



# Linear Regression

(ver. 6.0)

*Oscar Torres-Reyna*  
*Data Consultant*  
*otorres@princeton.edu*



We use regression to estimate the unknown effect 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.  $Y = x_1 + x_2 + \dots + x_N$ ).

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)]  
regress 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.

**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). [Click here](http://www.duxbury.com/highered/) to download the data or search for it at <http://www.duxbury.com/highered/>. Use the file `states.dta` (educational data for the U.S.).

# Regression: variables

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

use `http://dss.princeton.edu/training/states.dta`

`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	display format	value label	variable label
<code>csat</code>	int	%9.0g		Mean composite SAT score
<code>expense</code>	int	%9.0g		Per pupil expenditures prim&sec
<code>percent</code>	byte	%9.0g		% HS graduates taking SAT
<code>income</code>	double	%10.0g		Median household income, \$1,000
<code>high</code>	float	%9.0g		% adults HS diploma
<code>college</code>	float	%9.0g		% adults college degree
<code>region</code>	byte	%9.0g	region	Geographical region

`. summarize csat expense percent income high college region`

Variable	Obs	Mean	Std. Dev.	Min	Max
<code>csat</code>	51	944.098	66.93497	832	1093
<code>expense</code>	51	5235.961	1401.155	2960	9259
<code>percent</code>	51	35.76471	26.19281	4	81
<code>income</code>	51	33.95657	6.423134	23.465	48.618
<code>high</code>	51	76.26078	5.588741	64.3	86.6
<code>college</code>	51	20.02157	4.16578	12.3	33.3
<code>region</code>	50	2.54	1.128662	1	4

# Regression: what to look for

Lets run the regression:

```
regress csat expense, robust
```

Outcome variable (Y)

Predictor variable (X)

Robust standard errors (to control for heteroskedasticity)

1

This is the p-value of the model. It tests whether  $R^2$  is different from 0. Usually we need a p-value lower than 0.05 to show a statistically significant relationship between X and Y.

```
. regress csat expense, robust
```

Linear regression on

```
Number of obs = 51
F( 1, 49) = 36.80
Prob > F = 0.0000
R-squared = 0.2174
Root MSE = 59.814
```

2

R-square shows the amount of variance of Y explained by X. In this case expense explains 22% of the variance in SAT scores.

7

Root MSE: root mean squared error, is the sd of the regression. The closer to zero better the fit.

csat	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
expense	-.0222756	.0036719	-6.07	0.000	-.0296547 -.0148966
_cons	1060.732	24.35468	43.55	0.000	1011.79 1109.675

6

$csat = 1061 - 0.022 * expense$   
For each one-point increase in expense, SAT scores decrease by 0.022 points.

3

Adj  $R^2$  (not shown here) shows the same as  $R^2$  but adjusted by the # of cases and # of variables. When the # of variables is small and the # of cases is very large then  $Adj R^2$  is closer to  $R^2$ . This provides a more honest association between X and Y.

5

The t-values test the hypothesis that the coefficient is different from 0. To reject this, you need a t-value greater than 1.96 (for 95% confidence). You can get the t-values by dividing the coefficient by its standard error. The t-values also show the importance of a variable in the model.

4

Two-tail p-values test the hypothesis that each coefficient is different from 0. To reject this, the p-value has to be lower than 0.05 (you could choose also an alpha of 0.10). In this case, expense is statistically significant in explaining SAT.

# Regression: what to look for

Adding the rest of predictor variables:

```
regress csat expense percent income high college, robust
```

Robust standard errors (to control for heteroskedasticity)

Output variable (Y)

Predictor variables (X)

1

This is the p-value of the model. It indicates the reliability of X to predict Y. Usually we need a p-value lower than 0.05 to show a statistically significant relationship between X and Y.

```
. regress csat expense percent income high college, robust
```

Linear regression

```
Number of obs = 51
F( 5, 45) = 50.90
Prob > F = 0.0000
R-squared = 0.8243
Root MSE = 29.571
```

7

Root MSE: root mean squared error, is the sd of the regression. The closer to zero better the fit.

2

R-square shows the amount of variance of Y explained by X. In this case the model explains 82.43% of the variance in SAT scores.

csat	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
expense	.0033528	.004781	0.70	0.487	-.0062766 .0129823
percent	-2.618177	.2288594	-11.44	0.000	-3.079123 -2.15723
income	.1055853	1.207246	0.09	0.931	-2.325933 2.537104
high	1.630841	.943318	1.73	0.091	-.2690989 3.530781
college	2.030894	2.113792	0.96	0.342	-2.226502 6.28829
_cons	851.5649	57.28743	14.86	0.000	736.1821 966.9477

3

Adj R<sup>2</sup> (not shown here) shows the same as R<sup>2</sup> but adjusted by the # of cases and # of variables. When the # of variables is small and the # of cases is very large then Adj R<sup>2</sup> is closer to R<sup>2</sup>. This provides a more honest association between X and Y.

$$csat = 851.56 + 0.003 * expense - 2.62 * percent + 0.11 * income + 1.63 * high + 2.03 * college$$

6

5

The t-values test the hypothesis that the coefficient is different from 0. To reject this, you need a t-value greater than 1.96 (at 0.05 confidence). You can get the t-values by dividing the coefficient by its standard error. The t-values also show the importance of a variable in the model. In this case, percent is the most important.

4

Two-tail p-values test the hypothesis that each coefficient is different from 0. To reject this, the p-value has to be lower than 0.05 (you could choose also an alpha of 0.10). In this case, expense, income, and college are not statistically significant in explaining SAT; high is almost significant at 0.10. Percent is the only variable that has some significant impact on SAT (its coefficient is different from 0)

Region is entered here as dummy variable. The easy way to add dummy variables to a regression is using “xi” and the prefix “i.” (interpretation is the same as before). The first category is always the reference:

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

```
. xi: regress csat expense percent income high college i.region, robust
i.region      _Iregion_1-4      (naturally coded; _Iregion_1 omitted)
```

Linear regression

```
Number of obs =      50
F( 8,      41) =     69.82
Prob > F      =     0.0000
R-squared     =     0.9111
Root MSE     =     21.492
```

	csat	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
	expense	-.002021	.0035883	-0.56	0.576	-.0092676 .0052256
	percent	-3.007647	.2358047	-12.75	0.000	-3.483864 -2.53143
	income	-.1674421	1.196409	-0.14	0.889	-2.583638 2.248754
	high	1.814731	1.02694	1.77	0.085	-.2592168 3.888679
	college	4.670564	1.599798	2.92	0.006	1.439705 7.901422
Regions:	_Iregion_2	69.45333	17.99933	3.86	0.000	33.10295 105.8037
1 West	_Iregion_3	25.39701	12.52558	2.03	0.049	.101086 50.69293
2 N. East	_Iregion_4	34.57704	9.44989	3.66	0.001	15.4926 53.66149
3 South	_cons	808.0206	67.86418	11.91	0.000	670.9661 945.0751
4 Midwest						

**NOTE:** By default xi excludes the first value, to select a different value, before running the regression type:

```
char region[omit] 4
```

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

This will select Midwest (4) as the reference category for the dummy variables.

