

Shift-Share Instrumental Variables (SSIV)

Based on Peter Hull's Lecture Notes

Zhiyuan Chen

Empirical Methods

Renmin Business School

April 2025

Roadmap

- 1 What is Linear SSIV?
- 2 Shock Exogeneity
 - Motivation
 - Borusyak et al. (2022)
- 3 Share Exogeneity
 - Motivation
 - Goldsmith-Pinkham et al. (2020)
- 4 Choosing an Appropriate Framework
- 5 *Recentering Method by Borusyak and Hull (2023, ECTA)
 - Motivation & Intuition
 - Formal Framework
 - Applications

Roadmap

- 1 What is Linear SSIV?
- 2 Shock Exogeneity
 - Motivation
 - Borusyak et al. (2022)
- 3 Share Exogeneity
 - Motivation
 - Goldsmith-Pinkham et al. (2020)
- 4 Choosing an Appropriate Framework
- 5 *Recentering Method by Borusyak and Hull (2023, ECTA)
 - Motivation & Intuition
 - Formal Framework
 - Applications

What is a (Linear) SSIV?

A weighted sum of a common set of shocks, with weights reflecting heterogeneous exposure shares:

$$z_\ell = \sum_n s_{\ell n} g_n$$

What is a (Linear) SSIV?

A weighted sum of a common set of shocks, with weights reflecting heterogeneous exposure shares :

$$z_\ell = \sum_n s_{\ell n} g_n$$

- The shocks vary at a different “level” $n = 1, \dots, N$ than the shares $\ell = 1, \dots, L$, where we also observe an outcome y_ℓ & treatment x_ℓ

What is a (Linear) SSIV?

A weighted sum of a common set of shocks, with weights reflecting heterogeneous exposure shares :

$$z_\ell = \sum_n s_{\ell n} g_n$$

- The shocks vary at a different “level” $n = 1, \dots, N$ than the shares $\ell = 1, \dots, L$, where we also observe an outcome y_ℓ & treatment x_ℓ

We want to use z_ℓ to estimate parameter β_ℓ of the model $y_\ell = \beta_\ell x_\ell + \varepsilon_\ell$

What is a (Linear) SSIV?

A weighted sum of a common set of shocks, with weights reflecting heterogeneous exposure shares :

$$z_\ell = \sum_n s_{\ell n} g_n$$

- The shocks vary at a different “level” $n = 1, \dots, N$ than the shares $\ell = 1, \dots, L$, where we also observe an outcome y_ℓ & treatment x_ℓ

We want to use z_ℓ to estimate parameter β_ℓ of the model $y_\ell = \beta_\ell x_\ell + \varepsilon_\ell$

- Could be a “structural” equation or a potential outcomes model
- Could be misspecified, with heterogeneous treatment effects β_ℓ
- Could be a “reduced form” analysis, with $x_\ell = z_\ell$
- Could have other included controls w_ℓ

What is a (Linear) SSIV?

A weighted sum of a common set of shocks, with weights reflecting heterogeneous exposure shares :

$$z_\ell = \sum_n s_{\ell n} g_n$$

- The shocks vary at a different “level” $n = 1, \dots, N$ than the shares $\ell = 1, \dots, L$, where we also observe an outcome y_ℓ & treatment x_ℓ

We want to use z_ℓ to estimate parameter β_ℓ of the model $y_\ell = \beta_\ell x_\ell + \varepsilon_\ell$

- Could be a “structural” equation or a potential outcomes model
- Could be misspecified, with heterogeneous treatment effects β_ℓ
- Could be a “reduced form” analysis, with $x_\ell = z_\ell$
- Could have other included controls w_ℓ

Key question: Under what assumptions does this SSIV strategy “work”?

SSIV Examples

$$\text{Instrument } z_\ell = \sum_n \overset{\text{shares}}{\underbrace{s_{\ell n}}} \overset{\text{shocks}}{\underbrace{g_n}} \quad \text{for model } y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$$

Bartik (1991); Blanchard and Katz (1992):

- β = inverse local labor supply elasticity
- x_ℓ and y_ℓ = employment and wage growth in region ℓ
- Need a labor demand shifter as an IV

SSIV Examples

$$\text{Instrument } z_\ell = \sum_n \overset{\text{shares}}{\underbrace{s_{\ell n}}} \overset{\text{shocks}}{\underbrace{g_n}} \text{ for model } y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$$

Bartik (1991); Blanchard and Katz (1992):

- β = inverse local labor supply elasticity
- x_ℓ and y_ℓ = employment and wage growth in region ℓ
- Need a labor demand shifter as an IV
- g_n = national growth of industry n
- $s_{\ell n}$ = lagged employment shares (of industry in a region)
- z_ℓ = predicted employment growth due to national industry trends

SSIV Examples

$$\text{Instrument } z_\ell = \sum_n \overset{\text{shares}}{\boxed{s_{\ell n}}} \overset{\text{shocks}}{\boxed{g_n}} \quad \text{for model } y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$$

“Enclave instrument”, e.g., Card (2009)

- β = inverse elasticity of substitution between native and immigrant labor of some skill level (need a relative labor supply instrument)
- x_ℓ and y_ℓ = relative employment and wage in region ℓ
- g_n = national immigration growth from origin country n
- $s_{\ell n}$ = lagged shares of migrants from origin n in region ℓ
- z_ℓ = share of migrants predicted from enclaves & recent growth

SSIV Examples

$$\text{Instrument } z_\ell = \sum_n \overset{\text{shares}}{\underbrace{s_{\ell n}}} \overset{\text{shocks}}{\underbrace{g_n}} \quad \text{for model } y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$$

Hummels et al. (2014) on offshoring:

- β = effect of imports on wages
- x_ℓ = imports by Danish firm ℓ , y_ℓ = wages
- g_n = changes in transport costs by n = (product, country)
- $s_{\ell n}$ = lagged import shares
- z_ℓ = predicted change in firm inputs via transport costs

What Do We Do With This?

Of course, we can always run IV with such z_ℓ ... but what does the corresponding estimand *identify*?

What Do We Do With This?

Of course, we can always run IV with such z_ℓ ... but what does the corresponding estimand *identify*?

Recall IV validity condition: $E\left[\frac{1}{L}\sum_{\ell} z_{\ell}\epsilon_{\ell}\right] = 0$ for model residual ϵ_{ℓ}

- Looks a little different than normal because we're not assuming *i.i.d.* sampling, i.e.,

$$E\left[\frac{1}{L}\sum_{\ell} z_{\ell}\epsilon_{\ell}\right] = E[z_{\ell}\epsilon_{\ell}]$$

What Do We Do With This?

Of course, we can always run IV with such z_ℓ ... but what does the corresponding estimand *identify*?

Recall IV validity condition: $E\left[\frac{1}{L}\sum_{\ell} z_{\ell}\epsilon_{\ell}\right] = 0$ for model residual ϵ_{ℓ}

- Looks a little different than normal because we're not assuming *i.i.d.* sampling, i.e.,

$$E\left[\frac{1}{L}\sum_{\ell} z_{\ell}\epsilon_{\ell}\right] = E[z_{\ell}\epsilon_{\ell}]$$

What properties of shocks and shares make this condition hold?

- Is SSIV like a natural experiment? A diff-in-diff? Something new?
- Since z_{ℓ} combines multiple sources of variation, it can be difficult to think about it being randomly assigned across ℓ (unlike a lottery IV)

Roadmap

1 What is Linear SSIV?

2 Shock Exogeneity

- Motivation
- Borusyak et al. (2022)

3 Share Exogeneity

- Motivation
- Goldsmith-Pinkham et al. (2020)

4 Choosing an Appropriate Framework

5 *Recentering Method by Borusyak and Hull (2023, ECTA)

- Motivation & Intuition
- Formal Framework
- Applications

Exogenous Shocks in Industry-Level Regressions

Acemoglu-Autor-Dorn-Hanson-Price (AADHP, 2016) look at the effects of import competition with China on US manufacturing *industries*:

$$\Delta \log Emp_{nt} = \alpha + \beta \Delta IP_{nt} + \varepsilon_{nt},$$

where ΔIP_{nt} measures growth in import penetration from China in industry n , and ε_{nt} captures industry demand/productivity shocks

Exogenous Shocks in Industry-Level Regressions

Acemoglu-Autor-Dorn-Hanson-Price (AADHP, 2016) look at the effects of import competition with China on US manufacturing *industries*:

$$\Delta \log Emp_{nt} = \alpha + \beta \Delta IP_{nt} + \varepsilon_{nt},$$

where ΔIP_{nt} measures growth in import penetration from China in industry n , and ε_{nt} captures industry demand/productivity shocks

Two Key Problems with OLS estimation:

- ① Endogeneity of ΔIP_{nt} : OLS is not consistent for β
- ② GE spillovers: β does not capture aggregate effects

Problem 1: Endogeneity of ΔIP_{nt}

$$\Delta \log Emp_{nt} = \alpha + \beta \Delta IP_{nt} + \varepsilon_{nt}$$

ΔIP_{nt} is driven by productivity shocks in China, but also potentially by productivity and demand shocks in the US

- ε_{nt} captures productivity and demand shocks in the US

Problem 1: Endogeneity of ΔIP_{nt}

$$\Delta \log Emp_{nt} = \alpha + \beta \Delta IP_{nt} + \varepsilon_{nt}$$

ΔIP_{nt} is driven by productivity shocks in China, but also potentially by productivity and demand shocks in the US

- ε_{nt} captures productivity and demand shocks in the US

AADHP instrument ΔIP_{nt} with ΔIPO_{nt} , measuring average Chinese import penetration growth in 8 non-US countries

Problem 1: Endogeneity of ΔIP_{nt}

$$\Delta \log Emp_{nt} = \alpha + \beta \Delta IP_{nt} + \varepsilon_{nt}$$

ΔIP_{nt} is driven by productivity shocks in China, but also potentially by productivity and demand shocks in the US

- ε_{nt} captures productivity and demand shocks in the US

AADHP instrument ΔIP_{nt} with ΔIPO_{nt} , measuring average Chinese import penetration growth in 8 non-US countries

- **Relevance:** both ΔIP_{nt} and ΔIPO_{nt} are driven by the same Chinese productivity shocks
- **Validity:** local productivity/demand shocks in the US are uncorrelated with those of other countries (entering ΔIPO_{nt})

Identification from a Natural Experiment

Suppose ΔIPO_{nt} is as-good-as-randomly assigned, as in a RCT:

$$E[\Delta IPO_{nt} \mid \mathbb{I}] = \mu \quad \text{for all } n, t$$

where $\mathbb{I} = \{\epsilon_{nt}, \text{pre-trends, balance variables}, \dots\}$

Identification from a Natural Experiment

Suppose ΔIPO_{nt} is as-good-as-randomly assigned, as in a RCT:

$$E[\Delta IPO_{nt} \mid \mathbb{I}] = \mu \quad \text{for all } n, t$$

where $\mathbb{I} = \{\epsilon_{nt}, \text{pre-trends, balance variables}, \dots\}$

Consistent IV estimation then follows from many observations of nt , with sufficiently independent variation in ΔIPO_{nt}

Identification from a Natural Experiment

Can relax to add observables capturing systematic variation:

$$E[\Delta IPO_{nt} \mid \mathbb{I}] = q'_{nt}\mu \quad \text{for all } n, t$$

where q_{nt} may include:

- period FE, isolating within-period variation in the shocks
- FE of 10 broad sectors, isolating within-sector variation, etc.

Identification from a Natural Experiment

Can relax to add observables capturing systematic variation:

$$E[\Delta IPO_{nt} \mid \mathbb{I}] = q'_{nt}\mu \quad \text{for all } n, t$$

where q_{nt} may include:

- period FE, isolating within-period variation in the shocks
- FE of 10 broad sectors, isolating within-sector variation, etc.

