</sup> set  $T$ .
  - ▶  $\psi$  is a score function.
- Two-stage procedure:
    - Stage 1 estimating nuisance parameter  $\eta$ ,
    - Stage 2 making inference for the target parameter  $\theta$ .

---

<sup>1</sup>To ensure that  $\psi(\mathbf{W}; \theta_0, \eta_0 + r(\eta - \eta_0))$  is well defined for all  $r \in [0, 1)$ .

# Nuisance Parameters

- Specify the nuisance parameters  $\eta = (g, m, p)$  as follows,

$$g(Z, X) = \mathbf{E}_P[Y|Z, X] = Z\beta_{21} + X'\beta_{22} \quad \text{Reduced form}$$

$$m(Z, X) = \mathbf{E}_P[D|Z, X] = \Lambda(Z\beta_{11} + X'\beta_{12}) \quad \text{First stage}$$

$$p(X) = \mathbf{E}_P[Z|X] = \Lambda(X'\gamma) \quad \text{Propensity score}$$

- ▶ The logistic CDF  $\Lambda(t) = \frac{\exp(t)}{1+\exp(t)}$  for all  $t \in \mathbb{R}$
- ▶ In this example, the nuisance parameters  $\eta = (\beta_{11}, \beta_{12}, \beta_{21}, \beta_{22}, \gamma)$ .

## Properties of the Score $\psi$

- Moment condition:

$$\mathbf{E}_P[\underbrace{\psi(W_i; \theta_0, \eta_0)}_{a-\theta_0 \times b}] = 0.$$

- Neyman orthogonality condition:

- ▶ Path-wise (or Gateaux) derivative map  $D_r$

$$D_r[\eta - \eta_0] := \partial_r \{ \mathbf{E}_P[\psi(W; \theta_0, \eta_0 + r(\eta - \eta_0))] \} \text{ for } \eta \in T.$$

- ▶ The Neyman orthogonality condition holds at  $(\theta_0, \eta_0)$  if

$$D_0[\eta - \eta_0] = \partial_\eta \mathbf{E}_P \psi(W; \theta_0, \eta_0)[\eta - \eta_0] = 0$$

holds for all  $\eta \in \mathcal{T}_N$  for a nuisance realization set  $\mathcal{T}_N \subset T$ .

↪ The score function  $\psi$  is an **AR-type Neyman orthogonal score**.



# Algorithm Breakdown

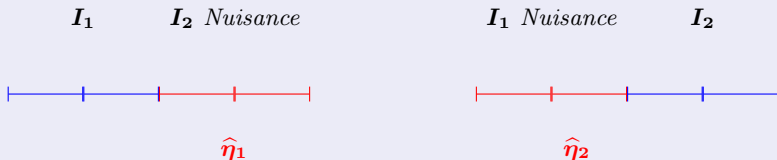
- estimate nuisance parameter  $\eta$   
Step 1-2: data splitting/ cross-validation
  
- make an inference for the target parameter  $\theta$  based on  $\hat{\eta}$   
Step 3-6: inversion of the condition QLR test statistics

# Estimate Nuisance Parameter $\eta$

Step 1: Randomly split the sample  $\{1, \dots, N\}$  into  $K$  folds  $\{I_1, \dots, I_K\}$ .

Step 2: For each  $k \in \{1, \dots, K\}$ , obtain  $\hat{\eta}_k$  by ML methods using only the subsample of those observations with indices  $i \in \{1, \dots, N\} \setminus I_k$ :

An illustration of  $K=2$ -fold cross-fitting.



## Estimate Nuisance Parameter, Step 2

(2.1)  $(\hat{\beta}_{21,k}, \hat{\beta}_{22,k})$  in reduce form: run ML (e.g., lasso) OLS regression to estimate,

$$(\hat{\beta}_{21,k}, \hat{\beta}_{22,k}) \in \arg \min_{\beta_{21}, \beta_{22}} \mathbb{E}_{I_k^c} [(Y_i - Z_i \beta_{21} - X_i' \beta_{22})^2] + \frac{\lambda_3}{|I_k^c|} \|(\beta_{21}, \beta_{22})\|_1,$$

(2.2)  $(\hat{\beta}_{11,k}, \hat{\beta}_{12,k})$  in first stage: run ML (e.g., lasso) logistic regression to estimate,

$$(\hat{\beta}_{11,k}, \hat{\beta}_{12,k}) \in \arg \min_{\beta_{11}, \beta_{12}} \left\{ \mathbb{E}_{I_k^c} [L_1(\mathbf{W}_i; \beta_{11}, \beta_{12})] + \frac{\lambda_1}{|I_k^c|} \|(\beta_{11}, \beta_{12})\|_1 \right\},$$

(2.3)  $\hat{\gamma}_k$  in the propensity score: run ML (e.g., lasso) logistic regression to estimate,

$$\hat{\gamma}_k \in \arg \min_{\gamma} \left\{ \mathbb{E}_{I_k^c} [L_2(\mathbf{W}_i; \gamma)] + \frac{\lambda_2}{|I_k^c|} \|\gamma\|_1 \right\},$$

▶  $\lambda_1, \lambda_2, \lambda_3$  : ▶ penalty parameter

▶  $L_1(\mathbf{W}_i; \beta_{11}, \beta_{12}) = D_i(Z_i \beta_{11} + X_i' \beta_{12}) - \log(1 + \exp(Z_i \beta_{11} + X_i' \beta_{12}))$ ,

▶  $L_2(\mathbf{W}_i; \gamma) = Z_i X_i' \gamma - \log(1 + \exp(X_i' \gamma))$ .

