2021 中级计量经济学作业 2

1. (关于分块回归) 已知 $y = \beta_0 + \beta_1 x_1 + \cdots + \beta_K x_K + \varepsilon$. 证明 OLS 回归系数估计 b_K 满足

$$b_K \to \beta_K + \frac{\operatorname{Cov}(\tilde{x}_K, \varepsilon)}{\operatorname{Var}(\tilde{x}_K)},$$

其中 $\tilde{x}_K = x_K - L(x_K \mid 1, x_1, \dots, x_{K-1}), L(x_K \mid 1, x_1, \dots, x_{K-1})$ 是 x_K 在 $1, x_1, \dots, x_{K-1}$ 上的 Linear Projection. (假设所需的秩条件和大数定律满足)

证明: 首先, 定义 $\mathbf{M}_{(K)}$ 为除 x_K 之外的自变量 $1, x_1, \ldots, x_{K-1}$ 的消灭矩阵.

$$egin{aligned} \widetilde{\mathbf{y}} &= \mathbf{M}_{(K)} \mathbf{y} \ &= \mathbf{M}_{(K)} (\mathbf{X} oldsymbol{eta} + oldsymbol{arepsilon}) \ &= \mathbf{M}_{(K)} \left[\mathbf{X}_{(K)} \quad \mathbf{X}_K
ight] \left[eta_{(K)}^{oldsymbol{eta}_{(K)}} + \mathbf{M}_{(K)} oldsymbol{arepsilon} \ &= \mathbf{M}_{(K)} \left(\mathbf{X}_{(K)} oldsymbol{eta}_{(K)} + \mathbf{X}_K eta_K
ight) + \mathbf{M}_{(K)} oldsymbol{arepsilon} \ &= \widetilde{\mathbf{X}}_K eta_K + \mathbf{M}_{(K)} oldsymbol{arepsilon} \end{aligned}$$

其次, 由分块回归知

$$b_{K} = \left(\widetilde{\mathbf{X}}_{K}'\widetilde{\mathbf{X}}_{K}\right)^{-1} \left(\widetilde{\mathbf{X}}_{K}'\widetilde{\mathbf{y}}\right)$$

$$= \left(\widetilde{\mathbf{X}}_{K}'\widetilde{\mathbf{X}}_{K}\right)^{-1} \left(\widetilde{\mathbf{X}}_{K}'(\widetilde{\mathbf{X}}_{K}\beta_{K} + \mathbf{M}_{(K)}\boldsymbol{\varepsilon})\right)$$

$$= \beta_{K} + \left(\widetilde{\mathbf{X}}_{K}'\widetilde{\mathbf{X}}_{K}\right)^{-1} \left(\widetilde{\mathbf{X}}_{K}'\mathbf{M}_{(K)}\boldsymbol{\varepsilon}\right)$$

$$\to \beta_{K} + \mathbf{E}[\widetilde{\mathbf{x}}_{K}\widetilde{\mathbf{x}}_{K}']^{-1} \mathbf{E}[\widetilde{\mathbf{x}}_{K}\boldsymbol{\varepsilon}].$$

注意这里 \tilde{x}_K 是标量, 因此 $b_K \to \beta_K + \mathrm{E}[\tilde{x}_K^2]^{-1}\mathrm{E}[\tilde{x}_K \varepsilon]$. 同时, 因为 $\mathbf{X}_{(K)}$ 中包含常数项, $\mathrm{E}[\tilde{x}_K] = 0$, 因此 $b_K \to \beta_K + \mathrm{Var}(\tilde{x}_K)^{-1}\mathrm{Cov}(\tilde{x}_K, \varepsilon)$.

2. 在同方差假设下, 比较 OLS 回归和 2SLS 回归的渐进方差, 说明哪一个更加有效 (more efficient).

OLS 和 2SLS 的渐进方差分别为:

$$\operatorname{Avar}(\sqrt{N}(\hat{\boldsymbol{\beta}}_{\text{OLS}} - \boldsymbol{\beta})) = \sigma^{2} \operatorname{E}[\mathbf{x}\mathbf{x}']^{-1}$$
$$\operatorname{Avar}(\sqrt{N}(\hat{\boldsymbol{\beta}}_{2\text{SLS}} - \boldsymbol{\beta})) = \sigma^{2} \operatorname{E}[\hat{\mathbf{x}}\hat{\mathbf{x}}']^{-1},$$

其中 $\hat{\mathbf{x}}$ 是第一阶段用工具变量对 \mathbf{x} 回归得到的 Linear projection, 因此 $\mathbf{x} = \hat{\mathbf{x}} + \mathbf{r}$, $\mathrm{E}(\hat{\mathbf{x}}\mathbf{r}') = \mathbf{0}$. 因此

$$E(\mathbf{x}\mathbf{x}') = E(\mathbf{\hat{x}}\mathbf{\hat{x}}' + \mathbf{\hat{x}}\mathbf{r}' + \mathbf{r}\mathbf{\hat{x}}' + \mathbf{r}\mathbf{r}')$$
$$= E(\mathbf{\hat{x}}\mathbf{\hat{x}}' + \mathbf{r}\mathbf{r}')$$

因此,

$$\sigma^2 \mathrm{E}[\hat{\mathbf{x}}\hat{\mathbf{x}}']^{-1} - \sigma^2 \mathrm{E}[\mathbf{x}\mathbf{x}']^{-1}$$

是半正定的.

