### Introduction to Stata

### Christopher F Baum

Faculty Micro Resource Center Boston College

August 2009



What is Stata? Stata is a full-featured statistical programming language for Windows, Macintosh, Unix and Linux. It can be considered a "stat package," like SAS, SPSS, RATS, or eViews. The number of variables is limited to 2,047 in standard Stata/IC, but can be much larger in Stata/SE or Stata/MP. The number of observations is limited only by memory.

Stata has traditionally been a command-line-driven package that operates in a graphical (windowed) environment. Stata version 11 (released July 2009) contains a graphical user interface (GUI) for command entry. Stata may also be used in a command-line environment on a shared system (e.g., Unix) if you do not have a graphical interface to that system.





What is Stata? Stata is a full-featured statistical programming language for Windows, Macintosh, Unix and Linux. It can be considered a "stat package," like SAS, SPSS, RATS, or eViews. The number of variables is limited to 2,047 in standard Stata/IC, but can be much larger in Stata/SE or Stata/MP. The number of observations is limited only by memory.

Stata has traditionally been a command-line-driven package that operates in a graphical (windowed) environment. Stata version 11 (released July 2009) contains a graphical user interface (GUI) for command entry. Stata may also be used in a command-line environment on a shared system (e.g., Unix) if you do not have a graphical interface to that system.





Stata is eminently portable, and its developers are committed to cross-platform compatibility. Stata runs the same way on Windows, Macintosh, Unix, and Linux systems. The only platform-specific aspects of using Stata are those related to native operating system commands: e.g. is that file

```
C:\Stata\StataData\myfile.dta
or
/users/baum/statadata/myfile.dta
```

And—perhaps unique among statistical packages—Stata's binary data files may be freely copied from one platform to any other, or even accessed over the Internet from any machine that runs Stata.



- data manipulation
- statistics
- graphics



- data manipulation
- statistics
- graphics



- data manipulation
- statistics
- graphics



- data manipulation
- statistics
- graphics



- data manipulation
- statistics
- graphics



In terms of **statistics**, Stata provides all of the standard univariate, bivariate and multivariate statistical tools, from descriptive statistics and t-tests through one-, two- and N-way ANOVA, regression, principal components, and the like. Stata's regression capabilities are full-featured, including regression diagnostics, prediction, robust estimation of standard errors, instrumental variables and two-stage least squares, seemingly unrelated regressions, vector autoregressions and error correction models, etc. It has a very powerful set of techniques for the analysis of limited dependent variables: logit, probit, ordered logit and probit, multinomial logit, and the like.



Stata's breadth and depth really shines in terms of its specialized statistical capabilities. These include environments for time-series econometrics (ARCH, ARIMA, VAR, VEC), model simulation and bootstrapping, maximum likelihood estimation, and nonlinear least squares. Families of commands provide the leading techniques utilized in each of several categories:

- "xt" commands for cross-section/time-series or panel (longitudinal) data
- "svy" commands for the handling of survey data with complex sampling designs
- "st" commands for the handling of survival-time data with duration models





Stata's breadth and depth really shines in terms of its specialized statistical capabilities. These include environments for time-series econometrics (ARCH, ARIMA, VAR, VEC), model simulation and bootstrapping, maximum likelihood estimation, and nonlinear least squares. Families of commands provide the leading techniques utilized in each of several categories:

- "xt" commands for cross-section/time-series or panel (longitudinal) data
- "svy" commands for the handling of survey data with complex sampling designs
- "st" commands for the handling of survival-time data with duration models





Stata's breadth and depth really shines in terms of its specialized statistical capabilities. These include environments for time-series econometrics (ARCH, ARIMA, VAR, VEC), model simulation and bootstrapping, maximum likelihood estimation, and nonlinear least squares. Families of commands provide the leading techniques utilized in each of several categories:

- "xt" commands for cross-section/time-series or panel (longitudinal) data
- "svy" commands for the handling of survey data with complex sampling designs
- "st" commands for the handling of survival-time data with duration models





Stata **graphics** are excellent tools for exploratory data analysis, and can produce high-quality 2-D publication-quality graphics in several dozen different forms. Every aspect of graphics may be programmed and customized, and new graph types and graph "schemes" are being continuously developed. The programmability of graphics implies that a number of similar graphs may be generated without any "pointing and clicking" to alter aspects of the graphs. Stata does not have 3-D graphics capabilities, but those are under development in the new graphics system.



For members of the Boston College community, Stata is available through ITS' applications server, http://apps.bc.edu. After downloading client software from this site, you may connect to the apps server from any BC-activated computer and run Stata in a window on your computer. It is actually running the Windows version of Stata 10.1, but the interface and commands is almost identical to Stata for Mac OS X or Stata for Linux. Up to 50 users may access Stata on the apps server simultaneously. Results from your analysis may be stored on MyFiles, as the m: disk is automatically mapped to your account on MyFiles. If you are working from off campus, you must use set up VPN on your computer; see http://www.bc.edu/help for details.





If you would like your own copy of Stata, it is quite inexpensive. The vendor's GradPlan program makes the full version of Stata version 10 software available to BC faculty and students for \$98.00 (one-year license for students) or \$179.00 (perpetual license for faculty). As a first with Stata 11, the entire documentation set is installed in PDF format when you install Stata, and hyperlinked to the on-line help for each command and feature.



If you would like your own copy of Stata, it is quite inexpensive. The vendor's GradPlan program makes the full version of Stata version 10 software available to BC faculty and students for \$98.00 (one-year license for students) or \$179.00 (perpetual license for faculty). As a first with Stata 11, the entire documentation set is installed in PDF format when you install Stata, and hyperlinked to the on-line help for each command and feature.

The "Small Stata" version is available to students for \$49.00 for a one-year license. It contains all of Stata's commands, but can only handle a limited number of observations and variables (thus not recommended for Ph.D. students or Senior Honors Thesis students). GradPlan orders are made direct to Stata, with delivery from on-campus inventory.

Stata is very well supported by telephone and email technical support, as well as the more informal support provided by other users on **StataList**, the Stata listserv. The manuals are useful—particularly the User's Guide—but full details of the command syntax are available online, and in hypertext form in the GUI environment, with hyperlinks to the appropriate pages of the full documentation set of over a dozen manuals. The command findit *keyword* can also be used to locate Stata materials, including descriptions of built-in commands, Stata FAQs, and hundreds of user-written routines.



One of Stata's great strengths is that it can be updated over the Internet. Stata is actually a web browser, so it may contact Stata's web server and enquire whether there are more recent versions of either Stata's executable (the kernel) or the ado-files. The kernel is updated relatively infrequently—once a month at most—but the ado-files may be modified every ten days or so. This enables Stata's developers to distribute bug fixes, enhancements to existing commands, and even entirely new commands during the lifetime of a given release. Updates during the life of the version you own are free. You need only have a licensed copy of Stata and access to the Internet (which may be by proxy server) to check for and, if desired, download the updates.





## But why should I type commands?

But before we discuss the specifics to back up these claims, let's consider a meta-issue: why would you want to learn how to use a command-line-driven package? Isn't that ever so 20th century?

Stata may be used in an interactive mode, and those learning the package may wish to make use of the menu system. But when you execute a command from a pull-down menu, it records the command that you could have typed in the Review window, and thus you may learn that with experience you could type that command (or modify it and resubmit it) more quickly than by use of the menus.

Let us consider a couple of reasons why a command-line-driven package makes for an effective and efficient research strategy.



### But why should I type commands?

But before we discuss the specifics to back up these claims, let's consider a meta-issue: why would you want to learn how to use a command-line-driven package? Isn't that ever so 20th century?

Stata may be used in an interactive mode, and those learning the package may wish to make use of the menu system. But when you execute a command from a pull-down menu, it records the command that you could have typed in the Review window, and thus you may learn that with experience you could type that command (or modify it and resubmit it) more quickly than by use of the menus.

Let us consider a couple of reasons why a command-line-driven package makes for an effective and efficient research strategy.



### But why should I type commands?

But before we discuss the specifics to back up these claims, let's consider a meta-issue: why would you want to learn how to use a command-line-driven package? Isn't that ever so 20th century?

Stata may be used in an interactive mode, and those learning the package may wish to make use of the menu system. But when you execute a command from a pull-down menu, it records the command that you could have typed in the Review window, and thus you may learn that with experience you could type that command (or modify it and resubmit it) more quickly than by use of the menus.

Let us consider a couple of reasons why a command-line-driven package makes for an effective and efficient research strategy.



## Reproducibility

First, the important issue of reproducibility. If you are conducting scientific research, you must be able to reproduce your results. Ideally, anyone with your programs and data should be able to do so without your assistance. If you cannot produce such reproducible research findings, it can be argued that you are not following the scientific method, nor is your work conforming to ethical standards of research.

A thorough discussion of this issue is covered in the webpage, http://fmwww.bc.edu/GStat/docs/pointclick.html.





## Reproducibility

First, the important issue of reproducibility. If you are conducting scientific research, you must be able to reproduce your results. Ideally, anyone with your programs and data should be able to do so without your assistance. If you cannot produce such reproducible research findings, it can be argued that you are not following the scientific method, nor is your work conforming to ethical standards of research.

A thorough discussion of this issue is covered in the webpage, http://fmwww.bc.edu/GStat/docs/pointclick.html.





In a computer program where all actions are point and click, such as a spreadsheet, who can say how you arrived at a certain set of results? Unless every step of your transformations of the data can be retraced, how can you find exactly how the sample you are employing differs from the raw data? A command-driven program is capable of this level of reproducibility, we should all instill this level of rigor in our research practices.

Reproducibility also makes it very easy to perform an alternate analysis of a particular model. What would happen if we added this interaction, or introduced this additional variable, or decided to handle zero values as missing? Even if many steps have been taken since the basic model was specified, it is easy to go back and produce a variation on the analysis if all the work is represented by a series of programs.

In a computer program where all actions are point and click, such as a spreadsheet, who can say how you arrived at a certain set of results? Unless every step of your transformations of the data can be retraced, how can you find exactly how the sample you are employing differs from the raw data? A command-driven program is capable of this level of reproducibility, we should all instill this level of rigor in our research practices.

Reproducibility also makes it very easy to perform an alternate analysis of a particular model. What would happen if we added this interaction, or introduced this additional variable, or decided to handle zero values as missing? Even if many steps have been taken since the basic model was specified, it is easy to go back and produce a variation on the analysis if all the work is represented by a series of programs.

