
Econometrics 2 - Lecture 6

Models Based on Panel Data

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Types of Data

Population of interest: individuals, households, companies, countries

Types of observations

- Cross-sectional data: observations of all units of a population, or of a representative subset, at one specific point in time
- Time series data: series of observations on units of the population over a period of time
- Panel data (longitudinal data): repeated observations of (the same) population units collected over a number of periods; data set with both a cross-sectional and a time series aspect; multi-dimensional data

Cross-sectional and time series data are one-dimensional, special cases of panel data

Pooling independent cross-sections: (only) similar to panel data

Example: Individual Wages

Verbeek's data set "males"

- Sample of
 - 545 full-time working males
 - each person observed yearly after completion of school in 1980 till 1987
- Variables
 - *wage*: log of hourly wage (in USD)
 - *school*: years of schooling
 - *exper*: $\text{age} - 6 - \text{school}$
 - dummies for union membership, married, black, Hispanic, public sector
 - others

Panel Data in GRET

Three types of data:

- Cross-sectional data: matrix of observations, units over the columns, each row corresponding to the set of variables observed for a unit
- Time series data: matrix of observations, each column a time series, rows correspond to observation periods (annual, quarterly, etc.)
- Panel data: matrix of observations with special data structure
 - Stacked time series: each column one variable, with stacked time series corresponding to observational units
 - Stacked cross sections: each column one variable, with stacked cross sections corresponding to observation periods
 - Use of index variables: index variables defined for units and observation periods

Stacked Data: Examples

Stacked time series

	unit	year	x_1	x_2
1:1	1	2009	1.197	252
2:1	2	2009	1.220	198
3:1	3	2009	1.173	167
...
1:2	1	2010	1.369	269
2:2	2	2010	1.397	212
3:2	3	2010	1.358	201
...

	unit	Year	x_1	x_2
1:1	1	2009	1.197	252
1:2	1	2010	1.369	269
1:3	1	2011	1.675	275
...
2:1	2	2009	1.220	198
2:2	2	2010	1.397	212
2:3	2	2011	1.569	275
...

Stacked cross sections

Panel Data Files

- Files with one record per observation
 - For each unit (individual, company, country, etc.) T records
 - Stacked time series or stacked cross sections
 - Allows easy differencing
 - Time-constant variable: on each record the same value
- Files with one record per unit
 - Each record contains all observations for all T periods
 - Time-constant variables are stored only once

Panel Data

Typically data at micro-economic level (individuals, households, firms), but also at macro-economic level (e.g., countries)

Notation:

- N : Number of cross-sectional units
- T : Number of time periods

Types of panel data:

- Large T , small N : “long and narrow”
- Small T , large N : “short and wide”
- Large T , large N : “long and wide”

Example: Data set “males”: short ($T = 8$) and wide ($N = 545$) panel ($N \gg T$)

Panel Data: Some Examples

Data set “males”: wages and related variables

- short and wide panel ($N = 545$, $T = 8$)
- rich in information (~40 variables)
- unobserved heterogeneity

Grunfeld investment data: investments in plant and equipment by

- $N = 10$ firms
- for each $T = 20$ yearly observations for 1935-1954

Penn World Table: purchasing power parity and national income accounts for

- $N = 189$ countries/territories
- for some or all of the years 1950-2009 ($T \leq 60$)

Use of Panel Data

Econometric models for describing the behaviour of cross-sectional units over time

Panel data models

- Allow controlling individual differences, comparing behaviour, analysing dynamic adjustment, measuring effects of policy changes
- More realistic models
- Allow more detailed or sophisticated research questions

Methodological implications

- Dependence of sample units in time-dimension
- Some variables might be time-constant (e.g., variable *school* in “males”, population size in the Penn World Table dataset)
- Missing values

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Example: Wages and Experience

Data set “males”

- Independent random samples for 1980 and 1987
- $N_{80} = N_{87} = 100$
- Variables: *wage* (log of hourly wage), *exper* (age – 6 – years of schooling)

		1980		1987	
		Full set	sample	Full set	sample
<i>wage</i>	mean	1.39	1.37	1.87	1.89
	st.dev.	0.558	0.598	0.467	0.475
<i>exper</i>	mean	3.01	2.96	10.02	9.99
	st.dev.	1.65	1.29	1.65	1.85
$\exp(wage)$		4.01		6.49	

Pooling of Samples

Independent random samples:

- Pooling gives an independently pooled cross section
- OLS estimates with higher precision, tests with higher power
- Requires
 - the same distributional properties of sampled variables
 - the same relation between variables in the samples

Example: Wages and Experience

Some wage equations (coefficients in bold letters: $p < 0.05$):

- 1980 data

$$wage = 1.315 + 0.026 * exper, R^2 = 0.006$$

- 1987 data

$$wage = 2.441 - \mathbf{0.057} * exper, R^2 = 0.041$$

- pooled 1980 and 1987 data

$$wage = 1.289 + \mathbf{0.052} * exper, R^2 = 0.128$$

- pooled data with dummy d_{87}

$$wage = 1.441 - \mathbf{0.016} * exper + \mathbf{0.583} * d_{87}, R^2 = 0.177$$

- pooled sample with dummy d_{87} and interaction

$$wage = 1.315 + 0.026 * exper + \mathbf{1.126} * d_{87} - \mathbf{0.083} * d_{87} * exper$$

d_{87} : dummy for observations from 1987

Wage Equations

Wage equations, dependent variable: *wage* (log of hourly wage)

