Seminar 9 Solutions

Giulio Rossetti* giuliorossetti94.github.io March 13, 2025

* email: giulio.rossetti.1@wbs.ac.uk

Disclaimer

Full solutions are available on my.wbs. All exercises are examinable material, not just the ones we covered in the seminars.

Roadmap

Exercise 1

Exercise 4

Part 1: Theory

Exercise 1

For T = 2, consider the standard panel data model:

$$y_{it} = x'_{it}\beta + \alpha_i + u_{it}, \quad t = 1, 2, \quad i = 1, \dots, n$$

where i denotes the cross-sectional unit and t denotes the time dimension. For simplicity, assume that in this model there is no intercept.

First-Difference Estimator

Show that the fixed-effects(FE) and first-difference (FD) estimators are identical (they deliver the same beta estimates.)

• FD: Remove unobs heterogeneity by differencing over time:

$$y_{i2} - y_{i1} = (x_{i2} - x_{i1})'\beta + (u_{i2} - u_{i1})$$

 $\Delta y_i = \Delta x'_i \beta + \Delta u_i.$

• Assuming independence of the error terms, β_{FD} :

$$\hat{\beta}_{FD} = \left(\sum_{i=1}^{n} \Delta x_i \Delta x'_i\right)^{-1} \sum_{i=1}^{n} \Delta x_i \Delta y_i.$$

Fixed-Effects Estimator

• FE : Remove unobs heterogeneity by demeaning:

$$\bar{y}_i = \frac{1}{2}(y_{i1} + y_{i2}), \quad \bar{x}_i = \frac{1}{2}(x_{i1} + x_{i2}), \quad \bar{u}_i = \frac{1}{2}(u_{i1} + u_{i2}).$$

• Then, we have:

$$y_{it} - \bar{y}_i = (x_{it} - \bar{x}_i)'\beta + u_{it} - \bar{u}_i, \quad t = 1, 2.$$

• β_{FE} :

$$\hat{\beta}_{FE} = \left(\sum_{i=1}^{n} \sum_{t=1}^{2} (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)'\right)^{-1} \sum_{i=1}^{n} \sum_{t=1}^{2} (x_{it} - \bar{x}_i)(y_{it} - \bar{y}_i).$$

Equivalence of FE and FD

Note that:

$$\sum_{t=1}^{2} (x_{it} - \bar{x}_i)(x_{it} - \bar{x}_i)' = \sum_{t=1}^{2} \left(x_{it} - \frac{x_{i1} + x_{i2}}{2} \right) \left(x_{it} - \frac{x_{i1} + x_{i2}}{2} \right)'$$
$$= \left(\frac{x_{i1} - x_{i2}}{2} \right) \left(\frac{x_{i1} - x_{i2}}{2} \right)' + \left(\frac{x_{i2} - x_{i1}}{2} \right) \left(\frac{x_{i2} - x_{i1}}{2} \right)'$$
$$= \frac{1}{2} \Delta x_i \Delta x'_i.$$

Equivalence of FE and FD

Similarly:

$$\sum_{t=1}^{2} (x_{it} - \bar{x}_i)(y_{it} - \bar{y}_i) = \frac{1}{2} \Delta x_i \Delta y_i.$$

Substituting into the FE estimator, we obtain:

$$\hat{\beta}_{FE} = \left(\frac{1}{2}\sum_{i=1}^{n} \Delta x_i \Delta x'_i\right)^{-1} \left(\frac{1}{2}\sum_{i=1}^{n} \Delta x_i \Delta y_i\right)$$
$$= \left(\sum_{i=1}^{n} \Delta x_i \Delta x'_i\right)^{-1} \sum_{i=1}^{n} \Delta x_i \Delta y_i = \hat{\beta}_{FD}.$$

.

Conclusion: The fixed-effects and first-difference estimators are identical when T = 2.

Including age as a Regressor

Suppose that we include the variable *age* as an additional regressor and use first differencing to estimate a fixed effects model.

- Requirements behind the FD estimator: Δx_{it} must have some variation across *i*.
- This fails if an explanatory variable such as age is included.
 - age changes by the same amount for each of the individuals over time

$$y_{i1} = \beta_1 x_{i1} + \beta_2 x_{i2} + \alpha_i + u_{i1}, \quad t = 1, \quad i = 1, \dots, n$$

 $y_{i2} = \beta_1 x_{i2} + \beta_2 x_{i2} + \alpha_i + u_{i2}, \quad t = 2, \quad i = 1, \dots, n.$

Differencing the Model

By subtracting the first equation from the second, we obtain:

$$\Delta y_i = \beta_1 \Delta x_{i1} + \beta_2 \Delta x_{i2} + \Delta u_i, \quad i = 1, \dots, n.$$

Since x_{i2} increases by the same amount *c* across individuals:

$$\Delta y_i = \beta_1 \Delta x_{i1} + \beta_2 c + \Delta u_i.$$

$$=\beta_1 \Delta x_{i1} + \delta + \Delta u_i.$$

where $\delta = \beta_2 c$ is a constant term.

Key issue: The constant term δ makes it problematic to identify β_2 .

