



## Panel Data

A panel, or longitudinal, data set is one where there are repeated observations on the same  $N$  units: individuals, households, firms, countries, or any set of entities that remain stable through time.

Repeated observations create a potentially very large panel data sets. With  $N$  units and  $T$  time periods, we have  $NT$  observations.

- Advantage: Large sample, great for estimation.
- Disadvantage: Dependence! Observations are, likely, not independent.

Modeling the potential dependence creates different models.







## Panel Data Models

With panel data we can study different issues:

- (1) *Cross sectional* variation (unobservable in time series data) vs *Time series* variation (unobservable in cross sectional data).
- (2) Heterogeneity (observable and unobservable heterogeneity).
- (3) Hierarchical structures (say zip codes, city, and state effects).
- (4) Dynamics in economic behavior.
- (5) Individual/group effects.
- (6) time effects.









$z_i$ : The variables responsible for unobserved heterogeneity (& dependence on the  $y_i$ 's). Usually, a nuisance component of the model.

The  $z_p$  variables are unobserved: Impossible to obtain information about each component in  $\sum_{p=1}^s z_{ip} \gamma_p$ . We define a term  $c_i$ , **the unobserved effect**, representing the joint impact of the  $z_p$  variables on  $y_i$  – like an index of unobservables for individual  $i$ :

$$C_i = \sum_{p=1}^S z_{ip} \gamma_p.$$

We can rewrite the regression model as:

$$y_{it} = \beta_1 + \sum_{j=2}^k x_{ij,t} \beta_j + c_i + \delta t + \varepsilon_{i,t}.$$

## Panel Data Models: Basic Model

Note: If the  $X_j$ 's are so comprehensive that they capture all relevant characteristics of individual  $i$ ,  $c_i$  can be dropped and, then, pooled OLS may be used. But this situation is very unlikely.

In general, dropping  $c_i$  leads to missing variables problem: omitted variable bias!

We usually think of  $c_i$  as **contemporaneously exogenous** to the conditional error, That is,  $E[\varepsilon_{it}|c_i] = 0$ ,  $t = 1, \dots, T$ .

A stronger assumption: **Strict exogeneity** can also be imposed. Then,  $E[\varepsilon_{it} | x_{i,1}, x_{i,2}, \dots, x_{i,T}, c_i] = 0$  for  $t = 1, \dots, T$ .

















# Pooled Model





## Pooled Model

In this context, OLS produces BLUE and consistent estimators. We refer to this as **pooled OLS estimation**.

Of course, if our assumption regarding the unobservable variables is wrong, we are in the presence of an omitted variable,  $c_i$ .

Then, we have potential bias and inconsistency of pooled OLS. The magnitude of these problems depends on how the true model behaves: “fixed” or “random”.











Disadvantage: We lose observations (and power!) since we have only  $N$

Remark: Under the usual assumptions, pooled OLS using the between







We have two periods: before and after the natural experiment (the treatment)

The number of groups,  $S$ , (treated or not treated) under consideration is





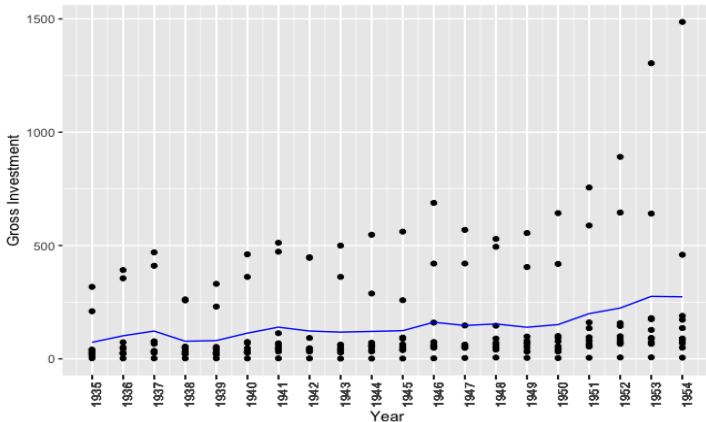


# Application





# Grunfeld Investment Study



**Figure 3:** Time heterogeneity. The blue line connects the mean values of invest for all firms across years.



## Grunfeld Investment Study

We achieve the same coefficient estimates by using function `plm()` from package `plm`. First, an index has to be supplied, corresponding to the entity and/or time dimension of the panel. The argument `model=` is set to ‘‘pooling’’.

```
pooled_ols_plm <- plm(inv ~ capital, data = Grunfeld,
index = c("firm", "year"),
effect = "individual", model = "pooling")
summary(pooled_ols_plm)
```

# Grunfeld Investment Study

## #Output

```
Call:plm(formula = inv ~ capital, data = Grunfeld, effect = "individual",
model = "pooling", index = c("firm", "year"))
```

Balanced Panel:  $n = 10$ ,  $T = 20$ ,  $N = 200$

Residuals:

Min.	1st Qu.	Median	3rd Qu.	Max.
-316.924	-96.450	-14.429	17.069	481.924

Coefficients:

Estimate	Std. Error	t-value	Pr(> t )
(Intercept)	14.236205	15.639266	0.9103
capital	0.477224	0.038339	12.4474

