

# Applied Data Analysis School

## Panel Data Models

Anabela Carneiro<sup>1</sup>   Miguel Portela<sup>2</sup>

<sup>1</sup>Universidade do Porto

<sup>2</sup>Universidade do Minho

November 17 & 19, 2020

- Outline

- ➊ Introduction to Panel Data Analysis
- ➋ Pooled OLS Model
- ➌ Random Effects Model
- ➍ Fixed Effects Model
- ➎ Test for the Presence of Fixed Effects
- ➏ Test for Random Effects: the Hausman Test
- ➐ Comparison of Estimators

# Introduction

- Econometric models aim to establish causal relationships.
- Endogeneity may confound the evaluation of the effects of interest.
- Laboratory or randomized field experiments are rare in economics and other social sciences research.
- Hence, researchers depend on observational data to make causal claims.
- One of the most important group of models that mitigates endogeneity bias is known as fixed effects models (FE).
- Fixed effects models allow you to deal with endogeneity due to omitted variable.

- A FE approach relies on panel data.
- The term panel data (or longitudinal data) is used for a wide variety of situations in econometrics.
- It refers to any data set with repeated observations over time for the same individuals (workers, households, firms, industries, regions, or countries).
- *Panel data* typically refers to situations in which the cross-sectional dimension ( $N$ ) is large relative to the time dimension ( $T$ ).

## Example

Employer-employee level data: wages, schooling, experience, location, industry, ...

## Example

Data on countries: GDP, average education, physical capital, ....

## Example

Hospital-doctor-patient data

## Example

School-class-teacher-student level data

# Panel Data - *Personnel Records* (Quadros de Pessoal): Matched Employer-employee Data Set

## Personnel Records - Firms

	FIRMID	LOCATION	INDUSTRY	NUMBER OF EMPLOYEES	SALES	...
2001	10000	1301	52111	10	€	
2002	10000	1301	52111	10	€	
2003	10000	1301	52111	12	€	
2001	20000	1602	17400	25	€	
2002	20000	1602	17400	20	€	
2003	20000	1602	17400	18	€	
2002	160000	1105	65121	2	€	
2003	160000	1105	65121	4	€	
	...	...	...	...	...	...

# Panel Data - *Personnel Records* (Quadros de Pessoal): Matched Employer-employee Data Set

## Personnel Records - Workers

	WORKERID	FIRMID	GENDER	AGE	EDUC	WAGE	...
2001	999999990	10000	1	21	9	€	
2002	999999990	10000	1	22	9	€	
2003	999999990	10000	1	23	9	€	
2001	999999991	10000	2	40	4	€	
2002	999999991	10000	2	41	6	€	
2001	999999992	20000	1	18	12	€	
2002	999999992	20000	1	19	12	€	
2003	999999992	20000	1	20	12	€	
2002	999999993	160000	1	50	16	€	
2003	999999993	160000	1	51	16	€	
2002	999999994	160000	2	48	9	€	
	...	...	...	...	...	...	...

# Advantages of Panel Data (Hsiao, 2003)

- They usually give the researcher a large number of data points, increasing the degrees of freedom and reducing the collinearity among explanatory variables - improved efficiency in estimators and gains in terms of identification.
- Longitudinal data allow a researcher to analyze a number of important economic questions that cannot be addressed using cross-sectional or time-series data sets alone: e. g., firm survival, worker mobility, patient survival, etc.
- Control for unobserved time-invariant variables potentially correlated with the error term.



$$y_{it} = \alpha_i + \mathbf{x}_{it}\boldsymbol{\beta} + u_{it}, \quad t = 1, \dots, T; i = 1, \dots, N \quad (1)$$

- where  $\mathbf{x}_{it}$  is a  $(1 \times K)$  vector of explanatory variables and  $\boldsymbol{\beta}$   $(K \times 1)$  the vector of unknown parameters.
- $\alpha_i$ : unobserved heterogeneity, individual effect, fixed effect.
- $u_{it}$ : is the random error term assumed to be independent and identically distributed  $\text{IID}(0, \sigma_u^2)$ .

- **Key issue:** is  $\alpha_i$  correlated with the observed explanatory variables  $\mathbf{x}_{it}$ ,  $t = 1, \dots, T$  ?
- If  $\alpha_i$  is referred as an “individual random effect”,  $\alpha_i$  is being assumed to be uncorrelated with the  $\mathbf{x}_{it}$ , i. e.,  $\text{Cov}(\mathbf{x}_{it}, \alpha_i) = \mathbf{0}$ ,  $t = 1, \dots, T$ .
- If  $\alpha_i$  is referred as a “fixed effect”,  $\alpha_i$  is allowed to be correlated with the  $\mathbf{x}_{it}$ , i. e.,  $\text{Cov}(\mathbf{x}_{it}, \alpha_i) \neq \mathbf{0}$ ,  $t = 1, \dots, T$ .

