

```

1  * * * * *
2  *
3  * PROGRAM:      Stata Tutorial
4  * PROGRAMMER:   Suanna Oh (SO)
5  * INSTITUTION:  Columbia University
6  * DATE:         November 30, 2016
7  *
8  * * * * *
9
10
11 capture cd "C:\Users\suannaoh\Dropbox\Programming TA"
12 capture cd "C:\Users\NewUser\..."
13 // capture: executes following commands, suppressing all output and ignoring error
14 // for instance, it will ignore line 14 if no such directory exists on the computer
15 clear all // clear all memory
16 set more off // run the entire do file without
17 displaying -more-
18 /* Before getting started
19
20 1. Set up a project folder
21
22 Create a folder for the project materials and change directory to this folder using [cd] as
23 above
24 Within this folder we can create subfolders such as:
25 - Data: for saving data output files
26 - Log files: for saving log files
27 - Raw data: for saving raw data files
28 - Do files: for saving do files
29
30 2. Do-file documenting tools:
31
32 // : comment out words following
33 * : comment out the entire line if placed at the beginning of a line
34 /* */ : comment out multiple lines between the two stars
35 /// : break long commands into multiple lines
36
37 */
38
39
40 *****
41 * Task 1
42 *****
43
44 /*
45
46 ** Note 1a: Merging two data sets
47
48 One should follow the steps below:
49 1. Decide which type of merging needs to be performed (1:1, m:1, or 1:m)
50 1:1 indicates there is one observation per ID in both data sets
51 m:1 indicates there are m observations per ID in the master data
52 1:m indicates there are m observations per ID in the using data
53 2. Check if there will be any conflicting information in the data sets and decide how to
54 deal with them
55 If the two data sets have variables with the same name, the master data information in
56 kept by default
57 3. In both the master and using data, check that the unit of observation is consistent
58 Check for any duplicate ID's
59 4. Make sure that the ID variables in both data sets are named the same way
60 5. Review _merge results
61
62 ** Note 1b: Using macros
63
64 Macros: A macro is simply a name associated with some text.
65 Local macros
66 - Created by typing either:
67 * local macro_name1 "macro content"
68 * local macro_name1 = expression to be evaluated

```

```

68 - Referenced by writing `macro_name1' in the same do-file
69 - Local macros are deleted after the current do-file finishes running
70 - If we reference a macro that does not exist, Stata returns an empty string for the
macro
71 - One can also temporarily save files under a local macro name using [tempfile]
72 Global macros:
73 - Global macros are only deleted after Stata is closed
74 - See help macro
75
76 */
77
78 ** 1. Importing data
79 import delimited using "Raw data/Respondent_list.csv", clear
80 label variable pid "Respondent ID"
81 lab var name "Respondent name"
82 lab var village "Village name"
83
84
85 ** 2. Merging data
86 // We first check if the individual level data
87 duplicates report pid // we see that the IDs are not actually unique. There
are 4 observations with duplicate ids.
88 duplicates tag pid, gen(dup) // this generate a var named dup, which equals 1 if id
is not unique
89 list if dup // we see that 1123 is a duplicate entry so we just
drop this
90 duplicates drop
91 // this only drops the complete duplicates (all variables are the same for two
observations)
92 // however, there are two different respondents with the same id 1491
93 // here we use the information Fake Name 394's id is actually 1495, not 1491.
94 replace pid=1495 if pid==1491 & name=="Fake Name 394"
95 isid pid // confirm pid uniquely identifies observations
96 drop dup // drop this variable since we no longer need it
97
98
99 // Then, we need to create dta version of the village data
100 // Instead of saving the previous data and opening new data, we will preserve the
individual data
101 preserve // preserve the current data set
102 import delimited using "Raw data/Village_info.csv", clear
103 isid village_code
104 rename village_code village // since the village id is named "village" in the other
data, we should rename this var
105 tempfile village // declare that [village] will be used as a macro name to
save a temporary file
106 save `village', replace // one uses the regular saving commands to save temporary
files
107 restore // restores the preserved data
108
109
110 // We can now merge the two data sets using the variable "village"
111 merge m:1 village using `village'
112 // since there are many respondents per village, we specify m:1 that this is the case
113 // we see that all respondents in our respondent list are matched to villages (there is
0 case where _merge==1)
114 // we also see that there are 110 villages with no respondents (_merge==2). We don't
need them in the data, so let's drop them
115 drop if _merge==2
116 drop _merge
117 order village, before(villagename) // this order the variable village before villagename
118 tabulate treatment
119
120
121 ** 3. Summarizing and saving
122 bysort subcounty: sum treatment
123 tab subcounty, sum(treatment) // more compact way to do the above
124 bys subcounty: egen subcounty_size=count(treatment) // calculate number of people in
each subcounty
125 bys subcounty: egen treatment_share=mean(treatment) // calculate the share of treated
people in each subcounty
126 // bysort: run commands on subsets of data after sorting data

