# **Linear Regression using Stata**

(v. 6.4)

# **Oscar Torres-Reyna**

otorres@princeton.edu

December 2007

http://www.princeton.edu/~otorres/

We use regression to estimate the <u>unknown effect</u> of changing one variable over another (Stock and Watson, 2003, ch. 4)

When running a regression we are making two assumptions, 1) there is a linear relationship between two variables (i.e. X and Y) and 2) this relationship is additive (i.e.  $Y = x1 + x2 + \dots + xN$ ).

Technically, linear regression estimates how much Y changes when X changes one unit.

In Stata use the command regress, type:

```
regress [dependent variable] [independent variable(s)]
```

reqress y x

In a multivariate setting we type:

regress y x1 x2 x3 ...

Before running a regression it is recommended to have a clear idea of what you are trying to estimate (i.e. which are your outcome and predictor variables).

A regression makes sense only if there is a sound theory behind it.

# Regression: a practical approach (setting)

**Example**: Are SAT scores higher in states that spend more money on education controlling by other factors?\*

- Outcome (Y) variable SAT scores, variable csat in dataset
- Predictor (X) variables
  - Per pupil expenditures primary & secondary (expense)
  - % HS graduates taking SAT (percent)
  - Median household income (income)
  - % adults with HS diploma (high)
  - % adults with college degree (college)
  - Region (region)

\*Source: Data and examples come from the book *Statistics with Stata (updated for version 9)* by Lawrence C. Hamilton (chapter 6). <u>Click here to download the data or search for it at <u>http://www.duxbury.com/highered/</u>. Use the file states.dta (educational data for the U.S.).</u>

## Regression: variables

It is recommended first to examine the variables in the model to check for possible errors, type:

describe csat expense percent income high college region summarize csat expense percent income high college region

. describe csat expense percent income high college region

variable name	storage type	di spl ay format	val ue l abel	variable label
csat	int	%9. 0g	regi on	Mean composite SAT score
expense	int	%9. 0g		Per pupil expenditures prim&sec
percent	byte	%9. 0g		% HS graduates taking SAT
i ncome	double	%10. 0g		Median household income, \$1,000
hi gh	float	%9. 0g		% adults HS diploma
col l ege	float	%9. 0g		% adults college degree
regi on	byte	%9. 0g		Geographical region

. summarize csat expense percent income high college region

Vari abl e	0bs	Mean	Std. Dev.	Mi n	Max
csat	51	944. 098 5235 061	66. 93497 1401 155	832	1093
percent	51	35. 76471	<b>26.</b> 19281	2900 4	9239 81
hi gh	51 51	33. 95657 76. 26078	6. 423134 5. 588741	23. 465 64. 3	48. 618 86. 6
college	51	20. 02157	4. 16578	12.3	33. 3
region	50	<b>Z. 54</b>	1. 128662	1	4

4

# Regression: what to look for



## Regression: what to look for



#### Regression: using dummy variables/selecting the reference category

If using categorical variables in your regression, you need to add n-1 dummy variables. Here 'n' is the number of categories in the variable. In the example below, variable 'industry' has twelve categories (type tab industry, or tab industry, nolabel)

To change the reference category to "Professional services" The easiest way to include a set of dummies in a regression is by using the prefix "i." By default, the first category (or lowest value) is (category number 11) instead of "Ag/Forestry/Fisheries" (category used as reference. For example: number 1), use the prefix "ib#." where "#" is the number of the reference category you want to use; in this case is 11. sysuse nlsw88.dta req wage hours i.industry, robust sysuse nlsw88.dta reg wage hours ib11.industry, robust Number of obs = 2228 Linear regression Linear regression Number of obs = 2228 F(12, 2215) = 24.96F(12, 2215) = 24.96Prob > F = 0.0000 Prob > F = 0.0000 R-squared = 0.0800 R-squared = 0.0800 Root MSE = 5.5454 = 5.5454 Root MSE Robust Robust Std. Err. P>|t| [95% Conf. Interval] wage Coef t Coef Std. Err. t P>|t| [95% Conf. Interval] wage .0723658 .0110213 0.000 .0507526 .093979 hours 6.57 .0723658 .0110213 6.57 0.000 .0507526 .093979 hours industrv industry Mining 9.328331 7.287849 1.28 0.201 -4.963399 23.62006 Ag/Forestry/Fisheries -2.094988 .8192781 -2.56 0.011 -3.701622-.4883548 4.371759 Construction 1.858089 1.281807 1.45 0.147 -.6555808 0.318 Mining 7.233343 7.245913 1 00 -6.97615 21 44284 Manufacturing 1.415641 .849571 1.67 0.096 -.2503983 3.081679 -.2368991 1.011309 -0.23 0.815 -2.220112 1.746314 Construction 5.432544 1.03998 5.22 0.000 3.393107 7.471981 Transport/Comm/Utility -.6793477 .3362365 -2.02 0.043 -1.338719 -.019976 Manufacturing .8548564 -1.218023 Wholesale/Retail Trade .4583809 0.54 0.592 2.134785 3.337556 .6861828 0.000 1.991927 4.683185 Transport/Comm/Utility 4.86 Finance/Ins/Real Estate 3.92933 .9934195 3.96 0.000 1.981199 5.877461 0.000 Wholesale/Retail Trade -1.636607 .3504059 -4.67 -2.323766 -.949449 1.990151 1.054457 Business/Repair Svc 1.89 0.059 -.0776775 4.057979 Finance/Ins/Real Estate 1.834342 .6171526 2.97 0.003 .6240837 3.0446 Personal Services -1.018771 .8439617 -1.21 0.228 -2.67381 .6362679 -.1048377 .7094241 -0.15 0.883 -1.496044 1.286368 Business/Repair Svc 1.111801 1.192314 0.93 0.351 -1.226369 3.449972 -3.113759 .3192289 -9.75 0.000 -3.739779 -2.48774 Entertainment/Rec Svc Personal Services Professional Services 2.094988 .8192781 2.56 0.011 .4883548 3.701622 -.983187 .9004471 0.275 -2.748996 .7826217 Entertainment/Rec Svc -1.09 Public Administration 3.232405 .8857298 3.65 0.000 1.495457 4.969352 Public Administration 1.137416 .4176899 2.72 0.007 .3183117 1.956521 cons 5.221617 .4119032 12.68 0.000 4.41386 6.029374 3.126629 .8899074 3.51 0.000 1.381489 4.871769 cons