Stata makes this reproducibility very easy through a log facility, the ability to generate a command log (containing only the commands you have entered: see help cmdlog), and a "do-file editor" which allows you to easily enter, execute and save "do-files": sequences of commands, or program fragments. There is also an elaborate hypertext-based help browser, providing complete access to commands' descriptions and examples of syntax, with links to the appropriate pages of the PDF manuals. Each of these components appears in a separate window on the screen in the GUI version of Stata.



## **Extensibility**

Another clear advantage of the command-line driven environment is its interaction with the continual expansion of Stata's capabilities. A command, to Stata, is a verb instructing the program to perform some action.

Commands may be "built in" commands—those elements so frequently used that they have been coded into the "Stata kernel." A relatively small fraction of the total number of official Stata commands are built in, but they are used very heavily.





## **Extensibility**

Another clear advantage of the command-line driven environment is its interaction with the continual expansion of Stata's capabilities. A command, to Stata, is a verb instructing the program to perform some action.

Commands may be "built in" commands—those elements so frequently used that they have been coded into the "Stata kernel." A relatively small fraction of the total number of official Stata commands are built in, but they are used very heavily.





The vast majority of Stata commands are written in Stata's own programming language—the "ado-file" language. If a command is not built in to the Stata kernel, Stata searches for it along the "adopath". Like the PATH in Unix, Linux or DOS, the adopath indicates the several directories in which an ado-file might be located. This implies that the "official" Stata commands are not limited to those coded into the kernel.

If Stata's developers tomorrow wrote a command named "concatenate", they would make two files available on their web site: concatenate.ado (the ado-file code) and concatenate.sthlp (the associated help file). Both are straight ASCII text. These files should be produced in a text editor, not a word processing program.



The vast majority of Stata commands are written in Stata's own programming language—the "ado-file" language. If a command is not built in to the Stata kernel, Stata searches for it along the "adopath". Like the PATH in Unix, Linux or DOS, the adopath indicates the several directories in which an ado-file might be located. This implies that the "official" Stata commands are not limited to those coded into the kernel.

If Stata's developers tomorrow wrote a command named "concatenate", they would make two files available on their web site: concatenate.ado (the ado-file code) and concatenate.sthlp (the associated help file). Both are straight ASCII text. These files should be produced in a text editor, not a word processing program.



The importance of this program design goes far beyond the limits of official Stata. Since the adopath includes both Stata directories and other directories on your hard disk (or on a server's filesystem), you may acquire new Stata commands from a number of web sites. The *Stata Journal* (SJ), a quarterly refereed journal, is the primary method for distributing user contributions. Between 1991 and 2001, the *Stata Technical Bulletin* played this role, and a complete set of issues of the STB are available on line at http://ideas.repec.org.

The SJ is a subscription publication (available at O'Neill Library: older issues online at IDEAS), but the ado- and sthlp-files may be freely downloaded from Stata's web site. The Stata command help accesses help on all installed commands; the Stata command findit will locate commands that have been documented in the STB and the SJ, and with one click you may install them in your version of Stata. Help for these commands will then be available in your own Stata.

The importance of this program design goes far beyond the limits of official Stata. Since the adopath includes both Stata directories and other directories on your hard disk (or on a server's filesystem), you may acquire new Stata commands from a number of web sites. The *Stata Journal* (SJ), a quarterly refereed journal, is the primary method for distributing user contributions. Between 1991 and 2001, the *Stata Technical Bulletin* played this role, and a complete set of issues of the STB are available on line at http://ideas.repec.org.

The SJ is a subscription publication (available at O'Neill Library: older issues online at IDEAS), but the ado- and sthlp-files may be freely downloaded from Stata's web site. The Stata command help accesses help on all installed commands; the Stata command findit will locate commands that have been documented in the STB and the SJ, and with one click you may install them in your version of Stata. Help for these commands will then be available in your own Stata.

# User extensibility: the SSC archive

But this is only the beginning. Stata users worldwide participate in the **StataList** listserv, and when a user has written and documented a new general-purpose command to extend Stata functionality, they announce it on the <code>StataList</code> listserv (to which you may freely subscribe: see Stata's web site). Since September 1997, all items posted to <code>StataList</code> (over 1,000) have been placed in the Boston College Statistical Software Components (SSC) Archive in **RePEc**, available from IDEAS (http://ideas.repec.org) and EconPapers (http://econpapers.repec.org).





Any component in the SSC archive may be readily inspected with a web browser, using IDEAS' or EconPapers' search functions, and if desired you may install it with one command from the archive from within Stata. For instance, if you know there is a module in the archive named "ivreset," you could use ssc install ivreset to install it. Anything in the archive can be accessed via Stata's ssc command: thus ssc describe ivreset will locate this module, and make it possible to install it with one click.

Windows users should not attempt to download the materials from a web browser; it won't work.





The command ssc new lists, in the Stata Viewer, all SSC packages that have been added or modified in the last month. You may click on their names for full details. The command ssc hot reports on the most popular packages on the SSC Archive.

The Stata command adoupdate checks to see whether all packages you have downloaded and installed from the SSC archive, the *Stata Journal*, or other user-maintained net from... sites are up to date. adoupdate alone will provide a list of packages that have been updated. You may then use adoupdate, update to refresh your copies of those packages, or specify which packages are to be updated.





The importance of all this is that Stata is **infinitely extensible**. Any ado-file on your adopath is a full-fledged Stata command. Stata's capabilities thus extend far beyond the official, supported features described in the Stata manual to a vast array of additional tools.

Since the current directory is on the adopath, if I create an ado-file hello.ado:

```
program define hello
display "hello from Stata!"
end
exit
```

Stata will now respond to the command hello. It's that easy.



The importance of all this is that Stata is **infinitely extensible**. Any ado-file on your adopath is a full-fledged Stata command. Stata's capabilities thus extend far beyond the official, supported features described in the Stata manual to a vast array of additional tools.

Since the current directory is on the adopath, if I create an ado-file hello.ado:

```
program define hello
display "hello from Stata!"
end
exit.
```

Stata will now respond to the command hello. It's that easy.



### **Transportability**

Stata binary files may be easily transformed into SPSS or SAS files with the third-party application Stat/Transfer. Stat/Transfer is available for Windows and Mac OS X systems as well as on various Unix systems on campus. Personal copies of Stat/Transfer version 9 (which handles Stata versions 6, 7, 8, 9, 10 and 11 datafiles) are available at a discounted academic rate of \$69.00 through the Stata GradPlan.

Stat/Transfer can also transfer SAS, SPSS and many other file formats into Stata format, without loss of variable labels, value labels, and the like. It can also be used to create a manageable subset of a very large Stata file (such as those produced from survey data) by selecting only the variables you need. It is a very useful tool.



# **Command syntax**

We now consider the form of Stata commands. One of Stata's great strengths, compared with many statistical packages, is that its command syntax follows strict rules: in grammatical terms, there are no irregular verbs. This implies that when you have learned the way a few key commands work, you will be able to use many more without extensive study of the manual or even on-line help. The search command will allow you to find the command you need by entering one or more keywords, even if you do not know the command's name.



The fundamental syntax of all Stata commands follows a *template*. Not all elements of the template are used by all commands, and some elements are only valid for certain commands. But where an element appears, it will appear in the same place, following the same grammar. Like Unix or Linux, Stata is case sensitive. Commands must be given in lower case. For best results, keep all variable names in lower case to avoid confusion.

Following the examples in the *Getting Started with Stata...* manual, we will make use of auto.dta, a dataset of 74 automobiles' characteristics.





### The general syntax of a Stata command is:

### where elements in square brackets are optional for some commands.

In some cases, only the cmdname itself is required. describe without arguments gives a description of the current contents of memory (including the identifier and timestamp of the current dataset), while summarize without arguments provides summary statistics for all (numeric) variables. Both may be given with a varlist specifying the variables to be considered.

What are the other elements?



### The general syntax of a Stata command is:

where elements in square brackets are optional for some commands.

In some cases, only the cmdname itself is required. describe without arguments gives a description of the current contents of memory (including the identifier and timestamp of the current dataset), while summarize without arguments provides summary statistics for all (numeric) variables. Both may be given with a varlist specifying the variables to be considered.

What are the other elements?



### The varlist

varlist is a list of one or more variables on which the command is to operate: the subject(s) of the verb. Stata works on the concept of a single set of variables currently defined and contained in memory, each of which has a name. As desc will show you, each variable has a data type (various sorts of integers and reals, and string variables of a specified maximum length). The varlist specifies which of the defined variables are to be used in the command.

The order of variables in the dataset matters, since you can use hyphenated lists to include all variables between first and last. (The order and move commands can alter the order of variables.) You can also use "wildcards" to refer to all variables with a certain prefix. If you have variables pop60, pop70, pop80, pop90, you can refer to them in a varlist as pop\* or pop?0.



# The exp clause

The *exp* clause is used in commands such as <code>generate</code> and <code>replace</code> where an algebraic expression is used to produce a new (or updated) variable. In algebraic expressions, the operators ==, &, | and ! are used as equal, AND, OR and NOT, respectively. The  $\land$  operator is used to denote exponentiation. The + operator is overloaded to denote concatenation of character strings.





#### The if and in clauses

Stata differs from several common programs in that Stata commands will automatically apply to all observations currently defined. You need not write explicit loops over the observations. You can, but it is usually bad programming practice to do so. Of course you may want not to refer to all observations, but to pick out those that satisfy some criterion. This is the purpose of the *if exp* and *in range* clauses. For instance, we might:

```
sort price
list make price in 1/5
```

to determine the five cheapest cars in auto.dta. The 1/5 is a *numlist*: in this case, a list of observation numbers.  $\ell$  is the last observation, thus *list make price in -5/\ell* will list the five most expensive cars in auto.dta

Even more commonly, you may employ the *if exp* clause. This restricts the set of observations to those for which the "exp", a Boolean expression, evaluates to true. Stata's missing value codes are greater than the largest positive number, so that the last command would avoid listing cars for which the price is missing.

```
list make price if foreign==1
```

lists only foreign cars, and

```
list make price if price > 10000 & price <.
```

lists only expensive cars (in 1978 prices!) Note the double equal in the *exp*. A single equal sign, as in the C language, is used for assignment; double equal for comparison.



# The using clause

Some commands access files: reading data from external files, or writing to files. These commands contain a *using* clause, in which the filename appears. If a file is being written, you must specify the "replace" option to overwrite an existing file of that name.

Stata's own binary file format, the .dta file, is cross-platform compatible, even between machines with different byte orderings (low-endian and high-endian). A .dta file may be moved from one computer to another using ftp (in binary transfer mode).





