

Econometric Analysis – Dr. Sobel

Econometrics Session 1:

1. Building a data set

Which software - usually best to use Microsoft Excel (XLS format) but CSV is also okay

Variable names (first row only, 15 character max in gretl software we will be using, 1st must be letter, no spaces but _ symbol is ok)

First column is sometimes dates or labels (state names) and it is okay if this has spaces, but give it a title too

EXAMPLE DATA SET IN EXCEL:

The screenshot shows an Excel spreadsheet with the following data:

Country	Entrepreneurship	GDP_PerCapita	GovSpend_PerCap	Percent Male	Median_Age
Australia	6.68	24000	4344.548191	49.94	36
Belgium	2.99	26100	10332.63878	49.06	40
Canada	8.82	27700	5108.758149	49.54	37.8
Denmark	6.53	28000	9583.742386	49.45	39.1
Finland	4.56	25800	5989.431937	48.76	40.3
France	3.2	25400	4030.143661	48.58	38.3
Germany	5.16	26200	9936.223177	49.10	41.3
Hungary	6.64	12000	1424.893704	47.70	38.4
Iceland	11.32	30200	11874.51872	50.71	34
Ireland	9.14	27300	7029.715911	49.48	33.1
Italy	5.9	24300	8963.272687	48.53	41
Japan	1.81	27200	5663.72633	48.89	42
Mexico	12.4	9000	1374.176867	48.63	23.8
Netherlands	4.62	25800	8384.709494	49.47	38.6
New Zealand	14.01	10500	1295.922033	48.89	33.1
Norway	8.69	30800	12790.22259	49.55	37.7
Poland	4.44	8800	1353.732959	48.58	36
South Korea	14.52	19600	1997.730061	50.35	33.2
Spain	4.59	18900	2722.414047	48.90	38.7
Sweden	4	24700	12394.29218	49.63	40.1
Switzerland	7.13	31100	4119.02669	49.63	40.2
United Kingdom	5.37	24700	9053.143461	49.14	38.4
United States	10.51	36300	6124.602077	48.91	35.8

Red annotations in the image highlight the variable names in the first row and a note on the right side of the spreadsheet:

Variable names must be 15 letters or less, begin with a letter, and cannot have any spaces (but underscores " _ " are okay). They must be on the first row, with data starting on the second row. A column of names that label the data is okay and it doesn't matter much what's in it, and it is okay to have spaces in it.

The number of “observations” (data points) matters. More is better. Always report the number of observations you use in a research paper.

Missing observations – usually best to have none, but most programs can deal with them – leave cell blank

Types of datasets: cross section, time series, panel (usually must specify, and label)

– in gretl software we will be using, go to Data menu, Dataset structure to change or specify

Source information – how to document it

Descriptions: what is included, what year, how measured, etc.

Keep track of units of measurement (data in millions of dollars? Percentages: would 5% be as 0.05 or 5.0?)

EXAMPLE OF HOW TO DOCUMENT YOUR DATA IN YOUR RESEARCH PAPER:

Appendix 1: List of Variables and Sources

Table 5

Variable name	Minimum	Maximum	Average	Obs.	Source	Definition
Dependent variables						
Average business start-up rate	0.114	0.229	0.145	50	a	Average business births as a percentage of total businesses with 1–9 employees (2002–2003)
Net business creation rate	0.001	0.052	0.018	50	a	Average net business creation as a percentage of total businesses with 1–9 employees (2002–2003)
Venture capital per capita	0	379.39	46.635	48	b	Venture capital invested (by destination) per capita (2005)
Patents per capita	0.044	1.086	0.249	50	c	Number of patents granted (all types) per 1,000 people (2005)
Productive entrepreneurship	3.2	43	23.583	48	d	Index of entrepreneurial activity from Sobel (2008) based on several measures of entrepreneurship
Independent variables						
Cultural diversity	0.021	0.44	0.132	50	e	Probability that 2 randomly chosen individuals from a state were born in different countries, calculated like a standard Herfindahl–Hirshman Index (2000)
Median age	27.1	37.9	35.536	50	f	Median age of state population (2000)
Percent male	88.8	107.6	94.114	50	f	Percent of 18 years and over state population that is male (2000)
Percent college degree	14.8	33.2	23.776	50	f	Percent of state population with a bachelor's degree or higher (2000)
Population density	0.0011	1.1344	0.1819	50	f	Population density per square mile of land area in state (2000), in thousands
Median household income	29.696	55.146	41.371	50	e	Median household income (1999), in thousands
Economic freedom	5.6	7.8	6.684	50	f	Economic freedom of North America (EFNA) index (2005)

Data Sources: *a* Company Statistics Division of the United States Census Bureau (<http://www.census.gov/csd/susb/susb.htm>), *b* Venture Capital Association's (NVCA) Yearbook, *c* the US Patent and Trademark Office (USPTO), *d* Index constructed in Sobel (2008), *e* US Census 2000 (<http://www.census.gov>), *f* Fraser Institute (Karabegovic and McMahon 2005)

Transformations (getting units right, like per capita, logs, changes, percent changes)

- generally try to not have the units be to different (billions for one variable, hundreds for another)
- use per capita to adjust for differences due to like the size of states or changes thru time
 - or sometimes per something else (school spending per student)
- sometimes the effects are “percentage” based (e.g., a 10% weight loss) – use logs or percent changes
- sometimes when we use “time series” data we do the year to year changes as the variables

2. Reading the data into “econometric” software that can run “regressions”

Everyone must obtain or use software called “gretl”, downloadable at <http://gretl.sourceforge.net/>

- will also be available on computers in Beaty 120 computer lab and Beatty Center Atrium
- if you know another program (like EViews) you are welcome to use it, but I can't help much

DOWNLOAD SOFTWARE FROM GRETL WEBSITE EXAMPLE (homepage):



URL: gretl.sourceforge.net

Gnu Regression, Econometrics and Time-series Library

gretl

Is a cross-platform software package for econometric analysis, written in the C programming language. It is free, open-source software. You may redistribute it and/or modify it under the terms of the GNU General Public License (GPL) as published by the [Free Software Foundation](#).

Warning: gretl on Ubuntu 11.04

The default gretl package for Ubuntu 11.04 (natty) is broken; it should be replaced with gretl version 1.9.5 (or higher), available via packages.debian.org.

gretl conference 2011

This took place at Nicolaus Copernicus University, Toruń, Poland on 16-17 June 2011; [details here](#).

Features

- Easy intuitive interface (now in French, Italian, Spanish, Polish, German, Basque, Portuguese, Russian, Turkish, Czech, Traditional Chinese, Albanian and Greek as well as English)
- A wide variety of estimators: least squares, maximum likelihood, GMM; single-equation and system methods
- Time series methods: ARMA, GARCH, VARs and VECMs, unit-root and cointegration tests, etc.
- Limited dependent variables: logit, probit, tobit, interval regression, models for count and duration data, etc.
- Output models as LaTeX files, in tabular or equation format
- Integrated scripting language: enter commands either via the gui or via script
- Command loop structure for Monte Carlo simulations and iterative estimation procedures
- GUI controller for fine-tuning Gnuplot graphs
- Links to [GNU R](#), [GNU Octave](#) and [Ox](#) for further data analysis

Data formats

Supported formats include: own XML data files; Comma Separated Values; Excel, Gnumeric and Open Document worksheets; Stata .dta files; SPSS .sav files; Eviews workfiles; Multi data files; own format binary databases (allowing mixed data frequencies and series lengths), RATS 4 databases and PC-Give databases. Includes a sample US macro database. See also the [gretl data page](#).

Download

gretl is available for both Windows and Mac, and can even be installed on computers on which you do not have administrative rights. There are even older versions available for older operating systems.

DOWNLOAD SOFTWARE FROM GRETL WEBSITE EXAMPLE (Windows download page):

System requirements

As of version 1.9.4, gretl requires Windows XP or higher and a processor that supports the SSE2 instruction set. SSE2 support is older systems, [see below](#).

Downloads

If you have the rights of a "power-user" or better on Windows, choose a self-installer from the first column below; just download

If you have no administrator rights on Windows choose a zip archive from the second column; unzip this in any location where you called "enable folders". The whole archive is in a directory called gretl. For example, if you unzip the archive into a directory name

The current "snapshot" of gretl is more up to date than the release: often it will contain bug-fixes but sometimes it will contain ne

	<i>self-installer</i>	OR	<i>zip archive (no admin rights)</i>
latest release (Dec 22, 2011)	gretl-1.9.7.exe		gretl-1.9.7-win32.zip
OR current snapshot	gretl-install.exe		gretl-win32.zip

The executables were cross-compiled under GNU/Linux using [mingw32](#) and [GTK for Windows](#) (thanks Tor Lillqvist!). The free insta

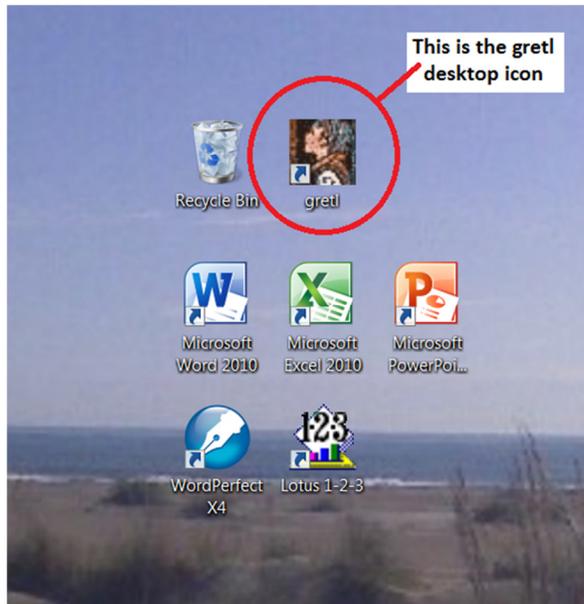
Optional extras you may wish to install

- X-12-ARIMA (seasonal adjustment, ARIMA)
- TRAMO/SEATS (seasonal adjustment, ARII)
- Datasets for Wooldridge, [Introductory Ec](#)
- Datasets for Gujarati, [Basic Econometrics](#)
- Datasets + scripts for Stock and Watson,
- Datasets for Davidson and MacKinnon, [Ec](#)
- Datasets for Marno Verbeek's [Guide to Mc](#)

choose to download “latest release” and use the one for “self-installer” if you do have administrative rights do not worry about downloading any of the “optional extras”

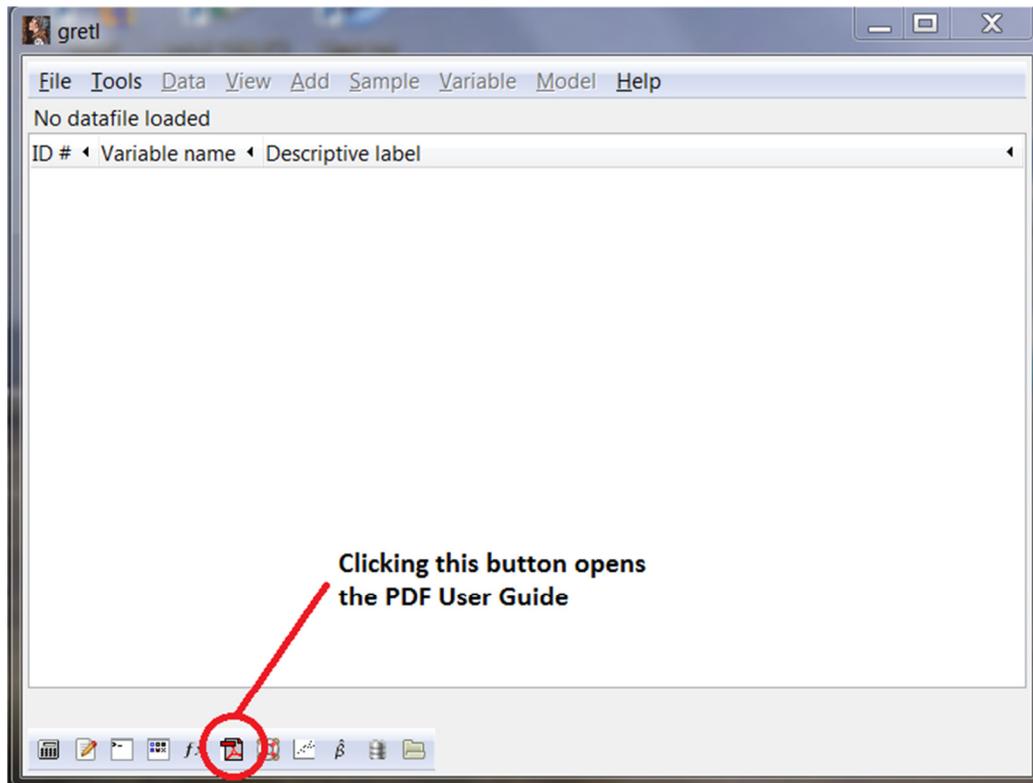
Once installed, open gretl by clicking on the gretl icon on your desktop:

