Stata

Popular program in Economics, Medicine, ???

The previous versions were command line based. The current version (8) has switched to more of a graphical menu driven interface.

However, everything can be done from a command line. Somethings have to be done from the command line.

Has very good Survey Sampling routines

More of a canned program, though you can write your own routines, and even link in your own C code (with version 8.1).

While you can do some programming with it, and a lot of people apparently do, I would not use Stata unless I was performing a "standard" analysis. Use S-Plus/R in those situations.

Has excellent backward compatibility with previous versions

Available for Windows, Macintosh, Unix (not available at Harvard that I can find).



The default setup will look similar to the above. However you can change the setup, font sizes, colours, etc. These will be kept from session to session.

👪 I	nterc	ooled	l Stat	a 8.0					
<u>F</u> ile	<u>E</u> dit	<u>P</u> refs	<u>D</u> ata	<u>Graphics</u>	Statistics	<u>U</u> ser	Window	<u>H</u> elp	
F	8 8	91	۵ 🔜			8			

👪 Intercooled Stat	a 8.0			
File Edit Prefs Data	Graphics Statistics User Window Help			
F B B M B	Easy graphs	Scatter plot		
Stata Results	Twoway graph (scatterplot, line, etc.) Overlaid twoway graphs	Connected scatter plot Line graph Area graph		
. insheet using "	Bar chart	Overlaid twoway graphs		
(26 vars, 93 obs)	Horizontal bar chart Dot chart Pie chart	Bar chart Horizontal bar chart Dot chart		
. insheet using "	Histogram Box plot	Pie chart (by variables) Pie chart (by category)		
(26 vars, 93 obs)	Horizontal box plot	Histogram Box plot		
(Scatterplot matrix			
 regress minpric 	Distributional graphs	Horizontal box plot		
Source	Smoothing and densities	Scatterplot matrix		
Model 5	Time series graphs	Regression fit Function graph		
Total 7	ROC analysis Quality control More statistical graphs	Adj R-squared Adj R-squar Root MSE		
minprice	Table of graphs	P> t [95% Con		
maxprice _cons	Manage graphs Change scheme/size Graph preferences	0.000 .6493481 0.1113239184		

🔣 Intercooled Stata 8.0		
File Edit Prefs Data Graphics	Statistics User Window Help	
FR B X B B C C	Summaries, tables, & tests	•
	Linear regression and related	٠.
🗖 Stata Results	Binary outcomes	•
drop _all	Ordinal outcomes	*
set obs `obs'	Count outcomes	1
tempvar z	Categorical outcomes	۲
gen z' = exp(z)	Selection models	9
summarize `z'	Generalized linear models (GLM)	1
return scalar m	Nonparametric analysis	۲
return scalar V	Time series	-
end	Multivariate time ceries	
	Cross-sectional time series	
. simulate "lnsim, o	Cross-sectional time series	25
	Survival analysis	•
make a dataset conta a log-normal distrib	Observational/Epi. analysis	•
7. Perform the expe	Survey data analysis	٠
simulate "lnsim o	ANOVA/MANOVA	
reps (10000)	Cluster analysis	•
	Other multivariate analysis	۲
Also see	Resampling & simulation	•
	General post-estimation	٠
Manual: [R] simulat	Other	•

Data Input:

👪 Intercooled	Stata 8.0	
File Edit Prefs	Data Graphics Statistics User	Window Help
🛎 🖬 🎒 🔗	Describe data	Describe variables in memory
Stata Resu	Data editor Data browser (read-only editor) Create or change variables	Describe variables in file Describe data contents (codebook) Inspect variables
> at 104\93c (26 vars, 93	Sort Combine datasets	 List data Compactly list variable names
. browse	Labels & notes	Summary statistics
	Variable utilities	•
 insheet us at 104\93c 	Matrices	ings\lrwin\My Documents\Har
(26 vars, 93	Other utilities	•

Stata can only deal with two file formats, its own, and text files for data.

