

Stata Introductory Series

Stata 202 – Data Manipulation

Our introductory Stata session #100 laid out the fundamentals of interacting with Stata; setting up DO and LOG files to easily save, replicate, modify, and share your analysis and its results; and how to detect and deal with missing and special values as you check your work.

Stata 202 will build off of last week's training on Data Cleaning (Stata 201), but do not worry if you did not attend last week or things are slightly confusing or a bit more difficult to follow along on. This week's session does not require attending prior trainings, and we will be around to answer questions and debug during and after the session. The goal of this session is to continue learning how to import, clean, and manipulate messy data that you may encounter in the real world from administrative, survey, or your own data collection.

What you will get out of this session:

- » Learn how to utilize messy, real-world data
- » Reshape, merge, and append your data sets
- » Review how to find and deal with duplicate records
- » Review how to clean and convert string to numeric and date data or vice versa
- » Review how to convert information stored in text strings into indicators

Basic command structure

<u>command</u>	<u>objects</u>	<u>conditions</u>	<u>, options</u>
<u>use</u>	<u>file.dta</u>		<u>, clear</u>
<u>generate</u>	<u>age = 15</u>	<u>if AGE2 == 15</u>	
<u>tabulate</u>	<u>state</u>	<u>if country == "US"</u>	<u>, missing</u>

Helpful resources

- » Stata manual: access by typing "help command" in the stata console
- » Statalist: <https://www.statalist.org/forums/forum/general-stata-discussion/general>
 - Often will come up if you google a question that isn't covered by the documentation
- » UCLA IDRE: <https://stats.idre.ucla.edu/stata/>
 - Provides helpful tips on how to use Stata as well as the statistics behind the programming
- » UNC CPC: http://www.cpc.unc.edu/research/tools/data_analysis/statatutorial
 - Guide to working with and analyzing data in Stata

Remember: Getting errors is a normal part of programming! The best way to debug is to read through every line carefully.

Next class:

- » Stata 301 – Automating Tasks and Exporting Output

```

1  clear
2  capture log close
3  capture ssc install savesome                /*Some commands have to be installed, but
4                                              only once, capture prevents the error from
5                                              popping up*/
6
7  set more off, permanently
8
9  *Set a working directory where all your project files will be located (replace my username
10 with your own):
11 *cd "D:\Users\CLou\Desktop\"
12 cd "C:\Users\urbanmeet\Desktop\"
13
14 log using "StataClass_$$_DATE.log", replace    /*$$_DATE saves the date in
15                                                  the name of your log file, so
16                                                  you save a log daily. This is
17                                                  helpful if you're changing
18                                                  things and need to go back*/
19
20
21 **This training builds off of the prior "Data Cleaning" session, though attending that
22 training first is not required or necessary to participate in this session.
23     *All materials and files necessary from the prior Data Cleaning training session are
24 provided, and some concepts are reviewed.
25
26 *Open up the "long.dta" Stata data set created at the end of the "Data cleaning" training
27 session:
28 use "long.dta", clear
29 /*Reshaping the data set*/
30 sort id arr_d
31 list in 1/10                                /*Data is currently in long format*/
32
33 *Right now, the data is setup in LONG format, meaning there are repeated
34 observations/rows/records
35     *of the same entity/unit of analysis, usually across time. Sometimes it is easier
36     *or better to work with data in WIDE format, where there is one observation/row/record
37     *for each entity or unit of analysis, as you do not have to worry about double counting.
38     *The code below will help with converting from LONG to WIDE format; the code is similar
39 for WIDE to LONG.
40
41 /*Need to generate a number that indicates how to order the wide format data*/
42 by id (arr_d), sort: egen arrest_n=seq() /*Creates a variable of the number of the
43                                           arrests by person id, the variable not in
44                                           parenthesis is the variable that it's created
45                                           by, the variable in parenthesis is sorted
46                                           within that*/
47 by id (arr_d), sort: gen arrest_n2=_n /* you can also use Stata's "_n" notation with just
48 regular generate.
49                                           The "_n" is essentially the index or (usually)
50 observation number in Stata and quite powerful.
51                                           When specifying groups with "by :", the "_n" index
52 actually resets for each group, which can be
53                                           use to your advantage if you want to know how many
54 observations you have in a group or to mark the
55                                           first or last observation within your group. "_N"
56 is the total # of observations overall
57                                           or within your group when combined with "by :".
58 Also, using [_n-1] and [_n+1] or +2,-2, etc.
59                                           indices directly or appending after a variable name
60 allows you to reference or use the values
61                                           or prior or subsequent observations in sorted data,
62 which is also quite useful.*/
63 list id arr_d arrest_n arrest_n2
64 assert arrest_n == arrest_n2
65 drop arrest_n2
66
67 sort id arr_d
68 list id arr_d arrest_n in 1/10
69