We would then just want to control for q_{nt} in the industry-level IV

Problem 2: GE Spillovers

Spillovers across different industries are likely important:

- When employment shrinks in industry n after a negative shock, aggregate employment may or may not respond

Problem 2: GE Spillovers

Spillovers across different industries are likely important:

- When employment shrinks in industry n after a negative shock, aggregate employment may or may not respond
- In a flexible labor market, comparing wages of similar workers across industries does not make sense

Problem 2: GE Spillovers

ADH Solution: specify the outcome equation for local labor markets

- Works if local economies are isolated “islands”
(simple model in Adao-Kolesar-Morales 2019; richer structure of spatial spillovers in Adao-Arkolakis-Esposito 2020)

Problem 2: GE Spillovers

ADH Solution: specify the outcome equation for local labor markets

- Works if local economies are isolated “islands”
(simple model in Adao-Kolesar-Morales 2019; richer structure of spatial spillovers in Adao-Arkolakis-Esposito 2020)

But correct specification is not the same as identification!

- **Key point:** the same industry-level natural experiment can be used to estimate a regional specification, via SSIV

Borusyak, Hull, and Jaravel (BHJ; 2022)

Consider the SSIV estimator of $y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$ instrumented by $z_\ell = \sum_n s_{\ell n} g_n$ and, for now, $\sum_n s_{\ell n} = 1$ for all ℓ

- Reduced-form allowed: $x_\ell = z_\ell$
- Only the shift-share structure of z_ℓ matters; x_ℓ can be anything
- Note: view g_n as stochastic, so can't assume z_ℓ is iid

Borusyak, Hull, and Jaravel (BHJ; 2022)

Consider the SSIV estimator of $y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$ instrumented by $z_\ell = \sum_n s_{\ell n} g_n$ and, for now, $\sum_n s_{\ell n} = 1$ for all ℓ

- Reduced-form allowed: $x_\ell = z_\ell$
- Only the shift-share structure of z_ℓ matters; x_ℓ can be anything
- Note: view g_n as stochastic, so can't assume z_ℓ is iid

E.g., $g_n = \Delta IPO_n$ aggregated w/mfg employment shares $s_{\ell n}$

- Can we leverage a natural experiment in g_n , as before?

Leveraging g_n

Shift-Share Estimand

Consider the SSIV estimator of $y_\ell = \beta x_\ell + \gamma' w_\ell + \varepsilon_\ell$ instrumented by $z_\ell = \sum_n s_{\ell n} g_n$ and, for now, $\sum_n s_{\ell n} = 1$ for all ℓ

First step: note that by the Frisch–Waugh–Lovell theorem, the estimator can be written

$$\hat{\beta} = \frac{\sum_\ell z_\ell \tilde{y}_\ell}{\sum_\ell z_\ell \tilde{x}_\ell} = \frac{\sum_\ell \sum_n s_{\ell n} g_n \tilde{y}_\ell}{\sum_\ell \sum_n s_{\ell n} g_n \tilde{x}_\ell}$$

where \tilde{v}_ℓ denotes sample residuals from regressing v_ℓ on w_ℓ . [Proof](#)

Leveraging g_n

BHJ Numerical Equivalence

BHJ show $\hat{\beta}$ can be obtained from a shock-level IV procedure that uses g_n to instrument for a shock-level “aggregate” of the treatment:

Leveraging g_n

BHJ Numerical Equivalence

BHJ show $\hat{\beta}$ can be obtained from a shock-level IV procedure that uses g_n to instrument for a shock-level “aggregate” of the treatment:

$$\hat{\beta} = \frac{\frac{1}{L} \sum_{\ell} \sum_n s_{\ell n} g_n \tilde{y}_{\ell}}{\frac{1}{L} \sum_{\ell} \sum_n s_{\ell n} g_n \tilde{x}_{\ell}} =$$

Leveraging g_n

BHJ Numerical Equivalence

BHJ show $\hat{\beta}$ can be obtained from a shock-level IV procedure that uses g_n to instrument for a shock-level “aggregate” of the treatment:

$$\hat{\beta} = \frac{\frac{1}{L} \sum_{\ell} \sum_n s_{\ell n} g_n \tilde{y}_{\ell}}{\frac{1}{L} \sum_{\ell} \sum_n s_{\ell n} g_n \tilde{x}_{\ell}} = \frac{\sum_n g_n \sum_{\ell} \frac{1}{L} s_{\ell n} \tilde{y}_{\ell}}{\sum_n g_n \sum_{\ell} \frac{1}{L} s_{\ell n} \tilde{x}_{\ell}} =$$

Leveraging g_n

BHJ Numerical Equivalence

BHJ show $\hat{\beta}$ can be obtained from a shock-level IV procedure that uses g_n to instrument for a shock-level “aggregate” of the treatment:

$$\hat{\beta} = \frac{\frac{1}{L} \sum_{\ell} \sum_n s_{\ell n} g_n \tilde{y}_{\ell}}{\frac{1}{L} \sum_{\ell} \sum_n s_{\ell n} g_n \tilde{x}_{\ell}} = \frac{\sum_n g_n \sum_{\ell} \frac{1}{L} s_{\ell n} \tilde{y}_{\ell}}{\sum_n g_n \sum_{\ell} \frac{1}{L} s_{\ell n} \tilde{x}_{\ell}} = \frac{\sum_n s_n g_n \bar{\tilde{y}}_n}{\sum_n s_n g_n \bar{\tilde{x}}_n},$$

where $s_n = \frac{1}{L} \sum_{\ell} s_{\ell n}$ are weights capturing the average importance of shock n , and $\bar{v}_n = \frac{\sum_{\ell} s_{\ell n} v_{\ell}}{\sum_{\ell} s_{\ell n}}$ is an exposure-weighted average of v_{ℓ}

Leveraging g_n

BHJ Numerical Equivalence

$$\hat{\beta} = \frac{\sum_n s_n g_n \tilde{y}_n}{\sum_n s_n g_n \tilde{x}_n}$$

The IV estimate from the original “location-level” IV procedure is equivalent to a “industry-level” IV regression with model $\tilde{y}_n = \alpha + \tilde{x}_n \beta + \bar{\epsilon}_n$ instrumented by g_n with weights s_n .

The residual $\bar{\epsilon}_n$ of this shock-level IV procedure is the average residual of observations with a high share of n

- E.g. in ADH, the average unobserved determinants of regional employment in regions most specialized in industry n

Leveraging g_n

BHJ Numerical Equivalence

$$\hat{\beta} = \frac{\sum_n s_n g_n \tilde{y}_n}{\sum_n s_n g_n \tilde{x}_n}$$

The IV estimate from the original “location-level” IV procedure is equivalent to a “industry-level” IV regression with model $\tilde{y}_n = \alpha + \tilde{x}_n \beta + \bar{\varepsilon}_n$ instrumented by g_n with weights s_n .

The residual $\bar{\varepsilon}_n$ of this shock-level IV procedure is the average residual of observations with a high share of n

- E.g. in ADH, the average unobserved determinants of regional employment in regions most specialized in industry n

It follows that $\hat{\beta}$ is consistent iff this shock-level IV procedure is...

BHJ Baseline Assumptions

A1 (Quasi-random shock assignment): $E[g_n \mid \bar{\epsilon}, s] = \mu$, for all n

- Each shock has the same expected value, conditional on the shock-level unobservables $\bar{\epsilon}_n$ and average exposure s_n

BHJ Baseline Assumptions

A1 (Quasi-random shock assignment): $E[g_n \mid \bar{\epsilon}, s] = \mu$, for all n

- Each shock has the same expected value, conditional on the shock-level unobservables $\bar{\epsilon}_n$ and average exposure s_n
- Implies SSIV exogeneity, as $z_\ell = \mu + \sum_n s_{\ell n}(g_n - \mu) = \mu + \text{“noise”}$

BHJ Baseline Assumptions

A2 (Many uncorrelated shocks):

- $E[\sum_n s_n^2] \rightarrow 0$: expected Herfindahl index of average shock exposure converges to zero (implies $N \rightarrow \infty$)
- $Cov(g_n, g_{n'} | \bar{\epsilon}, s) = 0$ for all $n' \neq n$: shocks are mutually uncorrelated given the unobservables

BHJ Baseline Assumptions

A2 (Many uncorrelated shocks):

- $E[\sum_n s_n^2] \rightarrow 0$: expected Herfindahl index of average shock exposure converges to zero (implies $N \rightarrow \infty$)
- $Cov(g_n, g_{n'} | \bar{\epsilon}, s) = 0$ for all $n' \neq n$: shocks are mutually uncorrelated given the unobservables
- Imply a shock-level law of large numbers: $\sum_n s_n g_n \bar{\epsilon}_n \xrightarrow{p} 0$

BHJ Baseline Assumptions

A2 (Many uncorrelated shocks):

- $E[\sum_n s_n^2] \rightarrow 0$: expected Herfindahl index of average shock exposure converges to zero (implies $N \rightarrow \infty$)
- $Cov(g_n, g_{n'} | \bar{\epsilon}, s) = 0$ for all $n' \neq n$: shocks are mutually uncorrelated given the unobservables
- Imply a shock-level law of large numbers: $\sum_n s_n g_n \bar{\epsilon}_n \xrightarrow{p} 0$

Both assumptions, while novel for SSIV, would be standard for a shock-level IV regression with weights s_n and instrument g_n

BHJ Extensions

Conditional Quasi-Random Assignment: $E[g_n \mid \bar{\varepsilon}, q, s] = q'_n \mu$ for some observed shock-level variables q_n

- Consistency follows when $w_\ell = \sum_n s_{\ell n} q_n$ is controlled for in the IV

BHJ Extensions

Conditional Quasi-Random Assignment: $E[g_n \mid \bar{\varepsilon}, q, s] = q'_n \mu$ for some observed shock-level variables q_n

- Consistency follows when $w_\ell = \sum_n s_{\ell n} q_n$ is controlled for in the IV

Weakly Mutually Correlated Shocks: $g_n \mid (\bar{\varepsilon}, q, s)$ are clustered or otherwise mutually dependent

- Consistency follows when mutual correlation is not too strong

BHJ Extensions

Conditional Quasi-Random Assignment: $E[g_n \mid \bar{\varepsilon}, q, s] = q'_n \mu$ for some observed shock-level variables q_n

- Consistency follows when $w_\ell = \sum_n s_{\ell n} q_n$ is controlled for in the IV

Weakly Mutually Correlated Shocks: $g_n \mid (\bar{\varepsilon}, q, s)$ are clustered or otherwise mutually dependent

- Consistency follows when mutual correlation is not too strong

Estimated Shocks: $g_n = \sum_\ell w_{\ell n} g_{\ell n}$ proxies for an infeasible g_n^*

- Consistency may require a “leave-out” adjustment: $z_\ell = \sum_n s_{\ell n} \tilde{g}_{\ell n}$ for $\tilde{g}_{\ell n} = \sum_{\ell' \neq \ell} \omega_{\ell' n} g_{\ell' n}$ (akin to JIVE solution to many-IV bias)

BHJ Extensions (cont.)

Panel Data: Have $(y_{\ell t}, x_{\ell t}, s_{\ell nt}, g_{nt})$ across $\ell = 1, \dots, L$, $t = 1, \dots, T$

- Consistency can follow from either $N \rightarrow \infty$ or $T \rightarrow \infty$
- Unit fixed effects “de-mean” the shocks, if $s_{\ell nt}$ are time-invariant

BHJ Extensions (cont.)