# Pooled OLS (POLS)

- The pooled model treats all observations as independent observations.
- The **pooled OLS (POLS) estimator** is consistent if the composite error term is uncorrelated with the regressors, i. e.,

$$\begin{aligned} \text{cov}(\mathbf{x}_{it}, u_{it}) &= \mathbf{0} \\ \text{cov}(\mathbf{x}_{it}, \alpha_i) &= \mathbf{0}, \quad t = 1, \dots, T \end{aligned} \tag{2}$$

- Even if assumption (2) holds, the composite errors will be serially correlated due to the presence of  $\alpha_i$  in each time period ( $\text{Cov}(v_{it}, v_{is}) = \sigma_\alpha^2$ ).
- Therefore, inference using pooled OLS requires the robust variance matrix estimator and robust test statistics.

# Random Effects Model (RE): Assumptions

- The RE model accommodates this correlation  $\implies$  apply GLS methods.
- Assumption RE.1:

$$y_{it} = \alpha_i + \mathbf{x}_{it}\boldsymbol{\beta} + u_{it}$$

- Assumption RE.2: We have a random sample in the cross-sectional dimension.
- Assumption RE.3: There are no perfect linear relationships among the explanatory variables.

# Random Effects Model (RE): Assumptions

- Assumption RE.4:

$$E(u_{it} | \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT}, \alpha_i) = 0, \quad t = 1, \dots, T$$

$$E(\alpha_i | \mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT}) = 0, \quad t = 1, \dots, T$$

- Assumption RE.5:

$$\text{var}(u_{it} | \mathbf{x}_i, \alpha_i) = \text{var}(u_{it}) = \sigma_u^2, \quad \text{for all } t = 1, \dots, T$$

$$\text{var}(\alpha_i | \mathbf{x}_i) = \sigma_\alpha^2, \quad t = 1, \dots, T$$

where  $\mathbf{x}_i \equiv (\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT})$ .

- Assumption RE.6:

$$\text{cov}(u_{it}, u_{is} | \mathbf{x}_i, \alpha_i) = 0, \quad \text{for all } t \neq s$$

# Random Effects Estimator

- Under these assumptions the RE estimator is consistent as  $N$  gets large for fixed  $T$ .
- The FGLS estimator of the RE model is the **random effects (RE) estimator**:

$$\hat{\beta}_{RE} = \left( \sum_{i=1}^N \mathbf{x}_i' \hat{\Omega}^{-1} \mathbf{x}_i \right)^{-1} \left( \sum_{i=1}^N \mathbf{x}_i' \hat{\Omega}^{-1} \mathbf{y}_i \right) \quad (3)$$

- The RE estimator is fully efficient under the RE assumptions, though the efficiency gain compared to POLS need not to be great.

- Where  $\Omega$  takes the special form

$$\Omega \equiv E(\mathbf{v}_i \mathbf{v}_i') =$$
$$= \begin{pmatrix} \sigma_\alpha^2 + \sigma_u^2 & \sigma_\alpha^2 & \cdots & \sigma_\alpha^2 \\ \sigma_\alpha^2 & \sigma_\alpha^2 + \sigma_u^2 & \cdots & \vdots \\ \vdots & \vdots & \ddots & \sigma_\alpha^2 \\ \sigma_\alpha^2 & \cdots & \cdots & \sigma_\alpha^2 + \sigma_u^2 \end{pmatrix}$$

# Random Effects Estimator

- $\hat{\sigma}_\alpha^2$  ?  $\hat{\sigma}_u^2$  ?
- These estimators can be based on the pooled OLS or fixed effects residuals.
- Let  $\hat{v}_{it}$  denote the pooled OLS residuals. A consistent estimator of  $\sigma_v^2$  under RE.1-RE.3 is:

$$\hat{\sigma}_v^2 = \frac{1}{(NT - K)} \sum_{i=1}^N \sum_{t=1}^T \hat{v}_{it}^2 \quad (4)$$

which is the usual variance estimator from the OLS regression on the pooled data.



