

# Canonical Research Designs V: Bartik, Simulated, and Granular Instruments

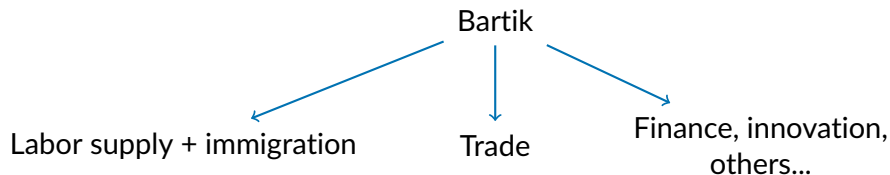
Paul Goldsmith-Pinkham

April 3, 2025

# Roadmap for Today

- In some cases, the source of exogenous variation (either in an IV setting, or just OLS) is straightforward
  - There is a single policy or source of variation
- However, in other settings, there are more complicated sources of variation exploited to identify effects. Today we'll focus on three:
  - Bartik (shift-share) instruments: three recent papers on commonly used identification approach
  - Simulated instruments: reframe an older literature in a new light using Borusyak and Hull (2022) paper
  - Granular instruments: identification approach in Gabaix and Koijen (2023) leveraging differences in the size distribution across firms
- Key historical feature of some of these approaches is that they had an “intuitive” feature of identification, but formal properties were not established for several decades
  - Analogous to staggered DiD lit!

# Bartik instruments are used everywhere



- Thread that links all Bartik applications:
  - local markets composed of many “categories”
  - need for identification
- Approach has been used since the early 90:
  - sometimes called “shift-share” or “industry mix” instruments

# Examples of Bartik instruments in many subfields

**Immigration:** Altonji and Card (1991), Card (2001)

**Bank Lending:** Amiti and Weinstein (2018), Greenstone, Mas and Nguyen (2015)

**Market Size + Demography:** Acemoglu and Linn (2004), Jaravel (2018)

**Labor Supply Elasticity:** Blanchard and Katz (1992), **Bartik** (1991)

**Fiscal Multipliers:** Nakamura and Steinsson (2014)

**Trade + Labor:** Autor, Dorn, and Hanson (2013), Autor, Dorn and Hanson (2018), etc.

**Foreign Aid:** Nunn and Qian (2014)

**Portfolio Allocation:** Calvet, Campbell, and Sodini (2009)

**Trade + Prices:** Piveteau and Smagghue (2017), de Roux et al. (2017)

**Automation:** Acemoglu and Restrepo (2017)

# Many paths lead to Bartik

- Diverse literature leads to many motivations and justifications for Bartik approach
- Two distinct approaches in the literature:
  1. Applied micro statistical approach: interested in a reduced form causal relationship; need an instrument that is uncorrelated with error term; make argument that Bartik instrument is defensible
  2. Structural approach: interested in particular parameters from model; assumptions of model motivate certain estimating equations
- So what is the Bartik approach anyway?

## Motivation: local labor market approaches + reduced form

Consider a local labor market regression like the following:

$$y_l = \beta_0 + \beta x_l + \epsilon_l$$

- $\mathbb{E}[x_l \epsilon_l] \neq 0 \Rightarrow$  need an instrument to estimate  $\beta$
- E.g. Autor, Dorn and Hanson (2013) setting:
  - $l$ : location (commuting zone)
  - $y_l$ : manufacturing employment *growth*
  - $x_l$ : import exposure to China *growth*
  - $\beta$ : effect of rise of China on manufacturing employment
  - an instrument for location-level exposure to trade with China

# The Bartik instrument

Accounting identity #1:

$$x_l = \sum_{k=1}^K z_{lk} g_{lk}$$

- $z_{lk}$ : location-industry shares ( $Z_l$ )
- $g_{lk}$ : location-industry growth (in imports) rates ( $G_l$ )

Accounting identity #2:

$$\underbrace{g_{lk}}_{\text{location-industry}} = \underbrace{g_k}_{\text{industry}} + \underbrace{\tilde{g}_{lk}}_{\text{idiosyncratic location-industry}}$$

Infeasible Bartik:

$$B_l = \sum_{k=1}^K z_{lk} g_k$$

## This gives us a simple 2SLS structure

$$y_l = \beta_0 + \beta x_l + \epsilon_l$$

$$x_l = \pi_0 + \pi_1 B_l + u_l$$

$$B_l = \sum_{k=1}^K z_{lk} g_k$$

$$g_{lk} = g_k + \tilde{g}_{lk}$$

Bank-lending relationships: e.g., Greenstone, Mas and Nguyen (2015)

- $z_{lk}$ : location ( $l$ ) share of loan origination from bank  $k$
- $g_{lk}$ : loan growth in location  $l$  by bank  $k$
- $g_k$ : part of loan growth due to bank supply shock



## Other instruments have this structure

$$y_I = \beta_0 + \beta x_I + \epsilon_I$$

$$x_I = \pi_0 + \pi_1 B_I + u_I$$

$$B_I = \sum_{k=1}^K z_{Ik} g_k$$

$$g_{Ik} = g_k + \tilde{g}_{Ik}$$

Immigrant enclave: e.g., Altonji and Card (1991)

- $z_{Ik}$ : share of people from foreign  $k$  living in  $I$  (in a base period)
- $g_{Ik}$ : growth in number of people from  $k$  to  $I$
- $g_k$ : growth in people from  $k$  nationally

## Other instruments have this structure

$$y_l = \beta_0 + \beta x_l + \epsilon_l$$

$$x_l = \pi_0 + \pi_1 B_l + u_l$$

$$B_l = \sum_{k=1}^K z_{lk} g_k$$

$$g_{lk} = g_k + \tilde{g}_{lk}$$

Market size and demography: e.g., Acemoglu and Linn (2004)