Panel Data: Have $(y_{\ell t}, x_{\ell t}, s_{\ell nt}, g_{nt})$ across $\ell = 1, \dots, L$, $t = 1, \dots, T$

- Consistency can follow from either $N \rightarrow \infty$ or $T \rightarrow \infty$
- Unit fixed effects “de-mean” the shocks, if $s_{\ell nt}$ are time-invariant

Heterogeneous Effects: LATE theorem logic goes through

- Under a first-stage monotonicity condition, SSIV identifies a convex weighted average of heterogeneous treatment effects

Practical Consideration 1: Incomplete Shares

The Problem

So far we have assumed a constant sum-of-shares: $S_\ell \equiv \sum_n s_{\ell n} = 1$

- But in some settings, S_ℓ varies across ℓ
- E.g. in ADH, S_ℓ is region ℓ 's share of non-manufacturing emp.

Practical Consideration 1: Incomplete Shares

The Problem

So far we have assumed a constant sum-of-shares: $S_\ell \equiv \sum_n s_{\ell n} = 1$

- But in some settings, S_ℓ varies across ℓ
- E.g. in ADH, S_ℓ is region ℓ 's share of non-manufacturing emp.

BHJ show that **A1/A2** are not enough for validity of z_ℓ in this case

- Now $z_\ell = \sum_n s_{\ell n}(\mu + (g_n - \mu)) = \mu S_\ell + \sum_n s_{\ell n}(g_n - \mu)$
- So z_ℓ is mechanically correlated with S_ℓ , which may be endogenous

E.g. in ADH, Comparing locations with larger and smaller z_ℓ could be comparing places with larger vs. smaller manufacturing employment (e.g. Midwest vs. South)

Practical Consideration 1: Incomplete Shares

The Solution

$$z_\ell = \sum_n s_{\ell n} (\mu + (g_n - \mu)) = \mu S_\ell + \underbrace{\sum_n s_{\ell n} (g_n - \mu)}_{\text{Clean Shock Variation}}$$

Controlling for the sum-of-shares S_ℓ isolates clean shock variation

Practical Consideration 1: Incomplete Shares

The Solution

$$z_{\ell} = \sum_n s_{\ell n} (\mu + (g_n - \mu)) = \mu S_{\ell} + \underbrace{\sum_n s_{\ell n} (g_n - \mu)}_{\text{Clean Shock Variation}}$$

Controlling for the sum-of-shares S_{ℓ} isolates clean shock variation

- Further controls are needed when **A1** only holds conditional on q_n ; e.g. in panels, S_{ℓ} should be interacted with time FE:

$$z_{\ell t} = \sum_n s_{\ell n} (\mu_t + (g_{nt} - \mu_t)) = \mu_t S_{\ell} + \underbrace{\sum_n s_{\ell n} (g_{nt} - \mu_t)}_{\text{Clean Shock Variation}}$$

Practical Consideration 2: Exposure Clustering

The Problem

Adão, Kolesar, and Morales (2019) study a novel inference challenge when SSIV identification leverages quasi-random shocks

- Observations with similar shares $s_{\ell 1}, \dots, s_{\ell N}$ are likely to have correlated z_{ℓ} , even when observations are not “clustered” in conventional ways (e.g., by distance)

Practical Consideration 2: Exposure Clustering

The Problem

Adão, Kolesar, and Morales (2019) study a novel inference challenge when SSIV identification leverages quasi-random shocks

- Observations with similar shares $s_{\ell 1}, \dots, s_{\ell N}$ are likely to have correlated z_{ℓ} , even when observations are not “clustered” in conventional ways (e.g., by distance)
- When ε_{ℓ} is similarly clustered (e.g. when $\varepsilon_{\ell} = \sum_n s_{\ell n} \mathbf{v}_n + \tilde{\varepsilon}_{\ell}$), large-sample distribution of $\hat{\beta}$ may not be well-approximated by standard central limit theorems (CLTs)

Practical Consideration 2: Exposure Clustering

The Problem

Adão, Kolesar, and Morales (2019) study a novel inference challenge when SSIV identification leverages quasi-random shocks

- Observations with similar shares $s_{\ell 1}, \dots, s_{\ell N}$ are likely to have correlated z_{ℓ} , even when observations are not “clustered” in conventional ways (e.g., by distance)
- When ε_{ℓ} is similarly clustered (e.g. when $\varepsilon_{\ell} = \sum_n s_{\ell n} \mathbf{v}_n + \tilde{\varepsilon}_{\ell}$), large-sample distribution of $\hat{\beta}$ may not be well-approximated by standard central limit theorems (CLTs)

They then derive a new CLT + SEs to address “exposure clustering”

- “Design-based”: leverage *iid*ness of shocks, not observations

Practical Consideration 2: Exposure Clustering

The Solution

BHJ use similar logic to show robust/clustered SEs can be valid when $\hat{\beta}$ is given by estimating the ‘industry-level’ regression

$$\bar{y}_n = \alpha + \beta \bar{x}_n + q'_n \tau + \bar{\varepsilon}_n^\perp,$$

instrumenting \bar{x}_n by g_n and weighting by s_n

Practical Consideration 2: Exposure Clustering

The Solution

BHJ use similar logic to show robust/clustered SEs can be valid when $\hat{\beta}$ is given by estimating the ‘industry-level’ regression

$$\bar{y}_n = \alpha + \beta \bar{x}_n + q'_n \tau + \bar{\varepsilon}_n^\perp,$$

instrumenting \bar{x}_n by g_n and weighting by s_n

- Numerically identical IV estimate, when controls include $\sum_n s_{\ell n} q_n$
- Clustering logic: valid SEs are obtained when estimating the IV at the level of identifying variation (here, shocks)

Practical Consideration 2: Exposure Clustering

The Solution

BHJ use similar logic to show robust/clustered SEs can be valid when $\hat{\beta}$ is given by estimating the ‘industry-level’ regression

$$\bar{y}_n = \alpha + \beta \bar{x}_n + q'_n \tau + \bar{\varepsilon}_n^\perp,$$

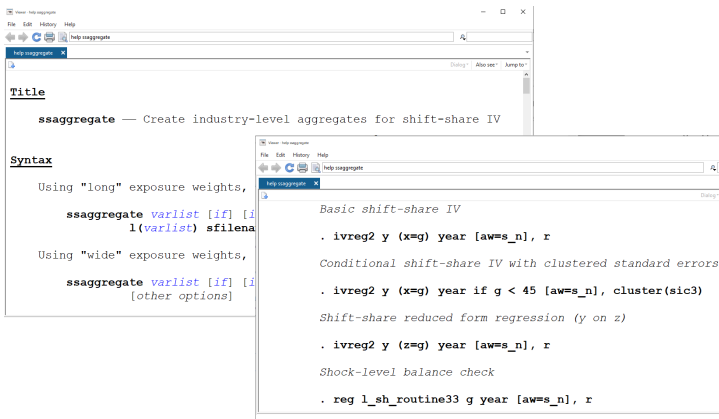
instrumenting \bar{x}_n by g_n and weighting by s_n

- Numerically identical IV estimate, when controls include $\sum_n s_{\ell n} q_n$
- Clustering logic: valid SEs are obtained when estimating the IV at the level of identifying variation (here, shocks)

Same logic applies to performing valid balance/pre-trend tests and evaluating first-stage strength of the instrument

SSIV with ssaggregate

Stata package **ssaggregate** leverages the BHJ equivalence result: it translates data to the shock level, after which researchers can proceed with familiar estimation commands (install w/ `ssc install ssaggregate`)



SSIV with ssaggregate...in R!

ssaggregate is now available in R!

kylebutts / ssaggregate Public

Watch 1 Fork 0 Star 0

Code Issues Pull requests Actions Projects Wiki Security Insights

main Go to file Add file Code About

kylebutts Fix publication year 3 days ago 3

R	Fix publication year	3 days ago
data-raw	Initial ssaggregate implementation	3 days ago
data	Initial ssaggregate implementation	3 days ago
inst	Initial ssaggregate implementation	3 days ago
man	Fix publication year	3 days ago

README.md

ssaggregate

ssaggregate converts "location-level" variables in a shift-share IV dataset to a dataset of exposure-weighted "industry-level" aggregates, as described in [Borusyak, Hull, and Jaravel \(2022\)](#).

Details

There are two ways to specify ssaggregate, depending on whether the industry exposure weights are saved in "long" format (unique rows for industry x location) in a separate dataset `shares` or in "wide" format (unique rows for location and columns for each industry) as part of `df`. In general `ssaggregate` will execute faster with "long" exposure weights. See the examples for proper syntax in both cases.

Create industry-level aggregates for shift-share IV following Borusyak, Hull, and Jaravel (2022)

Readme View license

0 stars 1 watching 0 forks

Releases

No releases published

Packages

No packages published

Languages

R 100.0%

Download at <https://github.com/kylebutts/ssaggregate>

Application: “The China Shock”

ADH study the effects of rising Chinese import competition on US commuting zones, 1991-2000 and 2000-2007

- Treatment x_ℓ : local growth of Chinese imports in \$1,000/worker (slightly different from AADHP and ADHS)
- Main outcome y_ℓ : local change in manufacturing emp. share

Application: “The China Shock”

ADH study the effects of rising Chinese import competition on US commuting zones, 1991-2000 and 2000-2007

- Treatment x_ℓ : local growth of Chinese imports in \$1,000/worker (slightly different from AADHP and ADHS)
- Main outcome y_ℓ : local change in manufacturing emp. share

To address endogeneity challenge, use a SSIV $z_{\ell t} = \sum_n s_{\ell nt} g_{nt}$

- n : 397 SIC4 manufacturing industries (\times 2 periods)
- g_{nt} : growth of Chinese imports in non-US economies per US worker
- $s_{\ell nt}$: lagged share of mfg. industry n in *total* emp. of location ℓ

ADH Revisited

BHJ show how ADH can be seen as leveraging quasi-random shocks

- *Ex ante* plausible: imagine random industry productivity shocks in China affecting imports in U.S. & elsewhere

ADH Revisited

Plausability of $A1/A2$

Evaluate **A1** by regional and industry-level balance tests

- Industry shocks are uncorrelated with observables

ADH Revisited

Plausability of A1/A2

Evaluate **A1** by regional and industry-level balance tests

- Industry shocks are uncorrelated with observables

Check sensitivity to adjusting for potential industry-level confounders:

- Control for $w_{\ell t} = \sum_n s_{\ell nt} q_{nt}$, where q_{nt} include period FE, sector FE, the Acemoglu et al. (2016) observables, ...

ADH Revisited

Plausability of A1/A2

Evaluate **A1** by regional and industry-level balance tests

- Industry shocks are uncorrelated with observables

Check sensitivity to adjusting for potential industry-level confounders:

- Control for $w_{\ell t} = \sum_n s_{\ell nt} q_{nt}$, where q_{nt} include period FE, sector FE, the Acemoglu et al. (2016) observables, ...

