

Some Stata Commands

Last modified: January 2, 2006 9:51AM

General Plotting Commands

1. Plot a histogram of a variable:
`histogram vname`
2. Plot a histogram of a variable using frequencies:
`histogram vname, freq`
`histogram vname, bin(xx) norm`
where `xx` is the number of bins.
3. Plot a boxplot of a variable:
`graph box vname`
4. Plot side-by-side box plots for one variable (`vone`) by categories of another variable `vtwo`. (`vtwo` should be categorical):
`graph box vone, over(vtwo)`
5. A scatter plot of two variables:
`scatter vone vtwo`
6. A matrix of scatter plots for three variables:
`graph matrix vone vtwo vthree`
7. A scatter plot of two variables with the values of a third variable used in place of points on the graph (`vthree` might contain numerical values or indicate categories, such as male ("m") and female ("f")):
`scatter vone vtwo, symbol([vthree])`
8. Normal quantile plot:
`qnorm vname`

General commands

1. To compute means and standard deviations of all variables:
`summarize`
or, using an abbreviation,
`summ`
2. To compute means and standard deviations of select variables:
`summarize vone vtwo vthree`
3. Another way to compute means and standard deviations that allows the `by` option:
`tabstat vone vtwo, statistics(mean, sd) by(vthree)`
4. To get more numerical summaries for one variable:
`summ vone, detail`

5. See `help tabstat` to see the numerical summaries available. For example:
`tabstat vone, statistics(min, q, max, iqr, mean, sd)`
6. Correlation between two variables:
`correlate vone vtwo`
7. To see all values (all variables and all observations, not recommended for large data sets):
`list`
 Hit the space bar to see the next page after "-more-" or type "q" to "break" (stop/interrupt the listing).
8. To list the first 10 values for two variables:
`list vone vtwo in 1/10`
9. To list the last 10 values for two variables:
`list vone vtwo in -10/1`
 (The end of this command is "minus 10" / "lowercase letter L".)
10. Tabulate categorical variable *vname*:
`tabulate vname`
 or, using an abbreviation,
`tab vname`
11. Cross tabulate two categorical variables:
`tab vone vtwo`
12. Cross tabulate two variables, include one or more of the options to produce column, row or cell percents and to suppress printing of frequencies:
`tab vone vtwo, column row cell`
`tab vone vtwo, column row cell nofreq`

Generating new variables

1. General.
 - a. Generate index of cases $1, 2, \dots, n$ (this may be useful if you sort the data, then want to restore the data to the original form without reloading the data):
`generate case= _n`
 or, using an abbreviation,
`gen case=_n`
 - b. Multiply values in *vx* by *b* and add *a*, store results in *vy*:
`gen vy = a + b * vx`
 - c. Generate a variable with values 0 unless *vtwo* is greater than *c*, then make the value 1:
`gen vone=0`
`replace vone=1 if vtwo>c`
 - d.
2. Random numbers.

- a. Set numbers of observations to n :
`set obs n`
- b. Set random number seed to $XXXX$, default is 1000:
`set seed $XXXX$`
- c. Generate n uniform random variables (equal chance of all outcomes between 0 and 1):
`gen $vname$ =uniform()`
- d. Generate n uniform random variables (equal chance of all outcomes between a and b):
`gen $vname$ = a + (b - a)*uniform()`
- e. Generate n discrete uniform random variables (equal chance of all outcomes between 1 and 6)
`gen $vname$ =1 + int(6*uniform())`
(These commands simulate rolling a six-sided die.)
- f. Generate normal data with mean 0 and standard deviation 1:
`gen $vname$ = invnorm(uniform())`
- g. Generate normal data with mean μ and standard deviation σ :
`gen $vname$ = μ + σ * invnorm(uniform())`

