

# Panel Data: Fixed and Random Effects

## 1 Introduction

In panel data, individuals (persons, firms, cities, ... ) are observed at several points in time (days, years, before and after treatment, ...). This handout focuses on panels with relatively few time periods (small  $T$ ) and many individuals (large  $N$ ).

This handout introduces the two basic models for the analysis of panel data, the fixed effects model and the random effects model, and presents consistent estimators for these two models. The handout does not cover so-called dynamic panel data models.

Panel data are most useful when we suspect that the outcome variable depends on explanatory variables which are not observable but correlated with the observed explanatory variables. If such omitted variables are constant over time, panel data estimators allow to consistently estimate the effect of the observed explanatory variables.

## 2 The Econometric Model

Consider the multiple linear regression model for individual  $i = 1, \dots, N$  who is observed at several time periods  $t = 1, \dots, T$

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + c_i + u_{it}$$

where  $y_{it}$  is the dependent variable,  $x'_{it}$  is a  $K$ -dimensional row vector of time-varying explanatory variables and  $z'_i$  is a  $M$ -dimensional row vector of time-invariant explanatory variables excluding the constant,  $\alpha$  is the intercept,  $\beta$  is a  $K$ -dimensional column vector of parameters,  $\gamma$  is a  $M$ -dimensional column vector of parameters,  $c_i$  is an *individual-specific effect* and  $u_{it}$  is an *idiosyncratic* error term.

We will assume throughout this handout that each individual  $i$  is observed in all time periods  $t$ . This is a so-called *balanced panel*. The treatment of unbalanced panels is straightforward but tedious.

The  $T$  observations for individual  $i$  can be summarized as

$$y_i = \begin{bmatrix} y_{i1} \\ \vdots \\ y_{it} \\ \vdots \\ y_{iT} \end{bmatrix}_{T \times 1} \quad X_i = \begin{bmatrix} x'_{i1} \\ \vdots \\ x'_{it} \\ \vdots \\ x'_{iT} \end{bmatrix}_{T \times K} \quad Z_i = \begin{bmatrix} z'_i \\ \vdots \\ z'_i \\ \vdots \\ z'_i \end{bmatrix}_{T \times M} \quad u_i = \begin{bmatrix} u_{i1} \\ \vdots \\ u_{it} \\ \vdots \\ u_{iT} \end{bmatrix}_{T \times 1}$$

and  $NT$  observations for all individuals and time periods as

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_N \end{bmatrix}_{NT \times 1} \quad X = \begin{bmatrix} X_1 \\ \vdots \\ X_i \\ \vdots \\ X_N \end{bmatrix}_{NT \times K} \quad Z = \begin{bmatrix} Z_1 \\ \vdots \\ Z_i \\ \vdots \\ Z_N \end{bmatrix}_{NT \times M} \quad u = \begin{bmatrix} u_1 \\ \vdots \\ u_i \\ \vdots \\ u_N \end{bmatrix}_{NT \times 1}$$

The data generation process (dgp) is described by:

*PL1: Linearity*

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + c_i + u_{it} \text{ where } E[u_{it}] = 0 \text{ and } E[c_i] = 0$$

The model is linear in parameters  $\alpha$ ,  $\beta$ ,  $\gamma$ , effect  $c_i$  and error  $u_{it}$ .

*PL2: Independence*

$$\{X_i, z_i, y_i\}_{i=1}^N \text{ i.i.d. (independent and identically distributed)}$$

The observations are independent across individuals but not necessarily across time. This is guaranteed by random sampling of individuals.

*PL3: Strict Exogeneity*

$$E[u_{it}|X_i, z_i, c_i] = 0 \text{ (mean independent)}$$

The idiosyncratic error term  $u_{it}$  is assumed uncorrelated with the explanatory variables of all past, current and future time periods of the same individual. This is a strong assumption which e.g. rules out lagged dependent variables. *PL3* also assumes that the idiosyncratic error is uncorrelated with the individual specific effect.