Evaluate **A2** by studying variation across industries

- Effective sample size (1/HHI of s_n weights): 58-192
- Shocks appear mutually uncorrelated across SIC3 sectors

BHJ do ADH: Shock-Level Balance

Table 3: Shock Balance Tests in the Autor et al. (2013) Setting

Balance variable	Coef.	SE
Production workers' share of employment, 1991	-0.011	(0.012)
Ratio of capital to value-added, 1991	-0.007	(0.019)
Log real wage (2007 USD), 1991	-0.005	(0.022)
Computer investment as share of total, 1990	0.750	(0.465)
High-tech equipment as share of total investment, 1990	0.532	(0.296)
# of industry-periods	794	

No significant correlations between shocks and industry observables,
controlling for year fixed effects

BHJ do ADH: Manufacturing Employment

Table 4: Shift-Share IV Estimates of the Effect of Chinese Imports on Manufacturing Employment

	(1)	(2)	(3)	(4)	(5)	(6)	(7)
Coefficient	-0.596	-0.489	-0.267	-0.314	-0.310	-0.290	-0.432
	(0.114)	(0.100)	(0.099)	(0.107)	(0.134)	(0.129)	(0.205)
<u>Regional controls</u>							
Autor et al. (2013) controls	✓	✓	✓		✓	✓	✓
Start-of-period mfg. share	✓						
Lagged mfg. share		✓	✓	✓	✓	✓	✓
Period-specific lagged mfg. share			✓	✓	✓	✓	✓
Lagged 10-sector shares					✓		✓
Local Acemoglu et al. (2016) controls						✓	
Lagged industry shares							✓
SSIV first stage F -stat.	185.6	166.7	123.6	272.4	64.6	63.3	27.6
# of region-periods	1,444	1,444	1,444	1,444	1,444	1,444	1,444
# of industry-periods	796	794	794	794	794	794	794

Roadmap

- 1 What is Linear SSIV?
- 2 Shock Exogeneity
 - Motivation
 - Borusyak et al. (2022)
- 3 Share Exogeneity
 - Motivation
 - Goldsmith-Pinkham et al. (2020)
- 4 Choosing an Appropriate Framework
- 5 *Recentering Method by Borusyak and Hull (2023, ECTA)
 - Motivation & Intuition
 - Formal Framework
 - Applications

The Mariel Boatlift as a Basic SSIV

Card (1990) leverages a big migration “push” of low-skilled workers from Cuba to Miami, a Cuban-enclave.

The Mariel Boatlift as a Basic SSIV

Card (1990) leverages a big migration “push” of low-skilled workers from Cuba to Miami, a Cuban-enclave. Imagine instrumenting immigrant inflows by the lagged share of Cuban workers $s_{\ell, \text{Cuba}}$ in a diff-in-diff setup

- Need parallel trends: regions with more/fewer Cuban workers on similar employment trends

This can be viewed as a simple shift-share instrument:

$$s_{\ell, \text{Cuba}} \equiv s_{\ell, \text{Cuba}} \cdot 1 + \sum_{n \neq \text{Cuba}} s_{\ell n} \cdot 0$$

The Mariel Boatlift as a Basic SSIV

Card (1990) leverages a big migration “push” of low-skilled workers from Cuba to Miami, a Cuban-enclave. Imagine instrumenting immigrant inflows by the lagged share of Cuban workers $s_{\ell, \text{Cuba}}$ in a diff-in-diff setup

- Need parallel trends: regions with more/fewer Cuban workers on similar employment trends

This can be viewed as a simple shift-share instrument:

$$s_{\ell, \text{Cuba}} \equiv s_{\ell, \text{Cuba}} \cdot 1 + \sum_{n \neq \text{Cuba}} s_{\ell n} \cdot 0$$

If several migration origins had a push shock, we can pool them together with a more traditional SSIV...

Goldsmith-Pinkham, Sorkin, and Swift (GPSS; 2020)

GPSS view the set of n and values of g_n as fixed, so $z_\ell = \sum_n s_{\ell n} g_n$ is a linear combination of shares

Goldsmith-Pinkham, Sorkin, and Swift (GPSS; 2020)

GPSS view the set of n and values of g_n as fixed, so $z_\ell = \sum_n s_{\ell n} g_n$ is a linear combination of shares

They then also establish a numerical equivalence: $\hat{\beta}$ can be obtained from an overidentified IV procedure that uses N share instruments $s_{\ell n}$ and a weight matrix based on the shocks g_n

Goldsmith-Pinkham, Sorkin, and Swift (GPSS; 2020)

Sufficient identifying assumption: shares $s_{\ell n}$ are exogenous for each n (like parallel trends when ε_{ℓ} are unobserved trends)

$$E[\varepsilon_{\ell} \mid s_{\ell n}] = 0, \forall n$$

Goldsmith-Pinkham, Sorkin, and Swift (GPSS; 2020)

Sufficient identifying assumption: shares $s_{\ell n}$ are exogenous for each n (like parallel trends when ε_{ℓ} are unobserved trends)

$$E[\varepsilon_{\ell} \mid s_{\ell n}] = 0, \forall n \implies E[\sum_{\ell} z_{\ell} \varepsilon_{\ell}] = \sum_{\ell} \sum_n g_n E[s_{\ell n}] E[\varepsilon_{\ell} \mid s_{\ell n}] = 0$$

This is N moment conditions at the level of observations, e.g. 38 for Card and 397 for ADH (vs. just 1 in BHJ, at the level of industries)

In other words, GPSS show that the SSIV estimator can be seen as pooling many Boatlift-style diff-in-diff IVs, one for each industry

Rotemberg Weights

How does SSIV pool different diff-in-diffs?

- GPSS propose “opening the black box” of overidentified IV by deriving the weights SSIV implicitly puts on each share instrument
- Builds on Rotemberg (1983), so they call these “Rotemberg weights”

$$\hat{\beta} = \sum_n \hat{\alpha}_n \hat{\beta}_n, \text{ where } \underbrace{\hat{\beta}_n = \frac{\sum_{\ell} s_{\ell n} \tilde{y}_{\ell}}{\sum_{\ell} s_{\ell n} \tilde{x}_{\ell}}}_{n\text{-specific IV estimate}} \text{ and } \underbrace{\hat{\alpha}_n = \frac{g_n \sum_{\ell} s_{\ell n} \tilde{x}_{\ell}}{\sum_{n'} g_{n'} \sum_{\ell} s_{\ell n'} \tilde{x}_{\ell}}}_{\text{Rotemberg weight}}$$

Rotemberg Weights

How does SSIV pool different diff-in-diffs?

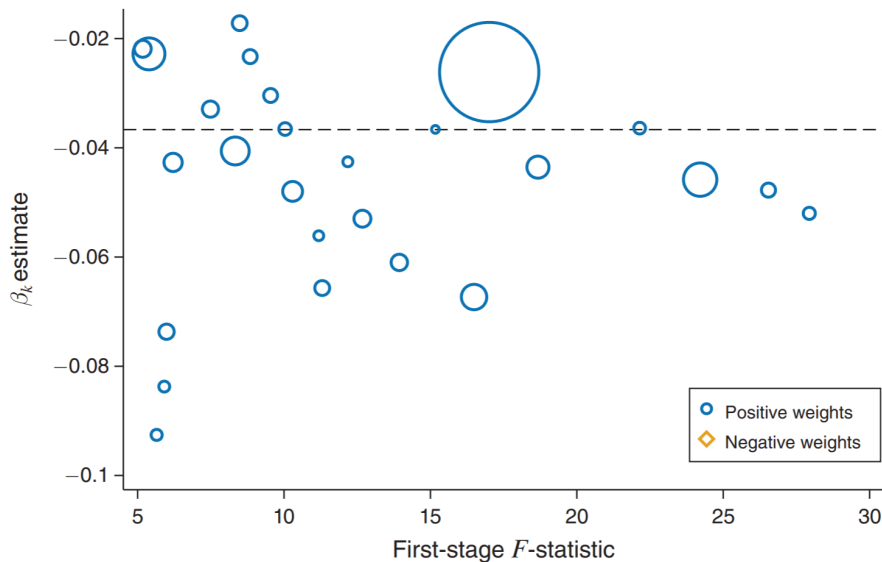
- GPSS propose “opening the black box” of overidentified IV by deriving the weights SSIV implicitly puts on each share instrument
- Builds on Rotemberg (1983), so they call these “Rotemberg weights”

$$\hat{\beta} = \sum_n \hat{\alpha}_n \hat{\beta}_n, \text{ where } \underbrace{\hat{\beta}_n = \frac{\sum_{\ell} s_{\ell n} \tilde{y}_{\ell}}{\sum_{\ell} s_{\ell n} \tilde{x}_{\ell}}}_{n\text{-specific IV estimate}} \quad \text{and} \quad \underbrace{\hat{\alpha}_n = \frac{g_n \sum_{\ell} s_{\ell n} \tilde{x}_{\ell}}{\sum_{n'} g_{n'} \sum_{\ell} s_{\ell n'} \tilde{x}_{\ell}}}_{\text{Rotemberg weight}}$$

Intuitively, more weight is given to share instruments with more extreme shocks g_n and larger first stages $\sum_{\ell} s_{\ell n} \tilde{x}_{\ell}$

- Weights can be negative (potential issue w/heterogeneous effects)

Rotemberg Weights in Card (2009)



Is Share Exogeneity Plausible?

Share exogeneity assumption is **not** that “shares don’t causally respond to the residual” (they can’t: shares are pre-determined)

- It’s: “all unobservables are uncorrelated with anything about the local share distribution”

Is Share Exogeneity Plausible?

This sufficient condition is typically violated when there are *any* unobserved shocks \mathbf{v}_n that affect $\boldsymbol{\varepsilon}_\ell$ via the same or correlated shares

- I.e. if $\boldsymbol{\varepsilon}_\ell = \sum_n s_{\ell n} \mathbf{v}_n + \tilde{\boldsymbol{\varepsilon}}_\ell$, then $s_{\ell n}$ and $\boldsymbol{\varepsilon}_\ell$ cannot be uncorrelated in large samples—even if \mathbf{v}_n are uncorelated with g_n
- E.g. in ADH, unobserved technology shocks across industries affect labor markets via lagged emp. shares, along with observed g_n
- Problem arises when shares are “generic” – predicting many things

Card and ADH Revisited

When share exogeneity is *ex ante* plausible, can test its assumptions *ex post* (focusing on high Rotemberg weight n):

- Balance/pre-trend tests
- Overidentification tests (under constant effects)
- Straightforward to implement; no different than any other IV

Card and ADH Revisited

When share exogeneity is *ex ante* plausible, can test its assumptions *ex post* (focusing on high Rotemberg weight n):

- Balance/pre-trend tests
- Overidentification tests (under constant effects)
- Straightforward to implement; no different than any other IV

GPSS find that balance/overidentification tests broadly pass for Card ... but fail badly for ADH, consistent with *ex ante* implausibility

Roadmap

- 1 What is Linear SSIV?
- 2 Shock Exogeneity
 - Motivation
 - Borusyak et al. (2022)
- 3 Share Exogeneity
 - Motivation
 - Goldsmith-Pinkham et al. (2020)
- 4 Choosing an Appropriate Framework
- 5 *Recentering Method by Borusyak and Hull (2023, ECTA)
 - Motivation & Intuition
 - Formal Framework
 - Applications

A Taxonomy of SSIV Settings

Case 1 the IV is based on a set of shocks which can be thought of as an instrument (i.e. many, plausibly quasi-randomly assigned)

- BHJ shows how this identifying variation can be mapped to estimate effects at a different “level” (i.e. industries \rightarrow local labor markets)

A Taxonomy of SSIV Settings

Case 1 the IV is based on a set of shocks which can be thought of as an instrument (i.e. many, plausibly quasi-randomly assigned)

- BHJ shows how this identifying variation can be mapped to estimate effects at a different “level” (i.e. industries \rightarrow local labor markets)

Case 2 the researcher does not directly observe many quasi-random shocks, but can estimate them in-sample

- Canonical setting of Bartik (1991), where g_n are average industry growth rates (thought to proxy for latent demand shocks)
- See also Card (2009), where national immigration rates are estimated

A Taxonomy of SSIV Settings

Case 1 the IV is based on a set of shocks which can be thought of as an instrument (i.e. many, plausibly quasi-randomly assigned)

- BHJ shows how this identifying variation can be mapped to estimate effects at a different “level” (i.e. industries \rightarrow local labor markets)

Case 2 the researcher does not directly observe many quasi-random shocks, but can estimate them in-sample

- Canonical setting of Bartik (1991), where g_n are average industry growth rates (thought to proxy for latent demand shocks)
- See also Card (2009), where national immigration rates are estimated

