

Post-Estimation Techniques in Statistical Analysis:

Introduction to Clarify and S-Post in Stata

PRISM Brownbag

November 16, 2004

By: Kevin Sweeney and Brandon Bartels

Presenters: Dave Darmofal and Corwin Smidt

(Note: If you're not in Political Science come talk to me to log in.)

Preliminaries

- We will be posting these Powerpoint presentations to the PRISM “Luncheons” webpage:
<http://psweb.sbs.ohio-state.edu/prism/luncheons.htm>
- Also, you will be logging—via a *.log* file—all of the *S-Post* and *Clarify* procedures you are about to run.
- *Bottom line:* Everything said and done here will be on the record, so there's less of a need to take extensive notes.
- Commands we'll be using in *S-Post*:
 - Open *Notepad*. Start → Programs → Accessories → Notepad
 - I: → general → Spost&Clarify → Post-estimation commands.txt

Introduction

- Three *S*'s of statistical analysis:
 - *Sign*
 - *Significance*
 - ***Strength, or substantive importance***
 - The *effect* of an independent variable on the dependent variable is “a change in an outcome for a change in an independent variable, holding all other variables constant” (Long 1997, 6).
- Most quantitative articles in leading journals contain post-estimation calculations of substantive effects of the independent variables of interest.

Effects in Linear and Non-Linear Models

- Examining marginal effects in OLS is easy: ***b***. A one-unit change in *X* produces a ***b***-unit change in *Y*, holding other variables constant.
- Examining effects in nonlinear models, such as logit, probit, ordered probit, and other ML models, is less straightforward. Marginal effects with respect to *X* are not constant (note: but not interactive).
- In nonlinear models, *the magnitude of the change in the probability of an event occurring, given a change in a particular independent variable, depends on the levels of the other independent variables.*
- *S-Post* (Long) and *Clarify* (Tomz, Wittenberg, and King) make post-estimation easy and offer a powerful means of presenting the substantive results from a statistical analysis.

What S-Post Can Do

- For a comprehensive presentation of S-Post, see:
 - Long, J. Scott, and Jeremy Freese. 2001. *Regression Models for Categorical Dependent Variables Using Stata*. College Station, TX: Stata Press.
 - Once S-Post is installed, the “help” files provide very good information on the commands. E.g, “help prchange”.
- For those who don’t usually use Stata, J. Scott Long and Simon Cheng also have Excel spreadsheets available to download which are easy to use and present nice graphs. Download at: <http://www.indiana.edu/~jslsoc/xpost.htm>

A Quick Look: S-Post Excel Spreadsheets

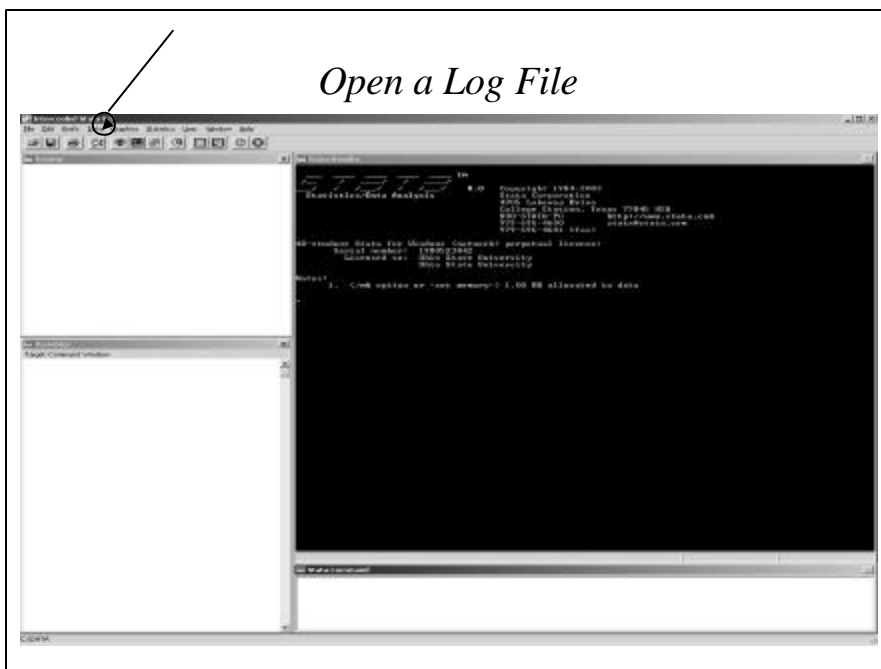
- The spreadsheets are easy to manage, come with clear instructions, and produce quick and easy graphs.
- However, they are read-only files and do not have much user flexibility.
- Still they serve as a nice way to save old results for easy reference.