**NOTE:** Another way to create dummy variables is to type:

```
tab region, gen(region)
```

This will create four new variables (or a many a categories in the variable), one for each region in this case.

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	
Model	(A) 200269.84	9	22252.2045 (D)	Number of obs = 50
Residual	(B) 12691.5396	40	317.28849 (E)	F( 9, 40) = 70.13
Total	(C) 212961.38	49	4346.15061 (F)	Prob > F = 0.0000

	R-squared = 0.9404
	Adj R-squared = 0.9270
	Root MSE = 17.813

$$F = \frac{\frac{MSS}{k-1}}{\frac{RSS}{n-k}} = \frac{\frac{200269.84}{9}}{\frac{12691.5396}{40}} = \frac{22252.2045}{317.28849} = \frac{D}{E} = 70.13$$

$$AdjR^2 = 1 - \frac{n-1}{n-k}(1-R^2) = 1 - \frac{49}{40}(1-0.9404) = 1 - \frac{E}{F} = 1 - \frac{317.28849}{4346.15061} = 0.9270$$

$$R^2 = \frac{MSS}{TSS} = 1 - \frac{\sum e_i^2}{\sum (y_i - \bar{y})^2} = \frac{200269.84}{12691.5396} = \frac{A}{C} = 0.9404$$

$$RootMSE = \sqrt{\frac{RSS}{(n-k)}} = \sqrt{\frac{12691.5396}{40}} = \sqrt{\frac{B}{40}} = 17.813$$

**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.



# Regression: eststo/esttab

To show the models side-by-side you can use the commands `eststo` 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) csat	(2) csat	(3) csat
expense	-0.0223*** (0.00367)	0.00335 (0.00478)	-0.00202 (0.00359)
percent		-2.618*** (0.229)	-3.008*** (0.236)
income		0.106 (1.207)	-0.167 (1.196)
high		1.631 (0.943)	1.815 (1.027)
college		2.031 (2.114)	4.671** (1.600)
_Iregion_2			69.45*** (18.00)
_Iregion_3			25.40* (12.53)
_Iregion_4			34.58*** (9.450)
percent2			
_cons	1060.7*** (24.35)	851.6*** (57.29)	808.0*** (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

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

Type `help eststo` and `help esttab` for more options.

## Regression: correlation matrix

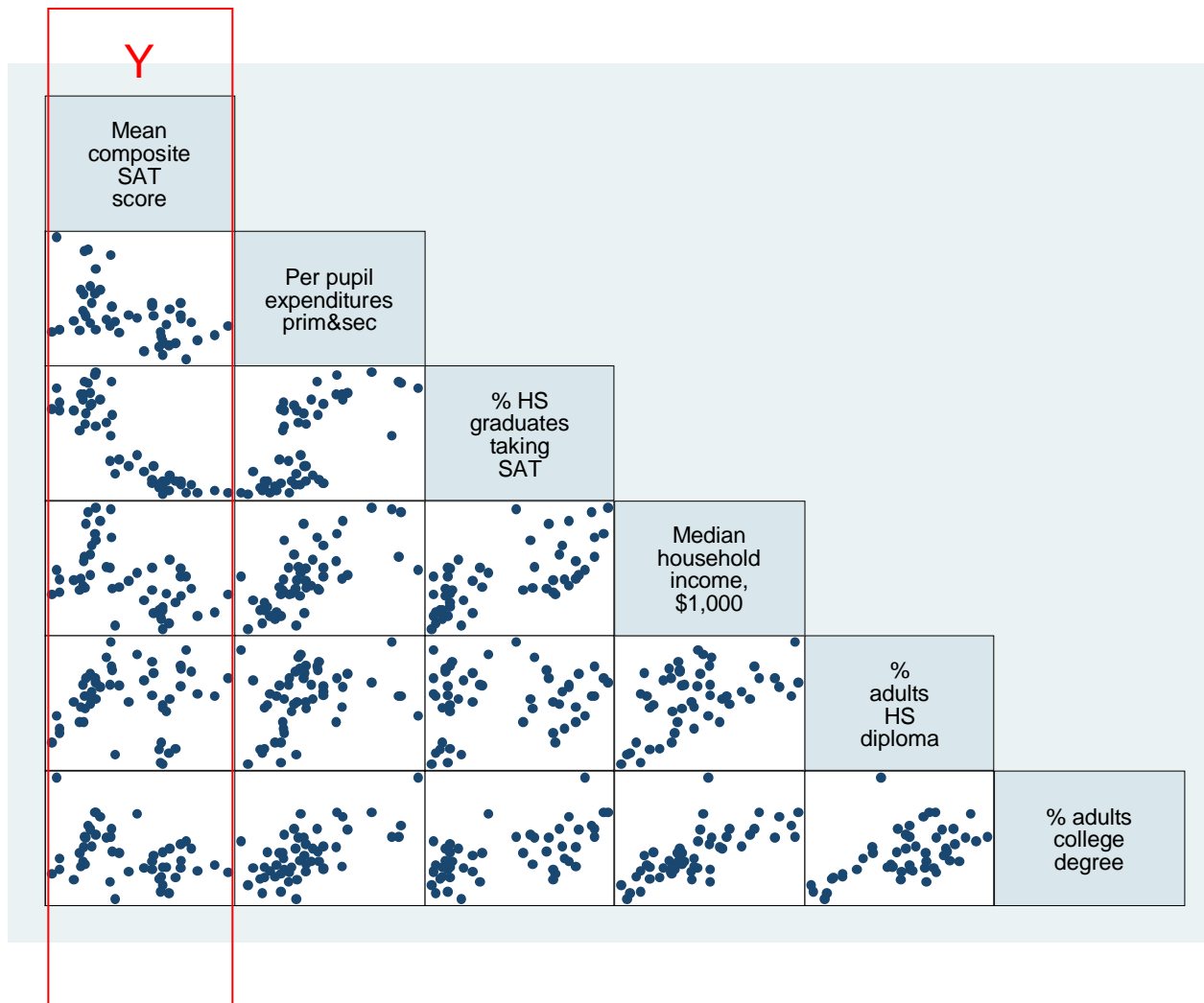
Below is a correlation matrix for all variables in the model. Numbers are Pearson correlation coefficients, go from -1 to 1. Closer to 1 means strong correlation. A negative value indicates an inverse relationship (roughly, when one goes up the other goes down).

