

Intro to Stata

James Ng, PhD
james.ng@nd.edu

what is Stata?

CENTER for
DIGITAL
SCHOLARSHIP

- statistical software package
- created in 1985 by economists

why bother when I can use Excel?

CENTER for
DIGITAL
SCHOLARSHIP

- documentation and reproducibility of data and results
- eases revision, collaboration
- integrates nicely with Word, Excel, LaTeX
- time and energy saver

steps in data analysis

CENTER for
DIGITAL
SCHOLARSHIP

- locate data
- load data into software package
- manipulate as needed (*bulk of your time*)
- analyze

"data"

- a set of numbers and/or text describing specific phenomena
 - economy, weather, traffic, pollution levels
- in social sciences, always rectangular:
 - columns contain "variables"
 - rows contain "observations"

example

Country	Population	GDP per capita (in USD)
USA	300,000,000	40,000
Malaysia	25,000,000	12,000
China	1,600,000,000	6,000
Vatican City	2,000	100,000

Data Editor (Browse) - [Untitled]

File Edit View Data Tools

Country[1] USA

	Country	Population	GDP_PC
1	USA	3.00e+08	40000
2	Malaysia	2.50e+07	12000
3	china	1.60e+09	6000
4	Vatican City	2000	100000

today's agenda

CENTER for
DIGITAL
SCHOLARSHIP

- how to load data and do basic manipulations and analysis
- on two widely-used, publicly-available datasets:
 - National Health Interview Survey (NHIS)
 - General Social Survey (GSS)

Stata environment

CENTER for
DIGITAL
SCHOLARSHIP

The screenshot shows the Stata/IC 12.0 software interface. The main window displays the Stata logo, version 12.0, and copyright information. Below this, it shows the 25-user Stata network perpetual license for the University of Notre Dame. The interface includes a menu bar (File, Edit, Data, Graphics, Statistics, User, Window, Help), a toolbar, and a command window at the bottom. Annotations with red callout boxes and black-bordered letters identify key components: 'History' (C) points to the left sidebar; 'Variables' (B) points to the top right sidebar; 'Results' (D) points to the main command window; and 'Command line interface' (A) points to the bottom command window.

Stata/IC 12.0 - [Results]

File Edit Data Graphics Statistics User Window Help

Review

Command _rc

There are no items to show.

Statistics/Data Analysis (R)

12.0 Copyright 1985-2011 StataCorp LP
StataCorp
4905 Lakeway Drive
College Station, Texas 77845 USA
800-STATA-PC http://www.stata.com
979-696-4600 stata@stata.com
979-696-4601 (fax)

25-user Stata network perpetual license:
Serial number: 30120547268
Licensed to: ND User
University of Notre Dame

Notes:
running C:\stata12\profile.do ...

Results

Command

Variables

Variable Label

There are no items to show.

Properties

Variables

Name
Label
Type
Format
Value Label
Notes

Data

Filename
Label
Notes
Variables 0
Observations 0
Size 0
Memory 64M

CAP NUM OVR

ways to use Stata

CENTER for
DIGITAL
SCHOLARSHIP

- point & click ← OK to start
- command line interface ← good
← today
- batch file (called a “do-file”) ← best

keeping records

CENTER for
DIGITAL
SCHOLARSHIP

- good practice:
keep a **log file** at start of each session
- Stata command:
log using anyfilename.log, **text**
replace

loading data into Stata [1]

- there are many ways

Command	File Type	File Extension
use	Stata format	.dta (always)
infix	Fixed-format ASCII	.dat, .raw, .fix, or simply nothing
infile (version 1)	Text-delimited ASCII	
infile (version 2)	Fixed-format ASCII, with a “dictionary”	
import delimited	Text-delimited ASCII	
import excel	Excel	.xls, .xlsx

loading data into Stata [1]

CENTER for
DIGITAL
SCHOLARSHIP

- today:
 - load a Stata-format dataset
 - load an ASCII dataset

loading data into Stata [2]

CENTER for
DIGITAL
SCHOLARSHIP

- before you start, must know: file path, file name

loading data into Stata [3]

- example 1: GSS data
- reading Stata-format data is trivial
- Stata command: use

```
use N:\Public\GSS\GSS2012.dta, clear
```

- good practice:

```
cd N:\Public\GSS\  
use GSS2012, clear
```

loading data into Stata [4]

CENTER for
DIGITAL
SCHOLARSHIP

- example 2: NHIS data
- fixed-format ASCII file

http://www.cdc.gov/nchs/nhis/nhis_2012_data_release.htm

- Stata command: `infix`
- script to load data already written by data provider – really helpful!
 - how to use it?

combining datasets

Merging

- adding variables to existing observations
- similar to SQL join, SPSS match files

id	sex		id	age		id	sex	age
001	M	+	001	21	=	001	M	21
002	F		002	23		002	F	23
data1.dta			data2.dta					

use data1, clear
merge 1:1 id using data2

Appending

- adding observations to existing variables
- similar to SPSS add files

id	sex		id	sex		id	sex
001	M	+	003	F	=	001	M
002	F		004	M		002	F
data1.dta			data3.dta			003	F
						004	M

use data1, clear
append using data3

inspecting your data [1]