The "ib#." option is available since Stata 11 (type help for arrlist for more options/details). For older Stata versions you need to use "xi:" along with "i." (type help xi for more options/details). For the examples above type (output omitted):

xi: req wage hours i.industry, robust

char industry[omit]11 /\*Using category 11 as reference\*/

xi: reg wage hours i.industry, robust

#### To create dummies as variables type

tab industry, gen(industry)

To include all categories by suppressing the constant type: reg wage hours bn.industry, robust hascons

## Regression: ANOVA table

If you run the regression without the 'robust' option you get the ANOVA table xi: regress csat expense percent income high college i.region

Source	SS	df	MS	Number of $obs = E(0, 0, 0, 0)$	50 70 12
 Model	(A) 200269. 84	9	22252. 2045 (D)	P(9, 40) = Prob > F =	0. 0000
 Resi dual	(B)12691.5396	40	317. 28849 (E)	R-squared = Adi R-squared =	0.9404
Total	<mark>(C)</mark> 212961. 38	49	4346. 15061 (F)	Root MSE =	17.813

$$F = \frac{\frac{MSS}{(k-1)}}{\frac{RSS}{n-k}} = \frac{\frac{200269.84}{9}}{\frac{12691.5396}{40}} = \frac{22252.2045}{317.28849} = \frac{D}{E} = 70.13 \qquad AdjR^2 = 1 - \frac{n-1}{n-k}(1-R^2) = 1 - \frac{49}{40}(1-0.9404) = 1 - \frac{E}{F} = 1 - \frac{317.28849}{4346.15061} = 0.9270$$

$$R^{2} = \frac{MSS}{TSS} = 1 - \frac{\sum e_{i}^{2}}{\sum (y_{i} - \overline{y})^{2}} = \frac{200269.84}{212961.38} = \frac{A}{C} = 0.9404 \qquad RootMSE = \sqrt{\frac{RSS}{(n-k)}} = \sqrt{\frac{12691.5396}{40}} = \sqrt{\frac{B}{40}} = 17.813$$

A = Model Sum of Squares (MSS). The closer to TSS the better fit.

B = Residual Sum of Squares (RSS)

C = Total Sum of Squares (TSS)

**D** = Average Model Sum of Squares = MSS/(k-1) where k = # predictors

**E** = Average Residual Sum of Squares = RSS/(n - k) where n = # of observations

**F** = Average Total Sum of Squares = TSS/(n-1)

 $R^2$  shows the amount of observed variance explained by the model, in this case 94%.

The *F*-statistic, F(9,40), tests whether  $R^2$  is different from zero.

Root MSE shows the average distance of the estimator from the mean, in this case 18 points in estimating SAT scores.

Source: Kohler, Ulrich, Frauke Kreuter, Data Analysis Using Stata, 2009

#### Regression: estto/esttab

To show the models side-by-side you can use the commands estto and esttab:

```
regress csat expense, robust
eststo model1
regress csat expense percent income high college, robust
eststo model2
xi: regress csat expense percent income high college i.region, robust
eststo model3
                                       . esttab, r2 ar2 se scalar(rmse)
esttab, r2 ar2 se scalar(rmse)
```

	(1)	(2)	(3)
	csat	csat	csat
expense	- 0. 0223***	0. 00335	- 0. 00202
1	( <b>0.00367</b> )	(0.00478)	(0.00359)
percent		- 2. 618***	- 3. 008***
L		( <b>0. 229</b> )	( <b>0. 236</b> )
income		0. 106	- 0. 167
		(1.207)	(1. 196)
hi ơh		1, 631	1, 815
8		(0.943)	(1.027)
college		2, 031	4.671**
8-		(2.114)	(1.600)
Iregion 2			69. 45***
			(18.00)
Iregion 3			25.40*
_iregron_o			(12.53)
Iregion 4			34 58***
_iregron_4			(9.450)
percent2			
cons	1060. 7***	851.6***	808. 0***
	(24.35)	(57.29)	(67.86)
N	51	51	50
R-sq	0. 217	0.824	0. 911
adj. R-sq	0. 201	0.805	0.894
rmse	59.81	29.57	21.49

Type help eststo and help esttab for more options.

Standard errors in parentheses \* p<0.05, \*\* p<0.01, \*\*\* p<0.001

# Regression: correlation matrix

Below is a correlation matrix for all variables in the model. Numbers are Pearson correlation coefficients, go from -1 to 1. Closer to 1 means strong correlation. A negative value indicates an inverse relationship (roughly, when one goes up the other goes down).

. pwcorr csat expense percent income high college, star(0.05) sig

	csat	expense	percent	income	hi gh	college
csat	1. 0000					
expense	- 0. 4663* 0. 0006	1.0000				
percent	- 0. 8758* 0. 0000	0. 6509* 0. 0000	1.0000			
income	- 0. 4713* 0. 0005	0. 6784* 0. 0000	0. 6733* 0. 0000	1.0000		
hi gh	0. 0858 0. 5495	0. 3133* 0. 0252	0. 1413 0. 3226	0. 5099* 0. 0001	1.0000	
col l ege	- 0. 3729* 0. 0070	0. 6400* 0. 0000	0. 6091* 0. 0000	0. 7234* 0. 0000	0. 5319* 0. 0001	1. 0000

#### Regression: graph matrix

Command graph matrix produces a graphical representation of the correlation matrix by presenting a series of scatterplots for all variables. Type:

graph matrix csat expense percent income high college, half
maxis(ylabel(none) xlabel(none))



## Regression: exploring relationships



scatter csat percent

There seem to be a curvilinear relationship between csat and percent, and slightly linear between csat and high. To deal with U-shaped curves we need to add a square version of the variable, in this case percent square

generate percent2 = percent^2

## Regression: functional form/linearity

The command acprplot (augmented component-plus-residual plot) provides another graphical way to examine the relationship between variables. It does provide a good testing for linearity. Run this command after running a regression

regress csat percent high /\* Notice we do not include percent2 \*/ acprplot percent, lowess acprplot high, lowess



The option lowess (locally weighted scatterplot smoothing) draw the observed pattern in the data to help identify nonlinearities. Percent shows a quadratic relation, it makes sense to add a square version of it. High shows a polynomial pattern as well but goes around the regression line (except on the right). We could keep it as is for now.