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

	csat	expense	percent	income	high	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		
high	0.0858 0.5495	0.3133* 0.0252	0.1413 0.3226	0.5099* 0.0001	1.0000	
college	-0.3729* 0.0070	0.6400* 0.0000	0.6091* 0.0000	0.7234* 0.0000	0.5319* 0.0001	1.0000

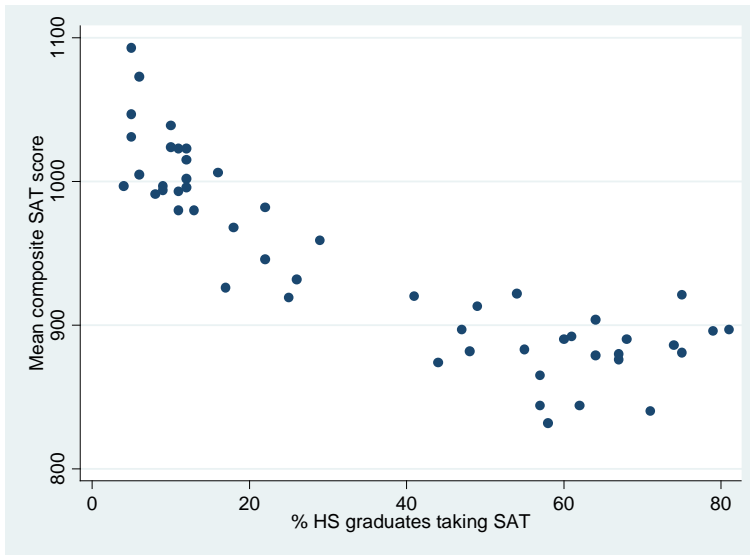
Command `graph matrix` produces a graphical representation of the correlation matrix by presenting a series of scatterplots for all variables. Type:

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

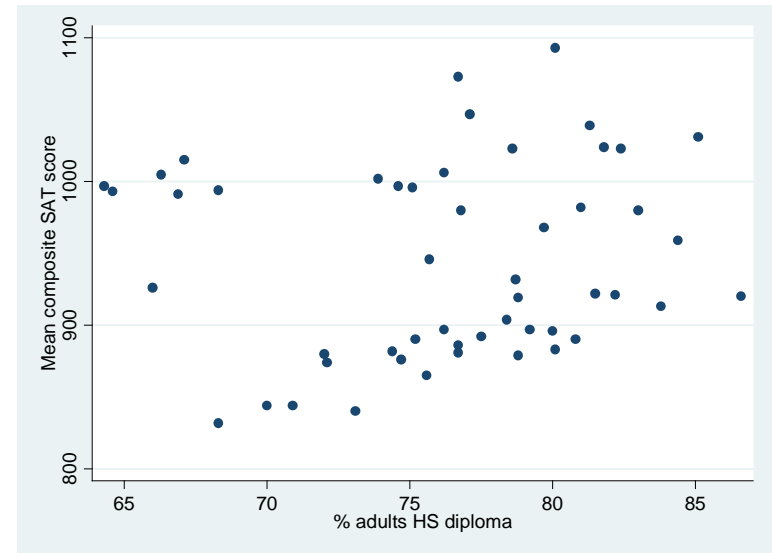


# Regression: exploring relationships

scatter csat percent



scatter csat high



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`

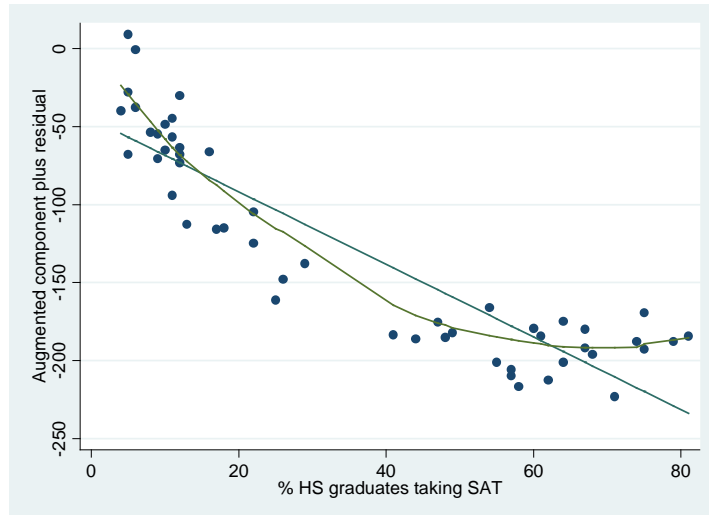
```
generate percent2 = percent^2
```

# Regression: functional form/linearity

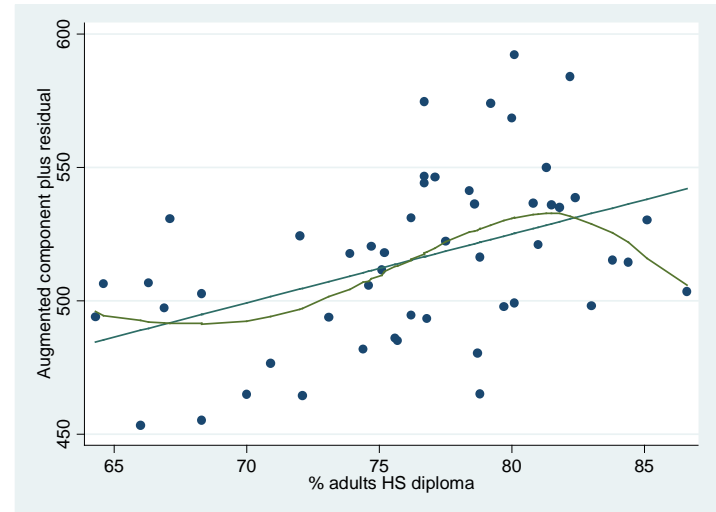
The command `acprplot` (augmented component-plus-residual plot) provides another graphical way to examine the relationship between variables. It does provide a good testing for linearity. Run this command after running a regression