3. 考虑如下关于吸烟是否影响新生儿体重的模型:

 $\log(bwght) = \beta_0 + \beta_1 male + \beta_2 parity + \beta_3 \log(faminc) + \beta_4 packs + \varepsilon,$

其中 log(bwght) 是新生儿体重的对数, male 是新生儿是否为男孩的虚拟变量, parity 是新生儿出生的顺序, faminc 是家庭收入, packs 是孕妇在怀孕期每天吸烟的平均数.

- 1. packs 会否和 ε 相关? 为什么?
 - 未观察到的其他影响健康的(从而影响婴儿体重)的因素可能和 packs 相关。比如,抽烟多的人可能也酗酒,或者饮食不健康,等等。这些遗漏变量会包括在 ε 中。
- 2. 假设你取得了样本中妇女所在省的香烟平均价格 cigprice 作为 packs 的工具变量. 讨论 cigprice 能否满足工具所需的两个条件.

相关性:需求理论会认为 packs 和 cigprice 负相关,但在实际数据中 cigprice 和 packs 相关性未必会大,因为仍需控制影响 packs 的其他变量,而这些变量不可得。例如,对香烟成瘾的程度决定了 packs 和 cigprice 的需求弹性,而成瘾程度数据不可得。同时,这里的香烟价格是省的平均价格,也会丢失了很多地区微观的变化。

外生性: cigprice 可能是外生的,因为香烟价格直觉上不会直接和婴儿体重相关。然而,香烟的价格受税率的影响很大,税率低的地区可能医疗条件等公共服务的水平也相对低,从而影响母亲的健康状况(包含在 ε 当中),从而也有可能不是外生的。

3. 用 BWGHT.DTA 数据估计模型. 先用 OLS. 再用 2SLS, 用 cigprice 作为 packs 的工具变量. 结果是否有很大的不同? (可直接附上 Stata 回归结果, 无需考虑格式)

Source	SS	d f	MS	Number of obs	=	1,38
				F(4, 1383)	=	12.55
Model	1.76664363	4	.441660908	Prob > F	=	0.000
Residual	48.65369	1,383	.035179819	R-squared	=	0.0350
				Adj R-squared	=	0.0322
Total	50.4203336	1,387	.036352079	Root MSE	=	.18756
lbwght	Coef.	Std. Err.	t	P> t [95% Co	nf.	Interval
			2 44	0.009 .006448	86	.0460328
male	.0262407	.0100894	2.60	0.009 .000446	, ,	
male parity	.0262407 .0147292	.0100894 .0056646		0.009 .003617	-	.0258414
			2.60		1	
parity	.0147292	.0056646	2.60	0.009 .003617	71	.0258414 .0290032 0501423

图 1

OLS 和 2SLS(这里是 IV) 回归的结果相差很大。OLS 中,多一包烟可减少新生儿体重 8.4%,且显著。IV 中,packs 的系数很大,是正数,但统计上不显著。这个结果不符合常识。

4. 作第一阶段回归, cigprice 和 packs 的相关性是否足够强? 这个结果如何影响了 2SLS 回归?

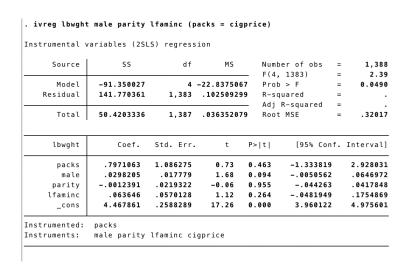


图 2

Source	SS	d f	MS	Number of obs	=	1,388
				F(4, 1383)	=	10.8
Model	3.76705108	4	.94176277	Prob > F	=	0.000
Residual	119.929078	1,383	.086716615	R-squared	=	0.030
				Adj R-squared	=	0.027
Total	123.696129	1,387	.089182501	Root MSE	=	.2944
packs	Coef.	Std. Err.	t	P> t [95% Con	1	intervat
	0047261	.0158539	-0.30	0.7660358264		.026374
male			2.04	0.041 .0007291		.035569
male parity	.0181491	.0088802	2.04	0.041 .000/231		
	.0181491 0526374	.0088802 .0086991		0.0000697023	-	035572
parity			-6.05			.035572 .002299

图 3

可以看到,*cigprice* 对 *packs* 的影响不显著。而且,符号也是正的,和我们预想的负号相反。因此,在这里 *cigprice* 是很弱的工具变量。同时,*cigprice* 也可能不是外生的。这造成了 2SLS 回归的不一致以及方差很大。

总结:这道题说明,虽然工具变量从理论上来说很好,但实际操作中,工具变量的选择是比较困难的。做实证研究涉及到工具变量选择时,需要谨慎讨论其合理性。