# Geneva Academics in Management and Economics

## *THE DO'S AND DON'TS OF A DO-FILE*

tips on how to do empirical work with Stata

Emmanuel Milet (emmanuel.milet@unige.ch)

$15^{th}$ Nov. 2016

# The do's and don'ts of a do-file

Why bother with a learning how to write a do-file?

- **THE** most important tool of your empirical analysis.
- You may share it with co-authors
- You will return to it several month after submitting your paper (and you'll possibly do this multiple times too)
- You may be asked by an editor to make your code and data publicly available

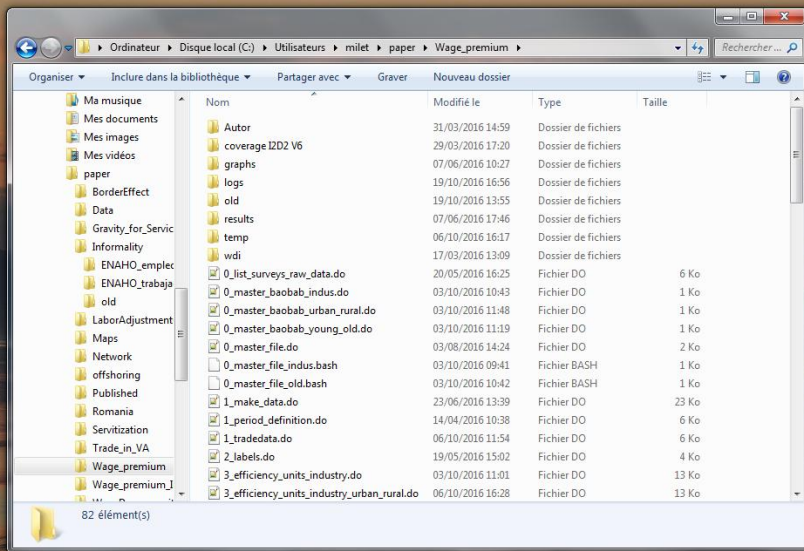# The do's and don'ts of a do-file

Some important steps:

1. Folder
2. Master do-file
3. Logs
4. Names and label variables
5. Regression output
6. Various commands

# The do's and don'ts of a do-file: FOLDER

You are going to generate TONS of output while doing your empirical analysis. To make things *easier* and *clearer*, I recommend to create a folder for each type of output:

- logs
- graphs
- regression results
- tables and other output for the stylized facts
- manuscript
- an "old" folder where you can store do-files and other output that you do not need anymore

Your data and do-file can be stored in the root folder.

Organiser ▾   Inclure dans la bibliothèque ▾   Partager avec ▾   Graver   Nouveau dossier

| Nom | Modifié le | Type | Taille |
|---|---|---|---|
| Autor | 31/03/2016 14:59 | Dossier de fichiers | |
| coverage I2D2 V6 | 29/03/2016 17:20 | Dossier de fichiers | |
| graphs | 07/06/2016 10:27 | Dossier de fichiers | |
| logs | 19/10/2016 16:56 | Dossier de fichiers | |
| old | 19/10/2016 13:55 | Dossier de fichiers | |
| results | 07/06/2016 17:46 | Dossier de fichiers | |
| temp | 06/10/2016 16:17 | Dossier de fichiers | |
| wdi | 17/03/2016 13:09 | Dossier de fichiers | |
| 0_list_surveys_raw_data.do | 20/05/2016 16:25 | Fichier DO | 6 Ko |
| 0_master_baobab_indus.do | 03/10/2016 10:43 | Fichier DO | 1 Ko |
| 0_master_baobab_urban_rural.do | 03/10/2016 11:48 | Fichier DO | 1 Ko |
| 0_master_baobab_young_old.do | 03/10/2016 11:19 | Fichier DO | 1 Ko |
| 0_master_file.do | 03/08/2016 14:24 | Fichier DO | 2 Ko |
| 0_master_file_indus.bash | 03/10/2016 09:41 | Fichier BASH | 1 Ko |
| 0_master_file_old.bash | 03/10/2016 10:42 | Fichier BASH | 1 Ko |
| 1_make_data.do | 23/06/2016 13:39 | Fichier DO | 23 Ko |
| 1_period_definition.do | 14/04/2016 10:38 | Fichier DO | 6 Ko |
| 1_tradedata.do | 06/10/2016 11:54 | Fichier DO | 6 Ko |
| 2_labels.do | 19/05/2016 15:02 | Fichier DO | 4 Ko |
| 3_efficiency_units_industry.do | 03/10/2016 11:01 | Fichier DO | 13 Ko |
| 3_efficiency_units_industry_urban_rural.do | 06/10/2016 16:28 | Fichier DO | 13 Ko |

