

Practice 2SLS with Artificial Data Part 1

Yona Rubinstein

July 2016

Practice with Artificial Data

- In this note we use artificial data to illustrate how to approach a selection problem using 2SLS.
- The "constructed" data set "MG4A4 2SLS DIET EXAMPLE PART 1.dta" can be found on my website under "Practice 2SLS using Artificial Data".
- The data contains 1000 individuals each individual observed over 100 periods (days, weeks).
- We have information regarding their weight and on whether they are dating.
- Specifically the data contains their (1) permanent weight, (2) change in weight if diet is not taken, (3) whether they are on diet and (iv) whether they received an invitation to date.
- People are on diet for three reasons: (1) if their weight exceeds 220 pounds; (2) if they gained 5 pounds; (3) if they receive an invitation to date.
- The latter is exogenous to fluctuations in their current weight.

The Model: The Production Function of Peorsons' Weight

- The casual model exhibits the following form:

$$Y_{it} = \beta_0 + \beta_D D_{it} + U_{it}. \quad (1)$$

- The variable Y_{it} is person's i weight (in pounds) in time t and D_{it} is a binary indicator which equals 1 if person i is on diet.
- The error term (U_{it}) is a composition of person's i permanent weight (relative to the population mean, that is θ_i) and person's i specific time varying fluctuations to his weight:

$$U_{it} = \theta_i + \varepsilon_{it}. \quad (2)$$

- The parameters that I used to impute weights (Y_{it}) are:

$$Y_{it} = 0 - 10 \cdot D_{it} + U_{it}. \quad (3)$$

Selection into Diet - the Simplest Case

- People are on diet for three reasons:
 - ① If their weight, without diet, exceeds 220 pounds (100kg):
 $\beta_0 + \theta_i + \varepsilon_{it} \geq 220$.
 - ② If they gained 5 pounds (or more): $\varepsilon_{it} \geq 5$
 - ③ If they receive an invitation to date: $Date_{it} = 1$
- We observe neither $vdate_{it}$ nor ε_{it} . Yet we observe whether the person is on diet and his weight. We know that person i is on diet if:

$$D_{it} = \max [(\beta_0 + \theta_i + \varepsilon_{it} - 220), Date_{it}, \varepsilon_{it}] > 0. \quad (4)$$

- Describe the data set

variable name	storage type	display format	value label	variable label
id	byte	%8.0g		person id number
time	byte	%8.0g		
PWi	float	%9.0g		Person's permanent weight
Eit	float	%9.0g		epsilon it
date	byte	%8.0g		shock to date value
vdate	float	%9.0g		date - PWi/190
Date	float	%9.0g		1 if date==1
Diet	float	%9.0g		1 if on diet
PWit	float	%9.0g		PWi + Eit - 10*Diet

- Summary statistics

```
> sum ;
```

Variable	Obs	Mean	Std. Dev.	Min	Ma
id	10,000	50.5	28.86751	1	10
time	10,000	50.5	28.86751	1	10
PWi	10,000	189.5	28.86751	140	23
Eit	10,000	-1.49e-07	13.81276	-60.54	63.4
date	10,000	.5004	.5000248	0	
vdate	10,000	-.4969684	.5217212	-1.257895	.263157
Date	10,000	.5004	.5000248	0	
Diet	10,000	.7144	.4517223	0	
PWit	10,000	182.356	30.97509	84.08	268.0

Estimating the Regression Model

- We next estimate the model in equation (1) using OLS

$$Y_{it} = b_0 + b_D D_{it} + e_{it}. \quad (5)$$

```
> eststo: reg PwIt Diet ;
```

Source	SS	df	MS	Number of obs	=	10,000
Model	240026.875	1	240026.875	F(1, 9998)	=	256.56
Residual	9353574.21	9,998	935.54453	Prob > F	=	0.0000
				R-squared	=	0.0250
				Adj R-squared	=	0.0249
Total	9593601.08	9,999	959.456054	Root MSE	=	30.587

PwIt	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
Diet	10.84626	.6771461	16.02	0.000	9.51892 12.17361
_cons	174.6074	.5723387	305.08	0.000	173.4855 175.7293

- According to the OLS estimate b_D^{OLS} diet leads to a gain of 10 pounds in weight!!

Estimating the Regression Model Controlling for Fixed Effects

- Next next turn to estimate the model using controlling for person fixed effects (θ_i):

$$Y_{it} = b_0 + b_D D_{it} + \theta_i + n_{it}. \quad (6)$$

```
eststo: areg PwIt Diet, absorb(id) ;
Linear regression, absorbing indicators
```

Number of obs	=	10,000
F(1, 9899)	=	106.28
Prob > F	=	0.0000
R-squared	=	0.8353
Adj R-squared	=	0.8336
Root MSE	=	12.6344

```
-----+-----
```

PwIt	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
Diet	2.946265	.285791	10.31	0.000	2.386056	3.506473
_cons	180.2512	.2400997	750.73	0.000	179.7805	180.7218

```
-----+-----
```

id	F(99, 9899) =	491.888	0.000	(100 categories)
----	---------------	---------	-------	------------------

- Accounting for person fixed effects matter! the bias is smaller, yet still large enough (3 rather than -10).

Estimating the Regression Model using 2SLS

- Next we turn to estimate the model using 2SLS using the following equations:

- 1 The first stage model:

$$D_{it} = a_0 + a_D \text{Date}_{it} + v_{it}, \quad (7)$$

where V_{it} is the error term.