- A consistent estimator of  $\sigma_\alpha^2$  under RE.1-RE.6 is:

$$\hat{\sigma}_\alpha^2 = \frac{1}{[NT(T-1)/2 - K]} \sum_{i=1}^N \sum_{t=1}^{T-1} \sum_{s=t+1}^T \hat{\hat{v}}_{it} \hat{\hat{v}}_{is} \quad (5)$$

$$\hat{\sigma}_u^2 = \hat{\sigma}_v^2 - \hat{\sigma}_\alpha^2 \quad (6)$$

# Testing for the Presence of an Unobserved Effect

- The absence of an unobserved effect is statistically equivalent to  $H_0: \sigma_\alpha^2 = 0$ .
- Notice that:

$$\text{Corr}(v_{it}, v_{is}) = \sigma_\alpha^2 / (\sigma_\alpha^2 + \sigma_u^2) \geq 0, \quad s \neq t.$$

- Thus, if  $\sigma_\alpha^2 = 0$  then  $\text{Corr}(v_{it}, v_{is}) = 0$ .
- Lagrange multiplier test for the random effects model based on the pooled OLS residuals ( $\hat{v}_{it}$ ) proposed by Breusch and Pagan (1980).  
For

$$H_0 : \sigma_\alpha^2 = 0$$

$$H_1 : \sigma_\alpha^2 \neq 0$$

- If  $\sigma_{\alpha}^2 = 0$  then  $\Omega$  takes the special form

$$\Omega \equiv E(\mathbf{v}_i \mathbf{v}_i') =$$
$$= \begin{pmatrix} \sigma_u^2 & 0 & \dots & 0 \\ 0 & \sigma_u^2 & \dots & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \dots & & \sigma_u^2 \end{pmatrix} = \sigma_u^2 \mathbf{I}_T$$

# Testing for the Presence of an Unobserved Effect

- The test statistic is

$$LM = \frac{NT}{2(T-1)} \left[ \frac{\sum_{i=1}^N \left[ \sum_{t=1}^T \hat{v}_{it} \right]^2}{\sum_{i=1}^N \sum_{t=1}^T \hat{v}_{it}^2} - 1 \right]^2$$

- Under the null hypothesis, LM is distributed as chi-squared with one degree of freedom. Large values of LM favor the RE model.

# Endogeneity Due to Omitted Variables

- Let  $v_{it} = \alpha_i + u_{it}$
- If  $\text{Cov}(\mathbf{x}_{it}, \alpha_i) \neq \mathbf{0}$ , then  $\text{Cov}(\mathbf{x}_{it}, v_{it}) \neq \mathbf{0} \Rightarrow$  endogeneity due to omitted variable.
- OLS and RE estimators are inconsistent.
- Use FE to control for individual unobserved heterogeneity.

# Fixed Effects (FE) Model: Assumptions

- Assumption FE.1:

$$y_{it} = \alpha_i + \mathbf{x}_{it}\boldsymbol{\beta} + u_{it}$$

- Assumption FE.2: We have a random sample in the cross-sectional dimension.
- Assumption FE.3: There are no perfect linear relationships among the explanatory variables.

# FE Model: Assumptions

- Assumption FE.4 (strict exogeneity assumption):

$$E(u_{it}|\mathbf{x}_i, \alpha_i) = 0, \quad t = 1, \dots, T \quad (7)$$

where  $\mathbf{x}_i \equiv (\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT})$ .

- Assumption FE.5:

$$\text{var}(u_{it}|\mathbf{x}_i, \alpha_i) = \text{var}(u_{it}) = \sigma_u^2, \quad \text{for all } t = 1, \dots, T$$

- Assumption FE.6:

$$\text{cov}(u_{it}, u_{is}|\mathbf{x}_i, \alpha_i) = 0, \quad \text{for all } t \neq s$$

# FE Estimator

- $\alpha_i$  allowed to be arbitrarily correlated with  $\mathbf{x}_i \implies$  control for unobserved heterogeneity.
- The idea for estimating  $\beta$  under FE.1 is to transform the equations to eliminate the unobserved effect  $\alpha_i \Rightarrow$  **FE transformation/within transformation**.
- The FE transformation is obtained by first averaging equation (1) over  $t = 1, \dots, T$  to get the cross section equation

$$\bar{y}_i = \bar{\mathbf{x}}_i \beta + \alpha_i + \bar{u}_i \quad (8)$$

where  $\bar{y}_i = T^{-1} \sum_{t=1}^T y_{it}$ ,  $\bar{\mathbf{x}}_i = T^{-1} \sum_{t=1}^T \mathbf{x}_{it}$  and  $\bar{u}_i = T^{-1} \sum_{t=1}^T u_{it}$ .