		1980	1987	80+87	80+87	80+87
Interc.	coeff	1.315	2.441	1.289	1.441	1.315
	s.e.	0.050	0.120	0.031	0.036	0.045
exper	coeff	0.026	-0.057	0.052	-0.016	0.026
	s.e.	0.014	0.012	0.004	0.009	0.013
d87	coeff				0.583	1.126
	s.e.				0.073	0.141
d87*exper	coeff					-0.083
	s.e.					0.019
	R ² (%)	0.6	4.1	12.8	17.7	19.2

At least the intercept changes from 1980 to 1987

Pooled Independent Cross-sectional Data

Pooling of two independent cross-sectional samples

$$y_{it} = \beta_1 + \beta_2 x_{it} + \varepsilon_{it} \text{ for } i = 1, \dots, N, t = 1, 2$$

- Implicit assumption: identical β_1, β_2 for $i = 1, \dots, N, t = 1, 2$
- OLS-estimation: requires
 - homoskedastic and uncorrelated ε_{it}
$$E\{\varepsilon_{it}\} = 0, \text{Var}\{\varepsilon_{it}\} = \sigma^2 \text{ for } i = 1, \dots, N, t = 1, 2$$
$$\text{Cov}\{\varepsilon_{i1}, \varepsilon_{j2}\} = 0 \text{ for all } i, j \text{ with } i \neq j$$
 - exogenous x_{it}

For the analysis of panel data, often a more realistic model is needed, taking into consideration

- changing coefficients
- correlated error terms
- exogenous regressors

Model with Time Dummy

Model for pooled independent cross-sectional data in presence of changes:

- Dummy variable d : indicator for $t = 2$ ($d_t=0$ for $t=1$, $d_t=1$ for $t=2$)

$$y_{it} = \beta_1 + \beta_2 x_{it} + \beta_3 d_t + \beta_4 d_t^* x_{it} + \varepsilon_{it}$$

allows changes (from $t = 1$ to $t = 2$)

- of intercept from β_1 to $\beta_1 + \beta_3$
- of coefficient of x from β_2 to $\beta_2 + \beta_4$
- Tests for constancy of (1) β_1 or (2) β_1, β_2 over time (cf. Chow test)
 $H_0^{(1)}: \beta_3 = 0$ or $H_0^{(2)}: \beta_3 = \beta_4 = 0$
- Similarly testing for constancy of σ^2 over time

Generalization to more than two time periods

Example: Wages and Experience

Wage equation

$$wage_{it} = \beta_1 + \beta_2 exper_{it} + \beta_3 d_t + \varepsilon_{it}$$

Wages might depend also on other variables; omitted variables are covered by the error terms

- *black*: time-constant variable, omission may cause autocorrelation of error terms; similar other time-constant factors like *hisp*
- *mar* (married): variable which is for many (not all) units time-constant, similar *rural*, *union*, *ne* (living in north east), etc.; omission may cause autocorrelation
- *school*: omission may cause endogeneity of *exper*
- Unobserved and unobservable variables can have similar effects, e.g., parental background, attitudes, etc.

Problems with Sample Pooling

The analysis of the data (y_{it}, x_{it}) , $i = 1, \dots, N$, $t = 1, 2$, by OLS estimation of the parameters of model

$$y_{it} = \beta_1 + \beta_2 x_{it} + \varepsilon_{it}$$

(or extensions based on a year dummy for $t=2$) may not fulfil usual requirements

- The independence assumption across time may be unrealistic
- Main reason is that effects of non-measured and non-measurable variables are only covered by the error terms
- Also exogeneity of regressors may be unrealistic

Consequences: OLS-estimates

- biased and inconsistent
- not efficient

Panel data models allow more adequate analyses

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Models for Panel Data

Model for y , based on panel data from N cross-sectional units and T periods

$$y_{it} = \beta_0 + x_{it}'\beta_1 + \varepsilon_{it}$$

$i = 1, \dots, N$: sample unit

$t = 1, \dots, T$: time period of sample

x_{it} and β_1 : K -vectors

- β_0 and β_1 : represent intercept and K regression coefficients; are assumed to be identical for all units and all time periods
- ε_{it} : represents unobserved factors that may affect y_{it}
 - Assumption that ε_{it} are uncorrelated over time not realistic
 - Standard errors of OLS estimates misleading, OLS estimation not efficient

Fixed Effects Model

The general model

$$y_{it} = \beta_0 + x_{it}'\beta_1 + \varepsilon_{it}$$

- Specification for the error terms: two components

$$\varepsilon_{it} = \alpha_i + u_{it}$$

- α_i unit-specific, time-constant factors, also called unobserved (individual) heterogeneity; may be correlated with x_{it}
- $u_{it} \sim \text{IID}(0, \sigma_u^2)$; uncorrelated over time; represents unobserved factors that change over time, also called idiosyncratic or time-varying error
- ε_{it} : also called composite error