GRETl DESKTOP ICON EXAMPLE:



When you open gretl you will get a main page that looks like this below:

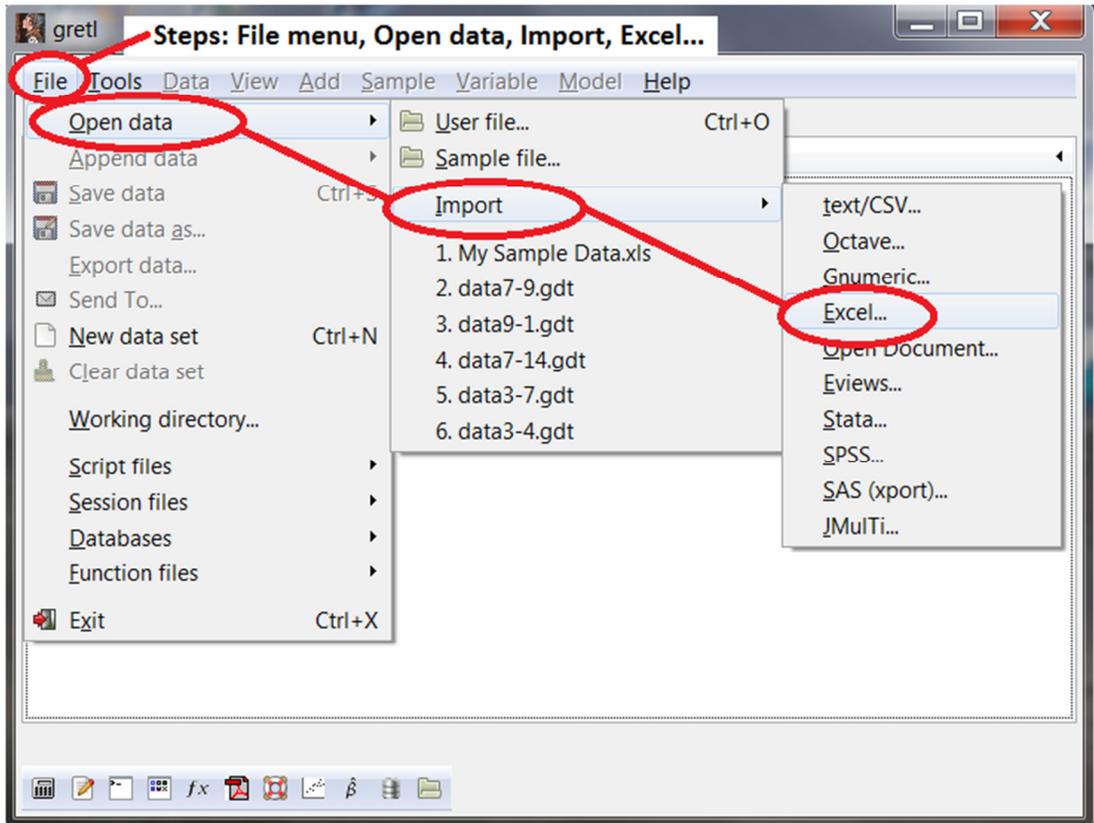
GRETl MAIN OPENING SCREEN EXAMPLE:



gretl has a built in user manual if you ever have questions or problems (PDF "users guide" button along bottom)

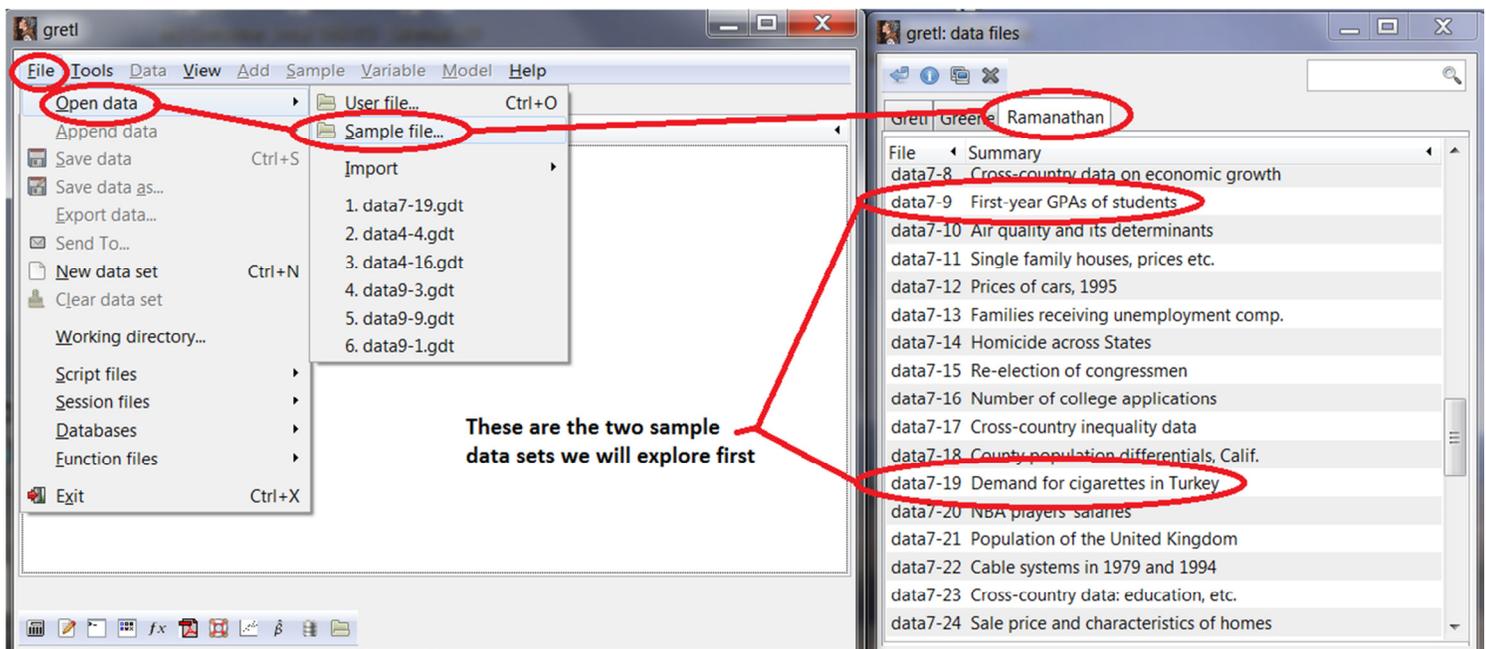
Reading data into gretl (File menu, Open Data, Import, then Excel if it's your data)

GRETl IMPORTING EXCEL DATA EXAMPLE:



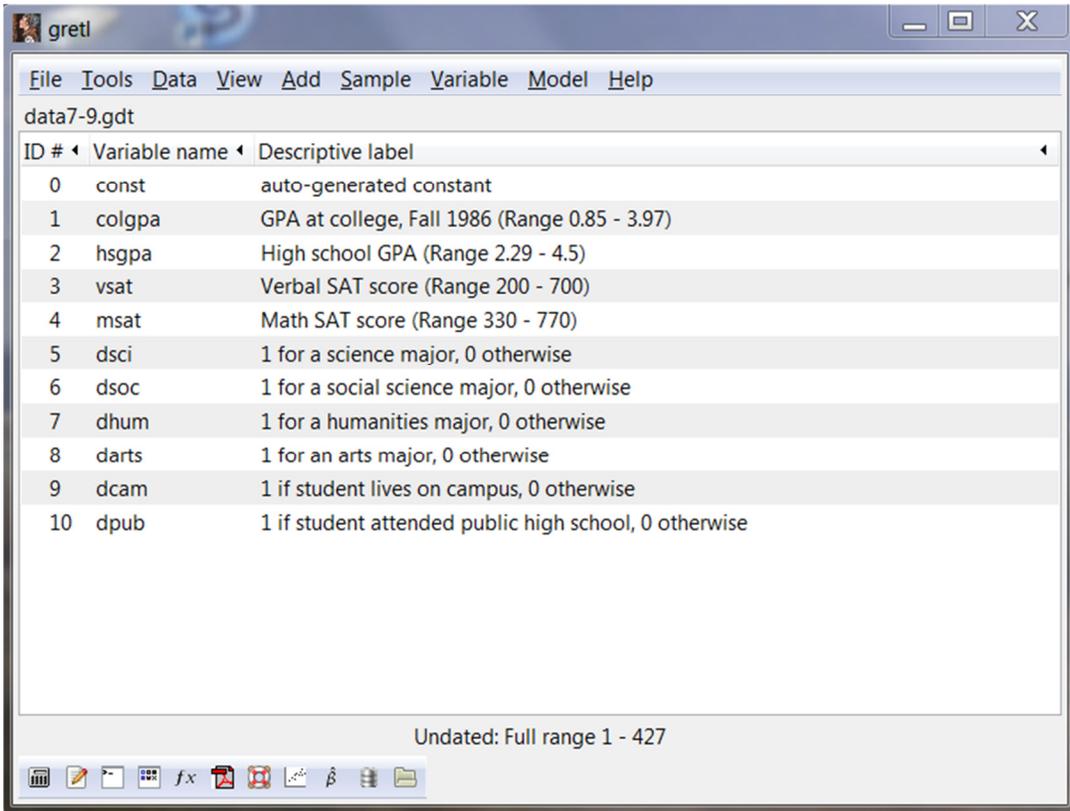
- For now, use sample dataset (File menu, Open Data, Sample file): Ramanathan data7-9 “First year GPAs of students”, we may also use data7-19 “Demand for cigarettes in Turkey” but we will start with the GPA data

GRETl OPENING RAMANATHAN GPA SAMPLE DATA EXAMPLE:



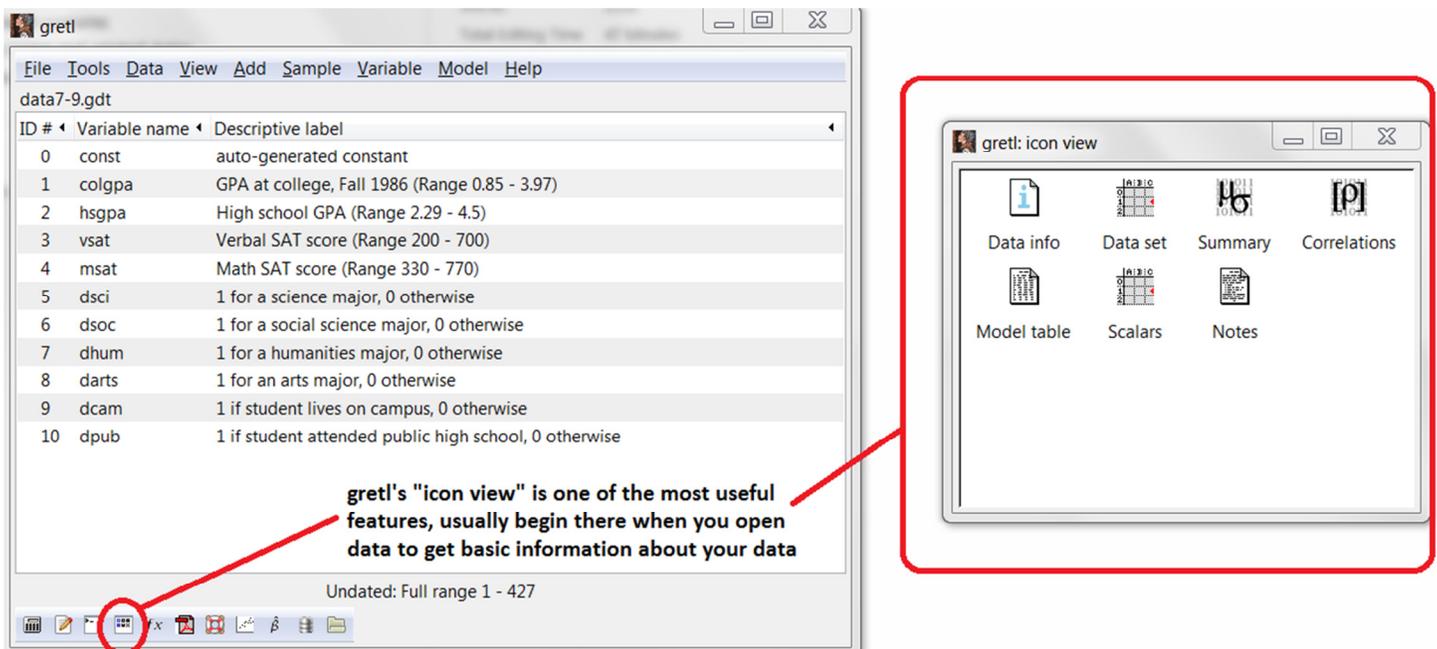
once you open the data, the main screen will look like this below:

GRET MAIN SCREEN AFTER OPENING SAMPLE GPA DATA:



one of the most used features is the "icon view", so let's open it. From this window you can look at your data and get basic statistical information about your data:

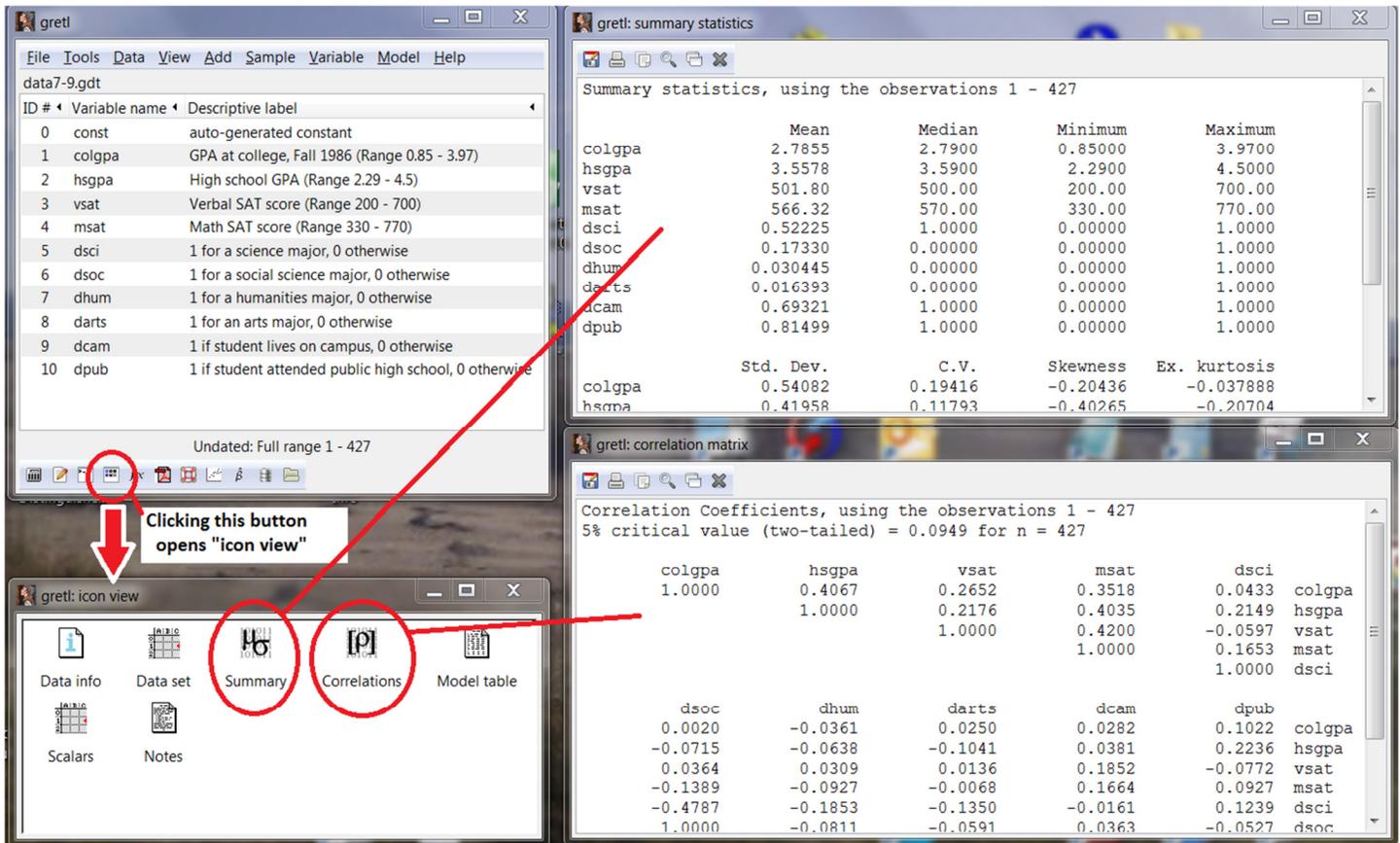
OPENING GRETL'S ICON VIEW:



3. First steps: Examine basic descriptive statistics, and examine correlations, creating new variables
 Note: these can be done in Excel too but are easier in gretl [Excel commands shown in brackets]

Examine “descriptive statistics” such as the mean (average), maximum, minimum of each variable
 - in gretl open “session icon view” and click “Summary” (or can use “View” menu then “Summary Statistics”)
 - [in Excel for data in cells A1 to A20 the commands are: =AVERAGE(A1:A20); =MAX(A1:A20); =MIN (A1:A20)]

USING GRETL'S ICON VIEW TO GET SUMMARY STATISTICS AND CORRELATIONS:



Examine “correlation coefficients”
 - in gretl either use “View” menu then “Correlation Matrix” or open “session icon view” and click “Correlations”
 - [in Excel for data in cells A1 to A20 and B1 to B20 the command is: =CORREL(A1:A20,B1:B20)]

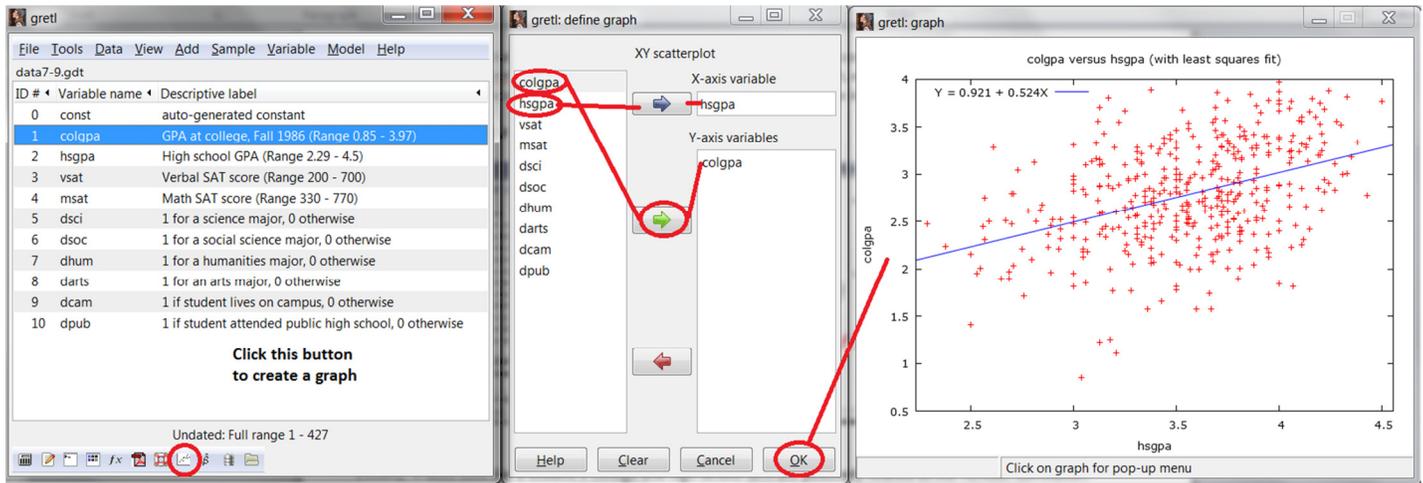
We will use these later, but for now just know the correlation coefficients measure whether two variables are positively or negatively related, and how strongly on a scale of zero to one. It shows the correlations between all of your variables, but most of the time you will only be interested in one or a few.

Example: what is the correlation between a student’s college GPA (colgpa) and highschool GPA (hsgpa)? +0.4067

Graph key variables against each other as XY plot – examine visual correlations

- in gretl click X-Y graph icon along bottom to create a chart, and choose the variables
- in Excel you have to create a “Scatter” type chart from the insert menu, after highlighting the data

CREATING A GRAPH OF TWO VARIABLES IN GRETL:

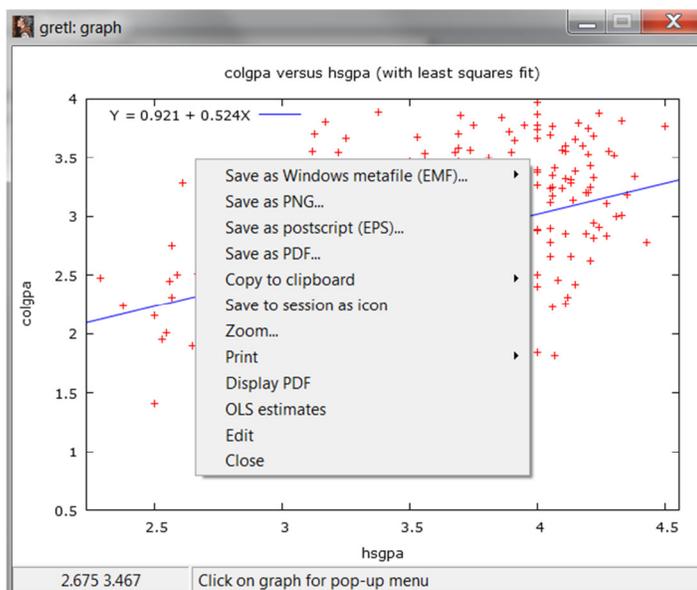


Visually, it does look like a student’s college and high school GPA are positively related as the +0.4067 correlation coefficient suggested.

Showing trend lines or “least squares fit” (which is an Ordinary Least Squares regression we are about to perform) in the graph is helpful to show you the relationship. The blue line in the graph above is a trend line. gretl usually adds a trend line automatically (in Excel you have right click any one data point in the scatter chart you create and choose “Add trendline”) In gretl it also gives you the equation of the trend line at the top of the chart.