- Subtracting equation (8) from equation (1) for each  $t$  gives FE transformed equation,

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\boldsymbol{\beta} + u_{it} - \bar{u}_i \quad (9)$$

- OLS estimator of equation (9) will be consistent, if  $\text{cov}[(\mathbf{x}_{it} - \bar{\mathbf{x}}_i), (u_{it} - \bar{u}_i)] = \mathbf{0}, t = 1, \dots, T$ .
- The **fixed effects (FE) estimator**, denoted by  $\hat{\beta}_{FE}$  is the OLS estimator from the regression

$$(y_{it} - \bar{y}_i) \text{ on } (\mathbf{x}_{it} - \bar{\mathbf{x}}_i) \quad (10)$$

# Identifying Assumption

- Identifying assumption: unobservable factors that might simultaneously affect the dependent variable and the independent variables of the regression are time-invariant.
- Identification of  $\beta$  is based on within group variation over time.
- Coefficients of variables that are time-invariant are not identified as  $(\mathbf{x}_{it} - \bar{\mathbf{x}}_i) = 0, \forall_i$ .

# Least Square Dummy Variable (LSDV) Estimator

- If the  $\alpha_i$  are parameters to estimate, how would we estimate each  $\alpha_i$  along with  $\beta$ ?
- One possibility is to define  $N$  dummy variables, one for each cross section observation:  $dn_i = 1$  if  $n = i$ ,  $dn_i = 0$  if  $n \neq i$ . Then, run the OLS regression

$$y_{it} \text{ on } d1_i, d2_i, \dots, dN_i, \mathbf{x}_{it} \quad (11)$$

- Then,  $\hat{\alpha}_1$  is the coefficient on  $d1_i$ ,  $\hat{\alpha}_2$  is the coefficient on  $d2_i$ , and so on. The estimator of  $\beta$  obtained from regression (11) is, in fact, the FE estimator.
- This is why  $\hat{\beta}_{FE}$  is sometimes referred to as the **least square dummy variable estimator (LSDV)**.

# Testing for Fixed Effects

- Run the OLS regression

$$y_{it} \text{ on constant, } d1_i, d2_i, \dots, dN-1_i, \mathbf{x}_{it} \quad (12)$$

- Let  $SSR_{ur}$  be the unrestricted sum of squared residuals from regression (12).
- Then, run the OLS regression

$$y_{it} \text{ on constant, } \mathbf{x}_{it} \quad (13)$$

- Let  $SSR_r$  be the restricted sum of squared residuals from regression (13).

# Testing for Fixed Effects

- $H_0 : \alpha_1 = \alpha_2 = \dots = \alpha_{N-1} = 0$

- Then, under  $H_0$

$$F = \frac{(SSR_r - SSR_{ur})}{SSR_{ur}} \cdot \frac{[N(T-1) - K]}{N-1} \sim F_{N-1, N(T-1)-K}$$

- If  $H_0$  is not rejected, using FE does not cause bias.
- Though, in this case, controlling for FE when all  $\alpha_i = 0$  will lead to larger standard errors.

# Testing for Fixed Effects

- If  $H_0$  is rejected  $\Rightarrow$  use FE to control for them.
- Notice, however, that the pooled OLS model and the FE model will not necessarily produce different results.
- Bias occurs when  $Cov(\mathbf{x}_{it}, \alpha_i) \neq \mathbf{0}$ .
- To cause bias, FE must exist and be correlated with the independent variables.

# Estimating the Fixed Effects

- Sometimes it is useful to obtain the  $\hat{\alpha}_i$  even when regression (11) is unfeasible ( $N$  is too large). Using the OLS first-order conditions, each  $\hat{\alpha}_i$  can be shown to be

$$\hat{\alpha}_i = \bar{y}_i - \bar{\mathbf{x}}_i \hat{\boldsymbol{\beta}}_{FE}, \quad i = 1, \dots, N$$

- Econometric software that employs fixed effects usually suppresses the "estimates" of the  $\alpha_i$ , although an overall intercept is often reported.



# RE vs FE - Hausman (1978) Test

- Comparing the FE and RE estimates can be a test for whether there is correlation between  $\alpha_i$  and  $\mathbf{x}_{it}$ , assuming that the idiosyncratic errors and explanatory variables are uncorrelated across all time periods.
- Given the random effects specification:

$$y_{it} = \alpha_i + \mathbf{x}_{it}\boldsymbol{\beta} + u_{it}$$

- and assuming that the orthogonality condition holds for the  $u_{it}$ , the null hypothesis that we want to test is:

$$H_0 : \text{Cov}(\mathbf{x}_{it}, \alpha_i) = \mathbf{0}$$

while the alternative hypothesis is:

$$H_1 : \text{Cov}(\mathbf{x}_{it}, \alpha_i) \neq \mathbf{0}$$

# RE vs FE - Hausman (1978) Test

- Under  $H_0$ :
  - the fixed effects estimator is consistent but inefficient;
  - the random effects estimator is consistent and efficient;
- Under  $H_1$ :
  - the fixed effects estimator remains consistent;
  - the random effects estimator becomes inconsistent.
- Therefore, under the null hypothesis the two estimators should not differ and this observation provides the basis for the test.

