

Applied Econometrics Homework

M2 FE

Khalil Janbek, Romain Jouhameau

01/01/2022

```
library(data.table)
library(dtplyr)
library(dplyr, warn.conflicts = FALSE)
library(readxl)
library(janitor)
library(tidyverse)
library(sf)
library(usmap)
library(tsibble)
library(DBI)

con <- dbConnect(RSQLite::SQLite(), "Data/DB.sqlite")
DB = tbl(con, "DB")

DB %>%
  head() %>%
  collect()

#> # A tibble: 6 x 43
#>   Loan_Seq_Number Reporting_Period Current_UPB Delinquency_Status Loan_Age
#>   <chr>             <dbl>         <dbl> <chr>                <int>
#> 1 F20Q10000703      18293         342000 0                      0
#> 2 F20Q10000703      18322         342000 0                      1
#> 3 F20Q10000703      18353         341000 0                      2
#> 4 F20Q10000703      18383         341000 0                      3
#> 5 F20Q10000703      18414         340000 0                      4
#> 6 F20Q10000703      18444         340000 0                      5
#> # ... with 38 more variables: Time_to_Maturity <int>, Zero_Balance_Code <int>,
#> #   Current_Interest_Rate <dbl>, E_LoanToValue <int>,
#> #   Delinquency_Due_To_Disaster <chr>, Borrower_Assistance_Status_Code <chr>,
#> #   Interest_Bearing_UPB <dbl>, Credit_score <int>, First_Payment_Date <int>,
#> #   FstTime_HB_Flag <chr>, Maturity_Date <int>, MSA <int>,
#> #   Mortgage_Insurance_pct <int>, Number_of_Units <int>,
#> #   Occupancy_Status <chr>, O_CombinedLoanToValue <int>, ...
```

PART A:

1. In your data set, which are the variables which are varying with respect to two indices (or more if you consider inflows and outflows from one individual or country to another individual or countries? Which are the variables which are varying only with respect to time? Which are the variables which are varying only with respect to individuals?

```
DB %>%
  count(Loan_Seq_Number) %>%
  summarise(Number_Loans = n(),
            Max_T = max(n, na.rm = TRUE),
            Min_T = min(n, na.rm = TRUE),
            Mean_T = mean(n, na.rm = TRUE),
            Median_T = median(n, na.rm = TRUE),
            ) %>%
  collect()
```

2. What is the largest number of period T for individuals? What is the number of individuals?

```
#> # A tibble: 1 x 5
#>   Number_Loans Max_T Min_T Mean_T Median_T
#>   <int> <int> <int> <dbl> <int>
#> 1      346724    17     1   14.0     15
```

```
DB %>%
  mutate(Type = case_when(
    Lender_Type == 1 ~ 'Bank',
    Lender_Type == 0 & Fintech == 0 ~ 'Shadow Bank',
    Lender_Type == 0 & Fintech == 1 ~ 'Fintech'
  )) %>%
  group_by(Type) %>%
  count(Loan_Seq_Number) %>%
  summarise(Number_Loans = n(),
            Max_T = max(n, na.rm = TRUE),
            Min_T = min(n, na.rm = TRUE),
            Mean_T = mean(n, na.rm = TRUE),
            Median_T = median(n, na.rm = TRUE)) %>%
  collect()
```

```
#> # A tibble: 3 x 6
#>   Type      Number_Loans Max_T Min_T Mean_T Median_T
#>   <chr>          <int> <int> <int> <dbl> <int>
#> 1 Bank           147506    17     1   14.3     15
#> 2 Fintech         52273    17     1   13.7     15
#> 3 Shadow Bank    146945    17     1   13.9     15
```

3. Comment on the structure of the unbalanced panel (how many (and which) countries have a single observation, discontinuities between observations, how many have at least 2 consecutive observations (which is useful to compute lags, autocorrelations, first difference and within estimators)?

4. Compute between transformed and within transformed variables for all variables. Present a table with the within, between and pooled variance for each variable. Compute the share of between and within variance in the total variance for each variable. Comment these results.

5. Plot the distribution of the within and between transformed dependent variable and of you key (preferred) explanatory variable (not all the explanatory variable) [in Burnside and Dollar: GDP growth and foreign aid EDA/GDP], using on the same graph an histogram, a normal law with same empirical mean and standard error and a kernel continuous approximation. Comment the between and within difference for each variable, and compare within/within for dependent and explanatory variable, and between/between for dependent and explanatory variable: kurtosis, skewness, non-normality, high leverage observation (far from the mean), several modes (mixture of distribution)?