\*\*\*---Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

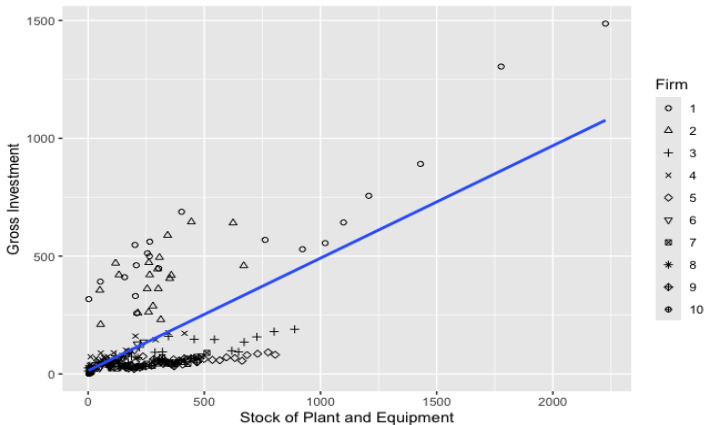
Total Sum of Squares: 9359900

Residual Sum of Squares: 5251000

R-Squared: 0.43899, Adj. R-Squared: 0.43616

F-statistic: 154.937 on 1 and 198 DF, p-value: < 2.22e-16





**Figure 4:** Scatter plot: Although firms could be distinguished by the variable firm, OLS estimation treats all observations as if they come from different entities and fits the regression line accordingly.

## Main Models: FEM and REM

There are two main approaches to fitting models using panel data:

- (1) Fixed effects regression.
- (2) Random effects regressions.

The key difference between these two approaches is how the unobservable characteristics—the **individual effects**—are modeled.

Fixed effects (FE): The individuals are fixed. The differences between them are not interest, only  $\beta$  is interesting. No intent on generalizing the results.

Random effects (RE): The individuals come from a random sample drawn from a larger population, and the variance between them is interesting and can be informative about the larger population.

# Fixed Effects Model

## Fixed Effects Model (FEM)

The fixed effects model:

$$(A1) \ y_{i,t} = x'_{i,t}\beta + c_i + \varepsilon_{i,t}, \quad (x'_{i,t} \text{ does not include an intercept}).$$

(A2)  $E[\varepsilon_{i,t}|X_{i,s}, c_{i,s}] = 0$ , for all  $t, s$ . ( $X_i$  and  $c_i$  are strictly exogenous).

The unobserved component,  $c_i$ , is arbitrarily correlated with  $x_{i,t}$ :  $E[c_i|X_i] = g(X_i) = \text{constant}_i$ , which implies  $\text{Cov}[x_{i,t}, c_i] \neq 0$ .

**Note 1:** Under the FEM, pooled OLS omits  $c_i$ . So, the estimators are biased and inconsistent.

We summarize (“**control for**”) these unobservable effects with  $\alpha_i$ , a constant. All time invariant characteristics of individual  $i$  (location, gender, nationality, etc.) are swept away under this formulation.

**Note 2:** In a FEM, individuals serve as their own controls.

## Estimation with Fixed Effects

Whatever effects the omitted variables have on the individual  $i$  at one time, they will also have the same effect at a later time, thus, their effects will be constant, or “**fixed**.”

For this, we need the omitted variables to have time-invariant values with time-invariant effects. Typical example, a CEO's IQ/gender. We expect this variable to have the same effect at  $t = 1$  or  $t = 10$ .

As we will see, FEM are estimated using the **within** transformation. Thus, if individuals do not change much (or at all) across time, a FEM may not work very well. We need within-individuals variability in the variables if we are to use individuals as their own controls.

## FEM: Estimation via LSDV

Model for individual  $i$  is:  $y_i = x_i' \beta + c_i + \varepsilon_i$ , where all notations are as before and  $c_i$  is a  $T_i \times 1$  vector. Recall, each individual has  $T_i$  observations.

Stacking over all individuals, the model can be expressed as,

$$y = X\beta + c + \varepsilon,$$

where all notations are as before and  $c, y$ , and  $\varepsilon$  are  $\sum_i T_i \times 1$  vectors.

The FEM can be represented with a set of dummy/indicator variables:

$$y_{i,t} = x_{i,t}' \beta + \sum_{j=1}^N c_j d_{ij,t} + \varepsilon_{i,t}, \quad \text{with } d_{ij,t} = 1, \text{ if } i = j.$$

## FEM: Estimation via LSDV

The FE model assumes  $c_j = \alpha_j$  (a constant, it does not vary with  $t$ ):

$$y_i = X_i\beta + d_i\alpha_i + \varepsilon_i, \quad \text{for each individual } i.$$

Staking the model, we have

$$\begin{aligned} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} &= \begin{pmatrix} X_1 & d_1 & 0 & \dots & 0 \\ X_2 & 0 & d_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ X_N & 0 & 0 & \dots & d_N \end{pmatrix} \begin{pmatrix} \beta \\ \alpha \end{pmatrix} + \varepsilon \\ &= [X \ D] \begin{pmatrix} \beta \\ \alpha \end{pmatrix} + \varepsilon = Z\delta + \varepsilon. \end{aligned}$$

The FEM is the CLM, but with many  $(k + N)$  independent variables. OLS estimators are unbiased, consistent, efficient, but impractical if  $N$  is large.