- $z_{lk}$ : spending share on drug  $l$  from age group  $k$
- $g_{lk}$ : growth in spending of group  $k$  on drug  $l$
- $g_k$ : growth in spending of group  $k$  (due to population aging)

# What's necessary for consistency?

$$y_I = \beta_0 + \beta x_I + \epsilon_I$$

$$x_I = \pi_0 + \pi_1 B_I + u_I$$

$$B_I = \sum_{k=1}^K z_{Ik} g_k$$

$$g_{Ik} = g_k + \tilde{g}_{Ik}$$

- We need  $B_I$  to be a valid instrument
- Requires two conditions with constant effects:
  1. Relevance:  $\pi_1 \neq 0$ , e.g.  $\text{Cov}(B_I, x_I) \neq 0$
  2. Exclusion:  $E(B_I \epsilon_I) = 0$
- Key flaw in this literature until recently: economic + statistical content of exclusion has been vague and sometimes confused

## Recent Literature on this topic

- Three papers addressed this question, and can be split into two grouping
- The division between papers can be split based on focus on  $z_{lk0}$  vs.  $g_{kt}$ 
  1. Goldsmith-Pinkham, Sorkin and Swift (2020) focus on  $z_{lk0}$  and make an analogy to difference-in-differences
  2. Adao, Kolesar and Morales (2019) and Borusyak, Hull and Jaravel (2020) focus on  $g_{kt}$ , and make a strong connection to the design based approach (e.g. these are as-if random shocks)
- Key problem, historically, in this literature, was the lack of a coherent defense of the identifying variation
  - These papers provide a way of doing this! But you have to pick one approach

# When is the estimator consistent for the estimand of interest?

## What is the identification condition?

$$\hat{\beta}_{Bartik} = \frac{\sum_{l=1}^L \sum_{t=1}^T \sum_{k=1}^K z_{lkt} g_{kt} y_{lt}^{\perp}}{\sum_{l=1}^L \sum_{t=1}^T \sum_{k=1}^K z_{lkt} g_{kt} x_{lt}^{\perp}}$$

Two ideas:

- “Shares” : talk about properties of  $z_{lkt}$ 
  - Conditional exogeneity
  - *model based* – diff-in-diff style approach
- “Shocks” (Borusyak, Hull and Jaravel (2018)): talk about properties of  $g_{kt}$ 
  - Random, and a large number (equivalent industry-level regression)
  - *design-based* (in spirit) – IV strategy

# Understanding the shares assumption in GPSS

1. One time period, two industries
2. T time periods, two industries
3. One time period, K industries

## Special case #1: One time period, two industries

- $z_{l2} = 1 - z_{l1}$
- Bartik:

$$\begin{aligned}B_l &= z_{l1}g_1 + z_{l2}g_2 = z_{l1}g_1 + (1 - z_{l1})g_2 \\&= g_2 + (g_1 - g_2)z_{l1}\end{aligned}$$

First-stage:

$$\begin{aligned}x_l &= \gamma_0 + \gamma B_l + \eta_l \\x_l &= \underbrace{\gamma_0 + \gamma g_2}_{\text{constant}} + \underbrace{\gamma(g_1 - g_2)}_{\text{coefficient}} z_{l1} + \eta_l\end{aligned}$$

The instrument is  $z_{l1}$ , while  $g_k$  affects relevance

► Why OLS is biased

## Special case #2: T time periods, two industries

Panel Bartik:

$$B_{lt} = Z_{l10}g_{1t} + Z_{l20}g_{2t} = g_{2t} + \underbrace{\Delta_{gt}}_{g_{1t}-g_{2t}} Z_{l10}$$

First stage:

$$\begin{aligned}x_{lt} &= \tau_l + \tau_t + \gamma B_{lt} + \eta_{lt} \\x_{lt} &= \tau_l + \underbrace{(\tau_t + \gamma g_{2t})}_{\tilde{\tau}_t} + \underbrace{\gamma \Delta_{gt}}_{\tilde{\gamma}_t} Z_{l10} + \eta_{lt}\end{aligned}$$

- Industry shares times time period is the instrument
- (Updated industry shares: similar)



## Special case #2: T time periods, two industries

- Analogy to continuous difference-in-differences
  - $\Delta_{gt}$  is size of policy
  - $z_{i10}$  is exposure to policy
- Sometimes a “pre-period” before policy: test for parallel pre-trends
  - E.g., in ADH, what happens from 1970 to 1990?

## Special case #3: One time period, K industries

- $G$ :  $K \times 1$  vector of  $g_k$
- $Z$ :  $L \times K$ , matrix of  $Z_l$
- $Y^\perp, X^\perp, B = (ZG)$ :  $L \times 1$ , vectors of  $y_l^\perp, x_l^\perp$  and  $B_l$
- $\Omega$ :  $K \times K$

$$\hat{\beta}_{Bartik} = \frac{B' Y^\perp}{B' X^\perp}$$
$$\hat{\beta}_{GMM} = \frac{(X^{\perp'} Z) \Omega (Z' Y^\perp)}{(X^{\perp'} Z) \Omega (Z' X^\perp)}$$

If  $\Omega = (GG')$ , then  $\hat{\beta}_{Bartik} = \hat{\beta}_{GMM}$

## Full general result with $T$ time periods and $K$ industries

Two estimators are numerically identical:

- TSLS with Bartik instrument
- GMM with industry shares  $\times$  time period as instruments and a particular weight matrix

$$\hat{\beta}_{Bartik} = \frac{\mathbf{B}'\tilde{\mathbf{Y}}^\perp}{\mathbf{B}'\tilde{\mathbf{X}}^\perp}$$
$$\hat{\beta}_{GMM} = \frac{(\mathbf{X}^\perp'\tilde{\mathbf{Z}})\Omega(\tilde{\mathbf{Z}}'\mathbf{Y}^\perp)}{(\mathbf{X}^\perp'\tilde{\mathbf{Z}})\Omega(\tilde{\mathbf{Z}}'\mathbf{X}^\perp)}$$

$\Omega = (\mathbf{G}\mathbf{G}')$ , and  $\tilde{\mathbf{Z}}$  is an  $LT \times KT$  stacked vector of  $Z_0$  interacted with time fixed effects and  $\mathbf{G}$  is a  $KT \times 1$  stacked vector of growth rates  $g_{kt}$ .