6. Plot boxplot of within distribution and between distribution for the dependent variable and the key explanatory variables. Comment that you find the same insights from question 5.
7. Compute univariate descriptive statistics (min, Q1, median, Q3, max, mean, standard error) for Within and Between transformed variables. Is the mean different from the median and why? How many standard errors from the mean are the min and max extremes (report (MAX-average)/standard error and (MIN-average)/standard error in the tables)?
8. Plot the boxplot of within transformed dependent variable and the key explanatory variable by a few individual (all of them if N around 50) and only the first 20 of them for larger data set. Comment on their differences of standard errors and means for each individuals
9. Compare and comment the within and between transformed bivariate correlation matrix for all variables (include a time trend 1,2,,T). Check poor simple correlation with the dependent variables and high correlation between explanatory variables.
10. Comment the bivariate auto-correlation and trend-correlations (check the number of observations).
11. Comment the bivariate graphs with linear, quadratic and Lowess fit for dependent and key explanatory variable (aid/gdp and growth of gdp): Within transformed, Between transformed.
12. Comment the results of estimations of Between, Within (fixed effects, (fe)) and Mundlak (random effects (re) including all X(i.) as regressors), two-way fixed effects (add year dummies in fe regression) and First differences, including all explanatory variables except the ones with high near-multicollinearity in their respective between or within space.
13. If one of your variable is time-invariant z(i) (Institutional quality ICRG for Burnside Dollar), run a baseline Hausman Taylor estimation including all X(i.) as instruments. Comment the results.
14. If one of your variable is time-invariant z(i) (Institutional quality ICRG for Burnside Dollar), run a between regression on z(i) explained by X(i.) and other time invariant variable (only with N observations). If the R² is low, this may signal X(i.) are weak instruments poorly correlated with the variable z(i) to be instrumented. Comment.
15. Optional: mention or propose improvements to the Python, STATA, SAS or R code (copy it here). Optional: propose improvements, additional insights, and you do not know how to code them.

PART B (update results)

1. Download 5 panel data variables from World Bank and/or IMF and/or FRED databases for the recent period (1990-2020) and for the largest coverage of emerging economies: GDP/head, GDP/head PPP-adjusted (very last update), Log(population), Foreign aid/GDP (ODA), of log an index of corruption (or good public sector governance) from the World Bank. From now on, consider as your sample only country-year observations which are available for ALL the 5 variables for at least TWO CONSECUTIVE years for a given country. The full class may coordinate for this updated database. In all the following questions except perhaps the last one, the PPP adjusted GDP is not used. So we consider 4 variables excluding GDP/head PPP adjusted.
2. Compute 2 growth rates using the difference of log: the growth of GDP/head (difference of log, denoted GDPg), the growth of foreign aid ODAg (but NOT the growth for foreign aid/GDP: remove the difference of log of GDP from the difference of log of foreign aid/GDP).
3. Compute the between average over time for the first period and for the second period for the 6 variables. Provide the top 10 of countries for ODA/GDP with average over time for each period.
4. Compute the proportion of country-years observations in your database such that $0 \leq \text{ODA/GDP} < 0.5\%$

5. Compute the between and within transformations of the 6 variables over the full period. Provide the 4 histograms for ODA/GPD, growth of ODA, growth of GDP/head, corruption index for both between and within transformed variables (hence 8 histograms). Comment.
6. Provide the 3 bivariate graphs (with acronyms for observations NIC12, for Nicaragua 2012) for between and within (hence 6 graphs) of growth of GDP/head (vertical axis) with (1) ODA/GDP, (2) the growth of ODA; of corruption index with ODA/GDP. Comment. ###7. Comment the between versus within correlation matrix for the 6 variables in this order
8. Run a one-way fixed effect foreign aid regression on ODA/GDP function of $\ln(\text{Population})$ and $\ln(\text{GDP/head})$. Comment.
9. Run a one-way fixed effect of Corruption Index function of $\ln(\text{GDP/head})$, of ODA/GDP and the growth of ODA. Comment.
10. Run a one-way fixed effect with the growth of GDP/head function of $\ln(\text{GDP/head})$, ODA/GDP, the growth of ODA and the Corruption index.
11. Propose an additional interesting estimation using this database.
12. Compute the between and within transformations of the 11 variables over the full period. Provide histograms for ODA/GPD, growth of ODA, growth of GDP/head for both between and within transformed. Comment.