

# Marginal Treatment Effects with Many Instruments

Matias D. Cattaneo, University of Michigan

Michael Jansson, UC-Berkeley and CREATES

Xinwei Ma, University of Michigan

July 2016

Draft prepared for NBER conference is here:

[http://www.umich.edu/~cattaneo/papers/Cattaneo-Jansson-Ma\\_2016\\_MTE.pdf](http://www.umich.edu/~cattaneo/papers/Cattaneo-Jansson-Ma_2016_MTE.pdf)

# Overview

## What this paper does:

- 1 Study estimation and inference for MTE and functionals thereof.
- 2 Develop new large-sample point estimation and distributional approximations allowing for many (included) instruments entering the probability of selection.
- 3 Results encompass standard ones, so we provide strict asymptotic improvements.
- 4 New estimation/inference methods are fully automatic (jackknife + wild bootstrap).
- 5 Results extend to a large class of two-step M-estimators with high-dimensional first step (e.g., propensity score weighting, generated regressors, etc.).

## What this paper does not do:

- 1 Discuss identification or interpretation of MTE, and functional thereof, with “few” or “many” instruments.
  - ▶ Our model/estimand is taken as given based on results in the literature.
- 2 Discuss model selection and/or ultra-high-dimensional methods for selecting in or out instruments.
  - ▶ Our methods are high-dimensional in the sense that the number of included instruments may be “large” (relative to the sample size).

## Setup and Estimation in Practice

- Random sample of units  $i = 1, 2, \dots, n$  from large population:

$$y_i = t_i y_i(1) + (1 - t_i) y_i(0), \quad t_i = \mathbf{1}[p_i \geq v_i], \quad v_i | \mathbf{x}_i \sim \text{Uniform}[0, 1].$$

- Parameter of interest:

$$\tau_{\text{mte}} = \tau_{\text{mte}}(a | \mathbf{x}_i) = \mathbb{E}[y_i(1) - y_i(0) | v_i = a, \mathbf{x}_i] = \frac{\partial}{\partial a} \mathbb{E}[y_i | p_i = a, \mathbf{x}_i]$$

- ▶ Identification assumption:  $(p_i, \mathbf{z}'_i) \perp\!\!\!\perp (y_i(1), y_i(0), v_i) \mid \mathbf{x}_i$
- ▶ Practical assumption:  $\mathbb{E}[y_i | p_i = a, \mathbf{x}_i] = m(\mathbf{x}_i, a, \boldsymbol{\theta}_0)$

- Estimation approach:

- ① Estimate propensity score using possibly “many” instruments  $\mathbf{z}_i \in \mathbb{R}^k$ :

$$\hat{p}_{n,i} = \mathbf{z}'_i \hat{\boldsymbol{\beta}}, \quad \hat{\boldsymbol{\beta}} \in \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (t_i - \mathbf{z}'_i \boldsymbol{\beta})^2,$$

- ② Estimate MTE:

$$\hat{\tau}_{\text{mte}} = \hat{\tau}_{\text{mte}}(a | \mathbf{x}) = \frac{\partial}{\partial a} m(\mathbf{x}, a, \hat{\boldsymbol{\theta}}), \quad \hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (y_i - m(\mathbf{x}_i, \hat{p}_{n,i}, \boldsymbol{\theta}))^2.$$

## Overview of Results

- 1 When  $k$  is “large” relative to  $n$ , first-order bias in distributional approximation:

$$\sqrt{n}(\hat{\tau}_{\text{mte}} - \tau_{\text{mte}}) = \frac{k}{\sqrt{n}}B + \frac{1}{\sqrt{n}} \sum_{i=1}^n \psi_i + o_{\mathbb{P}}(1) \rightarrow_d \mathcal{N}(\rho B, V)$$

where

$$\frac{k}{\sqrt{n}} \rightarrow \rho \in [0, \infty)$$

- 2 **Jackknife** bias and variance estimation works, even when  $\rho > 0$ :

$$\hat{B}_{\text{Jack}} - B = o_{\mathbb{P}}(1) \quad \text{and} \quad \hat{V}_{\text{Jack}} - V = o_{\mathbb{P}}(1)$$

- 3 **Wild bootstrap** works after jackknife estimation, even when  $\rho > 0$ :

$$\sup_{t \in \mathbb{R}} \left| \mathbb{P}^* \left[ \frac{\sqrt{n} \left( \hat{\tau}_{\text{mte}}^* - \hat{\tau}_{\text{mte}} - \frac{k}{n} (\hat{B}_{\text{Jack}}^* - \hat{B}_{\text{Jack}}) \right)}{\sqrt{\hat{V}_{\text{Jack}}^*}} \leq t \right] - \mathbb{P} \left[ \frac{\sqrt{n} \left( \hat{\tau}_{\text{mte}} - \tau_{\text{mte}} - \frac{k}{n} \hat{B}_{\text{Jack}} \right)}{\sqrt{\hat{V}_{\text{Jack}}}} \leq t \right] \right| = o_{\mathbb{P}}(1)$$

- 4 Results extend to functionals of  $\hat{\tau}_{\text{mte}} = \hat{\tau}_{\text{mte}}(a|\mathbf{x})$ , such as those used for extrapolation of treatment effects or to construct other treatment effects.