# The using clause

Some commands access files: reading data from external files, or writing to files. These commands contain a *using* clause, in which the filename appears. If a file is being written, you must specify the "replace" option to overwrite an existing file of that name.

Stata's own binary file format, the .dta file, is cross-platform compatible, even between machines with different byte orderings (low-endian and high-endian). A .dta file may be moved from one computer to another using ftp (in binary transfer mode).





To bring the contents of an existing Stata file into memory, the command:

use file [,clear]

is employed (clear will empty the current contents of memory). You must make sufficient memory available to Stata to load the entire file, since Stata's speed is largely derived from holding the entire data set in memory. Consult *Getting Started...* for details on adjusting the memory allocation on your computer, since it differs by operating system.





Reading and writing binary (.dta) files is much faster than dealing with text (ASCII) files (with the insheet or infile commands), and permits variable labels, value labels, and other characteristics of the file to be saved along with the file. To write a Stata binary file, the command

save file [,replace]

is employed. The compress command can be used to economize on the disk space (and memory) required to store variables.

Stata's version 10 and 11 datasets cannot be read by version 8 or 9; to create a compatible dataset, use saveold.





The amazing thing about "use filename" is that it is by no means limited to the files on your hard disk. Since Stata is a web browser,

webuse klein

or

use http://fmwww.bc.edu/ec-p/data/Wooldridge/crimel.dta will read these datasets into Stata's memory over the web.



# The type command can display any text file, whether on your hard disk or over the Web; thus

type http://fmwww.bc.edu/ec-p/data/Wooldridge/crime1.des

will display the codebook for this file, and

 $\verb"copy" http://fmwww.bc.edu/ec-p/data/Wooldridge/crime1.des crime.codebook" and the control of the control of$ 

will make a copy of the codebook on your own hard disk.





When you have used a dataset over the Web, you have loaded it into memory in your desktop Stata. You cannot save it to the Web, but can save the data to your own hard disk. The advantages of this feature for instructional and collaborative research should be clear. Students may be given a URL from which their assigned data are to be accessed; it matters not whether they are using Stata for Windows, Macintosh, Linux, or Unix.





# The options clause

Many commands make use of options (such as clear on use, or replace on save). All options are given following a single comma, and may be given in any order. Options, like commands, may generally be abbreviated (with the notable exception of replace).





#### **Prefix commands**

A number of Stata commands can be used as *prefix commands*, preceding a Stata command and modifying its behavior. The most commonly employed is the *by prefix*, which repeats a command over a set of categories. The *statsby:* prefix repeats the command, but collects statistics from each category. The *rolling:* prefix runs the command on moving subsets of the data (usually time series).

Several other command prefixes: *simulate:*, which simulates a statistical model; *bootstrap:*, allowing the computation of bootstrap statistics from resampled data; and *jackknife:*, which runs a command over jackknife subsets of the data. The *svy:* prefix can be used with many statistical commands to allow for survey sample design. See my separate slideshow on *Monte Carlo Simulation in Stata.* 



#### Prefix commands

A number of Stata commands can be used as *prefix commands*, preceding a Stata command and modifying its behavior. The most commonly employed is the by prefix, which repeats a command over a set of categories. The statsby: prefix repeats the command, but collects statistics from each category. The *rolling:* prefix runs the command on moving subsets of the data (usually time series).

Several other command prefixes: simulate:, which simulates a statistical model; bootstrap:, allowing the computation of bootstrap statistics from resampled data; and jackknife:, which runs a command over jackknife subsets of the data. The svy: prefix can be used with many statistical commands to allow for survey sample design. See my separate slideshow on *Monte Carlo Simulation in Stata*.

### The by prefix

You can often save time and effort by using the *by* prefix. When a command is prefixed with a *bylist*, it is performed repeatedly for each element of the variable or variables in that list, each of which must be categorical. For instance,

by foreign: summ price

will provide descriptive statistics for both foreign and domestic cars. If the data are not already sorted by the bylist variables, the prefix bysort should be used. The option , total will add the overall summary.





What about a classification with several levels, or a combination of values?

```
bysort rep78: summ price
```

bysort rep78 foreign: summ price

This is a very handy tool, which often replaces explicit loops that must be used in other programs to achieve the same end.





The by prefix should not be confused with the by *option* available on some commands, which allows for specification of a grouping variable: for instance

ttest price, by(foreign)

will run a t-test for the difference of sample means across domestic and foreign cars.

Another useful aspect of *by* is the way in which it modifies the meanings of the observation number symbol. Usually *\_n* refers to the current observation number, which varies from 1 to *\_N*, the maximum defined observation. Under a bylist, *\_n* refers to the observation within the bylist, and *\_N* to the total number of observations for that category. This is often useful in creating new variables.



The by prefix should not be confused with the by *option* available on some commands, which allows for specification of a grouping variable: for instance

```
ttest price, by(foreign)
```

will run a t-test for the difference of sample means across domestic and foreign cars.

Another useful aspect of *by* is the way in which it modifies the meanings of the observation number symbol. Usually \_n refers to the current observation number, which varies from 1 to \_N, the maximum defined observation. Under a bylist, \_n refers to the observation within the bylist, and \_N to the total number of observations for that category. This is often useful in creating new variables.



For instance, if you have individual data with a family identifier, these commands might be useful:

```
sort famid age
by famid: gen famsize = _N
by famid: gen birthorder = _N - _n +1
```

Here the famsize variable is set to \_N, the total number of records for that family, while the birthorder variable is generated by sorting the family members' ages within each family.





# Missing values

Missing value codes in Stata appear as the dot (.) in printed output (and a string missing value code as well: "", the null string). It takes on the largest possible positive value, so in the presence of missing data you do not want to say

```
generate hiprice = (price > 10000), but rather
generate hiprice = (price > 10000 & price <.)</pre>
```

which then generates a "dummy variable" for high-priced cars (for which price data are complete, with prices "less than missing").

As of version 8, Stata allows for multiple missing value codes (.a, .b, .c, ..., .z).





### **Display formats**

Each variable may have its own default display format. This does not alter the contents of the variable, but affects how it is displayed. For instance,  $\$9.2 \, \text{f}$  would display a two-decimal-place real number. The command

format varname %9.2f

will save that format as the default format of the variable, and

format date %tm

will format a Stata date variable into a monthly format (e.g., 1998m10).



### Variable labels

Each variable may have its own variable label. The variable label is a character string (maximum 80 characters) which describes the variable, associated with the variable via

label variable varname "text"

Variable labels, where defined, will be used to identify the variable in printed output, space permitting.





### Value labels

Value labels associate numeric values with character strings. They exist separately from variables, so that the same mapping of numerics to their definitions can be defined once and applied to a set of variables (e.g. 1=very satisfied...5=not satisfied may be applied to all responses to questions about consumer satisfaction). Value labels are saved in the dataset. For example:

```
label define sex1bl 0 male 1 female
label values sex sexibl
```

If value labels are defined, they will be displayed in printed output instead of the numeric values.





# **Generating new variables**

The command <code>generate</code> is used to produce new variables in the dataset, whereas <code>replace</code> must be used to revise an existing variable (and <code>replace</code> must be spelled out). The syntax just demonstrated is often useful if you are trying to generate indicator variables, or dummies, since it combines a <code>generate</code> and <code>replace</code> in a single command.

A full set of functions are available for use in the generate command including the standard mathematical functions, recode functions, string functions, date and time functions, and specialized functions (help functions for details). Note that generate's sum() function is a running sum.



# **Generating new variables**

The command <code>generate</code> is used to produce new variables in the dataset, whereas <code>replace</code> must be used to revise an existing variable (and <code>replace</code> must be spelled out). The syntax just demonstrated is often useful if you are trying to generate indicator variables, or dummies, since it combines a <code>generate</code> and <code>replace</code> in a single command.

A full set of functions are available for use in the <code>generate</code> command, including the standard mathematical functions, recode functions, string functions, date and time functions, and specialized functions (help functions for details). Note that <code>generate</code>'s <code>sum()</code> function is a running sum.

# The egen command

Stata is not limited to using the set of defined functions. The <code>egen</code> (extended <code>gen</code>erate) command makes use of functions written in the Stata ado-file language, so that  $\_gzap.ado$  would define the extended generate function zap(). This would then be invoked as

```
egen newvar = zap(oldvar)
```

which would do whatever zap does on the contents of oldvar, creating the new variable newvar.

A number of egen functions provide row-wise operations similar to those available in a spreadsheet: row sum, row average, row standard deviation, etc.



# **Time series operators**

The D., L., and F. operators may be used under a timeseries calendar (including in the context of panel data) to specify first differences, lags, and leads, respectively. These operators understand missing data, and numlists: e.g. L (1/4) .x is the first through fourth lags of x, while L2D.x is the second lag of the first difference of the x variable.

It is important to use the time series operators to refer to lagged or led values, rather than referring to the observation number (e.g.,  $\_n-1$ ). The time series operators respect the time series calendar, and will not mistakenly compute a lag or difference from a prior period if it is missing. This is particularly important when working with panel data to ensure that references to one individual do not reach back into the prior individual's data.

# **Time series operators**

The D., L., and F. operators may be used under a timeseries calendar (including in the context of panel data) to specify first differences, lags, and leads, respectively. These operators understand missing data, and numlists: e.g. L(1/4).x is the first through fourth lags of x, while L2D.x is the second lag of the first difference of the x variable.

It is important to use the time series operators to refer to lagged or led values, rather than referring to the observation number (e.g.,  $\_n-1$ ). The time series operators respect the time series calendar, and will not mistakenly compute a lag or difference from a prior period if it is missing. This is particularly important when working with panel data to ensure that references to one individual do not reach back into the prior individual's data.

### Mata: Matrix programming language

As of version 9, Stata contains a full-fledged matrix programming language, *Mata*, with all of the capabilities of MATLAB, Ox or GAUSS. Mata can be used interactively, or Mata functions can be developed to be called from Stata. A large library of mathematical and matrix functions is provided in Mata, including equation solvers, decompositions, eigensystem routines and probability density functions. Mata functions can access Stata's variables and can work with virtual matrices ("views") of a subset of the data in memory. Mata also supports file input/output.

Mata code is automatically compiled into bytecode, like Java, and ca be stored in object form or included in-line in a Stata do-file or ado-fil Mata code runs many times faster than the interpreted ado-file language, providing significant speed enhancements to many computationally burdensome tasks. See my separate slideshow Mata in Stata.