## Understanding the shocks assumption in BHJ (drop time for convenience)

$$\hat{\beta}_{Bartik} = \frac{\sum_{l=1}^L \sum_{k=1}^K z_{lk} g_k y_l^{\perp}}{\sum_{l=1}^L \sum_{k=1}^K z_{lk} g_k x_l^{\perp}}$$

- Now let's assume  $g_k$  is as good as randomly assigned (conditional on the observed data)
- Can view this as just-identified regression at industry level:

$$\hat{\beta}_{Bartik} = \frac{\sum_{k=1}^K g_k \sum_{l=1}^L z_{lk} y_l^{\perp}}{\sum_{k=1}^K g_k \sum_{l=1}^L z_{lk} x_l^{\perp}} = \frac{\sum_{k=1}^K g_k z_k \bar{y}_k^{\perp}}{\sum_{k=1}^K g_k z_k \bar{x}_k^{\perp}}$$

where  $z_k = L^{-1} \sum_{l=1}^L z_{lk}$  and  $\bar{y}_k^{\perp} = \sum_l z_{lk} y_l^{\perp} / \sum_l z_{lk}$  (same for  $\bar{x}_k^{\perp}$ )

## Understanding the shocks assumption in BHJ (drop time for convenience)

$$\hat{\beta}_{Bartik} = \frac{\sum_{k=1}^K g_k \sum_{l=1}^L z_{lk} y_l^{\perp}}{\sum_{k=1}^K g_k \sum_{l=1}^L z_{lk} x_l^{\perp}} = \frac{\sum_{k=1}^K g_k z_k \bar{y}_k^{\perp}}{\sum_{k=1}^K g_k z_k \bar{x}_k^{\perp}}$$

BHJ assumptions:

1.  $g_k$  is randomly assigned
2. You have *many* uncorrelated  $g_k$  shocks
  - Average shock exposure goes to zero
  - shocks are sufficiently uncorrelated (doesn't have to be completely independent, but need enough)
3. Key underappreciated assumption :  $g_k$  is often not observed. We use an empirical proxy.
  - When is this valid? We have  $\hat{g}_k = g_k + \psi_k$  - need  $\psi_k$  to not be correlated across  $k$ !
  - Important economic question of what drives these equilibrium measures we use (e.g. industry growth)

# When are these views plausible? What do they mean?

## Shares

*Conditional* exogeneity:

- Typically: exogenous to changes in error term, not levels of outcome
- Standard in diff-in-diff (exclusion): in a period, exposure to an industry matters for outcome only through  $x$

## Shocks

- Large number of industries (shares are misspecified, need it to average out)
- Random shocks across industries – need the shocks to be conditionally random

## How do we choose?

- The shocks approach is more design-based (which can be appealing), but requires an argument why shocks are randomly assigned
- The shares approach is model-based, so suffers from same issues as diff-in-diff, but may more naturally work in your setting.

# How to decide if a paper is about shocks or shares?

For Autor, Dorn and Hanson's China Shock: Shocks:

- Explains why  $g_{kt}^{high-income}$  rather than  $g_{kt}^{US}$  (hard to rationalize under shares)
- Natural in a trade model: why would imports from China rise (in a trade model)?  
Independent industry-specific shocks

Shares:

- Explains why  $z_{lkt-1}$  rather than  $z_{lkt}$  (hard to rationalize under shocks)
- Explains why it is important for identification to study local labor markets (as opposed to parameter of interest where we want to think about spillovers)

Bottom line: a little hard to tell what exactly ADH are assuming; ADH approach does not appear to satisfy testable assumptions under GPSS, but do appear to under BHJ.

## Key thing to remember from today

- Assuming independence or exogeneity on the basis of a model does not necessarily make it true
  - E.g. Hausman instruments in IO models – model may assume that exclusion restriction is satisfied, but not necessarily true in reality
- Assuming that two things are independent because they don't seem “related” doesn't make it true
  - Bartik literature many times argues that national nature of shocks “decouples” the instrument from local market conditions. However, it still exploits local characteristics. Need to make very specific arguments to validate claim (will come to this).
- When evaluating an identification strategy, you should be able to describe counterfactual claims using the measure. This is typically not concrete in Bartik – try to make it concrete! What is exactly changing in China? Why is it random?



# Decomposing Bartik (GPSS 2020)

(Special case of Rotemberg (1983), proposition 1)

$$\hat{\beta}_{Bartik} = \sum_k \hat{\alpha}_k \hat{\beta}_k, \quad \sum_k \hat{\alpha}_k = 1$$

IV estimate using only the  $k^{th}$  instrument:

$$\hat{\beta}_k = (Z_k' X)^{-1} Z_k' Y$$

“Rotemberg” weight:

$$\hat{\alpha}_k = \frac{g_k Z_k' X}{\sum_{k=1}^K g_k Z_k' X}$$

# Interpretation: sensitivity to misspecification elasticity

Conley, Hansen and Rossi (2012); Andrews, Gentzkow and Shapiro (2017)

Local misspecification:  $\epsilon_{lt} = L^{-1/2} V_{lt} + \tilde{\epsilon}_{lt}$ ,  $\text{Cov}(V_{lt}, Z_{lt}) \neq 0$ ,