#### The model is:

xi: regress csat expense percent percent2 income high college i.region, robust

Form more details see <a href="http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm">http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</a>, and/or type help acprplot and help lowess. 13

### **Regression: models**

xi: regress csat expense percent percent2 income high college i.region, robust eststo model4

esttab, r2 ar2 se scalar(rmse)

. esttab, r2 ar2 se scalar(rmse)

	(1)	(2)	(3)	(4)
	csat	csat	csat	csat
expense	-0.0223*** (0.00367)	0. 00335 (0. 00478)	-0.00202 (0.00359)	0. 00141 (0. 00372)
percent		-2.618*** (0.229)	-3.008*** (0.236)	- 5. 945*** (0. 641)
income		0. 106 (1. 207)	-0.167 (1.196)	-0.914 (0.973)
hi gh		1.631 (0.943)	1.815 (1.027)	1.869 (0.931)
col l ege		<b>2.031</b> ( <b>2.114</b> )	<b>4. 671**</b> ( <b>1. 600</b> )	<b>3. 418**</b> (1. 145)
_I regi on_2			69. 45*** (18. 00)	5. 077 (20. 75)
_I regi on_3			<b>25. 40*</b> (12. 53)	5. 209 (10. 42)
_I regi on_4			34. 58*** (9. 450)	19. 25* (8. 110)
percent2				0. 0460*** (0. 0102)
_cons	1060. 7*** (24. 35)	851. 6*** (57. 29)	808. 0*** (67. 86)	874. 0*** (58. 13)
N R-sq adj. R-sq rmse	51 0. 217 0. 201 59. 81	51 0. 824 0. 805 29. 57	50 0. 911 0. 894 21. 49	50 0. 940 0. 927 17. 81

Standard errors in parentheses \* p<0.05, \*\* p<0.01, \*\*\* p<0.001

#### Regression: getting predicted values

How good the model is will depend on how well it predicts *Y*, the linearity of the model and the behavior of the residuals.

There are two ways to generate the *predicted values of Y* (usually called *Yhat*) given the model:

<u>Option A</u>, using generate after running the regression:

```
xi: regress csat expense percent percent2 income high college i.region, robust
generate csat_predict = _b[_cons] + _b[percent]*percent + _b[percent]*percent + _b[percent2]*percent2
+ _b[high]*high + ...
```

Option B, using predict immediately after running the regression:

xi: regress csat expense percent percent2 income high college i.region, robust
predict csat\_predict
label variable csat\_predict "csat predicted"

. predict csat\_predict
(option xb assumed; fitted values)
(1 missing value generated)

. label variable csat\_predict "csat predicted"

#### Regression: observed vs. predicted values

For a quick assessment of the model run a scatter plot



scatter csat csat\_predict

We should expect a 45 degree pattern in the data. Y-axis is the observed data and x-axis the predicted data (*Yhat*).

In this case the model seems to be doing a good job in predicting csat

#### Regression: testing for homoskedasticity

An important assumption is that the variance in the residuals has to be homoskedastic or constant. Residuals cannot varied for lower of higher values of X (i.e. fitted values of Y since Y=Xb). A definition:

"The error term [e] is homoskedastic if the variance of the conditional distribution of  $[e_i]$  given  $X_i$  [var( $e_i|X_i$ )], is constant for i=1...n, and in particular does not depend on x; otherwise, the error term is heteroskedastic" (Stock and Watson, 2003, p.126)

When plotting residuals vs. predicted values (*Yha*t) we *should not observe* any pattern at all. In Stata we do this using rvfplot <u>right after</u> running the regression, it will automatically draw a scatterplot between residuals and predicted values.



#### Regression: testing for homoskedasticity

A non-graphical way to detect heteroskedasticiy is the Breusch-Pagan test. The null hypothesis is that residuals are homoskedastic. In the example below we fail to reject the null at 95% and concluded that residuals are homogeneous. However at 90% we reject the null and conclude that residuals are not homogeneous.

#### estat hettest

. estat hettest Breusch-Pagan / Cook-Weisberg test for heteroskedasticity Ho: Constant variance Variables: fitted values of csat chi 2(1) = 2.72Prob > chi 2 = 0.0993

The graphical and the Breush-Pagan test suggest the possible presence of heteroskedasticity in our model. The problem with this is that we may have the wrong estimates of the standard errors for the coefficients and therefore their t-values.

There are two ways to deal with this problem, one is using heteroskedasticity-robust standard errors, the other one is using weighted least squares (see Stock and Watson, 2003, chapter 15). WLS requires knowledge of the conditional variance on which the weights are based, if this is known (rarely the case) then use WLS. In practice it is recommended to use heteroskedasticity-robust standard errors to deal with heteroskedasticity.

By default Stata assumes homoskedastic standard errors, so we need to adjust our model to account for heteroskedasticity. To do this we use the option robust in the regress command.

xi: regress csat expense percent percent2 income high college i.region, robust

Following Stock and Watson, as a rule-of-thumb, you should always assume heteroskedasticiy in your model (see Stock and Watson, 2003, chapter 4).

### Regression: omitted-variable test

How do we know we have included all variables we need to explain Y?

Testing for omitted variable bias is important for our model since it is related to the assumption that the error term and the independent variables in the model are not correlated (E(e|X) = 0)

If we are missing variables in our model and

- "is correlated with the included regressor" and,
- "the omitted variable is a determinant of the dependent variable" (Stock and Watson, 2003, p.144),

...then our regression coefficients are inconsistent.