Saving your graph: in gretl, click anywhere on the graph, can save as PDF, copy to clipboard (to say paste in your paper), and you can save it so you can reopen it later by choosing “Save to session as icon”

RIGHT CLICK IN THE GRAPH ANYWHERE TO BRING UP A MENU OF OPTIONS TO SAVE OR PRINT YOUR GRAPH:

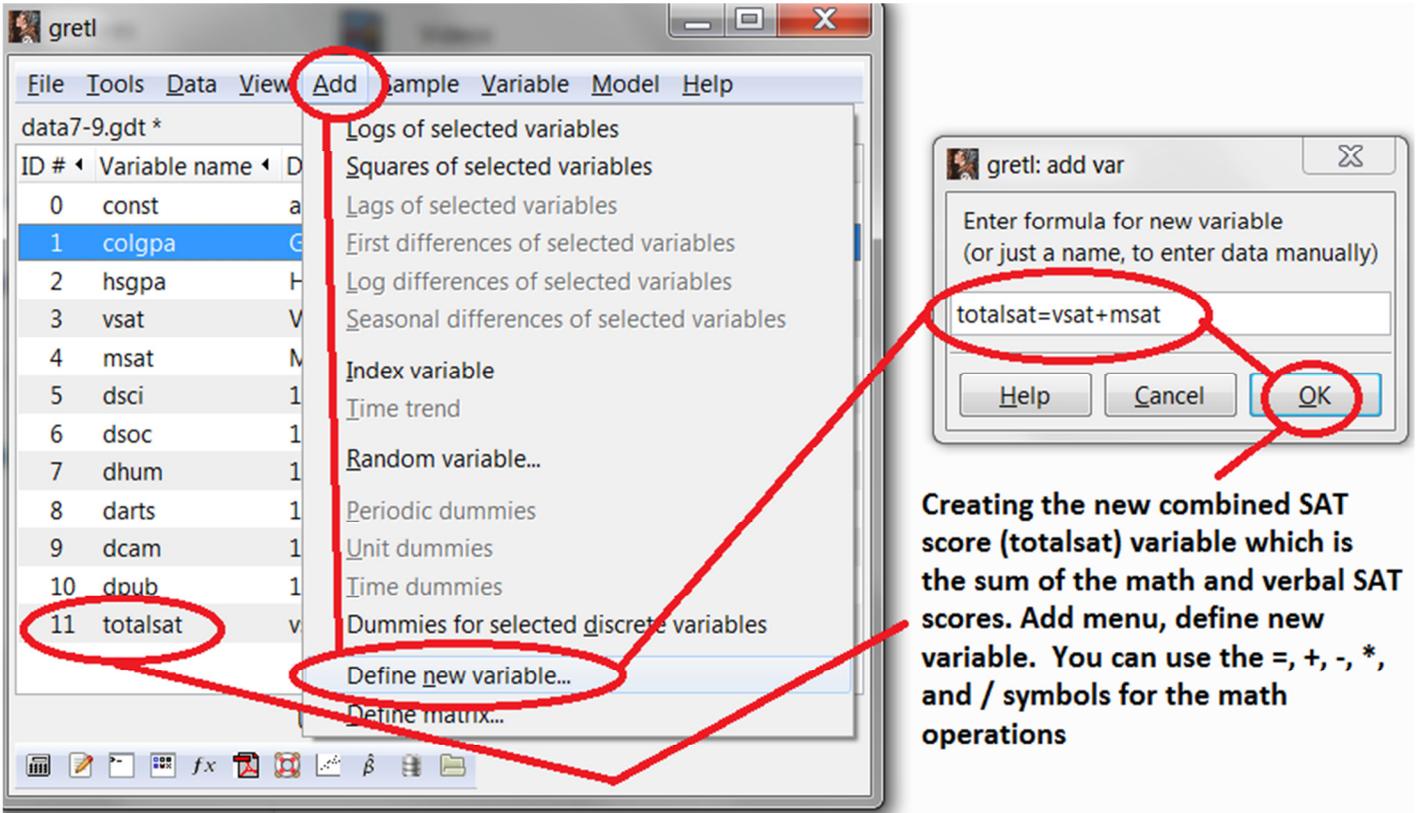


Creating a new variable in gretl from your existing variables

- Add menu, Define new variable, do as equation (e.g., $Z=X+Y$ if X and Y are in your data and Z is new)

In our sample data, there is both the student's verbal SAT (vsat) and math SAT (msat) score. But there is no variable for the student's combined math plus verbal score. So let's create a new variable named "totalsat" for that.

CREATING A NEW VARIABLE IN GRETL (COMBINED SAT SCORE EXAMPLE):



gretl can also create "logs" or "squared" versions of your variables automatically, these are the first two entries in the Add menu shown above where we picked "Define new variable".

gretl can also create random variables for you.

Econometric Analysis – Dr. Sobel

Econometrics Session 2:

4. Perform Ordinary Least Squares (OLS) Regression Analysis in gretl

What it is, what it does, and why we do it: Regression analysis is basically fitting and estimating the trend line in the X-Y graphs. There are many different types of regressions, but the most basic and easiest is Ordinary Least Squares (or OLS). This method fits the line through the data that minimizes the sum of the squared differences between all points and the trend line. Doing a regression allows us to include many different variables, so we can examine the relationship between say, X and Y, but we can also “control” for other things that might impact the relationship or variables.

The OLS regression model is only valid when certain assumptions are met. One is that the “errors” between the points and the trend line are “random” or “normally distributed”, and another is that the dependent variable is a “continuous” variable (meaning it can take on a wide range of numbers including the possibility of non-whole numbers). One example of when this is violated is when the dependent variable is a count variable (takes the values zero, one, two, three, etc.). For these cases you cannot use OLS regressions. We will examine some of these in part 3 of the econometric section.

A regression has a “dependent” variable and also “independent” variables (sometimes called “endogenous” and “exogenous” variables, respectively). The dependent variable is the one being explained or predicted and the independent variables are the ones we are using to explain it. The regression allows us to see how a change in any one of the independent variables impacts the dependent variable, holding constant the other variables. For every regression you must specify which one is the dependent variable and also which one or many to include as independent variables.

How to run a basic OLS regression in gretl: Click the button in the main window that looks like the greek letter beta (β). This will open a window where you can specify the “dependent” variable and “independent” variables. Highlight the variable you want to be the dependent variable and click the blue arrow, then highlight anything you want to use as independent (explanatory) variables and hit the green arrow to include them in the bottom list.

Let’s try to explain/predict a student’s college GPA (colgpa) using their high school GPA (hsgpa) using an OLS regression

RUNNING A BASIC OLS REGRESSION IN GRETL EXAMPLE (explaining college GPA with high school GPA):

Clicking this button is how you run a basic OLS regression

This is the “constant” and it should automatically be there (always include it)

	coefficient	std. error	t-ratio	p-value
const	0.920577	0.204631	4.499	8.83e-06 ***
hsgpa	0.524173	0.0571206	9.177	1.95e-18 ***

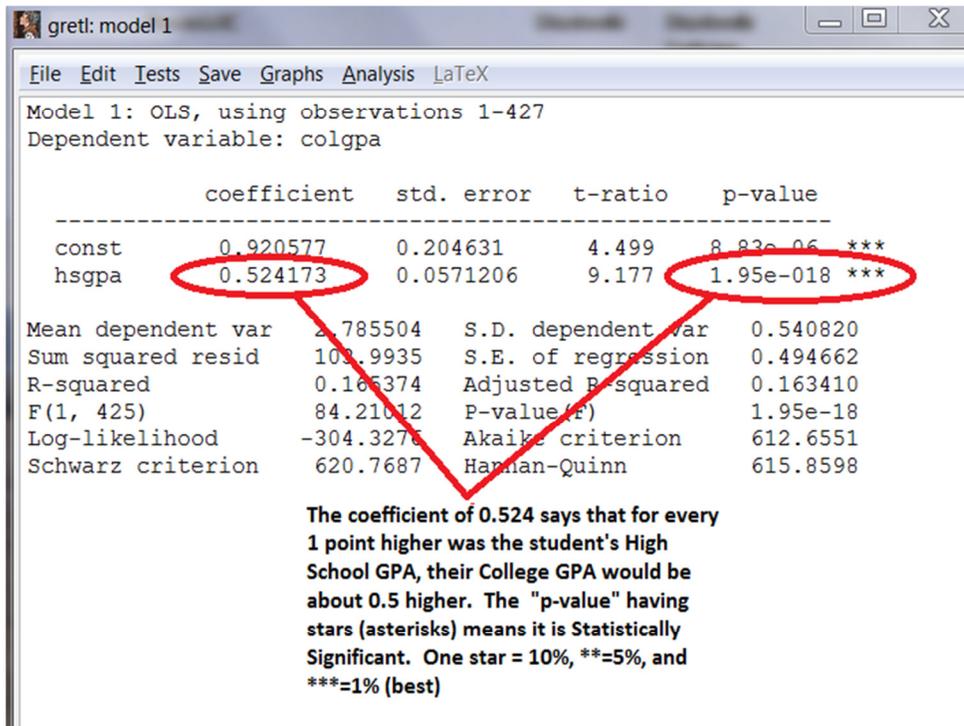
		S.D. dependent var	
Mean dependent var	2.785504	0.540820	
Sum squared resid	103.9935	S.E. of regression	0.494662
R-squared	0.165374	Adjusted R-squared	0.163410
F(1, 425)	84.21012	P-value(F)	1.95e-18
Log-likelihood	-304.3276	Akaike criterion	612.6551
Schwarz criterion	620.7687	Hannan-Quinn	615.8598

****NOTE:** You should ALWAYS include a “constant” as an independent variable. This is basically allowing the trend line to have a y-intercept. If you do not include it, it will force the trend line to go through the origin (0,0 point), which can mess up your regression. gretl will normally automatically include a constant for you in the list (but make sure it is indeed there as “const” in the independent variable list). But we never care if the constant is significant and never interpret it.

How to interpret the OLS regression results (“coefficient estimates”, “statistical significance”, & “goodness of fit”)

The “coefficient” estimate is the impact of a one unit change in that independent variable on the dependent variable. It is the slope of the line in the X-Y plot. Sometimes we only care if the coefficient estimate is positive or negative (does that variable positively or negatively impact the other variable). But we ALWAYS care if it is “statistically significant” (using p-value and the asterisks beside it). In research papers we report the t-ratio or standard (std.) error that gives us the p-value, and put the asterisks/stars beside the coefficient estimate. You can round it to fewer decimal places.

COEFFICIENT ESTIMATES FROM A REGRESSION, HOW TO INTERPRET THEM:



Model 1: OLS, using observations 1-427
Dependent variable: colgpa

	coefficient	std. error	t-ratio	p-value
const	0.920577	0.204631	4.499	8.83e-06 ***
hsgpa	0.524173	0.0571206	9.177	1.95e-018 ***

Mean dependent var	2.785504	S.D. dependent var	0.540820
Sum squared resid	103.9935	S.E. of regression	0.494662
R-squared	0.166374	Adjusted R-squared	0.163410
F(1, 425)	84.21012	P-value(F)	1.95e-18
Log-likelihood	-304.3276	Akaike criterion	612.6551
Schwarz criterion	620.7687	Hannan-Quinn	615.8598

The coefficient of 0.524 says that for every 1 point higher was the student's High School GPA, their College GPA would be about 0.5 higher. The "p-value" having stars (asterisks) means it is Statistically Significant. One star = 10%, **=5%, and *=1% (best)**

If one wanted to “predict” a student’s college GPA from their high school GPA one would use the equation:

$$\text{Predicted college GPA} = 0.920577 + 0.524173 * \text{High School GPA}$$

where the first number comes from the coefficient estimate for the “constant” (const)

Statistical significance is very important. This is shown by the stars or asterisks to the far right for each variable. Three stars is the “best”, two is “good”, and one is just “okay”. If there are no stars, the variable is said to be insignificant, and so the coefficient might as well be zero (meaning you could exclude it from the regression, it’s not an important predictor or doesn’t explain the other variable very well).

These stars/asterisks are based on the number in the column titled “p-value”. This is the probability value for the statistical test. A p-value of 1% (0.01) or lower gives three stars, a p-value greater than 1% (0.01) but less than 5% (0.05) gives two stars, and a p-value greater than 5% (0.05) but less than 10% (0.10) gives one star. If the p-value is higher than 10% (0.10) the variable is said to be “insignificant”. In writing up your results, for example, for a variable with two stars we usually would say that “the variable is significant at the 5% level”. Note the example has one using scientific notation, the “e-018” means there the true decimal is 18 places to the left (so the p-value is like 0.00000000000000000195)

Optional: The p-value is based on a “t-test”. It is calculated by dividing the coefficient estimate by the standard error. This gives the “t-ratio” reported in the results. If the t-ratio is larger (in absolute value) than roughly 1.645 the variable will be significant. T-ratios of 2 or more give statistical significance at the 5% level. This is called “Hypothesis testing” and it is important to specify your hypothesis and for some tests to know what the “null” is. Sometimes we use “Standard errors” (std. error) to construct a “confidence interval” for the estimate using the standard error times two (the impact is 0.524173 plus or minus $2 * 0.0571206$, that is 0.52 ± 0.114241 , or that the estimated impact of high school GPA on college GPA probably is in the range 0.409932 to 0.638414, with our best estimate being 0.524173).

How good the model is (goodness of fit) is measured by the R-squared or Adjusted R-squared (adjusted is better)

- R-squared ranges from zero to one (one would be a perfect fit)
- interpreted as “the percent of the variation explained”
- gretl will also save and report the “Log-likelihood” or “lnL” but you don’t have to show this

R-SQUARED FROM A REGRESSION, HOW GOOD THE REGRESSION FITS THE DATA:

Model 1: OLS, using observations 1-427
Dependent variable: colgpa

	coefficient	std. error	t-ratio	p-value
const	0.920577	0.204631	4.499	8.83e-06 ***
hsgpa	0.524173	0.0571206	9.177	1.95e-018 ***

Mean dependent var	2.785504	S.D. dependent var	0.540820
Sum squared resid	109.0935	S.E. of regression	0.504662
R-squared	0.165374	Adjusted R-squared	0.163410
F(1, 425)	81.21012	P-value (F)	1.95e-18
Log-likelihood	-304.3276	Akaike criterion	612.6551
Schwarz criterion	620.7687	Hannan-Quinn	615.8598

The "R-squared" tells us how good the model is at fitting the data. This shows that a student's high school GPA alone predicts about 16 percent of the variation in their college GPA. Normally (but not always) people report the "Adjusted R-squared" (which adjusts for the number of variables in the regression). You can report both if you like. Just make sure to tell which you report.