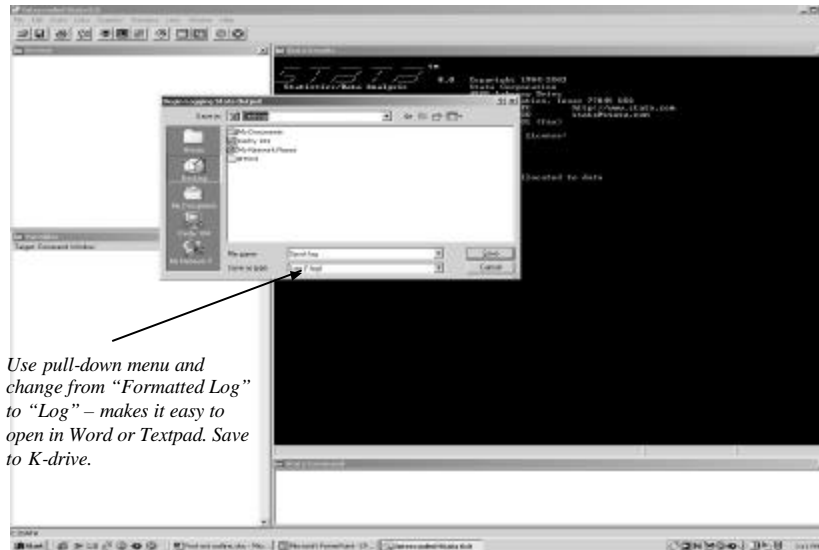
What S-Post Can Do – Key Commands

- `fitstat` – goodness-of-fit measures
- `listcoef` – odds ratios
- `prvalue` & `prtab` – predicted probabilities for particular covariate profiles.
- `prchange` – first differences
- `prgen` – setup for graphing
- Again, use the “help” files in Stata if you get stuck on these.
- These commands can be used in:
 - `regress`, `logit`, `probit`, `ologit`, `oprobit`,
`mlogit`, `clogit`, `cloglog`, `poisson`, `nbreg`,
`cnreg`, `intreg`, `tobit`, `zip`, `zinb`