Its own file format has changed over time, but the current versions can read older versions. However the other way doesn't work (e.g. Version 7 can read a version 8 file)

Reading in text files

Totercooled Stata 8.0		💻 insheet - Import ASCII data	×
File Edit Prefs Data Graphics	s Statistics User Window Help	ASCII dataset filename:	Prowney
Open Ctrl+O View Save Ctrl+S		New variable names: (optional)	Browse
Save As Shift+Ctrl+S Do Filename		Options Store variables as doubles	
Log		Specify value delimiter:	
Import	ASCII data created by a spreadsheet	Comma	
Open Graph	ASCII data in fixed format ASCII data in fixed format with a dictionary Unformatted ASCII data		
Print Graph Print Results	9 12 46 11.4 11.2 46.0	Replace data in memory	
Exit Alt+F4	14 10 60	OK Cancel	Submit

insheet using "C:\Documents and Settings\Irwin\My
Documents\Harvard\Courses\Stat 104\93cars.txt", clear

When creating commands with Stata, the command you would need to type is given in the Results window and the Review window.

The Review window can be used to either repeat commands (double click on the command) or as a start to create a new command (single click).

A single click will put the command into the Commands window and it can be edited

Also the previous version of a command is kept in the dialogue boxes so you can also edit things that way as well. The **I** button in a command dialogue box will clear the entries.

The underlying data structure is similar to a Spreadsheet. Rows correspond to observations and colums correspond to variables. You can also think of it like a S-Plus/R data frame.

In Stata you can only have one data file in use at a time. Its not like S-Plus/R where you can have many data frames available.

It is possible to combine datasets to get around this limitation.

All you can really read in is a table with r rows and c columns.

However the columns can be a mixture of numbers and strings. Just about any file that read.table() will work with, the Stata insheet command can handle.

```
. insheet using "filename.raw", clear
```

The default file extension for text files being read in is .raw, though any can be used.

The clear option is there to force the current data file to be deleted from memory and to be replace by the new file. If not, you will get an error message similar to

insheet using "Berkeley.txt"

you must start with an empty dataset

r(18);

While the expected column separators are either commas or tabs, any character can be used.

Exporting Data



outsheet using " Berkeley.raw", replace

outsheet using " Berkeley-men.raw" if gender == "Men", replace

Defining variables

👪 Intercooled	Stata 8.0	
File Edit Prefs	Data Graphics Statistics User	Window Help
	Describe data	
🗖 Stata Resu	Data editor Data browser (read-only editor) Create or change variables	Create new variable
21. F 22. F	Sort Combine datasets	Create new variable (extended) Other variable creation commands
24. F	Labels & notes Variable utilities	Change contents of variable Other variable transformation commands
. table maj	Matrices	count], chi2
option chi2	Other utilities	•

\blacksquare generate - Generate a new variable X	Expression builder	×
Main if/in Generate variable: Contents:	Category: Mathematical	OK Cancel
Generate variable as type: Attach value label:	Mathematical abs() 7 8 Probability distributions acos() asin() atan() 4 5 String atan() atan() 1 2 Programming atan() atan() 1 2 Date ceil() ceil() 0 0 Selecting time-spans cos() 0 0 Matrix Operators digamma() 1 8	9 / == 6 · > 3 · < 1 · < + >= 1 () !=
OK Cancel Submit	abs(x): The absolute value of x.	*

Going through the dialogues will generate a command like

generate float CityFuel= 100/HighMPG generate float CityFuel= 100/CityMPG HighFuel already defined generate can only be used for creating a variable for the first time. If you need to update a variable use replace

```
replace CityFuel= 100/CityMPG
```

```
(93 real changes made)
```

Stata has the property when creating a new variable to place results in every row, even when it may not make intuitive sense. For example

```
generate tcrit= invttail(24,0.025)
```

will place the value 2.063899 in every row of the data file. To store it only in the first row, use

gen tcrit= invttail(24,0.025) in 1

instead. Actually for this sort of calculation, you probably just want

```
disp invttail(24,0.025)
```

2.0638985

which just prints the result in the Results window.

Viewing Data

There are a number of ways that you can view your data in Stata.