- Fixed effects (FE) model

$$y_{it} = \sum_j \alpha_j d_{ij} + x_{it}'\beta_1 + u_{it}$$

d_{ij} : dummy variable for unit i : $d_{ij} = 1$ if $i = j$, otherwise $d_{ij} = 0$

- Overall intercept omitted; unit-specific intercepts α_i

Random Effects Model

Starting point is again the model

$$y_{it} = \beta_0 + x_{it}'\beta_1 + \varepsilon_{it}$$

with composite error $\varepsilon_{it} = \alpha_i + u_{it}$

- Specification for the error terms:
 - $u_{it} \sim \text{IID}(0, \sigma_u^2)$; uncorrelated over time
 - $\alpha_i \sim \text{IID}(0, \sigma_a^2)$; represents all unit-specific, time-constant factors; correlation of error terms over time only via the α_i
 - α_i and u_{it} are assumed to be mutually independent and independent of x_{js} for all j and s
- Random effects (RE) model
$$y_{it} = \beta_0 + x_{it}'\beta_1 + \alpha_i + u_{it}$$
- Unbiased and consistent ($N \rightarrow \infty$) estimation of β_0 and β_1
- Efficient estimation of β_0 and β_1 : takes error covariance structure into account; GLS estimation

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Fixed Effects (FE) Model

Model for y , based on panel data for T periods

$$y_{it} = \alpha_i + x_{it}'\beta + u_{it}, u_{it} \sim \text{IID}(0, \sigma_u^2)$$

$i = 1, \dots, N$: sample unit

$t = 1, \dots, T$: time period of sample

- α_i : fixed parameter, represents all unit-specific, time-constant factors, unobserved (individual) heterogeneity
- x_{it} : all K components are assumed to be independent of all u_{it} ; may be correlated with α_i

Model with dummies $d_{ij} = 1$ for $i = j$ and 0 otherwise:

$$y_{it} = \sum_j \alpha_i d_{ij} + x_{it}'\beta + u_{it}$$

- Number of coefficients: $N + K$
- Main interest: estimators for β

FE Model Parameters: Estimation

FE model with dummies $d_{ij} = 1$ for $i = j$ and 0 otherwise:

$$y_{it} = \sum_j \alpha_i d_{ij} + x_{it}'\beta + u_{it}$$

Number of coefficients: $N + K$

Various estimation procedures are available

- Least squares dummy variable (LSDV) estimator
- Within or fixed effects estimator
- First-difference estimator

A special case

- Differences-in-differences (DD or DID or D-in-D) estimator

Least Squares Dummy Variable (LSDV) Estimator

Estimation procedure for $N + K$ parameters β and α_i of the FE model

$$y_{it} = \sum_j \alpha_i d_{ij} + x_{it}'\beta + u_{it}$$

OLS estimation

- NT observations for estimating $N + K$ coefficients
- Numerically costly, not attractive
- Estimates for α_i usually not of interest

Fixed effects and first-difference estimators are more attractive

Fixed Effects Estimation

Within transformation: transforms y_{it} into time-demeaned \check{y}_{it} by subtracting the average $\bar{y}_i = (\sum_t y_{it})/T$:

$$\check{y}_{it} = y_{it} - \bar{y}_i$$

analogously \check{x}_{it} and \check{u}_{it} , for $i = 1, \dots, N$, $t = 1, \dots, T$

Model in time-demeaned variables

$$\check{y}_{it} = \check{x}_{it}'\beta + \check{u}_{it}$$

- Pooled OLS estimator b_{FE} for β
- b_{FE} : “fixed effects estimator”, also called “within estimator”
- Uses time variation in y and x within each cross-sectional observation; explains deviations of y_{it} from \bar{y}_i (not of \bar{y}_i from \bar{y}_j !)

GRETL: Model > Panel > Fixed or random effects ...

The Fixed Effects Estimator

FE model

$$y_{it} = \alpha_i + x_{it}'\beta + u_{it}, u_{it} \sim \text{IID}(0, \sigma_u^2)$$

x_{it} are assumed to be independent of all u_{it} but may be correlated with α_i

Estimation of β from the model in time-demeaned variables

$$\check{y}_{it} = \check{x}_{it}'\beta + \check{u}_{it}$$

gives

$$b_{FE} = (\sum_j \sum_t \check{x}_{it} \check{x}_{it}')^{-1} \sum_j \sum_t \check{x}_{it} \check{y}_{it}$$

- Time-demeaning differences away time-constant factors α_i
- Under the assumption that x_{it} are independent of all u_{it} : b_{FE} is unbiased
- b_{FE} coincides with LSDV estimator

Wage Equations

Wage equations, dependent variable: *wage* (log of hourly wage)