### Step 3

Step 3: Compute  $\hat{q}_N(\theta)$  and  $\hat{\Omega}(\theta_1, \theta_2)$  for later use,

$$\hat{q}_N(\theta) = \frac{1}{\sqrt{N}} \sum_{k=1}^K \sum_{i \in I_k} \psi(W_i; \theta, \hat{\eta}_k),$$
$$\hat{\Omega}(\theta_1, \theta_2) = \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \psi(W_i; \theta_1, \hat{\eta}_k) \psi(W_i; \theta_2, \hat{\eta}_k) \\ - \frac{1}{N^2} \sum_{k=1}^K \sum_{k'=1}^K \sum_{i \in I_k, i' \in I_{k'}} \psi(W_i; \theta_1, \hat{\eta}_k) \psi(W_{i'}; \theta_2, \hat{\eta}_{k'}).$$

An illustration of K=2-fold cross-fitting.

$I_1$  Score  $I_2$  Nuisance



$$\sum_{i \in I_1} \psi(W_i; \theta, \hat{\eta}_1)$$

$I_1$  Nuisance  $I_2$  Score



$$\sum_{i \in I_2} \psi(W_i; \theta, \hat{\eta}_2)$$

## Make Inference for Target Parameter $\theta$

**Step 4:** Take independent draws  $\xi \sim N(\mathbf{0}, \widehat{\Omega}(\theta_0, \theta_0))$  and calculate  $\mathbf{R} = \mathbf{R}(\xi, \mathbf{h}_N, \widehat{\Omega})$ , where ▶ null hypothesis  $H_0$

$$\mathbf{R}(\xi, \mathbf{h}_N, \widehat{\Omega}) = \xi^2 \widehat{\Omega}(\theta_0, \theta_0)^{-1} - \inf_{\theta} (V(\theta)\xi + \mathbf{h}_N)^2 \widehat{\Omega}(\theta, \theta)^{-1},$$

- ▶  $V(\theta) = \widehat{\Omega}(\theta, \theta_0) \widehat{\Omega}(\theta_0, \theta_0)^{-1}$ ,
- ▶  $\mathbf{h}_N(\theta) = \widehat{q}_N(\theta) - \widehat{\Omega}(\theta, \theta_0) \widehat{\Omega}(\theta_0, \theta_0)^{-1} \widehat{q}_N(\theta_0)$ .

**Step 5:** Calculate the conditional critical value  $c_{\alpha}(\tilde{\mathbf{h}})$  as

$$c_{\alpha}(\tilde{\mathbf{h}}) = \min\{c : P(\mathbf{R}(\xi, \mathbf{h}_N, \widehat{\Omega}) > c) \leq \alpha\}.$$

**Step 6:** Reject the null hypothesis  $H_0 : \mathbf{S}_N \in \mathcal{S}_0$  when  $\mathbf{R}(\xi, \mathbf{h}_N, \widehat{\Omega}) \geq c_{\alpha}(\mathbf{h}_N)$ , report the  $(1 - \alpha)$  confidence interval:  $CI_{\alpha} = \{\theta : \mathbf{R}(\xi, \mathbf{h}_N, \widehat{\Omega}) \leq c_{\alpha}(\mathbf{h}_N)\}$ .

# Table of Contents

- 1 Introduction
- 2 Overview
- 3 Theory**
- 4 Simulation
- 5 Applications
- 6 Takeaways

## Notations

- ▶ Let  $c > 0$ ,  $c_0 \geq 0$ ,  $c_1 \geq 0$ ,  $C_1 > 0$  be finite constants, and  $\mathbf{a}_N = \mathbf{p} \vee N$ .
- ▶ Let  $\{\Delta_N\}_{N \geq 1}$ ,  $\{\delta_N\}_{N \geq 1}$  (estimation errors) be sequences of positive constants that converges to zero such that  $\delta_N \geq N^{-1/2}$ .
- ▶ Let  $\|\delta\|_0$  represent the number of non-zero components of  $\delta$ .
- ▶ Let  $P \in \mathcal{P}_N$  be the probability law of  $\{W_i\}_{i=1}^N$ .
- ▶ Let  $\mathcal{P}_0$  be family of distribution consistent with the null.
- ▶ We use  $\mathbf{a} \lesssim \mathbf{b}$  to denote  $\mathbf{a} \leq \mathbf{cb}$  for some  $\mathbf{c} > \mathbf{0}$  that does not depends on  $N$ .
- ▶ The sequence  $\{M_N\}_{N \geq 1}$  be a set of positive constants such that  $M_N \geq (\mathbf{E}_P[(Z_i \vee \|X_i\|_\infty)^{2q}])^{1/2q}$ .

# Assumption: Regularity Conditions for the LATE

For  $P \in \mathcal{P}_N$ , the following conditions hold.

(i) The equations are satisfied with binary variables  $D$  and  $Z$ .

$$\left. \begin{aligned} D &= m_0(Z, X) + v, & \mathbf{E}_P[v|Z, X] &= \mathbf{0} \\ Y &= g_0(Z, X) + u, & \mathbf{E}_P[u|Z, X] &= \mathbf{0} \end{aligned} \right\} \rightarrow (Y, D) \perp\!\!\!\perp Z|X$$
$$Z = p_0(X) + e, \quad \mathbf{E}_P[e|X] = \mathbf{0}.$$

(ii) For some  $\varepsilon > 0$ ,  $\varepsilon \leq P(Z = 1|X) \leq 1 - \varepsilon$  almost surely.

(iii)  $\Theta$  is compact.

(iv)  $\mathbf{E}_P[D|Z = 1] \geq \mathbf{E}_P[D|Z = 0]$ . ▶ Assumption Comparison

(v)  $\|u\|_{P,2} \geq c_0$ , and  $\|\mathbf{E}_P[u^2|X]\|_{P,\infty} \leq c_1$ .

(vi)  $\|Y\|_{P,q} \leq c_1$ .