The first is the list command

list MinPrice MidPrice MaxPrice in 1/10, separator(5)

-	+		+
	MinPrice	MidPrice	MaxPrice
1.	12.9	15.9	18.8
2.	29.2	33.9	38.7
3.	25.9	29.1	32.3
4.	30.8	37.7	44.6
5.	23.7	30	36.2
6.	14.2	15.7	17.3
7.	19.9	20.8	21.7
8.	22.6	23.7	24.9
9.	26.3	26.3	26.3
10.	33	34.7	36.3
_	+		

The other common approaches are with the Data Editor (edit) and the Data Browser (browse)

They both appear the same, but only with the Editor can you make changes.

Also when either of these is active, it is not possible to run any analyses. You must close the window to proceed.

serve <u>R</u> es	tore <u>S</u> ort <<	>> <u>H</u> ide	Delete						
	M	anu[1] = Acura							
1	Manu	Model	Туре	MinPrice	MidPrice	MaxPrice	CityMPG	HighMPG	AirI
1	Acura	Integra	Small	12.9	15.9	18.8	25	31	
2	Acura	Legend	Midsize	29.2	33.9	38.7	18	25	
3	Audi	90	Compact	25.9	29.1	32.3	20	26	
4	Audi	100	Midsize	30.8	37.7	44.6	19	26	
5	BMW	535i	Midsize	23.7	30	36.2	22	30	
6	Buick	Century	Midsize	14.2	15.7	17.3	22	31	
7	Buick	LeSabre	Large	19.9	20.8	21.7	19	28	
8	Buick	Roadmaster	Large	22.6	23.7	24.9	16	25	
9	Buick	Riviera	Midsize	26.3	26.3	26.3	19	27	
10	Cadillac	DeVille	Large	33	34.7	36.3	16	25	
11	Cadillac	Seville	Midsize	37.5	40.1	42.7	16	25	
12	Chevrolet	Cavalier	Compact	8.5	13.4	18.3	25	36	
13	Chevrolet	Corsica	Compact	11.4	11.4	11.4	25	34	
14	Chevrolet	Camaro	Sporty	13.4	15.1	16.8	19	28	
15	Chevrolet	Lumina	Midsize	13.4	15.9	18.4	21	29	
16	Chevrolet	Lumina_APV	Van	14.7	16.3	18	18	23	
17	Chevrolet	Astro	Van	14.7	16.6	18.6	15	20	
18	Chevrolet	Caprice	Large	18	18.8	19.6	17	26	
19	Chevrolet	Corvette	Sporty	34.6	38	41.5	17	25	
20	Chrylser	Concorde	Large	18.4	18.4	18.4	20	28	
21	Chrysler	LeBaron	Compact	14.5	15.8	17.1	23	28	
22	Chrysler	Imperial	Large	29.5	29.5	29.5	20	26	
23	Dodge	Colt	Small	7.9	9.2	10.6	29	33	
24	Dodge	Shadow	Smal1	8.4	11.3	14.2	23	29	
25	Dodge	Spirit	Compact	11.9	13.3	14.7	22	27	
26	Dodge	Caravan	Van	13.6	19	24.4	17	21	
27	Dodge	Dynasty	Midsize	14.8	15.6	16.4	21	27	

With the data editor, you cannot give variable names within the editor. However you can change them later with the rename command, e.g.

rename Manu Manufacturer

Also many command names can be abbreviated. For example you can use ren for rename, reg for regress, etc. To see what they are, check the help pages. They can be displayed in the main results window, with a command like help regress or in Stata Viewer window (which you get if you click on the help button in the dialogue box)

Stata Viewer [help regress]			×
Back Befresh Search Helpt Contents What's New News			<u> </u>
Jommand: help regress			<i>6</i> 40
help for regress	manual:	[R] regr	ess
	dialogs:	regress	predict
Lincon regression			
hillear regression			
regress depvar [varlist] [weight] [if exp] [in range] [, level(#) beta robust score(newvar) hc2 hc3 hascons noconstant tsscons noheader eform(string plus]	<u>cluster(va</u> g) <u>dep</u> name()	rname) varname) <u>m</u>	sel
by : may be used with regress; see help \underline{by} .			
aweights, fweights, iweights, and pweights are allowed; see help weights.			
depvar and the varlist following depvar may contain time-series operators; see he	lp <u>varlist</u> .		
regress shares the features of all estimation commands; see help estcom.			
regress may be used with sw to perform stepwise estimation; see help $\frac{sw}{sw}$.			
The syntax of predict following regress is			
<pre>predict [type] newvarname [if exp] [in range] [, statistic]</pre>			
			-

The help facility isn't bad. One nice features is that it uses hyperlinks, allowing you to immediately go to another help page of interest.