Ma musique
Mes documents
Mes images
Mes vidéos
paper
  BorderEffect
  Data
  Gravity_for_Servic
  Informality
    ENAHO_emplec
    ENAHO_trabaja
    old
  LaborAdjustment
  Maps
  Network
  offshoring
  Published
  Romania
  Servitization
  Trade_in_VA
  Wage_premium
  Wage_premium_I

82 élément(s)

# The do's and don'ts of a do-file: MASTER DO FILE

A master do-file is a do-file which launches your other do-files in a specific order.

- It allows you to organize your do-files, and add comments to them.
- It makes sure that your final results (in your paper) are produced correctly.
- It allows you to declare a directory path, as well as other options that you may use across all your do-files.

```
 1                        * ----------------------------
 2                        * Wage Premium paper: master file
 3                        * ----------------------------
 4       *
 5       clear*
 6       set more off
 7       global sysdate=c(current_date)
 8       global path "C:\Users\milet\paper\Wage_premium"      // Directory Unige
 9       *global path=c(pwd)                                   // Directory Baobab
10       cd $path
11       *
12       *
13       *
14       *
15       * Data Creation
16       do 1_make_data.do          // loops overthe raw data, extract variables, and make the final dataset
17       do 1_period_definition.do  // selects the surveys that we need to define the various periods
18       do 2_labels.do             // puts the labels in the final dataset
19       do 1_tradedata.do          // gets us the imports of capital, R&D intensive and R&D un-intensive goods
20       *
21       *
22       * Compute the prices and quantities in terms of efficiency units
23       do 3_efficiency_units.do                      // Get the supply and the wages in efficiency units (needed for the descriptive
         statistics)
24       do 3_efficiency_units_industry.do             // efficiency units at the industry level (to get the elasticities of substitution)
25       *do 3_efficiency_units_young_old.do            // efficiency units for young and old wokers
26       do 3_panel_data.do                            // creates the dataset at the country*level with skill premium and skill supply
27       *
28       *
29       * Stylized Facts
30       do 4_share_of_labor_income.do                 // get the share of labor income of total income for each country
31       do 4_table_list_IncGroup_WrldRegion.do        // get the list of countries
32       do 4_table_years_period.do                    // get the years used for each country in each period
33       do 4_wage_distribution_percentiles.do         // Evolution of the 10th, 50th, 90th percentiles. All workers, women only, men only
34       do 4_wage_distribution_education_level.do     // Evolution of the composition-adjusted wages by education level
35       do 4_skillpremium.do                          // Evolution of the composition-adjusted wages by education level
36       *
37       * Elasticities of substitution
38       do 5_elasticities.do                                 // get the elasticities of substitution between skilled and unskilled workers
39       *
40       *
41       *
42       ** end of file
```

# The do's and don'ts of a do-file: LOGS

Log files are used to keep track of EVERYTHING your do-file does, and more:

- It saves what Stata erases in the result window.
- You should ALWAYS start your do-file by opening a log file, and end by closing it.
- Your log files should be stored in a specific log folder.
- Give the SAME name to your do-files and your logs.
- You can also add the date of the day in the log's name, so that you do not erase them each time your run your do-file.