# RE vs FE - Hausman (1978) Test

- Let  $\hat{\delta}_{RE}$  denote the vector of RE estimates for the coefficients on the regressors that change across both  $i$  and  $t$ , and let  $\hat{\delta}_{FE}$  denote the corresponding FE estimates; let these each be  $K \times 1$  vectors.
- The Hausman test statistic can be written as:

$$H = (\hat{\delta}_{FE} - \hat{\delta}_{RE})' [\hat{\mathbf{V}}(\hat{\delta}_{FE}) - \hat{\mathbf{V}}(\hat{\delta}_{RE})]^{-1} (\hat{\delta}_{FE} - \hat{\delta}_{RE}) \stackrel{H_0}{\sim} \chi_K^2$$

where the  $\hat{\mathbf{V}}$  denote the estimates of the true covariance matrices.

- If  $H_0$  is rejected, we have to conclude that the RE estimator is inconsistent because the orthogonality condition fails.

# Decomposing OLS Estimates

- Ignoring the individual effects, we have the standard model

$$y_{it} = \alpha + \mathbf{x}_{it}\boldsymbol{\beta} + u_{it} \quad (14)$$

- We can decompose the OLS estimator by considering two variations on this model.
- Deviations from means (within estimator):

$$y_{it} - \bar{y}_i = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\boldsymbol{\beta} + (u_{it} - \bar{u}_i) \quad (15)$$

- Group means only (between estimator):

$$\bar{y}_i = \alpha_i + \bar{\mathbf{x}}_i\boldsymbol{\beta} + \bar{u}_i$$

$$\bar{y}_i = \alpha + \bar{\mathbf{x}}_i\boldsymbol{\beta} + (\alpha_i - \alpha + \bar{u}_i) \quad (16)$$

# Decomposing OLS Estimates

- The **between estimator (BE)** is the OLS estimator from regression of  $\bar{y}_i$  on an intercept and  $\bar{x}_i$ .
- The between estimator is consistent if the regressors  $\bar{x}_i$  are independent of the composite error  $(\alpha_i - \alpha + \bar{u}_i)$  in (16).
- This will be the case for the random effects model, but the between estimator effectively discards the time series information in the data set (uses only variation between the cross section observations).
- It is more efficient to use the RE estimator.
- In contrast, for the fixed effects model the between estimator is inconsistent as  $\alpha_i$  is then assumed to be correlated with  $\mathbf{x}_{it}$  and hence  $\bar{\mathbf{x}}_i$ .

# Decomposing OLS Estimates

- In fact, model (14) implies

$$y_{it} - \bar{\bar{y}} = (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)\boldsymbol{\beta} + (\bar{\mathbf{x}}_i - \bar{\bar{\mathbf{x}}})\boldsymbol{\beta} + (u_{it} - \bar{\bar{u}})$$

$$\text{where } \bar{\bar{\mathbf{x}}} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T \mathbf{x}_{it} \text{ and } \bar{\bar{y}} = \frac{1}{NT} \sum_{i=1}^N \sum_{t=1}^T y_{it}.$$

# Decomposing OLS Estimates - Sum of Squares and Cross Products

Let

$$S_{xx}^t = \sum_{i=1}^N \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}})' (\mathbf{x}_{it} - \bar{\mathbf{x}})$$

$$S_{xy}^t = \sum_{i=1}^N \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}}) (y_{it} - \bar{y})$$

$$S_{xx}^w = \sum_{i=1}^N \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)' (\mathbf{x}_{it} - \bar{\mathbf{x}}_i)$$

$$S_{xy}^w = \sum_{i=1}^N \sum_{t=1}^T (\mathbf{x}_{it} - \bar{\mathbf{x}}_i) (y_{it} - \bar{y}_i)$$

total group variation

within group variation

# Decomposing OLS Estimates - Sum of Squares and Cross Products

$$S_{xx}^b = \sum_{i=1}^N T(\bar{\mathbf{x}}_i - \bar{\bar{\mathbf{x}}})'(\bar{\mathbf{x}}_i - \bar{\bar{\mathbf{x}}})$$

$$S_{xy}^b = \sum_{i=1}^N T(\bar{\mathbf{x}}_i - \bar{\bar{\mathbf{x}}})(\bar{y}_i - \bar{\bar{y}})$$

between group variation

It is easy to verify that

$$S_{xx}^t = S_{xx}^w + S_{xx}^b$$

$$S_{xy}^t = S_{xy}^w + S_{xy}^b$$



# Decomposing OLS Estimates - Sum of Squares and Cross Products

- There are, therefore, three possible least squares estimators of  $\beta$  corresponding to the decomposition.
- The least squares estimator is

$$\hat{\beta}^t = [S_{xx}^t]^{-1} S_{xy}^t = [S_{xx}^w + S_{xx}^b]^{-1} [S_{xy}^w + S_{xy}^b]$$