*PL4: Error Variance*

- a)  $V[u_i|X_i, z_i, c_i] = \sigma_u^2 I$ ,  $\sigma_u^2 > 0$  and finite  
(homoscedastic and no serial correlation)
- b)  $V[u_{it}|X_i, z_i, c_i] = \sigma_{u,it}^2 > 0$ , finite and  
 $Cov[u_{it}, u_{is}|X_i, z_i, c_i] = 0 \forall s \neq t$  (no serial correlation)
- c)  $V[u_i|X_i, z_i, c_i] = \Omega_{u,i}(X_i, z_i)$  is p.d. and finite

The remaining assumptions are divided into two sets of assumptions: the random effects model and the fixed effects model.

## 2.1 The Random Effects Model

In the random effects model, the individual-specific effect is a random variable that is uncorrelated with the explanatory variables.

*RE1: Unrelated effects*

$$E[c_i|X_i, z_i] = 0$$

*RE1* assumes that the individual-specific effect is a random variable that is uncorrelated with the explanatory variables of all past, current and future time periods of the same individual.

*RE2: Effect Variance*

- a)  $V[c_i|X_i, z_i] = \sigma_c^2 < \infty$  (homoscedastic)
- b)  $V[c_i|X_i, z_i] = \sigma_{c,i}^2(X_i, z_i) < \infty$  (heteroscedastic)

*RE2a* assumes constant variance of the individual specific effect.

*RE3: Identifiability*

a)  $\text{rank}(W) = K + M + 1 < NT$  and  $E[W_i'W_i] = Q_{WW}$  is p.d. and finite. The typical element  $w'_{it} = [1 \ x'_{it} \ z'_i]$ .

b)  $\text{rank}(W) = K + M + 1 < NT$  and  $E[W_i'\Omega_{v,i}^{-1}W_i] = Q_{WOW}$  is p.d. and finite.  $\Omega_{v,i}$  is defined below.

*RE3* assumes that the regressors including a constant are not perfectly collinear, that all regressors (but the constant) have non-zero variance and not too many extreme values.

The random effects model can be written as

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + v_{it}$$

where  $v_{it} = c_i + u_{it}$ . Assuming *PL2*, *PL4* and *RE1* in the special versions *PL4a* and *RE2a* leads to

$$\Omega_v = V[v|X, Z] = \begin{pmatrix} \Omega_{v,1} & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & & & \vdots \\ 0 & & \Omega_{v,i} & & 0 \\ \vdots & & & \ddots & \vdots \\ 0 & \cdots & 0 & \cdots & \Omega_{v,N} \end{pmatrix}_{NT \times NT}$$

with typical element

$$\Omega_{v,i} = V[v_i|X_i, z_i] = \begin{pmatrix} \sigma_v^2 & \sigma_c^2 & \cdots & \sigma_c^2 \\ \sigma_c^2 & \sigma_v^2 & \cdots & \sigma_c^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_c^2 & \sigma_c^2 & \cdots & \sigma_v^2 \end{pmatrix}_{T \times T}$$

where  $\sigma_v^2 = \sigma_c^2 + \sigma_u^2$ . This special case under *PL4a* and *RE2a* is therefore called the *equicorrelated random effects model*.

## 2.2 The Fixed Effects Model

In the fixed effects model, the individual-specific effect is a random variable that is allowed to be correlated with the explanatory variables.

*FE1: Related effects*

–

*FE1* explicitly states the absence of the unrelatedness assumption in *RE1*.

*FE2: Effect Variance*

–

*FE2* explicitly states the absence of the assumption in *RE2*.

*FE3: Identifiability*

$\text{rank}(\ddot{X}) = K < NT$  and  $E(\ddot{X}'_i \ddot{X}_i)$  is p.d. and finite

where the typical element  $\ddot{x}_{it} = x_{it} - \bar{x}_i$  and  $\bar{x}_i = 1/T \sum_t x_{it}$

*FE3* assumes that the time-varying explanatory variables are not perfectly collinear, that they have non-zero within-variance (i.e. variation over time for a given individual) and not too many extreme values. Hence,  $x_{it}$  cannot include a constant or any time-invariant variables. Note that only the parameters  $\beta$  but neither  $\alpha$  nor  $\gamma$  are identifiable in the fixed effects model.

### 3 Estimation with Pooled OLS

The *pooled OLS estimator* ignores the panel structure of the data and simply estimates  $\alpha$ ,  $\beta$  and  $\gamma$  as

$$\begin{pmatrix} \hat{\alpha}_{POLS} \\ \hat{\beta}_{POLS} \\ \hat{\gamma}_{POLS} \end{pmatrix} = (W'W)^{-1} W'y$$

where  $W = [\iota_{NT} \ X \ Z]$  and  $\iota_{NT}$  is a  $NT \times 1$  vector of ones.