- $\sqrt{L}(\hat{\beta} - \beta_0) \xrightarrow{d} \tilde{\beta}$ ,  $\mathbb{E}[\tilde{\beta}] = \text{bias (misspecification) of Bartik instrument}$
- $\sqrt{L}(\hat{\beta}_k - \beta_0) \xrightarrow{d} \tilde{\beta}_k$ ,  $\mathbb{E}[\tilde{\beta}_k] = \text{bias (misspecification) of } k\text{th instrument}$

Suppose  $\beta_0 \neq 0$ . Percentage bias:

$$\frac{\mathbb{E}[\tilde{\beta}]}{\beta_0} = \sum_k \alpha_k \frac{\mathbb{E}[\tilde{\beta}_k]}{\beta_0}$$

Industry with high  $\alpha_k$ :

- an industry where it matters whether it is misspecified (endogenous)
  - because it is “important” in the estimate

## Top five industries (out of 397)

	$\hat{\alpha}_k$	$g_k^{\text{high-income}}$	$\hat{\beta}_k$
Games and Toys	0.182	174.841	-0.151
Electronic Computers	0.182	85.017	-0.620
Household Audio and Video	0.130	118.879	0.287
Computer Equipment	0.076	28.110	-0.315
Telephone Apparatus	0.058	37.454	-0.305
	0.628/1.379		-0.230

*The **main source of variation in exposure** is within-manufacturing specialization in industries subject to different degrees of import competition...there is differentiation according to **local labor market reliance on labor-intensive industries**...By 2007, China accounted for over 40 percent of US imports in four four-digit SIC industries (**luggage, rubber and plastic footwear, games and toys, and die-cut paperboard**) and over 30 percent in 28 other industries, including **apparel, textiles, furniture, leather goods, electrical appliances, and jewelry**.*

— Autor, Dorn and Hanson (2013) , pg. 2123

## Three tests of the identifying condition (under GPSS (2020))

1. Confounds (or correlates)
2. Pre-trends
3. Alternative estimators and overidentification
  - There are *also* tests for BHJ – similar to assuming strict ignorability, you can test for balance on observables (like the confounds above) of industries and locations

# Alternative estimators and overidentification tests

Basic insight in GPSS: many instruments

- Estimators (maximum likelihood): LIML, Hausman, Newey, Woutersen, Chao and Swanson (2012) HFUL (heteroskedasticity-Fuller (1977))
- Estimators (two-step): TSLS (problematic), Bartik TSLS, MBTSLS (Anatolyev (2013), and Kolesar et al (2015))

Interpretation:

- Gap between maximum likelihood and two-step estimators is evidence of misspecification

Also, overidentification tests, which provides evidence of misspecification (but not robust to heterogeneous effects!)

# Switching gears: Economists have a nose for randomness

- Paraphrasing a Yale prof:  
*Economists are really good at doing almost the right thing in empirical work.*  
-Anonymous Yale Professor
- Economists are clever at finding things that look convincingly “random”
  - Sometimes, it is easy to know how to use this randomness



Andy Luttrell  
@AndyLuttrell15

...

Dear Economists, how do you hear about these natural experiments occurring in the world? This seems like a thing economists are very good at. Do you just have a Google alert for the words "at random" or something?

4:58 PM · Sep 14, 2021 · Twitter Web App

# Borusyak and Hull (2022) on exploiting randomness in IV

- Two key parts to this paper:
  1. Highlighting how seemingly complicated research designs can be framed as generalized propensity scores
  2. How complicated research designs can suffer from *interference*
- There are many interesting results that spiral out from these two insights, but these are the key kernels (third piece is thinking about uncertainty using randomization inference, but deeply tied to other pieces))
- Will first start with showing how complicated research designs → propensity scores

## Non-Random Exposure to Exogenous Shocks: Theory and Applications

Kirill Borusyak  
UCL and CEPR

Peter Hull  
UChicago and NBER\*

January 2021

### Abstract

We develop new tools for estimating the causal effects of treatments or instruments that combine multiple sources of variation according to a known formula. Examples include treatments capturing spillovers in social and transportation networks, simulated instruments for policy eligibility, and shift-share instruments. We show how exogenous shocks to some, but not all, determinants of such variables can be leveraged while avoiding omitted variables bias. Our solution involves specifying counterfactual shocks that may as well have been realized and adjusting for a summary measure of non-randomness in shock exposure: the average treatment (or instrument) across such counterfactuals. We further show how to use shock counterfactuals for valid finite-sample inference, and characterize the valid instruments that are asymptotically efficient. We apply this framework to address bias when estimating employment effects of market access growth from Chinese high-speed rail construction, and to boost power when estimating coverage effects of expanded Medicaid eligibility.

# Research designs from simple to complex

- Consider the trivial research design, following an RCT that randomly assigns  $x_i \in \{0, 1\}$ , and we want to estimate the effect of  $x_i$  on  $y_i$ :

$$y_i = \alpha + x_i\beta + \epsilon_i$$

- The research design is effectively a coin flip:  $E(x_i) = p$ , and each  $x_i$  is independent for each  $i$ 
  - $\beta$  is identified thanks to this coin flip design
- This is true even when we have covariates,  $w_i$  that stratify the experiment. We just need to control for  $w_i$  correctly:  $E(x_i|w_i) = p(w_i)$  and we can estimate the ATE directly
- Effectively, the (potentially) endogenous  $w_i$  affects treatment, but if we condition correctly, we can still identify a causal effect