		Pooled 80+87	FE 80+87	FE 80+87	FE 80+87	FE 80...87
Interc.	coeff	1.289	1.285	1.432	1.307	1.237
	s.e.	0.031	0.031	0.036	0.045	0.016
exper	coeff	0.052	0.053	-0.013	0.029	0.063
	s.e.	0.004	0.004	0.009	0.013	0.002
d87	coeff			0.564	1.107	
	s.e.			0.073	0.141	
d87*exper	coeff				-0.083	
	s.e.				0.019	
	adjR ² (%)	12.8	13.7	18.1	19.5	55.6

Properties of Fixed Effects Estimator

$$b_{FE} = (\sum_i \sum_t \ddot{x}_{it} \ddot{x}_{it}')^{-1} \sum_i \sum_t \ddot{x}_{it} \ddot{y}_{it}$$

- Unbiased if all x_{it} are independent of all u_{it}
- Normally distributed if normality of u_{it} is assumed
- Consistent (for $N \rightarrow \infty$) if x_{it} are strictly exogenous, i.e., $E\{x_{it} u_{is}\} = 0$ for all s, t
- Asymptotically normally distributed
- Covariance matrix

$$V\{b_{FE}\} = \sigma_u^2 (\sum_i \sum_t \ddot{x}_{it} \ddot{x}_{it}')^{-1}$$

- Estimated covariance matrix: substitution of σ_u^2 by

$$s_u^2 = (\sum_i \sum_t \tilde{u}_{it} \tilde{u}_{it}') / [N(T-1)]$$

with the residuals $\tilde{u}_{it} = \ddot{y}_{it} - \ddot{x}_{it}' b_{FE}$

- Attention! The standard OLS estimate of the covariance matrix underestimates the true values

Estimator for α_i

Time-constant factors $\alpha_i, i = 1, \dots, N$

Estimates based on the fixed effects estimator b_{FE}

$$a_i = \bar{y}_i - \bar{x}_i' b_{FE}$$

with averages over time \bar{y}_i and \bar{x}_i for the i -th unit

- Consistent (for $T \rightarrow \infty$) if x_{it} are strictly exogenous
- Potentially interesting aspects of estimates a_i
 - Distribution of the $a_i, i = 1, \dots, N$
 - Value of a_i for unit i of special interest

First-Difference Estimator

Elimination of time-constant factors α_i by differencing

$$\Delta y_{it} = y_{it} - y_{i,t-1} = \Delta x_{it}'\beta + \Delta u_{it}$$

Δx_{it} and Δu_{it} analogously defined as $\Delta y_{it} = y_{it} - y_{i,t-1}$

First-difference estimator: OLS estimation

$$b_{FD} = (\sum_i \sum_t \Delta x_{it} \Delta x_{it}')^{-1} \sum_i \sum_t \Delta x_{it} \Delta y_{it}$$

Properties

- Consistent (for $N \rightarrow \infty$) under slightly weaker conditions than b_{FE}
- Slightly less efficient than b_{FE} due to serial correlations of the Δu_{it}
- For $T = 2$, b_{FD} and b_{FE} coincide

Differences-in-Differences Estimator

Natural experiment or quasi-experiment:

- Exogenous event, e.g., a new law, changes in operating conditions
- Treatment group, control group
- Assignment to groups not (like in a true experiment) at random
- Data: before event, after event

Model for response y_{it}

$$y_{it} = \delta r_{it} + \mu_t + \alpha_i + u_{it}, \quad i = 1, \dots, N, \quad t = 1 \text{ (before), } 2 \text{ (after event)}$$

- Dummy $r_{it} = 1$ if i -th unit in treatment group, $r_{it} = 0$ otherwise
- δ : treatment effect
- Fixed effects model (for differencing away time-constant factors):

$$\Delta y_{it} = y_{i2} - y_{i1} = \delta \Delta r_{it} + \mu_0 + \Delta u_{it}$$

with $\mu_0 = \mu_2 - \mu_1$

Wage Differences 1980 - 1987

Effect of ethnicity

- *wage* (log of hourly wage) : increases from 1.419 (1980) to 1.892 (1987)
- i.e., increase of hourly wage from USD 4.13 (1980) to 6.63 (1987)

Does the wage increase depend on ethnicity?

- Dummy $black_{it} = 1$ if i -th person is afro-american, $black_{it} = 0$ otherwise
- Model for *wage*:

$$wage_{it} = \mu_t + \alpha_i + u_{it}, \quad i = 1, \dots, N, \quad t = 1980, 1987$$

- α_i : time-constant factors, e.g., schooling, rural, industry, etc.
- Model for differences with $\mu_0 = \mu_{1987} - \mu_{1980}$

$$\Delta wage_{it} = \mu_0 + \delta black_{it} + \Delta u_{it}$$

Wage Differences, cont'd

Increase of *wage* (log of hourly wage)

$$\Delta wage_{it} = \mu_0 + \delta black_{it} + \Delta u_{it}$$

OLS-estimation gives ($N = 545$, 63 afro-americans)

	μ_0	δ	adj R ²
Estimate	0.491	-0.154	0.47
Std.err.	0.027	0.081	

Differences in *wage* and in hourly wages

	μ_0	$\mu_0 + \delta$	all
	<i>black</i> = 0	<i>black</i> = 1	
<i>wage</i> (average)	0.491	0.337	0.473
hourly wages	1.634	1.401	1.605
Increase (%)	63.4	40.1	60.5

