

```

1  /*****
2  Title: STATA Tutorial 2
3  *** Purpose: Stata setup, basic data manipulation, Mean comparison test, Basic OLS regressions ***
4  ****
5  /
6  clear all
7
8  ****General setup**
9  ****
10 capture log close /* close the log file if any already open. capture: This command tells Stata to
    ignore any error messages and keep going*/
11
12 /*A macro is used as shorthand: you type a short macro name but are actually referring to some
    numerical value or a string of characters. Macros are of two types: local and global
13 1 local: work within the program or do-file in which they are created
14 2 global: work in all programs and do files
15 */
16
17 global computer "/Users/huihuaxie"
18
19 // if you wanna invoke, use $ as a pre-fix
20 global dropbox "$computer/Dropbox"
21 global datapath "$computer/Data"
22 global projectpath "$dropbox/Lectures/EC03211_2021Fall/Tutorials/Tutorial 1"
23
24 global rawdata "$datapath/Lectures"
25 global cleandata "$projectpath/Data"
26
27 global dopath "$projectpath/Do"
28 global result "$projectpath/Results"
29
30 cd "$cleandata"
31
32 ****
33 ***Put your notes in do-file ****
34 ****
35
36 use hlth hi age age2 angrist brooks educ empl famsize inc1 inc2 inc3 inc4 inc5 inc6 inc7 inc8
    marstat nwhite sex perweight using "$rawdata/NHIS2009_tutorial.dta", clear
37 /*use varlist using filename */
38
39 /*
40 /*It is good practice to keep extensive notes within your do-file
41 Thus, when you look back over it you know what you were trying to achieve with each command or
    set of commands
42 You can insert comments in several different ways */
43
44 // Stata will ignore a line if it starts with two consecutive forward slashes "/"
45
46 // cd "D:/Dropbox/lectures/EC03211_2021Fall/STATA/Stata 1"
47
48 /*cd "D:/Dropbox/lectures/EC03211_2021Fall/STATA/Stata 1" */
49
50 // blocking a whole set of commands
51
52 /*
53 rename sex gender /* example of renaming the variable*/
54
55 drop marstat /*drop a variable*/
56
57 keep if _n<=50 /*keep observations1-50 for all variables*/
58 */
59
60 * You can use three consecutive slashes "///" which will result in the rest of the line being
    ignored and the next line added at the end of the current line. This comment is useful for
    splitting a very long line of code over several lines.
61 It works in do-file window, but not in command window.

```

```

62 */
63
64 reg hlth hi age age2 angrist brooks educ empl famsize inc1 inc2 inc3 ///
65 inc4 inc5 inc6 inc7 inc8 marstat nwhite sex, robust
66
67 sum hlth if hi==0
68 local controlmean=r(mean)
69
70 outreg2 using Table1, excel dec(4) addstat(controlmean, `controlmean') replace //arrange
71 regression results into an excel table. "dec(4)" is used to fix decimals, "addstat" is used to
72 add additional statistics
73 * An additional row or column?
74
75 * Three ways to install a new package
76 net search outreg2 //1. search Internet for installable packages
77 ssc install outreg2, all replace //2. install or uninstall packages from the Boston College
78 Statistical Software Components (SSC) archive
79 help outreg2 //3. displays useful information about the how to use specified command or specific
80 topic
81
82 *****
83 **** Macros ****
84 *****
85 // firstly define Lcontrols and Gcontrols.
86 local Lcontrols age age2 educ empl famsize nwhite sex
87 global Gcontrols age age2 educ empl famsize nwhite sex
88
89 reg hlth hi `Lcontrols', robust //symbol` is located next to number 1
90
91 *Compare reg hlth hi `lcontrols', ro
92 // Stata is a case-sensitive, sensitive to Upper and lower cases.
93
94 outreg2 using Table1, excel dec(4) addstat(controlmean, `controlmean')
95 /*
96 . local controls age age2 educ empl famsize nwhite sex
97
98 . reg hlth hi `Lcontrols', robust
99
100 Linear regression
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125

```

	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
hi	.1424578	.0214263	6.65	0.000	.1004604	.1844553
age	-.0151922	.007768	-1.96	0.051	-.0304182	.0000339
age2	-.0000198	.0000918	-0.22	0.829	-.0001998	.0001601
educ	.0589874	.0021703	27.18	0.000	.0547333	.0632415
empl	.2437649	.0204933	11.89	0.000	.2035961	.2839336
famsize	.0093462	.0054408	1.72	0.086	-.0013182	.0200105
nwhite	-.1781165	.017185	-10.36	0.000	-.2118007	-.1444324
sex	.0010022	.0137301	0.07	0.942	-.0259101	.0279145
_cons	3.35254	.1602921	20.92	0.000	3.038352	3.666727

```

116
117
118 . outreg2 using Table1, excel dec(4) addstat(controlmean, `controlmean')
119 */
120
121 sum hlth if hi==0
122 local controlmean=r(mean)
123
124 outreg2 using Table1, excel dec(4) addstat(controlmean, `controlmean')
125