So, the data also includes the student’s math and verbal SAT scores. So which is the best predictor of a student’s college GPA? is it their high school GPA? their verbal SAT score? their math SAT score? To find out let’s do a regression for each and compare the adjusted R-squareds:

Model 1: OLS, using observations 1-427
Dependent variable: colgpa

	coefficient	std. error	t-ratio	p-value
const	1.62845	0.151348	10.76	4.85e-024 ***
msat	0.00204309	0.000263713	7.747	6.96e-014 ***

Mean dependent var	2.785504	S.D. dependent var	0.540820
Sum squared resid	109.1797	S.E. of regression	0.506846
R-squared	0.123752	Adjusted R-squared	0.121690
F(1, 425)	60.02224	P-value (F)	6.96e-14
Log-likelihood	-314.7177	Akaike criterion	633.4355
Schwarz criterion	641.5491	Hannan-Quinn	636.6402

Model 2: OLS, using observations 1-427
Dependent variable: colgpa

	coefficient	std. error	t-ratio	p-value
const	1.99740	0.141279	14.14	1.81e-037 ***
vsat	0.00157055	0.000277004	5.670	2.64e-08 ***

Mean dependent var	2.785504	S.D. dependent var	0.540820
Sum squared resid	115.8373	S.E. of regression	0.522071
R-squared	0.070319	Adjusted R-squared	0.068132
F(1, 425)	32.14621	P-value (F)	2.64e-08
Log-likelihood	-327.3551	Akaike criterion	658.7102
Schwarz criterion	666.8238	Hannan-Quinn	661.9150

- A) Adjusted R-squared from a regression of College GPA on just High School GPA: 0.1634
- B) Adjusted R-squared from a regression of College GPA on just Math SAT Score: 0.1217
- C) Adjusted R-squared from a regression of College GPA on just Verbal SAT Score: 0.0681

Based on this, the verbal SAT score is the worst predictor, and the high school GPA is the best predictor. However, each is statistically significant at the best 1% level (three stars for each one), so all are important predictors of college GPA, but the high school GPA just explains it more closely.

Regressions with multiple variables: Now let's see if using all three at the same time helps. That is, can we do even better at predicting college GPA if we use information on high school GPA and information on verbal and math SAT scores. So run a regression including all three as independent variables at the same time:

A REGRESSION WITH MULTIPLE INDEPENDENT VARIABLES:

gretl: model 3

Model 3: OLS, using observations 1-427
Dependent variable: colgpa

	coefficient	std. error	t-ratio	p-value	
const	0.423249	0.219749	1.926	0.0548	*
hsgpa	0.398349	0.0605865	6.575	1.44e-010	***
msat	0.00101521	0.000293603	3.458	0.0006	***
vsat	0.000737453	0.000280682	2.627	0.0089	***
Mean dependent var	2.785504	S.D. dependent var	0.540820		
Sum squared resid	97.16385	S.E. of regression	0.479272		
R-squared	0.220187	Adjusted R-squared	0.214657		
F(3, 423)	39.81265	P-value (F)	1.11e-22		
Log-likelihood	-289.8245	Akaike criterion	587.6490		
Schwarz criterion	603.8761	Hannan-Quinn	594.0584		

All three are still significant. The adjusted R-squared now rises to 0.214657, meaning we are doing an even better job of predicting/explaining college GPA (we now explain 21% of the variance in it using these three variables). If we did want to do an equation to predict a student's college GPA it would now be: COLGPA = 0.42 + .0398 HSGPA + 0.0010*msat + 0.0007*vsat. Note that because the coefficient on math SAT is bigger, this means a one unit higher math SAT score has a greater impact on college performance than a one unit higher verbal SAT score.

Note that the data also includes information about the student's major, whether they live on campus, and whether they went to a public (versus private) high school. Perhaps some majors are easier and some harder, or that living on campus matters, and we should control for this to do an even better job.

Note that these new variables are "dummy" or "indicator" variables that only take the value of zero/one. This is fine for an independent variable (if it was the dependent variable we would need to use something different from OLS, we will get to this in part 3). The coefficient on a dummy variable is easy to interpret, as it's the effect of going from zero to one, or rather the effect of having that characteristic versus not having it.

A REGRESSION WITH MULTIPLE INDEPENDENT VARIABLES (including "dummy" or "indicator" variables):

gretl: model 4

Model 4: OLS, using observations 1-427
Dependent variable: colgpa

	coefficient	std. error	t-ratio	p-value	
const	0.367296	0.224302	1.638	0.1023	
hsgpa	0.405914	0.0634178	6.401	4.17e-010	***
msat	0.00108585	0.000302752	3.587	0.0004	***
vsat	0.000725891	0.000289903	2.504	0.0127	**
dsci	-0.0273225	0.0573192	-0.4767	0.6338	
dsoc	0.0561481	0.0727785	0.7715	0.4409	
dhum	-0.00405912	0.141771	-0.02863	0.9772	
darts	0.228650	0.188921	1.210	0.2269	
dcam	-0.0407050	0.0521617	-0.7804	0.4356	
dpub	0.0294027	0.0630396	0.4664	0.6412	
Mean dependent var	2.785504	S.D. dependent var	0.540820		
Sum squared resid	96.20445	S.E. of regression	0.480319		
R-squared	0.227887	Adjusted R-squared	0.211223		

None of these indicator/dummy variables have stars by them. None of them are significant. So apparently living on campus, public vs. private high school, nor major really matters once you are already accounting for the student's High School GPA, and math and verbal SAT scores. Note that the adjusted R-squared is lower than the model with only those three variables.

Keeping/saving your results – What to report or show in a paper, how to copy it into a word file

FROM A WINDOW OF REGRESSION RESULTS:

- File menu, then “Save to session as icon” or you can “print” or save as a word document
- Save to session as icon is best, as you can then make a table of many results easily
it will save them as icons named “model 1”, “model 2”, etc.

SAVING YOUR REGRESSION RESULTS:

To save your results for later you can save them as an icon. Then these results will be an icon in the icon view

You can also copy and paste the results right into a document file using “Edit” menu, then “copy” and then paste

Generally, we run many different regressions and show them in a table all at the same time. In gretl this is very easy as long as you save each set of results as an icon. Once several (up to six maximum) models are saved as icons, open the icon view. Then you can either drag & drop the models onto the icon called “Model table” or you can right click the model and choose “Add to model table”. The order in which they appear in the table depends on the order in which you put them in the model table. Note that gretl’s model table will show everything you need to put in your research paper and more (we usually don’t need/show “lnL” but you can leave it in, and the full name for it is “Log-likelihood” if you want to put that in your table in the paper). Once created, click on the “Model table” icon to view the table of results.

CREATING A TABLE OF RESULTS FROM SEVERAL REGRESSIONS:

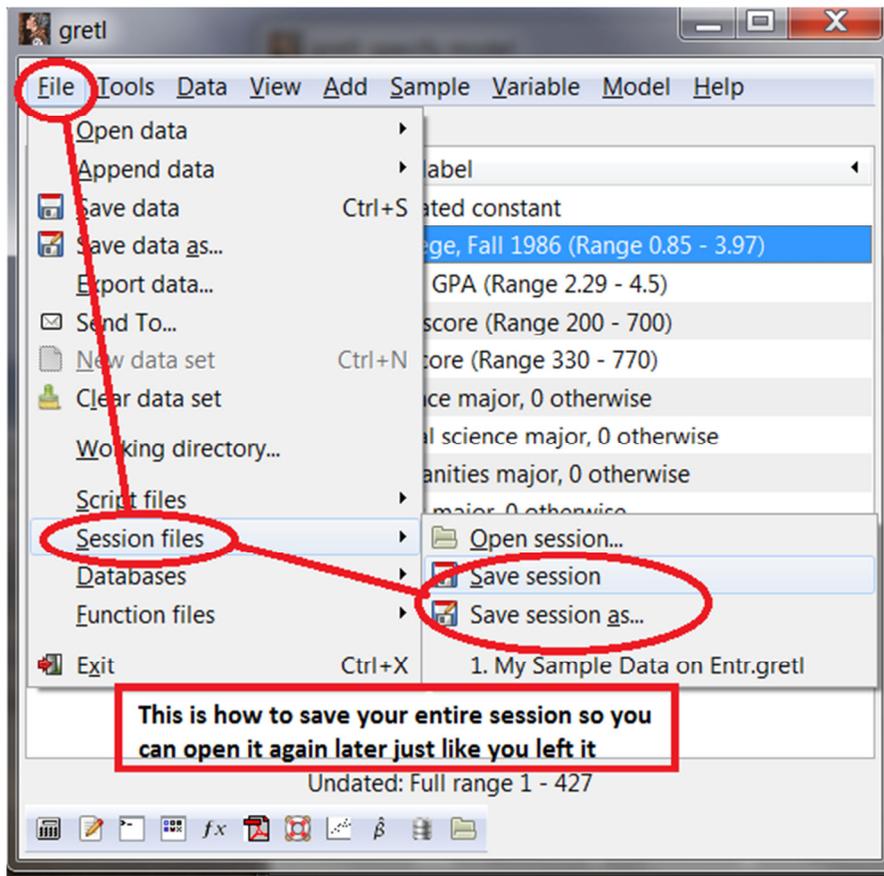
To create a table of results from several regressions using the same dependent variable, after saving each model result as an icon, then drag the icon for each model (in the order you want) onto the Model Table icon. Or you can right click the model's icon and select "add to model table". This table can be saved, or even copied and pasted into Microsoft Word or whatever software you are writing your paper using.

	(1)	(2)	(3)	(4)	(5)
const	0.9206** (0.2046)	1.628** (0.1513)	1.997** (0.1413)	1.402** (0.1696)	0.4232* (0.2197)
hsgpa	0.5242** (0.05712)				0.3983** (0.06059)
msat		0.002043** (0.0002637)		0.001695** (0.0002881)	0.001015** (0.0002936)
vsat			0.001571** (0.0002770)	0.0008444** (0.0002938)	0.0007375** (0.0002807)
Adj. R**2	0.1634	0.1217	0.0681	0.1364	0.2147
lnL	-304.3	-314.7	-327.4	-310.6	-289.8

Standard errors in parentheses
* indicates significance at the 10 percent level
** indicates significance at the 5 percent level

How to save the session to reopen later easily including any results or graphs you have done
“File” menu, “Session files”, “Save Session as”

SAVING YOUR ENTIRE SESSION SO YOU CAN REOPEN IT LATER WITH EVERYTHING THERE:



Frequently encountered basic problems in regressions:

One variable turns insignificant when you include additional variables – this can happen because the new variable simply is a better predictor (so given we have the new variable, the old one just doesn't matter in explaining the data), or it can happen because the two variables are measuring the same thing, a problem known as “multicollinearity”.

Multicollinearity – don't include two independent variables that are too closely correlated (meaning they measure roughly the same thing). If two variables have a high correlation coefficient (near 1, or let's say at least above 0.7 or 0.8), then don't include both variables. If you do, one or both will be insignificant, even if they matter. Check your correlation coefficients!

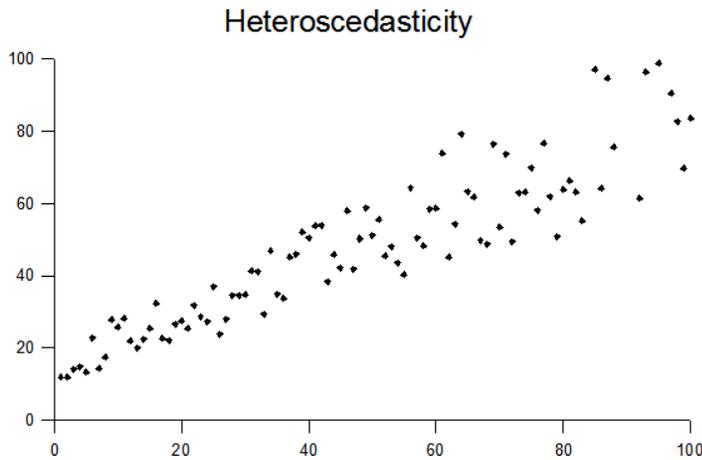
Outliers – one data point way off the line from the rest can really change your results. So always look at the graph of the data to see if there is a bad data point, then exclude that observation from the regression.