# The do's and don'ts of a do-file: LOGS

Example:

# The do's and don'ts of a do-file: NAMES & LABELS

1. Give explicit names to your variables.
2. Adopt a naming rule and stick to it! For instance:
   - log variables can be named: ln_*varname*
   - dummy variables created from variables: *varname*_dum
3. Attach labels to the variables.
4. Ideally, you should have a do-file with all the labels.
5. Labels will show up in any output that you produce.

```stata
106    gen rdoecd_incapoecd=lnm_rd_high_oecd-lnm_cap_high_oecd // share of R&D imports in imp
107    *
108    * Log of FDI inflows
109    gen lnfdi=ln(fdi_inflow)
110    *
111    * put labels
112    label var relwage            "Ln skill premium"
113    label var relsupply          "Ln relative supply"
114    label var captotal_intotal   "Share of capital imports in total imports"
115    label var caphigh_inhigh     "Share of capital imports in imports from high-income
116    label var capoecd_inoecd     "Share of capital imports in imports from OECD countr
117    label var rdtotal_incaptotal "Share of R&D imports in capital imports"
118    label var rdhigh_incaphigh   "Share of R&D imports in imports from high-income cou
119    label var rdoecd_incapoecd   "Share of R&D imports in imports from OECD countries"
120    label var lnfdi              "Ln FDI inflow"
121    label var year               "Year"
122    label var ccode              "Country code"
123    label var skilled_worker     "1=skilled worker"
124    label var lnwage_p           "Ln wage"
125    label var supply             "# workers"
126    label var lnsupply           "Ln # workers"
127    label var nworker_y          "# workers (sum of weights from the survey)"
128    label var fdi_inflow         "FDI inflow"
129    label var lnm_all            "Ln aggregate imports"
130    label var lnm_high           "Ln Imports from high income countries"
131    label var lnm_high_oecd      "Ln imports from OECD countries"
132    label var lnm_cap_all        "Ln capital imports"
133    label var lnm_cap_high       "Ln capital imports from high income countries"
134    label var lnm_cap_high_oecd  "Ln capital imports from OECD countries"
135    label var lnm_rd_all         "Ln R&D imports"
136    label var lnm_rd_high        "Ln R&D imports from high income countries"
137    label var lnm_rd_high_oecd   "Ln R&D imports from OECD countries"
138    label var lnm_notrd_all      "Ln capital non-R&D imports"
139    label var lnm_notrd_high     "Ln capital non-R&D imports form high income countrie
140    label var lnm_notrd_high_oecd "Ln capital non-R&D imports from OECD countries"
141    *
142    *
143    order ccode iso3 world_region inc_group_last skilled_worker supply lnsupply relwage r
144    compress
```

Filter variables here

| Name | Label |
|---|---|
| landlocked | 1 if landlocked |
| continent | Continent |
| year | Year |
| industry | |
| nworker_y | # workers (sum of weights from the survey) |
| m_all | Total imports (K USD) |
| m_high | Total imports from High-income countries (K USD) |
| m_high_oecd | Total imports from High-income OECD countries (K USD) |
| m_cap_all | Total capital imports (K USD) |
| m_cap_high | Capital imports from High-income countries (K USD) |
| m_cap_high_oecd | Capital imports from High-income OECD countries (K USD) |
| m_rd_all | Total R&D imports (K USD) |
| m_rd_high | R&D imports from High-income countries (K USD) |
| m_rd_high_oecd | R&D imports from High-income OECD countries (K USD) |
| m_notrd_all | Total R&D imports (K USD) |
| m_notrd_high | NOT R&D imports from High-income countries (K USD) |
| m_notrd_high_oecd | NOT R&D imports from High-income OECD countries (K USD) |
| fdi_inflow | FDI inflow |
| lnwage_p | Ln wage |
| lnm_all | Ln aggregate imports |
| lnm_high | Ln Imports from high income countries |
| lnm_high_oecd | Ln imports from OECD countries |
| lnm_cap_all | Ln capital imports |
| lnm_cap_high | Ln capital imports from high income countries |
| lnm_cap_high_oecd | Ln capital imports from OECD countries |
| lnm_rd_all | Ln R&D imports |
| lnm_rd_high | Ln R&D imports from high income countries |

Properties

# The do's and don'ts of a do-file: REGRESSION OUTPUT

Regression results are probably the most important output from your do-files.