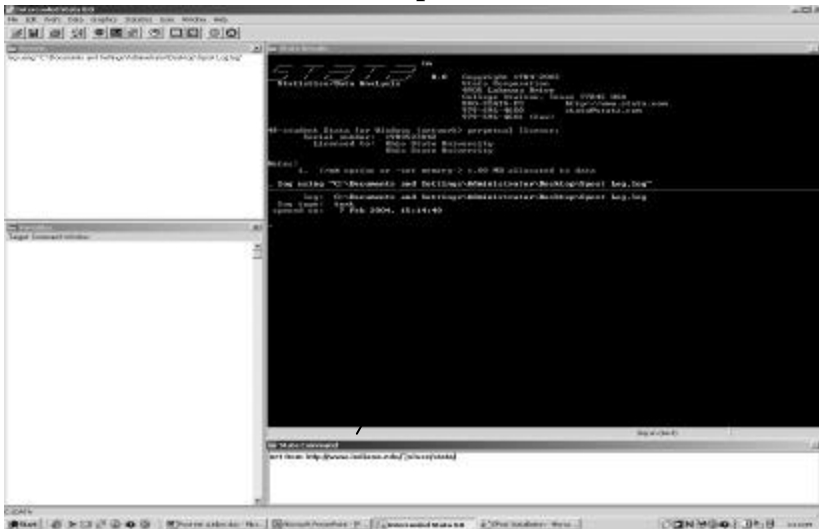
Open a Log File



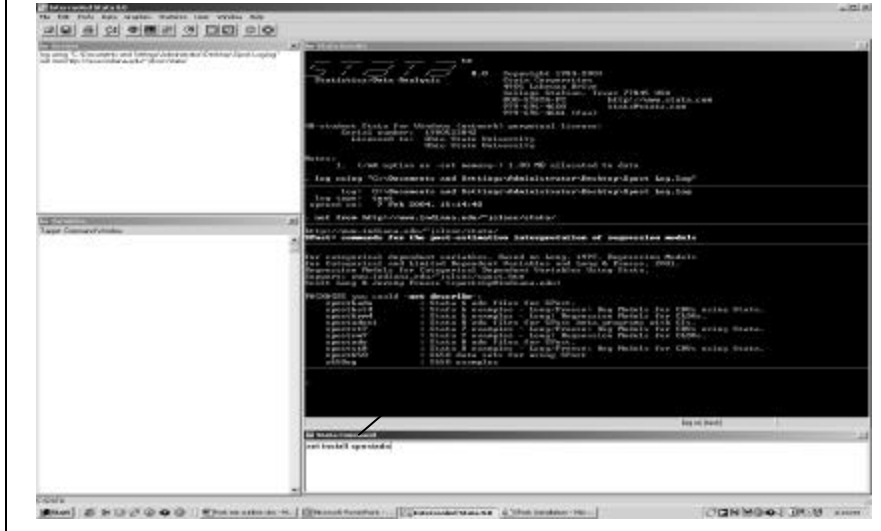
Open a Log File



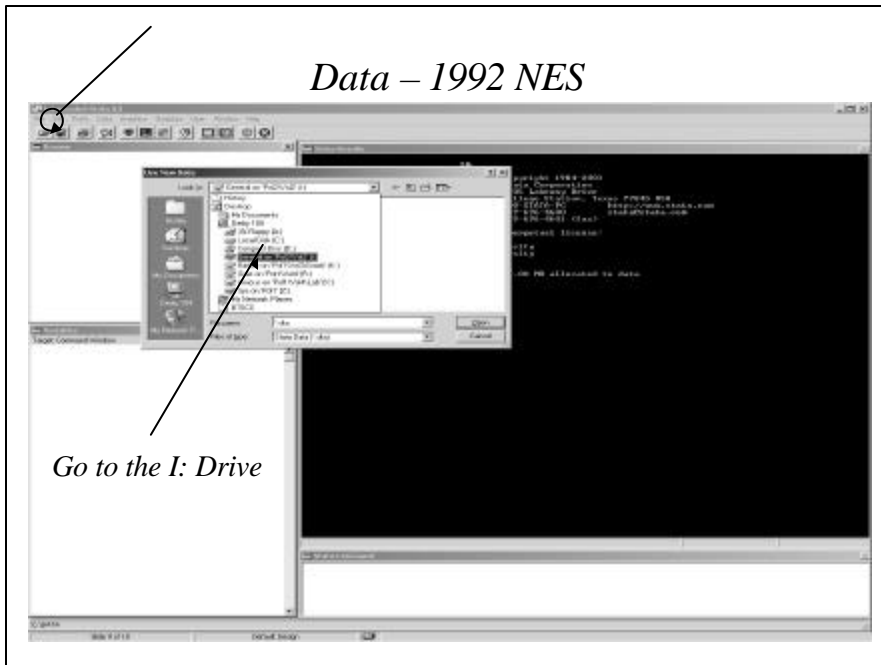
Installing S-Post Step 1



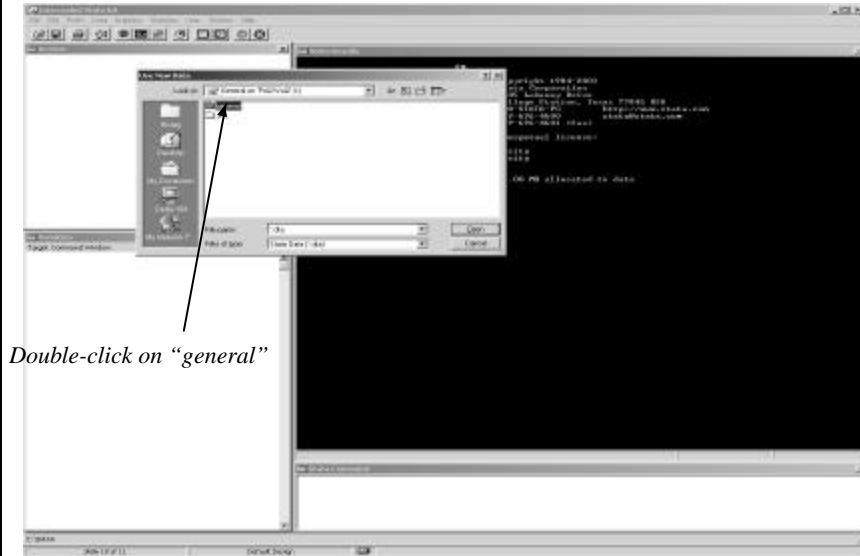
Installing S-Post Step 2



Data – 1992 NES

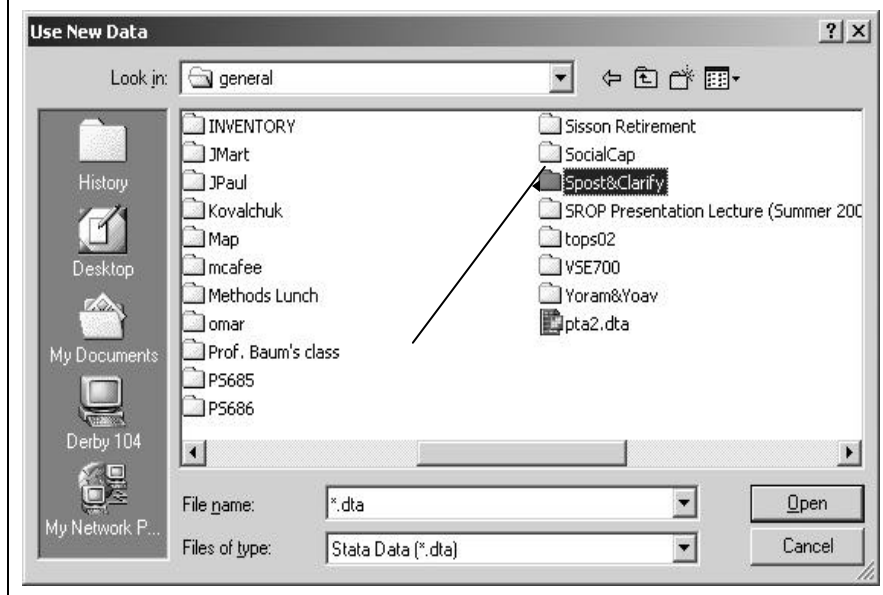


Data – 1992 NES

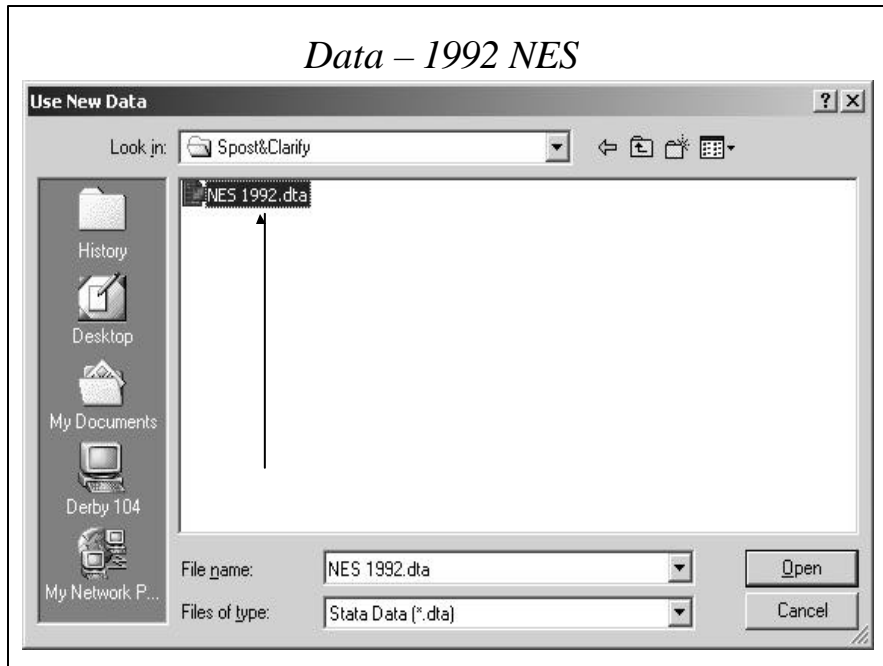


Double-click on "general"