```

```

127 // egen command can also calculate std, min, max, median, rank, total, etc.
128 // see help egen
129 save "Data/Respondent_village_data.dta", replace
130
131
132
133
134
135 *****
136 * Task 2
137 *****
138
139 /*
140
141 Note 2: Working with time variables
142
143 - Human readable forms (HRFs) is the readable date form usually stored as string variables.
144 While easy to read, Stata cannot work with these variables.
145 ex ) "2015-11-17", "2010 Jul 12"
146 - Stata internal form (SIF) is the date form Stata can understand.
147 These are numeric variables usually calculated with respect to 01jan1960 midnight.
148 One can convert string time variables into SIF forms using various datetime functions.
149 - To use any ts- commands, one must declare the time variable using [tsset]
150
151 */
152
153
154 ** 1. Importing data and converting time variables
155 import delimited using "Raw data/GDP.csv", clear
156 generate time_day = date(date, "YMD") // days since 01jan1960
157 format time_day %td // format the variable so that it is readable
158 generate time_quarter = qofd(time_day) // quaters since 1960q1
159 format time_quarter %tq
160 label var time_quarter "Quarter"
161
162
163 ** 2. Regress GDP on 1 period lagged var and 2 period lagged var.
164 tsset time_quarter // declare data to be time series
165 reg gdp L.gdp L2.gdp // L. (F.): create or use a lag (lead)
166 variable
167
168 ** 3. Detrend GDP using the BK filter and plot the series
169 tsfilter bk gdp_bk=gdp // de-trend gdp using a filter
170 label var gdp_bk "Detrended GDP"
171 tsline gdp_bk if time_quarter>tq(1990q1), lwidth(medthick) graphregion(color(white)) //
172 plot gdp_bk against time // I used thicker lines and white background to make the graph look nicer
173 graph export Graphs/fig_gdp.pdf, replace
174
175
176
177
178 *****
179 * Task 3
180 *****
181
182
183 ** 1. Manipulating string variables
184 import excel "Raw data\Data_Extract_From_World_Development_Indicators.xlsx", sheet("Data")
185 firstrow clear
186 drop if mi(CountryName)
187 gen ShortName=substr(SeriesName,1,strpos(SeriesName,"(")-1)
188 // substr helps us get a portion of a string
189 // strpos returns position of a particular string in another string
190 // go to 'help functions' -> 'string functions' to see other string functions available
191 split SeriesCode, gen(Code) parse(".")
192 // we can split a string based on a specific character
193 rename Code3 Code
194
195 ** 2. How many countries have St. in their names?

```

```

196 encode ShortName, gen(SeriesNum) // create a new variable that assigns a unique
197 number to each series
198 count if regexm(CountryName,"St.") & SeriesNum==1
199 // regexm returns 1 if the first argument (string) contains the second argument
200 (regular expression)
201
202 ** 3. Reshaping
203 keep SeriesNum CountryName CountryCode YR*
204 order SeriesNum CountryName CountryCode YR*
205 reshape long YR, i(SeriesNum CountryName) // converts wide data to long data
206 rename _j year
207 rename YR value
208 replace value="" if value==".."
209 destring value, replace
210 // destring converts a string variable to a numeric variable
211 // tostring performs the opposite function
212 // note the difference from encode and decode
213 sort SeriesNum CountryName year
214 by SeriesNum CountryName: gen growth_rate=(value-value[_n-1])/value[_n-1]
215 by SeriesNum CountryName: egen avg_growth=mean(growth_rate)
216 save "Data/Dev_indicators.dta", replace
217
218
219 *****
220 * Task 4
221 *****
222
223 /*
224
225 Note 4a: Using a loop
226
227 Using a loop, we can repeat a set of commands to each element in a specified list:
228 efficient coding!
229 Structure of a foreach loop:
230 foreach [lname] of(in) [list] {
231     commands referring to `lname'
232 }
233 "foreach" repeatedly sets local macro `lname' to each element of the [list] and executes
234 the commands enclosed in braces {}.
235 Allowed lists include:
236 - foreach lname in any_list {
237 - foreach lname of varlist varlist {
238 - foreach lname of numlist numlist {
239 - forvalues lname = range {
240     forvalues is the fastest way to loop over consecutive values
241
242 Note 4b: if programming command
243
244 "if" executes a set of commands if exp is true (nonzero). Otherwise the statement following
245 [else] is executed
246 - not to be confused with the if qualifier
247 Basic structure:
248 if exp {
249     multiple commands
250 }
251 else {
252     multiple commands
253 }
254
255 */
256
257 ** 1. Reshaping data
258 use "Data/Dev_indicators.dta", clear
259 bys CountryName SeriesNum year: keep if _n==1 // keep one observation per
260 country*series*year
261 keep SeriesNum CountryName CountryCode year growth_rate
262 reshape wide growth_rate, i(CountryName year) j(SeriesNum)
263 drop if year==2010