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

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 <http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm>, 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) csat	(2) csat	(3) csat	(4) csat
expense	<b>-0.0223***</b> (0.00367)	<b>0.00335</b> (0.00478)	<b>-0.00202</b> (0.00359)	<b>0.00141</b> (0.00372)
percent		<b>-2.618***</b> (0.229)	<b>-3.008***</b> (0.236)	<b>-5.945***</b> (0.641)
income		<b>0.106</b> (1.207)	<b>-0.167</b> (1.196)	<b>-0.914</b> (0.973)
high		<b>1.631</b> (0.943)	<b>1.815</b> (1.027)	<b>1.869</b> (0.931)
college		<b>2.031</b> (2.114)	<b>4.671**</b> (1.600)	<b>3.418**</b> (1.145)
_Iregion_2			<b>69.45***</b> (18.00)	<b>5.077</b> (20.75)
_Iregion_3			<b>25.40*</b> (12.53)	<b>5.209</b> (10.42)
_Iregion_4			<b>34.58***</b> (9.450)	<b>19.25*</b> (8.110)
percent2				<b>0.0460***</b> (0.0102)
_cons	<b>1060.7***</b> (24.35)	<b>851.6***</b> (57.29)	<b>808.0***</b> (67.86)	<b>874.0***</b> (58.13)
N	51	51	50	50
R-sq	0.217	0.824	0.911	0.940
adj. R-sq	0.201	0.805	0.894	0.927
rmse	59.81	29.57	21.49	17.81

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

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:

Option A, 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[percent2]*percent2
+ _b[high]*high + ...
```

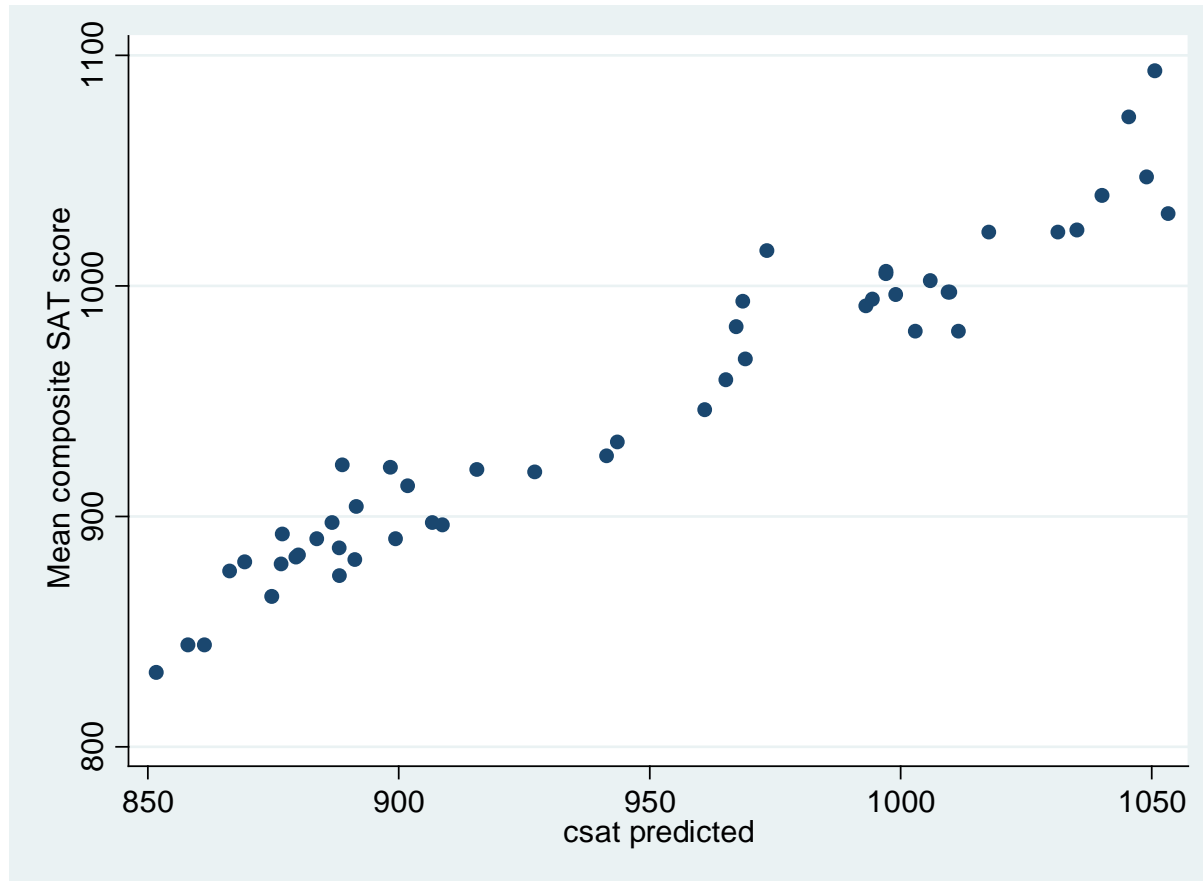
Option B, using `predict` immediately after running the regression:

```
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"

