

Stata .do Files

Why “.do files”?

Up until this point in your Stata use, most of you have probably used Stata as “dynamic”, “call and response” data analysis. In other words, for every command you enter at the command window you get some output in the results window, or a graph window opens. You can archive all of this resulting output in a log file – however, if you were interested in duplicating your analysis efforts at a different time, you would need to start from ground zero, re-entering each of the commands in the command window, one by one.

If this need to duplicate a previously used set of commands ever arises, wouldn’t it be nice to have all of the commands archived so that you just run them again quickly, without having to go re-enter things on a command by command basis?

Luckily, there is tool for archiving commands and re-running them: it is called a Stata .do file. Before we get into the details of the *.do file* itself, let’s talk about some situations where you may want to have a ready to go archive of your commands:

1. Sometimes, an analysis you have prepared for one set of data would be useful and appropriate to apply to other data sets as well. For example, if you were testing a hypothesis regarding factors linked to school performance in Boston high schools, you may be interested in testing this same set of hypotheses to schools in 3 other cities as well.
2. Sometime, especially in the exploratory stage, it is desirable to run a series of basic tabulations and univariate summary statistics on all variables in a dataset. (For example, all items collected in a survey). Usually, the exploration stage begins prior to the cleaning and collection of the entire data set, and therefore it may be desirable to re-run the same analysis as the data set gets updated.

What is a .do file?

It is essentially just a text file, consisting of a list of Stata commands, and (optional) user comments to document the commands.

A *.do file* can be

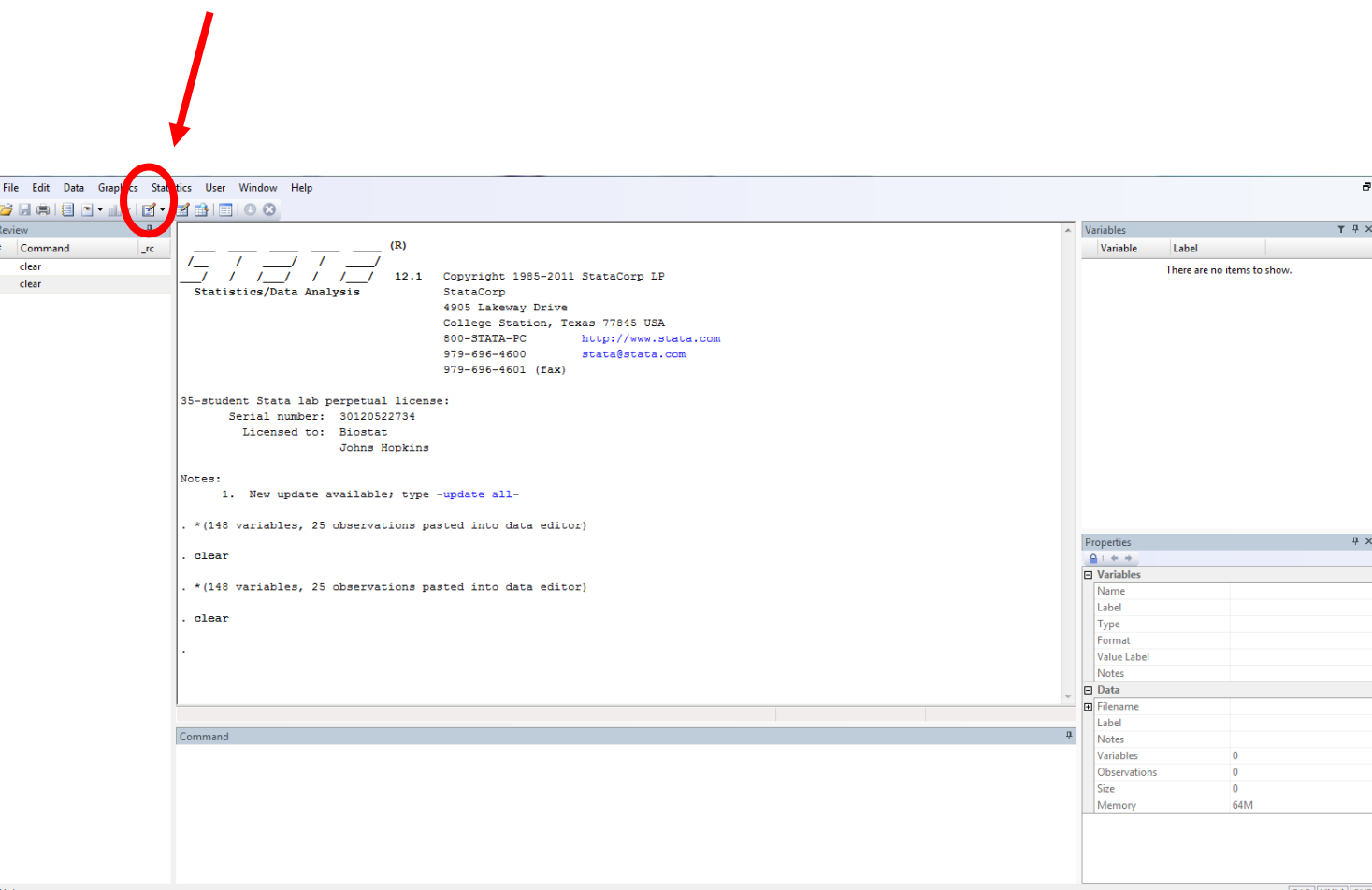
- 1) created from scratch, with the analyst hand typing each command into a text file, or
- 2) constructed by cutting and pasting command statements from a *.log file* or
- 3) created by saving the commands from the review window.

