

(FGLS)
Chapter 6

6.2.1

6.2.2

6.2.3

6.3

(LDV)
Chapter 10

10.1.2

10.1.3

10.2

10.3

10.4

10.5.1

(Dummy)
Chapter 7

7.1.2

7.2

7.4.1

AB

B.1.1

B.1.2

B.2

B.3

B.4

B.5

(IX)
Chapter 8

8.2

8.5

8.6

8.7.4

8.8

8.10

8.A

8.B.

CREATE
RECIPE

ALGORITHM
of Data Sc.

STATA
CHEATSHEETS

Start! you have data sets

- merging
- dropping vars
- generating vars
- labelling
- new refined data set

use .data, clear (1)

describe (2) summarize (3)

- N (largest observation)
- n (current)] → use

generate, replace, []
(new var) (modifies)
Transform. generate $x = y \times 100$
 $x = \frac{y}{z}$
sum x

gsort, sort (4)

ascending 1

given ascending 2

replace $x = \frac{y}{z} \times 100$

gsort y - z.
(descending)

gsort y + z
(asc.)

(sort)

list
by (5)

1/5 first 5
-5/1 last 5

sort by
list x y z in 1/5

gsort - z
 sort z list } gsort - z
 list x y z in 1/5

if ⑥ } gen z = y if x > 5000
 sort p
 list p x z, sep(0)

missing val ⑦ - ' .' (say failed cond's)

Sum z if x > 5000.

⑧ generate dummies

gen z = 0

replace z = 1 if x <= 5000

gen y = 0

replace y = 1 if x > 5000

list p x z y)

CAREFUL

gen y = 1 if x > 5000
 missing created
 not 0.

⑨ missing val as largest ~~max~~ num wher.

gen y = (x <= 5000) if x < . }
gen z = (x > 5000) if x < . }

TODAY

listen to ma'am recordings

Summary Statistics

Stata Baum IHDS practice play

Variables

if for subset of data analysis

} sum ~~for~~ r if x == 1

Varlist and if { Qualifiers

by x , sort: sum p, q, r } equivalent
sum p, q, r if $x == 1$ } (10)

by varlist: - for discrete categories

$x = \{ NE, SE \}$

by x , sort: sum p, q, r
then for each x value } automatic

sum p, q, r if $x == "NE"$

(2) log() - (11) (not by varlist:)

tabstat $p, \log(x)$ statistics (N mean
SD min max)

by() eg by(region) table with sum
for each region

byvarlist: - repeats command for each
value of region.

EXAMPLE

x and y to one dummy $z = 1$ for
 $x = 1$ and 2 for
 $y = 1$

$$\rightarrow \text{gen } z = x + 2 * y$$

by p z, sort: sum a b c

label (12)

label data "abc"

label variable x "population"

desc x

13) drop and keep

drop varlist if (condn, range)

keep varlist.

14) rename x y
(old) (new)

MB01
MB02

renprefix MB readability

15) Same xyz.dta, replace

16) outsheet
outreg (cs*)

DATA TRANSFORMS

Start from the first do files.

list all do files.

How data is transformed.

keep if $x == 1$

keep p q r s

by sort varlist : egen $z = \max(x)$

* Advanced [by sort egen]

Say create a var = 1 for just one obs in each ~~ex~~ var

by sort z : gen $y1 = (_n == 1)$ var

egen $y2 = \text{tag}(z)$

after new var (say using egen) - label
and drop duplicates.

(13) generate dummy var z
gen z = 0

replace z = 1 if x == 1 | y == 13

(14) merge m:1 using 2006a
matrix
drop _merge

(15) collapse(mean) [weight] by (varlist)
merge^{1:1} using gini (for ids)

we can use a gen for say max, min and use to get avg.

$$\begin{aligned} \text{gen } a &= \min(\text{inc}) \\ \text{gen } b &= \max(\text{inc}) \end{aligned}$$

So that create ordering.

gen $x = 1$ if ϕ

- N (for count) (16)

$$\begin{aligned} \text{inc} &= \text{inc} + b \\ &\text{avg} \end{aligned}$$

gen people-count = $-N$

(17) bys varlist : gen peoplecount = $-N$

bys varlist newvar : gen percent = $\left(\frac{-N}{\text{peoplecount}} \times 100 \right)$

(1HDS rank 1 (low) 2 (below avg) 3 (above avg))

[Indices Now] - index
govt index
post index

STATISTICAL ANALYSIS using STATA

Cheatsheet:

- use dataset
- `_N` (highest)
- `_n` (current)
- generate $x = \frac{y}{z}$
- summarize x (continuous)
- tabulate x (categorical)
- replace $x = \frac{x}{100}$
- sort $x - y$
- list $x\ y\ z$, sep by (x)

Q1) list 5 OS states w/ lowest popn
out of all $i, P(i)$

→ sort pop

→ list state region pop in 1/5

Q2) list 5 states w/ highest

→ gsort - pop # descending

→ list state region pop in 1/5

→ gen $X_2 = x$ if $y > 100$

* '01' - missing

Q3) Create a dummy off a
continuous.

→ gen $s = 0$

→ replace $s = 1$ if $pop \leq 5000$

→ gen $L = 0$

→ replace $L = 1$ if $pop > 5000$

g) why not ^{gen $s = 0$} gen $s = 1$ if $pop \leq 5000$
gen $L = 0$

A) for $pop > 5000$ s would be
set missing and not 0.

* (•) missing treated as ∞
∴ all values $< \bullet = \text{True!}$

ALTERNATE

gen $s = (pop \leq 5000)$ if $pop < \bullet$

gen $L = (pop > 5000)$ if $pop < \bullet$

→ replace new = 5 if old == 2.

→ recode old (2 = 5), gen(new)

STATA CHEATSHEETS

I SHORTCUTS

- 1) F2 - Describe data
- 2) Ctrl + 8 - open data editor
- 3) Ctrl D - highlighted execution of do

II SET UP

- pwd - dir
- cd PATH - dir *.dta
- log using FILE.txt, replace
- SSC install PACKAGE

IV IMPORT DATA

- use FILE.dta, clear
- import excel FILE.xlsx, *
* /sheet("Sheet1") cellrange(A2:H11)
first row

- import delimited FILE.csv, /*
*/ rowrange(2:11) colrange(1:8) varnames(

IV SYNTAX - Command.

[by var:] command [var2] [=exp] [ifexp]
[in range] [weight] [using FILE] [go/four]

foreg

" bysort RO4: Summarize price if foreign
== 0 & price <= 9000, detail "

Analysis

- bysort or just by - Application of
cmd across each
unique comp of
var in var1 (vector)

[var2] - Vector to apply command (f'n)
to

[=exp] - Same output as new variable

V LOGICAL OPERATORS

- +, -, *, /, ^
- & ! | == != < > =

VI DATA TYPES

missing data (no data)

true/false (byte)

String

int long float double

- gen var1 = string(var2) → say var1 = "1"

- to string var1, gen(var2)

decode var1, gen(var2)

VI DATA EXPLORE

- describe
- Count "count if var > 500"
- Codebook var gives overview
- Summarize or Sum Var
- inspect var (histogram)
- histogram var1, var2

missing val treated as largest number (positive) thus to exclude ask whether value < "."

levelsof var1 : unique vals of var1

browse

clist (compactify)

display (di) var1[4]

- list var1 var2 if var2 > 10,000
& !missing(var2)