# Assumption: Nuisance Parameter Estimators

- Sparse eigenvalue conditions: with probability  $1 - o(1)$ , for some  $l_N \rightarrow \infty$  slow enough, we have

$$\mathbf{1} \lesssim \phi_{\min}(l_N s_N) \leq \phi_{\max}(l_N s_N) \lesssim \mathbf{1}.$$

▶ sparse eigenvalue

- Sparsity:  $\|\beta_{12}^0\|_0 + \|\beta_{22}^0\|_0 + \|\gamma^0\|_0 \leq s_N$ .
- Parameters:  $\|\beta_{12}^0\| + \|\beta_{22}^0\| + \|\gamma^0\| \leq C_1$ .
- Covariates: for  $q > 4$ ,
  - ▶  $\inf_{\|\xi\|=1} \mathbf{E}_P[\left((Z_i, X_i')\xi\right)^2] \geq c$ .
  - ▶  $\sup_{\|\xi\|=1} \mathbf{E}_P[\left((Z_i, X_i')\xi\right)^2] \leq C_1$ .
  - ▶  $N^{-1/2+2/q} M_N^2 s_N \log^2 a_N \leq \Delta_N$ .

## Main Result

Propose an empirical process

$$\mathbb{G}_N(\cdot) = \underbrace{\frac{1}{\sqrt{N}} \sum_{i=1}^N (\psi(W_i; \cdot, \eta_0) - \mathbf{E}_P[\psi(W; \cdot, \eta_0)])}_{q_N(\cdot)},$$

and its estimator as

$$\widehat{\mathbb{G}}_N(\boldsymbol{\theta}) = \underbrace{\sqrt{N} \left( \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \psi(W_i; \boldsymbol{\theta}, \widehat{\eta}_k) - \mathbf{E}_P[\psi(W_i; \boldsymbol{\theta}, \widehat{\eta}_k)] \right)}_{\widehat{q}_N(\boldsymbol{\theta})}.$$

### Theorem 1

Suppose that the above assumptions are satisfied. Under the null, we have

$$\widehat{\mathbb{G}}_N(\boldsymbol{\theta}) = \mathbb{G}_N(\boldsymbol{\theta}) + O_P(N^{-1/2} \delta_N).$$

The process  $\widehat{\mathbb{G}}_N(\cdot)$  weakly converges to a centered Gaussian process  $\mathbb{G}(\cdot)$  for all  $P \in \mathcal{P}_0$  as  $N \rightarrow \infty$  with covariance function  $\Omega(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) = \mathbf{E}_P[(\psi(W; \boldsymbol{\theta}_1, \eta_0) - \mathbf{E}_P[\psi(W; \boldsymbol{\theta}_1, \eta_0)]) (\psi(W; \boldsymbol{\theta}_2, \eta_0) - \mathbf{E}_P[\psi(W; \boldsymbol{\theta}_2, \eta_0)])]$ .

# Variance Estimation

## Theorem 2

Under the same set of assumptions as above, the covariance function  $\Omega(\theta_1, \theta_2)$  can be consistently estimated uniformly for all  $\mathbf{P} \in \mathcal{P}_0$  by

$$\begin{aligned}\widehat{\Omega}(\theta_1, \theta_2) &= \frac{1}{N} \sum_{k=1}^K \sum_{i \in I_k} \psi(W_i; \theta_1, \widehat{\eta}_k) \psi(W_i; \theta_2, \widehat{\eta}_k) \\ &\quad - \frac{1}{N^2} \sum_{k, k'=1}^K \sum_{i \in I_k, i' \in I_{k'}} \psi(W_i; \theta_1, \widehat{\eta}_k) \psi(W_{i'}; \theta_2, \widehat{\eta}_{k'})\end{aligned}$$

and for any  $\varepsilon > 0$ ,

$$\lim_{N \rightarrow \infty} \sup_{\mathbf{P} \in \mathcal{P}_0} P \left\{ \sup_{\theta_1, \theta_2} \|\widehat{\Omega}(\theta_1, \theta_2) - \Omega(\theta_1, \theta_2)\| > \varepsilon \right\} = 0.$$

## How It Works

- Weak convergence over  $\Theta_I$ :
  - 1 the convergence of the finite dimensional distribution of  $\widehat{\mathbb{G}}_N(\boldsymbol{\theta})$  for  $\boldsymbol{\theta} \in \Theta_I$ .
  - 2 the stochastic equicontinuity of  $\widehat{\mathbb{G}}_N(\boldsymbol{\theta})$  over  $\Theta_I$ :

$$\lim_{\delta \rightarrow 0} \limsup_{N \rightarrow \infty} P \left( \sup_{|\boldsymbol{\theta}_1 - \boldsymbol{\theta}_2| \leq \delta} |\widehat{\mathbb{G}}_N(\boldsymbol{\theta}_1) - \widehat{\mathbb{G}}_N(\boldsymbol{\theta}_2)| > \varepsilon_1 \right) = 0,$$

for any  $\varepsilon_1 > 0$  and  $\boldsymbol{\theta}_1, \boldsymbol{\theta}_2 \in \Theta_I$ .

- 3 the boundedness of  $\Theta_I$ .
- The equivalence between testing  $\boldsymbol{\theta} \in \Theta_I$  and  $P \in \mathcal{P}_0$ .

$\hookrightarrow$  Uniformly consistent results for  $\widehat{\mathbb{G}}_N(\cdot)$ .

# Size Control

## Theorem 3

*Under the same set of assumptions above, the test that rejects the null hypothesis  $H_0 : \mathbf{S}_N \in \mathcal{S}_0$  when  $R(\mathbf{q}_N(\boldsymbol{\theta}_0), \mathbf{h}_N, \boldsymbol{\Omega})$  exceeds the  $(1 - \alpha)$  quantile  $c_\alpha(\mathbf{h}_N)$  of its conditional distribution given  $\mathbf{h}_N(\cdot)$  has uniformly correct asymptotic size.*

*Under the null, we have*

$$\lim_{N \rightarrow \infty} \sup_{P \in \mathcal{P}_0} P(R(\hat{\mathbf{q}}_N(\boldsymbol{\theta}_0), \mathbf{h}_N, \hat{\boldsymbol{\Omega}}) > c_\alpha(\mathbf{h}_N)) = \alpha.$$

# Table of Contents

- 1 Introduction
- 2 Overview
- 3 Theory
- 4 Simulation**
- 5 Applications
- 6 Takeaways

# Simulation Setup

- Primitive random vector  $\mathbf{X}'_i$  is constructed by

$$\mathbf{X}_i \sim \mathcal{N} \left( \mathbf{0}, \begin{pmatrix} U^0 & U^1 & \dots & U^{\dim(\mathbf{X})-2} & U^{\dim(\mathbf{X})-1} \\ U^1 & U^0 & \dots & U^{\dim(\mathbf{X})-3} & U^{\dim(\mathbf{X})-2} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ U^{\dim(\mathbf{X})-2} & U^{\dim(\mathbf{X})-3} & \dots & U^0 & U^1 \\ U^{\dim(\mathbf{X})-1} & U^{\dim(\mathbf{X})-2} & \dots & U^1 & U^0 \end{pmatrix} \right)$$

with  $U = 0.5$ .

