

Weak instrumental variables due to nonlinearities in panel data: A Super Learner Control Function estimator

Monika Avila Marquez

Methods and Data Analysis, University of Geneva, Switzerland

Motivation

- ▶ A model with **endogeneity** and fixed effects for $i \in [N], t \in [T_i]$:

$$y_{it} = x_{1it}\beta_{1o} + \alpha_{y,i} + \varepsilon_{it}, \quad (1)$$

$$x_{1it} = g_o(z_{it}) + \alpha_{x,i} + u_{it}. \quad (2)$$

Parameter of interest: β_{1o}

Source of Identification: Exogenous variation z_{it} .

But...

- ▶ Standard: W2SLS **ignoring the nonlinearity**
- ▶ Ignoring nonlinearity may cause a problem of **weak instrumental variables**

Setup

A.1 $O_i = (y_i, x_{1i}, \tilde{X}_i, z_i)$ are $T_i \times 1$ independent. $T_i < T_{max}$, $N_T = \sum_i T_i$.

A.2 Linear structural equation:

$$y_{it} = x_{1it}\beta_{1o} + \tilde{x}'_{it}\beta_{2o} + \alpha_{i,y} + \varepsilon_{it}, \quad (3)$$

$$\mathbb{E}[\varepsilon_{it} | \tilde{X}_i, z_i, \alpha_{i,y}] = \mathbb{E}[\varepsilon_{it}] = 0.$$

A.3 Nonlinear reduced form equation:

$$x_{1it} = g_o(\tilde{x}_{it}, z_{it}) + \alpha_{i,1x} + u_{it}, \quad (4)$$

$$\mathbb{E}[u_{it} | \tilde{X}_i, z_i, \alpha_{i,1x}] = 0.$$

A.4 Linear relationship between the structural error and the reduced form error:

$$\varepsilon_{it} = \rho u_{it} + \omega_{it}, \quad (5)$$

$$\mathbb{E}[\omega_{it} | u_i] = 0.$$

Identification: Control Function approach

- ▶ Identification of β_{1o}
- 1. Transform ($\tau : \mathbb{R} \rightarrow \mathbb{R}$) model:

$$\tau y_{it} = \tau x_{1it}\beta_{1o} + \sum_{k=2}^K \tau x_{kit}\beta_{2k_o} + \tau \varepsilon_{it}, \quad (6)$$

$$\tau x_{1it} = \tau g_o(\tilde{x}_{it}, z_{it}) + \tau u_{it}, \quad i \in [N], t \in \{t_a, \dots, T_i\}. \quad (7)$$

2. Augment transformed structural equation with τu_{it} :

$$\tau y_{it} = \tau x_{1it}\beta_{1o} + \sum_{k=2}^K \tau x_{kit}\beta_{2k_o} + \rho_o \tau u_{it} + \tau \omega_{it}. \quad (8)$$

Population moment conditions:

$$\mathbb{E}[\phi(O_i; \theta_o, \tau g_o)] = \mathbb{E}[(M_{\tau_i} H_i)' V_i^{-1} M_{\tau_i} \omega_i] = \mathbf{0}. \quad (9)$$

- ▶ τu_{it} is identified (Proposition 1)

Super Learner Control Function Estimation

Step 1:

1. Partition the set $\{1, 2, \dots, N\}$ in B subsets S_1, S_2, \dots, S_B with $n_{T,b} = \sum_{i \in S_b} T_i$.
2. Estimate $\tau g_o(\tilde{x}_{it}, z_{it}) = \mathbb{E}[\tau x_{1it} | I_t]$ using a Super Learner with partition $S_b^c = \{i \in S_j, j \neq b\}$, call the estimation $\widehat{\tau g_o}^{S_b^c}$.
3. Obtain the residuals $\widehat{\tau u}_{it}^{S_b} = \tau x_{1it} - \widehat{\tau g_o}^{S_b^c}(\tilde{x}_{it}, z_{it})$ for partition $S_b = \{i \in S_b\}$.
4. The estimator of θ_o for partition S_b is the solution of the sample moment conditions.