Graphics

👪 Intercooled Stat	a 8.0				
File Edit Prefs Data	Graphics Statistics User Window Help	2			
🗃 🖬 🍯 🕺 🚳 🔜	Easy graphs	Scatter plot			
Stata Results	Twoway graph (scatterplot, line, etc.) Overlaid twoway graphs	Connected scatter plot Line graph Area graph			
<pre>. insheet using " > at 104\93cars.t (26 vars, 93 obs) . browse</pre>	Bar chart Horizontal bar chart Dot chart Pie chart	Overlaid twoway graphs Bar chart Horizontal bar chart Dot chart			
. insheet using "	Histogram Box plot	Pie chart (by variables) Pie chart (by category)			
(26 vars, 93 obs)	Scatterplot matrix	Histogram Box plot			
 regress minpric 	Distributional graphs	Horizontal box plot			
Source	Smoothing and densities	Scatterplot matrix			
Model 5 Residual 1	Time series graphs	Regression fit Function graph R-squared			
Total 7	Quality control More statistical graphs	Adj R-square Root MSE			
minprice	Table of graphs	P> t [95% Con:			
maxprice _cons	Manage graphs Change scheme/size Graph preferences	0.000 .6493481 0.1113239184			

The quality of the graphs is ok. The resulting graphs are somewhat configurable (fonts, colours, etc), but there are some limitations.

For the common graphs, scatter plots, bar charts, etc, there are two approaches through the menus; Easy graphs and main menu entry.

The Easy graph dialogue boxes are a bit friendlier, but are not as powerful (missing options).

Easy Graph approach to a bar chart

🔲 graph bar - Bar cha	nrt			×
Main Over if/in Titles Y-Axis Statistic: Variable(s): sums count Percentile: 50 50	Options			
00		ок	Cancel	Submit

Main approach to a bar chart

🗖 g	raph	bar	- Bar	char	ts									×
Main	Over	By] if/in	Weights	Bars	Labels	Misc.	Y-Axis	Title	Caption	Legend Overall			1
Statis	tic:		E 60 - 2	New	name (d	optional)	- v	/ariable			New name (optional)	- V	/ariable	
sums		_					= c = [count		_		= = -		- 1
				ľ.			- - [- - [
none		•	50 -	Ξ										
none		•	50	3										
none		•	50	3										
00											ОК		Cancel	Submit

Both approaches could yield the command

graph bar (sum) count, over(major) ytitle(#
Applications) title(Berkeley Grade School
Applications in 1973)



One drawback in the Stata approach to graphics is that you can't add onto figures, as with S-Plus/R and Matlab.

However you can create figures that overlay plots on top of each other.

For example, one way to give a scatterplot with the regression line and 95% confidence intervals of the fit is

. regre	ss EngSize	Weig	ht							
Source	SS	df		MS		Number	r of	obs	=	93
+	+					F(1	,	91)	=	227.35
Model	70.7033831	1	70.70	33831		Prob :	> F		=	0.0000
Residual	28.2998401	91	.3109	87254		R-squa	ared		=	0.7142
	+					Adj R	-squ	ared	=	0.7110
Total	99.0032233	92	1.076	512199		Root I	MSE		=	.55766
EngSize	Coef.	Std.	Err.	t	P> t	[9]	5% C	onf.	Int	[erval]
	+									
Weight	.0014861	.0000	986	15.08	0.000	.00	0129	03	. (016819
_cons	-1.898924	.308	337	-6.16	0.000	-2.	5113	98	-1	.28645
. predi	ct ynat, xi	0								
. predi	ct residua	l, re	sidu	als						
mradi	at asfit	a + dm								

- . predict sefit, stdp
- . predict sepred, stdf
- . gen cilower = yhat invttail(91,0.025)*sefit
- . gen ciupper = yhat + invttail(91,0.025)*sefit