In Stata we test for omitted-variable bias using the ovtest command:

xi: regress csat expense percent percent2 income high college i.region, robust ovtest

. ovtest

Ramsey RESET test using powers of the fitted values of csat Ho: model has no omitted variables F(3, 37) = 1.25Prob > F = 0.3068

The null hypothesis is that the model does not have omitted-variables bias, the p-value is higher than the usual threshold of 0.05 (95% significance), so we fail to reject the null and conclude that we do not need more variables.

Another command to test model specification is linktest. It basically checks whether we need more variables in our model by running a new regression with the observed Y (csat) against Yhat (csat predicted or  $X\beta$ ) and Yhat-squared as independent variables<sup>1</sup>.

The thing to look for here is the significance of <u>hatsq</u>. The null hypothesis is that there is no specification error. If the p-value of \_hatsq is not significant then we fail to reject the null and conclude that our model is correctly specified. Type:

xi: regress csat expense percent percent2 income high college i.region, robust linktest

1	i	nkt	est

Source	SS	df	MS	Number of obs	= 50 = 370 90
Model Residual	200272. 359 12689. 0209	2 100 47 269.	136. 18 979169	Prob > F R-squared Adi P squared	$\begin{array}{rcrcr} - & 370.90 \\ = & 0.0000 \\ = & 0.9404 \\ - & 0.9379 \end{array}$
Total	212961. 38	49 4346	. 15061	Root MSE	= 16.431
csat	Coef.	Std. Err.	t P> t	[95% Conf.	Interval]
_hat _hatsq _cons	1. 144949 0000761 - 68. 69417	1. 50184 . 0007885 712. 388	0. 76 0. 45 -0. 10 0. 92 -0. 10 0. 92	0 - 1. 876362 3 0016623 4 - 1501. 834	4. 166261 . 0015101 1364. 446

<sup>1</sup> For more details see <a href="http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm">http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</a>, and/or type <a href="http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm">http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</a>, and/or type <a href="http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm">http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</a>, and/or type <a href="http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm">http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</a>, and/or type <a href="https://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm">https://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</a>, and/or type <a href="https://www.ats.ucla.edu/statareg2.htm">https://www.ats.ucla.edu/stat/statareg2.htm</a>, and a https://www.ats.ucla.edu/statareg2.htm</a>

# Regression: multicollinearity

An important assumption for the multiple regression model is that independent variables are *not perfectly multicolinear*. One regressor should not be a linear function of another.

When multicollinearity is present *standand errors may be inflated*. Stata will drop one of the variables to avoid a division by zero in the OLS procedure (see Stock and Watson, 2003, chapter 5).

The Stata command to check for multicollinearity is vif (variance inflation factor). Right after running the regression type:

Vari abl e	VI F	1/VIF
percent2	70. 80	0. 014124
percent	49. 52	0. 020193
_I regi on_2	8.47	0. 118063
income	4.97	0. 201326
_I regi on_3	4.89	0. 204445
hi gh	4.71	0. 212134
college	4.52	0. 221348
expense	3. 33	0. 300111
_I regi on_4	2.14	0. 467506
Mean VIF	17.04	

vi	f
V 1	1

Avif > 10 or a 1/vif < 0.10 indicates trouble.

We know that percent and percent2 are related since one is the square of the other. They are ok since percent has a quadratic relationship with *Y*, but this would be an example of multicolinearity.

The rest of the variables look ok.

## **Regression: outliers**

To check for outliers we use the avplots command (added-variable plots). Outliers are data points with extreme values that could have a negative effect on our estimators. After running the regression type:

avplot expense





These plots regress each variable against all others, notice the coefficients on each. All data points seem to be in range, no outliers observed.

For more details and tests on this and <u>influential</u> and <u>leverage</u> variables please check <u>http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</u>

Also type help diagplots in the Stata command window.

2000

1000

### **Regression: outliers**

#### avplots



# Regression: summary of *influence* indicators

DfBeta	Measures the influence of each observation on the <i>coefficient</i> of a particular independent variable (for example, x1). This is in standard errors terms. An observation is influential if it has a significant effect on the coefficient.	A case is an influential outlier if  DfBeta > 2/SQRT(N) Where N is the sample size. Note: Stata estimates standardized DfBetas.	<pre>In Stata type: reg y x1 x2 x3 dfbeta x1 Note: you could also type: predict DFx1, dfbeta(x1) To estimate the dfbetas for all predictors just type: dfbeta To flag the cutoff gen cutoffdfbeta = abs(DFx1) &gt; 2/sqrt(e(N)) &amp; e(sample)</pre>	In SPSS: Analyze-Regression- Linear; click Save. Select under "Influence Statistics" to add as a new variable (DFB1_1) or in syntax type REGRESSION /MISSING LISTWISE /STATISTICS COEFF OUTS R ANOVA /CRITERIA=PIN(.05) POUT(.10) /NOORIGIN /DEPENDENT Y /METHOD=ENTER X1 X2 X3 /CASEWISE PLOT(ZRESID) OUTLIERS(3) DEFAULTS DFBETA /SAVE MAHAL COOK LEVER DFBETA SDBETA DFFIT SDFIT COVRATIO .
DfFit	Indicator of leverage and high residuals. Measures how much an observation influences the regression model as a whole. How much the predicted values change as a result of including and excluding a particular observation.	High influence if  DfFIT  >2*SQRT(k/N) Where k is the number of parameters (including the intercept) and N is the sample size.	After running the regression type: predict dfits if e(sample), dfits To generate the flag for the cutoff type: gen cutoffdfit= abs(dfits)>2*sqrt((e(df_m)) +1)/e(N)) & e(sample)	Same as DfBeta above (DFF_1)
Covariance ratio	Measures the impact of an observation on the standard errors	High impact if $ COVRATIO-1  \ge 3*k/N$ Where k is the number of parameters (including the intercept) and N is the sample size.	In Stata after running the regression type predict covratio if e(sample), covratio	Same as DfBeta above (COV_1) 24 PU/DSS/OTR