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`



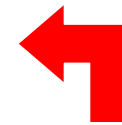
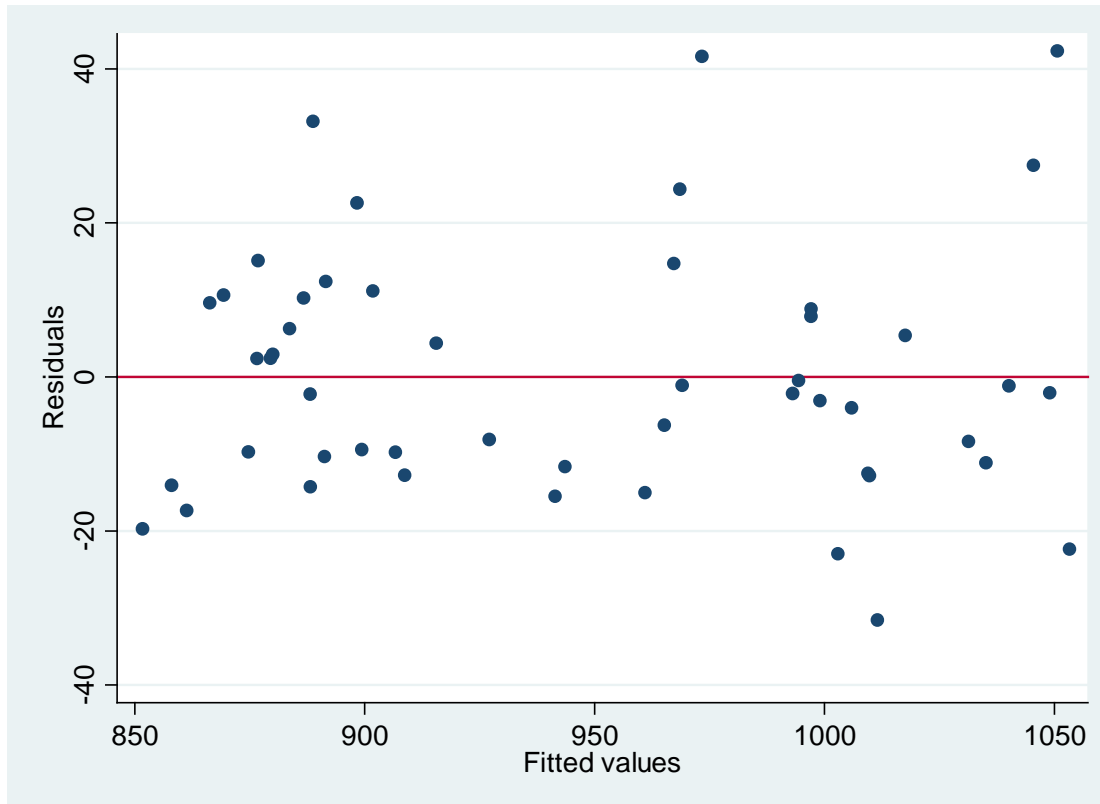
## Regression: testing for homoskedasticity

An important assumption is that the variance in the residuals has to be homoskedastic or constant. Residuals cannot varied for lower or 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  $X_i$   $[\text{var}(e_i|X_i)]$ , is constant for  $i=1\dots n$ , 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 ( $\hat{Y}$ ) we *should not observe* any pattern at all. In Stata we do this using `rvfplot` right after running the regression, it will automatically draw a scatterplot between residuals and predicted values.

`rvfplot, yline(0)`



Residuals seem to slightly expand at higher levels of  $\hat{Y}$ .

A non-graphical way to detect heteroskedasticity 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)      =      2.72
Prob > chi 2  =      0.0993
```

The graphical and the Breusch-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 heteroskedasticity 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) = 1.25  
Prob > 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.

## Regression: specification error

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  $\hat{Y}$  (`csat_predicted` or  $X\beta$ ) and  $\hat{Y}$ -squared as independent variables<sup>1</sup>.

The thing to look for here is the significance of `_hatsq`. 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
```

```
. linktest
```

Source	SS	df	MS			
Model	200272.359	2	100136.18	Number of obs =	50	
Residual	12689.0209	47	269.979169	F( 2, 47) =	370.90	
Total	212961.38	49	4346.15061	Prob > F =	0.0000	
				R-squared =	0.9404	
				Adj R-squared =	0.9379	
				Root MSE =	16.431	

csat	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
<code>_hat</code>	1.144949	1.50184	0.76	0.450	-1.876362	4.166261
<code>_hatsq</code>	-.0000761	.0007885	-0.10	0.923	-.0016623	.0015101
<code>_cons</code>	-68.69417	712.388	-0.10	0.924	-1501.834	1364.446

<sup>1</sup> For more details see <http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm>, and/or type `help linktest`.

## Regression: multicollinearity

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

When multicollinearity is present *standard 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:

```
. vif
```

Variable	VIF	1/VIF
percent2	70.80	0.014124
percent	49.52	0.020193
_Iregion_2	8.47	0.118063
income	4.97	0.201326
_Iregion_3	4.89	0.204445
high	4.71	0.212134
college	4.52	0.221348
expense	3.33	0.300111
_Iregion_4	2.14	0.467506
Mean VIF	17.04	

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

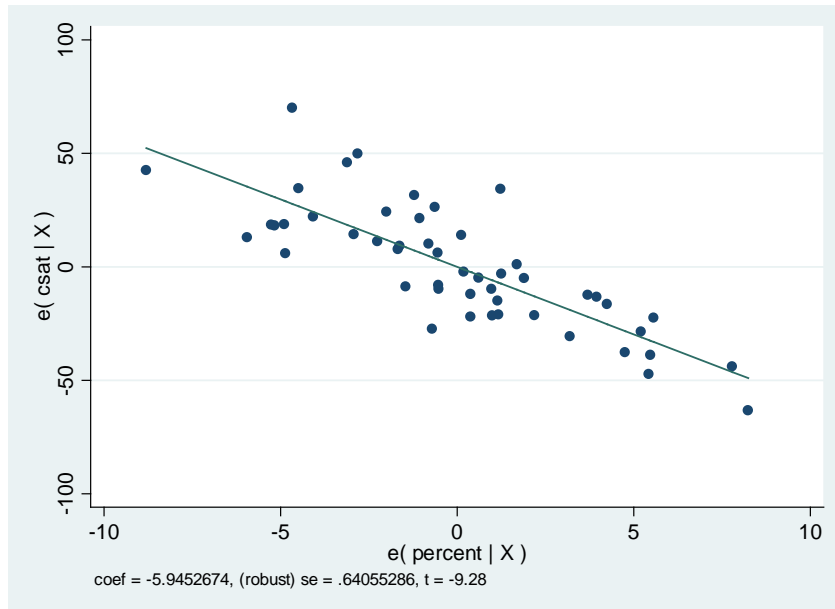
We know that `percent` and `percent2` are related since one is the square of the other. They are ok since `percent` has a quadratic relationship with  $Y$ , *but this would be an example of multicollinearity*.

The rest of the variables look ok.