Estimator of Treatment Effect

Effect of treatment (event) by comparing units

- with and without treatment
- before and after treatment

Model for panel data y_{it}

$$y_{it} = \delta r_{it} + \mu_t + \alpha_i + u_{it}, \quad i = 1, \dots, N, \quad t = 1 \text{ (before), } 2 \text{ (after event)}$$

Differences-in-differences (DD or DID or D-in-D) estimator of treatment effect δ

$$d_{DD} = \Delta \bar{y}^{\text{treated}} - \Delta \bar{y}^{\text{untreated}}$$

$\Delta \bar{y}^{\text{treated}}$: average difference $y_{i2} - y_{i1}$ of treatment group units

$\Delta \bar{y}^{\text{untreated}}$: average difference $y_{i2} - y_{i1}$ of control group units

- Treatment effect δ measured as difference between changes of y with and without treatment
- d_{DD} consistent if $E\{\Delta r_{it} \Delta u_{it}\} = 0$
- Allows correlation between time-constant factors α_i and r_{it}

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Random Effects Model

Model:

$$y_{it} = \beta_0 + x_{it}'\beta + \alpha_i + u_{it}, u_{it} \sim \text{IID}(0, \sigma_u^2)$$

- Time-constant factors α_i : stochastic variables with identical distribution for all units

$$\alpha_i \sim \text{IID}(0, \sigma_a^2)$$

- Attention! More information about α_i than in the fixed effects model
- $\alpha_i + u_{it}$: error term with two components
 - Unit-specific component α_i , time-constant
 - Remainder u_{it} , assumed to be uncorrelated over time
- α_i, u_{it} : mutually independent, independent of x_{js} for all j and s
- OLS estimators for β_0 and β are unbiased, consistent, not efficient (see next slide)

GLS Estimator

$\alpha_i i_T + u_i$: T -vector of error terms for i -th unit, T -vector $i_T = (1, \dots, 1)'$

$\Omega = \text{Var}\{\alpha_i i_T + u_i\}$: Covariance matrix of $\alpha_i i_T + u_i$

$$\Omega = \sigma_a^2 i_T i_T' + \sigma_u^2 I_T$$

Inverted covariance matrix

$$\Omega^{-1} = \sigma_u^{-2} \{ [I_T - (i_T i_T')/T] + \psi (i_T i_T')/T \}$$

with $\psi = \sigma_u^2 / (\sigma_u^2 + T\sigma_a^2)$

$(i_T i_T')/T$: transforms into averages

$I_T - (i_T i_T')/T$: transforms into deviations from average

GLS estimator

$$b_{\text{GLS}} = [\Sigma_i \Sigma_t \ddot{x}_{it} \ddot{x}_{it}' + \psi T \Sigma_i (\dot{x}_i - \bar{x})(\dot{x}_i - \bar{x})']^{-1} [\Sigma_i \Sigma_t \ddot{x}_{it} \ddot{y}_{it} + \psi T \Sigma_i (\dot{x}_i - \bar{x})(\bar{y}_i - \bar{y})]$$

with the average \bar{y} over all i and t , analogous \bar{x}

- $\psi = 0$: b_{GLS} coincides with b_{FE} ; b_{GLS} and b_{FE} equivalent for large T
- $\psi = 1$: b_{GLS} coincides with the OLS estimators for β_0 and β

Between Estimator

Model for individual means \bar{y}_i and \bar{x}_i :

$$\bar{y}_i = \beta_0 + \bar{x}_i' \beta + \alpha_i + \bar{u}_i, \quad i = 1, \dots, N$$

OLS estimator

$$b_B = [\sum_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})']^{-1} \sum_i (\bar{x}_i - \bar{x})(\bar{y}_i - \bar{y})$$

is called the between estimator

- Consistent if x_{it} strictly exogenous, uncorrelated with α_i
- GLS estimator can be written as

$$b_{GLS} = \Delta b_B + (I_K - \Delta) b_{FE}$$

Δ : weighting matrix, proportional to the inverse of $\text{Var}\{b_B\}$

- Matrix-weighted average of between estimator b_B and within estimator b_{FE}
- The more accurate b_B the more weight has b_B in b_{GLS}
- b_{GLS} : optimal combination of b_B and b_{FE} , more efficient than b_B and b_{FE}

GLS Estimator: Properties

$$b_{\text{GLS}} = [\Sigma_i \Sigma_t \ddot{x}_{it} \ddot{x}_{it}' + \psi T \Sigma_i (\dot{x}_i - \bar{x})(\dot{x}_i - \bar{x})']^{-1} [\Sigma_i \Sigma_t \ddot{x}_{it} \ddot{y}_{it} + \psi T \Sigma_i (\dot{x}_i - \bar{x})(\bar{y}_i - \bar{y})]$$