Data – 1992 NES



Data – 1992 NES



Data – 1992 NES

Variable	Description
<i>vote3</i>	<i>Vote 92: 0=Bush, 1=Clinton, 2=Perot</i>
<i>bush_app</i>	<i>Bush Approval, 1992 (4-pt)</i>
<i>ideology</i>	<i>Respondent's ideology (1=lib, 7=cons)</i>
<i>econworse</i>	<i>Economy WORSE than yr ago (5-pt)</i>
<i>militaryopp</i>	<i>Opposition to Use of Military Force (5-pt)</i>
<i>gulfwarworth</i>	<i>Gulf War Worth Cost (dichotomous)</i>
<i>pid</i>	<i>Party ID (SD=-3, SR=+3)</i>
<i>education</i>	<i>Years of School</i>
<i>govtemp</i>	<i>Government Employee (dichotomous)</i>
<i>union</i>	<i>Union Household (dichotomous)</i>
<i>income</i>	<i>Family Income, \$1,000</i>
<i>nonwhite</i>	<i>Nonwhite (dichotomous)</i>
<i>vote2</i>	<i>Vote 92: 0=non-Bush; 1=Bush</i>

Goodness-of-Fit Measures

fitstat

- Logit and Probit report one pseudo- R^2 measure: McFadden's R^2 : $(\text{init LL} - \text{final LL})/(\text{init LL})$.
- There are other pseudo- R^2 measures, too; see Long (1997, 104-113).
- Two statistics that are often reported in journal articles: *percent correctly predicted* (PCP; using a 0.5 threshold) and *proportional reduction in error* (PRE) – although see Train 2003, p.73 for why their theoretical basis is questionable.
- PRE is a measure comparing the predictive success of the estimated model to a null model, i.e., proportion of the DV in the modal category (PMC).
- $PRE = (PCP - PMC)/(1 - PMC)$
- The `fitstat` command can give these to you in an instant! Also available for models other than logit and probit.
- Let's estimate a simple vote choice model to check it out.
- `logit vote2 pid econworse militaryopp education nonwhite`

Goodness-of-Fit Measures

```

Notes:
1. (/## option or -set memory-) 1.00 MB allocated to data
use "I:\general\NES 1992.dta", clear
logit vote2 pid econworse militaryopp education nonwhite

Iteration 0: log likelihood = -254.0798
Iteration 1: log likelihood = -184.80258
Iteration 2: log likelihood = -157.16871
Iteration 3: log likelihood = -156.63857
Iteration 4: log likelihood = -156.63427
Iteration 5: log likelihood = -156.63427

logit estimates
Number of obs = 382
LR chi2(5) = 194.89
Prob > chi2 = 0.0000
Pseudo R2 = 0.3835

+-----+-----+-----+-----+-----+
| vote2 |      Coef. |   Std. Err. |      z |   P>|z| | [95% Conf. Interval] |
+-----+-----+-----+-----+-----+
| pid   |   -8191285 |   -8897104 |    9.13 |  0.000 |   -6432993 |   -9849577 |
| econworse | -2868483 |   -1638817 |   -1.75 |  0.080 |   -6080106 |   -949341 |
| militaryopp | -5852654 |    186129 |   -3.14 |  0.002 |   -9500714 |   -2204593 |
| education |  990359 |    705154 |    1.28 |  0.202 |   -0481717 |   2282436 |
| nonwhite |  -6128034 |    5694173 |   -0.02 |  0.981 |  -1.129841 |   1.102234 |
| _cons |  -6603643 |   1.314074 |   -0.49 |  0.626 |   -1.903573 |   2.114302 |
  
```

Goodness-of-Fit Measures

Iteration 3: log likelihood = -156.63427

logit estimates

vote2	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
pid	.8191285	.0897104	9.13	0.000	.6432993 .9949577
econworse	-.2868483	.1638817	-1.75	0.080	-.6080506 -.0343561
militaryopp	-.5852654	.186129	-3.14	0.002	-.9500714 -.2204593
education	-.0900359	.0705154	-1.28	0.202	-.4481717 .2282436
nonwhite	-.0138034	.5694173	-0.02	0.981	-1.129841 1.102234
_cons	-.6603643	1.354074	0.49	0.626	-1.993573 3.314301

Number of obs = 382
LR chi2(5) = 194.89
Prob > chi2 = 0.0000
Pseudo R2 = 0.3835

fitstat

Measures of Fit for logit of vote2

Log-Lik Intercept only:	-254.080	Log-Lik Full Model:	-156.634
D(376):	313.269	LR(5):	194.891
McFadden's R2:	0.384	Prob * LR:	0.000
Maximum Likelihood R2:	0.400	McFadden's Adj R2:	0.360
McKelvey and Zavoina's R2:	0.550	Cragg & Uhler's R2:	0.443
Variance of y*:	7.316	Efron's R2:	0.451
Count R2:	0.814	Variance of error:	3.290
AIC:	0.851	Adj Count R2:	0.514
BIC:	-1922.210	AIC(n):	325.269
		BIC:	-165.164

“Count R2” is PCP

“Adj Count R2” is PRE

Goodness-of-Fit Measures

logit vote2 p10 econworse tab vote2

nonwhite	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
nonwhite	-.0138034	.5694173	-0.02	0.981	-1.129841 1.102234
_cons	.6603643	1.354074	0.49	0.626	-1.993573 3.314301

fitstat

Measures of Fit for logit of vote2

Log-Lik Intercept only:	-254.080	Log-Lik Full Model:	-156.634
D(376):	313.269	LR(5):	194.891
McFadden's R2:	0.384	Prob * LR:	0.000
Maximum Likelihood R2:	0.400	McFadden's Adj R2:	0.360
McKelvey and Zavoina's R2:	0.550	Cragg & Uhler's R2:	0.443
Variance of y*:	7.316	Efron's R2:	0.451
Count R2:	0.814	Variance of error:	3.290
AIC:	0.851	Adj Count R2:	0.514
BIC:	-1922.210	AIC(n):	325.269
		BIC:	-165.164

tab vote2

Vote 02:	Freq.	Percent	Cum.
0=non-bush			
1=bush			
0	236	61.78	61.78
1	148	38.22	100.00
Total	382	100.00	

d1 (.814-.6178)/(1-.6178)
.5136054

PRE

Odds Ratios listcoef

- In logit, we can transform the LHS into the *log of the odds*, and the *log odds* are a linear combination of $X\mathbf{b}$. Exponentiating the betas—i.e., creating odds ratios—produces this transformation. $\exp(\mathbf{b})$ is an odds ratio.
- Odds ratios can then be interpreted as follows: *For a one-unit change in X_1 , the odds of voting for Bush change by a factor of $\exp(\mathbf{b}_1)$, holding all other variables constant.* (see Long 1997, 79-82)
- If $\exp(\mathbf{b}_1)$ is greater than 1, then the odds of voting for Bush increase as X_1 increases. If $\exp(\mathbf{b}_1)$ is less than 1, the odds of voting for Bush decrease as X_1 increases.
- To get odds ratios in S-Post, use:
 - *listcoef, help*
- [Note, you can also get odds ratios by using `logistic` instead of `logit`.]