It is important that they look nice!

- The "combo" esttab/estout commands is a good way to produce and export regression results.
- There is also outreg2 (personally not a fan)

|  | (1) | (2) |
|---|---|---|
|  | c_y_beta | c_y_beta |
| Ln FDI inflow | 0.002 | 0.005 |
|  | (0.191) | (0.389) |
| Ln aggregate i~s | -0.046*** | -0.066*** |
|  | (-3.241) | (-3.635) |
| Share of capit~a | 0.011 | 0.068 |
|  | (0.382) | (1.384) |
| Share of R&D i~ | 0.065 | 0.036 |
|  | (1.137) | (0.463) |
| Constant | 1.571*** | 1.218*** |
|  | (7.655) | (4.991) |
| Observations | 238 | 232 |

t statistics in parentheses
* p<0.1, ** p<0.05, *** p<0.01

. esttab using "./results/baseline.tex", compress nogaps label drop(c_dum*) ///
>                starlevels(($^{c}$) 0.1 ($^{b}$) 0.05 ($^{a}$) 0.01) b(%5.3f) t(%5.3f) replace

(output written to ./results/baseline.tex)

. estso clear

.

end of do-file

Do-file Editor - 3_panel_data.do

File Edit View Project Tools

5_elasticities.do × 3_panel_data.do* × 3_efficiency_units_industry.do × 0_master_file.do × Untitled1.do ×

```
196    *
197        * regress the country-level dummies on trade and fdi
198    * ---------------------------------------------
199    * 1) baseline:
200    estso: reg c_y_beta lnfdi lnm_all captotal_intotal rdtotal_incaptotal c_dum*
201    * 2) weighted regression:
202    estso: reg c_y_beta lnfdi lnm_all captotal_intotal rdtotal_incaptotal c_dum* [w=abs(c_y_tstat)]
203    *
204    esttab, compress nogaps label drop(c_dum*) starlevels(* 0.1 ** 0.05 *** 0.01) b(%5.3f) t(%5.3f)
205    esttab using "./results/baseline.tex", compress nogaps label drop(c_dum*) ///
206            starlevels(($^{c}$) 0.1 ($^{b}$) 0.05 ($^{a}$) 0.01) b(%5.3f) t(%5.3f) replace
207    estso clear
208
209    *
```

Line: 207, Col: 13  CAP  NUM  OVR

# The do's and don'ts of a do-file: REGRESSION OUTPUT

```
{
\def\sym#1{\ifmmode^{#1}\else\(^{#1}\)\ fi }
\begin{tabular}{l*{2}{c}}
\hline\hline
                &\multicolumn{1}{c}{(1)}&\multicolumn{1}{c}{(2)}\\
                &\multicolumn{1}{c}{Dem. shifters}&\multicolumn{1}{c}{Dem. shifters
\hline
Ln FDI inflow    &       0.002        &       0.005         \\
                 &      (0.191)       &      (0.389)        \\
Ln aggregate imports&      0.046{$^{a}$}&      0.066{$^{a}$}\\
                 &      (3.241)       &      (3.635)        \\
Share of capital imports in total imports&       0.011        &      0.068
\\
                 &      (0.382)       &      (1.384)        \\
Share of R\&D imports in capital imports&       0.065        &      0.036
\\
                 &      (1.137)       &      (0.463)        \\
Constant         &      1.571{$^{a}$}&      1.218{$^{a}$}\\
                 &      (7.655)       &      (4.991)        \\
\hline
Observations     &      238        &      232        \\
\hline\hline
\multicolumn{3}{l}{\footnotesize \textit{t} statistics in parentheses}\\
\multicolumn{3}{l}{\footnotesize {$^{c}$} p<0.1, {$^{b}$} p<0.05, {$^{a}$} p<0.01}
\end{tabular}
}
```