Selecting the sample / Working with subsamples / splitting the data – sometimes we may want to see if things are different among subsets of the data (men vs. women) or to exclude some observations that are outliers. To select to run the regression on only part (a subset) of your data: use the “Sample” menu from the top, “restrict based on criterion”, then use something like: $income > 50000$ or $year = 1995$ using one of your variables. To return to using the full sample, go to “Sample” menu, “restore full range”.

Note: After running a regression, in the results table there is a “tests” menu that has some tests for frequent problems including “influential observations” (outliers) and heteroskedasticity:

Correcting for “heteroskedasticity” - one of the assumptions of the basic OLS model is that the errors of the points around the line are uniform or random. Problems occur when the scatter of points looks like a cone:

HETEROSCEDASTICITY IN THE DATA:



Testing for this from the model results window: “Tests”, “Heteroskedasticity”, “White's test for heteroskedasticity”

White's test for heteroskedasticity
 OLS, using observations 1-427
 Dependent variable: uhat^2
 Omitted due to exact collinearity: sq_pub_hsgpa

	coefficient	std. error	t-ratio	p-value
const	-0.549304	0.945395		
hsgpa	0.459160	0.574973		
pub_hsgpa	0.0322970	0.106478		
sq_hsgpa	-0.0682941	0.0873609		
x2_x3	-0.00552510	0.0297412		

Unadjusted R-squared = 0.005434

Test statistic: $TR^2 = 2.320311$,
 with p-value = $P(\text{Chi-square}(4) > 2.320311) = 0.677074$

because this is not significant (not less than 0.10, we do NOT have a problem with heteroskedasticity)

If it is significant then you need to correct for it (e.g., a p-value smaller than 0.10) using the check box for “Robust standard errors” when you pick the variables for the regression:

Model 1: OLS, using observations 1-427
 Dependent variable: colgpa
 Heteroskedasticity-robust standard errors, variant HC1

	coefficient	std. error	t-ratio	p-value
const	0.920577	0.195798	4.702	3.49e-06 ***
hsgpa	0.524173	0.0544803	9.621	5.90e-20 ***

Mean dependent var 2.785504 S.D. dependent var 0.540820
 Sum squared resid 103.9935 S.E. of regression 0.494662
 R-squared 0.165374 Adjusted R-squared 0.163410
 F(1, 425) 92.56995 P-value (F) 5.90e-20
 Log-likelihood -304.3276 Akaike criterion 612.6551
 Schwarz criterion 620.7687 Hannan-Quinn 615.8598

****YOU SHOULD ALWAYS CHECK YOUR REGRESSION FOR THIS****

****SOME PEOPLE JUST ALWAYS CHECK THIS BOX****

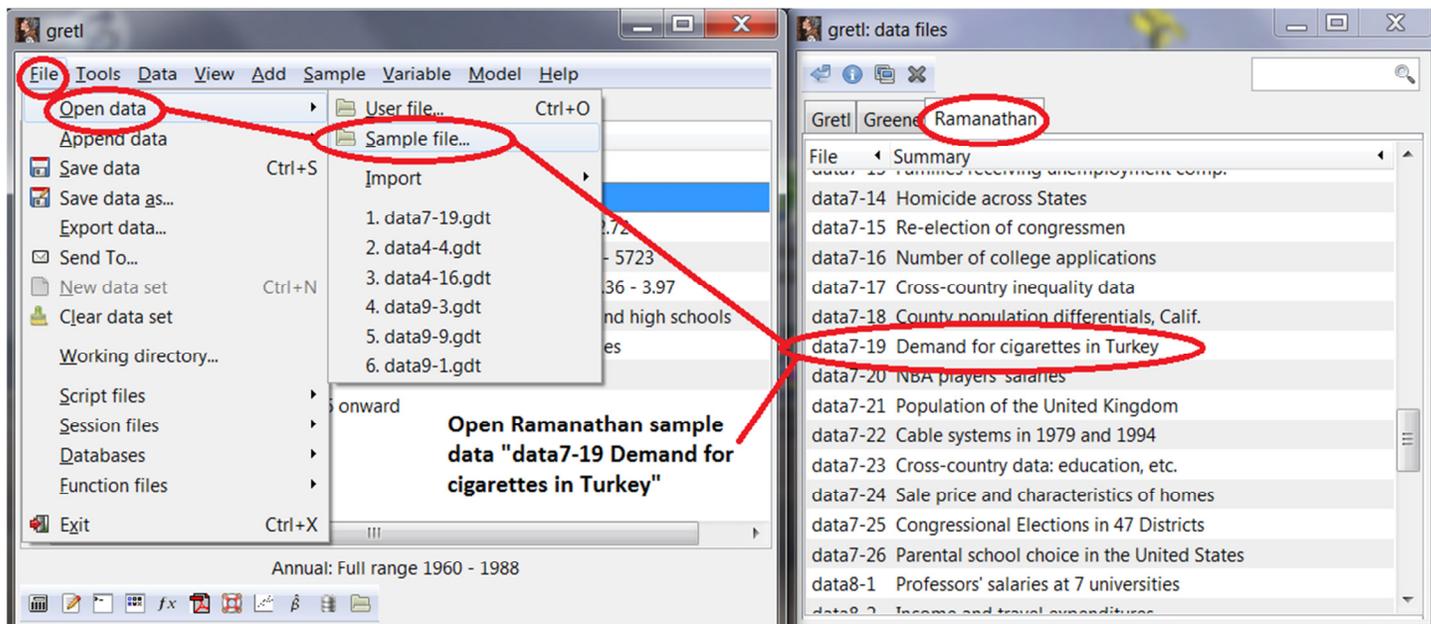
****IF YOU DO THIS, PLEASE NOTE IN YOUR RESULTS TABLE: “ROBUST STANDARD ERRORS”****

Econometric Analysis – Dr. Sobel

Econometrics Application to the Law of Demand:

1. Read in the Sample Data Set Ramanathan data7-19 “Demand for cigarettes in Turkey”

OPENING RAMANATHAN SAMPLE DATA 7-19:



2. Here is the data that will be in your data set:

The screenshot shows the gretl software interface displaying the data set 'data7-19.gdt'. The data set is summarized in the following table:

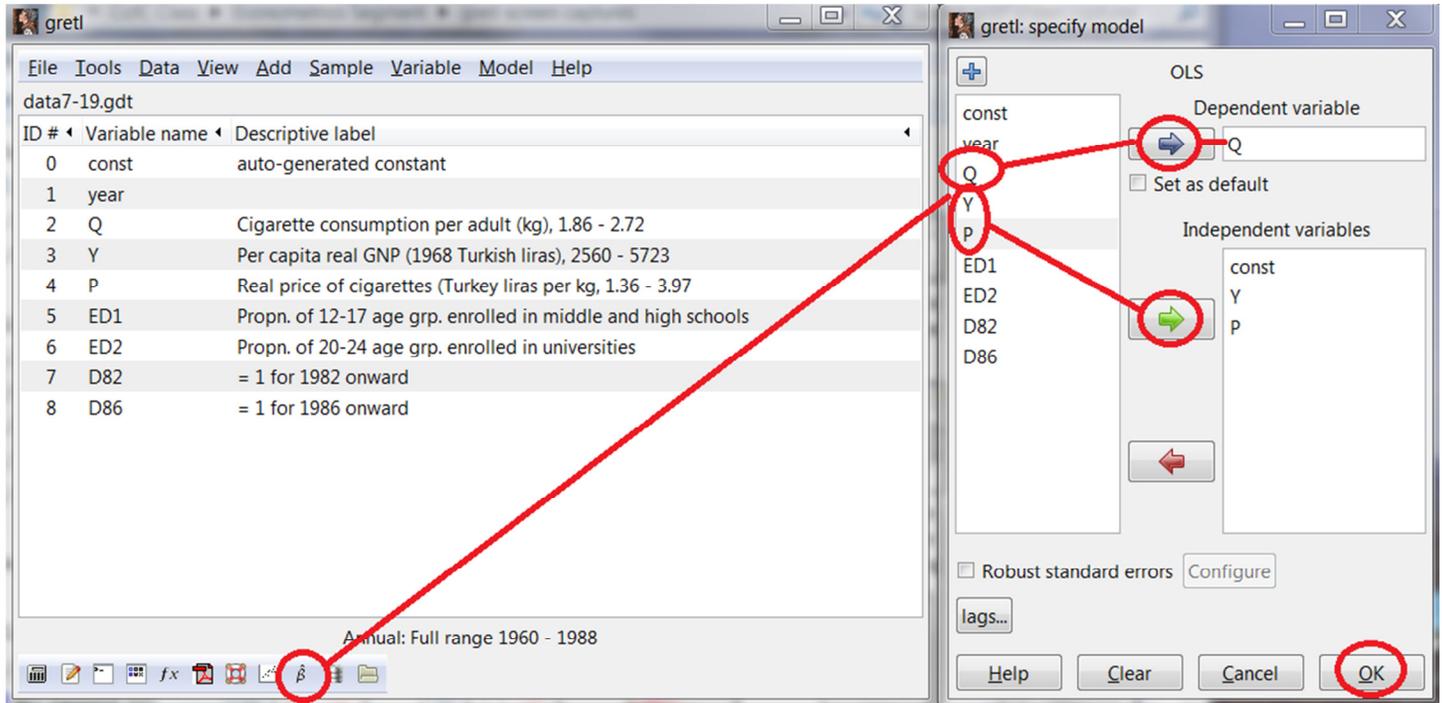
ID #	Variable name	Descriptive label
0	const	auto-generated constant
1	year	
2	Q	Cigarette consumption per adult (kg), 1.86 - 2.72
3	Y	Per capita real GNP (1968 Turkish liras), 2560 - 5723
4	P	Real price of cigarettes (Turkey liras per kg, 1.36 - 3.97
5	ED1	Propn. of 12-17 age grp. enrolled in middle and high schools
6	ED2	Propn. of 20-24 age grp. enrolled in universities
7	D82	= 1 for 1982 onward
8	D86	= 1 for 1986 onward

Annual: Full range 1960 - 1988

We have everything we need to estimate a demand curve: Q is the quantity demanded, Y is income, and P is price.

Based on economic theory, our hypothesis should be that the coefficient on price is NEGATIVE and significant (the law of demand), and the coefficient on income should be POSITIVE and significant (if cigarettes are a normal good).

3. Let's estimate the demand curve, select Q as the dependent variable, and as independent variables have the constant (const), income (Y), and price (P):



Model 1: OLS, using observations 1960-1988 (T = 29)
Dependent variable: Q

	coefficient	std. error	t-ratio	p-value
const	1.65654	0.123678	13.39	3.53e-013
Y	0.000344100	5.27935e-05	6.518	6.56e-07 ***
P	-0.423295	0.0969440	-4.366	0.0002 ***

Mean dependent var 2.204655 S.D. dependent var 0.243190
Sum squared resid 0.595167 S.E. of regression 0.151298
R-squared 0.640589 Adjusted R-squared 0.612942
F(2, 26) 23.17031 P-value (F) 1.67e-06
Log-likelihood 15.20081 Akaike criterion -24.40161
Schwarz criterion -20.29973 Hannan-Quinn -23.11695
rho 0.536727 Durbin-Watson 0.911596

Our hypotheses are confirmed, income (Y) has a positive and significant coefficient, and price (P) has a negative and significant coefficient.

If we wished to interpret the coefficients, they would say that for every 10,000 Turkish Liras of higher income, cigarette consumption is 3.4 kg higher, and that for every 1 Turkish Lira per kg the price increases, cigarette consumption falls by 0.42 kg.