*Random effects model:* The pooled OLS estimator of  $\alpha$ ,  $\beta$  and  $\gamma$  is unbiased under *PL1*, *PL2*, *PL3*, *RE1*, and *RE3a* in small samples. Additionally assuming *PL4* and normally distributed idiosyncratic and individual-specific errors, it is normally distributed in small samples. It is consistent and approximately normally distributed under *PL1*, *PL2*, *PL3*, *PL4*, *RE1*,

and *RE3a* in samples with a large number of individuals ( $N \rightarrow \infty$ ). However, the pooled OLS estimator is not efficient. More importantly, the usual standard errors of the pooled OLS estimator are incorrect and tests ( $t$ -,  $F$ -,  $z$ -, Wald-) based on them are not valid. Correct standard errors can be estimated with the so-called cluster-robust covariance estimator treating each individual as a cluster. Cluster-robust covariance matrix is consistent when the number of clusters  $N \rightarrow \infty$ . In practice we should have at least 50 clusters (see the handout on “Clustering in the Linear Model”).

*Fixed effects model:* The pooled OLS estimators of  $\alpha$ ,  $\beta$  and  $\gamma$  are biased and inconsistent, because the variable  $c_i$  is omitted and potentially correlated with the other regressors.

## 4 Random Effects Estimation

The *random effects estimator* is the feasible generalized least squares (GLS) estimator

$$\begin{pmatrix} \hat{\alpha}_{RE} \\ \hat{\beta}_{RE} \\ \hat{\gamma}_{RE} \end{pmatrix} = \left( W' \hat{\Omega}_v^{-1} W \right)^{-1} W' \hat{\Omega}_v^{-1} y.$$

where  $W = [\iota_{NT} \ X \ Z]$  and  $\iota_{NT}$  is a  $NT \times 1$  vector of ones.

The error covariance matrix  $\Omega_v$  is assumed block-diagonal with equicorrelated diagonal elements  $\Omega_{v,i}$  as in section 2.1 which depend on the two unknown parameters  $\sigma_v^2$  and  $\sigma_c^2$  only. There are many different ways to estimate these two parameters. For example,

$$\hat{\sigma}_v^2 = \frac{1}{NT} \sum_{t=1}^T \sum_{i=1}^N \hat{v}_{it}^2, \quad \hat{\sigma}_c^2 = \hat{\sigma}_v^2 - \hat{\sigma}_u^2$$

where

$$\hat{\sigma}_u^2 = \frac{1}{NT - N} \sum_{t=1}^T \sum_{i=1}^N (\hat{v}_{it} - \bar{\hat{v}}_i)^2$$

and  $\widehat{v}_{it} = y_{it} - \alpha_{POLS} - x'_{it}\widehat{\beta}_{POLS} - z'_i\widehat{\gamma}_{POLS}$  and  $\widehat{v}_i = 1/T \sum_{t=1}^T \widehat{v}_{it}$ . The degree of freedom correction in  $\widehat{\sigma}_u^2$  is also asymptotically important when  $N \rightarrow \infty$ .

*Random effects model:* We cannot establish small sample properties for the RE estimator. The RE estimator is consistent and asymptotically normally distributed under *PL1 - PL4*, *RE1*, *RE2* and *RE3b* when the number of individuals  $N \rightarrow \infty$  even if  $T$  is fixed. It can therefore be approximated in samples with many individual observations  $N$  as

$$\begin{pmatrix} \widehat{\alpha}_{RE} \\ \widehat{\beta}_{RE} \\ \widehat{\gamma}_{RE} \end{pmatrix} \overset{A}{\sim} N \left( \begin{pmatrix} \alpha \\ \beta \\ \gamma \end{pmatrix}, Avar \begin{bmatrix} \widehat{\alpha}_{RE} \\ \widehat{\beta}_{RE} \\ \widehat{\gamma}_{RE} \end{bmatrix} \right)$$