# Research designs from simple to complex- Medicaid eligibility

- Now imagine the eligibility rules for Medicaid were being randomly assigned
  - Drawn from a bag just like marbles, completely randomly
- We can now estimate the effect of Medicaid eligibility on things like child mortality
  - Issue: eligibility is also a function of many *endogenous* features
- We consider a known function,  $f_i$ , and eligibility rules,  $g_i$ , such that  $x_i = f(g_i, w_i)$  maps the  $w_i$  characteristics using the randomly drawn eligibility rules
  - Much like  $w_i$  strata case, but more complex b/c can be high-dimensional / non-linear



# Simulated instruments as a way to get a handle on this

- The challenge is that  $g_i$  is a complicated variable – it is a set of rules of that potentially complicated and hard to map to an “instrument” or “treatment”
- You don’t want to just use  $x_i$  because it contains endogenous  $w_i$
- Currie and Gruber (1996) solution: construct a variable  $z_i = \sum_j f(g_i, w_j)$  which takes  $w$  from a random population (outside the state) and uses it to construct a “predicted”  $x$ 
  - Intuitively, hold fixed the  $g_i$  and average over some distribution of  $w_j$

## **Saving Babies: The Efficacy and Cost of Recent Changes in the Medicaid Eligibility of Pregnant Women**

---

Janet Currie

*University of California, Los Angeles and National Bureau of Economic Research*

Jonathan Gruber

*Massachusetts Institute of Technology and National Bureau of Economic Research*

A key question for health care reform in the United States is whether expanded health insurance eligibility will lead to improvements in health outcomes. We address this question in the context of the dramatic changes in Medicaid eligibility for pregnant women that took place between 1979 and 1992. We build a detailed simulation model of each state’s Medicaid policy during this era and use this model to estimate (1) the effect of changes in the rules on the fraction of women eligible for Medicaid coverage in the event of pregnancy and (2) the effect of Medicaid eligibility changes on birth

# Simulated instruments as a way to get a handle on this

- The challenge is that  $g_i$  is a complicated variable – it is a set of rules of that potentially complicated and hard to map to an “instrument” or “treatment”
- You don’t want to just use  $x_i$  because it contains endogenous  $w_i$
- Currie and Gruber (1996) solution: construct a variable  $z_i = \sum_j f(g_i, w_j)$  which takes  $w$  from a random population (outside the state) and uses it to construct a “predicted”  $x$ 
  - Intuitively, hold fixed the  $g_i$  and average over some distribution of  $w_j$

To the extent that relevant state- and year-specific characteristics are not captured by state and year dummies (i.e., they are not constant within a state or across states within a year), the coefficient on the fraction eligible will be biased by omitted variables. Suppose, for example, that a state recession is associated with both increases in eligibility and a higher incidence of low birth weight. Then this source of variation in eligibility could induce a spurious positive correlation between Medicaid eligibility and low birth weight.

In order to overcome this potential problem, we instrument the actual fraction eligible with a measure of the generosity of Medicaid in a state and year that depends only on the state’s eligibility rules. To create our instrument, which we label the “simulated fraction eligible,” we first take a sample of 3,000 women from the CPS in each year. We then calculate the fraction of this sample of women who would be eligible for Medicaid in each state. By using the same group of women in each state simulation, we obtain an estimate of the fraction eligible that depends only on the legislative environment and is independent of other characteristics of states. This measure can be thought of as a convenient parameterization of legislative differences affecting women in different states and years: the generosity of state Medicaid policy can be naturally summarized in terms of the effect it would have on a given, nationally representative, population. Furthermore, we reduce the sampling variability in our estimates that derives from having relatively small cells for some states in the CPS.<sup>9</sup>

## This paper's approach vs. simulated instrument

- This is not the most efficient way to exploit this variation
- Remember our propensity score example: if we could just condition directly on  $w_i$ , then we would not worry about endogeneity
  - The solution, then, is to construct a propensity score and condition on that!
- Intuitively, “the eligibility rules for Medicaid were being randomly assigned”
  - In other words, we assert a counterfactual distribution over the policy rules  $Pr(g)$
  - This allows us to construct the propensity score for a given individual

$$p(w_i) = Pr(x_i | w_i) = \sum_g f_i(g, w_i) Pr(g)$$

- With pscore in hand, estimation is straightforward, and known to be semiparametrically efficient!
  - Either subtract this p-score off of the endogeneous variable (recentering) or control for it directly

# The return on propensity scores in an empirical example

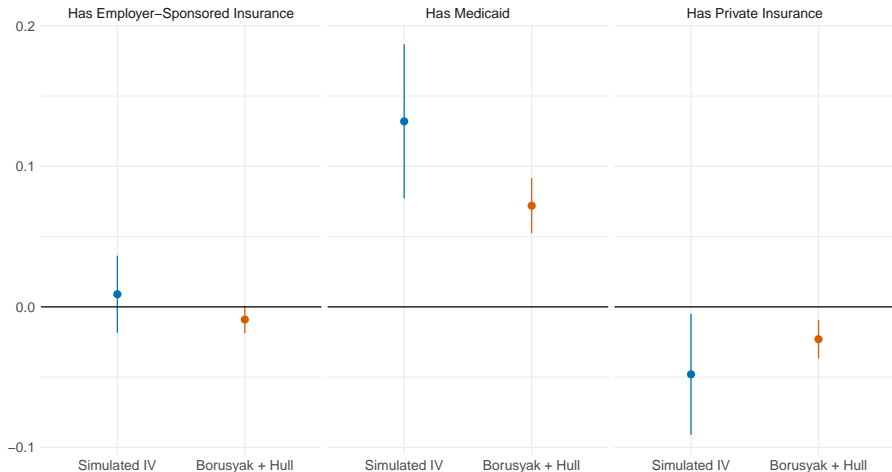
- Medicaid Empirical example in this paper: ACA medicaid expansion
- ACA expanded Medicaid in only some states thanks to NFIB v. Sebelius allowing choice by states
- Interested in understanding compositional shifts in health care across states
  - Use ACS micro data and consider structural equation

$$y_{it} = x_{it}\beta + \alpha_{s(i)} + \alpha_t + \epsilon_{it}$$

- $y_{it}$  are different health insurance take-up;  $x_{it}$  is Medicaid eligibility for individual  $i$
- “Simulated” IV: dummy for whether state expanded Medicaid  $z_{sim}$
- Borusyak and Hull IV: construct a person-level indicator for whether a person is eligible under their state’s law  $z_{bh}$ 
  - Also identify the  $p(w)$  that they are eligible on average across others states’ laws
  - They recenter ( $\tilde{z}_{bh} = z_{bh} - p(w)$ )

# The return on propensity scores in an empirical example

- **Much** more precise
- Makes sense!
- Seems like we should use it...



