

Bartik工具变量与模拟的工具变量

慧航

2024年5月

Bartik工具

Bartik工具变量，也称为移动-份额工具（shift-share instruments）最初由Bartik（1991）使用并讨论，在Blanchard和Katz（1992）之后更加流行，目前已近被应用于很多领域：

- 移民：Altonji和Card（1991）、Card（2001）
- 银行借贷：Amiti和Weinstein（2018）、Greenstone, Mas和Nguyen（2015）
- 市场规模与创新：Acemoglu和Linn（2004）
- 贸易：Autor, Dorn和Hanson（2013,2018）、Piveteau和Smagghue（2017）、de Roux等（2017）
- 国外援助：Nunn和Qian（2014）
- 自动化：Acemoglu和Restrepo（2017）

基本设定

基本（简单）设定

- 结构方程:

$$y_i = \beta_0 + \beta w_i + u_i$$

其中 $i = 1, \dots, N$ 为个体（通常为地区），

- y_i 为地区的被解释变量，为了保证接下来的外生性，通常被解释变量为相对值（比如差分、增长率等）
 - $\text{Cov}(w_i, u_i) \neq 0$ 从而有内生性
 - 假设 w_i 有如下的结构：

$$w_i = \sum_{m=1}^M z_{im} g_{im}$$

其中 z_{im} 为*i*地区第*m*个成分的份额, g_{im} 为增长率

基本设定

估计劳动供给弹性

- y : 工资增长率
- w : 就业增长率
- z_{im} : i 地区行业 m 的份额
- g_{im} : i 地区行业 m 的工资增长率

估计移民的影响

- y : 本地和移民工资差异
- w : 移民增长率
- z_{im} : i 地区来自国家 m 的移民份额
- g_{im} : i 地区自国家 m 的人口增长率

Bartik工具的构建

- 现在假定：

$$g_{im} = g_m + \tilde{g}_{im}$$

其中 g_m 为第 m 个成分（行业、国家等）的平均增长率

- 构建Bartik工具变量：

$$B_i = \sum_{m=1}^M z_{im} g_m$$

- 问题：何时 B_i 是一个合格的工具变量？两个条件：
 - $\text{Cov}(B_i, w_i) \neq 0$
 - $\text{Cov}(B_i, u_i) = 0$
- 两只文献：
 - Goldsmith-Pinkham, Sorkin和Swift (2020) 研究了份额 z_{im} 的外生性
 - Borusyak、Hull和Jaravel (2022) 则研究了冲击 g_m 的外生性

一个特例: $M = 2$

- 如果 $M = 2$, 那么:

$$B_i = z_{i1}g_1 + z_{i2}g_2$$

由于 $z_{i1} + z_{i2} = 1$, 从而:

$$B_i = g_2 + (g_1 - g_2)z_{i1}$$

- 第一阶段:

$$\begin{aligned} w_i &= \gamma_0 + \gamma B_i + \epsilon_i \\ &= \underbrace{\gamma_0 + \gamma g_2}_{\text{常数项}} + \underbrace{\gamma(g_1 - g_2)z_{i1}}_{\text{系数}} + \epsilon_i \end{aligned}$$

- 相关性要求 $g_1 - g_2 \neq 0$: $g_1 - g_2$ 看做是某一个“政策”变动,

一个特例: $M = 2$

- 真正的工具: z_{i1} !
 - z_{i1} 可以看做是政策的暴露水平
 - 问题: z_{i1} 会不会直接影响 y_i ?
- 一般而言, 无法排除 z_{i1} 作为份额与 y_i 的直接关系
 - 在劳动供给弹性的例子中, z_{i1} 为某个行业的就业份额, y_i 为收入, 显然这是长期均衡共同决定的
 - 在移民的例子中, 不同国家的收入不一样, 显然当地收入也会影响不同国家的移民比例
 - 事情不一定这么悲观: y_i 并不是收入, 而是收入增长率(变动, 而非水平), 外生性也许会满足