# Regression: summary of *distance* measures

Cook's distance	Measures how much an observation influences the overall model or predicted values. It is a summary measure of leverage and high residuals.	High influence if D > 4/N Where N is the sample size. A D>1 indicates big outlier problem	In Stata after running the regression type: predict D, cooksd	In SPSS: Analyze-Regression-Linear; click Save. Select under "Distances" to add as a new variable (COO_1) or in syntax type REGRESSION /MISSING LISTWISE /STATISTICS COEFF OUTS R ANOVA /CRITERIA=PIN(.05) POUT(.10) /NOORIGIN /DEPENDENT Y /METHOD=ENTER X1 X2 X3 /CASEWISE PLOT(ZRESID) OUTLIERS(3) DEFAULTS DFBETA /SAVE MAHAL COOK LEVER DFBETA SDBETA DFFIT SDFIT COVRATIO.
Leverage	Measures how much an observation influences regression coefficients.	High influence if leverage h > 2*k/N Where k is the number of parameters (including the intercept) and N is the sample size. A rule-of-thumb: Leverage goes from 0 to 1. A value closer to 1 or over 0.5 may indicate problems.	In Stata after running the regression type: predict lev, leverage	Same as above (LEV_1)
Mahalanobis distance	It is rescaled measure of leverage. M = leverage*(N-1) Where N is sample size.	Higher levels indicate higher distance from average values. The M-distance follows a Chi- square distribution with k-1 df and alpha=0.001 (where k is the number of independent variables). Any value over this Chi-square value may indicate problems	Not available	Same as above (MAH_1) 25 PU/DSS/OTR

Sources for the summary tables: influence indicators and distance measures

• Statnotes:

http://faculty.chass.ncsu.edu/garson/PA765/regress.htm#outlier2

- An Introduction to Econometrics Using Stata/Christopher F. Baum, Stata Press, 2006
- Statistics with Stata (updated for version 9) / Lawrence Hamilton, Thomson Books/Cole, 2006
- UCLA <u>http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm</u>

## Regression: testing for normality

Another assumption of the regression model (OLS) that impact the validity of all tests (p, t and F) is that residuals behave 'normal'. Residuals (here indicated by the letter "e") are the difference between the observed values (Y) and the predicted values (Yhat): e = Y - Yhat.

In Stata you type: predict e, resid. It will generate a variable called "e" (residuals).

Three graphs will help us check for normality in the residuals: kdensity, pnorm and qnorm.

#### kdensity e, normal



If residuals do not follow a 'normal' pattern then you should check for omitted variables, model specification, linearity, functional forms. In sum, you may need to reassess your model/theory. In practice normality does not represent much of a problem when dealing with really big samples. A kernel density plot produces a kind of histogram for the residuals, the option normal overlays a normal distribution to compare. Here residuals seem to follow a normal distribution. Below is an example using histogram.





### **Regression: testing for normality**

Standardize normal probability plot (pnorm) checks for non-normality in the middle range of residuals. Again, slightly off the line but looks ok.

Quintile-normal plots (qnorm) check for non-normality in the extremes of the data (tails). It plots quintiles of residuals vs quintiles of a normal distribution. Tails are a bit off the normal.



A non-graphical test is the Shapiro-Wilk test for normality. It tests the hypothesis that the distribution is normal, in this case the null hypothesis is that the distribution of the residuals is normal. Type

swilk e

swilk e

е	50	0. 95566	2. 085	1.567	0. 05855
Vari abl e	0bs	W	V	Ζ	Prob>z
	Shaj	piro-Wilk W	test for	normal data	

The null hypothesis is that the distribution of the residuals is normal, here the p-value is 0.06 we failed to reject the null (at 95%). We conclude then that residuals are normally distributed, with the caveat that they are not at 90%.

# Regression: joint test (F-test)

To test whether two coefficients are jointly different from 0 use the command test (see Hamilton, 2006, p.175).

xi: quietly regress csat expense percent percent2 income high college i.region, robust Note `quietly' suppress the regression output

To test the null hypothesis that both coefficients do not have any effect on csat ( $\beta_{high} = 0$  and  $\beta_{college} = 0$ ), type:

test high college

. test high college

(1) (2)	high = colleg	e = 0	
	F( 2,	40) =	17. 12
	P	rob > F =	0. 0000

The p-value is 0.0000, we reject the null and conclude that *both* variables have indeed a significant effect on SAT.

```
Some other possible tests are (see Hamilton, 2006, p.176):
test income = 1
test high = college
test income = (high + college)/100
```

#### Regression: saving regression coefficients

Stata temporarily stores the coefficients as \_b[varname], so if you type:

```
gen percent_b = _b[percent]
gen constant_b = _b[_cons]
```

You can also save the standard errors of the variables \_se[varname]

```
gen percent_se = _se[percent]
gen constant_se = _se[_cons]
```

summarize percent_	b percent_se cor	istant_b constant_se
--------------------	------------------	----------------------

Vari abl e	0bs	Mean	Std. Dev.	Mi n	Max
percent_b	51	- 5. 945267	0	- 5. 945267	- 5. 945267
percent_se	51	. 6405529	0	. 6405529	. 6405529
constant_b	51	873. 9537	0	873. 9537	873. 9537
constant_se	51	58. 12895	0	58. 12895	58. 12895

#### Regression: saving regression coefficients/getting predicted values

#### You can see a list of stored results by typing after the regression ereturn list:

. xi: quietly regress csat expense percent percent2 income high college i.region, robust \_I region\_1-4 (naturally coded; \_Iregion\_1 omitted) i.region

. ereturn list

scalars:

e(N)	=	50
e(df_m)	=	9
e(df_r)	=	40
e(F)	=	76. 92400040408057
e(r2)	=	. 9404045015031877
e(rmse)	=	17.81259357987284
e(mss)	=	200269. 8403983309
e(rss)	=	12691. 53960166909
$e(r2_a)$	=	. 9269955143414049
e(11)	=	- 209. 3636234584767
e(ll_0)	=	- 279. 8680043669825

macros:

e(cmdline)	:	"regress csat expense percent percent2 income high college _Iregion_*, robust"
e(title)	:	"Linear regression"
e(vce)	:	"robust"
e(depvar)	:	"csat"
e(cmd)	:	"regress"
e(properties)	:	"b V
e(predict)	:	"regres_p"
e(model)	:	"ol s"
e(estat_cmd)	:	"regress_estat"
e(vcetype)	:	"Robust"

matrices:

e(b) : 1 x 10 e(V) : 10 x 10

functions:

e(sample)