## Second kernel of the paper: interference

- Medicaid example is simple to think about, and clarifies idea that:
  1. Can convert high-dimensional variation into simple treatment effects
  2. Can be more *efficient* (e.g. smaller s.e.)
- However, you can take this much further.
- Consider the design of a railroad. Imagine the world in which a railroad designer randomly threw darts on a map to decide where to construct train lines
  - Similar to the analogy of “drawing” the Medicaid eligibility rules
  - But now, how do we think about the “random” piece interacting with different places?
  - Let’s start with something simple first



# Interference in network settings

- Consider a setting where the researcher want to measure the impact of a randomized experiment on a network
  - In other words, for a given person  $i$ , and observed network  $W$ , we randomly treat some subset of individuals on the network.
  - We want to know what the effect of having more treated individuals connected to you  $x_i$  is on  $y_i$
- Insight from paper: since the position in *network* affects probability of being connected to individuals, some individuals will inherently get more exposure!
  - Analogous to the friendship paradox
- Need to construct an analogous propensity score for the network setting, and control for that
  - Since we have a true RCT, this is not too hard!
  - (But we do have to make decisions about what the spillover is)
  - Aronow and Samii (2017) made serious progress on the network context



# Interference in spatial settings

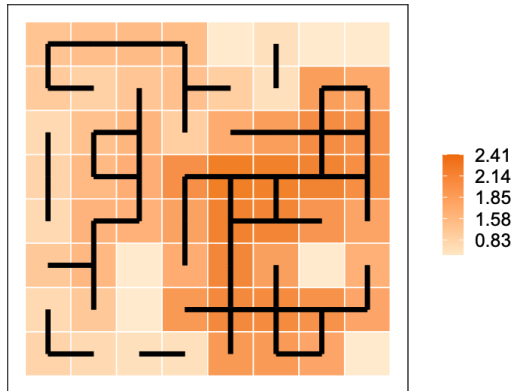
- Things can be more confusing than an RCT, but this same insight applies
  - Even with random shocks (darts on a board), some locations / people attract more treatment than others
  - Consider the application from the paper
- Estimate the impact of market access growth ( $MA$ ) on land values growth ( $V$ ) in China
  - $MA$  is influenced by transportation networks, and measures aggregated access to other populations

$$MA_{it} = \sum_j \tau(\mathbf{g}_t, loc_i, loc_j)^{-1} pop_j$$

- Want to estimate the effect of  $MA_{it}$  using “random” variation in network changes!
  - Can we just run the OLS? No!

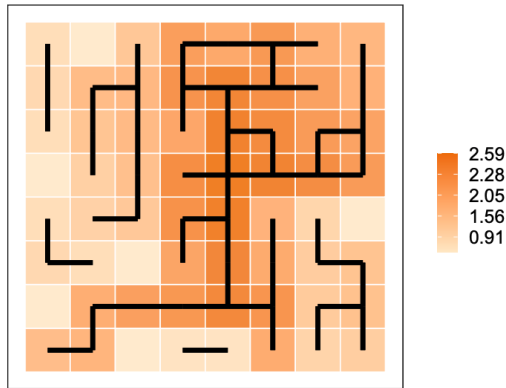
# Stylized Example of Market Access on a Square Island

- Take a square with square villages and randomly assign roads
- How does market access change?



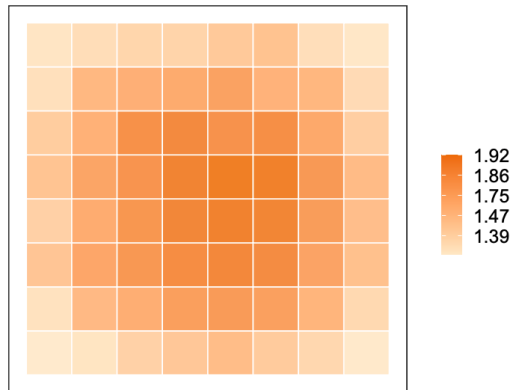
# Stylized Example of Market Access on a Square Island

- Take a square with square villages and randomly assign roads
- How does market access change?
- If we rerandomize, does it look different?



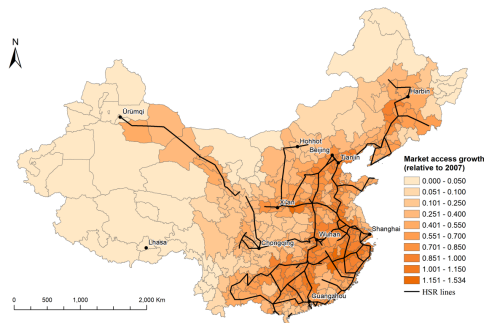
# Stylized Example of Market Access on a Square Island

- Take a square with square villages and randomly assign roads
- How does market access change?
- If we rerandomize, does it look different?
- As with the network, some places get more market access than others on average!
- Need to account for this propensity difference



# China: defining the counterfactual distribution

- In the stylized example, lines are laid randomly, making it easy to define the propensity scores
  - What about in China?
- What is the plausible counterfactual?