M不受限、截面

- 现在考虑控制变量：

$$y_i = \beta_0 + \beta w_i + x'_i \eta + u_i$$

记 $Y^\perp = Y - \mathbb{L}(Y|X)$, $W^\perp = W - \mathbb{L}(W|X)$

- 如果有 M 个行业，那么 Bartik 工具的 2SLS 估计：

$$\hat{\beta}^{Bartik} = \frac{B' Y^\perp}{B' W^\perp}$$

- 同时，如果使用 $M - 1$ 个份额作为工具变量，那么：

$$\hat{\beta}^{GMM} = \frac{W^{\perp'} Z \Xi Z' Y^\perp}{W^{\perp'} Z \Xi Z' W^\perp}$$

其中 Ξ 为权重矩阵

M不受限、截面

定理

定理：如果 $\Xi = GG'$, 那么 $\hat{\beta}^{Bartik} = \hat{\beta}^{GMM}$

启示

Bartik工具变量估计的结果等价于使用份额 z_{im} 作为工具变量，并没有使用 g_m 的外生性！

面板情形

- 对于面板设定（首先忽略控制变量）：

$$y_{it} = \beta_t + \beta w_{it} + u_{it}$$

其中 β_t 为时间固定效应

- Bartik工具：

$$B_{it} = \sum_{m=1}^M z_{im0} g_{mt}$$

其中 z_{im0} 为初始的份额，而增长率：

$$g_{imt} = g_{mt} + \tilde{g}_{imt}$$

特例: $M = 2, T = 2$

- 此时工具:

$$B_{it} = g_{1t}z_{i10} + g_{2t}z_{i20} = g_{2t} + (g_{1t} - g_{2t})z_{i10}$$

- 第一阶段:

$$\begin{aligned} w_{it} &= \gamma_t + \gamma B_{it} + \epsilon_{it} \\ &= \underbrace{\gamma_t + \gamma g_{2t}}_{\tilde{\tau}_t} + \gamma(g_{1t} - g_{2t})z_{i10} + \epsilon_{it} \end{aligned}$$

如果将其用两期的虚拟变量写出来:

$$w_{it} = \tilde{\tau}_t + \underbrace{\gamma(g_{11} - g_{21})}_{\tilde{\gamma}_1} \mathbf{1}\{t=1\} z_{i10} + \underbrace{\gamma(g_{12} - g_{22})}_{\tilde{\gamma}_2} \mathbf{1}\{t=2\} z_{i10} + \epsilon_{it}$$

其中 $\tilde{\gamma}_1, \tilde{\gamma}_2$ 为份额 z_{i10} 与时间固定效应的交互项的系数

- 等价关系: 特定权重矩阵下, Bartik工具的估计结果与使用“份额与时间固定效应的交互”作为工具变量是等价的。

Bartik与DID

- 如果将 $g_{1t} - g_{2t}$ 视为政策冲击，而 z_{i10} 视为政策暴露水平，以上设定与DID关系密切！
 - 具有高 ($m = 1$) 份额的地区与低 ($m = 1$) 份额的地区相比，对于 w 的增长率冲击 $g_{1t} - g_{2t}$ 的反应是不是不同的？
 - 如果高低地区的 w 反应不同、 y 反应也不同，那么推断 $w > y$ 有影响
- 问题：平行趋势？分组变量 (z_{i10}) 外生吗？
 - 如果存在政策冲击，比如 $g_{11} - g_{21} = 0$ ，那么可以检验 z_{i10} 是否对事前的 y 有影响：平行趋势检验
 - 直觉：如果 $g_{11} - g_{21} = 0$ ，且 z_{i10} 不直接影响 y ，那么 z_{i10} 由于不影响事前 w ，从而也不影响事前的 y ；如果有影响： z_{i10} 不外生！

正式的识别假设

- 不失一般性用面板设定：

$$y_{it} = \beta w_{it} + x'_{it}\eta + u_{it}$$

其中 x_{it} 包含时间、个体固定效应等。

假设1

对于所有的 s, m , 回归:

$$w_{it} = x'_{it}\zeta + \sum_m \sum_s \delta_{ms} \times 1\{t = s\} \times z_{im0} + e_{it}$$