## FEM: LSDV Estimator

Avoid dummy variable trap: If a constant is present in the model, the number of dummy variable should be  $N - 1$ . The omitted individual or group becomes the reference category.

However, the choice of reference category is often arbitrary and thus, the interpretation of the  $\alpha_j$  will not be particularly interesting.

Alternatively, we can drop the intercept  $\beta_1$  and define dummy variables for all of the individuals. This is the more common approach, as done above. The  $\alpha_i$ 's now become the intercepts for each of the  $i$ 's.



## FEM: LSDV Estimator

If  $E[\varepsilon_{i,t}|x_{i,s}, c_i] \neq 0$ , then LSDV cannot be used. It is inconsistent. In this case, we need to use IV's. Or a good natural experiment.

Note: Least Squares (LS) is an estimator, not a model. Given the formulation with a lot of dummy variables, this particular LS estimator is called Least Squares Dummy Variable (LSDV) estimator.



$$\begin{aligned} M_D^i &= I_{T_i} - d_i(d_i' d_i)^{-1} d_i = I_{T_i} - (1/T_i) d_i d_i' \\ X' M_D X &= \sum_{i=1}^N \sum_{t=1}^{T_i} (x_{i,t} - \bar{x}_i)(x_{i,t} - \bar{x}_i)' \\ X' M_D y &= \sum_{i=1}^N \sum_{t=1}^{T_i} (x_{i,t} - \bar{x}_i)(y_{i,t} - \bar{y}_i)' \end{aligned}$$

51 / 100

## FEM: Within Estimation

The within-groups method estimates the parameters using demeaned data. Since we do not include intercept  $\beta_1$ , we have,

$$y_{i,t} - \bar{y}_i = \sum_{j=2}^k (x_{ij,t} - \bar{x}_{ij}) \beta_j + \delta^* + (\varepsilon_{i,t} - \bar{\varepsilon}).$$

Recall: It is called **within-groups/individuals** method because it relies on variations **within** individuals rather than **between** individuals.

For the usual asymptotic results, we need:

$$(A2) \ E[\Delta \varepsilon_{i,t} | X_i] = 0$$

$$(A3') \ E[\Delta \varepsilon_i' \varepsilon_i | X_i, c_i] = \Sigma \quad (\text{different formulations are ok})$$

$$(A4) \ E[\Delta X_i' \Delta X_i] \text{ has full rank.}$$

## FEM: Within Transformation Removes Effects

There are costs in the simplicity of the within-groups estimation:

- (1) All **time invariant** variables (including constant) for each individual  $i$  drop out of the model. This eliminates all between-individuals variability (which may be contaminated by omitted variable bias) and leaves only within-subject variability to analyze.
- (2) Dependent variables are likely to have smaller variances than in the original specification (measured as deviations from the  $i$  mean).
- (3) The manipulation involves the loss of  $N$  degrees of freedom (we are estimating  $N$  means).

## FEM: Within Estimation

$\hat{\beta}$  is obtained by **within-groups** least squares (group mean deviations).

Then we use the normal equations to estimate  $\hat{\alpha}_{N \times 1}$

$$\begin{aligned} D'X\hat{\beta} + D'D\hat{\alpha} &= D'y \\ \hat{\alpha}_{N \times 1} &= (D'D)^{-1}D'(y - X\hat{\beta}), \quad \text{or,} \\ \hat{\alpha}_i &= \frac{1}{T_i} \sum_{t=1}^{T_i} (y_{i,t} - x'_{i,t}\hat{\beta}) = \bar{\varepsilon}. \end{aligned}$$

Note: This is simple algebra—the estimator is just OLS. Note what  $\hat{\alpha}_i$  is when  $T_i = 1$ . Follow this with  $y_{i,t} - \hat{\alpha} - x'_{i,t}\hat{\beta} = 0$  if  $T_i = 1$ .



## FEM: Estimation – FE or FD?

### Fixed-effects (or Within) Estimator

- Each variable is demeaned – i.e., subtracted by its average
- Dummy variable regression – i.e., put in a dummy variable for each cross-sectional unit, along with other explanatory variables. This may cause estimation difficulty when  $N$  is large.

## FD Estimator

- Each variable is differenced once over time, so we are effectively estimating the relationship between change of variables.



## FEM: Estimation – FE or FD?

When  $N$  is large and  $T$  is small but greater than 2 (for  $T = 2$ , FE=FD)

- FE is more efficient when  $\varepsilon_{i,t}$  are serially uncorrelated while FD is more efficient when  $\varepsilon_{i,t}$  follows a random walk ( $\rho = 1$ ).