```

```

57 *rename * *_ /*If I ran this line the variables generated in the reshape would be arr_d_1 */
58 /*arr_d_2 instead of arr_d1 arr_d2*/
59
60 reshape wide arr_*, i(id) j(arrest_n) /*Converts to wide format, i is the id, j
61 is the order of the dates, note arr_n
62 becomes suffix*/
63 list id arr_d* in 1/2
64 duplicates report id
65
66 gen t_arr1= arr_d2-arr_d1 /*Calculate the time between arrests*/
67 sum t_arr1
68
69 save "wide.dta", replace
70
71 import excel "Stata Class File_main.xlsx", ///
72 sheet("Sheet1") ///
73 firstrow ///
74 case(lower) ///
75 allstring ///
76 clear
77
78 describe
79 list in 1/10
80
81 duplicates report id
82
83 /*Clean birthdate*/
84 list dob in 1/10
85 gen birth_d=date(dob, "MDY") /*Need to input format of the date of birth,
86 "MDY" or "DMY"*/
87 label var birth_d "Birth Date"
88 format birth_d %td
89 assert birth_d!=. if dob!="" /*if date is not in the right format,
90 the new var will be missing when the
91 old var has a value always want to
92 check this*/
93 list dob birth_d if birth_d==. & dob!=""
94 replace dob="01/01/1985" if dob=="01/001/1985"
95 replace birth_d=date(dob, "MDY")
96 assert birth_d!=. if dob!=""
97 codebook birth_d, d
98 list birth_d dob if (birth_d<=td(01jan1915) | birth_d>td(01jan1996)) /*list new and
99 old var that
100 might not be
101 real*/
102 replace birth_d=. if birth_d<=td(01jan1915) /*replace with missing*/
103 replace birth_d=. if birth_d>td(01jan1996) /*replace with missing*/
104 drop dob
105
106
107 gen age=floor((td(01jan2014)-birth_d)/365) /*Generate age, floor rounds down to
108 the nearest interger*/
109 sum age, d
110 /*Generate categorical age variable with labels*/
111 assert age<=100 | age==. /*Note missing is considered infinity be
112 careful with open ended greater/less than*/
113 recode age ///
114 (18/25=1 "18-25 Years Old") ///
115 (26/30=2 "26-30 Years Old") ///
116 (30/100=3 "30+ Years Old"), ///
117 gen(age_cat)
118 label var age_cat "Age Category"
119 tab age_cat, m
120
121 /*Cleaning string variables, generating sex categorical variables*/
122 tab sex, m
123 replace sex=lower(sex) /*makes string lower case*/
124 replace sex=trim(sex) /*removes leading/trailing spaces*/
125 tab sex
126 gen gender=(substr(sex,1,1)=="f") /*If the first letter is f, syntax is