- The **within-groups** estimator is

$$\hat{\beta}^w = [S_{xx}^w]^{-1} S_{xy}^w$$

- An alternative estimator would be the **between-groups** estimator,

$$\hat{\beta}^b = [S_{xx}^b]^{-1} S_{xy}^b$$

# Decomposing OLS Estimates - Sum of Squares and Cross Products

- The standard OLS estimator of model (14) can be decomposed into

$$\hat{\beta}^t = [S_{xx}^t]^{-1} S_{xy}^t = F^w \hat{\beta}^w + F^b \hat{\beta}^b$$

where

$$F^w = ([S_{xx}^w + S_{xx}^b])^{-1} S_{xx}^w$$

$$F^b = (1 - F^w)$$

- which shows that the OLS estimator can be interpreted as weighted average of the within and between estimators with weights that depend on the “within” versus “between” variability of the explanatory factors.

# Decomposing OLS Estimates - Sum of Squares and Cross Products

However:

- in general this is not the most efficient way to exploit jointly the within and between variability;
- it leads to biased and inconsistent estimates of the true causal effect  $\beta$  if the individual specific effects are correlated with the regressors.

# Decomposing RE Estimates

- Like the OLS estimator, also the random effects estimator can be interpreted as a weighted average of the within and between estimator:

$$\hat{\beta}^{re} = F^w \hat{\beta}^w + (F^b) \hat{\beta}^b$$

where

$$F^w = ([S_{xx}^w + \theta S_{xx}^b])^{-1} S_{xx}^w$$

$$\theta = \sigma_u^2 / (\sigma_u^2 + T\sigma_\alpha^2)$$

# Decomposing RE Estimates

- The RE estimator is obtained by estimating the transformed equation

$$y_{it} - \lambda \bar{y}_i = (\mathbf{x}_{it} - \lambda \bar{\mathbf{x}}_i) \boldsymbol{\beta} + (v_{it} - \lambda \bar{v}_i) \quad (17)$$

where

$$\lambda = 1 - [\sigma_u^2 / (\sigma_u^2 + T \sigma_\alpha^2)]^{1/2}$$

- The GLS estimator is simply the pooled OLS estimator of equation (17).

# Decomposing RE Estimates

Note that:

- If  $\sigma_{\alpha}^2 = 0 \rightarrow \lambda = 0$ : the RE estimator is identical to the OLS estimator because there is no individual heterogeneity.
- If  $\sigma_u^2 = 0 \rightarrow \lambda = 1$ : in which case the only existing ignorance would be the individual-specific one captured by  $\alpha_i$  and the RE estimator would be identical to the FE estimator.
- If  $T \rightarrow \infty \rightarrow \lambda = 1$ : if  $T$  goes to infinity the unobserved  $\alpha_i$  "becomes observable" and the RE estimator would be identical to the FE estimator.

# First-Difference (FD) Estimator

- Assumption FD.1:

$$y_{it} = \alpha_i + \mathbf{x}_{it}\boldsymbol{\beta} + u_{it} \quad (18)$$

- Assumption FD.2: We have a random sample in the cross-sectional dimension.
- Assumption FD.3: There are no perfect linear relationships among the explanatory variables.

# First-Difference (FD) Estimator

- Assumption FD.4:

$$E(u_{it} | \mathbf{x}_i, \alpha_i) = 0, t = 1, 2, \dots, T$$

where  $\mathbf{x}_i \equiv (\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots, \mathbf{x}_{iT})$ .

- Lagging the model (18) one period and subtracting gives

$$\Delta y_{it} = \Delta \mathbf{x}_{it} \boldsymbol{\beta} + \Delta u_{it}, \quad t = 2, 3, \dots, T$$

where  $\Delta y_{it} = y_{it} - y_{i,t-1}$ ,  $\Delta \mathbf{x}_{it} = \mathbf{x}_{it} - \mathbf{x}_{i,t-1}$ ,  
and  $\Delta u_{it} = u_{it} - u_{i,t-1}$ .