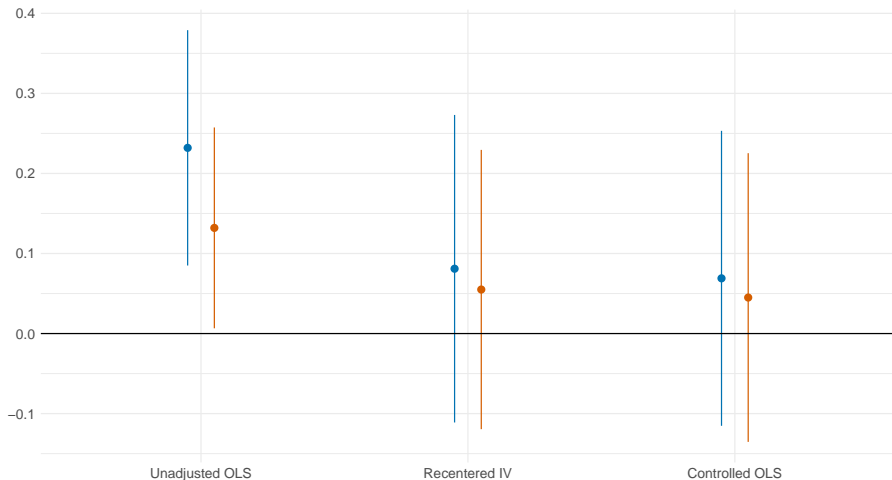
# China: defining the counterfactual distribution

- In the stylized example, lines are laid randomly, making it easy to define the propensity scores
  - What about in China?
- What is the plausible counterfactual?
- Paper proposes an idea, and analogous to other examples
  - Use *planned* lines are randomized between unbuilt but planned, and built lines
  - Calculate distribution of propensity score by constructing  $MA_i$  under each counterfactual scenario



# China railroads: the punchline

- There was substantial bias from using OLS!
- Makes sense – geography is king...
- No effect in randomized setting



## Defining the counterfactual distribution

- If one takes issue with the counterfactuals, that is reasonable (but of course, challenging to prove one way or the other)
- Key issue: this paper is just making **text** what was already **subtext**
  - There was always an assumption about some counterfactual comparison in these designs!
- The issue is that many of these paper do not understand how to describe the randomization aspect of their research design
  - Consequentially, they cannot describe the “as-if random” component coherently
  - If a researcher has an alternative proposal, they should try that and see what estimates are available!
- Also suggests that reserchers can show a “range” of estimates under different scenarios



## Key takeaways from paper

- Provide a toolbox for contexts when economists have found good “as-if” random variation (and can describe the counterfactual distribution)
- Show that in cases where treatment is not influenced by others’ treatment status, approach maps very tightly with traditional propensity methods, and can be much more efficient
- In spatial and network cases where treatment spillovers exist, show how to adjust for bias arising from units location on network or graph (or relevant characteristic)

# Caveats

- We focused on Currie and Gruber (1996a, b) cases of simulated instruments, but there are other cases of “simulated instruments”
  - Gruber and Saez is all about characteristics (income) responding endogenously
- These are tax elasticities and not clear that exclusion restriction holds
- This tax literature is strongly tied to functional form or additional assumptions (see Blomquist, Newey, Kumar and Liang (2021))

# Granular Instruments (Gabaix and Koijen)

- Note that the previous approach heavily leaned on considering a “random” distribution of exogenous shocks, independent of the model
- Granular IV (Gabaix and Koijen (2022)) similarly relies on random shocks, but uses model structure to generate instruments for demand and supply estimation

## Granular Instrumental Variables\*

Xavier Gabaix and Ralph S.J. Koijen

June 6, 2022

### Abstract

We propose a new way to construct instruments in a broad class of economic environments. In the economies we study, a few large firms, industries or countries account for an important share of economic activity. As the idiosyncratic shocks from these large players affect aggregate outcomes, they are valid and often powerful instruments. We provide a methodology to extract idiosyncratic shocks from the data and create “granular instrumental variables” (GIVs), which are size-weighted sums of idiosyncratic shocks. These GIVs allow us to then estimate parameters of interest, including causal elasticities and multipliers. We illustrate the idea in a basic supply and demand framework. GIVs provide a novel approach to identify both supply and demand elasticities based on idiosyncratic shocks to either supply or demand. We then show how to extend the basic procedure to cover a range of empirically relevant situations. As an application, we measure how “sovereign yield shocks” transmit across countries in the Eurozone. We sketch how GIVs could be useful to estimate a host of other causal parameters in economics.

# Simplified Setup from Granular Instruments

- Goal of the paper is to estimate a structural model of the following setup:

$$y_{it} = \phi^d p_t + \gamma^y X_{1,it} + \underbrace{\lambda_i \eta_t}_{\text{Unobserved}} + u_{it} \quad (1)$$

$$p_t = \psi y_{St} + \gamma^p X_{2,t} + \epsilon_t \quad (2)$$

$$y_{St} = \sum_i \underbrace{S_i}_{\text{Size Weights}} y_{it}. \quad (3)$$

- Effectively, we want to estimate the elasticities for this system:

# Setup from Granular Instruments

- Goal of the paper is to estimate a structural model of the following setup:

$$y_{it} = \phi^d p_t + \gamma^y X_{1,it} + \underbrace{\lambda_i \eta_t}_{\text{Unobserved}} + u_{it} \quad (4)$$

$$p_t = \psi y_{St} + \gamma^p X_{2,t} + \epsilon_t \quad (5)$$

- A few things to note:
  1. OLS is biased for either regression (can solve and show endogeneity bias)
  2. Equations are estimated in different aggregation levels (key!)
- Key assumptions for the paper:
  - $u_{it}$  are completely random:  $\mathbf{u}_t \perp \eta_t, \gamma^y X_{1,it}, \gamma^p X_{2,t}$
  - Models are correctly specified