# Mata: Matrix programming language

As of version 9, Stata contains a full-fledged matrix programming language, *Mata*, with all of the capabilities of MATLAB, Ox or GAUSS. Mata can be used interactively, or Mata functions can be developed to be called from Stata. A large library of mathematical and matrix functions is provided in Mata, including equation solvers, decompositions, eigensystem routines and probability density functions. Mata functions can access Stata's variables and can work with virtual matrices ("views") of a subset of the data in memory. Mata also supports file input/output.

Mata code is automatically compiled into bytecode, like Java, and can be stored in object form or included in-line in a Stata do-file or ado-file. Mata code runs many times faster than the interpreted ado-file language, providing significant speed enhancements to many computationally burdensome tasks. See my separate slideshow *Mata in Stata*.



All estimation commands share the same syntax. Multiple equation estimation commands use a list of equations, rather than a *varlist*, where equations are defined in parenthesized *varlist*s. Most estimation commands allow the use of various kinds of weights.

Estimation commands display confidence intervals for the coefficients, and tests of the most common hypotheses. More complex hypotheses may be analyzed with the test and lincom commands; for nonlinear hypothesis, testnl and nlcom may be applied, making use of the delta method.

Robust (Huber/White) estimates of the covariance matrix are available for almost all estimation commands by employing the robust option.

All estimation commands share the same syntax. Multiple equation estimation commands use a list of equations, rather than a *varlist*, where equations are defined in parenthesized *varlist*s. Most estimation commands allow the use of various kinds of weights.

Estimation commands display confidence intervals for the coefficients, and tests of the most common hypotheses. More complex hypotheses may be analyzed with the test and lincom commands; for nonlinear hypothesis, testnl and nlcom may be applied, making use of the delta method.

Robust (Huber/White) estimates of the covariance matrix are available for almost all estimation commands by employing the robust option.

All estimation commands share the same syntax. Multiple equation estimation commands use a list of equations, rather than a *varlist*, where equations are defined in parenthesized *varlist*s. Most estimation commands allow the use of various kinds of weights.

Estimation commands display confidence intervals for the coefficients, and tests of the most common hypotheses. More complex hypotheses may be analyzed with the test and lincom commands; for nonlinear hypothesis, testnl and nlcom may be applied, making use of the delta method.

Robust (Huber/White) estimates of the covariance matrix are available for almost all estimation commands by employing the robust option.

Predicted values and residuals may be obtained after any estimation command with the <code>predict</code> command. For nonlinear estimators, <code>predict</code> will produce other statistics as well (e.g. the log of the odds ratio from logistic regression). The <code>mfx</code> command may be used to generate marginal effects, including elasticities and semi–elasticities, for any estimation command.

All estimation commands "leave behind" results of estimation in the e() array, where they may be inspected with ereturn list. Any item here, including scalars such as  $R^2$  and RMSE, the coefficient vector, and the estimated variance-covariance matrix, may be saved for use in later calculations.





Predicted values and residuals may be obtained after any estimation command with the <code>predict</code> command. For nonlinear estimators, <code>predict</code> will produce other statistics as well (e.g. the log of the odds ratio from logistic regression). The <code>mfx</code> command may be used to generate marginal effects, including elasticities and semi–elasticities, for any estimation command.

All estimation commands "leave behind" results of estimation in the e() array, where they may be inspected with  $ereturn\ list.$  Any item here, including scalars such as  $R^2$  and RMSE, the coefficient vector, and the estimated variance-covariance matrix, may be saved for use in later calculations.





The estimates suite of commands allow you to store the results of a particular estimation for later use in a Stata session. For instance, after the commands

regress price mpg length turn
estimates store model1
regress price weight length displacement
estimates store model2
regress price weight length gear\_ratio foreign
estimates store model3





#### the command

estimates table model1 model2 model3

will produce a nicely-formatted table of results. Options on estimates table allow you to control precision, whether standard errors or t-values are given, significance stars, summary statistics, etc.

### For example

```
estimates table model1 model2 model3, b(%10.3f)
se(%7.2f) stats(r2 rmse N) title(Some models of auto
price)
```





#### the command

estimates table model1 model2 model3

will produce a nicely-formatted table of results. Options on estimates table allow you to control precision, whether standard errors or t-values are given, significance stars, summary statistics, etc.

#### For example:

```
estimates table model1 model2 model3, b(%10.3f) se(%7.2f) stats(r2 rmse N) title(Some models of auto price)
```





Although estimates table can produce a summary table quite useful for evaluating a number of specifications, we often want to produce a publication-quality table for inclusion in a word processing document. Ben Jann's estout command processes stored estimates and provides a great deal of flexibility in generating such a table.

Programs in the estout suite can produce tab-delimited tables for MS Word, HTML tables for the web, and—my favorite—LATEX tables for professional papers. In the LATEX output format, estout can generate Greek letters, sub- and superscripts, and the like. estout is available from SSC, with extensive on-line help, and was described in the *Stata Journal*, 5(3), 2005 and 7(2), 2007. It has its own website at http://repec.org/bocode/e/estout.



From the example above, rather than using estimates save and estimates table we use Jann's eststo (store) and esttab (table) commands:

```
eststo clear
eststo: reg price mpg length turn
eststo: reg price weight length displacement
eststo: reg price weight length gear_ratio foreign
esttab using autol.tex, stats(r2 bic N) ///
subst(r2 \$R^2$) title(Models of auto price) ///
replace
```



Publication-quality tables

Table 1: Models of auto price			
	(1)	(2)	(3)
	price	price	price
mpg	$-186.7^{*}$		
	(-2.13)		
length	52.58	-97.63*	-88.03*
	(1.67)	(-2.47)	(-2.65)
turn	-199.0		
	(-1.44)		
weight		4.613**	5.479***
		(3.30)	(5.24)
displacement		0.727	
		(0.10)	
gear_ratio			-669.1
0			(-0.72)
foreign			3837.9***
			(5.19)
_cons	8148.0	10440.6*	7041.5
	(1.35)	(2.39)	(1.46)
$R^2$	0.251	0.348	0.552
bic	1387.2	1377.0	1353.5
N	74	74	74

t statistics in parentheses



 $<sup>^*</sup>$   $p < 0.05, \,^{**}$   $p < 0.01, \,^{***}$  p < 0.001

### File handling

File extensions usually employed (but not required) include:

```
.ado
       automatic do-file (defines a Stata command)
.dct
       data dictionary, optionally used with infile
. do
       do-file (user program)
.dta
       Stata binary dataset
.gph
       graphics output file (binary)
.loa
       text log file
.smcl
       SMCL (markup) log file, for use with Viewer
       ASCII data file
.raw
.sthlp Stata help file
```

These extensions need not be given (except for .ado). If you use other extensions, they must be explicitly specified.



Comma-separated (CSV) files or tab-delimited data files may be read very easily with the <code>insheet</code> command—which despite its name does not read spreadsheet files. If your file has variable names in the first row that are valid for Stata, they will be automatically used (rather than default variable names). You usually need not specify whether the data are tab- or comma-delimited. Note that <code>insheet</code> cannot read space-delimited data (or character strings with embedded spaces, unless they are quoted).



If the file extension is .raw, you may just use

insheet using filename

to read it. If other file extensions are used, they must be given:

insheet using filename.csv
insheet using filename.txt





A free-format ASCII text file with space-, tab-, or comma-delimited data may be read with the infile command. The missing-data indicator (.) may be used to specify that values are missing.

The command must specify the variable names. Assuming auto.raw contains numeric data,

infile price mpg displacement using auto

will read it. If a file contains a combination of string and numeric values in a variable, it should be read as string, and <code>encode</code> used to convert it to numeric with string value labels.



If some of the data are string variables without embedded spaces, they must be specified in the command:

infile str3 country price mpg displacement using auto2

would read a three-letter country of origin code, followed by the numeric variables. The number of observations will be determined from the available data





The infile command may also be used with fixed-format data, including data containing undelimited string variables, by creating a dictionary file which describes the format of each variable and specifies where the data are to be found. The dictionary may also specify that more than one record in the input file corresponds to a single observation in the data set.

If data fields are not delimited—for instance, if the sequence '102' should actually be considered as three integer variables. A dictionary must be used to define the variables' locations. The byvariable() option allows a variable-wise dataset to be read, where one specifies the number of observations available for each series.





The infile command may also be used with fixed-format data, including data containing undelimited string variables, by creating a dictionary file which describes the format of each variable and specifies where the data are to be found. The dictionary may also specify that more than one record in the input file corresponds to a single observation in the data set.

If data fields are not delimited—for instance, if the sequence '102' should actually be considered as three integer variables. A dictionary must be used to define the variables' locations. The byvariable() option allows a variable-wise dataset to be read, where one specifies the number of observations available for each series.



An alternative to infile with a dictionary is the infix command, which presents a syntax similar to that used by SAS for the definition of variables' data types and locations in a fixed-format ASCII data set: that is, a data file in which certain columns contain certain variables. The \_column() directive allow contents of a fixed-format data file to be retrieved selectively.

infix may also be used for more complex record layouts where one individual's data are contained on several records in an ASCII file.



A logical condition may be used on the infile or infix commands to read only those records for which certain conditions are satisfied: i.e.

```
infix using employee if sex=="M"
infile price mpg using auto in 1/20
```

where the latter will read only the first 20 observations from the external file. This might be very useful when reading a large data set, where one can check to see that the formats are being properly specified on a subset of the file.





If your data are already in the internal format of SAS, SPSS, Excel, GAUSS, MATLAB, or a number of other packages, the best way to get it into Stata is by using the third-party product Stat/Transfer. Stat/Transfer will preserve variable labels, value labels, and other aspects of the data, and can be used to convert a Stata binary file into other packages' formats. It can also produce subsets of the data (selecting variables, cases or both) so as to generate an extract file that is more manageable. This is particularly important when the 2.047-variable limit on standard Stata data sets is encountered. Stat/Transfer is well documented, with on-line help available in both Windows. Mac OS X and Unix versions, and an extensive manual.





### Combining data sets

In many empirical research projects, the raw data to be utilized are stored in a number of separate files: separate "waves" of panel data, timeseries data extracted from different databases, and the like. Stata only permits a single data set to be accessed at one time. How, then, do you work with multiple data sets? Several commands are available, including append, merge, and joinby.

The append command combines two Stata-format data sets that possess variables in common, adding observations to the existing variables. The same variables need not be present in both files, as long as a subset of the variables are common to the "master" and "using" data sets. It is important to note that "PRICE" and "price" are different variables, and one will not be appended to the other.



### Combining data sets

In many empirical research projects, the raw data to be utilized are stored in a number of separate files: separate "waves" of panel data, timeseries data extracted from different databases, and the like. Stata only permits a single data set to be accessed at one time. How, then, do you work with multiple data sets? Several commands are available, including append, merge, and joinby.