Assuming the equicorrelated model (*PL4a* and *RE2a*),  $\widehat{\sigma}_v^2$  and  $\widehat{\sigma}_c^2$  are consistent estimators of  $\sigma_v^2$  and  $\sigma_c^2$ , respectively. Then  $\widehat{\alpha}_{RE}$ ,  $\widehat{\beta}_{RE}$  and  $\widehat{\gamma}_{RE}$  are asymptotically efficient and the asymptotic variance can be consistently estimated as

$$\widehat{Avar} \begin{bmatrix} \widehat{\alpha}_{RE} \\ \widehat{\beta}_{RE} \\ \widehat{\gamma}_{RE} \end{bmatrix} = \left( W' \widehat{\Omega}_v^{-1} W \right)^{-1}$$

Allowing for arbitrary conditional variances and for serial correlation in  $\Omega_{v,i}$  (*PL4c* and *RE2b*), the asymptotic variance can be consistently estimated with the so-called cluster-robust covariance estimator treating each individual as a cluster (see the handout on “Clustering in the Linear Model”). In both cases, the usual tests ( $z$ -, Wald-) for large samples can be performed.

In practice, we can rarely be sure about equicorrelated errors and better always use cluster-robust standard errors for the RE estimator.

*Fixed effects model:* Under the assumptions of the fixed effects model (*FE1*, i.e. *RE1* violated), the random effects estimators of  $\alpha$ ,  $\beta$  and  $\gamma$  are biased and inconsistent, because the variable  $c_i$  is omitted and potentially correlated with the other regressors.

## 5 Fixed Effects Estimation

Subtracting time averages  $\bar{y}_i = 1/T \sum_t y_{it}$  from the initial model

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + c_i + u_{it}$$

yields the *within model*

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

where  $\ddot{y}_{it} = y_{it} - \bar{y}_i$ ,  $\ddot{x}_{itk} = x_{itk} - \bar{x}_{ik}$  and  $\ddot{u}_{it} = u_{it} - \bar{u}_i$ . Note that the individual-specific effect  $c_i$ , the intercept  $\alpha$  and the time-invariant regressors  $z_i$  cancel.

The *fixed effects estimator* or *within estimator* of the slope coefficient  $\beta$  estimates the within model by OLS

$$\hat{\beta}_{FE} = \left( \ddot{X}'\ddot{X} \right)^{-1} \ddot{X}'\ddot{y}$$

Note that the parameters  $\alpha$  and  $\gamma$  are not estimated by the within estimator.

*Random effects model and fixed effects model:* The fixed effects estimator of  $\beta$  is unbiased under *PL1*, *PL2*, *PL3*, and *FE3* in small samples. Additionally assuming *PL4* and normally distributed idiosyncratic errors, it is normally distributed in small samples. Assuming homoscedastic errors with no serial correlation (*PL4a*), the variance  $V\left[\hat{\beta}_{FE}|X\right]$  can be unbiasedly estimated as

$$\hat{V}\left[\hat{\beta}_{FE}|X\right] = \hat{\sigma}_u^2 \left( \ddot{X}'\ddot{X} \right)^{-1}$$

where  $\hat{\sigma}_u^2 = \hat{\ddot{u}}'\hat{\ddot{u}}/(NT - N - K)$  and  $\hat{\ddot{u}}_{it} = \ddot{y}_{it} - \ddot{x}'_{it}\hat{\beta}_{FE}$ . Note the non-usual degrees of freedom correction. The usual  $z$ - and  $F$ -tests can be performed.

The FE estimator is consistent and asymptotically normally distributed under *PL1 - PL4* and *FE3* when the number of individuals  $N \rightarrow \infty$  even if  $T$  is fixed. It can therefore be approximated in samples with many individual observations  $N$  as

$$\hat{\beta}_{FE} \overset{A}{\sim} N\left(\beta, Avar\left[\hat{\beta}_{FE}\right]\right)$$



Assuming homoscedastic errors with no serial correlation (*PL4a*), the asymptotic variance can be consistently estimated as

$$\widehat{Avar} [\widehat{\beta}_{FE}] = \widehat{\sigma}_u^2 (\ddot{X}' \ddot{X})^{-1}$$

where  $\widehat{\sigma}_u^2 = \widehat{u}' \widehat{u} / (NT - N)$ .