- Consider  $N = 500$ ,  $\dim(\mathbf{X}) = 5$ , 200, 400, and 600.  
high-dimensional controls

## Simulation Setup, Continued

- Consider the threshold crossing model:

- ▶ The latent tendency to receive treatment  $\delta_i \sim \mathcal{N}(0, 1)$ .
- ▶ The treatment assignment is given by  $Z_i = \mathbb{1}\{\delta_i \geq 0\}$ .
- ▶ The potential treatment indicators  $D_i(Z_i)$  are given by

$$D_i(0) = \mathbb{1}\{\Phi(\delta_i) < P_{AT}\}, \quad D_i(1) = \mathbb{1}\{\Phi(\delta_i) < 1 - P_{NT}\},$$

with  $\Phi(\cdot)$  denotes the CDF of a standard normal distribution.

- ▶ The target parameter is set to  $\theta_0 = 1$ .
  - ▶ The outcome  $Y_i = D_i + X_i + \varepsilon_i$  with  $\varepsilon_i \sim \mathcal{N}(0, 1)$ .
- Scenarios:
    - ▶ Strongly identified case:  $(P_{AT}, P_{NT}) = (0.25, 0.25) \rightarrow P_C = 0.5$
    - ▶ Weakly identified case:  $(P_{AT}, P_{NT}) = (0.45, 0.45) \rightarrow P_C = 0.1$ 
      - ▶ Completely unidentified



# Results

I compare the proposed method **HD-QLR** (this paper) with

- the conditional QLR test (AM16<sup>2</sup>) : robust against weak identification but not against high dimensionality.
- ML methods (CCDDHNR18<sup>3</sup> and BCFH17<sup>4</sup>): robust against high dimensionality but not against weak identification.

---

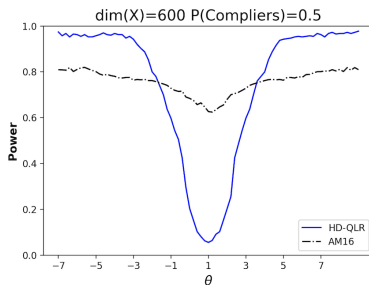
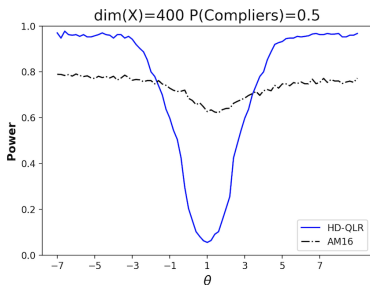
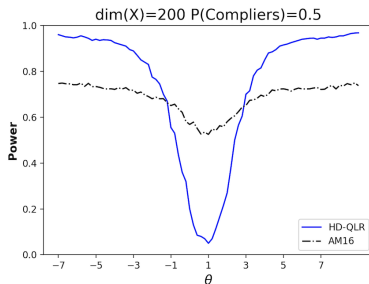
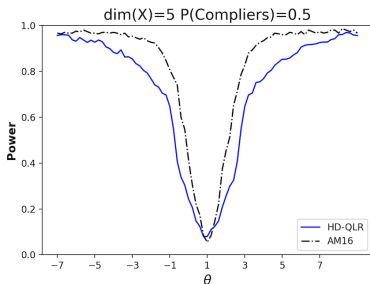
<sup>2</sup>Andrews and Mikusheva (2016).

<sup>3</sup>Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins (2018).

<sup>4</sup>Belloni, Chernozhukov, Fernandez-Val, and Hansen (2017).

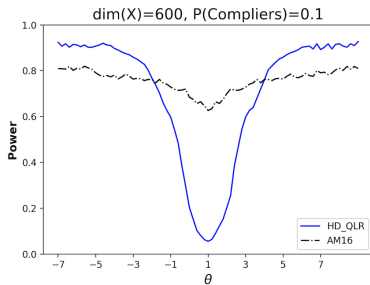
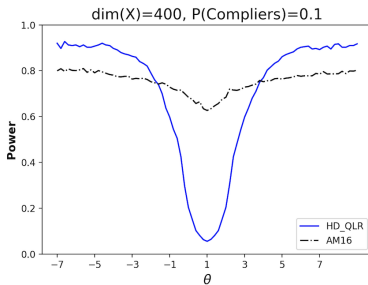
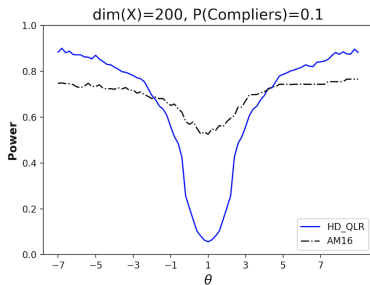
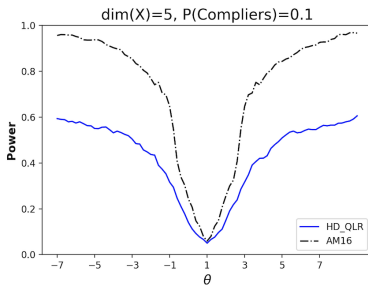
# Comparison: Strong Identification ▶ Power Curve

- Power curve of nominal 5%: AM16 and HD-QLR (this paper)



## Comparison: Weak Identification

- Power curve of nominal 5%: AM16 and HD-QLR (this paper)



# Results

I compare the proposed method **HD-QLR** (this paper) with

- the conditional QLR test (AM16<sup>2</sup>) : robust against weak identification but not against high dimensionality.
- ML methods (**CCDDHNR18**<sup>3</sup> and **BCFH17**<sup>4</sup>): robust against high dimensionality but not against weak identification.