## Motivation and Discussion: Simulations

- Standard  $(1 - \alpha)\%$  confidence intervals for  $\tau_{\text{mte}}$ :

$$\left[ \hat{\tau}_{\text{mte}} - \Phi_{1-\alpha/2}^{-1} \cdot \sqrt{\frac{\hat{V}}{n}} \quad , \quad \hat{\tau}_{\text{mte}} - \Phi_{\alpha/2}^{-1} \cdot \sqrt{\frac{\hat{V}}{n}} \right]$$

where:

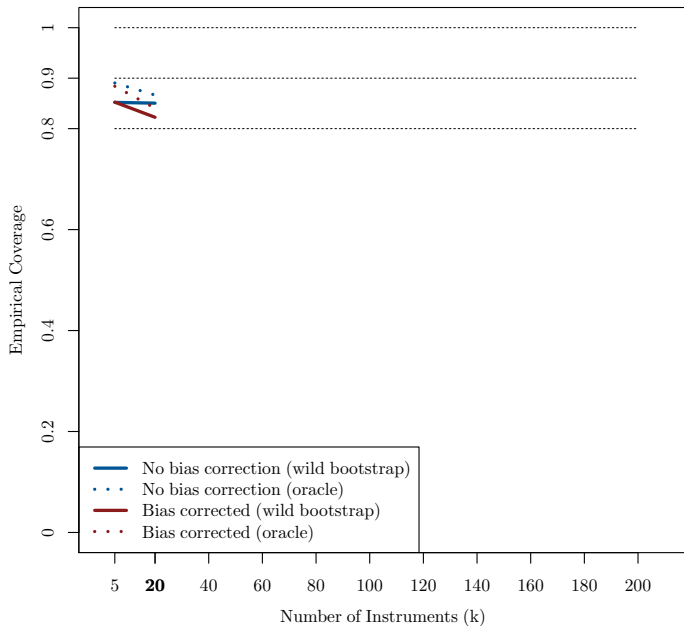
- ▶  $\Phi_{\alpha}^{-1}$  is standard normal quantile (e.g.,  $\Phi_{1-\alpha/2}^{-1} = \Phi_{\alpha/2}^{-1} = 1.65$  if  $\alpha = 0.90$ ).
  - ▶  $\hat{V}/n$  is some standard error estimator... Because  $V$  is complex and involves many unknown quantities, automatic methods such as  $\hat{V}_{\text{Jack}}$  are preferred (if valid).
- Our proposed method:  $(1 - \alpha)\%$  confidence intervals for  $\tau_{\text{mte}}$ :

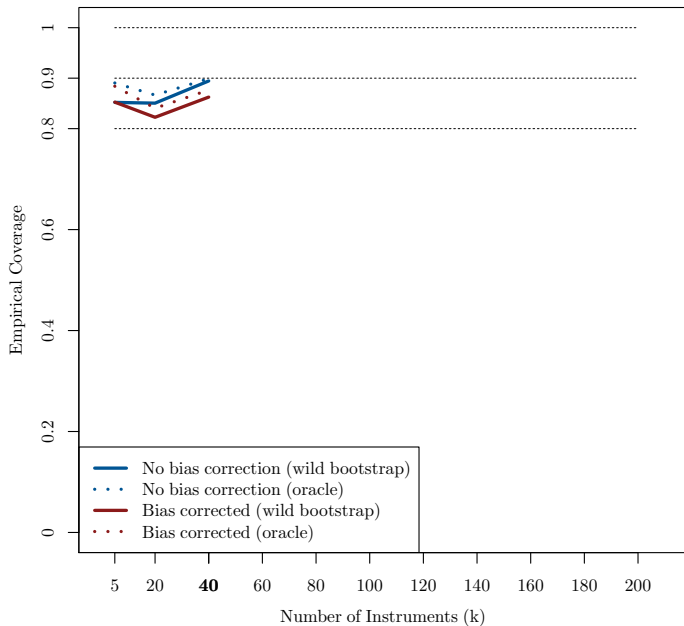
$$\left[ \left( \hat{\tau}_{\text{mte}} - \frac{k}{n} \hat{B}_{\text{Jack}} \right) - q_{1-\alpha/2}^* \cdot \sqrt{\frac{\hat{V}_{\text{Jack}}}{n}} \quad , \quad \left( \hat{\tau}_{\text{mte}} - \frac{k}{n} \hat{B}_{\text{Jack}} \right) - q_{\alpha/2}^* \cdot \sqrt{\frac{\hat{V}_{\text{Jack}}}{n}} \right]$$

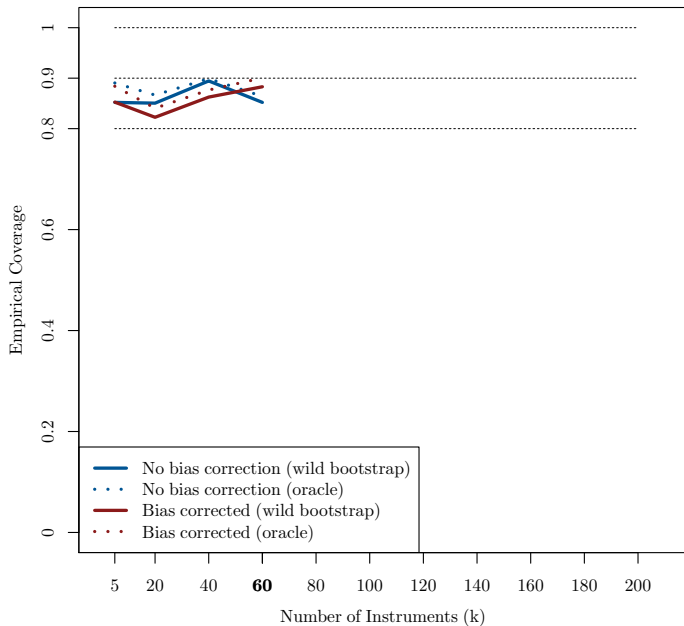
where:

- ▶  $\hat{B}_{\text{Jack}}$  and  $\hat{V}_{\text{Jack}}$  jackknife estimators, shown consistent even when  $\frac{k}{\sqrt{n}} \neq 0$ .
- ▶  $q_{\alpha}^*$  is quantile from wild bootstrap distribution of jackknife-based t-test statistic:

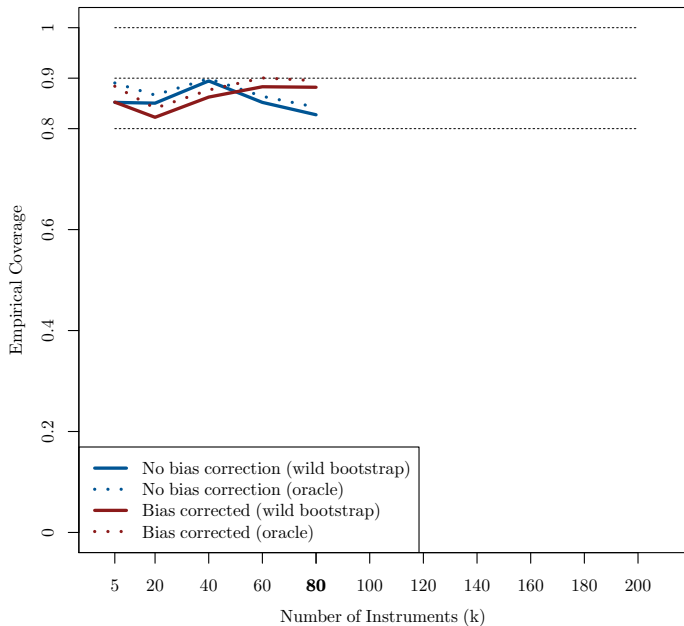
$$T^* = \frac{\sqrt{n} \left( \hat{\tau}_{\text{mte}}^* - \hat{\tau}_{\text{mte}} - \frac{k}{n} (\hat{B}_{\text{Jack}}^* - \hat{B}_{\text{Jack}}) \right)}{\sqrt{\hat{V}_{\text{Jack}}^*}}$$