When  $T$  is large and  $N$  is small

- FD has advantage for processes with large positive auto-correlation. (If  $\rho$  is near 1, FD solves the nonstationary problem!)
- FE is more sensitive to nonnormality, heteroscedasticity, and serial correlation in  $\varepsilon_{i,t}$ .
- On the other hand, FE is less sensitive to violation of the strict exogeneity assumption. Then, FE is preferred when the processes are weakly dependent over time.

## FEM: Calculation of $Var[\hat{\beta}|X]$

Since we have assumed strict exogeneity:  $Cov[\varepsilon_{i,t}, (x_{j,t}, c_j)] = 0$ , we have OLS in the CLM. That is,

$$\text{Asy.Var}[\hat{\beta}|X] = (\sigma_\varepsilon^2 / \sum_{i=1}^N T_i) \text{plim} [\sigma_\varepsilon^2 / \sum_{i=1}^N T_i \sum_{i=1}^N X_i' M_D^i X_i]^{-1} \quad \text{where,}$$

$$\hat{\sigma}_\varepsilon^2 = \frac{(\sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \hat{\alpha}_i - x_{it}' \hat{\beta}))^2}{(\sum_{i=1}^N T_i - N - k)}.$$

PCSE Remark: All previous remarks apply to FEM.

We build the SE according to the type of data we have:

- If we do not suspect auto-correlated errors—not a strange situation—, we can rely on clustered White SE's.
- Otherwise, we use Driscoll and Kraay SE.

## FEM: Testing for Fixed Effects

Under  $H_0 : \alpha_i = \alpha$  for all  $i$  (i.e., no fixed-effects)

—Here, we test whether to pool or not to pool the data.

Different tests:

- F-test based on the LSDV dummy variable model: constant or zero coefficients for D. Test follows an  $F_{N-1, NT-N-k}$  distribution.
- F-test based on FEM (the unrestricted model) vs pooled model (the restricted model). Test follows an  $F_{N-1, NT-N-k}$  distribution.

An LR can also be done, usually, assuming normality. Test follows a  $\chi^2_{N-1}$  distribution.

# FEM: Hypothesis Testing

Based on estimated residuals of the fixed effects model.

(1) Estimate FEM:  $y_{i,t} = x'_{i,t}\beta + \alpha_i + \varepsilon_{i,t}$ , and store the residuals  $\hat{\varepsilon}_{FE,i,t}$ .

(2) Test as usual:

- For heteroscedasticity, use Breusch and Pagan (1980) test.
- For auto-correlation in AR(1) model, use Breusch and Godfrey (1981) test.

# Application

# Grunfeld Investment Study

## (1) Least Squares Dummy Variable Estimation

### (a) Fixed Effect Model using the `lm()` function.

With function `lm()` a FE model can be estimated by including dummy variables for all firms. This is the so called least squares dummy variable (LSDV) approach. Similarly to the pooled OLS model, I am regressing `inv` on `capital`. If there is a large number of individuals, the LSDV method is expensive from a computational point of view.

```
fe_model_lm <- lm(inv ~ capital + factor(firm), data = Grunfeld)
summary(fe_model_lm)
```

NOTE: One firm dummy variable is dropped to avoid the dummy variable trap.

Residual standard error: 63.57 on 189 degrees of freedom  
Multiple R-squared: 0.9184, Adjusted R-squared: 0.9141  
F-statistic: 212.7 on 10 and 189 DF, p-value: < 2.2e-16

# Grunfeld Investment Study

(b) Fixed Effect Model using the `lm()` function, excluding the intercept.

Next up, I calculate the same model but drop the constant (intercept) by adding `-1` to the formula, so that no coefficient (level) of firm is excluded. Note that this does not alter the coefficient estimate of capital!

```
fe_model_lm_nocons <- lm(inv ~ capital + factor(firm) -1, data = Grunfeld)
summary(fe_model_lm_nocons)
```



```
Call:lm(formula = inv ~ capital + factor(firm) - 1, data = Grunfeld)
```

```
Residuals:
```

```
Min      1Q      Median      3Q      Max
-190.715 -20.835   -0.459   21.383  293.687
```

```
Coefficients:
```

	Estimate	Std. Error	t value	Pr(> t )
capital	0.37075	0.01937	19.143	< 2e-16***
factor(firm)1	367.61297	18.96710	19.382	< 2e-16***
factor(firm)2	301.15762	15.31806	19.660	< 2e-16***
factor(firm)3	-46.06917	16.18939	-2.846	0.00492**
factor(firm)4	41.17196	14.40645	2.858	0.00474**
factor(firm)5	-118.66544	17.05605	-6.957	5.52e-11***
factor(firm)6	16.74738	14.35657	1.167	0.24487
factor(firm)7	-69.17024	15.46733	-4.472	1.33e-05***
factor(firm)8	11.14050	14.31023	0.778	0.43725
factor(firm)9	-68.55731	15.34015	-4.469	1.35e-05***
factor(firm)10	0.88169	14.21425	0.062	0.95061

```
---Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 63.57 on 189 degrees of
freedomMultiple R-squared:  0.9439, Adjusted R-squared:  0.9407
F-statistic: 289.3 on 11 and 189 DF,  p-value: < 2.2e-16
```