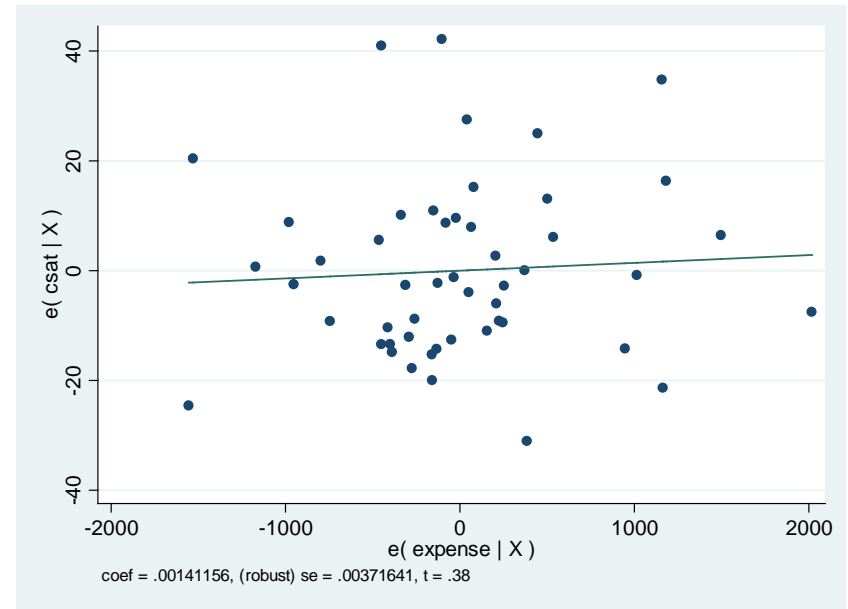
# Regression: outliers

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 percent`



`avplot expense`



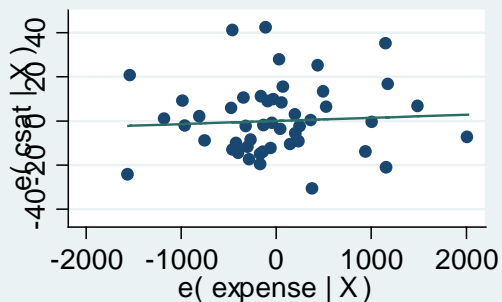
These plots regress each variable against all others, notice the coefficients on each. All data points seem to be in range, no outliers observed.

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

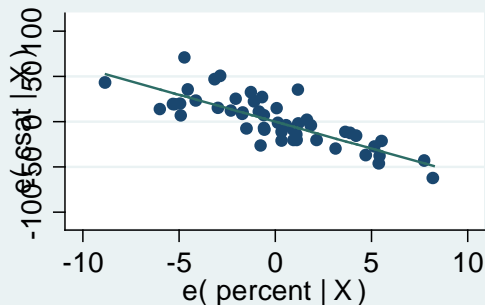
Also type `help diagplots` in the Stata command window.

# Regression: outliers

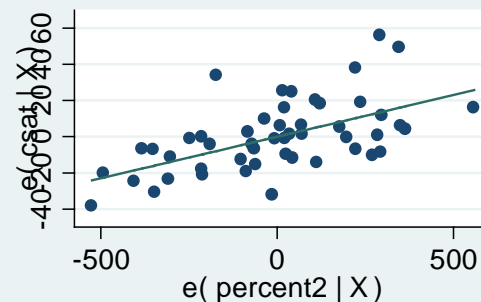
avplots



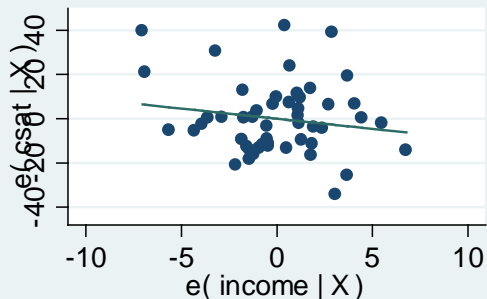
coef = .00141156, (robust) se = .00371641, t =



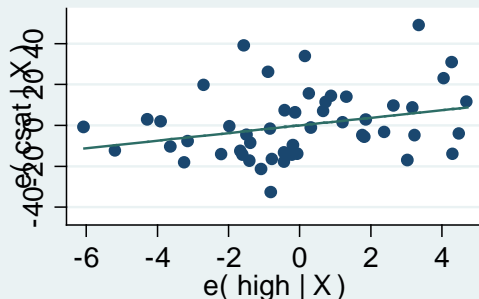
coef = -5.9452674, (robust) se = .64055286, t =



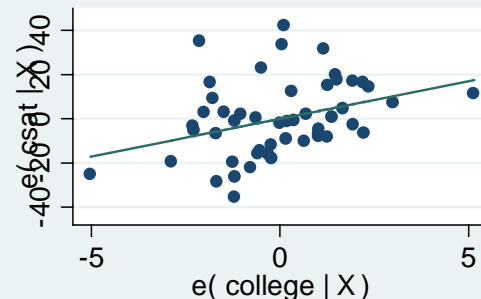
coef = .0460468, (robust) se = .01019105, t = 4.5



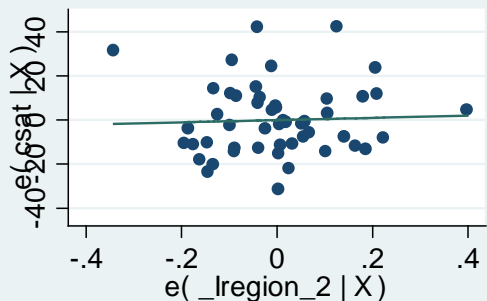
coef = -.9143708, (robust) se = .97326373, t =



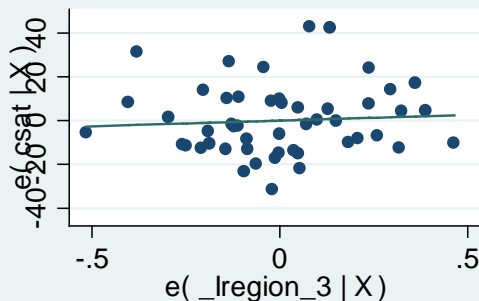
coef = 1.8691679, (robust) se = .93111302, t =



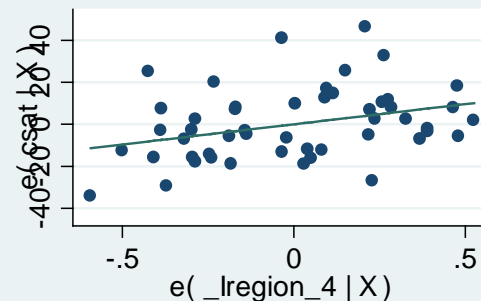
coef = 3.4175732, (robust) se = 1.1450333, t = 2.



coef = 5.0765963, (robust) se = 20.753948, t =



coef = 5.2088169, (robust) se = 10.422781, t =



coef = 19.245404, (robust) se = 8.1097615, t = 2.