The append command combines two Stata-format data sets that possess variables in common, adding observations to the existing variables. The same variables need not be present in both files, as long as a subset of the variables are common to the "master" and "using" data sets. It is important to note that "PRICE" and "price" are different variables, and one will not be appended to the other.



We now describe the merge command. Its syntax has changed considerably in Stata version 11. As you may not have access yet to Stata version 11, I describe the version 10 syntax. To invoke the older version of merge, use

version 10: merge ...

For details of the new version of merge, see help merge in version 11.





The merge command is very powerful. Like append, it works on a "master" data set—the current contents of memory—and one or more "using" data sets. One or more merge variables are specified, and both master and using data sets must be sorted on those variables.

The distinction between "master" and "using" is important. When the same variable is present in each of the files, Stata's default behavior is to hold the master data inviolate and discard the using dataset's copy of that variable. This may be modified by the update option, which specifies that non-missing values in the using dataset should replace missing values in the master, and update replace, which specifies that non-missing values in the using dataset should take precedence.





The merge command is very powerful. Like append, it works on a "master" data set—the current contents of memory—and one or more "using" data sets. One or more merge variables are specified, and both master and using data sets must be sorted on those variables.

The distinction between "master" and "using" is important. When the same variable is present in each of the files, Stata's default behavior is to hold the master data inviolate and discard the using dataset's copy of that variable. This may be modified by the <code>update</code> option, which specifies that non-missing values in the using dataset should replace missing values in the master, and <code>update replace</code>, which specifies that non-missing values in the using dataset should take precedence.



August 2009

A "one-to-one" merge specifies that each record in the using data set is to be combined with one record in the master data set. This would be appropriate if you acquired additional variables for the same observations. A new variable, \_merge, takes on integer values indicating whether an observation appears in the master only, the using only, or appears in both. This may be used to determine whether the merge has been successful, or to remove those observations which remain unmatched (e.g. merging a set of households from different cities with a comprehensive list of ZIP codes; one would then discard all the unused ZIP code records). The \_merge variable must be dropped before another merge is performed on this data set.

The unique option should be used if you believe that both data sets should have unique values of the *merge key*.



A "one-to-one" merge specifies that each record in the using data set is to be combined with one record in the master data set. This would be appropriate if you acquired additional variables for the same observations. A new variable, \_merge, takes on integer values indicating whether an observation appears in the master only, the using only, or appears in both. This may be used to determine whether the merge has been successful, or to remove those observations which remain unmatched (e.g. merging a set of households from different cities with a comprehensive list of ZIP codes; one would then discard all the unused ZIP code records). The merge variable must be dropped before another merge is performed on this data set.

The unique option should be used if you believe that both data sets should have unique values of the *merge key*.



The merge command can also do a "**match merge**", or "**one-to-N**" merge, in which each record in the using data set is matched with a number of records in the master data set. If a number of the households lived in the same ZIP code, then the match would place variables from the ZIP code file on the household records, repeating where necessary. This is a very useful technique to combine aggregate data with disaggregate data without dealing with the details.

Although "one-to-N" and "N-to-one" merges are commonplace and very useful, you probably never want to do a "N-to-N" merge, which will yield seemingly random results. To ensure that one data set has unique identifiers, specify the uniquaster or uniquaing options, or use the isid command to ensure that a dataset has a unique identifier.



The merge command can also do a "match merge", or "one-to-N" merge, in which each record in the using data set is matched with a number of records in the master data set. If a number of the households lived in the same ZIP code, then the match would place variables from the ZIP code file on the household records, repeating where necessary. This is a very useful technique to combine aggregate data with disaggregate data without dealing with the details.

Although "one-to-N" and "N-to-one" merges are commonplace and very useful, you probably never want to do a "N-to-N" merge, which will vield seemingly random results. To ensure that one data set has unique identifiers, specify the unique identifiers, specify the unique identifiers, or use the isid command to ensure that a dataset has a unique identifier.

### Writing external data

If you want to transfer data to another package, Stat/Transfer is very useful. But if you just want to create an ASCII file from Stata, the outfile command may be used. It takes a *varlist*, and the if or in clauses may be used to control the observations to be exported. Applying sort prior to outfile will control the order of observations in the external file. You may specify that the data are to be written in comma-separated format.

The outsheet command can write a comma-delimited or tab-delimited ASCII file, optionally placing the variable names in the first row. Such a file can be easily read by a spreadsheet program such as Excel. Note that outsheet does not write spreadsheet files.

For customized output, the file command can write out information (including scalars, matrices and macros, text strings, etc.) in any ASCII or binary format of your choosing.



### Writing external data

If you want to transfer data to another package, Stat/Transfer is very useful. But if you just want to create an ASCII file from Stata, the outfile command may be used. It takes a *varlist*, and the if or in clauses may be used to control the observations to be exported. Applying sort prior to outfile will control the order of observations in the external file. You may specify that the data are to be written in comma-separated format.

The outsheet command can write a comma-delimited or tab-delimited ASCII file, optionally placing the variable names in the first row. Such a file can be easily read by a spreadsheet program such as Excel. Note that outsheet does *not* write spreadsheet files.

For customized output, the file command can write out information (including scalars, matrices and macros, text strings, etc.) in any ASCII or binary format of your choosing.



#### Writing external data

If you want to transfer data to another package, Stat/Transfer is very useful. But if you just want to create an ASCII file from Stata, the outfile command may be used. It takes a *varlist*, and the if or in clauses may be used to control the observations to be exported. Applying sort prior to outfile will control the order of observations in the external file. You may specify that the data are to be written in comma-separated format.

The outsheet command can write a comma-delimited or tab-delimited ASCII file, optionally placing the variable names in the first row. Such a file can be easily read by a spreadsheet program such as Excel. Note that outsheet does *not* write spreadsheet files.

For customized output, the file command can write out information (including scalars, matrices and macros, text strings, etc.) in any ASCII or binary format of your choosing.



A very useful capability is provided by the postfile and post commands, which permit a Stata data set to be created in the course of a program. For instance, you may be simulating the distribution of a statistic, fitting a model over separate samples, or bootstrapping standard errors. Within the looping structure, you may post certain numeric values to the postfile. This will create a separate Stata binary data set, which may then be opened in a later Stata run and analysed. Note, however, that only numeric expressions may be written to the postfile, and the parens () given in the documentation, surrounding each exp., are required.





### Reconfiguring data

Data are often provided in a different orientation than that required for statistical analysis. The most common example of this occurs with panel, or longitudinal, data, in which each observation conceptually has both cross-section (i) and time-series (t) subscripts. Often one will want to work with a "pure" cross-section or "pure" time-series. If the microdata themselves are the objects of analysis, this can be handled with sorting and a loop structure. If you have data for N firms for T periods per firm, and want to fit the same model to each firm, one could use the statsby command, or if more complex processing of each model's results was required, a foreach block could be used. If analysis of a cross-section was desired, a bysort would do the job.





But what if you want to use average values for each time period, averaged over firms? The resulting dataset of T observations can be easily created by the collapse command, which permits you to generate a new data set comprised of summary statistics of specified variables. More than one summary statistic can be generated per input variable, so that both the number of firms per period and the average return on assets could be generated. collapse can produce counts, means, medians, percentiles, extrema, and standard deviations.



Different models applied to longitudinal data require different orientations of those data. For instance, seemingly unrelated regressions (sureg) require the data to have T observations ("wide"), with separate variables for each cross—sectional unit. Fixed—effects or random-effects regression models  $\mathtt{xtreg}$ , on the other hand, require that the data be stacked or "vec"'d in the "long" format. It is usually much easier to generate transformations of the data in stacked format, where a single variable is involved.

The reshape command allows you to transfer the data from the former ("wide") format to the latter ("long") format or vice versa. It is a complicated command, because of the many variations on this process one might encounter, but it is very powerful.





Different models applied to longitudinal data require different orientations of those data. For instance, seemingly unrelated regressions (sureg) require the data to have T observations ("wide"), with separate variables for each cross—sectional unit. Fixed—effects or random-effects regression models  $\mathtt{xtreg}$ , on the other hand, require that the data be stacked or "vec"d in the "long" format. It is usually much easier to generate transformations of the data in stacked format, where a single variable is involved.

The reshape command allows you to transfer the data from the former ("wide") format to the latter ("long") format or vice versa. It is a complicated command, because of the many variations on this process one might encounter, but it is very powerful.



As an example, a dataset from the World Bank, provided as a spreadsheet, has rows labelled by both country (ccode) and variable (vcode), and columns labelled by years. Two applications of reshape were needed to transfer the data to the desired long format, where the observations have both country and year subscripts, and the columns are variables:

```
reshape long d, i(ccode vcode) j(year)
reshape wide d, i(ccode year) j(vcode) string
```

The resulting data set is in the appropriate format for xtreg modelling. If it were to be used in sureg—type models, a further reshape wide could be applied to transform it into that format.

See Stata Tip 45, Baum and Cox, Stata Journal 7:2, 2007.



As an example, a dataset from the World Bank, provided as a spreadsheet, has rows labelled by both country (ccode) and variable (vcode), and columns labelled by years. Two applications of reshape were needed to transfer the data to the desired long format, where the observations have both country and year subscripts, and the columns are variables:

```
reshape long d, i(ccode vcode) j(year) reshape wide d, i(ccode year) j(vcode) string
```

The resulting data set is in the appropriate format for xtreg modelling. If it were to be used in sureg—type models, a further reshape wide could be applied to transform it into that format.

See Stata Tip 45, Baum and Cox, Stata Journal 7:2, 2007.





#### Repeating commands

One of Stata's great strengths is the ability to perform repetitive tasks without spelling out the details (e.g. the  $\mathtt{by}$  prefix). However, the  $\mathtt{by}$  prefix can only execute a single command; so that while you may run a regression for each country in your sample, you cannot also save the residuals or predicted values for those country-specific regressions.

Stata provides two commands that allow construction of a true block structure or loop: foreach and forvalues. These commands permit a delimited block of commands to be repeated over elements of a varlist or numlist. Indeed, the target of foreach may be any list of names, and can be a list of new variables to be created. The forvalues numlist may include an increment, so that it could for instance count from 10 to 100 in steps of 10, or count down from 10 to 1.

#### Repeating commands

One of Stata's great strengths is the ability to perform repetitive tasks without spelling out the details (e.g. the by prefix). However, the by prefix can only execute a single command; so that while you may run a regression for each country in your sample, you cannot also save the residuals or predicted values for those country-specific regressions.