```

```

126 /*Local macros are "private"
127 If you use several programs within a single do-file, you need not worry about whether some other
program has been using local macros with the same names*/
128 reg hlth hi `lcontrols', robust
129
130 /* See result below
131 . reg hlth hi `lcontrols', robust
132
133 Linear regression                                Number of obs      =       18,790
134                                                    F(1, 18788)        =       268.90
135                                                    Prob > F            =       0.0000
136                                                    R-squared           =       0.0159
137                                                    Root MSE           =       .94507
138
139 -----
140          |               Robust
141          |               Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
142 -----+-----
143          | hi |   .3288268   .0200527    16.40   0.000   .2895217   .3681318
144          | _cons |  3.655683   .018636   196.16   0.000   3.619154   3.692211
145 -----
146 */
147
148 // Global macros are "public"
149 /*Gcontrols refers to exactly the same list of variables irrespective of the program that uses
it, global macros are prefixed by the dollar sign: $
150 You should refrain from using global macros when a local macro suffices
151 This is good programming practice as it forces you to define these macro variables explicitly
instead of defining them in some hard-to-find place in your code
152 If you use global macros you should make sure that you define them at the beginning of your code.
153 */
154
155 reg hlth hi $Gcontrols, robust
156
157 sum hlth if hi==0
158 local controlmean=r(mean)
159 outreg2 using Table1, excel dec(4) addstat(controlmean, `controlmean')
160
161 /*. reg hlth hi $Gcontrols, robust
162
163 Linear regression                                Number of obs      =       18,790
164                                                    F(8, 18781)        =       228.29
165                                                    Prob > F            =       0.0000
166                                                    R-squared           =       0.0984
167                                                    Root MSE           =       .90476
168
169 -----
170          |               Robust
171          |               Coef.   Std. Err.      t    P>|t|     [95% Conf. Interval]
172 -----+-----
173          | hi |   .1424578   .0214263     6.65   0.000   .1004604   .1844553
174          | age |  -.0151922   .007768    -1.96   0.051  -.0304182   .0000339
175          | age2 | -.0000198   .0000918    -0.22   0.829  -.0001998   .0001601
176          | educ | .0589874   .0021703    27.18   0.000   .0547333   .0632415
177          | empl | .2437649   .0204933    11.89   0.000   .2035961   .2839336
178          | famsize | .0093462   .0054408     1.72   0.086  -.0013182   .0200105
179          | nwhite | -.1781165   .017185   -10.36   0.000  -.2118007  -.1444324
180          | sex | .0010022   .0137301     0.07   0.942  -.0259101   .0279145
181          | _cons |  3.35254   .1602921    20.92   0.000   3.038352   3.666727
182 -----
183 */
184
185
186 *****
187 ***Organize data***
188 *****
189 //compare egen & gen
190 egen inc_sexmean = mean(inc), by(sex)

```

```

191 *egen age_mean = mean(age), by(sex) remember by followed with parentheses
192 //egen typically creates new variables based on summary measures, such as sum, mean, min and max.
    Use function mean to get mean income for each gender
193
194 egen educ_sexmax = max(educ), by(sex)
195
196 egen inc_count = count(inc)
197
198 egen inc_diff = diff(inc inc1) //An indicator. generate a variable indicating whether variables
    inc and inc1 are different or not
199
200 //sort the data by age first, generate the mean income for each age group
201 bysort age: egen inc_serial = mean(inc)
202 *by age, sort: egen inc_serial2 = mean(inc)
203
204 //numeric or string variables
205 /*
206 Stata stores or formats data in either of two ways-numeric or string. Numeric will store numbers
    while string will store text. Numeric variables are in black/blue color and string variables are
    in red color. String variable can also be used to store numbers, but you will not be able to
    perform numerical analysis on those numbers.
207 */
208 tostring year, replace //change variable to the form of string. either replace or generate
209
210 destring year, gen(year1) //change variable from string to numeric
211 //at the same time, gen a new variable. parentheses
212
213 gen yr=substr(year, 3, 2)
214 * 2021 to 21
215 /*substr: Divide up a variable or to extract part of a variable to create a new one
216 The first term in parentheses is the string variable that you are extracting from
217 The second term (3) is the position of the first character you want to extract
218 The third term (2) is the number of characters to be extracted*/
219 // substr only works for string type not float type. If not tostring first, substr cannot be
    used.
220
221 gen yr2=substr(year,-2,2)
222 * 2021 to 21
223 /*Alternatively, you can select your starting character by counting from the end (2 positions
    from the end instead of 3 positions from the start)*/
224
225
226 collapse (mean) mean_inc=inc (max) max_inc=inc (count) count_inc=inc, by (age)
227                                     *(median) median_inc = inc
228 //Dangerous!! The change could be eternal.
229 * age is the first column
230 //collapse: This command converts the data into a dataset of summary statistics, such as sums,
    means, medians, and so on. And eternally leaves out all original information. This command is
    useful only if you want to aggregate dataset from individual level.
231 // definitely see the results. (browse)
232 // ctrl+shift+s 'save to other' to carefully save the result and protect the original data.
233
234 compress // Different var types take up different sizes of memories. compress attempts to reduce
    the amount of memory used by your data
235
236 /* Example:
237 compress
238 variable mpg was float now byte
239 variable price was long now int
240 variable yenprice was double now long
241 variable weight was double now int
242 variable make was str26 now str17
243 See help compress for detailed explanation
244 */
245
246 save "$cleandata/temp.dta", replace
247 * here we invoke a global defined string again.
248

```