# First-difference Estimator

- The **first-difference (FD) estimator**,  $\hat{\beta}_{FD}$ , is the OLS estimator from the regression

$$\Delta y_{it} \text{ on } \Delta \mathbf{x}_{it}, \quad t = 2, 3, \dots, T \quad (19)$$

- OLS estimation of the first-differenced equations will be consistent because the strict exogeneity holds in the first-differenced equation:

$$E(\Delta u_{it} | \Delta \mathbf{x}_{i2}, \Delta \mathbf{x}_{i3}, \dots, \Delta \mathbf{x}_{iT}) = 0, \quad t = 2, 3, \dots, T$$

which means the FD estimator is actually unbiased conditional on  $\mathbf{X}$ .

- Under FE.1-FE.6, the FE estimator is asymptotically efficient in the class of estimators using the strict exogeneity assumption FE.4.
- Therefore, the first difference estimator is less efficient than fixed effects under FE.1-FE.6.
- FE.5-FE.6 assumes homoskedasticity and no serial correlation in  $u_{it}$ .
- Assuming that the  $\{u_{it} : t = 1, 2, \dots, T\}$  are serially uncorrelated may be too strong.

# FD or FE?

- An alternative assumption is that the first difference of the idiosyncratic errors,  $\{e_{it} \equiv \Delta u_{it} : t = 2, \dots, T\}$ , are serially uncorrelated (and have constant variance):
- Assumption FD.5:

$$\text{var}(e_{it} | \mathbf{x}_i, \alpha_i) = \sigma_e^2, \quad \text{for all } t = 2, \dots, T$$

- Assumption FD.6:

$$\text{cov}(e_{it}, e_{is} | \mathbf{x}_i, \alpha_i) = 0, \quad \text{for all } t \neq s$$

- Under FD.6 we can write  $u_{it} = u_{i,t-1} + e_{it}$ , so that no serial correlation in the  $e_{it}$  implies that  $u_{it}$  is a random walk.

- When  $T = 2$ , fixed effects estimation and first differencing produce identical estimates and inference.
- When  $T > 2$ , the choice between FD and FE hinges on the assumptions about the idiosyncratic errors,  $u_{it}$ .
- The FE estimator is more efficient under FE.6 - the  $u_{it}$  are serially uncorrelated - while the FD estimator is more efficient when  $u_{it}$  follows a random walk.

# FE vs FD - Testing for Serial Correlation

- Under FD.6, the errors  $e_{it} \equiv \Delta u_{it}$  should be serially uncorrelated.
- As suggested by Wooldridge (2010), we can easily test this assumption given the pooled OLS residuals from regression (16).
- Since the strict exogeneity assumption holds, we can use the  $t$  statistic on  $\hat{\rho}$  in the regression

$$\hat{e}_{it} = \rho \hat{e}_{i,t-1} + error_{it}, \quad t = 3, 4, \dots, T \quad (20)$$

- If the idiosyncratic errors  $\{u_{it} : t = 1, 2, \dots, T\}$  are uncorrelated to begin with,  $\{e_{it} : t = 2, 3, \dots, T\}$  will be autocorrelated.
- In fact, under FE.6 it is easily shown that  $Corr(e_{it}, e_{i,t-1}) = -0.5$ .
- In any case, a finding of significant serial correlation in the  $e_{it}$  warrants computing the robust variance matrix for the FD estimator.

## PANEL DATA MODELS

Assumed Model

	POLS	RE	FE
POLS	Consistent	Consistent	Inconsistent
RE	Consistent	Consistent	Inconsistent
BE	Consistent	Consistent	Inconsistent
FE	Consistent	Consistent	Consistent
FD	Consistent	Consistent	Consistent

# FE Model Pros and Cons: Wrap-up

PROS	CONS
Accounts for permanent observed and unobserved heterogeneity	Needs panel data
Mitigates endogeneity issues due to omitted variables	Identification of the parameters of the regressors is based on intragroup variation over time
Always consistent whether the fixed effect is correlated or uncorrelated with the regressors	Less efficient than POLS and RE if the fixed effect is not correlated with the regressors
	Coefficients of time-invariant covariates are not identified

*"Besides the advantage that panel data allow us to construct and test more complicated behavioral models than purely cross-sectional or time-series data, the use of panel data also provides a means of resolving or reducing the magnitude of a key econometric problem that often arises in empirical studies, namely, the often-heard assertion that the real reason one finds (or does not find) certain effects is because of omitted (mismeasured, not observed) variables that are correlated with explanatory variables."*

Hsiao (2003)



# Estimation of a High-Dimensional FE Model

- Increasingly large data sets make it possible to control for different sources of unobserved heterogeneity
- Examples are:
  - Employer-employee level data
- With high-dimensional models explicit introduction of dummy variables to account for fixed effects is not an option.