Stata provides two commands that allow construction of a true block structure or loop: foreach and forvalues. These commands permit a delimited block of commands to be repeated over elements of a varlist or numlist. Indeed, the target of foreach may be any list of names, and can be a list of new variables to be created. The forvalues numlist may include an increment, so that it could for instance count from 10 to 100 in steps of 10, or count down from 10 to 1

This code fragment loops over a varlist, calculates (but does not display) the descriptives of each variable, and then summarizes the observations of that variable that exceed its mean. Note the use of 'var', in particular the backtick (') on the left of the word. This syntax is mandatory when referring to the placeholder.

```
foreach var of varlist pri-rep t* {
         quietly summarize 'var'
         summarize 'var' if 'var' > r(mean)
}
```

Generally a forvalues or foreach loop is the best way to solve any programming problem that involves repetition. It is usually much faster, in the long run, to figure out how to place a problem in this context. Nested loops may also be defined with these commands.

This code fragment loops over a varlist, calculates (but does not display) the descriptives of each variable, and then summarizes the observations of that variable that exceed its mean. Note the use of 'var', in particular the backtick (') on the left of the word. This syntax is mandatory when referring to the placeholder.

```
foreach var of varlist pri-rep t* {
        quietly summarize 'var'
        summarize 'var' if 'var' > r(mean)
}
```

Generally a forvalues or foreach loop is the best way to solve any programming problem that involves repetition. It is usually much faster, in the long run, to figure out how to place a problem in this context.

Nested loops may also be defined with these commands.

A loop structure may also be explicitly defined by the while command, which is akin to the "do while" construct in other programming languages. A while structure often will make use of an if command—not to be confused with the if clause on other commands—which will create conditional logic. The if command may also use an else clause to express conditional logic. For many purposes, it is more efficient (in terms of your time) to employ foreach or forvalues, since those commands handle the logic of repetition without explicit detail. Programs written with these commands are easier to maintain and modify.





#### Local macros and scalars

In programming terms, **local macros** and **scalars** are the "variables" of Stata programs (not to be confused with the variables of the data set). The distinction: a local macro can contain a string, while a scalar can contain a single number (at maximum precision). You should use these constructs whenever possible to avoid creating variables with constant values merely for the storage of those constants. This is particularly important when working with large data sets.

When you want to work with a scalar object—such as a counter in a foreach or forvalues command—it will involve defining and accessing a local macro. In addition, all Stata commands that compute results or estimates generate one or more objects to hold those items, which are scalars (numbers), local macros (strings) or matrices.



#### Local macros and scalars

In programming terms, **local macros** and **scalars** are the "variables" of Stata programs (not to be confused with the variables of the data set). The distinction: a local macro can contain a string, while a scalar can contain a single number (at maximum precision). You should use these constructs whenever possible to avoid creating variables with constant values merely for the storage of those constants. This is particularly important when working with large data sets.

When you want to work with a scalar object—such as a counter in a foreach or forvalues command—it will involve defining and accessing a local macro. In addition, all Stata commands that compute results or estimates generate one or more objects to hold those items, which are scalars (numbers), local macros (strings) or matrices.



This behavior of Stata's computational commands allows you to write a do-file that makes use of these quantities. We saw one example of this above. The command  $\mathtt{summarize}$  generates a number of scalars, such as  $\mathtt{r}(\mathtt{N})$ , the number of observations;  $\mathtt{r}(\mathtt{mean})$ , the mean;  $\mathtt{r}(\mathtt{Var})$ , the variance; etc. The available items are shown by  $\mathtt{return}$  list for a "r-class" command. The contents of these scalars may be used in expressions. In the example above, the mean of a variable was used to govern a following  $\mathtt{summarize}$  command:

```
quietly summarize price
summarize price if price > r(mean)
```





In this example, the scalar  $r \pmod{m}$  may be used directly. But what if you wanted to issue another command that generated results, which would wipe out all of the  $r \pmod{m}$  returns? Then you use the local statement to preserve the item in a macro of your choosing:

local mu r(mean)

Later in the program, you could use

regress mpg weight length if price > 'mu'

Note the use of the backtick (') on the left of the local macro. This syntax is mandatory, as it makes the *dereference* clear: you are referring to the *value* of the local macro mu rather than the contents of the variable mu.



In this example, the scalar r (mean) may be used directly. But what if you wanted to issue another command that generated results, which would wipe out all of the r() returns? Then you use the local statement to preserve the item in a macro of your choosing:

```
local mu r (mean)
```

Later in the program, you could use

```
regress mpg weight length if price > 'mu'
```

Note the use of the backtick (') on the left of the local macro. This syntax is mandatory, as it makes the dereference clear: you are referring to the value of the local macro mu rather than the contents of the variable mil.





#### Returned results

Stata commands are either *r-class* commands like <code>summarize</code>, that return results, or *e-class* commands, that return estimates. You may examine the set of results from a r-class command with the command <code>return list</code>. For an e-class command, use <code>ereturn list</code>. An e-class command will return <code>e()</code> scalars, macros and matrices: for instance, after <code>regress</code>, the local macro <code>e(N)</code> will contain the number of observations, <code>e(r2)</code> the  $R^2$  value, <code>e(depvar)</code> will contain the name of the dependent variable, and so on.

Commands may also return matrices. For instance, regress (like all estimation commands) will return the matrix e(b), a row vector of point estimates, and the matrix e(V), the estimated variance—covariance matrix of the estimated parameters.



#### Returned results

Stata commands are either *r-class* commands like <code>summarize</code>, that return results, or *e-class* commands, that return estimates. You may examine the set of results from a r-class command with the command <code>return list</code>. For an e-class command, use <code>ereturn list</code>. An e-class command will return <code>e()</code> scalars, macros and matrices: for instance, after <code>regress</code>, the local macro <code>e(N)</code> will contain the number of observations, <code>e(r2)</code> the  $R^2$  value, <code>e(depvar)</code> will contain the name of the dependent variable, and so on.

Commands may also return matrices. For instance, regress (like all estimation commands) will return the matrix e(b), a row vector of point estimates, and the matrix e(V), the estimated variance—covariance matrix of the estimated parameters.



Use  $\mathtt{display}$  to examine the contents of a scalar or local macro. For the latter, you must use the backtick and apostrophe to indicate that you want to access the contents of the macro: contrast  $\mathtt{display}$  r (mean) with  $\mathtt{display}$  "The mean is ` mu' ". The contents of matrices may be displayed with the matrix list command.

Since items are accessible in local macros, it is very easy to write a program that makes use of results in directing program flow. Local macros can be created by the local statement, and used as counters (e.g. in foreach).

For more information, see my separate slideshow *Why should you become a Stata programmer?* 



Use  $\mathtt{display}$  to examine the contents of a scalar or local macro. For the latter, you must use the backtick and apostrophe to indicate that you want to access the contents of the macro: contrast  $\mathtt{display}$  r (mean) with  $\mathtt{display}$  "The mean is ` mu' ". The contents of matrices may be displayed with the matrix list command.

Since items are accessible in local macros, it is very easy to write a program that makes use of results in directing program flow. Local macros can be created by the local statement, and used as counters (e.g. in foreach).

For more information, see my separate slideshow *Why should you become a Stata programmer?* 





August 2009

## Some useful Stata commands

help: online help on a specific command

findit: online references on a keyword or topic ssc: access routines from the SSC Archive

log: log output to an external file

tsset: define the time indicator for timeseries or panel data

compress: economize on space used by variables

pwd : print the working directorycd : change the working directory

clear : clear memory

quietly: do not show the results of a command

update query: see if Stata is up to date

adoupdate: see if user-written commands are up to date

exit: exit the program (,clear if dataset is not saved)





# Data manipulation commands

generate : create a new variable

replace: modify an existing variable

rename: rename variable

renvars: rename a set of variables

sort : change the sort order of the dataset

drop: drop certain variables and/or observations

keep: keep only certain variables and/or observations

append: combine datasets by stacking

merge: merge datasets (one-to-one or match merge)

encode: generate numeric variable from categorical variable

recode: recode categorical variable

destring: convert string variables to numeric

foreach: loop over elements of a list, performing a block of code

forvalues : loop over a numlist, performing a block of code

local: define or modify a local macro (scalar variable)



describe: describe a data set or current contents of memory

use : load a Stata data set

save : write the contents of memory to a Stata data set

insheet: load a text file in tab- or comma-delimited format

infile: load a text file in space-delimited format or as defined in a

dictionary

outfile: write a text file in space- or comma-delimited format outsheet: write a text file in tab- or comma-delimited format

contract: make a dataset of frequencies

collapse : make a dataset of summary statistics

tab: abbreviation for tabulate: 1- and 2-way tables

table: tables of summary statistics





## Statistical commands

summarize : descriptive statistics

correlate: correlation matrices

ttest: perform 1-, 2-sample and paired t-tests

anova: 1-, 2-, n-way analysis of variance

regress: least squares regression

predict: generate fitted values, residuals, etc.

test: test linear hypotheses on parameters lincom: linear combinations of parameters

cnsreg: regression with linear constraints

testnl: test nonlinear hypothesis on parameters

mfx: marginal effects (elasticities, etc.)

ivregress: instrumental variables regression

prais: regression with AR(1) errors

sureg: seemingly unrelated regressions

reg3: three-stage least squares

qreg: quantile regression



89 / 132

## Limited dependent variable estimation commands

logit, logistic : logit model, logistic regression

probit: binomial probit model

tobit: one- and two-limit Tobit model

cnsreg: Censored normal regression (generalized Tobit)

ologit, oprobit : ordered logit and probit models

mlogit: multinomial logit model poisson: Poisson regression heckman: selection model





### **Time series estimation commands**

arima: Box-Jenkins models, regressions with ARMA errors

arch: models of autoregressive conditional heteroskedasticity

dfgls: unit root tests

corrgram : correlogram estimation

var: vector autoregressions (basic and structural)

irf: impulse response functions, variance decompositions

vec : vector error-correction models (cointegration)

rolling: prefix permitting rolling or recursive estimation over subsets





### Panel data estimation commands

xtreg,fe: fixed effects estimator

xtreg,re: random effects estimator

xtgls: panel-data models using generalized least squares

xtivreg: instrumental variables panel data estimator

xtlogit : panel-data logit models

xtprobit : panel-data probit models

xtpois: panel-data Poisson regression

xtgee: panel-data models using generalized estimating equations

xtmixed : linear mixed (multi-level) models

xtabond : Arellano-Bond dynamic panel data estimator





## **Nonlinear estimation commands**

The nl command may be used to estimate a nonlinear model, while ml supports maximum likelihood estimation with a user-specified likelihood function. See my separate slideshow on *Maximum Likelihood Estimation and Nonlinear Least Squares in Stata*.