其中 δ_{ms} 有限且 $\sum_s \sum_m g_{ms} \delta_{ms} \neq 0$ 。

假设2

对于所有的 $g_m \neq 0$ 的 m , 有: $\mathbb{E}(u_{it} z_{im0} | x_{it}) = 0$ 。

打开Bartik工具的黑箱

- 根据以上介绍的等价关系，虽然 B_{it} 是一个工具变量，但是实际上等价于很多份额作为工具变量
 - 截面：份额作为工具变量
 - 面板：份额与时间固定效应的乘积作为工具变量
- 所以Bartik工具综合了很多很多工具变量的结果
 - 可以分解！
 - 可以直接使用这些多工具的结果进行估计，同时进行过度识别检验

Bartik工具变量估计的分解

Rotemberg权重（截面）

Bartik工具变量的结果可以分解为：

$$\hat{\beta}^{Bartik} = \sum_m \hat{\alpha}_m \hat{\beta}^m$$

其中 $\hat{\beta}^m$ 为使用份额 m 作为工具的 2SLS 估计结果， $\hat{\alpha}_m$ 为 Rotemberg 权重，且

$$\sum_m \hat{\alpha}_m = 1$$

Bartik工具变量估计的分解

注意！

权重 $\hat{\alpha}_m$ 有可能为负！从而：

- 如果同质性影响， $\hat{\beta}^m$ 都几乎相同，无影响
- 如果有异质性影响， $\hat{\beta}^{Bartik}$ 可能有问题
 - 使用bartik_weight (<https://github.com/paulgp/bartik-weight>) 计算Rotemberg权重
 - 查看Rotemberg权重比较大的 m 的估计结果
 - 根据 $\hat{\alpha}_m$ 的正负分组，并观察两个组别权重的和

Rotemberg权重

- 可以计算 $\hat{\alpha}_m$ 与 g_m 的相关性： $\hat{\alpha}_m$ 多大程度上可以由 g_m 解释
- 可以计算 $\hat{\alpha}_m$ 与 z_{ik} 的方差的相关性，其中的方差按照：

$$\frac{1}{N-1} \sum_i (z_{ik} - \bar{z}_k)^2$$

计算

- 同时可以计算 $\hat{\alpha}_m$ 与第一阶段F统计量 \hat{F}_k 的相关性：第一阶段F大的工具不一定权重大
- 对于面板，同理，不过由于此时工具变量有 $(M-1) \times T$ 个，所以需要加总：

$$\hat{\alpha}_m = \sum_t \hat{\alpha}_{m,t}$$

Bartik工具所需要做的一些检验

为了保证Bartik工具能够正确识别，一系列检验需要进行：

- 检验1：检查份额 z_{im0} 与初始期的特征的相关性
 - 真正重要的是 z_{im0} 能否直接预测 y （的变化、增长率等）
 - 如果 z_{im0} 与一些能够预测 y 变化的因素相关，可能会有遗漏变量的偏误
- 检验2：如果有政策变化——平行趋势检验
- 检验3：
 - 其他的IV估计量（LIML、MB2SLS、HFUL等）结果是否大致相同
 - 过度识别检验（直接使用份额做IV）
 - 异质性影响：同样会使得以上两个检验不通过：直接查看 $\hat{\beta}^m$

Bartik工具：另一种视角

- 以上介绍Bartik工具聚焦份额的外生性，而Borusyak、Hull和Jaravel (2022) 考虑了 g_m 的外生性
- 假设：
 - g_m 外生
 - g_m 之间不相关
 - $M \rightarrow \infty$
- 根据大数定律，Bartik工具仍然是一致的
- 然而标准误可能有问题 (Adao, Kolesar和Morales, 2019)
 - BHJ重新提出了稳健估计量
 - Stata: ssaggregate

合成工具变量

- 工具变量（解释变量）需要满足外生性，比如一些自然实验的冲击
- 然而有时，有些冲击不得不与一些内生的变量合在一起才能被观察到：
 - Bartik工具：可能 g_m 是外生的，但是 z_{im} 的内生的
 - 交通设施：哪些地方开通高铁可能不是外生的，但是开通高铁的时间可能是外生的
 - 美国每个州的（给穷人的）健康保险政策是外生的，但是人口特征不是（比如可能跟健康状况有关）
- 问题：这些工具变量是否可以使用？如何使用？Borusyak和Hull (2021) 回答了这一问题。