Among these options, the last (number 3) is the most efficient!

First Approach: How do I get to a text editor to create a do-file from scratch?

Stata has a *.do file* editor built right in – to access this “blank canvas”, you just need click the *.do file* editor icon on the command icon bar – the *.do file* editor icon looks like a ring notebook. You can type all of your commands into this editor, and then save the resulting *.do file*.

.do file editor icon



Third, one then “right-clicks” on the highlighted commands and chooses “Send to Do-File Editor”

The screenshot displays the Stata software interface. The main window shows the 'Review' tab for a regression model. The model summary includes the following statistics:

	Model	698684.887	2	349342.443	F(2, 4288) = 41116.19
Residual	36432.8575	4288	8.49646863	Prob > F = 0.0000	
				R-squared = 0.9504	
				Adj R-squared = 0.9504	
				Root MSE = 2.9149	
Total	735117.744	4290	171.356118		

The coefficient table for the regression model is as follows:

	height	weight	age	_cons
Coef.	2.843894	.2418318	.0056806	44.81849
Std. Err.	.0313833	.0056806	.1811823	
t	90.62	42.57	247.37	
P> t	0.000	0.000	0.000	
[95% Conf. Interval]	2.782367 2.905422	.2306949 .2529687	44.46328 45.17371	

A context menu is open over the command window, with the option 'Send to Do-file Editor' highlighted in red. The command window shows the following commands:

```
. predict resid, resid
option resid not allowed
r(198);

. predict resid, resid
option resid not allowed
r(198);

. predict resid yhat, title(Residuals from wt=ht age) yline(0)
option resid yhat, yline(0) title(Residuals from wt=ht age)
```

The Do-file Editor window at the bottom shows the command:

```
predict resid, resid
```

The Variables window on the right lists the variables: id, age, sex, weight, height, armcirc, resid, yhat, and fitted values. The Properties window on the right shows the file name 'nepali_anthropometry.dta' and other file details.

The commands are then copied into a do-file which can be edited and saved:

The screenshot displays the Stata software interface. A central window titled "Do-file Editor: Untitled1.do" contains the following code, with lines 1 through 6 circled in red:

```

1 use "C:\Users\JMCGRAD\Desktop\JMCGRAD\Statistical Reasoning 1 Makeover 2013\lecture
2 regress height weight age
3 predict resid, resi
4 predict yhat
5 twoway (scatter resid yhat), title(Residuals from wt=ht age) yline(0)
6 twoway (scatter resid yhat), yline(0) title(Residuals from wt=ht age)
7

```

Below the editor, the Command window shows the command `predict resid, resid`. To the right, the Variables list shows variables: `id`, `age`, `sex`, `weight`, `height`, `armcirc`, `resid` (labeled Residuals), and `yhat` (labeled Fitted values). The Properties panel on the far right shows details for the variable `id`, including its name, label, type (long), format (%12.0g), and value label.

