

Data Processing

with Stata

For more info, see Stata's reference manual (stata.com)

Useful Shortcuts

- F2** — keyboard buttons **Ctrl** + **9** open a new do-file describe data
- Ctrl** + **8** open the data editor highlight text in do-file, then **ctrl** + **d** executes it in the command line
- delete** delete data in memory

AT COMMAND PROMPT

- PgUp** **PgDn** scroll through previous commands
- Tab** autocompletes variable name after typing part
- cls** clear the console (where results are displayed)

Set up

- pwd** print current (working) directory
- cd "C:\Program Files\Stata16"** change working directory
- dir** display filenames in working directory
- dir *.dta** List all Stata data in working directory

- capture log close** close the log on any existing do-files use "capture" or "cap"
- log using "myDoFile.txt", replace** create a new log file to record your work and results
- search** find the package mdesc packages contain extra commands that expand Stata's toolkit
- ssc install mdesc** install the package mdesc; needs to be done once

Import Data

- sysuse auto, clear** load system data (auto data) for many examples, we use the auto dataset.
- use "yourStataFile.dta", clear** load a dataset from the current directory frequently used commands are highlighted in yellow
- import excel "yourSpreadsheet.xlsx", /*** ***/ sheet("Sheet1") cellrange(A2:H1) firstrow**
- import delimited "yourFile.csv", /*** ***/ rowrange(2:11) colrange(1:8) varnames(2)**
- import sas "yourSASfile.sas7bdat", bcat("value labels file")**
- import spss "yourSPSSfile.sav"** see help import for more options
- webuse set "https://github.com/GeoCenter/StataTraining/raw/master/Day2/Data"**
- webuse "wb_indicators_long"** set web-based directory and load data from the web

Basic Syntax

All Stata commands have the same format (syntax):



To find out more about any command—like what options it takes—type **help command**

Basic Data Operations

Arithmetic

- +** add (numbers) & combine (strings)
- subtract
- *** multiply
- /** divide
- ^** raise to a power

Logic

- ==** tests if something is equal
- =** assigns a value to a variable
- !=** not equal
- <** less than
- <=** less than or equal to
- >** greater than
- >=** greater or equal to

if foreign != 1 & price >= 10000

make	foreign	price
Chevy_Colt	0	3,984
Buick_Riviera	0	10,372
Honda_Civic	1	4,499
Volvo_260	1	11,995

Explore Data

VIEW DATA ORGANIZATION

- describe** make price display variable type, format, and any value/variable labels
- count** number of rows (observations) can be combined with logic
- count if** price > 5000
- inspect** mpkg show histogram of data and number of missing or zero observations
- lookfor** "in." search for variable types, variable name, or variable label
- isid** mpkg check if mpkg uniquely identifies the data
- histogram** mpkg, frequency plot a histogram of the distribution of a variable

BROWSE OBSERVATIONS WITHIN THE DATA

- browse** or **Ctrl** + **8** open the data editor
- list** make price if price > 10000 & **missing**(price)
- display** price[4] display the 4th observation in price; only works on single values
- gsort** price mpkg (ascending) **gsort** -price -mpkg (descending)
- duplicates report** sort in order, first by price then miles per gallon
- levelsof** rep78 finds all duplicate values in each variable verify truth of claim **assert** price!=.
- display the unique values for rep78

Change Data Types

Stata has 6 data types, and data can also be missing:

- no data** true/false
- missing** words
- byte** numbers
- string** int long float double

To convert between numbers & strings:

- gen** foreignString = string(foreign)
- tostring** foreign, gen(foreignString)
- dencode** foreign, gen(foreignString)
- gen** foreignNumeric = real(foreignString)
- tostring** foreignString, gen(foreignNumeric)
- dencode** foreignString, gen(foreignNumeric)
- foreign**

recast double mpkg generic way to convert between types

Summarize Data

- include missing values create binary variable for every rep78
- tabulate** rep78, **mi** **gen**(repairRecord) one-way table: number of rows with each value of rep78
- tabulate** rep78 foreign, **mi** two-way table: cross-tabulate number of observations for each combination of rep78 and foreign
- by**sort rep78: **tabulate** foreign for each value of rep78, apply the command tabulate foreign
- tabstat** price weight mpkg, **by**(foreign) **stat**(mean sd n) create compact table of summary statistics displays stats formats numbers for all data
- table** foreign, **contents**(mean price sd price) **f**(%9.2fc) **row** create a flexible table of summary statistics
- collapse** (mean) price (max) mpkg, **by**(foreign) - replaces data calculate mean price & max mpkg by car type (foreign)

Create New Variables

- generate** mpgSq = mpg^2 **gen** byte lowPr = price < 4000 create a new variable. Useful also for creating binary variables based on a condition (**generate** byte)
- generate** id = _n **by**sort rep78: **gen** repairIdx = _n _n creates a running index of observations in a group
- generate** totRows = _N **by**sort rep78: **gen** repairTot = _N _N creates a running count of the total observations per group
- ptile** mpg Quartile = mpg, **nq** = 4 create quartiles of the mpg data
- egen** meanPrice = **mean**(price), **by**(foreign) calculate mean price for each group in foreign see help egen for more options