Odds Ratios

logit (n=382): Factor change in odds

Odds of: 1 vs 0

vote2	b	z	P> z	e^b	e^b*SD	SD
pid	0.81913	9.131	0.000	2.2685	5.5529	2.0929
econworse	-0.28685	-1.750	0.080	0.7506	0.7714	0.9047
militaryppp	-0.38327	-3.244	0.002	0.5570	0.8279	0.7952
education	0.09004	1.277	0.202	1.0942	1.2140	2.1533
nonwhite	-0.01380	-0.024	0.981	0.9862	0.9861	0.3451

b = raw coefficient
 z = z-score for test of b=0
 P>|z| = p-value for z-test
 e^b = exp(b) = factor change in odds for unit increase in X
 e^b*SD = exp(b*SD of X) = change in odds for SD increase in X
 SD = standard deviation of X

Holding the other variables constant, a one-unit increase in PID increases the odds of voting for Bush by a factor of 2.27.

Odds Ratios – Percent Change

logit (N=882): Percentage Change in odds

	b	z	P> z	%	%stdx	sdofx
pid	0.81913	9.131	0.000	325.2	451.3	2.0928
econworse	-0.28685	-1.750	0.080	-24.0	-32.9	0.9047
militaryopp	-0.58527	-3.144	0.002	-44.3	-37.2	0.7952
education	0.09004	1.277	0.202	9.4	21.4	2.1535
nonwhite	-0.01380	-0.024	0.981	-1.4	-0.4	0.2851

b = raw coefficient
 z = z-score for test of b=0
 P>|z| = p-value for z-test
 % = percent change in odds for unit increase in X
 %stdx = percent change in odds for SD increase in X
 sdofx = standard deviation of X

Holding other variables constant, a one unit increase in negative economic perceptions decreases the odds of voting for Bush by about 25%. [25% is simply calculated as 1 minus the odds ratio, i.e., 1 - .751]

Predicted Probabilities

- In ML models, you can convert equations into probability statements.
- Recall, in Logit and Probit: $\Pr(Y = 1 | X) = F(Xb)$
 - F is the cumulative distribution function (CDF).
 - For *logit*, we use the logistic CDF, and get:

$$\Pr(Y = 1 | X) = \frac{\exp(Xb)}{1 + \exp(Xb)}$$

- For *probit*, we use the standard normal CDF:

$$\Pr(Y = 1 | X) = \Phi(Xb)$$

- *Ordered probit*:

$$\Pr(Y = m | X) = \Phi(t_m - Xb) - \Phi(t_{m-1} - Xb)$$

Predicted Probabilities

- In our logit model of vote choice:

$$\Pr(Y = 1|X) = \frac{\exp(\mathbf{b}_0 + \mathbf{b}_1PID + \mathbf{b}_2Econ + \mathbf{b}_3Milit + \mathbf{b}_4Educ + \mathbf{b}_5Nonwhite)}{1 + \exp(\mathbf{b}_0 + \mathbf{b}_1PID + \mathbf{b}_2Econ + \mathbf{b}_3Milit + \mathbf{b}_4Educ + \mathbf{b}_5Nonwhite)}$$

- Present results in terms of probability to draw conclusions about the substantive importance of variables.
- A number of ways to do this:
 - *Predicted probabilities* for various *covariate profiles*.
 - *First differences*. Change in the probability of an event occurring given a particular change in an IV, holding other variables constant at baseline values.
 - *Graphing* the probability of an event occurring as a function of an IV of interest, holding other variables constant at baseline values.

Predicted Probabilities for Covariate Profiles

`prvalue`

- Let's say I wanted to know the probability that a highly educated individual who is strongly opposed to military force voted for Bush.
- Use `prvalue` to set these two particular variables to the desired values, and set the other variables to baseline values (I will set everything to mean levels).
- Basic syntax:
`prvalue, x(education max militaryopp max) rest(mean)`

Predicted Probabilities for Covariate Profiles

Stata Command Window Output:

```

logit vote2 pid econworse militaryopp education nonwhite
-----+-----
             b         z         P>|z|         %         %stdx         SbofX
-----+-----
             +-----+-----+-----+-----+-----+-----+
pid          0.81913    9.111    0.000    126.9    451.1    2.8929
econworse   -0.28655   -1.750    0.080    -24.9    -22.9    0.9047
militaryopp -0.58527   -1.144    0.252    -44.3    -37.2    0.7952
education    0.09004    1.277    0.202     9.4     21.4    2.1535
nonwhite    -0.01360   -0.028    0.981     -1.4     -0.4    0.2891
-----+-----

             b = raw coefficient
             z = z-score for test of b=0
             P>|z| = p-value for z-test
             % = percent change in odds for unit increase in x
             %stdx = percent change in odds for SD increase in x
             SbofX = standard deviation of x

             prvalue, x(education max militaryopp max) rest(mean)
logit: Predictions for vote2
             Pr(Y=1|X):    0.1399    95% CI: (0.0611,0.2754)
             Pr(Y=0|X):    0.8601    95% CI: (0.7246,0.9345)

             +-----+-----+-----+-----+-----+-----+
             pid          econworse    militaryopp    education    nonwhite
             +-----+-----+-----+-----+-----+-----+
             .07591623    3.9790576             5             17             .08900124
    
```