Case 3 the g_n cannot be naturally viewed as an instrument

- Either too few or implausibly exogenous, even given some q_n .
- Identification may (or may not) instead follow from share exogeneity

Ex Ante vs. Ex Post Validity

BHJ emphasize that the decision to pursue a “shocks” vs. “shares” identification strategy must be made *ex ante*

- Undesirable to base identifying assumptions on *ex post* tests, though balance/pre-trend tests can be used to falsify assumptions
- The two identification strategies have different economic content

Ex Ante vs. Ex Post Validity

BHJ emphasize that the decision to pursue a “shocks” vs. “shares” identification strategy must be made *ex ante*

- Undesirable to base identifying assumptions on *ex post* tests, though balance/pre-trend tests can be used to falsify assumptions
- The two identification strategies have different economic content

They suggest thinking about whether shares are “tailored” to the economic question/treatment, or are “generic”

- Generic shares (e.g. ADH): unobserved \mathbf{v}_n are likely to enter $\boldsymbol{\varepsilon}_\ell$ via the same or similar shares, violating share exogeneity
- Tailored shares have a diff-in-diff feel; don’t even need the shocks, except to possibly improve power or avoid many-IV bias

Roadmap

- 1 What is Linear SSIV?
- 2 Shock Exogeneity
 - Motivation
 - Borusyak et al. (2022)
- 3 Share Exogeneity
 - Motivation
 - Goldsmith-Pinkham et al. (2020)
- 4 Choosing an Appropriate Framework
- 5 *Recentering Method by Borusyak and Hull (2023, ECTA)
 - Motivation & Intuition
 - Formal Framework
 - Applications

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- 1 Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- 1 Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Who got selected for the intervention & who neighbors whom

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- ① Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Who got selected for the intervention & who neighbors whom

- ② Regional growth of market access from transportation upgrades:

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- ① Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Who got selected for the intervention & who neighbors whom

- ② Regional growth of market access from transportation upgrades:

Location + timing of upgrades & location and size of markets

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- 1 Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Who got selected for the intervention & who neighbors whom

- 2 Regional growth of market access from transportation upgrades:

Location + timing of upgrades & location and size of markets

- 3 An individual's eligibility for a public program, e.g. Medicaid:

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- ① Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Who got selected for the intervention & who neighbors whom

- ② Regional growth of market access from transportation upgrades:

Location + timing of upgrades & location and size of markets

- ③ An individual's eligibility for a public program, e.g. Medicaid:

State-level policy & individual income and demographics

Motivation

Many treatments/instruments are SSIV-like: combining multiple sets of variation, w/ some as-good-as-randomly assigned, but not all:

- ① Spatial/network/GE spillover treatments: e.g. the number of neighbors selected for a randomized intervention:

Who got selected for the intervention & who neighbors whom

- ② Regional growth of market access from transportation upgrades:

Location + timing of upgrades & location and size of markets

- ③ An individual's eligibility for a public program, e.g. Medicaid:

State-level policy & individual income and demographics

How can we just leverage the exogenous shocks to such z_i ?

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”
 - ▶ Specify many counterfactual sets of shocks

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”
 - ▶ Specify many counterfactual sets of shocks
 - ▶ Compute μ_i , the average z_i across counterfactuals, by simulation

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”
 - ▶ Specify many counterfactual sets of shocks
 - ▶ Compute μ_i , the average z_i across counterfactuals, by simulation
 - *the key confounder (similar to a propensity score)*

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”
 - ▶ Specify many counterfactual sets of shocks
 - ▶ Compute μ_i , the average z_i across counterfactuals, by simulation
 - *the key confounder (similar to a propensity score)*
 - ▶ “Recenter” z_i by μ_i (i.e. instrument with $\tilde{z}_i = z_i - \mu_i$) or control for μ_i

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”
 - ▶ Specify many counterfactual sets of shocks
 - ▶ Compute μ_i , the average z_i across counterfactuals, by simulation
 - *the key confounder (similar to a propensity score)*
 - ▶ “Recenter” z_i by μ_i (i.e. instrument with $\tilde{z}_i = z_i - \mu_i$) or control for μ_i
 - ▶ Conventional solutions (e.g. directly instrumenting with shocks or controlling for all features of exposure) are often infeasible

Borusyak and Hull (BH, 2023): Main Points

- 1 Non-random exposure to exogenous shocks generates systematic variation which can lead to omitted variable bias.
 - ▶ Randomizing roads \nrightarrow random market access growth from them
- 2 The systematic variation can be removed via novel “recentering”
 - ▶ Specify many counterfactual sets of shocks
 - ▶ Compute μ_i , the average z_i across counterfactuals, by simulation
 - *the key confounder (similar to a propensity score)*
 - ▶ “Recenter” z_i by μ_i (i.e. instrument with $\tilde{z}_i = z_i - \mu_i$) or control for μ_i
 - ▶ Conventional solutions (e.g. directly instrumenting with shocks or controlling for all features of exposure) are often infeasible
- 3 Recentering solution also can have attractive efficiency properties
 - ▶ Leverages non-random exposure to best predict shock effects

(Some) Other Settings where these Points are Relevant

Linear shift-share IV (Autor et al. 2013, Borusyak et al. 2022)

Nonlinear shift-share IV (Boustan et al. 2013, Berman et al. 2015, Chodorow-Reich and Wieland 2020, Derenoncourt 2021)

IV based on centralized school assignment mechanisms (Abdulkadiroğlu et al. 2017, 2019, Angrist et al. 2020)

Model-implied optimal IV (Adão-Arkolakis-Esposito 2021)

Weather instruments (Gomez et al. 2007, Madestam et al. 2013)

“Free space” instruments for mass media access (Olken 2009, Yanagizawa-Drott 2014)

Example 1: Market Access Effects via RCT

Theory suggests transportation upgrades affect local outcomes (e.g. land value) of regions i by increasing their market access (MA):

$$\Delta \log V_i = \beta \Delta \log MA_i + \varepsilon_i,$$

$$\text{where } MA_{it} = \sum_j \tau(g_t, loc_i, loc_j)^{-1} pop_j,$$

for road network g_t in periods $t=1,2$, region locations loc_j (co-determining travel cost τ), and regional population pop_j

Example 1: Market Access Effects via RCT

Theory suggests transportation upgrades affect local outcomes (e.g. land value) of regions i by increasing their market access (MA):

$$\Delta \log V_i = \beta \Delta \log MA_i + \varepsilon_i,$$

where $MA_{it} = \sum_j \tau(g_t, loc_i, loc_j)^{-1} pop_j,$

for road network g_t in periods $t = 1, 2$, region locations loc_j (co-determining travel cost τ), and regional population pop_j

Imagine an experiment randomly connecting adjacent regions by road

Example 1: Market Access Effects via RCT

Theory suggests transportation upgrades affect local outcomes (e.g. land value) of regions i by increasing their market access (MA):

$$\Delta \log V_i = \beta \Delta \log MA_i + \varepsilon_i,$$

where $MA_{it} = \sum_j \tau(g_t, loc_i, loc_j)^{-1} pop_j,$

for road network g_t in periods $t=1,2$, region locations loc_j (co-determining travel cost τ), and regional population pop_j

Imagine an experiment randomly connecting adjacent regions by road

Example 1: Market Access Effects via RCT

Theory suggests transportation upgrades affect local outcomes (e.g. land value) of regions i by increasing their market access (MA):

$$\Delta \log V_i = \beta \Delta \log MA_i + \varepsilon_i,$$

where $MA_{it} = \sum_j \tau(g_t, loc_i, loc_j)^{-1} pop_j,$

for road network g_t in periods $t=1,2$, region locations loc_j (co-determining travel cost τ), and regional population pop_j

Imagine an experiment randomly connecting adjacent regions by road

- MA only grows because of the random transportation shocks
- So can we view variation in MA growth as random and just run OLS?

Example 1: Market Access Effects via RCT

Theory suggests transportation upgrades affect local outcomes (e.g. land value) of regions i by increasing their market access (MA):

$$\Delta \log V_i = \beta \Delta \log MA_i + \varepsilon_i,$$

where $MA_{it} = \sum_j \tau(g_t, loc_i, loc_j)^{-1} pop_j,$

for road network g_t in periods $t=1,2$, region locations loc_j (co-determining travel cost τ), and regional population pop_j

Imagine an experiment randomly connecting adjacent regions by road

- MA only grows because of the random transportation shocks
- So can we view variation in MA growth as random and just run OLS?

No. Randomizing roads \nrightarrow randomizing MA due to them!