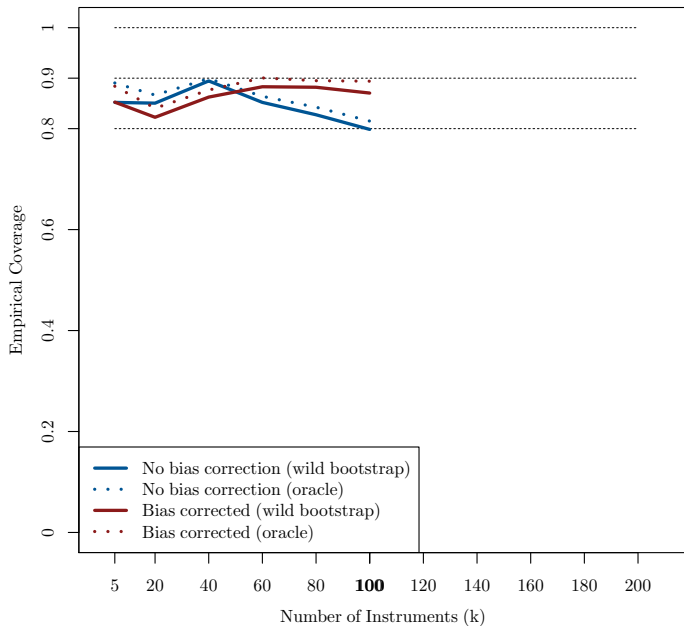


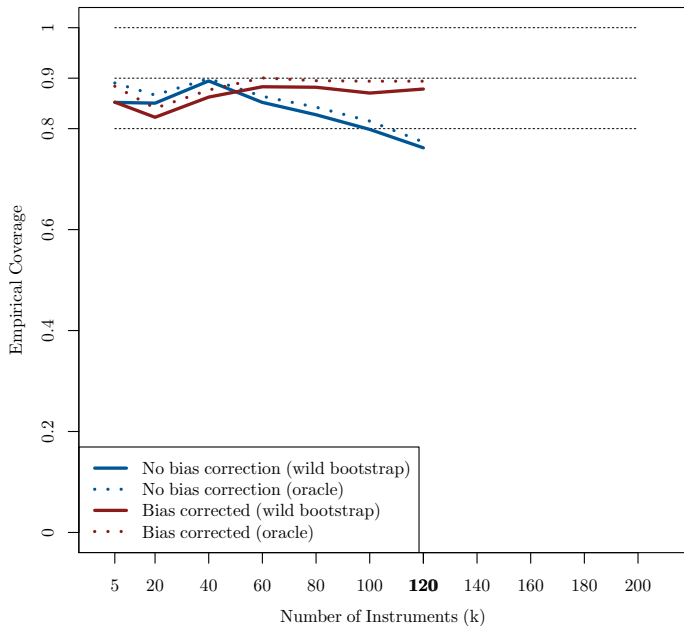


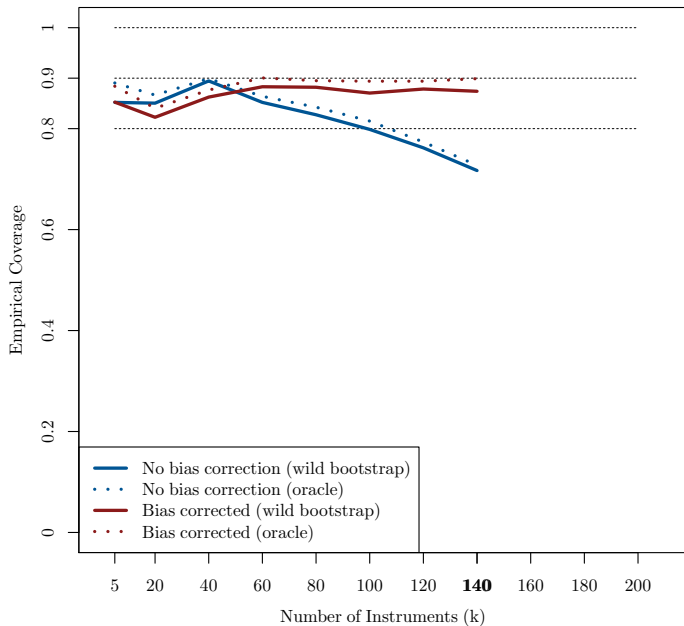


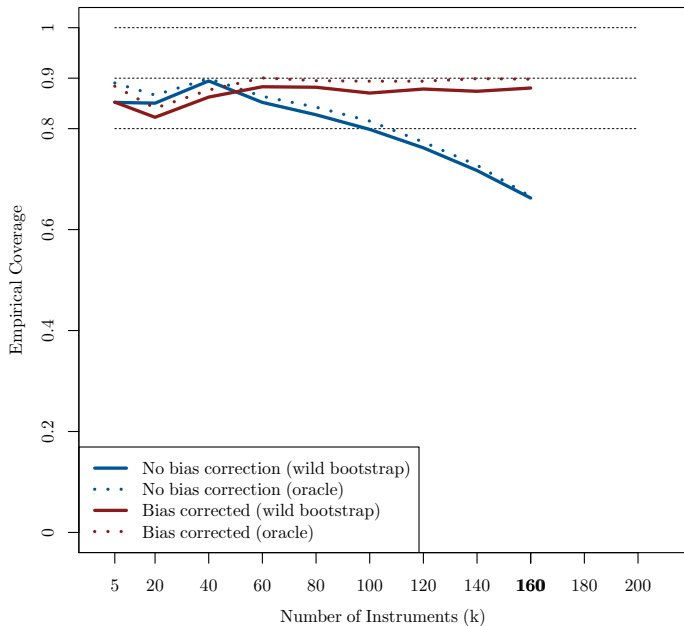


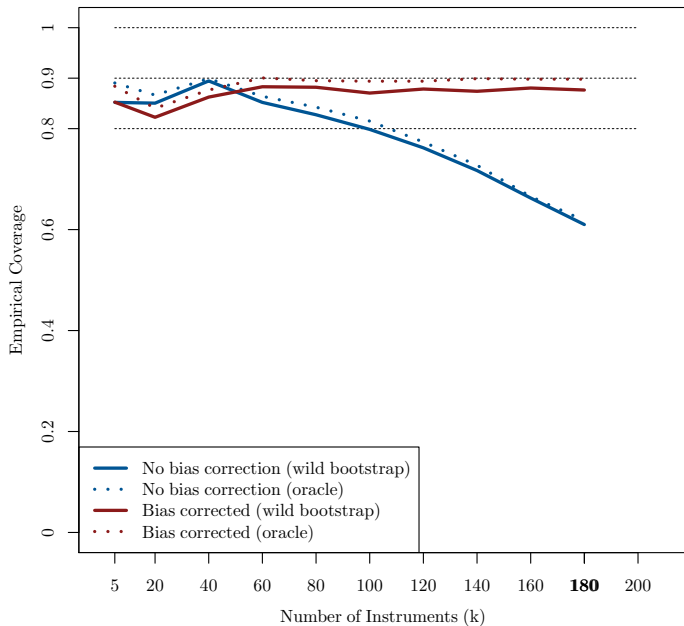


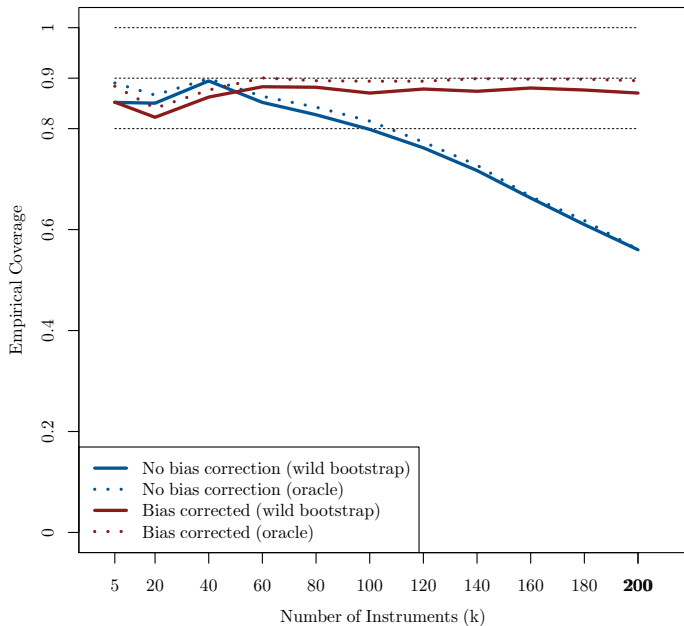












## Intuition Behind Results

- Recall that

$$\hat{\tau}_{\text{mte}} = \frac{\partial}{\partial a} m(\mathbf{x}, a, \hat{\boldsymbol{\theta}}), \quad \hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^n (y_i - m(\mathbf{x}_i, \hat{p}_{n,i}, \boldsymbol{\theta}))^2,$$

$$\hat{p}_{n,i} = \mathbf{z}'_i \hat{\boldsymbol{\beta}}, \quad \hat{\boldsymbol{\beta}} \in \arg \min_{\boldsymbol{\beta}} \sum_{i=1}^n (t_i - \mathbf{z}'_i \boldsymbol{\beta})^2,$$