---

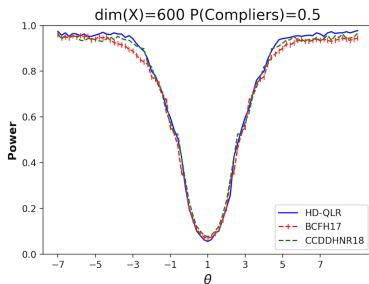
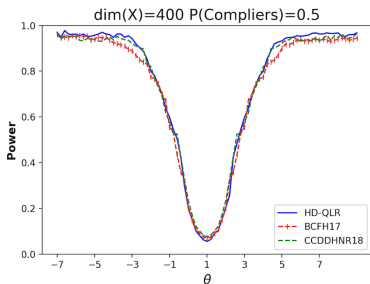
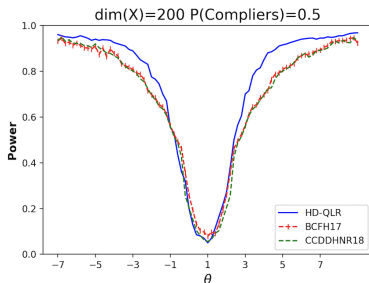
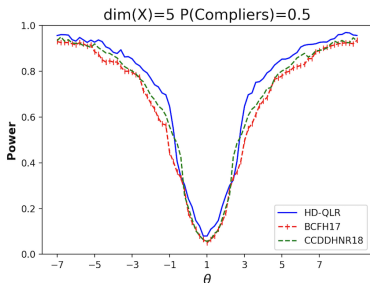
<sup>2</sup>Andrews and Mikusheva (2016).

<sup>3</sup>Chernozhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins (2018).

<sup>4</sup>Belloni, Chernozhukov, Fernandez-Val, and Hansen (2017).

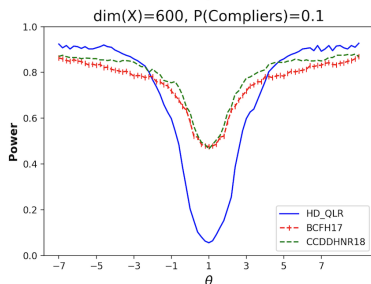
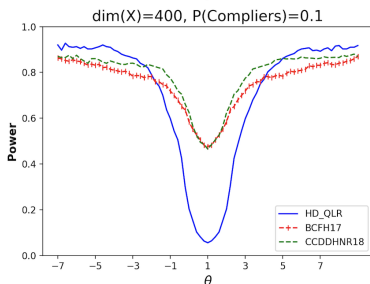
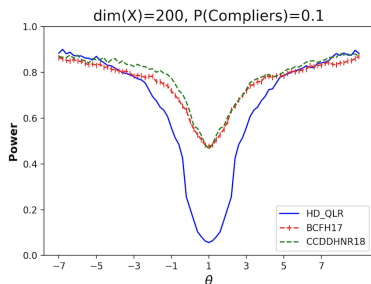
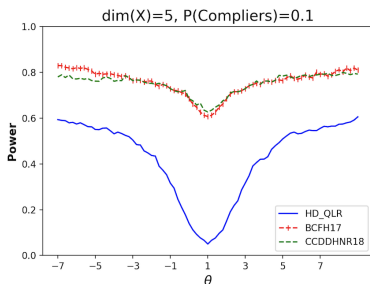
# Comparison: Strong Identification

- Power curve: CCDDHNR18, BCFH17 and HD-QLR (this paper)



# Comparison: Weak Identification

- Power curve: CCDDHNR18, BCFH17 and HD-QLR (this paper)



# Comparison Across Four Approaches

## AM16

- Neyman orthogonal score  $\psi$
- test statistics  $R$
- use traditional methods to handle  $X_i$

## CCDDHNR18 BCFH17

- Neyman orthogonal score  $\psi$
- normal t-test
- use ML to handle  $X_i$

## Drawbacks

- overfitting
- multicollinearity
- cannot perform well under weakly identified scenarios

The proposed HD-QLR takes advantage from both methods:

- Neyman orthogonal score  $\psi$
- test statistics  $R$
- use ML to handle the high-dimensional  $X_i$

# Comparison Across Four Approaches

## AM16

- Neyman orthogonal score  $\psi$
- test statistics  $R$
- use traditional methods to handle  $X_i$

## CCDDHNR18 BCFH17

- Neyman orthogonal score  $\psi$
- normal t-test
- use ML to handle  $X_i$

## Drawbacks

- overfitting
- multicollinearity
- cannot perform well under weakly identified scenarios

The proposed HD-QLR takes advantage from both methods:

- Neyman orthogonal score  $\psi$
- test statistics  $R$
- use ML to handle the high-dimensional  $X_i$



# Comparison Across Four Approaches

## AM16

- Neyman orthogonal score  $\psi$
- test statistics  $R$
- use traditional methods to handle  $X_i$

## CCDDHNR18 BCFH17

- Neyman orthogonal score  $\psi$
- normal t-test
- use ML to handle  $X_i$

## Drawbacks

- overfitting
- multicollinearity
- cannot perform well under weakly identified scenarios

The proposed **HD-QLR** takes advantage from both methods:

- Neyman orthogonal score  $\psi$
- test statistics  $R$
- use ML to handle the high-dimensional  $X_i$