What does a *.do* file look like?

It is really just a file with a list of commands!

Great, now how do I get this .do file to run?

This icon is the “run current file” icon: click on this to run the .do file.



The screenshot shows the Stata software interface. The main window is the 'Do-file Editor - Untitled1.do', which contains the following code:

```

1 use "C:\Users\JMCGRAD\Desktop\JMCGRAD\Statistical Reasoning 1 Makeover 2013\lecture
2 regress height weight age
3 predict resid, resi
4 predict yhat
5 twoway (scatter resid yhat), title(Residuals from wt=ht age) yline(0)
6 twoway (scatter resid yhat), yline(0) title(Residuals from wt=ht age)
7

```

The 'Run' icon (a blue square with a white play button) in the toolbar is circled in red. Below the Do-file Editor, the Command window shows the command:

```

predict resid, resid

```

On the right side of the interface, there are two panels: 'Variables' and 'Properties'. The 'Variables' panel lists variables: id, age, sex, weight, height, armcirc, resid (Residuals), and yhat (Fitted values). The 'Properties' panel shows details for the variable 'id', including its name, label, type (long), format (%12.0g), value label, and notes. It also shows data summary information: Filename (nepali_anthropometry.dta), Label, Notes, Variables (8), Observations (4,386), Size (111.36K), and Memory (64M).

A less exiting, but equally effective way to run your do file is to type:

run filename

at the Stata command line. (see upcoming section about directory paths)

What's the Catch? The Directory Path Issue!

The only potential “sticky issue” when creating a *.do file* is concerned with directory path recognition issues. Many times when doing a particular analysis, your data files are stored in directories specific to your project, and it would be desirable to place any resulting output (graphs, log files) in the same directory path. Without being specifically told, Stata has no idea where to find your data, or to place your output files. Stata has its own “default” directory, a place where it looks to find things, and posts output unless being specifically told to do otherwise.

For example, suppose you have your data stored in a file called “hw2data” in folder called “bio623” on your disk drive A.

The full name of your file is:

A: \ bio622\hw2data.dta

In order to access this data in Stata, you may be tempted to type:

use “hw2data”

However, chances are good that Stata is not looking in your “A:\bio622” folder for this file. *(in fact, if you wish to know what Stata uses as it's “go to”, default director, look in lower right hand corner of the Stata window to see the default working directory)* You will probably get an error message which says “file hw2data not found”.

One way to remedy this problem is to modify your statement to:

use “A:\bio622\hw2data”

That is to tell Stata EXACTLY where to find the data file.

Another option is to change the default directory where Stata looks for files if a directory is not specified: This can be done using the “cd” command:

cd “A:\bio622\”

Now when you type “use hw2data”, Stata will be looking for the file in the “A:\bio622” directory. Similarly, if you wanted to begin a log file called “mylog”, and typed

log using mylog

Stata would open the file in the A:\bio622 directory.

Stata would continue using A:\bio622 as its default until you change it at the command line, or log out of Stata.

How do I make sure that Stata knows where to look for my files, and direct any output I want to save?

There are 3 options:

- 1) *Every time you reference a file, include the full name (directory path and all)*

use A:\studentname\dofile.dta"
log using "C:\myfiles\term2\bio622\hw1.log"

This option is really only appropriate if you are using files from multiple directories in one analysis – if there is one directory serves as your “base of operations” this approach would be very inefficient (and lead to unnecessary extra typing)

- 2) *Prior to running your .do file, change your default directory at the Stata command line.*

Somewhat useful if you are doing your work from one directory – a pain, though, because you must consciously remember to do this every time you want to run your .do file.

- 3) *Include a change directory command at the beginning of your .do file. This will automatically change the default directory path EVERY TIME you run the .do file, on any computer.*

The best, most efficient option if you are running everything out of one directory (getting your data etc...).

Of course, you can “mix and match” the above options: for example if you are doing most of your work from the directory “A:\bio622\hw1” but want to save your log file in “C:\mystuff\log1.log”, you could change your default directory to “A:\bio622\hw1”, and then explicitly type the full filename when opening the .log file.