- $\hat{\tau}_{\text{mte}} \rightarrow_{\mathbb{P}} \tau_{\text{mte}}$ , even when  $\frac{k}{\sqrt{n}} \not\rightarrow 0$ , if  $\max_{1 \leq i \leq n} |\hat{p}_{n,i} - p_i| = o_{\mathbb{P}}(1)$ .

- Distribution theory is about “local” behavior:

$$\begin{aligned} & \sqrt{n}(\hat{\tau}_{\text{mte}} - \tau_{\text{mte}}) \\ & \approx \frac{1}{\sqrt{n}} \sum_{i=1}^n g(\mathbf{x}_i, p_i) + \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial g(\mathbf{x}_i, p_i)}{\partial p_i} \cdot [\hat{p}_{n,i} - p_i] + \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial^2 g(\mathbf{x}_i, p_i)}{\partial p_i^2} \cdot [\hat{p}_{n,i} - p_i]^2 \\ & \quad \downarrow \qquad \qquad \qquad \downarrow \qquad \qquad \qquad \downarrow \\ & \mathcal{N}(0, V_1) \qquad \qquad \mathcal{N}(0, V_2) + \frac{k}{\sqrt{n}} B_1 \qquad \qquad \frac{k}{\sqrt{n}} B_2 \end{aligned}$$



## Application to Treatment Effect Extrapolation

- Large class of treatment effects can be recovered from MTE:

$$\tau_{\omega}(\mathbf{x}) = \int_0^1 \tau_{\text{mte}}(a|\mathbf{x})\omega(a|\mathbf{x})da$$

- (Policy, Extrapolation, etc.) Estimator:

$$\hat{\tau}_{\omega} = \hat{\tau}_{\omega}(\mathbf{x}) = \int_0^1 \hat{\tau}_{\text{mte}}(a|\mathbf{x})\omega(a|\mathbf{x})da$$

- Our results apply directly (using the “delta method”):

$$\sqrt{n}(\hat{\tau}_{\omega} - \tau_{\omega}) = \frac{k}{\sqrt{n}}B_{\omega} + \frac{1}{\sqrt{n}}\sum_{i=1}^n \phi_i + o_{\mathbb{P}}(1) \rightarrow_d \mathcal{N}(\rho B_{\omega}, V_{\omega})$$

where

$$\frac{k}{\sqrt{n}} \rightarrow \rho \in [0, \infty)$$

- **Jackknife** bias and variance estimation works, even when  $\rho > 0$ .
- **Wild bootstrap** works after jackknife estimation, even when  $\rho > 0$ .

## (Backup) Empirical Illustration: Card (1993)

- We try out a preliminary (backup) empirical application, using Card (1993)'s returns to college education paper.
- We will soon have access to the NLSY Geocode data files, to revisit Carneiro, Heckman and Vytlačil (2011)'s approach to estimate marginal returns to education, with a much richer set of raw and expanded instruments.
- Card (1993) Empirical Illustration:
  - ▶ Subsample of National Longitudinal Survey of Young Men (NLSYM).
  - ▶ Initial survey 1966, sample survey 1976; for example,  $n = 2,307$  white males.
  - ▶ Treatment:  $t_i = \mathbb{1}(\text{education in 1976} \geq 13)$ .
  - ▶ Consider specifications from 25 to 59 instruments.
  - ▶ Consider 2SLS, ATE, ATE with bias-correction.
  - ▶ Excluded instrument: college proximity measures.
  - ▶ **Main empirical finding:** results are robust to “many instruments” bias.

**Table:** Returns to College Education, Whites Subsample,  $n = 2,307$

	(1)	(2)	(3)	(4)	(5)	(6)	(7)	(8)	(9)	(10)	(11)	(12)
2SLS	0.125 (0.018)	0.129 (0.017)	0.138 (0.016)	0.131 (0.017)	0.136 (0.016)	0.143 (0.015)	0.144 (0.062)	0.170 (0.046)	0.186 (0.032)	0.192 (0.055)	0.177 (0.034)	0.186 (0.027)
ATE	0.128 (0.030)	0.136 (0.026)	0.142 (0.022)	0.139 (0.029)	0.137 (0.022)	0.145 (0.018)	0.106 (0.067)	0.142 (0.045)	0.158 (0.034)	0.162 (0.053)	0.148 (0.031)	0.164 (0.027)
ATE (bc)	0.143 (0.035)	0.164 (0.033)	0.156 (0.026)	0.148 (0.034)	0.171 (0.030)	0.158 (0.024)	0.118 (0.075)	0.188 (0.061)	0.172 (0.041)	0.170 (0.064)	0.213 (0.047)	0.170 (0.037)
<b>Outcome Eqn.</b>												
Baseline and Geographic <sup>(i)</sup>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Parental Education <sup>(ii)</sup>							✓	✓	✓	✓	✓	✓
<b>Selection Eqn.</b>												
Baseline and Geographic <sup>(i)</sup>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Parental Education <sup>(ii)</sup>	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓	✓
Family Structure <sup>(iii)</sup>			✓			✓			✓			✓
4-yr College Proximity <sup>(iv)</sup>	✓	✓	✓				✓	✓	✓			
Interactions <sup>(v)</sup>		✓	✓					✓	✓			
Three College Proximities <sup>(vi)</sup>				✓	✓	✓				✓	✓	✓
Interactions <sup>(v)</sup>					✓	✓					✓	✓
$k$	25	35	37	27	57	59	25	35	37	27	57	59
$k/n$	0.01	0.02	0.02	0.01	0.02	0.03	0.01	0.02	0.02	0.01	0.02	0.03
$k/\sqrt{n}$	0.52	0.73	0.77	0.56	1.19	1.23	0.52	0.73	0.77	0.56	1.19	1.23