# Regression: summary of *influence* indicators

<b>DfBeta</b>	<p>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.</p> <p>An observation is influential if it has a significant effect on the coefficient.</p>	<p>A case is an influential outlier if</p> $ DfBeta  > 2/\sqrt{N}$ <p>Where N is the sample size.</p> <p>Note: Stata estimates standardized DfBetas.</p>	<p>In Stata type:</p> <pre>reg y x1 x2 x3</pre> <pre>dfbeta x1</pre> <p>Note: you could also type:</p> <pre>predict DFx1, dfbeta(x1)</pre> <p>To estimate the dfbetas for all predictors just type:</p> <pre>dfbeta</pre> <p>To flag the cutoff</p> <pre>gen cutoffdfbeta = abs(DFx1) &gt; 2/sqrt(e(N)) &amp; e(sample)</pre>	<p>In SPSS: Analyze-Regression-Linear; click Save. Select under “Influence Statistics” to add as a new variable (DFB1_1) or in syntax type</p> <pre>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 .</pre>
<b>DfFit</b>	<p>Indicator of leverage and high residuals.</p> <p>Measures how much an observation influences the regression model as a whole.</p> <p>How much the predicted values change as a result of including and excluding a particular observation.</p>	<p>High influence if</p> $ DfFIT  > 2*\sqrt{k/N}$ <p>Where k is the number of parameters (including the intercept) and N is the sample size.</p>	<p>After running the regression type:</p> <pre>predict dfits if e(sample), dfits</pre> <p>To generate the flag for the cutoff type:</p> <pre>gen cutoffdfit= abs(dfits)&gt;2*sqrt((e(df_m) +1)/e(N)) &amp; e(sample)</pre>	<p>Same as DfBeta above (DFF_1)</p>
<b>Covariance ratio</b>	<p>Measures the impact of an observation on the standard errors</p>	<p>High impact if</p> $ COVRATIO-1  \geq 3*k/N$ <p>Where k is the number of parameters (including the intercept) and N is the sample size.</p>	<p>In Stata after running the regression type</p> <pre>predict covratio if e(sample), covratio</pre>	<p>Same as DfBeta above (COV_1)</p>



# Regression: summary of *distance* measures

<p><b>Cook's distance</b></p>	<p>Measures how much an observation influences the overall model or predicted values.</p> <p>It is a summary measure of leverage and high residuals.</p>	<p>High influence if</p> <p><math>D &gt; 4/N</math></p> <p>Where N is the sample size.</p> <p>A <math>D &gt; 1</math> indicates big outlier problem</p>	<p>In Stata after running the regression type:</p> <p><code>predict D, cooks</code></p>	<p>In SPSS: Analyze-Regression-Linear; click Save. Select under "Distances" to add as a new variable (COO_1) or in syntax type</p> <pre>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.</pre>
<p><b>Leverage</b></p>	<p>Measures how much an observation influences regression coefficients.</p>	<p>High influence if</p> <p>leverage <math>h &gt; 2*k/N</math></p> <p>Where k is the number of parameters (including the intercept) and N is the sample size.</p> <p>A rule-of-thumb: Leverage goes from 0 to 1. A value closer to 1 or over 0.5 may indicate problems.</p>	<p>In Stata after running the regression type:</p> <p><code>predict lev, leverage</code></p>	<p>Same as above (LEV_1)</p>
<p><b>Mahalanobis distance</b></p>	<p>It is rescaled measure of leverage.</p> <p><math>M = \text{leverage} * (N-1)</math></p> <p>Where N is sample size.</p>	<p>Higher levels indicate higher distance from average values.</p> <p>The M-distance follows a Chi-square distribution with k-1 df and <math>\alpha=0.001</math> (where k is the number of independent variables).</p> <p>Any value over this Chi-square value may indicate problems.</p>	<p>Not available</p>	<p>Same as above (MAH_1)</p>

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 <http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm>

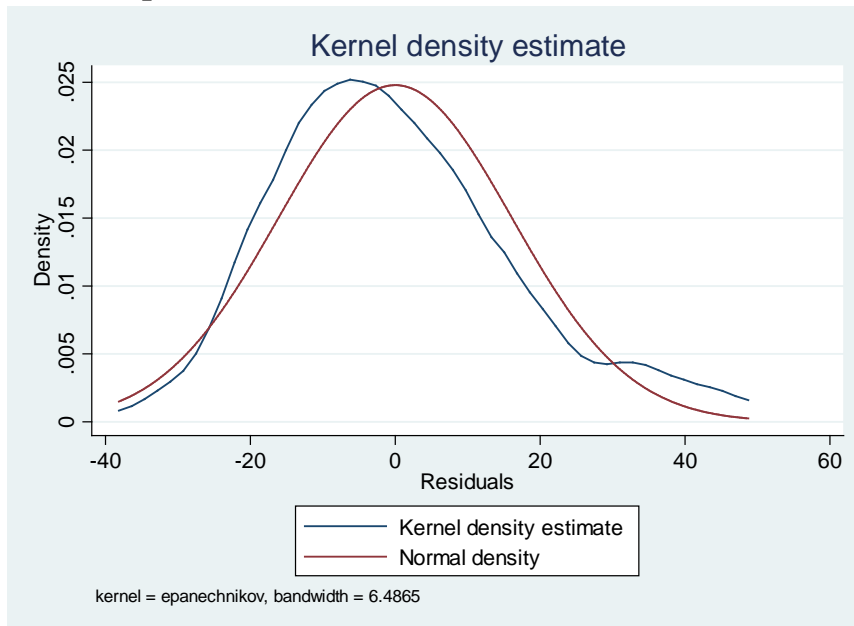
# Regression: testing for normality

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 ( $\hat{Y}$ ):  $e = Y - \hat{Y}$ .