- δ does not represent the intercept (since there was no intercept in the original model).
- It also does not represent any change in the intercept by definition:
 - Since we allow α_i to be correlated with x_{i2} , we cannot separate the effect of α_i on y_i from the effect of any other variable that does not change over time.

Implications of $Cov(x_{it}, \alpha_i) = 0$

Suppose that $Cov(x_{it}, \alpha_i) = 0$. What does this imply for the FE and FD estimators?

- When we assume that $Cov(x_{it}, \alpha_i) = 0$, the original model becomes a random effects model.
- The random effects assumptions include all of the fixed effects assumptions plus the additional requirement that α_i is independent of all explanatory variables in all time periods.
- WNote that given $\text{Cov}(x_{it}, \alpha_i) = 0$, β can be consistently estimated by Pooled OLS.

Composite Error Term

However, this ignores a key feature of the model. If we define the composite error term as:

 $v_{it} = \alpha_i + u_{it},$

we can show that:

$$\operatorname{corr}(v_{it}, v_{is}) = \frac{\sigma_{\alpha}^2}{\sigma_{\alpha}^2 + \sigma_u^2}, \quad t \neq s,$$

where:

$$\sigma_{\alpha}^2 = \operatorname{Var}(\alpha_i), \quad \sigma_u^2 = \operatorname{Var}(u_{it}).$$

Implications for Estimation

- positive serial correlation in the error term makes pooled OLS standard errors incorrect.
- We must:
 - Either correct the OLS SE, or
 - Use the GLS random effects estimator

Roadmap

Exercise 1

Exercise 4

Exercise 4: Rental Prices and Student Presence

The data for the years 1980 and 1990 include rental prices and other variables for college towns. The goal is to determine whether a stronger presence of students affects rental rates. The model is:

 $\log(\operatorname{rent}_{it}) = \beta_0 + \delta_0 y 90_t + \beta_1 \log(\operatorname{pop}_{it}) + \beta_2 \log(\operatorname{avginc}_{it}) + \beta_3 \operatorname{pctstu}_{it} + e_{it},$

where:

- pop is city population,
- avginc is average income,
- pctstu is student population as a percentage of city population (during the school year).

Pooled OLS Estimation Results

You estimate the model with pooled OLS and obtain the following results:

Source	33	df	MS		Number of obs	=	128
Model Residual Total	12.1080112 1.9501234 14.0581346	4 123 127	3.02700281 .015854662 .110693974		F(4, 123) Prob > F R-squared Adj R-squared Root M3E	= = =	190.92 0.0000 0.8613 0.8568 .12592
lrent	Coef.	Std. E	rr. t	P> t	[95% Conf.	Int	erval]
y 90	.2622267	.03476	32 7.54	0.000	.1934151	. 3	310384
lpop	.0406863	.02251	54 1.81	0.073	0038815	. 0	852541
lavginc	.5714461	.05309	81 10.76	0.000	.4663417	. 6	5765504
petstu	.0050436	.00101	92 4.95	0.000	.0030262		007061
_cons	5688069	.53488	08 -1.06	0.290	-1.627571	.4	899568

Figure: Pooled OLS Estimation Results for Rental Prices and Student Presence

Interpreting the Regression Results

- Almost all regressors are statistically significant.
- City population is borderline significant.
- However, population per se is not a strong driving factor:
 - The number of inhabitants affects rents only if land size is limited.
 - This constraint is not explicitly considered in the model.
- There is a clear omitted variable bias:
 - City size is not constant and may depend on the city itself.
 - Example: London and Coventry do not have the same size.
- This leads to the so-called heterogeneous bias.
- To address this issue:
 - A fixed effects model can be used if regressors are correlated with city-specific effects.
 - A random effects model can be used if regressors are uncorrelated with city-specific effects.

Pooled OLS Estimation Results

Now you estimate the model with fixed effect and obtain the following results:

corr(u_i, Xb)	0.1297			F(4,60) Frob > I	-	624.15 0.0000
lrenb	Coef.	Std. Err.	6	Po o	[95% Conf.	Interval]
290	.3055219	.0365245	10.47	0.000	.3110615	. 4591813
lpop	.0722456	.0882426	0.82	0.417	104466	.2489571
lauging	.2099605	.0664771	4.66	0.000	.1769865	.4429246
petota	.0112033	.0041319	2.71	0.009	.0029382	.0194684
_cons	1.409384	1.167298	1.21	0.292	9254994	9.744208
pigma_u	.15905877					
sigma_e	.06372673					
rho	.0616755	(fraction (of varia	nce due to	u_i)	
I test that al	11 u_1=0:	r(63, 60) =	10.2	0	Frob > 1	r = 0.0000

Figure: FE Estimation Results for Rental Prices and Student Presence

Fixed Effects and Model Selection

- By fully acknowledging unobservable fixed effects, the impact of *lpop* disappears.
- From the output, we see that:

 $\operatorname{corr}(\alpha_i, x_{it}) = -0.129,$

which is relatively small.

- Given this small correlation, it might be sensible to use a *random effects model* instead.
- However, determining the appropriate model is difficult without first implementing a **Hausman test**.