```

```

127                                     substr(varname, position, number)* /
128     replace gender=. if sex=="
129     replace gender=. if sex=="u"
130     label var gender "Gender"
131     capture label define gender 0 "Male" 1 "Female" /*Capture allows the program to
132                                                         continue even if there's an error*/
133     label values gender gender
134     tab sex gender, m
135     drop sex
136
137     save "main.dta", replace
138     use "main.dta", clear
139
140     /*Merge Data Sets*/
141     /*Don't need to sort either data set prior to merging*/
142     merge 1:m id using "long.dta" /*need to specify variable on which to merge and the
143                                     extent of duplicates in each file no duplicates is a 1, any duplicates
144                                     is an m, can do 1:1, 1:m or m:1 never merge m:m (it probably doesn't do
145                                     what you think-- see the HELP file), use joinby instead,
146                                     which creates all pairwise combinations of the data based on the linking
147                                     variables specified, which is what the m:m merge sounds like it should
148                                     do.*/
149                                     /*if same varnames will override variables in the using data set*/
150                                     /*generates variable called _merge, indicating how well each
151                                     observation merged*/
152     keep if _merge==3 /*Keep the ones that matched*/
153     drop _merge
154     sort id
155     list in 1/10
156     save "data_long.dta", replace
157
158     /*Or you can merge 1:1 using the wide data set*/
159     use "main.dta", clear
160     merge 1:1 id using "wide.dta"
161     keep if _merge==3
162     drop _merge
163     sort id
164     list in 1/10
165
166     save "data_wide.dta", replace
167     keep id gender birth_d age arr_d1 arr_y1
168
169     savesome                                     ///
170     id gender birth_d age arr_d1 arr_y1 ///
171     if arr_y1==2014                                     ///
172     using "arr_2014.dta", replace /*Saves a portion of the data set*/
173
174     sum arr_y1
175     drop if arr_y1==2014 /*drops that same portion*/
176     sum arr_y1
177
178     *Append essentially stacks one data set on top of another
179     *as opposed to merge which places them side-by-side.
180     append using "arr_2014.dta" /*adds that data set back in*/
181     sum arr_y1
182     clear
183
184
185
186     ***EXERCISES***
187     *(It will be helpful to run the code above 1st,
188     * as the questions below use some of the data files that are created.)
189
190
191     /*Open the main.dta data set. How many
192     duplicates are there on id?*/
193
194     *Answer:
195     use "main.dta", clear

```

```

196     duplicates report id
197
198     /*Merge in the data_long.dta (Hint: type "help merge" into the command line
199     to find out the syntax for this kind of _merge)*/
200
201     merge 1:m id using "data_long.dta"
202
203     /*Keep only those observations in both files, that matched.
204     Save this data set in the temp folder using the name main_long.dta*/
205
206     keep if _merge==3
207     save "main_long.dta", replace
208
209     /*Reshape you data so that there is only one record or observation per id. Save your new
210     data file under the name "main_wide.dta":*/
211     drop _merge
212     sort id arr_d
213     by id: gen arrest_n3=_n
214     reshape wide arr_* age* birth* gender , i(id) j(arrest_n3)
215
216     /*import the csv file called main_names.csv (Hint: you can either search help import
217     or use the interface to determine the code for importing a csv)*/
218
219     import delimited "main_names.csv", clear
220
221     /*How many true duplicates are there? How many duplicates on id? Drop any true duplicates*/
222
223     duplicates report
224     duplicates report id
225     duplicates drop
226
227     /*Remove any true duplicates in terms of name (accounting for different letter casing and
228     any leading or trailing spaces*/
229
230     replace full_name=trim(full_name)
231     replace full_name=lower(full_name)
232     duplicates drop
233
234     /*Try to merge in the main_long.dta data set that you saved in the temp folder*/
235
236     merge 1:m id using "main_long.dta"
237
238     /*Did you get an error message? How can you address the error?*/
239
240     describe id
241     tostring id, replace
242
243     merge 1:m id using "main_long.dta"
244
245     /*Did you get another error message? How can you address the error?*/
246     merge 1:m id using "main_long.dta", gen(new_merge)
247
248
249     capture log close
250     exit
251

```