

```
1 *****
2 * Penguin Analysis
3 * Demonstration Do File
4 *****
5
6 * So many projects have the same, or similar, workflow.
7 * DO YOUR THINKING IN CODE!!!
8
9 * have a question ->
10 * get data ->
11 * process and clean data ->
12 * analyze data ->
13 * visualize data ->
14 * make conclusions
15
16 /* do files are useful to preserve
17 a record of your work. They help
18 to keep an audit trail of the
19 decisions that you have made. */
20
21 /* do files thus serve as a way of creating an
22 automated, replicable and documented workflow
23 as well as finding and minimizing errors */
24
25 * A `*` character at the beginning of a line makes that
26 line a comment
27
28 /* You can also use asterisk slash to denote multiple
29 lines of comment */
30
31 *****
32 * get data
33 *****
34
35 * a good workflow habit is to
36 * always--or at least frequently--
37 * work from your raw data.
38
39 * i.e. run your script so you are always--
40 * or at least often--opening your raw data,
41 * cleaning the data, creating new variables,
42 * and then running analyses.
```

```
41
42 clear all // clear the workspace
43
44 * get data from web
45
46 use "penguins.dta", clear
47
48 *****
49 * take a look at the data
50 *****
51
52 * NB if you have a lot of variables, the commands below
  will produce a lot of (too much) output
53
54 * you may need to `describe` or `codebook` specific
  variables
55
56 describe // describe the variables
57
58 codebook // full descriptions of all the variables;
  produces a lot of output
59
60 *****
61 * descriptive statistics
62 *****
63
64 summarize // descriptive statistics for all variables
65
66 summarize body_mass_g // descriptive statistics for this
  variable
67
68 tabulate species // tabulate this categorical variable
69
70 * dtable is a useful new command
71 * for producing tables of descriptive statistics
72 * be sure to denote indicator variables with an `i.`
73
74 dtable culmen_length_mm body_mass_g i.species
75
76 *****
77 * data wrangling
78 *****
```

```
79
80 * find variables of interest
81
82 lookfor mass // look for a variable w a particular keyword
83
84 * sometimes it is useful to `keep` only the variables in
85   which you have an interest
86 * to reduce the size of the data set
87
88 * recode variables
89
90 generate big_penguin = body_mass_g > 4000 // create a
91   big penguin variable
92
93 tabulate big_penguin
94
95 *****
96 * ANOVA
97 *****
98
99 oneway body_mass_g species, tabulate
100
101 *****
102 * regression
103 *****
104
105 regress culmen_length_mm body_mass_g
106
107 est store M1 // store regression estimates
108
109 regress culmen_length_mm body_mass_g i.species
110
111 est store M2 // store regression estimates
112
113 * /// indicates that a command spans multiple lines
114
115 etable, estimates(M1 M2) /// nicely formatted table of
116   regression estimates
117 cstat(_r_b) /// beta's only
118 showstars showstarsnote // show stars and note
119
120 *****
```

```
118 * graph
119 *****
120
121 graph bar body_mass_g, over(species) // bar graph
122
123 twoway scatter culmen_length_mm body_mass_g // scatterplot
124
125
126
127
128
129
```