## Discussion and Conclusion

- 1 Motivation: Elucidate the finite-sample sensitivity of distributional approximations
- 2 Some examples:
  - ▶ HAC Estimation:  $b_n \propto n$
  - ▶ Weak IV:  $\pi \propto n^{-1/2}$
  - ▶ Many IV:  $\dim(\pi) \propto n$
  - ▶ FE Panels:  $T \propto n$
  - ▶ NN-based (“matching”):  $M_n \propto 1$
  - ▶ Kernel-based semiparametrics:  $h_n \propto n^\gamma$
  - ▶ Series-based semiparametrics:  $K_n \propto n^\gamma$
  - ▶ High-dimensional linear regression models:  $K_n \propto n$
- 3 Today:
  - ▶ Two-step estimation with high-dimensional first step:  $k \propto \sqrt{n}$
  - ▶ Gaussian distributional approximation not valid: first-order bias!
  - ▶ Jackknife bias and variance estimation works: fully automatic!
  - ▶ Wild bootstrap applied to jackknife-based t-statistic works very well: fully automatic!

THANK YOU !

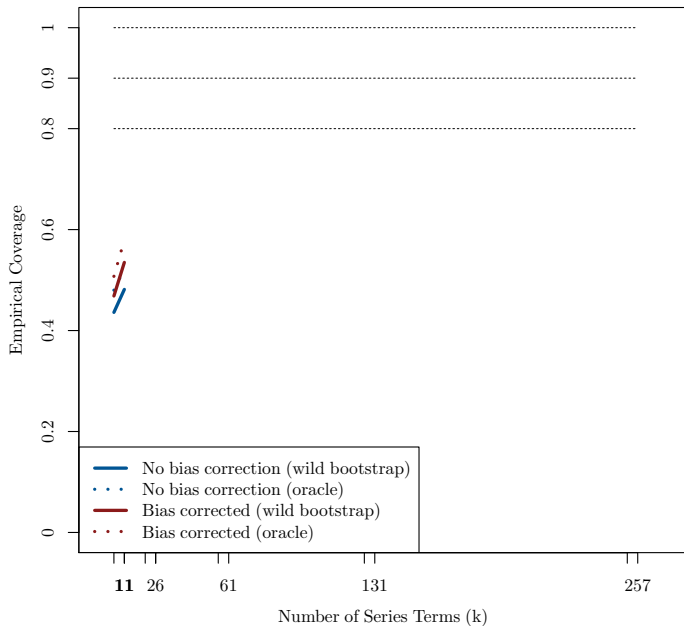
## Appendix: Simulations with Misspecification Error

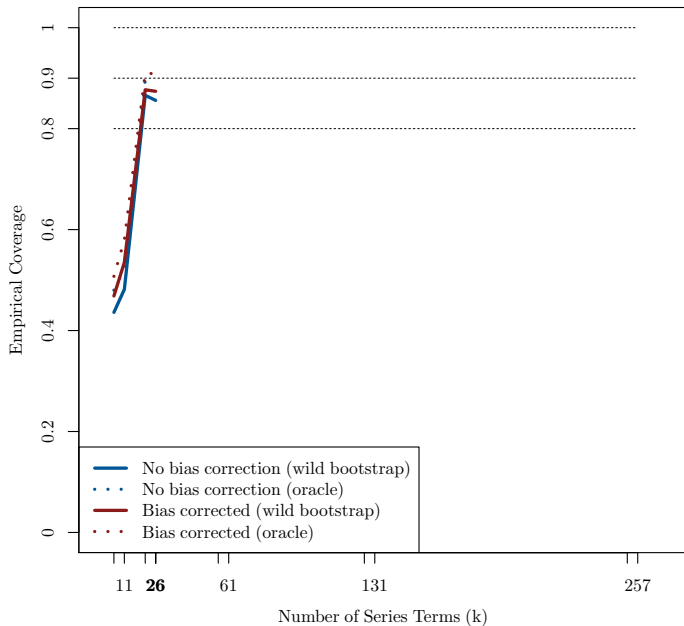
- $n = 2,000$

Table: Polynomial Basis Expansion.