The predict command are available under General Post-Estimation

twoway - Twoway	/ graphs			×
Plot 1 Plot 2 Plot 3 Plot 4 By Required Type: scatter	Y-Axis X-Axis R-Axis Title Caption Lege	ind Overall		
Markers Symbol: Default Size: Default Color: Default Marker labels Variable:	if:	reate		
Size: Default Color: Default Position: Default	Additional graph options:	ОК	Cancel	Submit

🗖 twoway - Twoway	/ graphs				×
Plot 1 Plot 2 Plot 3 Plot 4 By Required Type: line	Y-Axis X-Axis R-Axis Title 1 X: Weight (I⊄ S	Caption Legend Overall	[]		
Plot on right axis	if: Line Color: Default Pattem: Default Width: Default Type: Default	Create			
2	Additional graph options:				
00			ОК	Cancel	Submit

Required			
Plot on right axis	if: Line Color: Red Pattem: Dash Width: Default Type: Default ▼	Create	
	Additional graph options:		

twoway (scatter EngSize Weight) (line yhat Weight, sort) (line cilower Weight, sort clcolor(red) clpat(dash)) (line ciupper Weight, sort clcolor(red) clpat(dash)), ytitle(Engine Size) title(Regression of Engine Size on Weight with 1993 Cars)



Actually there is an easier approach that doesn't involve running the regression ahead of time

. twoway (lfitci EngSize Weight, ciplot(rline) blcolor(red)) (scatter EngSize Weight, sort), ytitle(Engine Size) title(Regression of Engine Size on Weight with 1993

Cars)



Graphs in Stata can be saved in the following formats:

Stata graph (.gph), Windows metafile (.wmf), Enhanced metafile (.emf), Portable Network Graphics (.png), Postscript (.ps), and Encapsulated Postscript (.eps)

Multiway tables

Relationships between categorical variables

Example: 1973 graduate applications at Berkeley for the 6 largest programs

. table major status gender [fweight=count],row col

			Gender	and Status		
		Men			Women	
Major	Admitted	Rejected	Total	Admitted	Rejected	Total
	+					
А	511	314	825	89	19	108
В	353	207	560	17	8	25
С	120	205	325	202	391	593
D	138	279	417	131	244	375
E	54	137	191	94	299	393
F	22	351	373	24	317	341
Total	1,198	1,493	2,691	557	1,278	1,835

When dealing with data like this, it can be entered a couple of ways. For the above table, it was entered as

list

-	+			+
	major	gender	status	count
1.	 A	Men	Admitted	511
2.	A	Men	Rejected	314
3.	A	Women	Admitted	89
4.	A	Women	Rejected	19
5.	В	Men	Admitted	353
6.	B	Men	Rejected	207
7.	B	Women	Admitted	17
19.	E	Women	Admitted	94
20.	E	Women	Rejected	299
21.	F	Men	Admitted	22
22.	F	Men	Rejected	351
23.	F	Women	Admitted	24
24.	F	Women	Rejected	317
-				+

Instead, there could be a line for each person. For example you would need 551 lines with A Men Admitted, 314 lines with A Men Rejected, etc, for a total of 4526 entries.

The fweight (frequency weight) entry is saying that each line should be treated as existing count times.

Note that fweight can be used with many commands. Also there are other types of weights (sampling, analytic, importance)



tabulate gender status [fweight=count], chi2
expected row



	Sta	tus	
Gender	Admitted	Rejected	Total
Men	1,198	1,493	2,691
	1,043.5	1,647.5	2,691.0
	44.52	55.48	100.00
Women	557	1,278	1,835
	711.5	1,123.5	1,835.0
	30.35	69.65	100.00
Total	1,755	2,771	4,526
	1,755.0	2,771.0	4,526.0
	38.78	61.22	100.00

Pearson chi2(1) = 92.2053 Pr = 0.000

So it appears that women had a much harder time to get into Berkeley in 1973. However, lets look at each Major separately by major, sort: tabulate gender status
[fweight=count], chi2 expect exact

-> major = A	7		
+ Key	++		
freque freque expected f	ency requency		
	Stat	cus	
Gender	Admitted	Rejected	Total
Men	511 530.5	314 294.5	825 825.0
Women	89 69.5	19 38.5	108 108.0
Total	600 600.0	333 333.0	933 933.0