# Table of Contents

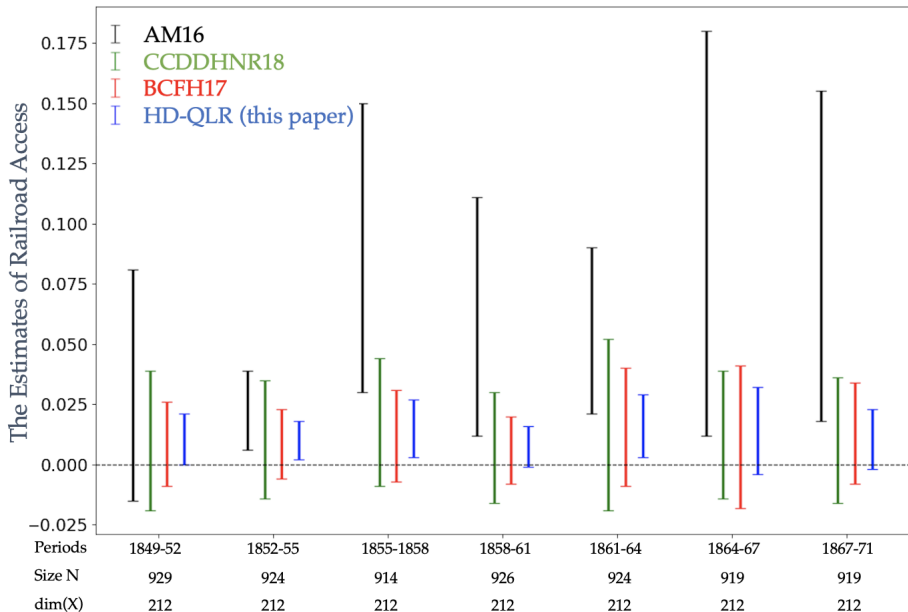
- 1 Introduction
- 2 Overview
- 3 Theory
- 4 Simulation
- 5 Applications**
- 6 Takeaways

## Example One: Hornung (2015) “Railroads and growth in Prussia”

- ★ Data: highly detailed city-level data from the historical German state of Prussia.
- $Y_{it}$  : urban population growth rate for city  $i$  during time period  $t$ .
- $D_i$  : whether the city  $i$  was connected to the railroad in 1848.
- $Z_i$  : whether the city  $i$  was located within a straight-line corridor between two important cities.
- $X_i$  : whether the city had access to the main roads, whether the city had waterway access, military population, age composition, school enrollment rate, etc.

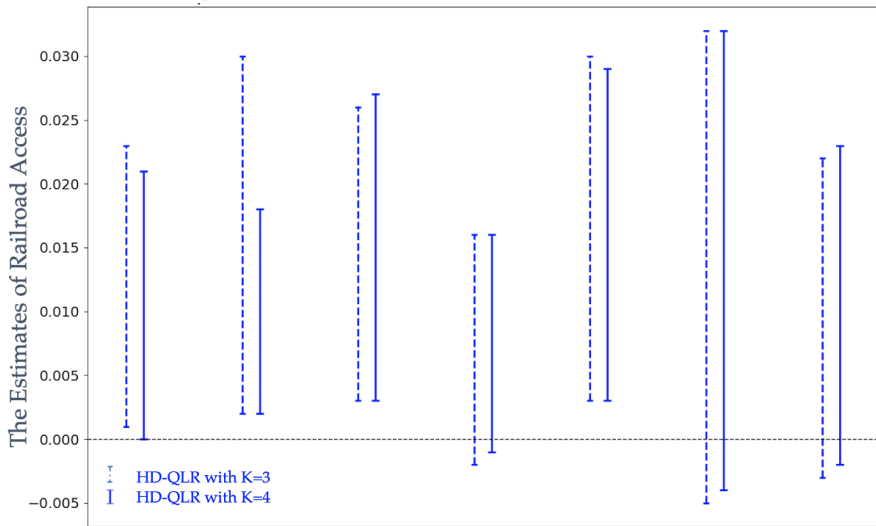
# Results

## Comparison of 95% Confidence Intervals Across Four Approaches



# Results, Continued

## Comparison of Results for Different Numbers of Folds



Periods	1849-52	1852-55	1855-1858	1858-61	1861-64	1864-67	1867-71
Size N	929	924	914	926	924	919	919
dim(X)	212	212	212	212	212	212	212

## Example Two: Ambrus, Field, and Gonzalez (2020) “The Impact on Housing Prices of A Cholera Epidemic”

- ★ *“In August 1854, St. James experienced a sudden outbreak of cholera when one of the 13 shallow wells that serviced the parish, the Broad Street pump, became contaminated with cholera bacteria.”*
- $Y_i$  : the log rental price of house  $i$  in 1864.
- $D_i$  : whether house  $i$  had at least one cholera death.
- $Z_i$  : whether house  $i$  fell inside the contaminated areas.
- $X_i$  : distance to the closest pump, distance to the fire station, distance to the urinal, sewer access, among a total of 23 variables.

## Results

	AM16	CCDDHNR18	BCFH17	HD-QLR (this paper)
95% CI	[-2.160, -0.670]	[-1.132, 0.406]	[-1.291, 0.576]	[-1.080, 0.035]
length of CI	1.490	1.538	1.866	1.115

**Table:** Displayed are the CIs and the length of CI. Inference results are based on 10 iterations of resampled cross fitting with  $K = 4$  folds for cross fitting. The number of observations  $N = 467$ .

# Table of Contents

- 1 Introduction
- 2 Overview
- 3 Theory
- 4 Simulation
- 5 Applications
- 6 Takeaways**



# Takeaways

- I develop an inference method for the **high-dimensional LATE**, independent of **the strength of identification**.

	Low-dimensional Model	High-dimensional Model
Strong Identification	z-test	CCDDHNR18, BCFH17
Weak Identification	AM16	HD-QLR (my paper)

- The proposed method has uniformly correct asymptotic size.
- The proposed test is robust against weak identification and high dimensionality, outperforming other conventional methods.
- The proposed method yields narrower confidence intervals than conventional methods, as demonstrated in applications.

# Takeaways

- I develop an inference method for the **high-dimensional LATE**, independent of **the strength of identification**.

	Low-dimensional Model	High-dimensional Model
Strong Identification	z-test	CCDDHNR18, BCFH17
Weak Identification	AM16	HD-QLR (my paper)

- The proposed method has uniformly correct asymptotic **size**.
- The proposed test is robust against **weak identification** and **high dimensionality**, outperforming other conventional methods.
- The proposed method yields **narrower confidence intervals** than conventional methods, as demonstrated in applications.