Regression

1. Compute simple regression line (vy is response, vx is explanatory variable):
`regress vy vx`
2. Compute predictions, create new variable $yhat$:
`predict $yhat$`
3. Produce scatter plot with regression line added:
`graph twoway lfit vy vx || scatter vy vx`
4. Compute residuals, create new variable $residuals$:
`predict $residuals$, resid`
5. Produce a residual plot with horizontal line at 0:
`graph $residuals$, yline(0)`
6. Identify points with largest and smallest residuals:
`sort $residuals$`
`list in 1/5`
`list in -5/1`
(The last command is "minus 5" / "lowercase letter L".)
7. Compute multiple regression equation (vy is response, $vthree$, $vtwo$, and $vvthree$ are explanatory variables):
`regress vy $vone$ $vtwo$ $vthree$`

Important Notes on the "stem" command

In some versions of Stata, there is a potential glitch with Stata's `stem` command for stem-and-leaf plots. The `stem` function seems to permanently reorder the data so that they are sorted according to the variable that the stem-and-leaf plot was plotted for. The best way to avoid this problem is to avoid doing any stem-and-leaf plots (do histograms instead). However, if you really want to do a stem-and-leaf plot you should always create a variable containing the original observation numbers (called *index*, for example). A command to do so is:

```
generate index = _n
```

If you do this, then you can re-sort the data after the stem-and-leaf plot according to the *index* variable:

```
sort index.
```

Then, the data are back in the original order.

Summary of These and Other Commands

Here is a list of the commands demonstrated above and some other commands that you may find useful (this is by no means an exhaustive list of all Stata commands):

<code>anova</code>	general ANOVA, ANCOVA, or regression
<code>by</code>	repeat operation for categories of a variable
<code>ci</code>	confidence intervals for means
<code>clear</code>	clears previous dataset out of memory
<code>correlate</code>	correlation between variables
<code>describe</code>	briefly describes the data (# of obs, variable names, etc.)
<code>diagplot</code>	distribution diagnostic plots
<code>drop</code>	eliminate variables from memory
<code>edit</code>	better alternative to <code>input</code> for Macs
<code>exit</code>	leave Stata
<code>generate</code>	creates new variables (e.g., <code>generate years = last - first</code>)
<code>graph</code>	general graphing command (this command has many options)
<code>help</code>	online help
<code>histogram</code>	create a histogram graphic
<code>if</code>	lets you select a subset of observations (e.g., <code>list if radius >= 3000</code>)
<code>infile</code>	read non-Stata-format dataset (ASCII or text file)
<code>input</code>	type in raw data

<code>insheet</code>	read non-Stata-format spreadsheet with variable names on first line
<code>list</code>	lists the whole dataset in memory (you can also list only certain variables)
<code>log</code>	save or print Stata output (except graphs)
<code>lookup</code>	keyword search of commands, often precursor to <code>help</code>
<code>oneway</code>	oneway analysis of variance
<code>pcorr</code>	partial correlation coefficients
<code>plot</code>	text-mode (crude) scatterplots
<code>predict</code>	calculated predicted values (\hat{y}), residuals (ordinary, standardized and studentized), leverages, Cook's distance, standard error of predicted individual y , standard error of predicted mean y , standard error of residual from regression
<code>qnorm</code>	create a normal quantile plot
<code>regress</code>	regression
<code>replace</code>	lets you change individual values of a variable
<code>save</code>	saves data and labels in a Stata-format dataset
<code>scatter</code>	create a scatter plot of two numerical variables
<code>set</code>	set Stata system parameters (e.g., <code>obs</code> and <code>seed</code>)
<code>sebarr</code>	standard error-bar chart
<code>sort</code>	sorts observations from smallest to largest
<code>stem</code>	stem and leaf display
<code>summarize</code>	produces summary statistics (# obs, mean, sd, min, max) (has a <code>detail</code> option)
<code>tabstat</code>	produces summary statistics of your choice
<code>tabulate</code>	produces counts/frequencies for categorical data
<code>test</code>	conducts various hypothesis tests (refers back to most recent model fit (e.g., <code>regress</code> or <code>anova</code>) (see <code>help</code> function for info and examples))
<code>ttest</code>	one and two-sample t-tests
<code>use</code>	retrieve previously saved Stata dataset