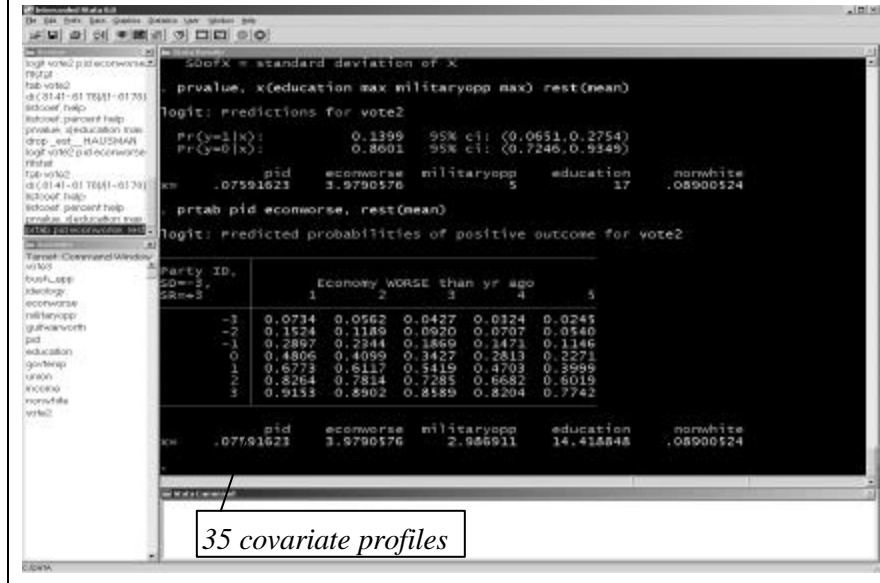
Probability of a highly educated individual who is strongly opposed to military force (holding the other variables constant at their means) voting for Bush is about 0.14.

Cross-Tabs of Predicted Probabilities

prtab

- Another way to present probabilities for particular covariate profiles is by creating a cross-tab of probabilities.
- S-Post's `prtab` command computes a table of predicted probabilities for all combinations of as many as 4 categorical variables.
- Let's say I wanted to examine the probability of voting Bush given all possible covariate profiles of PID and negative economic perceptions, holding other variables constant at their mean values.
- Syntax:
`prtab pid econworse, rest(mean)`

Cross-Tabs of Predicted Probabilities



First Differences prchange

- *First difference*: The change in the probability of voting for Bush given a particular change in an independent variable, holding all other variables constant at some baseline value (e.g., means or modes for dichotomous variables). I recommend using *Clarify* to do these, but I'll quickly run through the syntax.
- *prchange* calculates these for you.
- First, do a `set mat` command to expand the matrix in Stata.

set mat 60

prchange, help

First Differences prchange

```

set mat 60
prchange, help

logit: changes in predicted probabilities for vote2

      min--max      0--x1      -x1/2      -+sd/2      MargEfect
      pid           0.7887      0.1893      0.1893      0.3458      0.1704
      econworse    -0.2580      -0.0713      -0.0598      -0.0540      -0.0597
      militaryopp   -0.4585      -0.1337      -0.1214      -0.0966      -0.1218
      education     0.1817      0.0086      0.0187      0.0403      0.0187
      nonwhite      -0.0029      -0.0029      -0.0029      -0.0008      -0.0029

Pr(y|x) = 0.7047 0.2953

      x=      pid      econworse      militaryopp      education      nonwhite
sd(x)= 2.09286      .904687      .79522      2.15354      .089005

Pr(y|x): probability of observing each y for specified x values
Avg|Chg|: average of absolute value of the change across categories
Min--Max: change in predicted probability as x changes from its minimum to
its maximum
0--x1: change in predicted probability as x changes from 0 to 1
-x1/2: change in predicted probability as x changes from 1/2 unit below
base value to 1/2 unit above
-+sd/2: change in predicted probability as x changes from 1/2 standard
dev below base to 1/2 standard dev above
MargEfect: the partial derivative of the predicted probability/rate with
respect to a given independent variable

```

First Differences prchange, fromto

```

prchange, fromto help

logit: changes in predicted probabilities for vote2

      from:      to:      dif:      from:      to:      dif:
      econworse  0.0326      0.8213      0.7887      0.2825      0.4715      0.1893
      militaryopp 0.4961      0.2182      -0.2780      0.5674      0.4961      -0.0713
      education   0.5727      0.1142      -0.4585      0.7064      0.3727      -0.1337
      nonwhite    0.1661      0.3458      0.1817      0.1626      0.1312      0.0086
      nonwhite    0.2955      0.2926      -0.0029      0.2955      0.2926      -0.0029

      from:      to:      dif:      from:      to:      dif:
      pid        x=-1/2      x=1/2      -x1/2      x=1/2sd      +x1/2sd      -+sd/2
      econworse  0.2176      0.3809      0.1633      0.1510      0.4965      0.3458
      militaryopp 0.3505      0.3823      0.0318      0.1250      0.2090      -0.0840
      education   0.2860      0.3047      0.0187      0.1459      0.2402      -0.0966
      nonwhite    0.2967      0.2938      -0.0029      0.2957      0.2948      -0.0008

      MargEfect
      pid           0.1704
      econworse    -0.0597
      militaryopp   -0.1218
      education     0.0187
      nonwhite      -0.0029

Pr(y|x) = 0.7047 0.2953

      x=      pid      econworse      militaryopp      education      nonwhite
sd(x)= 2.09286      .904687      .79522      2.15354      .089005

Pr(y|x): probability of observing each y for specified x values
Avg|Chg|: average of absolute value of the change across categories
Min--Max: change in predicted probability as x changes from its minimum to
its maximum
0--x1: change in predicted probability as x changes from 0 to 1
-x1/2: change in predicted probability as x changes from 1/2 unit below
base value to 1/2 unit above
-+sd/2: change in predicted probability as x changes from 1/2 standard
dev below base to 1/2 standard dev above
MargEfect: the partial derivative of the predicted probability/rate with
respect to a given independent variable