# Grunfeld Investment Study

## Scatter Plot

Due to the introduction of firm dummy variables each firm has its own intercept with the y axis! For comparison, I plot the fitted values from the pooled OLS model (blue dashed line). Its slope is more steep compared to the LSDV approach as influential observations of firm 1 lead to an upward bias.

```
ggplot(data = broom::augment(fe_model_lm),
  aes(x = capital, y = .fitted)) +
  geom_point(aes(color = 'factor(firm)')) +
  geom_line(aes(color = 'factor(firm)')) +
  geom_line(data=broom::augment(pooled_ols_lm),
  aes(x = capital, y =.fitted),
  color = "blue", lty="dashed", linewidth = 1) +
  labs(x = "Stock of Plant and Equipment", y = "Fitted Values (inv ~ capital)",
  color = "Firm")
```

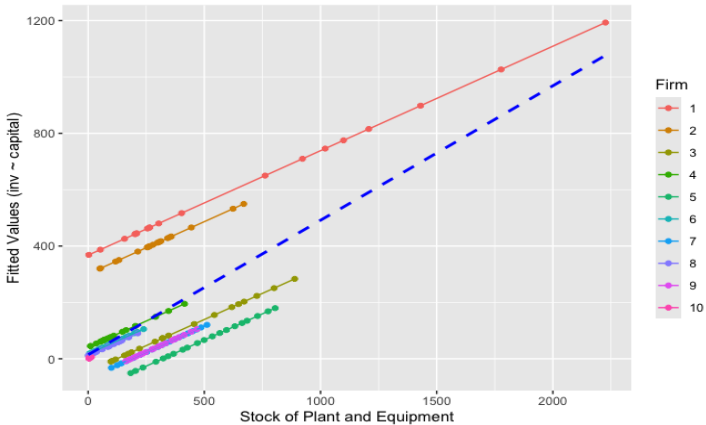


Figure 5: Fitted values from fixed-effects and pooled OLS models.

# Grunfeld Investment Study

## (2) Within Groups Estimator

### (a) Fixed Effects Model using `plm()` function.

The same coefficient estimates as with the LSDV approach can be computed with function `plm()`. The argument `model=` is now set to "within". This is the within estimator with  $n$  entity-specific intercepts.

```
fe_model_plm <- plm(inv ~ capital, data = Grunfeld,
  index = c("firm", "year"),
  effect = "individual", model = "within")
summary(fe_model_plm)
```

The coefficient of capital indicates how much `inv` changes over time, on average per firm, when capital increases by one unit.

# Grunfeld Investment Study

```
Call:plm(formula = inv ~ capital, data = Grunfeld, effect = "individual",
model = "within", index = c("firm", "year"))
```

Balanced Panel: n = 10, T = 20, N = 200

Residuals:

Min.	1st Qu.	Median	3rd Qu.	Max.
-190.71466	-20.83474	-0.45862	21.38262	293.68714

Coefficients:

	Estimate	Std. Error	t-value	Pr(> t )
capital	0.370750	0.019368	19.143	< 2.2e-16***

Total Sum of Squares: 2244400, Residual Sum of Squares: 763680

R-Squared: 0.65973, Adj. R-Squared: 0.64173

F-statistic: 366.446 on 1 and 189 DF, p-value: < 2.22e-16

# Grunfeld Investment Study

```
fixef(fe_model_plm)
```

With function `fixef()` the fixed effects, i.e. the constants for each firm can be extracted. Compare them with the coefficients of the LSDV approach (w/o the constant)- they must be identical.

1	2	3	4	5	6	7
367.61297	301.15762	-46.06917	41.17196	-118.66544	16.74738	-69.17024
8	9	10				
11.14050	-68.55731	0.88169				

## Testing for Fixed-Effects

With the function `pFtest()` one can test for fixed effects with the null hypothesis that pooled OLS is better than FE.

```
pFtest(fe_model_plm, pooled_ols_plm)
```

```
# Output
```

```
F test for individual effects
```

```
data: inv ~ capital
```

```
F = 123.39, df1 = 9, df2 = 189, p-value < 2.2e-16
```

```
alternative hypothesis: significant effects
```

Alternatively, this test can be carried out by jointly assessing the significance of dummy variables in the LSDV approach. The results are identical.

## Testing for Fixed-Effects

Joint significane test with LSDV approach

```
car::linearHypothesis(fe_model_lm, hypothesis.matrix =  
matchCoefs(fe_model_lm, "firm"))
```

# Output

Hypothesis:

```
factor(firm)2 = 0, factor(firm)3 = 0, factor(firm)4 = 0,  
factor(firm)5 = 0, factor(firm)6 = 0, factor(firm)7 = 0  
factor(firm)8 = 0, factor(firm)9 = 0, factor(firm)10 = 0
```

Model 1: restricted model

Model 2:  $\text{inv} \sim \text{capital} + \text{factor}(\text{firm})$

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	198 5250996				
2	189 763680	9	4487316	123.39	< 2.2e-16***

---Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

In both cases the null hypothesis is rejected in favor of the alternative that there are significant fixed effects.