# The do's and don'ts of a do-file: REGRESSION OUTPUT

|  | (1) | (2) |
| --- | :---: | :---: |
|  | Dem. shifters | Dem. shifters |
| Ln FDI inflow | 0.002 | 0.005 |
|  | (0.191) | (0.389) |
| Ln aggregate imports | -0.046$^a$ | -0.066$^a$ |
|  | (-3.241) | (-3.635) |
| Share of capital imports in total imports | 0.011 | 0.068 |
|  | (0.382) | (1.384) |
| Share of R&D imports in capital imports | 0.065 | 0.036 |
|  | (1.137) | (0.463) |
| Constant | 1.571$^a$ | 1.218$^a$ |
|  | (7.655) | (4.991) |
| Observations | 238 | 232 |

$t$ statistics in parentheses

[c] $p < 0.1$, [b] $p < 0.05$, [a] $p < 0.01$

# The do's and don'ts of a do-file: USEFUL COMMANDS

The **COLLAPSE** command:

- It reduces the dimensionality of your dataset and calculates many statistics (count, median, mean, standard deviation, percentiles, first obs. ...) base on certain dimensions (useful when using panel data, or data with more than 2 dimensions: firm×product×destination×year for instance)
- BEWARE: Weight normalization impacts only the sum, count, sd, semean, and sebinomial statistics
- A weighted average cannot be obtained directly from collapse, despite the fact that you can ask for the mean of a variable, and specify weights.

# The do's and don'ts of a do-file: USEFUL COMMANDS

The **TAG** command:

- It creates a dummy variable taking the value 1 for each occurrence of a variable.
- It is useful when you have nested dimensions (say individual/household/county/region), and want to quickly get statistics at a specific level.

# The do's and don'ts of a do-file: USEFUL COMMANDS

The **TAG** command:

- It creates a dummy variable taking the value 1 for each occurrence of a variable.
- It is useful when you have nested dimensions (say individual/household/county/region), and want to quickly get statistics at a specific level.

The **GROUP** command:

- Creates a variable taking integer values from 1 to $\mathbb{N}$ for each occurrence of a given variable or list of variables (i.e. an occurrence is therefore a combination of various variables).
- This is handy to create fixed-effect variables.
- Turns a string variable into a numeric variable (the command **encode** does this too, but only for 1 variable).

# The do's and don'ts of a do-file: USEFUL COMMANDS

The **LEVELSOF** command:

- Especially useful for loops
- IT lists all occurrences of a variable, and stores the list into a local variable.
- You can then loop over the elements of this list.
- The nice feature is that Stata does not create any variable, the elements of the list (i.e. the occurrences) are *local* elements.

The **#delimit** command:

- This is a line breaker.
- Essentially the same as using three forward slash bars: ///
- It is more convenient though (personal opinion).

```
levelsof ccode, local(ccode_local)
]foreach c of local ccode_local{
preserve
keep if ccode=="`c'"
    *
        * 1a) make the graph
        * ------------------
    #delimit ;
    twoway (scatter relsupply year, msymbol(Dh) mcolor(emerald))
           (line relsupply year, lcolor(emerald))
           (scatter relwage year, msymbol(Oh) mcolor(dkorange) yaxis(2))
           (line relwage year, lcolor(dkorange) yaxis(2)),
           scheme(s1color) xtitle("") ytitle("") ytitle("", axis(2))
           legend(order(2 "Relative supply index (left axis)" 3 "Skill premium (right axis)"
           region(lpattern(blank)));
    #delimit cr
    graph export "./graphs/skillpremium_`c'.pdf", as(pdf) replace
    *
    *
```

- That's it for today.
- You can find a written version fo all this on my webpage:
  `http://emmanuelmilet.weebly.com/`
- Thank you for your attention