Illustration: Market Access on a Square Island

Start from no roads, assume equal population everywhere

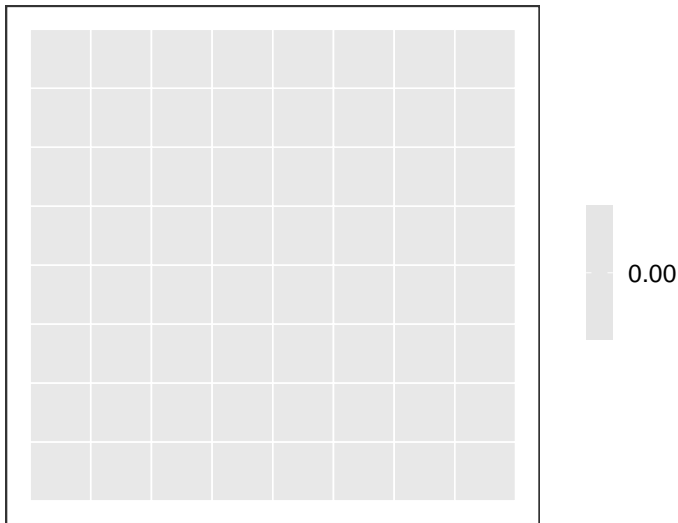


Illustration: Market Access on a Square Island

Randomly connect adjacent regions by road

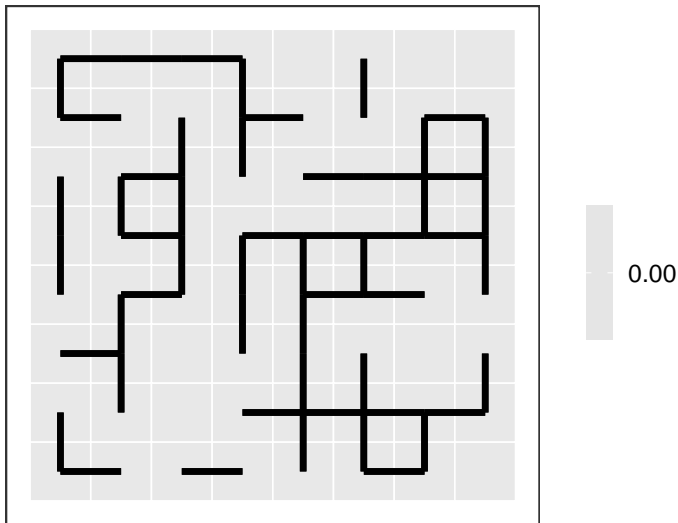


Illustration: Market Access on a Square Island

Randomly connect adjacent regions by road and compute MA growth

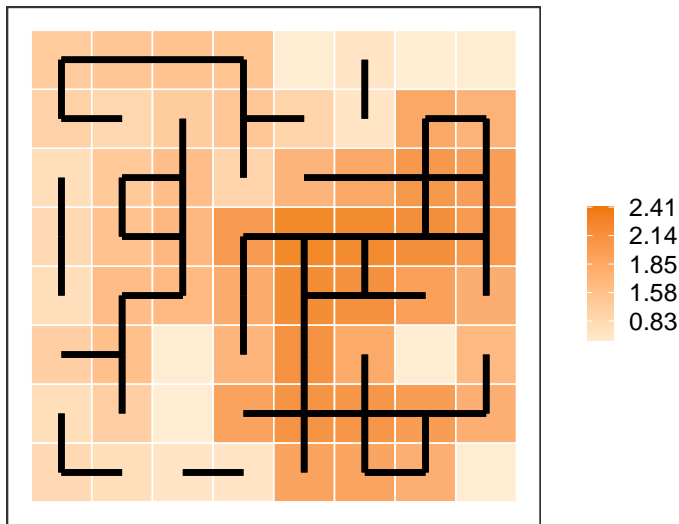


Illustration: Market Access on a Square Island

Randomly connect adjacent regions by road and compute MA growth

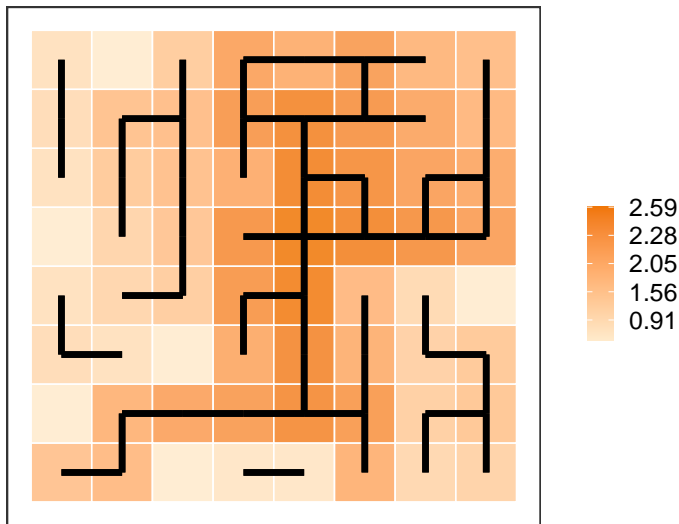
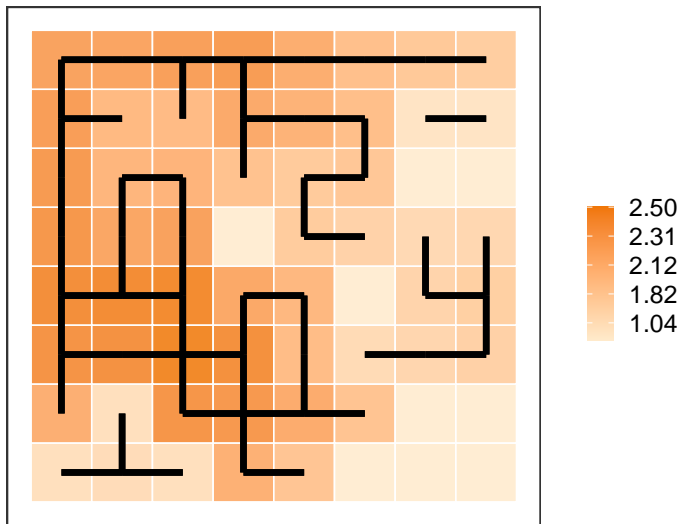


Illustration: Market Access on a Square Island

Randomly connect adjacent regions by road and compute MA growth



Expected Market Access Growth μ_i

Some regions get systematically more MA

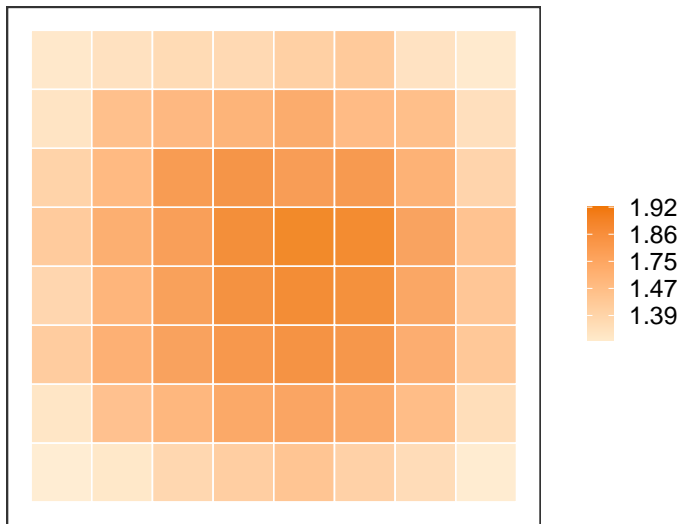


Illustration: High-Speed Rail in China

149 lines were built or planned (as of April 2019)



Illustration: High-Speed Rail in China

The 83 lines actually built by 2016. Suppose timing is random

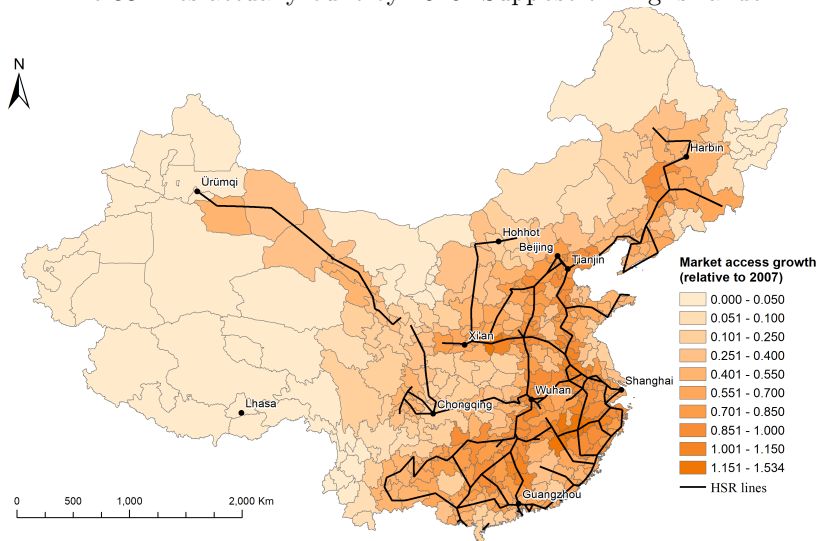


Illustration: High-Speed Rail in China

A counterfactual draw of 83 lines by 2016

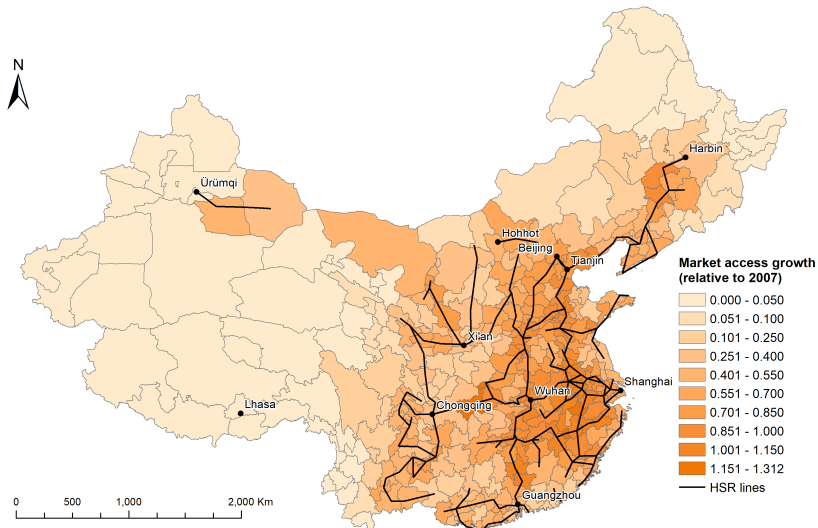
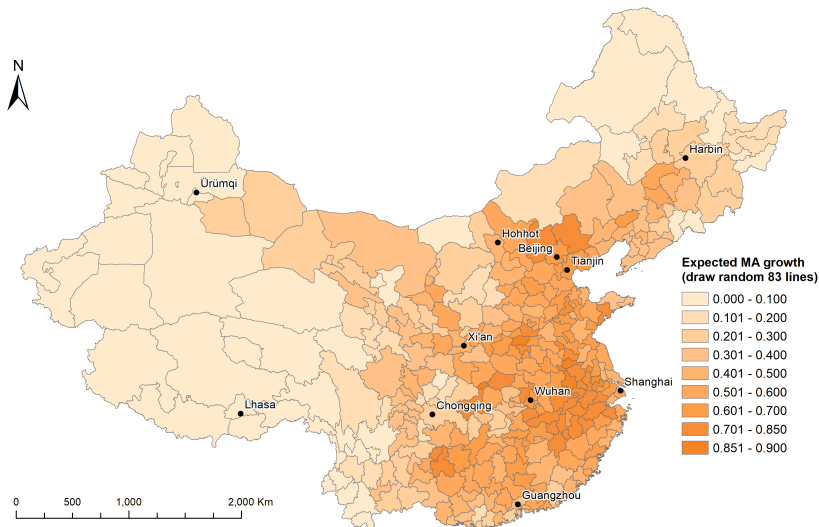


Illustration: High-Speed Rail in China

Expected MA growth, μ_i



OVB and Recentering Solution

Systematic variation in MA growth can generate OVB

- E.g. land values fall in the periphery because of rising sea levels
- More vs less developed Chinese regions may be on different trends

OVB and Recentering Solution

Systematic variation in MA growth can generate OVB

- E.g. land values fall in the periphery because of rising sea levels
- More vs less developed Chinese regions may be on different trends

Systematic variation can be removed via “recentering”:

$$\begin{array}{c} \text{Recentered} \\ \text{MA growth} \end{array} = \begin{array}{c} \text{Realized} \\ \text{MA growth} \end{array} - \begin{array}{c} \text{Expected} \\ \text{MA growth} \end{array}$$

OVB and Recentering Solution

Systematic variation in MA growth can generate OVB

- E.g. land values fall in the periphery because of rising sea levels
- More vs less developed Chinese regions may be on different trends

Systematic variation can be removed via “recentering”:

$$\text{Recentered MA growth} = \text{Realized MA growth} - \text{Expected MA growth}$$

Recentered MA is a valid instrument for realized MA growth

- Compares MA from actual and counterfactual shocks
- As it turns out, we can also control for expected MA growth

Linear SSIV Redux

Classic SSIV is a special case where $z_i = \sum_n \textcolor{violet}{s}_i \textcolor{brown}{g}_n$ is linear in the exogenous shocks

Linear SSIV Redux

Classic SSIV is a special case where $z_i = \sum_n s_{in} g_n$ is linear in the exogenous shocks

The expected instrument is $\mu_i = E[\sum_n s_{in} g_n \mid s] = \sum_n s_{in} E[g_n \mid s]$

Linear SSIV Redux

Classic SSIV is a special case where $z_i = \sum_n s_{in} g_n$ is linear in the exogenous shocks

The expected instrument is $\mu_i = E[\sum_n s_{in} g_n \mid s] = \sum_n s_{in} E[g_n \mid s]$

- If $E[g_n \mid s] = \gamma$, we need to adjust for $\gamma(\sum_n s_{in})$

Linear SSIV Redux

Classic SSIV is a special case where $z_i = \sum_n s_{in} g_n$ is linear in the exogenous shocks

The expected instrument is $\mu_i = E[\sum_n s_{in} g_n \mid s] = \sum_n s_{in} E[g_n \mid s]$

- If $E[g_n \mid s] = \gamma$, we need to adjust for $\gamma(\sum_n s_{in})$
- Linear in the sum-of-shares $S_i = \sum_n s_{in}$; it turns out controlling for this observable is enough (recall FWL theorem!)

