

Problem Set 5

Due by 3/25

A researcher is interested in the relationship between a state's mortality rate and its spending on hospitals (and other health services). The zip file `PS5.zip` contains data sets on state mortality rate, state spending, state education level and state per capita income for the years 1993 through 2015. The data come from different sources such as the US Census, the Bureau of Economic Analysis and the US Mortality Database. More specifically, the data files are as follow:

- a. `mortality_data.csv`,
 - b. `income_data.csv`,
 - c. `education_0715.csv`,
 - d. education data for the years 1993 through 2006 are in the folder *education*, one file per year: `education_1993.csv`, ..., `education_2006.csv`,
 - e. expenditure data for the years 1993 through 2015 are in the folder *expenditure*, one file per year: `expnd_1993.csv`, ..., `expnd_2015.csv`.
- (1) Import the *mortality* data set and name it `mort_data`. Keep only the observations for the years 1993 through 2015.
 - (2) Change the column names for columns 4 through 11, to, `[mort_rate, prob_death, ave_length_surv, num_of_surv, num_of_deaths, num_years_lived, num_years_left, life_expec]`.
 - (3) `Age` column is a character type and needs to be changed to a numeric type. As such, first generate a new column, say `Age2`, by locating the "-" in the character string, and then slicing the string from the first character to the chracter just before "-". To this end, you can use `str_locate()`, `str_length()`, `str_sub()` functions from `stringr` package. Then, convert this to a numeric type using `as.numeric()`.
 - (4) Generate a new column, say `age_group` by cutting `Age2` to three intervals: `[0,18)`, `[18,64)`, `[65,)`. Use `cut()` function and assign the labels "`<18`", "`18-64`", "`>64`".
 - (5) Drop `Age` and `Age2` columns, and reorder the columns as `[state, year, age_group, 7 mortality variables]`.
 - (6) Using `aggregate()` function calculate the sum for the 7 mortality variables by `[state, year, age_group]`. When using `sum()` in `aggregate()`, don't forget to specify `na.rm=TRUE`.
 - (7) Import the *income* data set and name it `inc_data`. Note that the data set is in the *wide* form and needs to be converted to the *long* form. To do so, use `reshape()` and set argument `varying` to the column names corresponding to multiple years in the wide form, and set argument `sep = "."`.
 - (8) Drop the last column, and sort the income data by as `[state, year]`.
 - (9) Append the education data sets from 1993 through 2006 and `education_0715`, name it `educ_data`. Rename columns 3 and 4 as `[phs, pcoll]`.
 - (10) Append the expenditure data sets from 1993 to 2015, and name it `expnd_data`. Note that the columns may have been named slightly different for some years.
 - (11) Merge `inc_data` and `educ_data` by `state` and `year`, and name the merged data set `data`. Notice that this is a *one-to-one* merge.

- (12) Merge `data` and `expnd.data` by `state` and `year`, and name the merged data set again `data`. Notice again that this is a *one-to-one* merge.
- (13) Merge `mort.data` and `data` by `state` and `year`, and name the merged data set again `data`. Notice that this is a *many-to-one* merge.
- (14) Remove `mort.data`, `educ.data`, `expnd.data`.
- (15) Change the measurement of `pinc`, `tot_revenue`, `taxes`, `tot_expnd`, `education`, `public_welfare`, `hospital`, `health` to in 10,000 dollars, i.e., divide each by 10,000.
- (16) Change the measurement of `phs`, `pcoll` to ratios, i.e., divide each by 100.
- (17) Generate a table of descriptive statistics for your data set using `stargazer()`. You can export the table by setting arguments `type = "html"`, `out = "descriptives.doc"`.
- (18) Regress `mort_rate` for the age group 65 and older, on an intercept, `health`, `hospital`, `log pinc`, `phs` and `pcoll`. Name the results `reg1`.
- (19) Regress `mort_rate` for the age group 65 and older, on an intercept, `health`, `hospital`, `log pinc`, `phs`, `pcoll` and state dummies. Name the results `reg2`.
- (20) Regress `mort_rate` for the age group 65 and older, on an intercept, `health`, `hospital`, `log pinc`, `phs`, `pcoll` and state and year dummies. Name the results `reg3`.
- (21) Generate a table for the regression results using `stargazer()`. You can export the table by setting arguments `type = "html"`, `out = "regressions.doc"`.