## Simplified Setup from Granular Instruments

- Simplified application:

$$y_{it} = \phi^d p_t + \eta_t + u_{it} \quad (6)$$

$$p_t = \psi y_{St} + \epsilon_t \quad (7)$$

- Aggregate risk exposure is identical ( $\eta_t$ )
- What happens if aggregate  $y_{it}$  within-time period?
  - $y_{St} = \sum_i S_i y_{it}$ . Everything is identical across firms except  $u_{it}$ .
  - Define  $z_t = y_{St} - y_{Et} = \sum_i S_i y_{it} - n^{-1} \sum_i y_{it} = u_{St} - u_{Et}$
  - By assumption,  $z_t$  is independent of  $\epsilon_t$ !
  - By assumption,  $E(z_t \eta_t) = 0$  AND  $E(z_t u_{Et}) = 0$

## Simplified Setup from Granular Instruments

$$y_{it} = \phi^d p_t + \eta_t + u_{it} \quad (8)$$

$$p_t = \psi y_{St} + \epsilon_t \quad (9)$$

- Conceptually, this model has no external instruments, but if it is correctly specified, it can identify residuals which are assumed to be independent
- E.g. consider the reduced forms:

$$y_{St} = (1 - \phi^d \psi) \phi^d \epsilon_t + (1 - \phi^d \psi) \eta_t + (1 - \phi^d \psi) u_{St} \quad (10)$$

$$p_t = (1 - \phi^d \psi)^{-1} \psi \eta_t + (1 - \phi^d \psi)^{-1} \psi u_{St} + (1 - \phi^d \psi)^{-1} \epsilon_t \quad (11)$$

## Simplified Setup from Granular Instruments

$$y_{it} = \phi^d p_t + \eta_t + X_{it}\delta + u_{it} \quad (12)$$

$$p_t = \psi y_{St} + \epsilon_t \quad (13)$$

- How can this break? What if there is a time-varying characteristic we do not control for?
- Then,  $z_t$  will capture this size weighted component – if it's correlated with the average shocks on either side, that will be problematic.



# Application for Granular Instruments (Gabaix and Koijen 2022b)

- Shocks to mutual fund flows as a source of variation in demand

## In Search of the Origins of Financial Fluctuations: The Inelastic Markets Hypothesis

Xavier Gabaix and Ralph S.J. Koijen\*

May 12, 2022

### Abstract

We develop a framework to theoretically and empirically analyze the fluctuations of the aggregate stock market. Households allocate capital to institutions, which are fairly constrained, for example operating with a mandate to maintain a fixed equity share or with moderate scope for variation in response to changing market conditions. As a result, the price elasticity of demand of the aggregate stock market is small, and flows in and out of the stock market have large impacts on prices.

Using the recent method of granular instrumental variables, we find that investing \$1 in the stock market increases the market's aggregate value by about \$5. We also develop a new measure of capital flows into the market, consistent with our theory. We relate it to prices, macroeconomic variables, and survey expectations of returns.

We analyze how key parts of macro-finance change if markets are inelastic. We show how general equilibrium models and pricing kernels can be generalized to incorporate flows, which makes them amenable to use in more realistic macroeconomic models and to policy analysis.

Our framework allows us to give a dynamic economic structure to old and recent datasets comprising holdings and flows in various segments of the market. The mystery of apparently random movements of the stock market, hard to link to fundamentals, is replaced by the more manageable problem of understanding the determinants of flows in inelastic markets. We delineate a research agenda that can explore a number of questions raised by this analysis, and might lead to a more concrete understanding of the origins of financial fluctuations across markets.

# Application for Granular Instruments (Gabaix and Koijen 2022b)

- Shocks to mutual fund flows as a source of variation in demand

**GIV: Requirements and threats to identification** For the GIV to be consistent, we need  $\mathbb{E}[u_{it}\eta_t] = 0$  to hold: the idea is that there are random “bets” or “shocks” to various fund managers, institutions and sectors, that are orthogonal to all reasonable common macro factors such as GDP, TFP, and so forth. For the GIV to be a powerful instrument, we need large idiosyncratic shocks, and a few large sectors, so that the market is “granular” in the sense that the idiosyncratic shocks to a few large sectors meaningfully affect the aggregate.<sup>44</sup> Fortunately, this is verified in our setting, as it is in related settings in macro (Gabaix (2011), Carvalho and Grassi (2019)), trade (Di Giovanni and Levchenko (2012)) or finance (Amiti and Weinstein (2018), Herskovic et al. (forthcoming), Galaasen et al. (2020)). Ben-David et al. (forthcoming) and Ghysels et al. (2021) study the impact of investor granularity on the cross-section of US stock returns.

The main threats to identification with GIV are that we do not properly control for common factors, or that the loadings on the omitted factor are correlated with size, such that  $\lambda_S - \lambda_E \neq 0$ . To mitigate the risk of omitted factors, we extract additional factors and explore the stability of the estimates as we add extra factors.

# Application for Granular Instruments (Gabaix and Koijen 2022b)

- Shocks to mutual fund flows as a source of variation in demand

Intuitively, we use the sector-specific, or idiosyncratic, demand shocks of one sector as a source of exogenous price variation to estimate the demand elasticity of another sector. Viewed this way, the GIV estimator generalizes the idea behind the index inclusion literature to estimate the micro elasticity. In the index inclusion literature, a demand shock to the group of index investors (assuming the inclusion of a stock into the index is random) can be used to estimate the slope of the demand curve of the non-index investors.