Linear SSIV Redux

Classic SSIV is a special case where $z_i = \sum_n s_{in} g_n$ is linear in the exogenous shocks

The expected instrument is $\mu_i = E[\sum_n s_{in} g_n \mid s] = \sum_n s_{in} E[g_n \mid s]$

- If $E[g_n \mid s] = \gamma$, we need to adjust for $\gamma(\sum_n s_{in})$
- Linear in the sum-of-shares $S_i = \sum_n s_{in}$; it turns out controlling for this observable is enough (recall FWL theorem!)
- If g_n is only exogenous conditional on q_n , with $E[g_n \mid s, q] = q_n' \gamma$, we need to adjust for $\sum_n s_{in} E[g_n \mid s, q] = \gamma(\sum_n s_{in} q_n)$

Linear SSIV Redux

Classic SSIV is a special case where $z_i = \sum_n s_{in} g_n$ is linear in the exogenous shocks

The expected instrument is $\mu_i = E[\sum_n s_{in} g_n \mid s] = \sum_n s_{in} E[g_n \mid s]$

- If $E[g_n \mid s] = \gamma$, we need to adjust for $\gamma(\sum_n s_{in})$
- Linear in the sum-of-shares $S_i = \sum_n s_{in}$; it turns out controlling for this observable is enough (recall FWL theorem!)
- If g_n is only exogenous conditional on q_n , with $E[g_n \mid s, q] = q_n' \gamma$, we need to adjust for $\sum_n s_{in} E[g_n \mid s, q] = \gamma(\sum_n s_{in} q_n)$
- Controlling for $\sum_n s_{in} q_n$ is enough (sound familiar?)

Example 2: Effects of Program Eligibility

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

Example 2: Effects of Program Eligibility

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

- Suppose state policies are as-good-as-random

Example 2: Effects of Program Eligibility

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

- Suppose state policies are as-good-as-random
- But pre-determined demographics are endogenous \Rightarrow OLS biased

Example 2: Effects of Program Eligibility

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

- Suppose state policies are as-good-as-random
- But pre-determined demographics are endogenous \Rightarrow OLS biased

Standard “simulated instruments” solution (Currie and Gruber (1996)):
use state-level variation (average policy generosity across a “simulated” group of individuals) as a single IV for x_i

- This works, but is likely inefficient: the policy shocks likely have heterogeneous effects across individuals w/different demos

Gaining Power from Recentering

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

Gaining Power from Recentering

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

The BH approach:

- Formalize the policy experiment as “all permutations of g across states are equally likely”

Gaining Power from Recentering

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

The BH approach:

- Formalize the policy experiment as “all permutations of g across states are equally likely”
- Compute μ_i = the share of states in which i would be eligible

Gaining Power from Recentering

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

The BH approach:

- Formalize the policy experiment as “all permutations of g across states are equally likely”
- Compute μ_i = the share of states in which i would be eligible
- Leverage all variation in x_i but recenter by μ_i (or control for μ_i)

Gaining Power from Recentering

Consider the effects of individual eligibility x_i for Medicaid:

$$y_i = \beta x_i + \varepsilon_i$$

where x_i is determined by i 's state policy g_{state_i} and demographics

The BH approach:

- Formalize the policy experiment as “all permutations of g across states are equally likely”
- Compute μ_i = the share of states in which i would be eligible
- Leverage all variation in x_i but recenter by μ_i (or control for μ_i)
- Yields efficiency gain by better first-stage prediction, e.g. by removing i who are always or never eligible

General Setup

We have a model of $y_i = \beta x_i + \varepsilon_i$ for a fixed population $i = 1 \dots N$

- In the paper: extensions to heterogeneous effects, other controls, multiple treatments, nonlinear outcome models, panel data...

General Setup

We have a model of $y_i = \beta x_i + \varepsilon_i$ for a fixed population $i = 1 \dots N$

- In the paper: extensions to heterogeneous effects, other controls, multiple treatments, nonlinear outcome models, panel data...

We have a candidate instrument $z_i = f_i(g, w)$, where g is a vector of shocks; w collects predetermined variables; $f_i(\cdot)$ are known mappings

- Applies to any z_i which can be constructed from observed data
- Nests reduced-form regressions: $x_i = z_i$
- Allows $g = (g_1, \dots, g_K)$ to vary at a different level than i

General Setup

We have a model of $y_i = \beta x_i + \varepsilon_i$ for a fixed population $i = 1 \dots N$

- In BH: extensions to heterogeneous effects, other controls, multiple treatments, nonlinear outcome models, panel data...

We have a candidate instrument $z_i = f_i(g, w)$, where g is a vector of shocks; w collects predetermined variables; $f_i(\cdot)$ are known mappings

Assumptions:

- 1 Shocks are exogenous: $g \perp \varepsilon \mid w$
- 2 Conditional distribution $G(g \mid w)$ is known (e.g. via randomization protocol or uniform across permutations of g)

Main Results

The expected instrument, $\mu_i = E[f_i(g, w) \mid w] \equiv \int f_i(g, w) dG(g \mid w)$, is the sole confounder generating OVB:

$$E \left[\frac{1}{N} \sum_i z_i \varepsilon_i \right] = E \left[\frac{1}{N} \sum_i \mu_i \varepsilon_i \right] \neq 0, \text{ in general}$$

Main Results

The expected instrument, $\mu_i = E[f_i(g, w) \mid w] \equiv \int f_i(g, w) dG(g \mid w)$, is the sole confounder generating OVB:

$$E \left[\frac{1}{N} \sum_i z_i \varepsilon_i \right] = E \left[\frac{1}{N} \sum_i \mu_i \varepsilon_i \right] \neq 0, \text{ in general}$$

The *recentered instrument* $\tilde{z}_i = z_i - \mu_i$ is a valid instrument for x_i :

$$E \left[\frac{1}{N} \sum_i \tilde{z}_i \varepsilon_i \right] = 0$$

Main Results

The expected instrument, $\mu_i = E[f_i(g, w) \mid w] \equiv \int f_i(g, w) dG(g \mid w)$, is the sole confounder generating OVB:

$$E\left[\frac{1}{N}\sum_i z_i \varepsilon_i\right] = E\left[\frac{1}{N}\sum_i \mu_i \varepsilon_i\right] \neq 0, \text{ in general}$$

The *recentered instrument* $\tilde{z}_i = z_i - \mu_i$ is a valid instrument for x_i :

$$E\left[\frac{1}{N}\sum_i \tilde{z}_i \varepsilon_i\right] = 0$$

Regressions which control for μ_i also identify β (implicitly recenter, by the FWL theorem)

Extensions

Consistency: follows when \tilde{z}_i is weakly mutually dependent across i

Robustness to heterogeneous treatment effects: \tilde{z}_i identifies a convex avg. of β_i under appropriate first-stage monotonicity

Randomization inference provides exact confidence intervals for β (under constant effects) and falsification tests

BH also characterize the **asy. efficient** recentered IV among all $f_i(\cdot)$

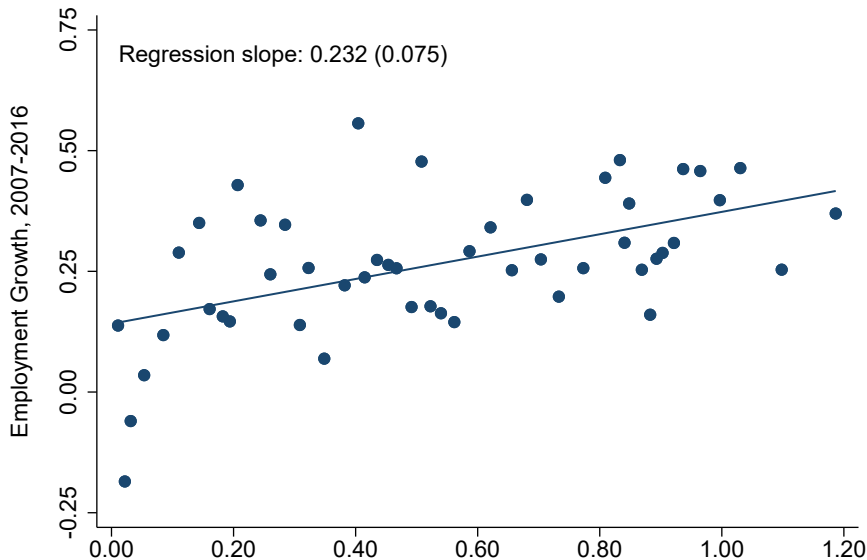
App. 1: Market Access from Chinese High-Speed Rail

BH first show how instrument recentering can address OVB when estimating the effects of market access growth

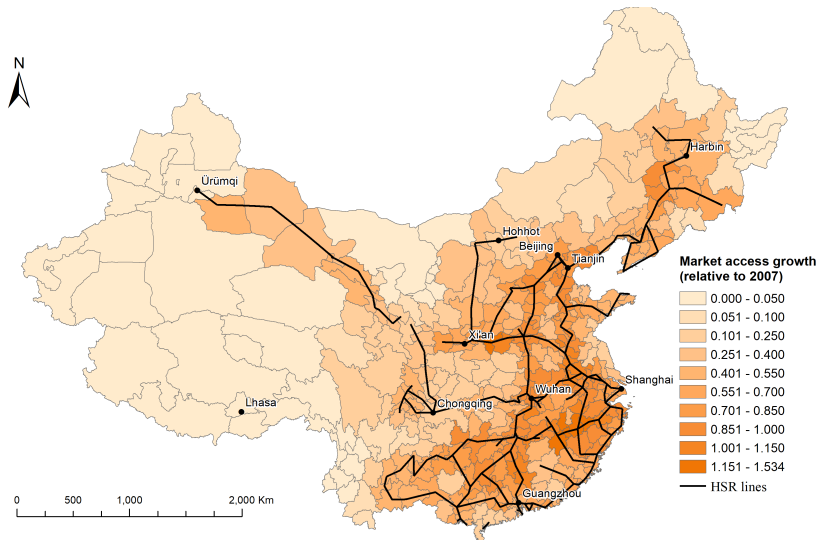
Setting: Chinese HSR; 83 lines built 2008–2016, 66 yet unbuilt

- Market access: $MA_{it} = \sum_k \exp(-0.02\tau_{ikt}) p_{k,2000}$, where τ_{ikt} is HSR-affected travel time between prefecture capitals (Zheng and Kahn, 2013) and $p_{i,2000}$ is prefecture i 's population in 2000
- Relate to employment growth in 274 prefectures, 2007-2016

Simple OLS Regressions Suggest a Large MA Effect



High vs. Low MA Growth is Not a Convincing Contrast!



How to Find Valid Treatment-Control Contrasts?

Add controls (province FE, longitude, etc...)

- Hard to justify *ex ante* since MA is a variable constructed based on a structural model
- No experimental analog

How to Find Valid Treatment-Control Contrasts?

Add controls (province FE, longitude, etc...)