- Unbiased, if x_{it} are independent of all α_i and u_{it}
- Consistent for N or T or both tending to infinity if
 - $E\{\ddot{x}_{it} u_{it}\} = 0$
 - $E\{\dot{x}_i u_{it}\} = 0, E\{\ddot{x}_{it} \alpha_i\} = 0$
 - These conditions are required also for consistency of b_B
- More efficient than the between estimator b_B and the within estimator b_{FE} ; also more efficient than the OLS estimator

Random Effects Estimator

EGLS or Balestra-Nerlove estimator: Calculation of b_{GLS} from model

$$y_{it} - \vartheta \bar{y}_i = \beta_0(1 - \vartheta) + (x_{it} - \vartheta \bar{x}_i)' \beta + v_{it}$$

with $\vartheta = 1 - \psi^{1/2}$, $v_{it} \sim \text{IID}(0, \sigma_v^2)$

quasi-demeaned $y_{it} - \vartheta \bar{y}_i$ and $x_{it} - \vartheta \bar{x}_i$

Two step estimator:

1. Step 1: Transformation parameter ψ calculated from (method by Swamy & Arora)

- within estimation: $s_u^2 = (\sum_i \sum_t \tilde{v}_{it} \tilde{v}_{it}) / [N(T-1)]$

- between estimation: $s_B^2 = (1/N) \sum_i (\bar{y}_i - b_{0B} - \bar{x}_i' b_B)^2 = s_a^2 + (1/T) s_u^2$

- $s_a^2 = s_B^2 - (1/T) s_u^2$

2. Step 2:

- Calculation of $1 - [s_u^2 / (s_u^2 + T s_a^2)]^{1/2}$ for parameter ϑ

- Transformation of y_{it} and x_{it}

- OLS estimation gives the random effect estimator b_{RE} for β

Random Effects Estimator: Properties

b_{RE} : EGLS estimator of β from

$$y_{it} - \vartheta \bar{y}_i = \beta_0(1 - \vartheta) + (x_{it} - \vartheta \bar{x}_i)' \beta + v_{it}$$

with $\vartheta = 1 - \psi^{1/2}$, $\psi = \sigma_u^2 / (\sigma_u^2 + T\sigma_a^2)$

- Asymptotically normally distributed under weak conditions
- Covariance matrix

$$\text{Var}\{b_{RE}\} = \sigma_u^2 [\Sigma_i \Sigma_t \bar{x}_{it} \bar{x}_{it}' + \psi T \Sigma_i (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})']^{-1}$$

- More efficient than the within estimator b_{FE} (if $\psi > 0$)

Wage Equations, 1980-1987

Dependent variable: *wage* (log of hourly wage)

	Between	Fixed Effects	Random Effects	Pooled OLS
Intercept	0.511	1.053	-0.079	0.049
<i>school</i>	0.089***	--	0.100***	0.095***
<i>exper</i>	-0.032	0.118***	0.111***	0.087***
<i>exper2</i>	0.004	-0.004***	-0.004***	-0.003***
<i>union</i>	0.262***	0.082***	0.109***	0.179***
<i>mar</i>	0.184***	0.045**	0.064***	0.126***
<i>black</i>	-0.141***	--	-0.149***	-0.150***
<i>rural</i>	0.188***	0.049*	-0.026	-0.138***

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Summary of Estimators

- Between estimator
- Fixed effects (within) estimator
- Combined estimators
 - OLS estimator
 - Random effects (EGLS) estimator
- First-difference estimator

Estimator		Consistent, if
Between	b_B	x_{it} strictly exog, x_{it} and α_i uncorr
Fixed effects	b_{FE}	x_{it} strictly exog
OLS	b	x_{it} and α_i uncorr, x_{it} and u_{it} contemp. uncorr
Random effects	b_{RE}	conditions for b_B and b_{FE} are met
First-difference	b_{FD}	$E\{\dot{x}_{it} \ddot{u}_{it}\} = 0$

Fixed Effects or Random Effects?

Random effects model

$$E\{y_{it} | x_{it}\} = x_{it}'\beta$$

- Large values N ; of interest: population characteristics (β), not characteristics of individual units (α_i)
- More efficient estimation of β , given adequate specification of the time-constant model characteristics

Fixed effects model

$$E\{y_{it} | x_{it}\} = x_{it}'\beta + \alpha_i$$

- Of interest: besides population characteristics (β), also characteristics of individual units (α_i), e.g., of countries or companies; rather small values N
- Large values of N , if x_{it} and α_i correlated: consistent estimator b_{FE} in case of correlated x_{it} and α_i

Diagnostic Tools

- Test of common intercept of all units
 - Applied to pooled OLS estimation: Rejection indicates preference for fixed or random effects model
 - Applied to fixed effects estimation: Non-rejection indicates preference for pooled OLS estimation
- Hausman test (of correlation between x_{it} and α_i):
 - Null-hypothesis that GLS estimates are consistent
 - Rejection indicates preference for fixed effects model
- Test of non-constant variance of the error terms, Breusch-Pagan test
 - Rejection indicates preference for fixed or random effects model
 - Non-rejection indicates preference for pooled OLS estimation