However, we need to check to be sure our model doesn't suffer from any common problems, most importantly heteroskedasticity. So run the White's test:

gretl: model 1

	d. error	t-ratio	p-value
const	123678	13.39	3.53e-013 ***
Y	27935e-05	6.518	6.56e-07 ***
P	0969440	-4.366	0.0002 ***

White's test for heteroskedasticity
OLS, using observations 1960-1988 (T = 29)
Dependent variable: uhat^2

	coefficient	std. error	t-ratio	p-value
const	0.224376	0.134559	1.667	0.1090
Y	-0.000288480	0.000118342	-2.438	0.0229 **
P	0.297077	0.210234	1.413	0.1710
sq_Y	6.36452e-08	3.22655e-08	1.973	0.0607 *
X2_X3	-0.000101817	8.49743e-05	-1.198	0.2430
sq_P	0.0310950	0.0499291	0.6228	0.5396

Warning: data matrix close to singularity!
Unadjusted R-squared = 0.475305
Test statistic: $TR^2 = 13.783847$,
with p-value = $P(\text{Chi-square}(5) > 13.783847) = 0.017042$

Yes, we do have a problem with heteroskedasticity

The test p-value is smaller than 0.10 (10%), so yes we have a problem with heteroskedasticity. We need to re-run the regression using "Robust Standard Errors":

gretl: specify model

OLS

Dependent variable: Q

Independent variables: const, Y, P

Robust standard errors

gretl: model 2

Model 2: OLS, using observations 1960-1988 (T = 29)
Dependent variable: Q
HAC standard errors, bandwidth 2 (Bartlett kernel)

	coefficient	std. error	t-ratio	p-value
const	1.65654	0.126315	13.11	5.71e-013 ***
Y	0.000344100	7.19664e-05	4.781	5.98e-05 ***
P	-0.423295	0.0970912	-4.360	0.0002 ***

Mean dependent var 2.204655 S.D. dependent var 0.243190
Sum squared resid 0.595167 S.E. of regression 0.151298
R-squared 0.640589 Adjusted R-squared 0.612942
F(2, 26) 11.48910 P-value(F) 0.000266
Log-likelihood 15.20081 Akaike criterion -24.40161
Schwarz criterion -20.29973 Hannan-Quinn -23.11695
rho 0.536727 Durbin-Watson 0.911596

Even after adjusting for heteroskedasticity, our hypotheses are both confirmed again, but these are better estimates. So now, let's save this model as an icon:

gretl: model 2

File Edit Tests Save Graphs Analysis LaTeX

Save as...
Save to session as icon
Save as icon and close

	std. error	t-ratio	p-value
const	0.126315	13.11	5.71e-013 ***
Y	7.19664e-05	4.781	5.98e-05 ***
P	0.0970912	-4.360	0.0002 ***

Mean dependent var 2.204655 S.D. dependent var 0.243190
Sum squared resid 0.595167 S.E. of regression 0.151298
R-squared 0.640589 Adjusted R-squared 0.612942
F(2, 26) 11.48910 P-value(F) 0.000266
Log-likelihood 15.20081 Akaike criterion -24.40161
Schwarz criterion -20.29973 Hannan-Quinn -23.11695
rho 0.536727 Durbin-Watson 0.911596

Now, let's try including two of the other variables in the data set, measuring educational enrollment, ED1 and ED2 which are the proportions of the 12-17 age group enrolled in middle and high schools, and the proportion of the 20-24 age group enrolled in universities.

The screenshot shows two windows from the gretl software. The left window, titled 'gretl: specify model', shows the OLS specification process. The dependent variable is 'Q'. The independent variables list includes 'const', 'Y', 'P', 'ED1', and 'ED2'. Red circles highlight the 'ED1' and 'ED2' variables in the list and the 'OK' button at the bottom. The right window, titled 'gretl: model 3', displays the OLS regression results for Model 3, using observations 1960-1988 (T = 29). The dependent variable is 'Q'. The results table is as follows:

	coefficient	std. error	t-ratio	p-value
const	0.707979	0.454836	1.557	0.1327
Y	0.000957520	0.000291661	3.283	0.0031 ***
P	-0.313591	0.104349	-3.005	0.0061 ***
ED1	-5.86691	2.64112	-2.221	0.0360 **
ED2	-3.22065	3.56141	-0.9043	0.3748

Below the table, summary statistics are provided:

Mean dependent var	2.204655	S.D. dependent var	0.243190
Sum squared resid	0.493419	S.E. of regression	0.143385
R-squared	0.702033	Adjusted R-squared	0.652372
F(4, 24)	14.13646	P-value(F)	4.62e-06
Log-likelihood	17.91932	Akaike criterion	-25.83863
Schwarz criterion	-19.00216	Hannan-Quinn	-23.69753
rho	0.501147	Durbin-Watson	0.962300

At the bottom of the results window, it states: 'Excluding the constant, p-value was highest for variable 6 (ED2)'.

Again our hypotheses are confirmed, income (Y) has a positive and significant coefficient, and price (P) has a negative and significant coefficient. The new variables on educational enrollment suggest that a higher proportions of young individuals enrolled in school does REDUCE the demand for cigarettes (they are both negative), although the second one for university enrollment is not significant.

But again, we need to check for problems with heteroskedasticity:

The screenshot shows the 'gretl: LM test (heteroskedasticity)' window. It displays the results of White's test for heteroskedasticity, using observations 1960-1988 (T = 29). The dependent variable is 'uhat^2'. The results table is as follows:

	coefficient	std. error	t-ratio	p-value
const	2.20922	1.37947	1.602	0.1316
Y	-0.00254965	0.00165336	-1.542	0.1453
P	-1.03264	0.672014	-1.537	0.1467
ED1	21.6174	15.8213	1.366	0.1934
ED2	35.4459	22.9572	1.544	0.1449
sq_Y	6.64943e-07	4.86999e-07	1.365	0.1937
X2_X3	0.000761843	0.000422984	1.801	0.0933 *
X2_X4	-0.0105546	0.00887014	-1.190	0.2539
X2_X5	-0.0258571	0.0125935	-2.053	0.0592 *
sq_P	0.120376	0.0584879	2.058	0.0587 *
X3_X4	-8.20766	3.93072	-2.088	0.0555 *
X3_X5	-5.29263	4.23482	-1.250	0.2319
sq_ED1	42.4348	40.7316	1.042	0.3152
X4_X5	251.700	97.2115	2.589	0.0214 **
sq_ED2	106.559	88.1117	1.209	0.2466

Below the table, a warning message is displayed: 'Warning: data matrix close to singularity!'. The Unadjusted R-squared is 0.749925. The test statistic is $TR^2 = 21.747822$, with a p-value = $P(\text{Chi-square}(14) > 21.747822) = 0.083947$. The p-value 0.083947 is circled in red. A red callout bubble points to this p-value with the text: 'This model also has a problem with heteroskedasticity'.

So, now re-run this extended model with “Robust Standard Errors”:

The screenshot shows two windows from the gretl software. The left window is titled "gretl: specify model" and shows the OLS specification interface. The dependent variable is "Q". The independent variables are "const", "Y", "P", "ED1", and "ED2". The "Robust standard errors" checkbox is checked and circled in red. A red arrow points from this checkbox to the right window.

The right window is titled "gretl: model 4" and displays the results for Model 4: OLS, using observations 1960-1988 (T = 29). The dependent variable is Q, and HAC standard errors with a bandwidth of 2 (Bartlett kernel) are used. The results are as follows:

	coefficient	std. error	t-ratio	p-value	
const	0.707979	0.400761	1.767	0.0900	*
Y	0.000957520	0.000268534	3.566	0.0016	***
P	-0.313591	0.117486	-2.669	0.0134	**
ED1	-5.86691	2.76959	-2.118	0.0447	**
ED2	-3.22065	4.05229	-0.7948	0.4345	

Additional statistics shown include: Mean dependent var: 2.204655, S.D. dependent var: 0.243190, Sum squared resid: 0.493419, S.E. of regression: 0.143385, R-squared: 0.702033, Adjusted R-squared: 0.652372, F(4, 24): 28.58676, P-value(F): 8.11e-09, Log-likelihood: 17.91932, Akaike criterion: -25.83863, Schwarz criterion: -19.00216, Hannan-Quinn: -23.69753, rho: 0.501147, Durbin-Watson: 0.962300. A note at the bottom states: "Excluding the constant, p-value was highest for variable 6 (ED2)".

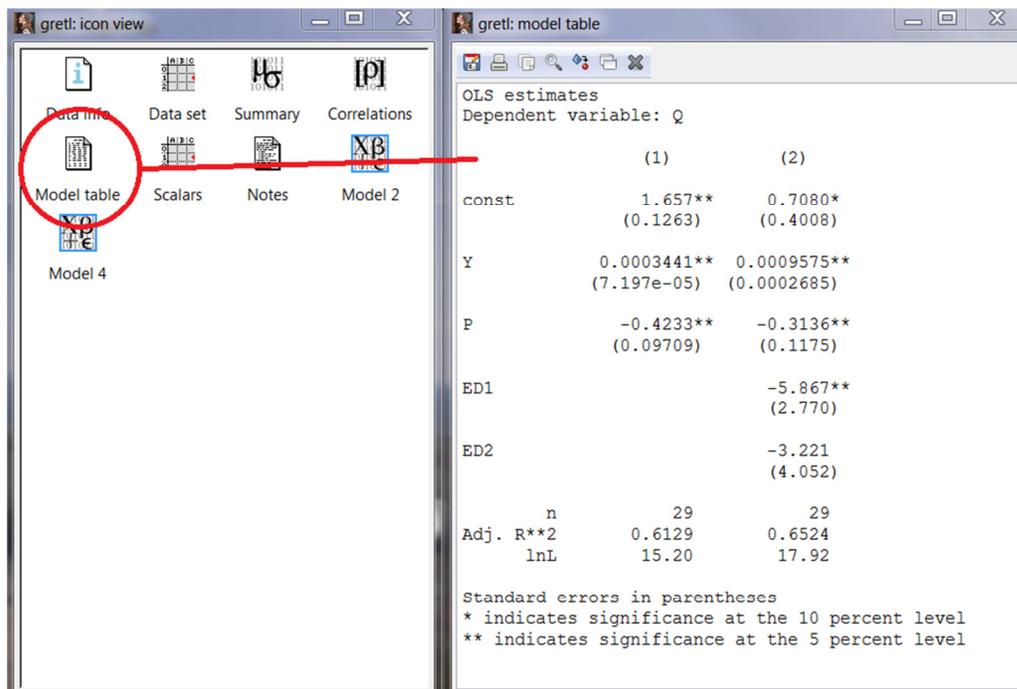
And again, our hypotheses are confirmed, these are just better estimates. So now, let's save this model as an icon as well.

Now, open the icon view by clicking on the button on the lower part of the main screen:

The screenshot shows the gretl main window on the left and the "gretl: icon view" window on the right. The main window displays a list of variables with their descriptive labels. The "Icon view" window shows a grid of icons for different views: Data info, Data set, Summary, Correlations, Model table, Scalars, Notes, and Model 2. The "Model table" and "Model 2" icons are circled in red. A red arrow points from the "Model table" icon in the icon view back to the main window's toolbar, where the "Model table" icon is also circled in red.

Now, let's create a model table that we could put in a paper with these two sets of results. First drag and drop the Model 2 icon onto the Model Table icon. Then next, drag and drop the Model 4 icon onto the Model Table icon.

Now, double click the Model Table icon to open up the table of results we created:



These results are ready to be put into a paper, or you can print them. In the table in your paper, however, you do need to make a note that the standard errors that are reported in parenthesis are Robust Standard Errors. You should also be sure to better label the variables with longer descriptions. I saved it as RTF (Word) format and then opened it and pasted it below and added descriptions and updated the standard error note. I did add borders to the table.