- Hard to justify *ex ante* since MA is a variable constructed based on a structural model
- No experimental analog

Find valid contrasts for *one* source of variation—a natural experiment

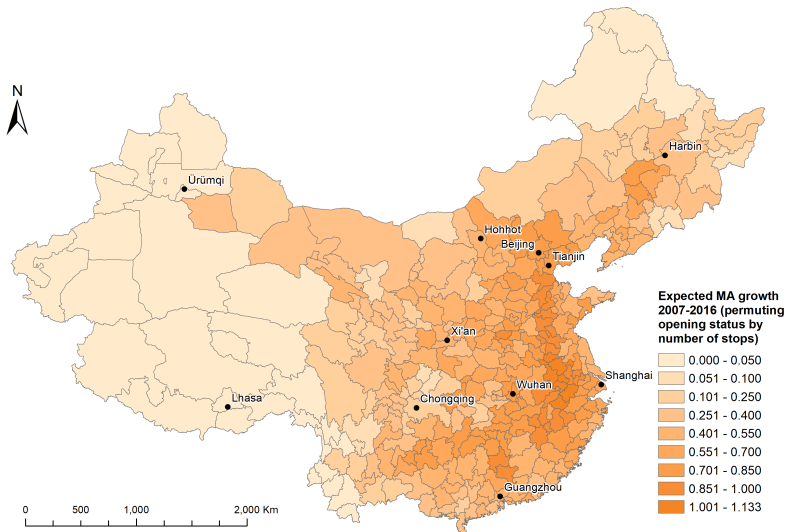
- Bartelme (2018): shocks affecting market size
- Donaldson (2018): built vs unbuilt lines
- BH application: assume random timing of observably similar lines

Built and Planned HSR Lines

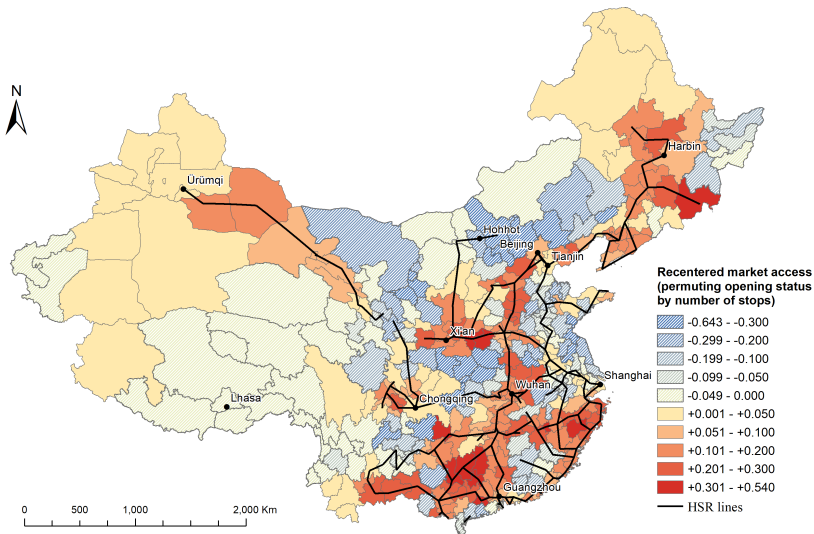
BH reshuffle built & planned lines connecting the same # of regions



Expected Market Access Growth (2007–2016), μ_i



Recentered Market Access Growth (2007–2016), \tilde{z}_i

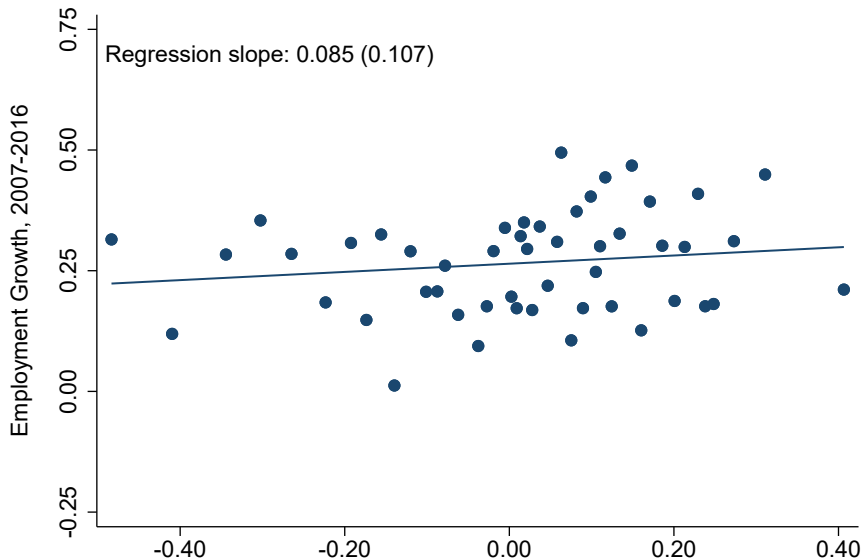


Market Access Balance Regressions

	Unadjusted	Recentered		
	(1)	(2)	(3)	(4)
Distance to Beijing	-0.292 (0.063)	0.069 (0.040)		0.089 (0.045)
Latitude/100	-3.323 (0.648)	-0.325 (0.277)		-0.156 (0.320)
Longitude/100	1.329 (0.460)	0.473 (0.239)		0.425 (0.242)
Expected Market Access Growth			0.027 (0.056)	0.056 (0.066)
Constant	0.536 (0.030)	0.014 (0.018)	0.014 (0.020)	0.014 (0.018)
Joint RI p-value		0.489	0.807	0.536
R^2	0.823	0.079	0.007	0.082
Prefectures	274	274	274	274

Regressions of unadjusted and recentered market access growth on geographic features.
Spatial-clustered standard errors in parentheses.

Recentered MA Doesn't Predict Employment Growth!



Adjusted Estimates of Market Access Effects

	Unadjusted OLS (1)	Recentered IV (2)	Controlled OLS (3)
<i>Panel A. No Controls</i>			
Market Access Growth	0.232 (0.075)	0.081 (0.098) [-0.315, 0.328]	0.069 (0.094) [-0.209, 0.331]
Expected Market Access Growth			0.318 (0.095)
<i>Panel B. With Geography Controls</i>			
Market Access Growth	0.132 (0.064)	0.055 (0.089) [-0.144, 0.278]	0.045 (0.092) [-0.154, 0.281]
Expected Market Access Growth			0.213 (0.073)
Recentered Prefectures	No 274	Yes 274	Yes 274

Regressions of log employment growth on log market access growth in 2007–2016.
Spatial-clustered standard errors in parentheses; permutation-based 95% CI in brackets

App. 2: Efficient Estimation of Medicaid Effects

Setting: U.S. Medicaid, partially expanded in 2014 under the ACA

- 19 of 43 states with low Medicaid coverage expanded to 138% FPL
- View expansion decisions as random across states with same-party governors, but not household demographics or pre-2014 policy
- Outcomes: Medicaid takeup and private insurance crowdout

App. 2: Efficient Estimation of Medicaid Effects

Setting: U.S. Medicaid, partially expanded in 2014 under the ACA

- 19 of 43 states with low Medicaid coverage expanded to 138% FPL
- View expansion decisions as random across states with same-party governors, but not household demographics or pre-2014 policy
- Outcomes: Medicaid takeup and private insurance crowdout

We compare two estimators, both valid under the same assumptions:

- Simulated IV: use state-level variation only (i.e. expansion dummy)
- Recentered IV: predict eligibility from expansion decisions & non-random demographics, and recenter

App. 2: Efficient Estimation of Medicaid Effects

Setting: U.S. Medicaid, partially expanded in 2014 under the ACA

- 19 of 43 states with low Medicaid coverage expanded to 138% FPL
- View expansion decisions as random across states with same-party governors, but not household demographics or pre-2014 policy
- Outcomes: Medicaid takeup and private insurance crowdout

We compare two estimators, both valid under the same assumptions:

- Simulated IV: use state-level variation only (i.e. expansion dummy)
- Recentered IV: predict eligibility from expansion decisions & non-random demographics, and recenter

Via non-random variation, recentered IV has ≈ 3 times smaller SEs

Estimates with Simulated vs. Recentered IV

	Has Medicaid		Has Private Insurance		Has Employer-Sponsored Insurance	
	Simulated IV (1)	Recentered IV (2)	Simulated IV (3)	Recentered IV (4)	Simulated IV (5)	Recentered IV (6)
<i>Panel A. Eligibility Effects</i>						
Eligibility	0.132 (0.028) [0.080,0.216]	0.072 (0.010) [0.051,0.093]	-0.048 (0.023) [-0.110,0.009]	-0.023 (0.007) [-0.040,-0.007]	0.009 (0.014) [-0.034,0.052]	-0.009 (0.005) [-0.021,0.004]
<i>Panel B. Enrollment Effects</i>						
Has Medicaid			-0.361 (0.165) [-0.813,0.082]	-0.321 (0.092) [-0.566,-0.108]	0.068 (0.111) [-0.232,0.421]	-0.125 (0.061) [-0.263,0.070]
P-value: SIV=RIV			0.719		0.104	
Exposed Sample	N	Y	N	Y	N	Y
States	43	43	43	43	43	43
Individuals	2,397,313	421,042	2,397,313	421,042	2,397,313	421,042

1% ACS sample of non-disabled adults in 2013–14, diff-in-diff IV regressions using one of the two instruments. Controls include state and year fixed effects and an indicator for Republican governor interacted with year. State-clustered standard errors in parentheses; wild score bootstrap 95% CI in brackets

Conclusions

In both linear SSIV and more elaborate settings, the most important thing is to decide *ex ante* what identifying variation you want to use

Conclusions

In both linear SSIV and more elaborate settings, the most important thing is to decide *ex ante* what identifying variation you want to use

- When leveraging a natural experiment, recentering (e.g. controlling for sum-of-shares, in linear SSIV) can help

Conclusions

In both linear SSIV and more elaborate settings, the most important thing is to decide *ex ante* what identifying variation you want to use

- When leveraging a natural experiment, recentering (e.g. controlling for sum-of-shares, in linear SSIV) can help
- Non-experimental assumptions (e.g. parallel trends) typically require other approaches

Conclusions

In both linear SSIV and more elaborate settings, the most important thing is to decide *ex ante* what identifying variation you want to use

- When leveraging a natural experiment, recentering (e.g. controlling for sum-of-shares, in linear SSIV) can help
- Non-experimental assumptions (e.g. parallel trends) typically require other approaches
- The source of variation can (should?) guide inference

Conclusions

In both linear SSIV and more elaborate settings, the most important thing is to decide *ex ante* what identifying variation you want to use

- When leveraging a natural experiment, recentering (e.g. controlling for sum-of-shares, in linear SSIV) can help
- Non-experimental assumptions (e.g. parallel trends) typically require other approaches
- The source of variation can (should?) guide inference

After deciding + appropriately adjusting the analysis, try to falsify the identifying variation (*ex post*) – via balance or pre-trend tests

Conclusions

In both linear SSIV and more elaborate settings, the most important thing is to decide *ex ante* what identifying variation you want to use

- When leveraging a natural experiment, recentering (e.g. controlling for sum-of-shares, in linear SSIV) can help
- Non-experimental assumptions (e.g. parallel trends) typically require other approaches
- The source of variation can (should?) guide inference

After deciding + appropriately adjusting the analysis, try to falsify the identifying variation (*ex post*) – via balance or pre-trend tests

Much more work to be done on the various econometrics here!

Appendix: Proof of IV Estimator using FWL Theorem

- The first-stage regression is $x_\ell = \hat{\alpha}z_\ell + \hat{\lambda}'w_\ell + \tilde{x}_\ell$. By the FWL theorem, the estimator $\hat{\alpha}$ can be expressed as

$$\hat{\alpha} = \frac{\text{Cov}(\tilde{x}_\ell, \tilde{z}_\ell)}{V(\tilde{z}_\ell)} = \frac{\text{Cov}(\tilde{x}_\ell, z_\ell)}{V(\tilde{z}_\ell)}$$

- Plugging \hat{x}_ℓ into the second-stage regression, we have

$$y_\ell = \beta(\hat{\alpha}z_\ell + \hat{\lambda}'w_\ell) + \gamma'w_\ell + \varepsilon_\ell$$

By the FWL theorem, the estimator $\beta\hat{\alpha}$ can be expressed as

$$\hat{\beta}\hat{\alpha} = \frac{\text{Cov}(y_\ell, \tilde{z}_\ell)}{V(\tilde{z}_\ell)} = \frac{\text{Cov}(y_\ell, z_\ell)}{V(\tilde{z}_\ell)} \Rightarrow \hat{\beta} = \frac{\text{Cov}(z_\ell, \tilde{y}_\ell)}{\text{Cov}(z_\ell, \tilde{x}_\ell)}$$