

# W271 Live Session 12: Analysis of Panel Data 2

Devesh Tiwari

8/1/2017

## Main topics covered in Week 12 (Async Unit 12)

- Fixed Effect Model
- A Digression: differencing when there are more than 2 time periods
- Random effect model
- Fixed effect vs. random effect models

## Readings:

**W2016:** Jeffrey Wooldridge. *Introductory Econometrics: A Modern Approach*. 6th edition. Cengage Learning

- Ch. 14.1 - 14.2
- [package plm] (<https://cran.r-project.org/web/packages/plm/plm.pdf>)
- [plm vignettes] (<https://cran.r-project.org/web/packages/plm/vignettes/plm.pdf>)

## Breakout Session 1

Imagine you are on a data science team working for a company that is interested in expanding its market in a developing country, Tiwaristan. The company's main source of revenue is its mobile application, which is why they are keen on understanding why users are happy or unhappy. By some miracle, they were able to conduct 4 rounds of surveys (once every month). They have survey data on 30 users from 6 different provinces. They have collected the following data:

- (a) Customer satisfaction score (0 - 100)
- (b) Number of hours the phone was not able to connect to the data network due to bandwidth related issues (connectivity).
- (c) Gender
- (d) Annual income
- (e) Education level
- (f) Province
- (g) Date / round of survey

1. Describe this dataset. How many observations are there in this dataset? What is the unit of analysis?
2. What estimation challenges would you face if you examined the relationship between customer satisfaction and connectivity? How would you overcome them?
3. Assume that you no longer have access to which province a user lives in. Can you use the same strategy you used to estimate #2 above to estimate the relationship with the other covariates and the DV? What are the benefits of using this approach? Costs?

## Breakout Session 2

In this exercise, we will analyze a simple panel dataset from the World Bank. This dataset has country level data from the years 2005 and 2010. In particular, it has data on a country's GDP per capita, the amount of money received from remittances abroad, and population. We are interested in estimating the relationship between GDP per capita, remittances, and population. BE SURE TO CONDUCT AN EDA AND MAKE ANY NECESSARY TRANSFORMATIONS TO THE DATA!!

1. Load this dataset as a *plm* package.
2. Estimate the pooling estimator using the *plm* function in R, and then examine the residuals. What do you notice?
3. Estimate the same data using a “within” or “fixed effect” estimator. Why do you think that this estimator is called a “within” estimator? What are we measuring? Examine these residuals as well.
4. Using the *lm* function in R, estimate a model that contains a dummy variable for each country. Compare these results to the fixed effects model you just estimated. What do you notice?
5. Compare your within estimate to the pooled-OLS estimate. Do they seem different? If so, what does that tell you?

## Group Discussion: Random effects

- Fixed effects estimator is strict in the sense that we can only reliably estimate time-variant variables that have a decent amount of within-unit variation.
- We might want to estimate a time-variant variable using more information than the fixed effects model, or we might want to estimate time-invariant variables.
- Doing so means that we have to assume that our variables of interest are exogenous with respect to the error term. If they are not exogenous, we have to interpret the results with a lot of caution because the coefficients will be biased.
- So, even though we are “OK” with generating biased estimates, we still should not use the pooled-OLS estimator because the errors within cross-sectional units are likely correlated.
- The random effects model helps us deal with this by “removing” this correlation from the observations.

## Breakout Session 3: Random effects model

1. Run a random effects model using the *plm* package. What do the “effects” portion of the output refer to? What is theta? Can you describe theta in plain English?
2. Examine and compare the regression outcome of the pooled-OLS, within estimator, and the random effects model. Pay attention to both the coefficients and standard errors. What do you notice? Which model produces the smallest standard errors? Largest? If your coefficients are different across models, what does that tell you about the relationship between the covariate and the error?
3. Run a Hausman test to determine whether or not a fixed effects model or a random effects model is appropriate.

## Open questions

Remainder of time