Hausman Test

Tests of correlation between x_{it} and α_i

H_0 : x_{it} and α_i are uncorrelated

Test statistic:

$$\xi_H = (b_{FE} - b_{RE})' [\tilde{V}\{b_{FE}\} - \tilde{V}\{b_{RE}\}]^{-1} (b_{FE} - b_{RE})$$

with estimated covariance matrices $\tilde{V}\{b_{FE}\}$ and $\tilde{V}\{b_{RE}\}$

- b_{RE} : consistent if x_{it} and α_i are uncorrelated
- b_{FE} : consistent also if x_{it} and α_i are correlated

Under H_0 : $\text{plim}(b_{FE} - b_{RE}) = 0$

- ξ_H asymptotically chi-squared distributed with K d.f.
- K : dimension of x_{it} and β

Hausman test may indicate also other types of misspecification

Robust Inference

Consequences of heteroskedasticity and autocorrelation of the error term:

- Standard errors and related tests are incorrect
- Inefficiency of estimators

Robust covariance matrix for estimator b of β from $y_{it} = x_{it}'\beta + \varepsilon_{it}$

$$b = (\sum_i \sum_t x_{it} x_{it}')^{-1} \sum_i \sum_t x_{it} y_{it}$$

- Adjustment of covariance matrix similar to Newey-West: assuming uncorrelated error terms for different units ($E\{\varepsilon_{it} \varepsilon_{js}\} = 0$ for all $i \neq j$)

$$V\{b\} = (\sum_i \sum_t x_{it} x_{it}')^{-1} \sum_i \sum_t \sum_s e_{it} e_{is} x_{it} x_{is}' (\sum_i \sum_t x_{it} x_{it}')^{-1}$$

e_{it} : OLS residuals

- Allows for heteroskedasticity and autocorrelation within units
- Called panel-robust estimate of the covariance matrix

Analogous variants of the Newey-West estimator for robust covariance matrices of random effects and fixed effects estimators

Testing for Autocorrelation and Heteroskedasticity

Tests for heteroskedasticity and autocorrelation in random effects model error terms

- Computationally cumbersome

Tests based on fixed effects model residuals

- Easier case
- Applicable for testing in both fixed and random effects case

Test for Autocorrelation

Durbin-Watson test for autocorrelation in the fixed effects model

- Error term $u_{it} = \rho u_{i,t-1} + v_{it}$
 - Same autocorrelation coefficient ρ for all units
 - v_{it} iid across time and units
- Test of $H_0: \rho = 0$ against $\rho > 0$
- Adaptation of Durbin-Watson statistic

$$dw_p = \frac{\sum_{i=1}^N \sum_{t=2}^T (\hat{u}_{it} - \hat{u}_{i,t-1})^2}{\sum_{i=1}^N \sum_{t=1}^T \hat{u}_{it}^2}$$

- Tables with critical limits d_U and d_L for K , T , and N ; e.g., Verbeek's Table 10.1

Test for Heteroskedasticity

Breusch-Pagan test for heteroskedasticity of fixed effects model residuals

- $V\{u_{it}\} = \sigma^2 h(z_{it}'\gamma)$; unknown function $h(\cdot)$ with $h(0)=1$, J -vector z
- $H_0: \gamma = 0$, homoskedastic u_{it}
- Auxiliary regression of squared residuals on intercept and regressors z
- Test statistic: $N(T-1)$ times R^2 of auxiliary regression
- Chi-squared distribution with J d.f. under H_0

Wage Equations, 1980-1987

Fixed effects estimation, standard and HAC standard errors

	Coeff.	s.e.	HAC s.e.	Δ
Intercept	1.053	0.0276	0.0384	1.39
<i>exper</i>	0.118	0.0084	0.0108	1.29
<i>exper2</i>	-0.004	0.0006	0.0007	1.17
<i>union</i>	0.082	0.0193	0.0227	1.18
<i>mar</i>	0.045	0.0183	0.0210	1.15
<i>rural</i>	0.049	0.0290	0.0391	1.35

Δ : ratio of HAC s.e. to s.e.

Goodness-of-Fit

Goodness-of-fit measures for panel data models: different from OLS estimated regression models

- Focus may be on within or between variation in the data
- The usual R^2 measure relates to OLS-estimated models

Definition of goodness-of-fit measures: squared correlation coefficients between actual and fitted values

- R^2_{within} : squared correlation between within transformed actual and fitted y_{it} ; maximized by within estimator
- R^2_{between} : based upon individual averages of actual and fitted y_{it} ; maximized by between estimator
- R^2_{overall} : squared correlation between actual and fitted y_{it} ; maximized by OLS

Corresponds to the decomposition

$$[1/TM]\sum_i\sum_t(y_{it} - \bar{y})^2 = [1/TM]\sum_i\sum_t(y_{it} - \bar{y}_i)^2 + [1/M]\sum_i(\bar{y}_i - \bar{y})^2$$

Goodness-of-Fit, cont'd

Fixed effects estimator b_{FE}

- Explains the within variation
- Maximizes R^2_{within}

$$R^2_{within}(b_{FE}) = \text{corr}^2\{\hat{y}_{it}^{FE} - \hat{y}_i^{FE}, y_{it} - \bar{y}_i\}$$