# Thank you!

feel free to email me any comments  
[yukun.ma@vanderbilt.edu](mailto:yukun.ma@vanderbilt.edu)

## Motivation: Lee et al. (2022)

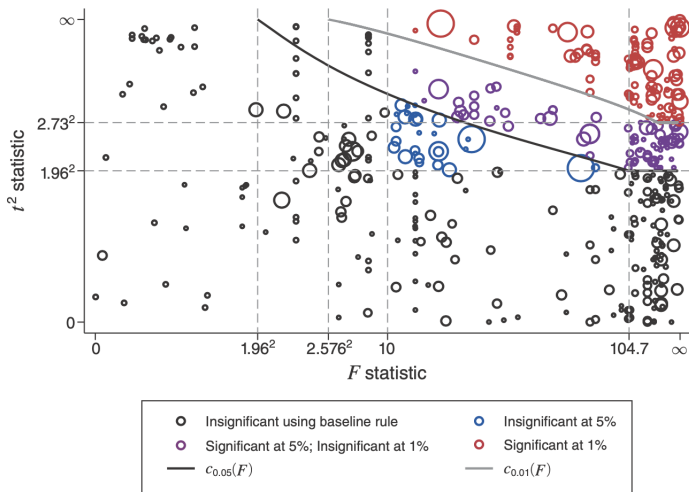


Figure: *American Economic Review* 2013-2019

## Tuning Parameters

### Lemma (Convergence rate for Lasso with logistic model)

Suppose some regularity assumptions hold. In addition, suppose that the penalty choice  $\lambda_1 = K_1 \sqrt{N \log(pN)}$  and  $\lambda_2 = K_2 \sqrt{N \log(pN)}$  for  $K_1, K_2 > 0$ . Then with probability  $1 - o(1)$ ,

$$\|(\hat{\beta}_{11}, \hat{\beta}_{12}) - (\beta_{11}^0, \beta_{12}^0)\| \vee \|\hat{\gamma} - \gamma^0\| \lesssim \sqrt{\frac{s_N \log(pN)}{N}}.$$

### Lemma (Convergence rate for Lasso with OLS)

Suppose some regularity assumptions hold. Moreover, suppose that the penalty choice  $\lambda_3 = K_3 \sqrt{N \log(pN)}$  for  $K_3 > 0$ . Then with probability  $1 - o(1)$ ,

$$\|(\hat{\beta}_{21}, \hat{\beta}_{22}) - (\beta_{21}^0, \beta_{22}^0)\| \lesssim \sqrt{\frac{s_N \log(pN)}{N}}.$$

# Null Hypothesis

- Define  $S_N(\cdot) = \mathbf{E}_P[N^{-1/2} \sum_{i=1}^N \psi(W_i; \cdot, \eta_0)]$ .

Case 1:  $H_0 : \theta = \theta_0$  with assuming  $\theta$  is point-identified

$$\hookrightarrow S_N(\theta_0) = \mathbf{0}.$$

Case 2:  $H_0 : \theta \in \Theta_I$  with the identified set  $\Theta_I \subset \Theta$  when point identification fails

$$\hookrightarrow S_N(\theta) = \mathbf{0} \text{ for } \forall \theta \in \Theta_I.$$

- Let  $\mathcal{S}_0$  be the collection of function  $S_N(\cdot)$  satisfying  $S_N(\theta) = \mathbf{0}$ .

$$\hookrightarrow H'_0 : S_N(\cdot) \in \mathcal{S}_0.$$

## Spare eigenvalue

For any  $\mathbf{T} \subset [p + 1]$ ,  $\boldsymbol{\delta} = (\delta_1, \dots, \delta_{p+1})' \in \mathbb{R}^{p+1}$  with  $\delta_{\mathbf{T},j} = \delta_j$  if  $j \in \mathbf{T}$  and  $\delta_{\mathbf{T},j} = 0$  if  $j \notin \mathbf{T}$ . Define the minimum and maximum sparse eigenvalue by

$$\begin{aligned}\phi_{\min}(m) &= \inf_{\|\boldsymbol{\delta}\|_0 \leq m} \frac{\|(\mathbf{Z}_i, \mathbf{X}'_i)\boldsymbol{\delta}\|_{2,N}}{\|\boldsymbol{\delta}_{\mathbf{T}}\|_1} \\ \phi_{\max}(m) &= \sup_{\|\boldsymbol{\delta}\|_0 \leq m} \frac{\|(\mathbf{Z}_i, \mathbf{X}'_i)\boldsymbol{\delta}\|_{2,N}}{\|\boldsymbol{\delta}_{\mathbf{T}}\|_1}.\end{aligned}$$

▸ back

## Identification Assumption Comparison

- In my paper:

$$\mathbf{E}_P[D|Z = 1] \geq \mathbf{E}_P[D|Z = 0].$$

- In weak identification literature:

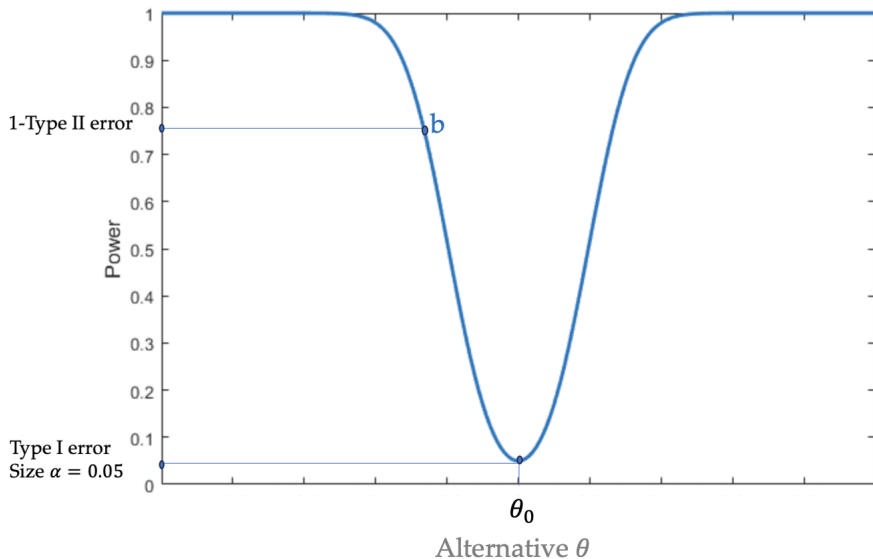
$$\mathbf{E}_P[D|Z = 1] - \mathbf{E}_P[D|Z = 0] = \frac{C_1}{\sqrt{N}} \quad \text{with } C_1 > 0.$$