OLS estimates
Dependent variable: Q (Cigarette Consumption Per Adult)

	(1)	(2)
Constant	1.657** (0.1263)	0.7080* (0.4008)
Y (Income)	0.0003441** (7.197e-05)	0.0009575** (0.0002685)
P (Price)	-0.4233** (0.09709)	-0.3136** (0.1175)
ED1 (School Enrollment Middle/High)	--	-5.867** (2.770)
ED2 (School Enrollment University)	--	-3.221 (4.052)
n (Observations)	29	29
Adj. R ²	0.6129	0.6524
lnL	15.2	17.92

Robust standard errors in parentheses
* indicates significance at the 10 percent level
** indicates significance at the 5 percent level

Econometric Analysis – Dr. Sobel

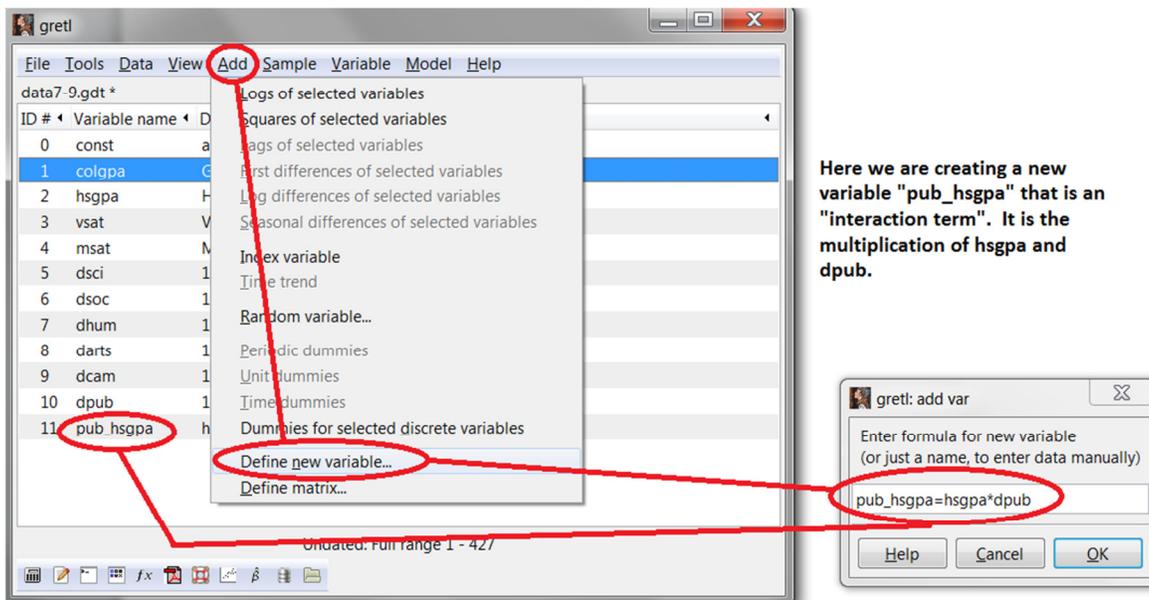
Econometrics Session 3 (time permitting and mostly for reference if your paper needs these adjustments):

5. Advanced Adjustments & Special Cases (time permitting and mostly for reference if your paper needs these adjustments)

Interaction terms – when the effect of one variable depends on another variable (e.g. different for men/women)
- for two variables X and Y, we do this by including both X, Y, and also a new variable X*Y (X times Y)

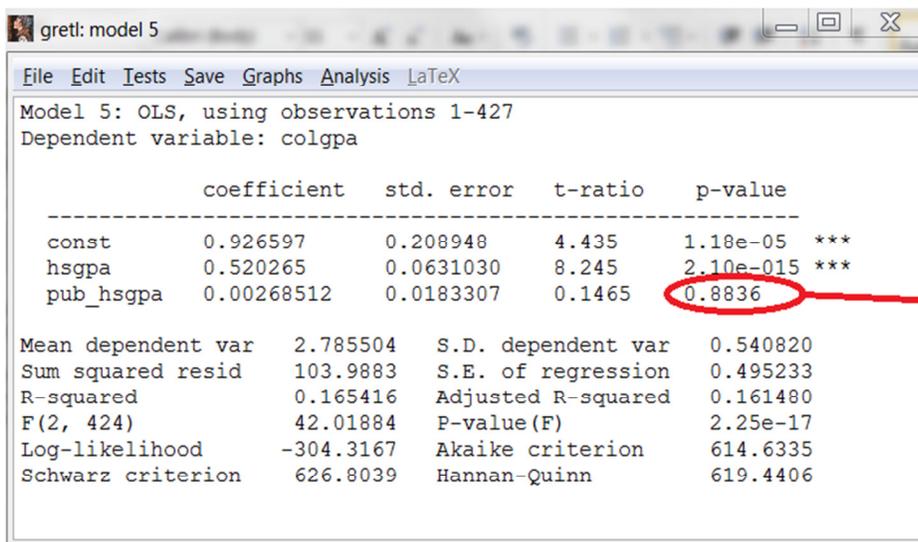
For example, let's ask whether the impact of a person's high school GPA on their college GPA depends on whether they went to a private or public high school. In other words, if you went to a private school your high school GPA might matter more or less than it does if you went to a public high school. We would do this by including both the high school GPA as a variable AND including a new variable that we have to create that is the high school GPA times the zero/one variable indicating they went to a public high school. So let's first create the new variable $\text{pub_hsgpa} = \text{hsgpa} * \text{dpub}$

CREATING A NEW VARIABLE THAT IS AN INTERACTION TERM:



The screenshot shows the gretl software interface. The 'Add' menu is open, and the 'Define new variable...' option is selected. A red circle highlights this option. The 'add var' dialog box is also open, showing the formula $\text{pub_hsgpa} = \text{hsgpa} * \text{dpub}$ entered in the text field. A red circle highlights the formula. A red arrow points from the 'Define new variable...' option in the menu to the dialog box. A text box on the right explains: 'Here we are creating a new variable "pub_hsgpa" that is an "interaction term". It is the multiplication of hsgpa and dpub.'

Now run the regression including this new variable as well:



The screenshot shows the gretl software interface displaying the results for Model 5: OLS, using observations 1-427. The dependent variable is colgpa. The regression results are as follows:

	coefficient	std. error	t-ratio	p-value	
const	0.926597	0.208948	4.435	1.18e-05	***
hsgpa	0.520265	0.0631030	8.245	2.10e-015	***
pub_hsgpa	0.00268512	0.0183307	0.1465	0.8836	

Additional statistics:

Mean dependent var	2.785504	S.D. dependent var	0.540820
Sum squared resid	103.9883	S.E. of regression	0.495233
R-squared	0.165416	Adjusted R-squared	0.161480
F(2, 424)	42.01884	P-value (F)	2.25e-17
Log-likelihood	-304.3167	Akaike criterion	614.6335
Schwarz criterion	626.8039	Hannan-Quinn	619.4406

A red circle highlights the p-value for the pub_hsgpa variable, which is 0.8836.

Because this new variable is insignificant (no stars), it means there is NOT a difference. Had this been significant, there would be a difference in the effect of high school GPA depending on whether you went to a public or private high school.

“Discrete”, “Count”, or Dummy Variables as the dependent variable (probit / logit / tobit)

- When the dependent variable is a zero/one variable we use “Probit” or “Logit” (usually Probit) but you can see which gives the best results. Run using “Model” menu, “Nonlinear models”, “Probit” (or “Logit”), then “Binary”. Note that the coefficients cannot be interpreted the same. Choosing “Show slopes at mean” will give you the “marginal effects” coefficients you can interpret the same (how a one unit change effects the PROBABILITY it is a one, and note these are in decimal form so 0.35 means it increases the probability of a one by 35%). But you will also need to run it the other way selecting “show p-values” to get the significance to show in your table you report as well.

Let’s try to see whether a student’s math SAT score impacts the PROBABILITY they choose to major in science (dsci). We have a dummy/indicator variable dsci that is equal to one if they live on campus. Let’s run a Probit model using this as the dependent variable.

RUNNING A PROBIT MODEL FOR A DEPENDENT VARIABLE THAT IS ZERO/ONE:

Here's how to run a probit model for a variable that is zero/one only (binary). I'm trying to see whether a student's Math SAT score (msat) helps to predict whether they choose to major in science in college (dsci).

There are two options, one gives you the "p-values" (statistical significance) and the other gives you the "slopes at mean" (the "marginal effects") that you generally also need to report. So you may need to run each model twice.

	coefficient	std. error	z	slope
const	-1.22924	0.381213	-3.225	
msat	0.00227212	0.000665825	3.412	0.000904946

	coefficient	std. error	z	p-value
const	-1.22924	0.381213	-3.225	0.0013 ***
msat	0.00227212	0.000665825	3.412	0.0006 ***

The answer is “yes” students with higher math SAT scores are more likely to major in science disciplines (we know this from the stars to the right in the version that came from “Show p-values” to the right above).

The “show slopes at mean” coefficients are the ones to interpret. The coefficient on msat is 0.00227212. So for every 1 point higher on the math SAT they are 0.00227212 (or 0.227212%) more likely to major in science. Let’s multiply those so they make sense. For every 100 points higher on the math SAT, they are roughly 22.7% (or 23%) more likely to major in science. This is really “percentage points” so for example if a student’s odds of majoring in science were 60% normally, a 100 point increase in their math SAT would increase their probability of majoring in science to 83%.

- When the dependent variable is a count (0, 1, 2, 3...) variable we use "Poisson" or "Negative binomial" model. Run using "Model" menu, "Nonlinear models", "Count data" and then you can pick Poisson or one of two negative binomial models. These are harder to interpret and if you need to do this you should see me.
- When the dependent variable is censored (or "truncated") by being cut off at some high or low value we use a "Tobit" model ("Model" menu, "Nonlinear models", "Tobit") but this is unusual. For this you should also see me.

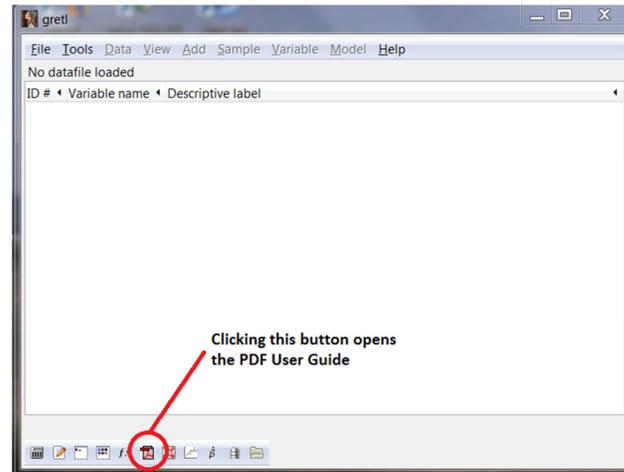
Panel Data Analysis & Fixed Effects – when you have both cross sectional and time series data (e.g., all states for a number of years) you generally include "fixed effects" which are simply dummy variables for each time period and/or for each state. This controls for anything specific to all states in a given year (like a recession) or for anything specific to one state in all years (like Hawaii is just always different). You can create and include these dummy variables yourself (although you must EXCLUDE one of the years or states as that will be in the constant). Alternatively, if you make sure you tell gretl your data is "Panel Data" (go to Data menu, Dataset structure to change or specify) you can then have it make these up for you by going to "Model" menu, selecting "Panel" then "Fixed and random effects".

Special problems/issues with time series data (you would again have to see me for help with these)

- The problem of "nonstationary" data (will give significant results even when not)
 - The EASIEST way to overcome this is to run your model in "Change" form where you use the year to year changes in ALL of the variables for the model rather than the raw data. You can, however, test for this as a problem using "Variable menu" then "Unit Root Tests".
- In time series models we generally want to predict the future based on the past, which is called "autoregressive" or "AR" model, use "Model" then "Time series" then "Autoregressive estimation".
- For some sets of two time series variables we just want to see if they tend to hang together through time, which is called a "Cointegration test" (found under "Model" menu, then "Time series")
- Granger Causality (Chicken Egg paper we will be covering) – can only be done on time series data, then go to "Model", "Time series", "Vector autoregression", here you put BOTH variables as "Endogenous" (dependent variables), you can run your model without any "Exogenous" (independent) variables. Again, this is one you would need to see me for help with.

List of gretl User Manual References & Pages for Important Commands

GRETL USER GUIDE (PDF “users guide” button along bottom):



- 1) How to run a basic OLS regression and details on all the menu commands for OLS regressions is in Chapter 2, Pages 5-12
- 2) Saving sessions, creating a model table, and the icon view is in Chapter 3, Section 3.4 on pages 16-18
- 3) Data files, how to read them in, rules on variable names, types of data is in Chapter 4, on pages 19-23
- 4) Dealing with missing values in the data is in Chapter 4, Section 4.6, on pages 27-28
- 5) Using subsamples of data, setting the sample is in Chapter 6, pages 36-38
- 6) Creating graphs and graph options is in Chapter 7, pages 40-41
- 7) Heteroskedasticity, robust standard errors, White's correction is in Chapter 15, pages 113-114
- 8) Using "panel" data and fixed effects is in Chapter 16, pages 121-125
- 9) Time series models are in Chapters 23-25, on pages 173-210. See Section 24.2 on Page 188 for Granger causality
- 10) Probit, Logit, and Tobit models are in Chapter 19, pages 234-253