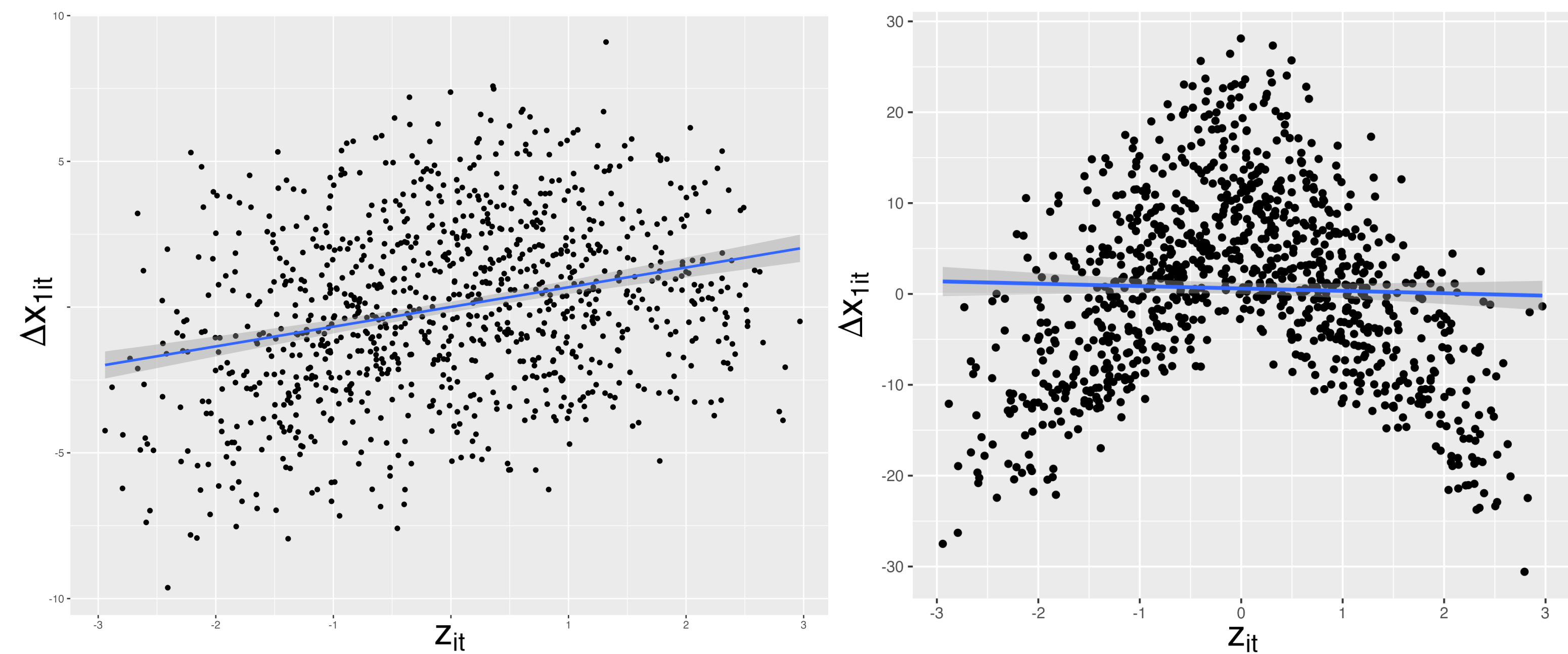
Step 2:

Average $\hat{\theta}_b$ to obtain the final estimator of θ_o as:

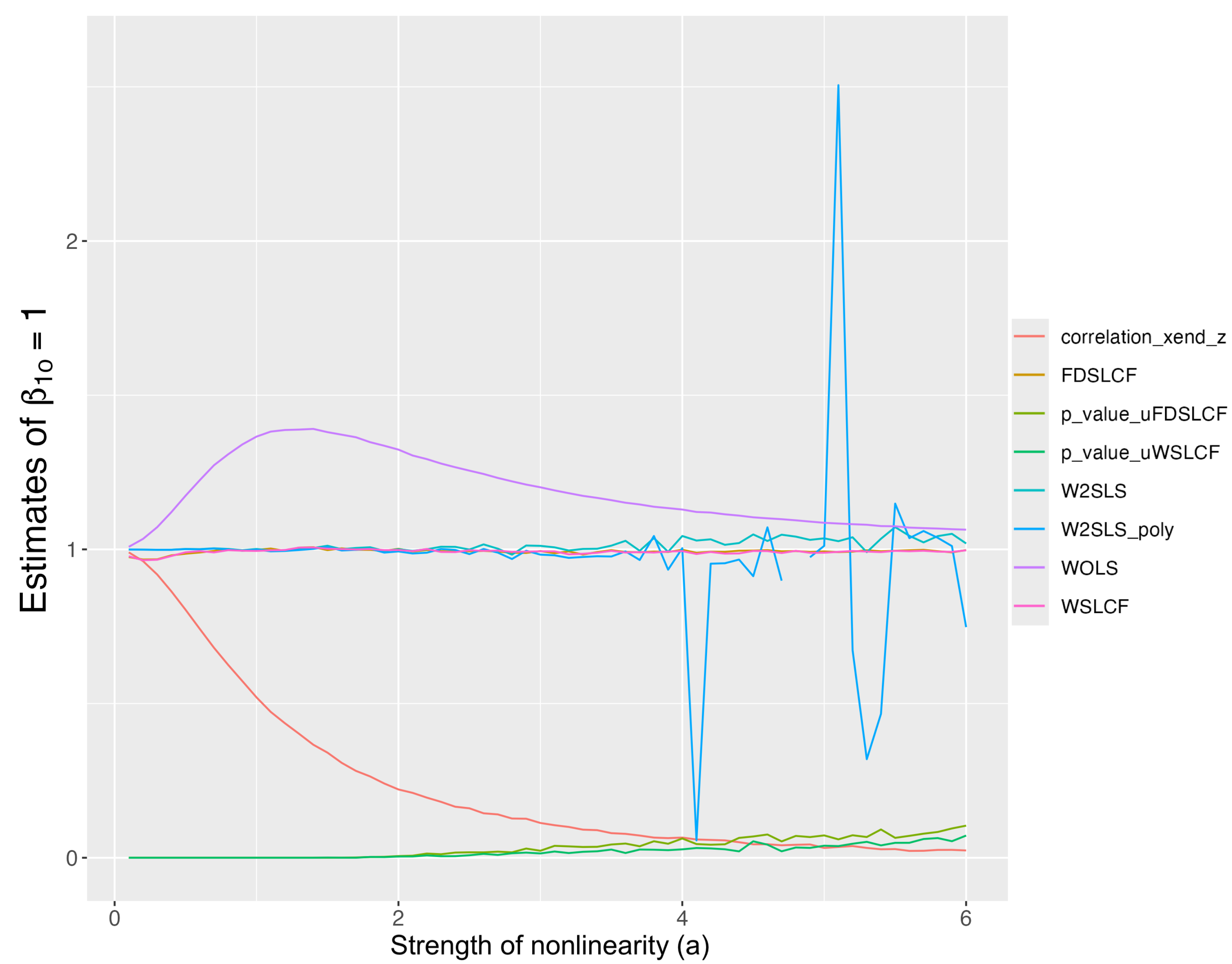
$$\hat{\theta}_o = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b.$$

Results

- ▶ Score $\phi(O_i; \theta_o, \tau g_o)$ is an Orthogonal score (Proposition 2).
- ▶ **Large sample properties**
SLCFE is consistent and asymptotically normal with parametric convergence rate $\sqrt{N_T}$ (Theorems 1 and 2).
- ▶ **Monte Carlo simulation**
The design: Relationship between the EEV and the instrument



The results: Average estimated β_{1o}



Average estimates of β_{1o} for 100 samples simulated with different values of a , $N = 1000, T = 2$

Empirical Application

Causal effect of air pollution on educational outcomes:

Table 1: Estimated effect of air pollution (PM 2.5 concentration) on student performance

	WOLS	W2SLS	FDSLFCF	WSLCF
Estimate	-0.072	-0.844	-0.979	-0.859
Standard Error	(0.011)	(0.762)	(0.101)	(0.085)

Note: US counties, N = 787, 2009-2013

Conclusion

- ▶ τu_{it} is identified (Proposition 1)
- ▶ Score is an orthogonal score (Proposition 2)
- ▶ SLCFE is $\sqrt{N_T}$ -consistent (Theorems 1 and 2)

Scan for paper & references

An R package is available on request

monika.avila@unige.ch

monikaavila.com