- 2 The second stage model:

$$Y_{it} = b_0 + b_D \hat{D}_{it} + e_{it}, \quad (8)$$

where $\hat{D}_{it} = a_0^{OLS} + a_D^{OLS} \text{Date}_{it}$.

- We present the results in next slide.

First Stage

$$D_{it} = a_0 + a_D \text{Date}_{it} + v_{it}, \quad (9)$$

```
> reg Diet Date ;
```

Source	SS	df	MS	Number of obs	=	10,000
Model	816.979723	1	816.979723	F(1, 9998)	=	6676.90
Residual	1223.34668	9,998	.12235914	Prob > F	=	0.0000
				R-squared	=	0.4004
				Adj R-squared	=	0.4004
Total	2040.3264	9,999	.204053045	Root MSE	=	.3498

Diet	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
Date	.5716573	.006996	81.71	0.000	.5579438	.5853708
_cons	.4283427	.0049489	86.55	0.000	.4186419	.4380435

Second Stage (without correcting SE).

$$Y_{it} = b_0 + b_D \hat{D}_{it} + e_{it}, \quad (10)$$

```
> eststo: reg PWit Diethat ;
```

Source	SS	df	MS	Number of obs	=	10,000
Model	65389.6299	1	65389.6299	F(1, 9998)	=	68.61
Residual	9528211.45	9,998	953.011748	Prob > F	=	0.0000
				R-squared	=	0.0068
				Adj R-squared	=	0.0067
Total	9593601.08	9,999	959.456054	Root MSE	=	30.871

PWit	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
Diethat	-8.94641	1.080049	-8.28	0.000	-11.06352	-6.829296
_cons	188.7473	.8310522	227.12	0.000	187.1183	190.3763

Estimating the 2SLS using "ivreg"

- Note that we obtain identical point estimates. The standard errors were corrected to account for using a projected variable \hat{D}_{it} .
- Using dates as an instrument allows to correct of selection on "LHS" variable.

```
> eststo: ivreg PWit (Diet=Date) ;
Instrumental variables (2SLS) regression
```

Source	SS	df	MS	Number of obs	=	10,000
Model	-559270.763	1	-559270.763	F(1, 9998)	=	64.39
Residual	10152871.8	9,998	1015.49028	Prob > F	=	0.0000
				R-squared	=	.
				Adj R-squared	=	.
Total	9593601.08	9,999	959.456054	Root MSE	=	31.867

PWit	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
Diet	-8.94641	1.114891	-8.02	0.000	-11.13182 -6.761
_cons	188.7473	.8578613	220.02	0.000	187.0657 190.4289

Controlling for Fixed Effects

- We can take advantage of the panel data and further control for omitted person fixed effects.

- 1 The first stage model controls for person fixed effects (γ_i):

$$D_{it} = a_0 + a_D \text{Date}_{it} + \gamma_i + v_{it}, \quad (11)$$

where V_{it} is the error term.

- 2 The second stage model also controls for person fixed effects (δ_i):

$$Y_{it} = b_0 + b_D \hat{D}_{it} + \delta_i + e_{it}, \quad (12)$$

where $\hat{D}_{it} = a_0^{OLS} + a_D^{OLS} \text{Date}_{it} + \gamma_i^{OLS}$

- We present the results in next slide.

Estimating the 2SLS controlling for person fixed effects using "xtivreg"

- Person fixed effects are part of the causal model. While our instrument provides a source of exogenous variation - controlling for persons' fixed effects does not hurt.

```
> eststo: xtivreg Pwit (Diet=Date) ;
G2SLS random-effects IV regression
Group variable: id
R-sq:
    within = 0.0106
    between = 0.3901
    overall = 0.0250
corr(u_i, X)      = 0 (assumed)
Number of obs     = 10,000
Number of groups  = 100
Obs per group:
    min = 100
    avg = 100.0
    max = 100
Wald chi2(1)     = 385.69
Prob > chi2      = 0.0000
```

Pwit	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
Diet	-9.548751	.486211	-19.64	0.000	-10.50171	-8.595795
_cons	189.1776	2.624326	72.09	0.000	184.034	194.3212

Estimating the Model using OLS, FE and 2SLS

	OLS (1)	FE (2)	OLS2ND (3)	2SLS (4)	2SLS w/FE (5)
Dependent variable: "Weight"					
Diet	10.846*** (0.677)	2.946*** (0.286)	-8.946*** (1.080)	-8.946*** (1.115)	-9.549*** (0.486)
Constant	175*** (0.572)	180*** (0.240)	189*** (0.831)	189*** (0.858)	189*** (2.624)

First Stage: Dependent variable "Diet"

Date			0.572*** (0.007)	0.572*** (0.007)	0.571*** (0.007)
			0.428*** (0.005)	0.428*** (0.005)	0.429*** (0.005)
R-square			0.400	0.400	0.390
Observation:	10000	10000	10000	10000	10000

Take Home Message

- We can identify the causal impact of a treatment on an outcome of interest – **accounting for selection into treatment** – if we have a variable that
 - 1 Is uncorrelated with the **unobserved component** in the outcome equation (the error term);
 - 2 Does **not** affect directly the outcome of interest
- Controlling for subjects fixed effects might eliminate some of the bias but not all as long as subject self-sort into treatment on **time varying unobservables**.