# Regression: general guidelines

The following are general guidelines for building a regression model\*

- 1. Make sure all relevant predictors are included. These are based on your research question, theory and knowledge on the topic.
- 2. Combine those predictors that tend to measure the same thing (i.e. as an index).
- 3. Consider the possibility of adding interactions (mainly for those variables with large effects)
- 4. Strategy to keep or drop variables:
  - 1. Predictor not significant and has the expected sign -> Keep it
  - 2. Predictor not significant and does not have the expected sign -> Drop it
  - 3. Predictor is significant and has the expected sign -> Keep it
  - 4. Predictor is significant but does not have the expected sign -> Review, you may need more variables, it may be interacting with another variable in the model or there may be an error in the data.

# Regression: publishing regression output (outreg2)

The command outreg2 gives you the type of presentation you see in academic papers. It is important to notice that outreg2 is not a Stata command, it is a user-written procedure, and you need to install it by typing (only the first time)

#### ssc install outreg2

#### (1)Follow this example (letters in italics you type) VARIABLES Model 1 4.950e+08 x1use "H:\public html\Stata\Panel101.dta", clear (6.902e+08)1.524e+09\*\* req y x1, r Constant (6.636e+08)outreg2 using myreg.doc, replace ctitle(Model 1) Observations 70 . outreg2 using myreq.doc, replace ctitle (Model 1) R-squared 0.006 myreg.doc 📢 Windows users click here to open the file myreg.doc in Word (you Robust standard errors in parentheses dir : seeout can replace this name with your own). Otherwise follow the Mac \*\*\* p<0.01, \*\* p<0.05, \* p<0.1 instructions. Mac users click here to go to the directory where myreq.doc is saved, open it with Word (you can replace this name with your own) (1)(2)You can add other model (using variable x2) by using the option append VARIABLES Model 1 Model 2

x1

 $x^2$ 

Constant

Observations

R-squared

(NOTE: make sure to close myreg.doc)

```
reg y x1 x2, r
outreg2 using myreg.doc, append ctitle(Model 2)
```

```
. outreg2 using myreg.doc, append ctitle(Model 2)
myreg.doc
dir : seeout
```

You also have the option to export to Excel, just use the extension \*.xls.

For older versions of outreg2, you may need to specify the option word or excel (after comma)

33

5.513e+08 (6.869e+08)

3.808e+07

(2.478e+08)

1.483e+09\*\*

(6.595e+08)

70

0.006

4.950e+08

(6.902e+08)

1.524e+09\*\*

(6.636e+08)

70

0.006

Robust standard errors in parentheses

\*\*\* p<0.01, \*\* p<0.05, \* p<0.1

# **Regression:** publishing regression output (outreg2)

You can use outreg2 for almost any regression output (linear or no linear). In the case of logit models with odds ratios, you need to add the option eform, see below

```
use "H:\public html\Stata\Panel101.dta", clear
logit y bin x1
outreq2 using mymod.doc, replace ctitle(Logit coeff)
           . outreg2 using mymod.doc, replace ctitle(Logit coeff)
           mymod.doc
           dir : seeout
logit y bin x1, or
outreg2 using mymod.doc, append ctitle(Odds ratio) eform
          . outreg2 using mymod.doc, append ctitle(Odds ratio) eform
          mymod.doc
                            Windows users click here to open the file mymod.doc in Word (you
          dir : seeout
                            can replace this name with your own). Otherwise follow the Mac
                            instructions.
```

Mac users click here to go to the directory where mymod.doc is saved, open it with Word (you can replace this name with your own)

	EQUATION	VARIABLES	(1) Logit coeff	(2) Odds ratio
	y_bin	x1	0.493	1.637
		Constant	(0.042) 1.082** (0.482)	2.952** (1.422)
For more details/options type		Observations	70	70
help outreg2 Standar *** p<0			s in parentheses p<0.05, * p<0.1	

# Regression: publishing regression output (outreg2)

For predicted probabilities and marginal effects, see the following document :

https://www.princeton.edu/~otorres/Margins.pdf

## Regression: interaction between dummies

Interaction terms are needed whenever there is reason to believe that the effect of one independent variable depends on the value of another independent variable. We will explore here the interaction between two dummy (binary) variables. In the example below there could be the case that the effect of student-teacher ratio on test scores may depend on the percent of English learners in the district\*.

- Dependent variable (Y) Average test score, variable testscr in dataset.
- Independent variables (X)
  - Binary hi\_str, where '0' if student-teacher ratio (str) is lower than 20, '1' equal to 20 or higher.
    - In Stata, first generate hi\_str = 0 if str<20. Then replace hi\_str=1 if str>=20.
  - Binary hi\_el, where '0' if English learners (el\_pct) is lower than 10%, '1' equal to 10% or higher
    - In Stata, first generate hi\_el = 0 if el\_pct<10. Then replace hi\_el=1 if el\_pct>=10.
  - Interaction term str\_el = hi\_str \* hi\_el. In Stata: generate str\_el = hi\_str\*hi\_el

#### We run the regression

```
regress testscr hi_el hi_str str_el, robust
```

. regress te	stscr hi_ell	hi_str_str_el,	robust			
Linear regres:	sion				Number of obs F( 3, 416) Prob > F R-squared Root MSE	= 420 = 60.20 = 0.0000 = 0.2956 = 16.049
testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf.	Interval]
hi_el hi_str str_el _cons	-18.16295 -1.907842 -3.494335 664.1433	2.345952 1.932215 3.121226 1.388089 4	-7.74 -0.99 -1.12 78.46	0.000 0.324 0.264 0.000	-22.77435 -5.705964 -9.629677 661.4147	-13.55155 1.890279 2.641006 666.8718

The equation is testscr\_hat = 664.1 - 18.1\*hi\_el - 1.9\*hi\_str - 3.5\*str\_el

The effect of hi\_str on the tests scores is -1.9 but given the interaction term (and assuming all coefficients are significant), the net effect is -1.9 -3.5\*hi\_el. If hi\_el is 0 then the effect is -1.9 (which is hi\_str coefficient), but if hi\_el is 1 then the effect is -1.9 -3.5 = -5.4. In this case, the effect of student-teacher ratio is more negative in districts where the percent of English learners is higher.