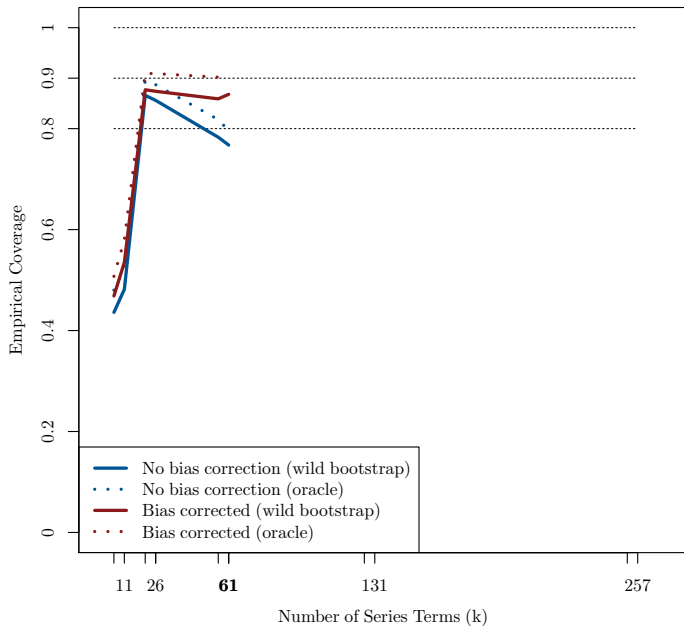
Series terms ( $k$ )	$\mathbf{s}^k(\mathbf{z}_i)$
6	1 and $\mathbf{z}_i$
11	1, $\mathbf{z}_i$ and $[z_{1i}^2, z_{2i}^2, \dots, z_{5i}^2]$
21	$\mathbf{s}^{11}(\mathbf{z}_i)$ and 2 <sup>nd</sup> -order interactions
26	$\mathbf{s}^{21}(\mathbf{z}_i)$ and $[z_{1i}^3, z_{2i}^3, \dots, z_{5i}^3]$
56	$\mathbf{s}^{26}(\mathbf{z}_i)$ and 3 <sup>rd</sup> -order interactions
61	$\mathbf{s}^{56}(\mathbf{z}_i)$ and $[z_{1i}^4, z_{2i}^4, \dots, z_{5i}^4]$
126	$\mathbf{s}^{61}(\mathbf{z}_i)$ and 4 <sup>th</sup> -order interactions
131	$\mathbf{s}^{126}(\mathbf{z}_i)$ and $[z_{1i}^5, z_{2i}^5, \dots, z_{5i}^5]$
252	$\mathbf{s}^{131}(\mathbf{z}_i)$ and 5 <sup>th</sup> -order interactions
257	$\mathbf{s}^{252}(\mathbf{z}_i)$ and $[z_{1i}^6, z_{2i}^6, \dots, z_{5i}^6]$

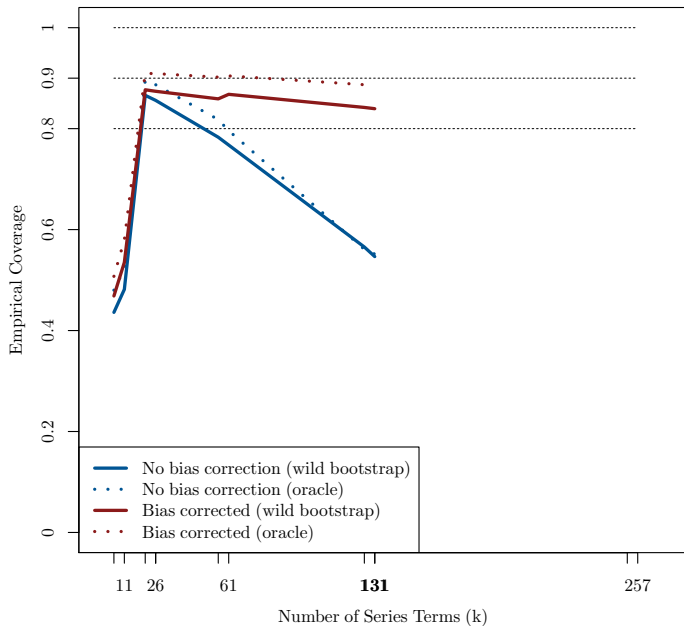
**Notes.** In DGP 3, the propensity score is a nonlinear function of five raw instruments, denoted by  $\mathbf{z}_i \in \mathbb{R}^5$ . We use the series method to estimate the propensity score, by including interactions and higher moments of the raw instruments, where the series basis is denoted by  $\mathbf{s}^k(\mathbf{z}_i)$ . An intercept is also included.