# Grunfeld Investment Study

## (3) First Difference Estimator

There is another way of estimating a FE model by specifying `model = "fd"` in the function `plm()`.

```
fe_model_fd<- plm(inv ~ capital -1, data = Grunfeld,
index = c("firm", "year"),
effect = "individual", model = "fd")
summary(fe_model_fd)
```

The coefficient of capital is now different compared to the LSDV approach and within-groups estimator. This is because the coefficients and standard errors of the first-difference model are only identical to the previously obtained results when there are two time periods. For longer time series, both the coefficients and standard errors will be different.

## Oneway (individual) effect First-Difference Model

```
Call:plm(formula = inv ~ capital - 1, data = Grunfeld, effect =
"individual",      model = "fd", index = c("firm", "year"))
```

Balanced Panel: n = 10, T = 20, N = 200

Observations used in estimation: 190

### Residuals:

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-240.4	-11.7	0.1	3.5	12.6	333.2

### Coefficients:

	Estimate	Std. Error	t-value	Pr(> t )
capital	0.230780	0.059639	3.8696	0.00015***

---Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Total Sum of Squares: 584410

Residual Sum of Squares: 561210

R-Squared: 0.04476, Adj. R-Squared: 0.04476

F-statistic: 14.9739 on 1 and 189 DF, p-value: 0.00014998

## With Two Periods

Let's verify that the results from LSDV and within approaches are identical to the FD approach when there are two time periods only. Let's drop all years except 1935 and 1936 from the Grunfeld dataset and estimate the model again.

```
# Within estimation (two periods)
fe_model_plm_check <- plm(inv ~ capital, data =Grunfeld,
subset = year %in% c(1935, 1936),
index = c("firm", "year"), effect = "individual", model = "within")
lmtest::coeftest(fe_model_plm_check)
```

# Output

t test of coefficients:

	Estimate	Std. Error	t value	Pr(> t )
capital	0.91353	0.85333	1.0705	0.3122

## With Two Periods

Let's verify that the results from LSDV and within approaches are identical to the FD approach when there are two time periods only. Let's drop all years except 1935 and 1936 from the Grunfeld dataset and estimate the model again.

```
# FD estimation (two periods)
fe_model_fd_check<- plm(inv ~ capital -1, data = Grunfeld,
subset = year %in% c(1935, 1936),
index = c("firm", "year"), effect = "individual", model = "fd")
lmtest::coeftest(fe_model_fd_check)
```

## # Output

t test of coefficients:

	Estimate	Std. Error	t value	Pr(> t )
capital	0.91353	0.85333	1.0705	0.3122

# Random Effects Model

## The Random Effects Model (REM)

DGP:  $y_{i,t} = x'_{i,t}\beta + z'_i\gamma + \varepsilon_{i,t}$ , – obs. for individual  $i$  at time  $t$ .

When the observed characteristics are constant for each  $i$  (e.g., gender), a FEM is not an effective tool because such variables cannot be included.

An alternative approach, known as random effects model (REM) that, subject to two conditions, provides a solution to this problem.

(1) It is possible to treat each of the unobserved  $z_p$  variables as being drawn randomly from a given distribution.

The  $z_p$  variables are distributed independently of all the  $X_i$  variables, i.e.,  $E[z'_i X_j] = 0$ .

## (1) Randomly drawn unobserved $Z_p$ variables

The  $c_i = z_i' \gamma$  may be treated as a RV (thus, the name of this approach) drawn from a given distribution. Let,  $u_i = c_i - \alpha$ , where  $E[c_i|X] = \alpha$ . Then,

$$\begin{aligned} y_{i,t} &= \beta_1 (= \alpha) + \sum_{j=2}^k x_{ij,t} \beta_j + u_i + \delta t + \varepsilon_{i,t} \\ &= \beta_1 (= \alpha) + \sum_{j=2}^k x_{ij,t} \beta_j + \delta t + w_{i,t}, \quad \text{where } w_{i,t} = u_i + \varepsilon_{i,t}. \end{aligned}$$

We deal with the unobserved effect by subsuming it into a compound disturbance term  $w_{i,t}$ . We assume that  $u_i \sim D(0, \sigma_u^2)$ . Then,

$$E[w_{i,t}] = E[u_i] + E[\varepsilon_{i,t}] = 0.$$

The zero mean assumption ( $E[u_i] = 0$ ) is not crucial, since any nonzero component will be absorbed by the intercept  $\beta_1$ .

## The Random Effects Model (REM)

**(2)  $z_p$  is independent of all  $x_j$  variables**

Otherwise,  $u_i$  (&  $w_{i,t}$ ) will not be uncorrelated with  $x_j$ . The RE estimation will be biased and inconsistent.

Note: We would have to use the FEM, even if the first condition seems to be satisfied.

If conditions (1) and (2) are satisfied, we can use the REM, and OLS will work, but there is a complication:  $w_{i,t}$  is heteroscedastic.



## REM: *aka* Error Components Model

Model:  $y_{i,t} = x'_{i,t}\beta + u_i + \varepsilon_{i,t} = x'_{i,t}\beta + w_{i,t}$  ( $x'_{i,t}$  incl. a col. of 1's).