# Estimation of a High-Dimensional FE Model

- Increasingly large data sets make it possible to control for different sources of unobserved heterogeneity
- Examples are:
  - Employer-employee level data
  - Hospital-doctor-patient data
- With high-dimensional models explicit introduction of dummy variables to account for fixed effects is not an option.

# Estimation of a High-Dimensional FE Model

- Increasingly large data sets make it possible to control for different sources of unobserved heterogeneity
- Examples are:
  - Employer-employee level data
  - Hospital-doctor-patient data
  - School-class-teacher-student level data
- With high-dimensional models explicit introduction of dummy variables to account for fixed effects is not an option.

# Estimation of a High-Dimensional FE Model

- There are a few commands that allow to estimate a high-dimensional fixed effects model
- The gold standard!
  - the Stata command **reghdfe** - absorbs any number of fixed effects and their interactions, implements IV estimation, much faster and takes advantage of multiple cores, excellent support (github) - by Sergio Correia
  - the R package **lfe** uses the same algorithm as in the Stata module **reg2hdfe** developed by Paulo Guimarães

- Bailey, M. A. (2016), Real Econometrics: The Right Tools to Answer Important Questions, Oxford University Press.
- Cameron, A. C. and P. K. Trivedi (2005), Microeconometrics – Methods and Applications, Cambridge University Press.
- Correia, S. (2016), "REGHDFE: Stata Module to Perform Linear or Instrumental-variable Regression Absorbing any Number of High-dimensional Fixed Effects". Statistical Software Components s457874, Boston College Department of Economics, revised 25 Jul 2015.
- Gaure, S. (2013), "LFE: Linear Group Fixed Effects", The R Journal, 5(2) 104-116.

# References

- Guimarães, P. and Portugal, P. (2010). "A Simple Feasible Procedure to Fit Models with High-dimensional Fixed Effects", the Stata Journal, 10(4) 628-649.
- Hausman, J. A. (1978). "Specification Tests in Econometrics", Econometrica 46(6): 1251-1271.
- Hsiao, C. (2003). Analysis of Panel Data, Cambridge University Press, Cambridge.
- Verbeek, M. (2017), A Guide to Modern Econometrics, John Wiley & Sons.
- Wooldridge, J. M. (2010), Econometric Analysis of Cross Section and Panel Data, MIT Press, Cambridge, MA.

# Appendix

$$\mathbf{y}_i = \begin{bmatrix} \mathbf{y}_{i1} \\ \mathbf{y}_{i2} \\ \vdots \\ \mathbf{y}_{iT} \end{bmatrix}$$

$$\mathbf{X}_i = \begin{bmatrix} \mathbf{X}_{i1}^1 & \mathbf{X}_{i1}^2 & \cdots & \mathbf{X}_{i1}^k \\ \mathbf{X}_{i2}^1 & \mathbf{X}_{i2}^2 & \cdots & \mathbf{X}_{i2}^k \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{X}_{iT}^1 & \mathbf{X}_{iT}^2 & \cdots & \mathbf{X}_{iT}^k \end{bmatrix}$$

$$\mathbf{v}_i = \begin{bmatrix} \mathbf{v}_{i1} \\ \mathbf{v}_{i2} \\ \vdots \\ \mathbf{v}_{iT} \end{bmatrix}$$

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_N \end{bmatrix}$$

$$\mathbf{v} = \begin{bmatrix} \mathbf{v}_1 \\ \mathbf{v}_2 \\ \vdots \\ \mathbf{v}_N \end{bmatrix}$$

- Wald test
- Any hypothesis involving the coefficients of a regression equation can be expressed as one or more restrictions on the coefficient vector:

$$R\beta = r$$

where  $R$  is a  $(q \times k)$  matrix and  $r$  is a  $q$ -element column vector; each row of  $R$  imposes one restriction on the coefficient vector.

- The Wald test uses the point and variance estimates from the unrestricted model to evaluate whether there is evidence that the restrictions are false.



- Given an hypothesis expressed as  $H_0 : R\beta = r$ , we can construct the Wald statistic as:

$$W = \left( \mathbf{R}\hat{\beta} - \mathbf{r} \right)' \left\{ \mathbf{R}\widehat{\text{Var}}(\hat{\beta})\mathbf{R}' \right\}^{-1} \left( \mathbf{R}\hat{\beta} - \mathbf{r} \right) \sim \chi^2_{(h)}$$

where  $\widehat{\text{Var}}(\hat{\beta})$  is a consistent estimator of  $\text{Var}(\hat{\beta})$  and  $h$  the number of restrictions under the null.