Mata now contains a full-featured set of optimization commands as  ${\tt optimize}$  ( ). These commands are now the preferred method to implement optimization in Stata.





# **Graphics commands:**

twoway produces a variety of graphs, depending on options listed histogram rep78 histogram of this categorical variable twoway scatter price mpg a Y vs X scatterplot twoway line price mpg a Y vs X line plot tsline GDP a Y vs time time-series plot twoway area price mpg an Y vs X area plot twoway rline price mpg a Y vs X range plot (hi-lo) with lines The command twoway may be omitted in most cases.





The flexibility of Stata graphics allows any of these plot types (including many more that are available) to be easily combined on the same graph. For instance, using the auto.dta dataset,

```
twoway (scatter price mpg) (lfit price mpg)
```

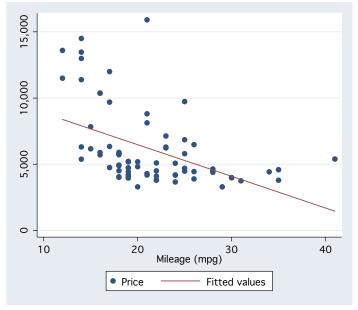
will generate a scatterplot, overlaid with the linear regression fit, and

```
twoway (lfitci price mpg) (scatter price mpg)
```

will do the same with the confidence interval displayed.

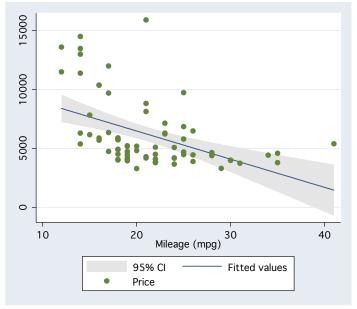
















A nonparametric fit of a bivariate relationship can be readily overlaid on a graph via

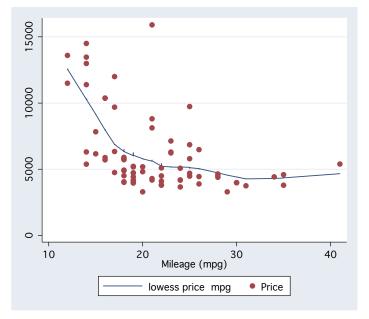
```
twoway (lowess price mpg) (scatter price mpg)
```

Twoway graphs may also represent mathematical functions, without explicit data:

```
twoway (function y=log(x)*sin(x)) (function y=x*cos(x))
```

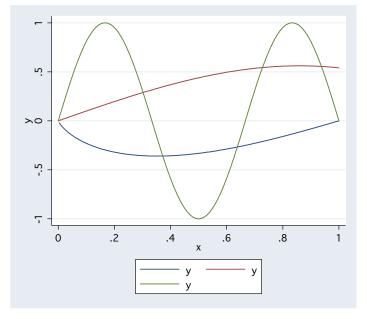














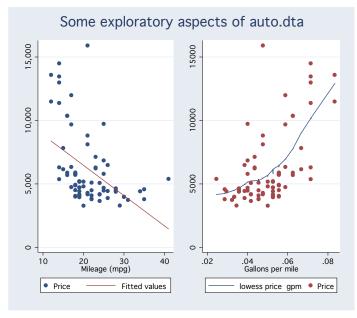


# Graphs may also be readily combined into a single graphic for presentation. For instance,

```
twoway (scatter price mpg) (lfit price mpg), name(auto1)
gen gpm = 1/mpg
label var gpm "Gallons per mile"
twoway (lowess price gpm) (scatter price gpm),
name(auto2)
graph combine auto1 auto2, saving(myauto, replace) ///
ti("Some exploratory aspects of auto.dta")
```

where the "///" is a continuation of the line.









### Instructional data sets

A list of over 100 datasets suitable for instructional use is available on the economics web pages as

http://fmwww.bc.edu/ec-p/data/ecfindata.html#teach

#### Sample Stata do-files

Consider the data Zvi Griliches used in his 1976 article on the wages of young men (*Journal of Political Economy*, 84, S69-S85). These are cross-sectional data on 758 individuals collected over several survey years.

do http://fmwww.bc.edu/ec-p/software/stata/stataintrol



```
* StataIntro: cross-section example
log using introl, replace
use http://fmwww.bc.edu/ec-p/data/hayashi/griliches76
describe
summarize
label define ur 0 rural 1 urban
label values smsa ur
tab smsa
tab mrt smsa, chi2
ttest med, by (smsa)
anova lw mrt smsa
anova lw mrt smsa mrt*smsa
```



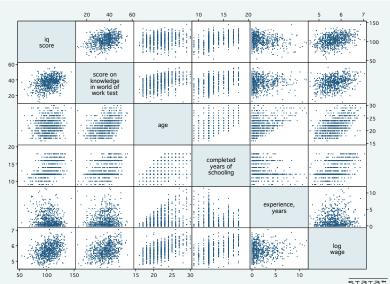


anova, regress

```
regress lw tenure kww smsa
predict lweps,resid
scatter lweps kww
bysort year: regress lw tenure kww smsa
graph matrix iq kww age s expr lw, msize(tiny)
gen medrural = med*(smsa==0)
gen medurban = med*(smsa==1)
regress lw tenure kww medurban medrural
test medurban=medrural
log close
```











The following example reads some daily Dow-Jones Averages data, graphs daily returns, then performs Dickey-Fuller tests for unit roots on the DJIA, its log, and its returns (log price relatives), and on their first differences. AR(3) models are then estimated on the series, and the Box–Pierce portmanteau test is then performed on the residuals.

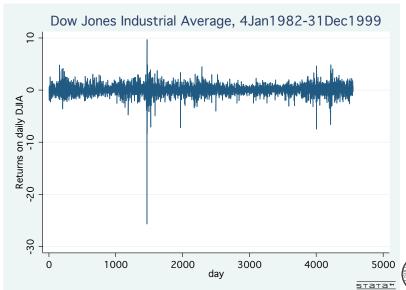
In this example, we make use of "local macros" (with values 'v'), which enable us to perform the same operations on several named variables without having to write out the commands for each variable. This facility may be used with <code>varlists</code> of any length, and makes it very easy to generate parallel analyses, produce graphs, etc. for an arbitrary set of variables or time periods.

do http://fmwww.bc.edu/ec-p/software/stata/stataintro2



```
* StataIntro: time-series example
log using intro2, replace
use http://fmwww.bc.edu/ec-p/data/micro/ddjia.dta
desc
summ
tsset
tsline ret
foreach v of varlist djia ldjia ret {
      dfgls 'v', maxlag(12)
      dfgls D.'v', maxlag(12)
      regress v' L(1/3) v', robust
      predict eps_'v', resid
      wntesta eps 'v'
log close
```







### **Examples of Stata programming**

Let us form a "rolling forecast" of volatility from a moving-window regression (we had not learned that Baum's rolling command (or Stata's rolling: prefix) could do this job for us). Assume that we have 120 time-series observations which have been tsset:

```
gen volfc=.
local win 12
forv i=13/120 {
    local first = 'i'-'win'+4
    quietly regress y L(1/4).y in 'first'/'i'
    quietly replace volfc = e(rmse) in 'i'/'i'
}
```

This program will generate the series volfc as the RMS error of an AR(4) model fit to a window of 12 observations for the y series.



The use of local macros and the appropriate loop constructs make it possible to write a Stata program that is fairly general, and requires little modification to be reused on different series, or with different parameters. This makes your work with Stata very productive, since much of the code is reusable and adaptable to similar tasks. Let us consider how this approach might be pursued in the context of the volatility forecast example.

[ For more information, see my separate slideshow *Why should you become a Stata programmer?* and my book *An Introduction to Stata Programming*, available in O'Neill Library. ]



We show here a complete Stata program, volfc, which is stored in the file volfc.ado on the adopath. Since this is a personally-authored program, it should be placed in the personal subdirectory of the ado directory (not the Stata directory's ado subdirectory!) For more information, see adopath.

This program makes use of Stata's syntax parsing capabilities to allow this user-written command to emulate all Stata commands' syntax. It does not make use of many of the features that might be useful in such a command: handling if and in clauses, providing more specific error messages for inappropriate option values, and so on.



We show here a complete Stata program, volfc, which is stored in the file volfc.ado on the adopath. Since this is a personally-authored program, it should be placed in the personal subdirectory of the ado directory (not the Stata directory's ado subdirectory!) For more information, see adopath.

This program makes use of Stata's syntax parsing capabilities to allow this user-written command to emulate all Stata commands' syntax. It does not make use of many of the features that might be useful in such a command: handling if and in clauses, providing more specific error messages for inappropriate option values, and so on.



The program generalizes the do-file shown above by allowing the moving—window volatility estimate to be generated from a specified variable, and placed in a new variable specified in the vol() option. The window width (option win()) and AR length (option AR()) take on default values 12 and 4, but may be overridden by the user. The program automatically calculates the first and last observations to be used in the loop from the data and specified options. It could readily be generalized to use a different volatility measure from the rolling regression (e.g. mean absolute error).

To be complete, we should provide a help file for volfc in the file volfc.sthlp. The help file would specify the syntax of the command, explain its purpose, define each of the options, and provide any references to other Stata commands that might be useful.



August 2009

```
program define volfc, rclass
version 10.0
syntax varname(numeric) , Vol(string) [Win(integer 12)
quietly tsset
if 'win' < 'ar' {
   di "You must have a longer window than AR length!"
   error 198
quietly gen 'vol' =.
local start = 'win'+'ar'
quietly summ 'varlist', meanonly
local last = r(N)
dis _n "'vol': volatility forecast for 'varlist' with
         window='win', AR('ar')"
```

(continues...)





## This program defines the volfc command, which will appear like any other Stata command on your machine. It may be executed as

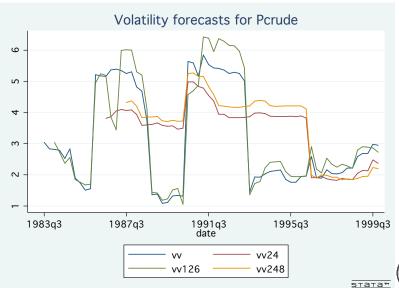
```
use http://fmwww.bc.edu/ec-p/data/macro/bdh, clear
volfc pcrude, vol(vv)
volfc pcrude, vol(vv24) win(24)
volfc pcrude, vol(vv126) ar(6)
volfc pcrude, vol(vv248) win(24) ar(8)
```

# The volatility series might then be graphed (presuming a time variable date which is the variable that has been tsset) with

```
tsline vv vv24 vv126 vv248 if tin(1983q1,), ti(Volatility forecasts for Pcrude)
```











This illustrates the relative simplicity of developing a quite general tool in Stata's programming language. Although you may use Stata without ever authoring an "ado-file", much of the productivity enhancement that a Stata user may enjoy is likely to be tied to this sort of development. Many research tasks are quite repetitive in some context, and developing a general-purpose tool to implement that repetition is likely to be a very good investment in terms of time and effort.