Between estimator b_B

- Explains the between variation
- Maximizes $R^2_{between}$

$$R^2_{between}(b_B) = \text{corr}^2\{\hat{y}_i^B, \bar{y}_i\}$$

Wage Equations, 1980-1987

Dependent variable: *wage* (log of hourly wage)

	Between	F.E.	R.E.	OLS
Intercept	0.511	1.053	-0.079	0.049
<i>school</i>	0.089***	--	0.100***	0.095***
<i>exper</i>	-0.032	0.118***	0.111***	0.087***
<i>exper2</i>	0.004	-0.004***	-0.004***	-0.003***
<i>union</i>	0.262***	0.082***	0.109***	0.179***
<i>mar</i>	0.184***	0.045**	0.064***	0.126***
<i>black</i>	-0.141***	--	-0.149***	-0.150***
<i>rural</i>	0.188***	0.049*	-0.026	-0.138***
R ² (%)	16.07	5.66	18.42	19.70

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

Panel Data and GRET

Estimation of panel models

Pooled OLS

- `Model > Ordinary Least Squares ...`
- Special diagnostics on the output window: `Tests > Panel diagnostics`

Fixed and random effects models

- `Model > Panel > Fixed or random effects...`
- Provide diagnostic tests
 - Fixed effects model: Test for common intercept of all units
 - Random effects model: Breusch-Pagan test, Hausman test

Further estimation procedures

- Between estimator
- Weighted least squares
- Instrumental variable panel procedure

Your Homework

1. Use Verbeek's data set MALES which contains panel data for 545 full-time working males over the period 1980-1987. Estimate a wage equation which explains the individual log wages by the variables years of schooling, years of experience and its squares, and dummy variables for union membership, being married, black, Hispanic, and working in the public sector. Use (i) pooled OLS, (ii) the between and (iii) the within estimator, and (iv) the random effects estimator.

Contents

- Panel Data
- Pooling Independent Cross-sectional Data
- Panel Data: Pooled OLS Estimation
- Panel Data Models
- Fixed Effects Model
- Random Effects Model
- Analysis of Panel Data Models
- Panel Data in **GRETL**

A Model for Two-period Panel Data

Model for y , based on panel data for two periods:

$$\begin{aligned}y_{it} &= \beta_0 + \delta_1 d_t + \beta_1 x_{it} + \varepsilon_{it} \\ &= \beta_0 + \delta_1 d_t + \beta_1 x_{it} + \alpha_i + u_{it}\end{aligned}$$

$i = 1, \dots, N$: sample units of the panel

$t = 1, 2$: time period of sample

d_t : dummy for period $t = 2$

- $\varepsilon_{it} = \alpha_i + u_{it}$: composite error
- α_i : represents all unit-specific, time-constant factors; also called unobserved (individual) heterogeneity
- u_{it} : represents unobserved factors that change over time, also called idiosyncratic or time-varying error
 - u_{it} (and ε_{it}) may be correlated over time for the same unit

Model is called unobserved or fixed effects model

Estimation of the Parameters of Interest

Parameter of interest is β_1

Estimation concepts:

1. Pooled OLS estimation of β_1 from $y_{it} = \beta_0 + \delta_1 d_t + \beta_1 x_{it} + \varepsilon_{it}$ based on the pooled dataset, $\varepsilon_{it} = \alpha_i + u_{it}$
 - ❑ Inconsistent, if x_{it} and α_i are correlated
 - ❑ Incorrect standard errors due to correlation of u_{it} (and ε_{it}) over time; typically too small standard errors
2. First-difference estimator: OLS estimation of β_1 from the first-difference equation
$$\Delta y_i = y_{i1} - y_{i2} = \delta_1 + \beta_1 \Delta x_i + \Delta u_i$$
 - ❑ α_i are differenced away; correlation of x_{it} and α_i not relevant
 - ❑ Correlation of u_{it} (and ε_{it}) over time not relevant
3. Fixed effects estimation (see below)

Wage Equations

Data set “males”, cross-sectional samples for 1980 and 1987

(1): OLS estimation in pooled sample

(2): OLS estimation in pooled sample
with interaction dummy

		(1)	(2)
interc.	coeff	1.045	1.241
	s.e.	0.048	0.056
exper	coeff	0.160	0.073
	s.e.	0.017	0.021
exper ²	coeff	-0.008	-0.006
		0.001	0.001
d87	coeff		0.479
	s.e.		0.076
	R ² (%)	16.2	19.0

Pooled OLS Estimation

Model for y , based on panel data from T periods:

$$y_{it} = x_{it}'\beta + \varepsilon_{it}$$

Pooled OLS estimation of β

- Assumes equal unit means α_i
- Consistent if x_{it} and ε_{it} (at least contemporaneously) uncorrelated
- Diagnostics of interest:
 - Test whether panel data structure to be taken into account
 - Test whether fixed or random effects model preferable

In **GRET**L: output window of OLS estimation applied to panel data structure offers a special test: `Test > Panel diagnostics`

- Tests H_0 : pooled model preferable to fixed effects and random effects model
- Hausman test (H_0 : random effects model preferable to fixed effects model)