一个示例：市场进入

- 考虑交通基础设施（市场进入， MA_i ）对土地价格 V_i 的影响：

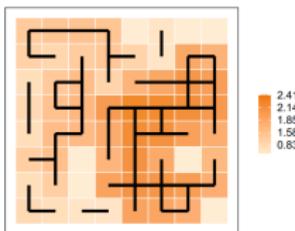
$$\Delta \ln V_i = \beta_0 + \beta \Delta \ln MA_i + u_i$$

- 假设在一个棋盘上随机取点修建铁路
- 虽然哪个点修建铁路是外生的，但是地理位置是内生的！
 - 棋盘中间天然具有更高的市场进入
- 解决办法：
 - 如果棋盘上的点是随机被选取出来修建铁路的，计算一个期望的市场进入 $\mu_i = \mathbb{E}(\Delta \ln MA_i)$
 - 使用 $\Delta \ln MA_i - \mu_i$ 作为解释变量
 - μ_i : 可以通过模拟计算得到

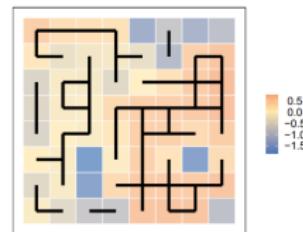
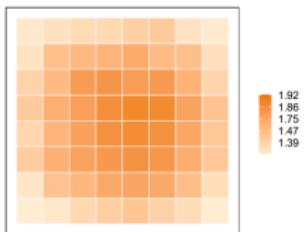
市场进入示例

Figure 1: Market Access Growth in the Motivating Example

A. Line Construction and Market Access Growth



B. Expected Market Access Growth C. Recentered Market Access Growth



设定

- 假设：

$$z_i = f_i(g, x)$$

其中：

- $f_i(\cdot), i = 1, \dots, N$ 为已知函数
- g 为 $N \times 1$ 的外生冲击
- x 为一些可观测变量

- 回归方程：

$$y_i = \beta_0 + \beta w_i + u_i$$

重要引理

假设1

假设冲击是外生的，即 $g \perp\!\!\!\perp u|x$

在以上条件下，由于：

$$\begin{aligned}\mathbb{E} \left(\frac{1}{N} \sum_i z_i u_i \right) &= \mathbb{E} \left(\frac{1}{N} \sum_i f_i(g, x) u_i \right) \\&= \mathbb{E} \left(\mathbb{E} \left[\frac{1}{N} \sum_i f_i(g, x) u_i | x \right] \right) \\&= \mathbb{E} \left(\frac{1}{N} \sum_i u_i \mathbb{E} [f_i(g, x) | x] \right) \\&= \mathbb{E} \left(\frac{1}{N} \sum_i u_i \mu_i \right)\end{aligned}$$

其中 $\mu_i = \mathbb{E} [f_i(g, x) | x]$ 为工具变量的期望。

新的工具

- 根据以上结论，有：

$$\mathbb{E} \left(\frac{1}{N} \sum_i (z_i - \mu_i) u_i \right) = 0$$

- 两种做法：

- 计算得到 μ_i 使用 $z_i - \mu_i$ 作为工具变量：重新中心化
(recentering)
- 直接将 μ_i 作为控制变量、 z_i 作为工具变量
 - 启示：如果 μ_i 被其他控制变量、固定效应吸收了，不用做额外处理，使用 z_i 做工具就够了。

工具期望的计算

问题是：如何计算 μ_i ？

- 假设我们知道给定 x , g 的分布 $G(g|x)$, 那么可以直接计算:

$$\mu_i = \int f_i(g; x) dG(g|x)$$

- 然而 $G(\cdot|\cdot)$ 可能不好设定: 直接将 g 向量进行重新排列组合, 得到新的组合 $\pi(g)$, 然后计算 $f_i(\pi(g); x)$, 重复多次后计算均值即可。
 - 比如对于高铁开通, 如果使用2016年的数据, 假设开通时间是随机的, 那么将2016已经开通的城市和有规划的城市进行一个排列组合, 取MA的均值即可。
 - 对于保险资格, 假设某个个体随机属于某一个州, 按照这个州的标准判断其在这个州是不是有资格, 然后求平均