See the next slide for more detailed computations.

## Regression: interaction between dummies (cont.)

You can compute the expected values of test scores given different values of hi\_str and hi\_el. To see the effect of hi\_str given hi\_el type the following right after running the regression in the previous slide.

```
. predict yhat1 if hi_str==0 & hi_el==0
(option xb assumed; fitted values) h
(271 missing values generated)
. predict yhat2 if hi_str==1 & hi_el==0
(option xb assumed; fitted values)
(341 missing values generated)
. predict yhat3 if hi_str==0 & hi_el==1
(option xb assumed; fitted values)
(331 missing values generated)
. predict yhat4 if hi_str==1 & hi_el==1
(option xb assumed; fitted values)
(317 missing values generated)
```

These are different scenarios holding constant hi\_el and varying hi\_str. Below we add some labels

- . label variable yhat1 "Low str/Low el"
- . label variable yhat2 "High str/Low el"
- . label variable yhat3 "Low str/High el"
- . label variable yhat4 "High str/High el"

We then obtain the average of the estimations for the test scores (for all four scenarios, notice same values for all cases).

. s	ummarize yh	at1 yhat2 yh	hat3 yhat4			
	variable	obs	Mean	Std. Dev.	Min	Max
	yhat1 yhat2 yhat3 yhat4	149 79 89 103	664.1433 662.2355 645.9803 640.5782	0 0 0	664.1433 662.2355 645.9803 640.5782	664.1433 662.2355 645.9803 640.5782

```
. display 664.1 - 662.2
1.9
. display 645.9 - 640.5
5.4
. display 5.4 - 1.9
3.5
```

Here we estimate the net effect of low/high student-teacher ratio holding constant the percent of English learners. When hi\_el is 0 the effect of going from low to high student-teacher ratio goes from a score of 664.2 to 662.2, a difference of 1.9. From a policy perspective you could argue that moving from high str to low str improve test scores by 1.9 in low English learners districts.

When hi\_el is 1, the effect of going from low to high student-teacher ratio goes from a score of 645.9 down to 640.5, a decline of 5.4 points (1.9+3.5). From a policy perspective you could say that reducing the str in districts with high percentage of English learners could improve test scores by 5.4 points.

#### Regression: interaction between a dummy and a continuous variable

Lets explore the same interaction as before but we keep student-teacher ratio continuous and the English learners variable as binary. The question remains the same\*.

- Dependent variable (Y) Average test score, variable testscr in dataset.
- Independent variables (X)
  - Continuous str, student-teacher ratio.
  - Binary hi\_el, where '0' if English learners (el\_pct) is lower than 10%, '1' equal to 10% or higher
  - Interaction term str\_el2 = str \* hi\_el. In Stata: generate str\_el2 = str\*hi\_el

We will run the regression

regress testscr str hi\_el str\_el2, robust

. regress testscr str hi_el str_el2, robust							
Linear regress	sion				Number of obs F( 3, 416) Prob > F R-squared Root MSE	= 420 = 63.67 = 0.0000 = 0.3103 = 15.88	
testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf.	Interval]	
str hi_el str_el2 _cons	9684601 5.639141 -1.276613 682.2458	.5891016 19.51456 .9669194 11.86781	-1.64 0.29 -1.32 57.49	0.101 0.773 0.187 0.000	-2.126447 -32.72029 -3.17727 658.9175	.1895268 43.99857 .6240436 705.5742	

The equation is testscr\_hat = 682.2 - 0.97\*str + 5.6\*hi\_el - 1.28\*str\_el2

The effect of str on testscr will be mediated by hi\_el.

- If hi\_el is 0 (low) then the effect of str is 682.2 0.97\*str.
- If hi\_el is 1 (high) then the effect of str is 682.2 0.97\*str + 5.6 1.28\*str = 687.8 2.25\*str

Notice that how hi\_el changes both the intercept and the slope of str. Reducing str by one in low EL districts will increase test scores by 0.97 points, but it will have a higher impact (2.25 points) in high EL districts. The difference between these two effects is 1.28 which is the coefficient of the interaction (Stock and Watson, 2003, p.223).

#### Regression: interaction between two continuous variables

Lets keep now both variables continuous. The question remains the same\*.

- Dependent variable (Y) Average test score, variable testscr in dataset.
- Independent variables (X)
  - Continuous str, student-teacher ratio.
  - Continuous el\_pct, percent of English learners.
  - Interaction term str\_el3 = str \* el\_pct. In Stata: generate str\_el3 = str\*el\_pct

We will run the regression

regress testscr str el\_pct str\_el3, robust

. regress test	scr str el_p	ct str_el3,	robust			
Linear regress	sion				Number of obs F( 3, 416) Prob > F R-squared Root MSE	= 420 = 155.05 = 0.0000 = 0.4264 = 14.482
testscr	Coef.	Robust Std. Err.	t	P>[t]	[95% Conf.	Interval]
str el_pct str_el3 _cons	-1.117018 6729116 .0011618 686.3385	.5875135 .3741231 .0185357 11.75935	-1.90 -1.80 0.06 58.37	0.058 0.073 0.950 0.000	-2.271884 -1.408319 0352736 663.2234	.0378468 .0624958 .0375971 709.4537

The equation is testscr\_hat = 686.3 - 1.12\*str - 0.67\*el\_pct + 0.0012\*str\_el3

The effect of the interaction term is very small. Following Stock and Watson (2003, p.229), algebraically the slope of str is

-1.12 + 0.0012\*el\_pct (remember that str\_el3 is equal to str\*el\_pct). So:

- If el\_pct = 10, the slope of str is -1.108
- If el\_pct = 20, the slope of str is -1.096. A difference in effect of 0.012 points.

In the continuous case there is an effect but is very small (and not significant). See Stock and Watson, 2003, for further details.