- In ML literature:

$$\mathbf{E}_P[D|Z = 1] - \mathbf{E}_P[D|Z = 0] \geq C_2 \quad \text{with } C_2 > 0.$$



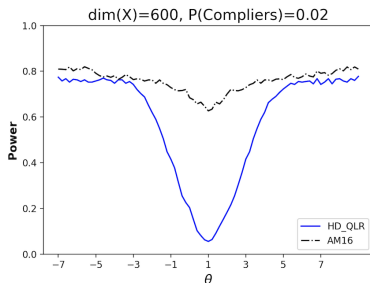
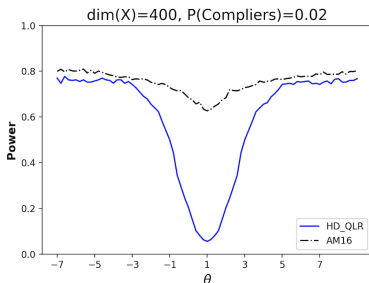
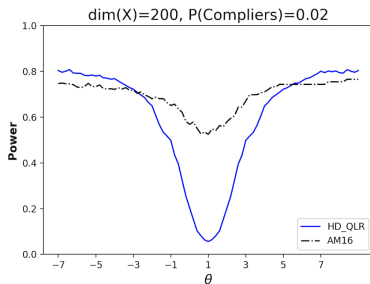
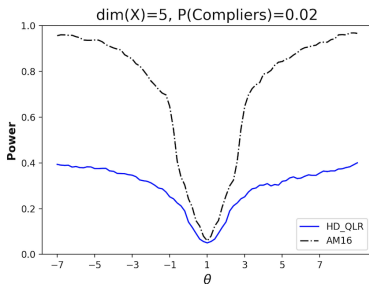
# Power Curve



▶ back

# Comparisons: Unidentified case

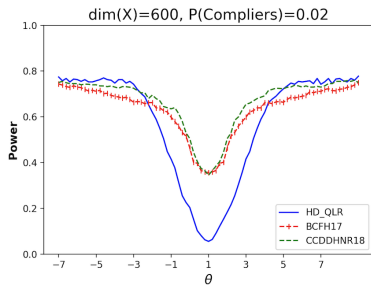
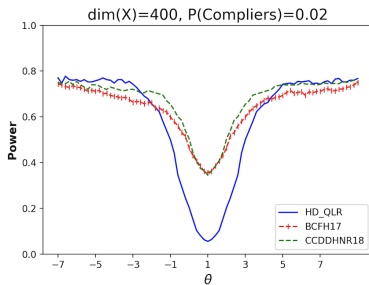
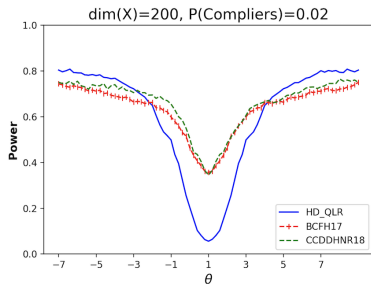
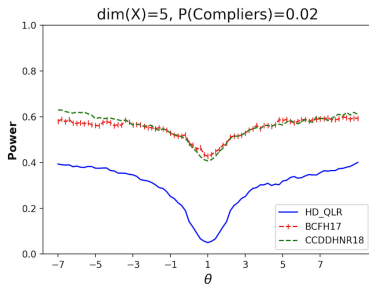
- Power curve: AM16 with HD-QLR (this paper)



# Comparisons: Unidentified Case

▶ back

- Power curve: **CCDDHNR18**, **BCFH17** with **HD-QLR** (this paper)



- BELLONI, A., V. CHERNOZHUKOV, AND K. KATO (2015): “Uniform post-selection inference for least absolute deviation regression and other Z-estimation problems,” *Biometrika*, 102, 77–94.
- BLANDHOL, C., J. BONNEY, M. MOGSTAD, AND A. TORGOVITSKY (2022): “When is TSLS actually late?” Tech. rep., National Bureau of Economic Research.
- CHERNOZHUKOV, V., D. CHETVERIKOV, AND K. KATO (2013): “Gaussian approximations and multiplier bootstrap for maxima of sums of high-dimensional random vectors,” *The Annals of Statistics*, 41, 2786–2819.
- (2016): “Empirical and multiplier bootstraps for suprema of empirical processes of increasing complexity, and related Gaussian couplings,” *Stochastic Processes and their Applications*, 126, 3632–3651.
- (2017): “Central limit theorems and bootstrap in high dimensions,” *The Annals of Probability*, 45, 2309–2352.
- KLEIBERGEN, F. (2005): “Testing parameters in GMM without assuming that they are identified,” *Econometrica*, 73, 1103–1123.
- MOREIRA, M. J. (2003): “A conditional likelihood ratio test for structural models,” *Econometrica*, 71, 1027–1048.
- STOCK, J. H. AND J. H. WRIGHT (2000): “GMM with weak identification,” *Econometrica*, 68, 1055–1096.
- TAN, Z. (2006): “Regression and weighting methods for causal inference using instrumental variables,” *Journal of the American Statistical Association*, 101, 1607–1618.
- VAN DER VAART, A. AND J. A. WELLNER (1996): “WEAK CONVERGENCE AND EMPIRICAL PROCESSES,” .