REM Assumptions:

$$E[\varepsilon_{i,t}|X_i] = 0, \quad E[\varepsilon_{i,t}^2|X_i] = \sigma_\varepsilon^2$$

$$E[u_i|X_i] = 0, \quad E[u_i^2|X_i] = \sigma_u^2$$

$$E[u_{j \in j,t} | X_j] = 0, \quad (u \text{ \& } \varepsilon \text{ are independent})$$

$$E[u_i u_j | X_i] = 0, \quad (i \neq j, \text{no cross-correlation of RE})$$

$$E[\varepsilon_{i,t}\varepsilon_{j,t}|X_i] = 0, \quad (i \neq j, \text{no cross-correlation of errors, } \varepsilon_{i,t})$$

$$E[\varepsilon_{i,t}\varepsilon_{i,s}|X_i] = 0, \quad (t \neq s, \text{there is no auto-correlation for } \varepsilon_{i,t})$$

$$\sigma_{w_{it}}^2 = \sigma_{u_i + \varepsilon_{it}}^2 = \sigma_{u_i}^2 + \sigma_{\varepsilon_{it}}^2 + 2\sigma_{u_i, \varepsilon_{it}} = \sigma_u^2 + \sigma_\varepsilon^2$$

$$\sigma_{w_{it1}, w_{it2}} = \sigma_{(u_j + \varepsilon_{jt1}), (u_j + \varepsilon_{jt2})} = \sigma_u^2.$$

## REM: Notation (Greene)

$$\begin{aligned} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} &= \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_N \end{pmatrix} \beta + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{pmatrix} + \begin{pmatrix} u_{1t_1} \\ u_{2t_2} \\ \vdots \\ u_{Nt_N} \end{pmatrix} \quad \begin{matrix} T_1 \text{ obs} \\ T_2 \text{ obs} \\ \vdots \\ T_N \text{ obs} \end{matrix} \\ &= X\beta + \varepsilon + u, \quad \sum_{i=1}^N T_i \text{ observations} \\ &= X\beta + w. \end{aligned}$$

In all that follows, except where explicitly noted,  $X$ ,  $X_i$ , and  $x'_{it}$  contain a constant term as the first element. To avoid notational clutter, in those cases,  $x'_{it}$  etc. will simply denote the counterpart without the constant term. Use of the symbol  $k$  for the number of variables will thus be context specific but will usually include the constant term.



9. *Journal of the American Medical Association*, 2000; 283: 2686-2692.

1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100
1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32	33	34	35	36	37	38	39	40	41	42	43	44	45	46	47	48	49	50	51	52	53	54	55	56	57	58	59	60	61	62	63	64	65	66	67	68	69	70	71	72	73	74	75	76	77	78	79	80	81	82	83	84	85	86	87	88	89	90	91	92	93	94	95	96	97	98	99	100



## REM: FGLS - Estimators for the Variances

Naturally,  $\sigma_\varepsilon^2$  and  $\sigma_u^2$  are unknown, which calls for FGLS to estimate the RE model. We first need to estimate the variance components. Assuming a balanced panel, the usual steps are:

- (1) Start with a consistent estimator of  $\beta$ . For example, pooled OLS,  $\hat{\beta}$ .
- (2) Compute  $\sum_{i=1}^N \sum_{t=1}^T (y_{i,t} - x'_{i,t} \hat{\beta})^2$  - estimates  $\sum_{i=1}^N \sum_{t=1}^T (\sigma_\varepsilon^2 + \sigma_u^2)$ .
- (3) Divide by a function of  $NT$ . For example:  $NT - k - 1$ 
  - We estimate  $\sigma^2$  as,  $\hat{\sigma}_{pooled}^2 = (\hat{\varepsilon}'_{pooled} \hat{\varepsilon}_{pooled}) / (NT - k - 1)$ .
  - We will use  $\hat{\sigma}_{pooled}^2$  to estimate the sum:  $\sigma_\varepsilon^2 + \sigma_u^2$
- (4) Use LSDV estimation to get  $\hat{\alpha}_i$  and  $\hat{\beta}_{LSDV}$ . Keep residuals  $\hat{\varepsilon}_{FE,i,t}$ .
- (5) Compute  $\sum_i \sum_t (y_{i,t} - \hat{\alpha}_i - x'_{i,t} \hat{\beta}_{LSDV})^2$ , - estimates  $\sum_{i=1}^N \sum_{t=1}^T \sigma_\varepsilon^2$ .
- (6) To estimate  $\sigma_\varepsilon^2$ , divide by  $NT - k - N$ :  $\hat{\sigma}_\varepsilon^2 = \sum_{i=1}^N \sum_{t=1}^T (\hat{\varepsilon}_{FE,i,t})^2 / (NT - k - N)$ .
- (7) Estimate  $\sigma_u^2$  as  $\hat{\sigma}_u^2 = \hat{\sigma}_{pooled}^2 - \hat{\sigma}_\varepsilon^2$ .



## REM: Practical Problems with FGLS

All of the preceding estimators regularly produce negative estimates of  $\sigma_u^2$ . Estimation is even more complicated in unbalanced panels.