Pearson chi2(1) = 17.4307 Pr = 0.000 Fisher's exact = 0.000 -> major = B

	Stat	us	
Gender	Admitted	Rejected	Total
Men	353	207	560
	354.2	205.8	560.0
Women	17 15.8	9.2	25 25.0
Total	370	215	585
	370.0	215.0	585.0
Pea	arson chi2(1)	= 0.2537	Pr = 0.614
Fi	sher's exact	=	0.677

-> major = C

	Stat	us	
Gender	Admitted	Rejected	Total
Men	120	205	325
	114.0	211.0	325.0
Women	202	391	593
	208.0	385.0	593.0
Total	322	596	918
	322.0	596.0	918.0
Pea:	rson chi2(1)	= 0.7535	Pr = 0.385
Fi	sher's exact		0.387

->	major = I	C		
		Stat	us	
	Gender	Admitted	Rejected	Total
	Men	138 141.6	279 275.4	417 417.0
	Women	131 127.4	244 247.6	375 375.0
	Total Pea F:	269 269.0 arson chi2(1)	523 523.0 = 0.2980	792 792.0 Pr = 0.585 0.599

-> major = E

	Stat	us	
Gender	Admitted	Rejected 	Total
Men	54	137	191
	48.4	142.6	191.0
Women	94	299	393
	99.6	293.4	393.0
Total	148	436	584
	148.0	436.0	584.0
Pea	arson chi2(1)	= 1.2877	Pr = 0.256
F:	isher's exact	=	0.266

-> major = F

	Stat	us	
Gender	Admitted	Rejected	Total
Men	22	351	373
	24.0	349.0	373.0
Women	24	317	341
	22.0	319.0	341.0
Total	46	668	714
	46.0	668.0	714.0
Pea	arson chi2(1)	= 0.3841	Pr = 0.535
Fi	sher's exact	=	0.546

Summary of tests

Major	Chi2	P-value
А	17.4307	0.000
В	0.2537	0.614
С	0.7535	0.385
D	0.2980	0.585
E	1.2877	0.256
F	0.3841	0.535

So in every major, expect A, the data is consistent with the acceptance probabilities being the same for men and women. In major A, women appear to be accepted at a higher rate than men 82.4% vs 61.9

The data set gives an example of Simpson's paradox, where aggregating over third factor can switch the direction of association between to factors.

For a further description of Simpson's Paradox, see Moore and McCabe or a short article I wrote for the OSU Biostat Center's newsletter a few years ago (link on the overheads page)

Pearson Chi-square test

Can be used to check for independence between two categorical variables (or homogeneity of proportions between one factor, for each level of a second factor)

Compares observed counts within each cell with the expected counts assuming independence (or homogeneity)

These are gotten by

expected count = $\frac{\text{row total} \times \text{column total}}{n}$

The common test statistic for examining whether there is association between two variables in a $r \times c$ table is

$$X^{2} = \sum_{\text{all cells}} \frac{(\text{observed count - expected count})^{2}}{\text{expected count}}$$

If the null hypothesis isn't true, there should be some cells in the table where the observed and expected counts are very different. This would lead to a large value for X^2 .

If the number of observations is large, the sampling distribution of X^2 is approximately Chi-square (χ^2) with (r-1)(c-1) degrees of freedom.

What to do if the sample sizes are small.

Similar to the one sample binomial problems, there is an exact procedure for two way tables that should be used when the normal approximation doesn't hold. In the 2×2 table case, its known as Fisher's exact test. It can be easily performed in Stata. Idea of Fisher's exact test

Assume that the margins for each table are fixed

```
-> major = B
```

	Stat	tus	
Gender	Admitted	Rejected	Total
Men	353	207	560
Women	17	8	25
Total	370	215	585
P	earson chi2(1) = 0.253	7 $Pr = 0.614$

Fisher's exact =

In a 2 by 2 table, with fixed margins, once you set one cell in the table, the others are determined exactly. For example lets set the value for the Rejected Women cell at 6. Then the other cells must be 351, 209, and 19.

Under the null hypothesis, the probability of any possible table is given by the hypergeometric distribution

0.677

a	b	U
С	d	V
S	Т	Ν

The probability of this table is given by



Fisher's exact tests sums these probabilities over all tables considered as or more extreme than the observed table.