# Creating dummies

You can create dummy variables by either using recode or using a combination of tab/gen commands:

tab major, generate(major\_dum)

. t	ab major,	generate(major_dum)						
	Major	Freq.	Percent	⊂um.				
	Econ Math Politics	10 10 10	33.33 33.33 33.33	33.33 66.67 100.00				
	Total	30	100.00					

Check the 'variables' window, at the end you will see three new variables. Using tab1 (for multiple frequencies) you can check that they are all 0 and 1 values

Variables		×
Name	Label	^
city	City	
state	State	
gender	Gender	
status	Status: grad or undergad	
major	Major	
country	Country	
age	Age	
sat	SAT	
score	Average score (grade)	
height	Height (in)	
readnews	Newspaper read / week	
score2	Score in decimals	
readnews2	Monthly readership	
agegroups	Age by groups	
sex	Gender	
major_dum1	major==Econ	
major_dum2	major==Math	
major_dum3	major==Politics	
<		>

. tabl major	dum1 major_dum	12 major_d	um3	
-> tabulation of major_dum1				
major==Econ	Freq.	Percent	⊂um.	
0 1	20 10	66.67 33.33	66.67 100.00	
Total	30	100.00		
-> tabulation of major_dum2				
major==Math	Freq.	Percent	Cum.	
0 1	20 10	66.67 33.33	66.67 100.00	
Total	30	100.00		
-> tabulation of major_dum3				
major==Poli tics	Freq.	Percent	Cum.	
0 1	20 10	66.67 33.33	66.67 100.00	
Total	30	100.00	40	

# Creating dummies (cont.)

#### Here is another example:

tab agregroups, generate(agegroups\_dum)

. tab agegrou	ups, generate	generate(agegrups_dum)	
Age by groups	Freq.	Percent	⊂um.
18 to 19 20 to 29 30 to 39	10 9 11	33.33 30.00 36.67	33.33 63.33 100.00
Total	30	100.00	

Check the 'variables' window, at the end you will see three new variables. Using tab1 (for multiple frequencies) you can check that they are all 0 and 1 values

	Variables		×
	Name	Label	^
	status	Status: grad or undergad	
	major	Major	
	country	Country	
	age	Age	
	sat	SAT	
	score	Average score (grade)	_
	height	Height (in)	
	readnews	Newspaper read / week	
	score2	Score in decimals	
	readnews2	Monthly readership	
	agegroups	Age by groups	
	sex	Gender	
	major_dum1	major==Econ	
	major_dum2	major==Math	
	major_dum3	major==Politics	
1	agegrups_dum1	agegroups==18 to 19	
	agegrups_dum2	agegroups==20 to 29	
	agegrups_dum3	agegroups==30 to 39	~
			-
	5		

. tab1 agegr	ups_dum1 a	igegrups_dum2	2 agegrups_dum3	
-> tabulation of agegrups_dum1				
agegroups== 18 to 19	Fred	4. Percer	nt Cum.	
0 1	2 1	0 66.0 .0 33.3	57 66.67 33 100.00	
Total	3	0 100.0	00	
-> tabulation	-> tabulation of agegrups_dum2			
agegroups== 20 to 29	Fred	ą. Percer	nt Cum.	
0 1	2	1 70.0 9 30.0	00 70.00 00 100.00	
Total	3	0 100.0	00	
-> tabulation of agegrups_dum3				
agegroups== 30 to 39	Fred	ą. Percer	nt Cum.	
0 1	1 1	.9 63.3 1 36.6	63.33 67 100.00	
Total	3	0 100.0	00 41	

# Frequently used Stata commands

Category	Stata commands
Getting on-line help	help
	search
Operating-system interface	pwd
	cd
	sysdir
	mkdir
	dir / ls
	erase
	сору
	type
Using and saving data from disk	use
	clear
	save
	append
	merge
	compress
Inputting data into Stata	input
	edit
	infile
	infix
	insheet
The Internet and Updating Stata	update
	net
	ado
	news

	Basic data reporting	describe
		codebook
6		inspect
lo		list
rce		browse
: ht		count
ttp:/		assert
/w/		summarize
WW.		Table (tab)
ats.		tabulate
.ucl	Data manipulation	generate
a.e		replace
du/		egen
sta		recode
t/sta		rename
ata/		drop
'not		keep
es2		sort
2/cc		encode
nm		decode
nan		order
ds.		by
htm		reshape
	Formatting	format
	Kaaning track of your work	
	Reeping track of your work	notos
	Convenience	display purper/ottp
	convenience	aispiay PU/DSS/OTR

# Useful links / Recommended books

- ESS https://economics.princeton.edu/undergraduate-program/ess/#
- UCLA Resources to learn and use STATA <u>http://www.ats.ucla.edu/stat/stata/</u>
- Introduction to Stata (PDF), Christopher F. Baum, Boston College, USA. "A 67-page description of Stata, its key features and benefits, and other useful information." <a href="http://fmwww.bc.edu/GStat/docs/StataIntro.pdf">http://fmwww.bc.edu/GStat/docs/StataIntro.pdf</a>
- STATA FAQ website <u>http://stata.com/support/faqs/</u>

#### Books

- Introduction to econometrics / James H. Stock, Mark W. Watson. 2nd ed., Boston: Pearson Addison Wesley, 2007.
- Data analysis using regression and multilevel/hierarchical models / Andrew Gelman, Jennifer Hill. Cambridge ; New York : Cambridge University Press, 2007.
- Econometric analysis / William H. Greene. 6th ed., Upper Saddle River, N.J. : Prentice Hall, 2008.
- Designing Social Inquiry: Scientific Inference in Qualitative Research / Gary King, Robert O. Keohane, Sidney Verba, Princeton University Press, 1994.
- Unifying Political Methodology: The Likelihood Theory of Statistical Inference / Gary King, Cambridge University Press, 1989
- Statistical Analysis: an interdisciplinary introduction to univariate & multivariate methods / Sam Kachigan, New York : Radius Press, c1986
- Statistics with Stata (updated for version 9) / Lawrence Hamilton, Thomson Books/Cole, 2006