```

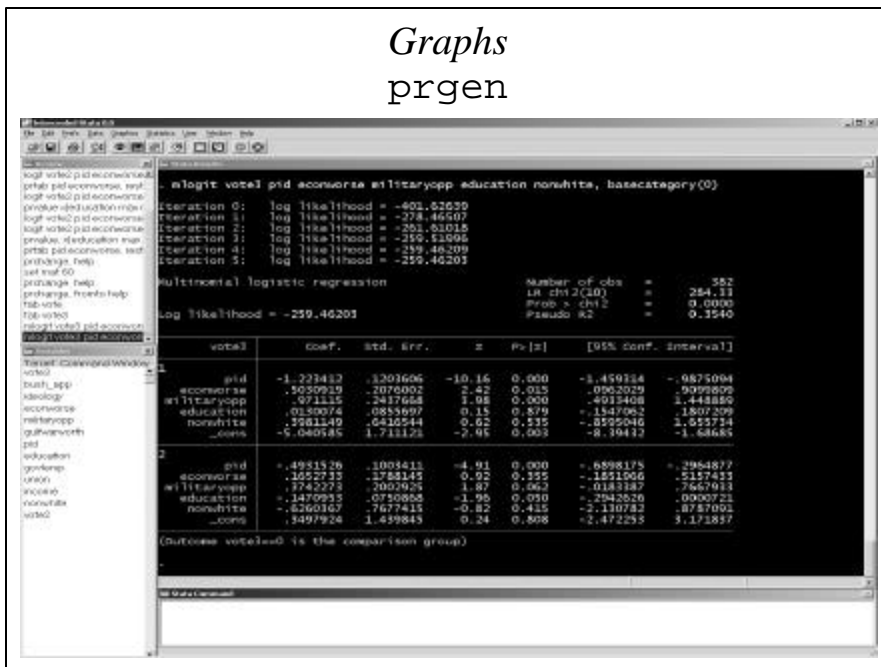

Graphs prgen

- Graphing is a very powerful way to present the results of a statistical model.
- Let's say we wanted to graph the probability of an event occurring as a function of an independent variable of interest.
- To show you that I can use S-Post outside of logit, let's run a multinomial logit model of vote choice, with Bush, Clinton, and Perot as the three nominal outcomes of the DV.

```
mlogit vote3 pid econworse militaryopp education nonwhite, basecategory(0)
```

- [Note: One can test whether this model violates the I.I.A. assumption using S-Post's `mlogtest` command. Do "help mlogtest" for more info.]

Graphs prgen



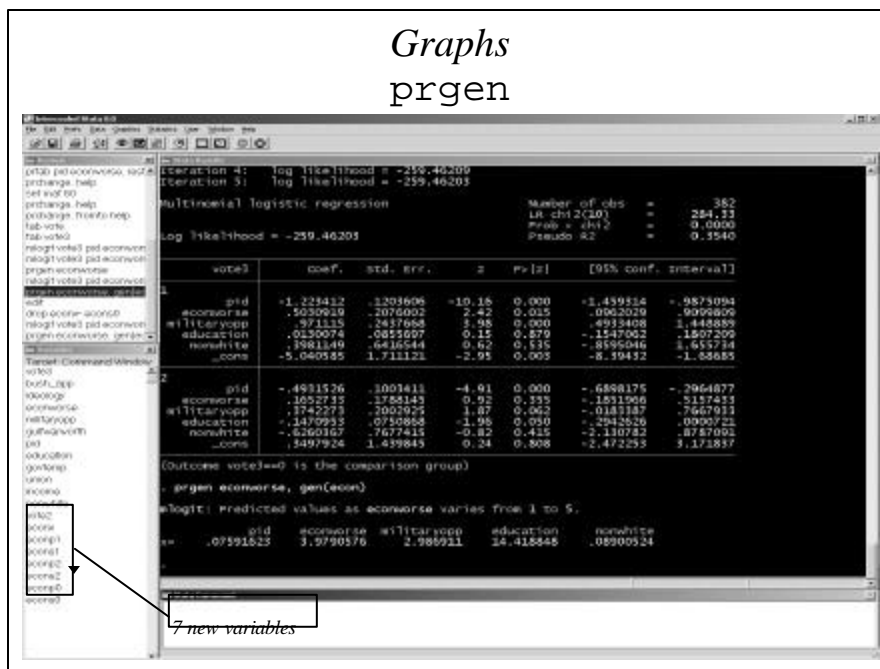
Graphs prgen

- Effect of economic perceptions on vote choice? Graph the probabilities of voting for Bush, Clinton, and Perot as a function of economic perceptions, holding other variables constant at a baseline value.

prgen econworse, gen(econ)

- Note:* the default settings generate predicted probabilities of voting for each of the three candidates as “econworse” ranges from its minimum (1) to maximum (5) value (holding other variables at their mean levels).
- See “help prgen” for more options.