Allowing for heteroscedasticity and serial correlation of unknown form (*PL4c*), the asymptotic variance  $Avar[\widehat{\beta}_k]$  can be consistently estimated with the so-called cluster-robust covariance estimator treating each individual as a cluster (see the handout on “Clustering in the Linear Model”). In both cases, the usual tests ( $z$ -, Wald-) for large samples can be performed.

In practice, the idiosyncratic errors are often serially correlated (violating *PL4a*) when  $T > 2$ . Bertrand, Duflo and Mullainathan (2004) show that the usual standard errors of the fixed effects estimator are drastically understated in the presence of serial correlation. It is therefore advisable to always use cluster-robust standard errors for the fixed effects estimator.

## 6 Random Effects vs. Fixed Effects Estimation

The random effects model can be consistently estimated by both the RE estimator or the FE estimator. We would prefer the RE estimator if we can be sure that the individual-specific effect really is an unrelated effect (*RE1*). This is usually tested by a (Durbin-Wu-)Hausman test. However, the Hausman test is only valid under homoscedasticity and cannot include time fixed effects.

The unrelatedness assumption (*RE1*) is better tested by running an auxiliary regression (Wooldridge 2010, p. 332, eq. 10.88, Mundlak, 1978):

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + \bar{x}'_i\lambda + \delta_t + u_{it}$$

where  $\bar{x}_i = 1/T \sum_t x_{it}$  are the time averages of all time-varying regressors. Include time fixed  $\delta_t$  if they are included in the RE and FE estimation.

A joint Wald-test on  $H_0: \lambda = 0$  tests *RE1*. Use cluster-robust standard errors to allow for heteroscedasticity and serial correlation.

Note: Assumption *RE1* is an extremely strong assumption and the FE estimator is almost always much more convincing than the RE estimator. Not rejecting *RE1* does not mean accepting it. Interest in the effect of a time-invariant variable is no sufficient reason to use the RE estimator.

## 7 Least Squares Dummy Variables Estimator (LSDV)

The least squares dummy variables (LSDV) estimator is pooled OLS including a set of  $N - 1$  dummy variables which identify the individuals and hence an additional  $N - 1$  parameters. Note that one of the individual dummies is dropped because we include a constant. Time-invariant explanatory variables,  $z_i$ , are dropped because they are perfectly collinear with the individual dummy variables.

The LSDV estimator of  $\beta$  is numerically identical with the FE estimator and therefore consistent under the same assumptions. The LSDV estimators of the additional parameters for the individual-specific dummy variables, however, are inconsistent as the number of parameters goes to infinity as  $N \rightarrow \infty$ . This so-called *incidental parameters* problem generally biases all parameters in *non-linear* fixed effects models like the probit model.

## 8 First Difference Estimator

Subtracting the lagged value  $y_{i,t-1}$  from the initial model

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + c_i + u_{it}$$

yields the *first-difference model*

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

where  $\dot{y}_{it} = y_{it} - y_{i,t-1}$ ,  $\dot{x}_{it} = x_{it} - x_{i,t-1}$  and  $\dot{u}_{it} = u_{it} - u_{i,t-1}$ . Note that

the individual-specific effect  $c_i$ , the intercept  $\alpha$  and the time-invariant regressors  $z_i$  cancel. The *first-difference estimator* (FD) of the slope coefficient  $\beta$  estimates the first-difference model by OLS.

$$\hat{\beta}_{FD} = \left( \dot{X}'\dot{X} \right)^{-1} \dot{X}'\dot{y}$$

Note that the parameters  $\alpha$  and  $\gamma$  are not estimated by the FD estimator. In the special case  $T = 2$ , the FD estimator is numerically identical to the FE estimator.

*Random effects model and fixed effects model:* The FD estimator is a consistent estimator of  $\beta$  under the same assumptions as the FE estimator. It is less efficient than the FE estimator if  $u_{it}$  is not serially correlated (*PL4a*).

## 9 Fixed Effects vs. First Difference Estimation

Given the fixed effects model (*PL1, PL2, PL3, FE3*), both the fixed effects and the first difference estimator of  $\beta$  are consistent. Hence, the two estimators should be similar in large samples. In practice, however, the two estimator often differ substantially. The reason for this is typically a misspecification of the timing in the linear model. *PL1* assumes that changes in  $x_{it}$  have only an instantaneous effect on  $y_{it}$  at time  $t$ . In practice, effects often need several periods to materialize. Such patterns are called *dynamic treatment effects*. In this situation, the first difference estimator will only pick up the instantaneous effect at time  $t$  while the fixed effects estimator picks up an average of the dynamic treatment effects.

