

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

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

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

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