```

```

262
263
264 ** 2. Looping
265 foreach x of numlist 1/3{ // let `x` take the value 1, 2,
and then 3 for each series
266     bys CountryName: egen count`x`=count(growth_rate`x') // calculate the number of
observations per country for each series
267 }
268 egen min_count=rowmin(count*) // calculate the minimum number
of observations
269
270 capture log close // close any log files open
271 log using "Log files/correlation.log", replace // start a log file
272 levelsof CountryName // get the list of all country
names for the country loop
273
274 foreach x in `r(levels)'{ // here `x` becomes each
country name
275     di "`x'" // print country name
276     quietly sum min_count if CountryName=="`x'" // count the minimum number of
observations for this country
277     if r(mean)<4 di "Missing values" // if the number is less than
4, print "Missing values"
278     else pwcorr growth_rate? if CountryName=="`x'" // otherwise calculate
correlations
279     di "" // creates a line gap between
each country result for readability
280 }
281 log close
282
283
284 *****
285 * Task 5
286 *****
287
288
289 // First we can clean the animals data
290 use "Raw data/animals.dta", clear
291 replace type="NA" if mi(type)
292 isid village hhn id animal type
293
294
295 // need to impute values of animals by animal type
296 destring number price, replace force
297 egen category=group(animal type)
298 levelsof category if mi(price)
299 foreach y in `r(levels)'{
300     qui sum price if category==`y'
301     replace price=`r(mean)' if mi(price) & category==`y'
302 }
303 gen total_value=number*price
304 collapse (sum) total_value, by(village hhn id animal)
305 rename animal asset_type
306 tempfile animal
307 save `animal', replace
308
309
310 // Next we can clean other asset data
311 use "Raw data/cash.dta", clear
312 foreach x of varlist balance balance2 cash jewelry cloth{
313     if "`x'"=="jewelry" gen `x'_DK=`x'=="DK" | `x'=="0"
314     else gen `x'_DK=`x'=="DK"
315     destring `x', replace force
316     sum `x'
317     replace `x'=`r(mean)' if `x'_DK==1
318     drop `x'_DK
319 }
320 egen saving=rowtotal(balance balance2), m
321
322
323 // To create an asset level data, we can append partial data sets instead of reshaping
324 local k=1 // we can increase the value of this macro by 1 for each loop iteration

```

```
325 foreach x of varlist cash jewelry cloth saving{
326     preserve
327     keep village hhn id `x'
328     drop if mi(`x')
329     rename `x' total_value
330     gen asset_type=`x'"
331     tempfile data`k'
332     save `data`k'', replace
333     restore
334     local `k++'      // increase the value of macro `k' by 1
335 }
336
337 use `data1', clear
338 append using `data2' `data3' `data4'
339 append using `animal'
340
341 sort village hhn id asset_type total_value
342 save "Data/asset_data.dta", replace
343
344 /*
345
346 Some useful websites for learning Stata:
347 http://data.princeton.edu/stata/
348 http://www.ats.ucla.edu/stat/stata/
349 http://www.cpc.unc.edu/research/tools/data\_analysis/statatutorial/
350 http://staskolenikov.net/stata/
351
352 */
353 exit
354
355
356
```