CENTER for
DIGITAL
SCHOLARSHIP

- read the manual / codebook / user guide

- some essential commands:

sort

order

browse

describe

lookfor

sum

tab

selecting variables

```
keep id happy abpoor age race sex health1 region
```

- see also: drop

- save your work data in a new file:

```
save temp_gss2012
```

- or overwrite existing file:

```
save temp_gss2012, replace
```

- be careful not to unintentionally overwrite dataset

creating a new variable [1]

CENTER for
DIGITAL
SCHOLARSHIP

- create a variable to indicate unhappiness based on an existing variable

- don't be misled by "value labels"

```
tab happy, nolabel
```

- watch out for missing values!

```
tab happy, nolabel missing
```

creating a new variable [2]

CENTER for
DIGITAL
SCHOLARSHIP

- here's how

```
gen unhappy = .
```

```
replace unhappy = 1 if happy == 3
```

```
replace unhappy = 0 if happy == 1 | happy == 2
```

- cross-check:

- `tab unhappy happy, nolabel missing`

creating a new variable [3]

CENTER for
DIGITAL
SCHOLARSHIP

- good practice: label all variables

```
label var unhappy "Is respondent unhappy? 1=yes 0-  
no"
```

creating a new variable [3]

CENTER for
DIGITAL
SCHOLARSHIP

- create a variable indicating whether a person feels poor

```
gen poor = abpoor==1
```

```
replace poor = . if missing(abpoor)
```

```
label var poor "Does respondent feel poor? 1=yes 0-  
no"
```

creating a new variable [4]

- you can also label a variable's values
- let's label values of unhappy
- 2-step process:

- define labels for variable's values:

```
label define labels_for_unhappy 0 "happy" 1 "unhappy"
```

- assign value labels to variable:

```
label values unhappy labels_for_unhappy
```


basic analysis [1]

CENTER for
DIGITAL
SCHOLARSHIP

- descriptive statistics

```
sum
```

```
sum age
```

```
tab race
```

```
tab race, nolabel
```

```
tab poor
```

```
tab unhappy if race==1
```

```
tab unhappy poor
```

```
tab unhappy poor, row column
```

basic analysis [2]

CENTER for
DIGITAL
SCHOLARSHIP

- distribution of a variable

histogram age, normal

- comparison of means

ttest unhappy, by (poor)

basic analysis [3]

- what is the relationship between poverty and unhappiness?

```
corr unhappy poor
```

```
reg unhappy poor
```

- what is this relationship controlling for some other factors?

```
recode sex (2=0), gen(male)
```

```
xi: reg unhappy poor male age i.health1, cl(region)
```

- does it vary by gender?

```
xi: reg unhappy i.poor*male age i.health1, cl(region)
```

```
xi: reg unhappy poor age i.health1 if male==1, cl(region)
```

```
xi: reg unhappy poor age i.health1 if male==0, cl(region)
```

basic analysis [4]

- how did average happiness change over time?

- use data compiled across years

```
use combined1972_2012, clear
```

```
browse
```

```
collapse (mean) ave_unhappiness=unhappy [pw=wtssall], by(year)
```

```
label var ave_unhappiness "fraction of respondents who felt unhappy"
```

- we can now graph it:

```
scatter ave_unhappiness year, xlabel(1972 1982 1991 2002 2012, grid)
```

- map Census regions according to level of unhappiness
- Command: `spmap`
- not part of basic installation; download and install from Stata server in one easy step:

```
ssc install spmap
```


using a “do-file”

- send commands to Stata through a batch file (.do)
 - “do-file”
- all commands in this session can be found in a do-file available on [Box](https://notredame.box.com/s/vs4aq0x64ovdk4zsoat6):
<https://notredame.box.com/s/vs4aq0x64ovdk4zsoat6>
- Stata reads each line as an executable statement
 - ignores lines beginning with an asterisk, * ← documentation, good practice!

if you get stuck

- Stata has an extensive internal help system

- need help with how to load data?

help loading data

- need help with `regress` command?

help regress

- WWW is your friend

- <http://www.ats.ucla.edu/stat/stata/>

- Google

ending your session

CENTER for
DIGITAL
SCHOLARSHIP

```
log close
```

```
Exit
```

or simply close Stata with your mouse

accessing workshop materials

CENTER for
DIGITAL
SCHOLARSHIP

- This PowerPoint is on CDS website:
 - <http://library.nd.edu/cds/workshops.shtml#DataAnalysis>
- Stata datasets and do-files are on Box:
 - <https://notredame.box.com/s/vs4aq0x64ovdk4zsoat6>

other resources on campus

CENTER for
DIGITAL
SCHOLARSHIP

- Center for Social Research: <http://csr.nd.edu>

rate this workshop

CENTER for
DIGITAL
SCHOLARSHIP

<http://library.nd.edu/cds/WorkshopFeedbackForm.shtml>