## 10 Time Fixed Effects

We often also suspect that there are time-specific effects  $\delta_t$  which affect all individuals in the same way

$$y_{it} = \alpha + x'_{it}\beta + z'_i\gamma + \delta_t + c_i + u_{it}.$$

---

We can estimate this extended model by including a dummy variable for  $T - 1$  time periods with one period serving as the reference period. Assuming a fixed number of time periods  $T$  and the number of individuals  $N \rightarrow \infty$ , both the RE estimator and the FE estimator are consistent using time dummy variables under above conditions.

## 11 Implementation in Stata 17

Stata provides a series of commands that are especially designed for panel data. See `help xt` for an overview.

Stata requires panel data in the so-called *long form*: there is one line for every individual and every time observation. The very powerful Stata command `reshape` helps transforming data into this format. Before working with panel data commands, we have to tell Stata the variables that identify the individual and the time period. For example, load data and define individuals (variable *idcode*) and time periods (variable *year*)

```
webuse nlswork.dta
xtset idcode year
```

Stata provides descriptive statistics for panel data with the commands

```
xtdescribe
xtsum
```

The pooled OLS estimator with corrected standard errors is calculated with the standard ols command `regress`:

```
generate ttl_exp2 = ttl_exp^2
reg ln_wage grade ttl_exp ttl_exp2, vce(cluster idcode)
```

where the `vce` option was used to report correct cluster-robust standard errors. This command multiplies  $\widehat{Avar}$  with  $(NT - 1)/(NT - M - K - 1) \cdot N/(N - 1)$  as a small sample correction and uses  $N - 1$  degrees of freedom for t- and F-tests.

The random effects estimator is calculated by the Stata command `xtreg` with the option `re`:

```
xtreg ln_wage grade ttl_exp ttl_exp2, re
```

Stata reports asymptotic *z*- and Wald-tests with random effects estimation. Cluster-robust standard errors are reported with:

```
xtreg ln_wage grade ttl_exp ttl_exp2, re vce(cluster idcode)
```

Since version 10, Stata always assumes clustering with robust standard errors in random and fixed effects estimations. So we could also just use

```
xtreg ln_wage grade ttl_exp ttl_exp2, re vce(robust)
```

The fixed effects estimator is calculated by the Stata command `xtreg` with the option `fe`:

```
xtreg ln_wage ttl_exp ttl_exp2, fe
```

Note that the effect of time-constant variables like `grade` is not identified by the fixed effects estimator. The parameter reported as `_cons` in the Stata output is the average fixed effect  $1/N \sum_i c_i$ . This command uses  $NT - N - K$  degrees of freedom for  $t$ - and  $F$ -tests. Cluster-robust standard errors are reported with the `vce` option:

```
xtreg ln_wage ttl_exp ttl_exp2, fe vce(cluster idcode)
```

This command multiplies  $\widehat{Avar}$  with  $(NT - 1)/(NT - N - K) \cdot N/(N - 1)$  as a small correction and reports reports cluster-robust  $t$ - and  $F$ -tests with  $N - 1$  degrees of freedom. The latter is particularly useful with large  $T$  (see Stock and Watson, 2008).

The Hausman test is calculated by

```
xtreg ln_wage grade ttl_exp ttl_exp2, re
estimates store b_re
xtreg ln_wage ttl_exp ttl_exp2, fe
estimates store b_fe
hausman b_fe b_re, sigmamore
```

and the auxiliary regression version by

```
regress ln_wage grade ttl_exp ttl_exp2
tegen ttl_exp_mean = mean(ttl_exp) if e(sample), by(idcode)
egen ttl_exp2_mean = mean(ttl_exp2) if e(sample), by(idcode)
regress ln_wage grade ttl_exp ttl_exp2 ///
      ttl_exp_mean ttl_exp2_mean, vce(cluster idcode)
test ttl_exp_mean ttl_exp2_mean
```

## 12 Implementation in R 4.2.2

The R package `plm` provides a series of functions and data structures that are especially designed for panel data.