Graphs prgen

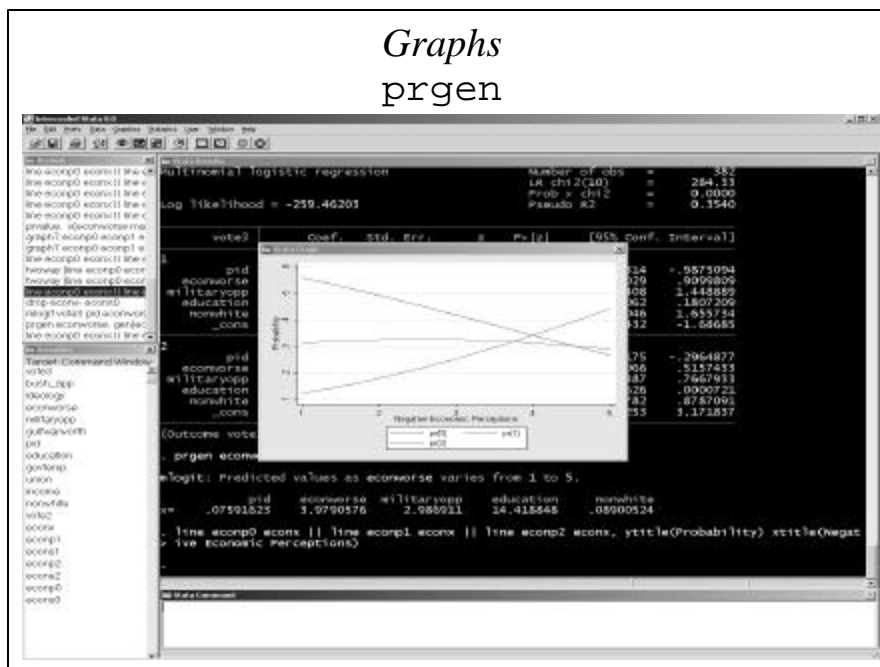


Graphs prgen

- To produce a simple graph of the effect of economic perceptions, use the `line` function:


```
line econp0 econp1 econp2 econx, ytitle(Probability) xtitle(Negative Economic Perceptions)
```
- One always lists the y-axis variables upfront with the x-axis variable (`econx`) coming last.

Graphs prgen



Graphing Hints in Stata

- The problem with the `line` command is that it is only useful if you have a color printer since it doesn't allow line symbols for differentiation.
- While some individuals resort to using Excel for graphs, it would be well worth your time to learn the graphing capabilities of Stata, especially `scatter` in this case.
- `scatter` is Stata's scatterplot graph command which allows for symbols to mark the lines. One only needs to make a few simple adjustments to it for one to produce a clean and highly informative graph. Since Stata's graphing help file is huge, I thought I would highlight a few, key commands to implement.

Scatter Commands

- First, graph adjustment commands are all made after the comma:
`scatter var1 var2 varx, adjustment commands....`
- One should first use `msymbol()` to assign symbols to each y-variable and then connect the symbols with lines using `connect()`:
`scatter var1 var2 varx, msymbol(d x) connect(l l)`
where `msymbol` options include:
d – diamond, x – x-mark, s – square, + – plus, o – circle, . – dot, Dh – large hollow diamond,
Sh – large hollow square, Oh – large hollow circle, and many more...
and `l` within `connect` connects each variable by a line.
- Note: the symbol and connect commands apply in order to each y-variable, where one uses a space to distinguish between variables.
- To title each axis use the `xtitle(...)` and `ytitle(...)` commands, with the titles typed within the parentheses. Use the `label var` command to label the variables so the legend is more descriptive.

Final Product

- I've posted a text document which walks you through a number of other easy and useful adjustments when graphing.

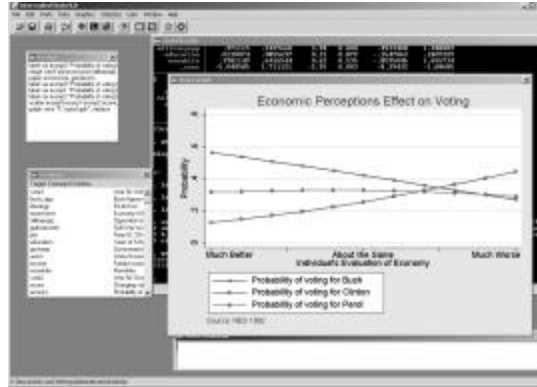
I:\General\Spst&Clarify\graphing.txt

- The final product can be seen by first labeling the variables by typing:

```
label var econp0 "Probability of voting for Bush"
```

```
label var econp1 "Probability of voting for Clinton"
```

```
label var econp2 "Probability of voting for Perot"
```



- Then entering the following:

```
scatter econp0 econp1 econp2 econx, msymbol(* Oh d) connect(1 1) legend(cols(1) post(7)) ytitle(Probability) xtitle(Individual's Evaluation of Economy)
yscale(r(0 .8)) xlabel(#4) xlabel(1.25 "Much Better" 3 "About the Same" 4.75 "Much Worse", noticks) stick(#5) title(Economic Perceptions Effect on Voting)
note(Source: NES 1992)
```

- Which gives you a paper-worthy graph...

Saving Graphs in Stata

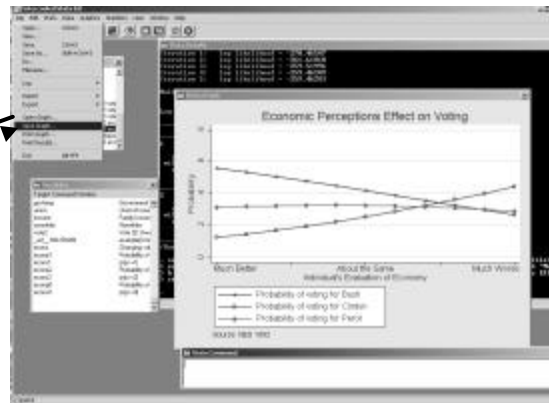
- After producing the graph make sure to save it as a windows metafile for use in other windows programs.

- Click on:

File → Save Graph...

- Then select "Windows Metafile (*.wmf)" under the "Save as type:" box.

- This should enable the figure to be used in other windows programs.



Conclusion

- *Get this book!*
 - Long, J. Scott, and Jeremy Freese. 2001. *Regression Models for Categorical Dependent Variables Using Stata*. College Station, TX: Stata Press.
- Use “help” files. They provide very good information on the commands, e.g, “help prchange”.