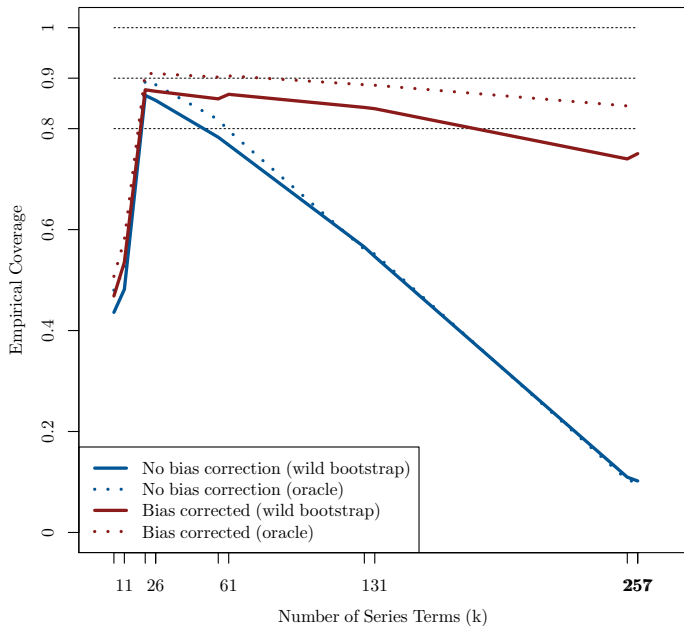




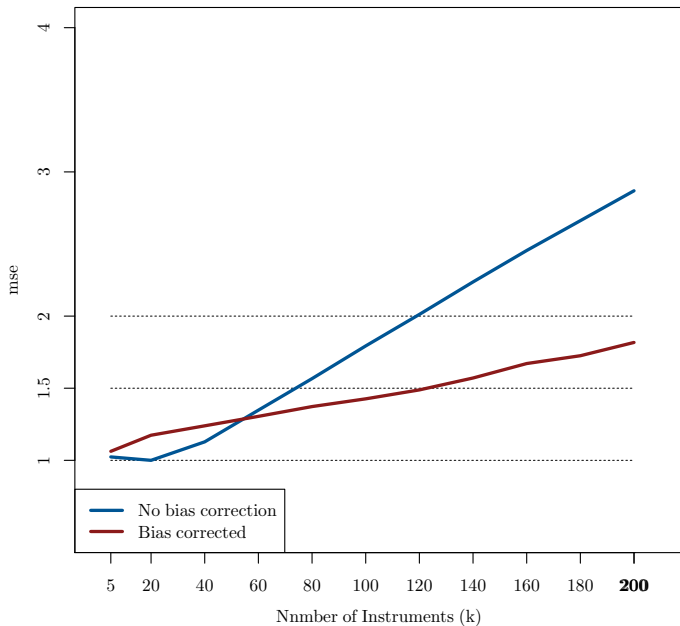




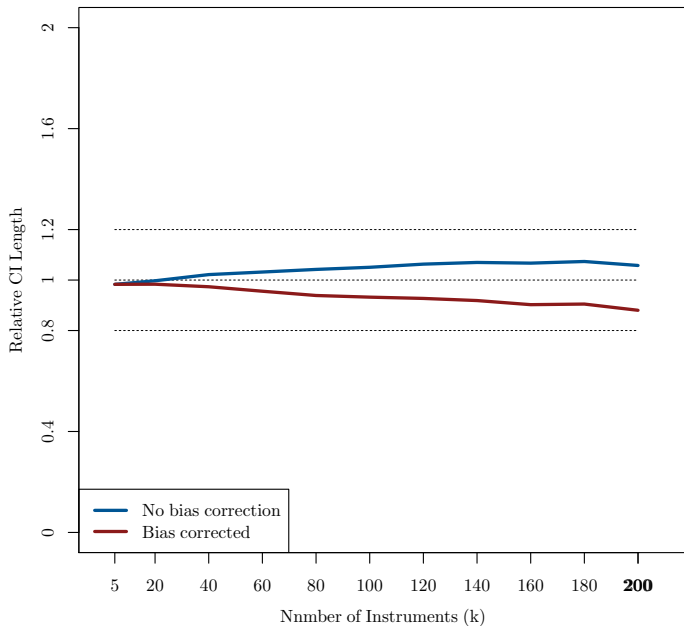




## Appendix: Simulations, MSE



## Appendix: Simulations, Interval Length



## Appendix: Jackknife Estimation

- **Step 1.** For each observation  $j = 1, 2, \dots, n$  estimate the propensity score without using the  $j$ -th observation, denoted by  $\hat{p}_{n,i}^{(j)}$ , and compute the leave- $j$ -out estimator by solving

$$\hat{\boldsymbol{\theta}}^{(j)} = \arg \min_{\boldsymbol{\theta}} \frac{1}{n-1} \sum_{i=1, i \neq j}^n \left( y_i - m(\mathbf{x}_i, \hat{p}_{n,i}^{(j)}, \boldsymbol{\theta}) \right)^2.$$

Define  $\hat{\boldsymbol{\theta}}^{(\cdot)} = \frac{1}{n} \sum_{j=1}^n \hat{\boldsymbol{\theta}}^{(j)}$ .

- **Step 2.** The jackknife bias estimator is defined as

$$\hat{\mathbf{B}} = (n-1) \cdot \sqrt{n} \left( \hat{\boldsymbol{\theta}}^{(\cdot)} - \hat{\boldsymbol{\theta}} \right) = \frac{n-1}{n} \sum_{j=1}^n \sqrt{n} \left( \hat{\boldsymbol{\theta}}^{(j)} - \hat{\boldsymbol{\theta}} \right),$$

and the bias corrected estimator is  $\hat{\boldsymbol{\theta}}_{bc} = \hat{\boldsymbol{\theta}} - \hat{\mathbf{B}}/\sqrt{n}$ .

- **Step 3.** The jackknife variance estimator is

$$\hat{\mathbf{V}} = (n-1) \sum_{j=1}^n \left( \hat{\boldsymbol{\theta}}^{(j)} - \hat{\boldsymbol{\theta}}^{(\cdot)} \right) \left( \hat{\boldsymbol{\theta}}^{(j)} - \hat{\boldsymbol{\theta}}^{(\cdot)} \right)'$$

**RECALL.** The MTE and functionals thereof are

$$\hat{\tau}_{\text{mte}} = \frac{\partial}{\partial a} m(\mathbf{x}, a, \hat{\boldsymbol{\theta}}) \quad \text{and} \quad \hat{\tau}_{\omega} = \int_0^1 \hat{\tau}_{\text{mte}}(a|\mathbf{x}) \omega(a|\mathbf{x}) da$$