Many of the modules available from the SSC Archive were first conceived by individuals looking to ease the burden of their own work. Stata's unique extensibility makes it trivial to incorporate user-written additions—including those which you author—into your copy of Stata, and to share it with collaborators or the Stata user community if desired.

As should be evident from this programming example, the program define command is used to declare a program. The program name must match the name of the ado-file in which it is stored. Most user-written programs are r-class. This program could be modified to return its parameters to the calling program with the return statement:

```
return local vol 'vol'
return local win 'win'
return local ar 'ar'
return local first 'start'
return local last 'last'
```

With these statements added to the end of the routine, the local macros are defined, and their values stored.



The second element to be noted is the <code>syntax</code> statement, which defines the allowable syntax for a user-written command. One may specify that the command allows a single variable, with <code>varname</code>; a set of variables, with <code>varlist</code>, optionally specifying how many are allowed. For instance, a statistical technique that operates on a pair of variables could specify that exactly two existing variables are to be provided. Likewise, one may specify that a new variable (or set of variables) are the <code>newvarlist</code> of the command, and syntax will check that they are indeed new variables.

Although not illustrated above, the syntax command will often specify that if and in clauses are optional elements. Optional elements of syntax (such as the options Win and AR above) are placed in brackets ([]).



The second element to be noted is the syntax statement, which defines the allowable syntax for a user-written command. One may specify that the command allows a single variable, with varname; a set of variables, with varlist, optionally specifying how many are allowed. For instance, a statistical technique that operates on a pair of variables could specify that exactly two existing variables are to be provided. Likewise, one may specify that a new variable (or set of variables) are the newvarlist of the command, and syntax will check that they are indeed new variables.

Although not illustrated above, the syntax command will often specify that if and in clauses are optional elements. Optional elements of syntax (such as the options Win and AR above) are placed in brackets ([]).



This programming example illustrates a "required option"—the vol option, which must be used on this command to specify the output of the command. The other two options are indeed optional, and take on default values if they are not specified. The argument of the vol option is meant to be a new variable name; that will be trapped when the generate statement attempts to create the variable if it is already in use, or is not a valid variable name.

Most user-written programs could be improved by adding code to trap errors in users' input. If the program is primarily for your own use, you may eschew extensive development of error trapping: for instance, checking the options for sensibility (although one test is applied here to prevent nonsensical results).



This programming example illustrates a "required option"—the vol option, which must be used on this command to specify the output of the command. The other two options are indeed optional, and take on default values if they are not specified. The argument of the vol option is meant to be a new variable name; that will be trapped when the generate statement attempts to create the variable if it is already in use, or is not a valid variable name.

Most user-written programs could be improved by adding code to trap errors in users' input. If the program is primarily for your own use, you may eschew extensive development of error trapping: for instance, checking the options for sensibility (although one test is applied here to prevent nonsensical results).



Local macros are exactly that: objects with local scope, defined within the program in which they are used, disappearing when that program terminates. This is generally the desired outcome, preventing a clutter of objects from being retained when a program calls numerous others in the course of execution. At times, though, it is necessary to have objects that can be passed from one subprogram to another. The return logic above would not really serve, since although it passes local macros from a program to its caller, they would then have to be passed as arguments to a second program.

To deal with the need for persistent objects, Stata contains *global macros*. These objects, once defined, live for the duration of your Stata session, and may be read or written within any Stata program. They are defined with the global command, rather than local, and referred to as \$macroname. Global macros should only be used where they are required.

Local macros are exactly that: objects with local scope, defined within the program in which they are used, disappearing when that program terminates. This is generally the desired outcome, preventing a clutter of objects from being retained when a program calls numerous others in the course of execution. At times, though, it is necessary to have objects that can be passed from one subprogram to another. The return logic above would not really serve, since although it passes local macros from a program to its caller, they would then have to be passed as arguments to a second program.

To deal with the need for persistent objects, Stata contains *global macros*. These objects, once defined, live for the duration of your Stata session, and may be read or written within any Stata program. They are defined with the global command, rather than local, and referred to as \$macroname. Global macros should only be used where they are required.

We now present an example of a Stata program that operates on panel, or longitudinal data. When you use panel data, you must use the panel data form of tsset in which both a unit variable and a time variable are specified.

Assume that you have a panel data set, properly identified as such, containing several time series for each unit in the panel: for instance, investment or population measures for several countries. We would like to generate a new series containing the deviations from a constant growth path (exponential trend) or, alternatively, the constant growth values themselves (the predicted values from the exponential trend line).



We now present an example of a Stata program that operates on panel, or longitudinal data. When you use panel data, you must use the panel data form of tsset in which both a unit variable and a time variable are specified.

Assume that you have a panel data set, properly identified as such, containing several time series for each unit in the panel: for instance, investment or population measures for several countries. We would like to generate a new series containing the deviations from a constant growth path (exponential trend) or, alternatively, the constant growth values themselves (the predicted values from the exponential trend line).





This program, pangrodev, performs this task for each unit of a panel, automatically identifying the observations belonging to each unit, taking the logarithm of the specified variable, running the appropriate regression and prediction commands, and assembling the results in the specified new variable.

The program makes use of Stata's tempname and tempvar commands to create non-scalar objects (in this case the matrix VV and variables lvar and pvar which, like local macros, will exist only for the duration of the ado-file). These temporary facilities, like the associated tempfile which allows temporary files to be specified, help reduce clutter and guarantee that objects' names will not conflict with other items in the user's namespace.



This program, pangrodev, performs this task for each unit of a panel, automatically identifying the observations belonging to each unit, taking the logarithm of the specified variable, running the appropriate regression and prediction commands, and assembling the results in the specified new variable.

The program makes use of Stata's tempname and tempvar commands to create non-scalar objects (in this case the matrix VV and variables lvar and pvar which, like local macros, will exist only for the duration of the ado-file). These temporary facilities, like the associated tempfile which allows temporary files to be specified, help reduce clutter and guarantee that objects' names will not conflict with other items in the user's namespace.



```
*! pangrodev 1.1.0 CFBaum 21Jan2006
* generate deviations from constant growth in panel
* 1.1.0: promote to v9, use levelsof
program define pangrodev, rclass
version 10.0
syntax varname, Gen(string) [xb]
local togens "deviations from constant growth"
if "'xb'" != "" {
local togens "predicted growth"
qui tsset
local ivar = r(panelvar)
local timevar = r(timevar)
tempname VV
tempvar lvar pvar
```

(continues...)



```
qui gen double 'lvar' = log('varlist')
* get list of units
qui levelsof 'ivar', local(vals)
local nuals: word count 'vals'
qui gen double 'gen' =.
local xc 0
local tbar 0
local rsqr 0
```



```
foreach v of local vals {
summ 'lvar' if 'ivar' == 'v', meanonly
if r(N) > 2 {
qui regress 'lvar' 'timevar' if 'ivar' == 'v'
capt drop 'pvar'
qui predict double 'pvar' if e(sample), xb
qui replace 'qen' = exp('pvar') if e(sample)
if "'xb'" =="" {
qui replace 'gen' = 'varlist'-'gen' if e(sample)
local xc = 'xc' + 1
local tbar = 'tbar' + e(N)
local rsqr = 'rsqr' + e(r2)
```



```
local tbar = int(100*'tbar' / 'xc')/100.0
local rsqr = int(1000*'rsqr' / 'xc')/1000.0
di in gr _n "'gen' : 'togens' for 'xc' of " ///
                 "'nvals' units"
di in qr "tbar = 'tbar' rsq-bar = 'rsqr'"
exit.
end
```



### This program defines the pangrodev command, which will appear like any other Stata command on your machine. It may be executed as

```
. use http://fmwww.bc.edu/ec-p/data/macro/cap797wa
(World Bank Database for Sectoral Investment, 1948-199
. pangrodev TotSECap, g(totcapdev)
```

```
totcapdev: deviations from constant growth for
                 57 of 63 units
tbar = 25.94 rsq-bar = .673
```

. pangrodev TotSECap, q(totcaphat) xb (output omitted)

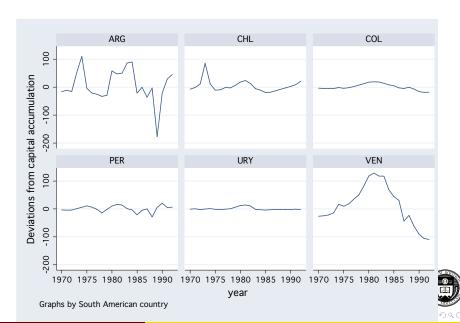


Selected series computed by pangrodev can now be graphed by the tsline command, which accepts a by (*varlist*) option:

```
replace totcapdev=totcapdev/10^9
keep if (ccode=="ARG" | ccode=="CHL" | ccode=="COL"
ccode=="PER" | ccode=="URY" | ccode=="VEN")
label var totcapdev "Deviations from capital accum"
label var ccode "South American country"
tsline totcapdev if year>1969, by(ccode)
```

will demonstrate how many countries followed the same pattern of below-trend growth of the capital stock (curtailed investment) during the 1980s.





Whether or not you use Stata's programming facilities to write your own ado-files, a "reading knowledge" of the programming language is very useful in case you want to adapt an existing Stata command (official or user-contributed) in a do-file you are writing.

Since the code for all Stata commands that are implemented as ado-files (as the command which... will show) are available on your hard disk, Stata itself is a fertile source of programming techniques that may be adapted to solve any programming problem.

For a thorough treatment of the subject, see my book *An Introduction to Stata Programming* (2009) in O'Neill Library.





Whether or not you use Stata's programming facilities to write your own ado-files, a "reading knowledge" of the programming language is very useful in case you want to adapt an existing Stata command (official or user-contributed) in a do-file you are writing.

Since the code for all Stata commands that are implemented as ado-files (as the command  $\mbox{which...}$  will show) are available on your hard disk, Stata itself is a fertile source of programming techniques that may be adapted to solve any programming problem.

For a thorough treatment of the subject, see my book *An Introduction to Stata Programming* (2009) in O'Neill Library.