A bulletproof solution (originally used in TSP, now LIMDEP and others):

(1) From robust LSDV estimator:  $\hat{\sigma}_\varepsilon^2 = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \hat{\alpha}_i - x'_{it} \hat{\beta}_{LSDV})^2}{\sum_{i=1}^N T_i}$ .

(2) From pooled OLS estimator:  $\hat{\sigma}_\varepsilon^2 + \hat{\sigma}_u^2 = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \hat{\alpha}_{OLS} - x'_{it} \hat{\beta}_{OLS})^2}{\sum_{i=1}^N T_i} \geq \hat{\sigma}_\varepsilon^2$ .

$$\hat{\sigma}_u^2 = \frac{\sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \hat{\alpha}_{OLS} - x'_{it} \hat{\beta}_{OLS})^2 - \sum_{i=1}^N \sum_{t=1}^{T_i} (y_{it} - \hat{\alpha}_i - x'_{it} \hat{\beta}_{LSDV})^2}{\sum_{i=1}^N T_i}.$$

Bullet proof solution: Do not correct by degrees of freedom. Then, given that the unrestricted RSS (LSDV) will be lower than the restricted (pooled OLS) RSS,  $\hat{\sigma}_{\mu}^2$  will be positive.



## Testing for Random Effects: LM Test

The most common test for the RE model vs a pooled OLS model is an LM test based on Breusch & Pagan (1980).  $H_0 : \sigma_u^2 = 0$  (i.e., no common effects). As with all LM test, it is based on the constrained model (here, pooled OLS model). Let  $\hat{\varepsilon}$  be pooled OLS residuals.

Breusch and Pagan Lagrange Multiplier statistic: Assuming normality (and for convenience now, a balanced panel)

$$LM = \frac{NT}{2(T-1)} \left[ \frac{\sum_{i=1}^N T \bar{\hat{\varepsilon}}^2}{\sum_{i=1}^n \hat{\varepsilon}_{it}^2} - 1 \right]^2 = \frac{NT}{2(T-1)} \left[ \frac{\sum_{i=1}^N T \bar{\hat{\varepsilon}}^2 - \sum_{i=1}^n \hat{\varepsilon}_{it}^2}{\sum_{i=1}^n \hat{\varepsilon}_{it}^2} \right]^2 \stackrel{H_0}{\sim} \chi_1^2.$$

For unbalanced panels, the scale in front becomes  $\frac{(\sum_{i=1}^N T_i)^2}{2 \sum_{i=1}^N T_i(T_i-1)}$ .



**Case for RE:**

- Under no omitted variables –or if the omitted variables are uncorrelated with  $x_{i,t}$  in the model– then a REM is probably best: It produces unbiased and efficient estimates, & uses all the data available.
- RE can deal with observed characteristics that remain constant for each individual. In FE, they have to be dropped from model.
- In contrast with FE, RE estimates a small number of parameters – We do not lose  $N$  degrees of freedom.
- Philosophically speaking, a REM is more attractive: Why should we assume one set of unobservables fixed and the other random?

## FE vs RE

### Case against RE:

- If either of the conditions for using RE is violated, we should use FE.

Condition (1): Randomly drawn unobserved  $z_p$  variables. This is a reasonable assumption in many cases: Many of the panels are designed to be a random sample (for example, NLSY).

But, it would not be a reasonable assumption if the units of observation in the panel data set were data from the S&P 500 firms.

Condition (2):  $z_p$  is independent of all of the  $x_j$  variables.

A violation of condition (2) causes inconsistency in the RE estimation.

## FE vs RE

FE estimation is always consistent. On the other hand, a violation of condition (2) causes inconsistency in the RE estimation.

In summary, the FE model always produce consistent estimates (under the null and alternative hypotheses), while the RE model is more efficient under the null hypothesis, but inconsistent under the alternative hypothesis.

To decide between FE and RE, we utilize the Hausman test.

*Journal of Management Education* 36(8) 907-924

# Application





```
re_model_plm <- plm(inv ~ capital, data = Grunfeld,
index = c("firm", "year"), effect = "individual", model = "random")
summary(re_model_plm)
```

## Oneway (individual) effect Random Effect Model

```
Call:plm(formula = inv ~ capital, data = Grunfeld, effect = "individual",
model = "random", index = c("firm", "year"))
```

Balanced Panel:  $n = 10$ ,  $T = 20$ ,  $N = 200$

Effects:

	var	std.dev	share
idiosyncratic	4040.63	63.57	0.135
individual	25949.52	161.09	0.865



# Grunfeld Investment Study

A decision between a fixed and random effects model can be made with the Hausman test, which checks whether the individual error terms are correlated with the regressors. The null hypothesis states that there is no such correlation (RE). The alternative hypothesis is that a correlation exists (FE). The test is implemented in function `phptest()`

```
phptest(fe_model_plm, re_model_plm)
```

## # Output

## Hausman Test

```
data:  inv ~ capitalchisq = 0.93423, df = 1, p-value = 0.3338
```

alternative hypothesis: one model is inconsistent

The null hypothesis cannot be rejected here, hence we should use a RE model.

# Thank you!