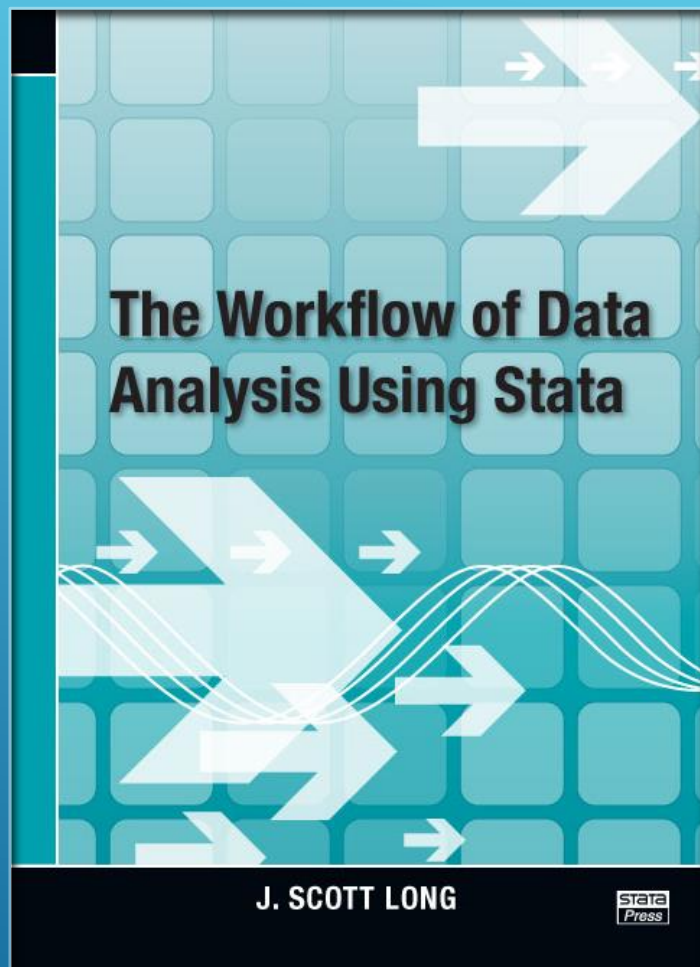# RESEARCH WORKFLOW USING STATA

How to Be an Effective Researcher

CCPR Workshop

# THE WORKFLOW OF DATA ANALYSIS USING STATA

J. Scott Long

# THREE WAYS TO EXECUTE COMMANDS

# THREE WAYS TO EXECUTE COMMANDS

# THREE WAYS TO EXECUTE COMMANDS

# DO-FILES

# ADVANTAGES OF USING DO-FILES

1. You have a record of the commands you ran.

   You can rerun them in the future to replicate your results
   You can quickly modify your code

2. You can use the features of the text editor

   i.e. copy/paste, find and replace, select all

# TWO RULES OF USING DO-FILES

1. Do-files must be robust

   Robust do-files produce exactly the same result when run at a later time or on another computer

2. Do-files must be legible

   Legible do-files are documented and formatted so that it is easy to understand what is being done

# MAKING DO-FILES ROBUST

Robust do-files are self-contained

# MAKING DO-FILES ROBUST

Exclude directory information

```
1
2    ********* load dataset using directory location    **************
3    use "C:\Users\ahicks\Documents\Stata_data\lafans\parent1.dta", clear
4
5    ********* load dataset without directory location    **************
6    cd C:\Users\ahicks\Documents\Stata_data\lafans
7    use parent1.dta
```

# MAKING DO-FILES ROBUST

Use version control

```
. version 14
this is version 13.0 of Stata; it cannot run version 14.0 programs
    You can purchase the latest version of Stata by visiting http://www.stata.com.
r(9);
```

Include seeds for random numbers

```
. set seed 90049
```

# MAKING DO-FILES LEGIBLE

Legible do-files are internally documented and formatted

Use comments!!

```
1  **************************************************************
2  *           Use Comments!                                    *
3  **************************************************************
4
5  * Stata treats the entire line as a comment if the line starts with a *
6
7  gen wave1date=date(pdate, "MDY")
8  gen hsattend =3
9  *replace hsattend =1 if (pg40==1)
10 replace hsattend =2 if (pg68==1)
11
12
13
14 /* You can create comments across multiple lines by using an opening
15    '/*' and a closing'*/'.  You can also use this type of comment to
16    temporarly stop entire sections of code from being executed  */
17
18 gen wave1date=date(pdate, "MDY")
19 gen hsattend =3
20 /*
21 replace hsattend =1 if (pg40==1)
22 replace hsattend =2 if (pg68==1)
23 label var hsattend "current or ever headstart attend"
24 */
25 recode hsattend (1 2=1) (3=0)
26 tab hsattend
27
28
29
30 // Any thing that comes after a double slash is treat as a comment
31
32 logit lfp wc hc // this analysis only includes education, add wages next
33
```

# MAKING DO-FILES LEGIBLE

Use alignment and indentation

```
*******************************************************************
* Use Alignment and indentation                                  *
*******************************************************************


* This code:
logit foreign price mpg trunk weight length turn displacement gear_ratio



* Is the same as this code:
logit foreign    price         mpg          trunk ///
                 weight        length       turn ///
                 displace      gear_ratio

* Use short lines
mi estimate, post nois saving(lwi_sub5_1, replace): regress W2_lwi_ss concdis_average childmale birth
```

# MAKING DO-FILES LEGIBLE

Limit your abbreviations

```
****************************************************************
* Be careful with abbreviations                               *
****************************************************************
* three commands that give the same results:
summarize displacement
su d
sum displace

* What on earth does this mean?
l ma p f in 1/3
```

# SAVING YOUR SESSION TO A LOG