The `plm` package works with data stored in a dataframe in the so-called *long form*. Long form data means that there is one line for every individual and every time observation. For example, load data

```
library(haven)
nlswork <- read_dta("https://www.stata-press.com/data/r17/nlswork.dta")
```

where individuals are defined by `idcode` and time periods by `year`.

Pooled OLS with cluster-robust standard errors can be estimated with a standard regression and the packages `lmtest` and `sandwich`

```
pols1 <- lm(ln_wage~grade+ttl_exp+I(ttl_exp^2), data = nlswork)
library(lmtest)
library(sandwich)
coefest(pols1, vcov = vcovCL, cluster = ~idcode)
```

This command multiplies  $\widehat{Avar}$  with  $(NT-1)/(NT-M-K-1) \cdot N/(N-1)$  as a small sample correction.

Alternatively, pooled OLS with corrected standard errors is estimated by the package `plm` with the function `plm` and its model option `pooling`:

```
library(plm)
pols2 <- plm(ln_wage~grade+ttl_exp+I(ttl_exp^2), model="pooling",
             data = nlswork, index=c("idcode", "year"))
summary(pols2)
coefest(pols2, vcov=vcovHC(pols2, cluster="group", type="HC1"))
```

where `coefest` reports cluster-robust standard errors. `cluster="group"` defines the clusters by the individual identifier set by the option `index` in `plm`, i.e. the variable `idcode` in the example. This command multiplies  $\widehat{Avar}$  with  $(NT-1)/(NT-M-K-1)$  but not with  $N/(N-1)$  and uses  $NT-M-K-1$  degrees of freedom for t- and F-tests.

The random effects estimator is calculated by `plm` option `random`:

```
re <- plm(ln_wage~grade+ttl_exp+I(ttl_exp^2), model="random",
          data = nlswork, index=c("idcode", "year"))
summary(re)
```

Cluster-robust standard errors are reported with

```
coeftest(re, vcov=vcovHC(re, cluster="group", type="HC1"))
```

The fixed effects estimator is calculated by plm option `within`

```
fe <- plm(ln_wage ~ grade + ttl_exp + I(ttl_exp^2), model="within",  
          data=nlswork, index=c("idcode", "year"))  
summary(fe)
```

Note that effects of time-constant variables like `grade` are not identified by the fixed effects estimator. This command uses  $NT - N - K$  degrees of freedom for t- and F-tests. Cluster-robust standard errors are given by:

```
coeftest(fe, vcov=vcovHC(fe, cluster="group", type="HC1"))
```

This command multiplies  $\widehat{Avar}$  with  $(NT - 1)/(NT - K - 1)$  as a small sample correction and uses  $NT - K - 1$  degrees of freedom for t- and F-tests.

The Hausman test is calculated by estimating RE and FE and then comparing the estimates:

```
phptest(fe, re)
```



## References

### Introductory textbooks

Stock, James H. and Mark W. Watson (2020), Introduction to Econometrics, 4th Global ed., Pearson. Chapter 10.

Wooldridge, Jeffrey M. (2009), Introductory Econometrics: A Modern Approach, 4th ed., South-Western Cengage Learning. Ch. 13 and 14.

### Advanced textbooks

Cameron, A. Colin and Pravin K. Trivedi (2005), Microeconometrics: Methods and Applications, Cambridge University Press. Chapter 21.

Wooldridge, Jeffrey M. (2010), Econometric Analysis of Cross Section and Panel Data, MIT Press. Chapter 10.

### Companion textbooks

Angrist, Joshua D. and Jörn-Steffen Pischke (2009), Mostly Harmless Econometrics: An Empiricist's Companion, Princeton University Press. Chapter 5.

### Articles

Manuel Arellano (1987), Computing Robust Standard Errors for Within-Group Estimators, Oxford Bulletin of Economics and Statistics, 49, 431-434.

Bertrand, M., E. Duflo and S. Mullainathan (2004), How Much Should We Trust Differences-in-Differences Estimates?, Quarterly Journal of Economics, 119(1), 249-275.

Mundlak, Y. (1978), On the pooling of time series and cross section data, Econometrica, 46, 69-85.

Stock, James H. and Mark W. Watson (2008), Heteroskedasticity-Robust Standard Errors for Fixed Effects Panel Data Regression, Econometrica, 76(1), 155-174. [advanced]