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

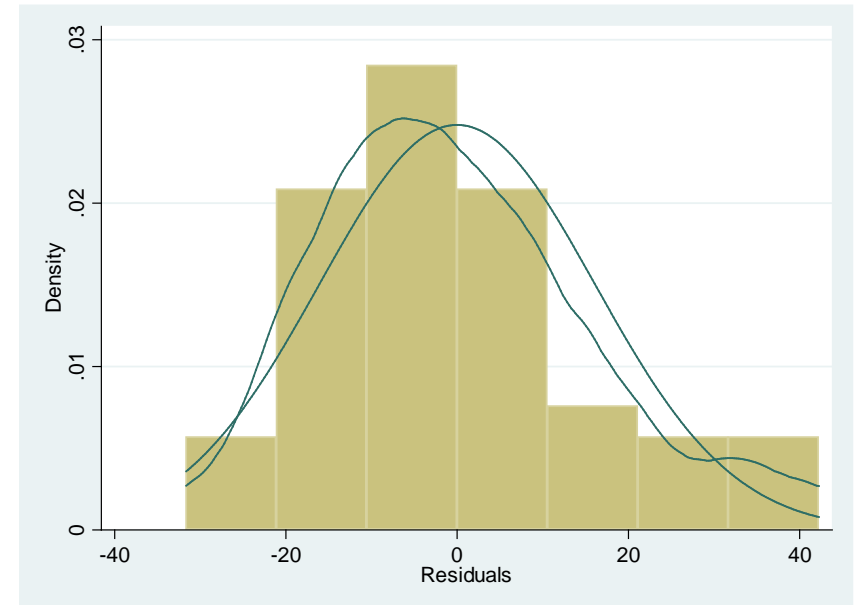
Three graphs will help us check for normality in the residuals: `kdensity`, `pnorm` and `qnorm`.

`kdensity e, normal`



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`.

`histogram e, kdensity normal`

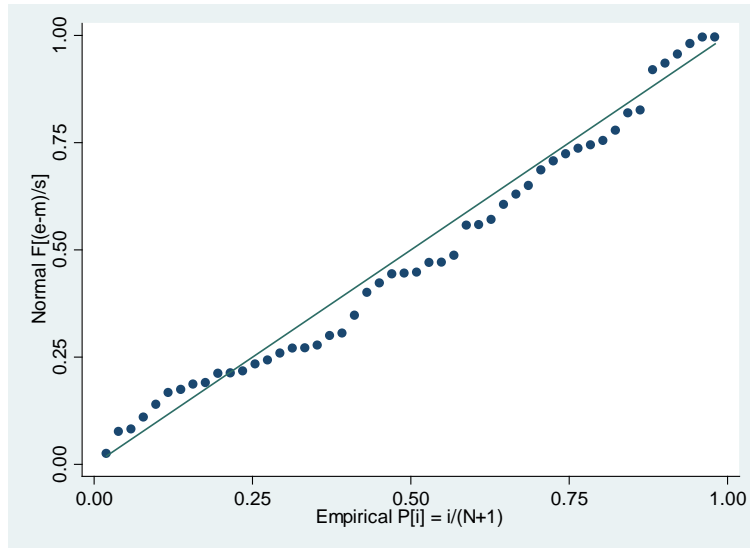


If residuals do not follow a 'normal' pattern then you should check for omitted variables, model specification, linearity, functional forms. In sum, you may need to reassess your model/theory. In practice normality does not represent much of a problem when dealing with really big samples.

# Regression: testing for normality

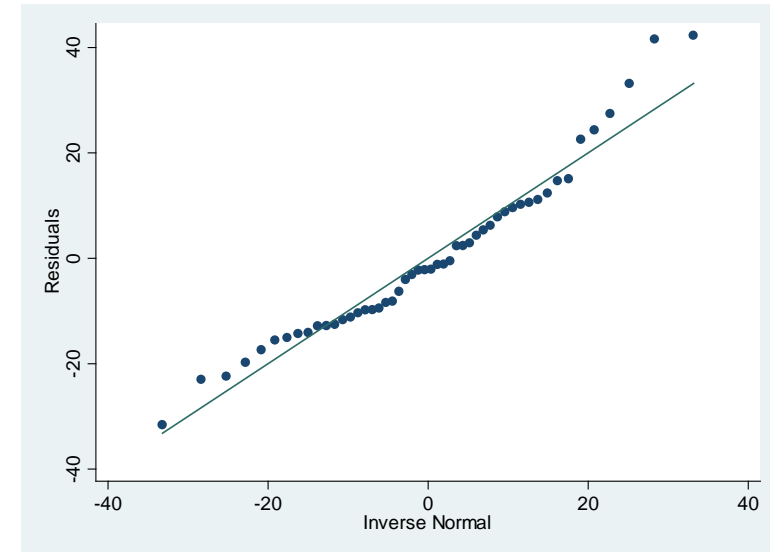
Standardize normal probability plot (`pnorm`) checks for non-normality in the middle range of residuals. Again, slightly off the line but looks ok.

`pnorm e`



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.

`qnorm e`



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`

Shapiro-Wilk W test for normal data

Variable	Obs	W	V	z	Prob>z
e	50	0.95566	2.085	1.567	0.05855

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)  high = 0
( 2)  college = 0
```

```
      F( 2,      40) =      17.12
      Prob > 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
```

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 constant_b constant_se
```

Variable	Obs	Mean	Std. Dev.	Min	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      _Iregion_1-4      (naturally coded; _Iregion_1 omitted)

. 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(ll) = -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) : "ols"
      e(estat_cmd) : "regress_estat"
      e(vcetype) : "Robust"

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

functions:
      e(sample)
```

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.

\*Gelman, Andrew, Jennifer Hill, *Data Analysis Using Regression and Multilevel/Hierarchical Models*, 2007, p. 69



# Regression: publishing regression output (outreg2)

The command `outreg2` gives you the type of presentation you see in published papers. If `outreg2` is not available you need to install it by typing `ssc install outreg2`

Let's say the regression is `regress csat percent percent2 high, robust`

The basic syntax for `outreg2` is: `outreg2 using [pick a name], [type either word or excel]`

After the regression type the following if you want to export the **results to excel\***

```
outreg2 using results, excel
```

```
. outreg2 using results, excel  
"results.xml"  
seeout
```

Click here to see the file



Or this if you want to **export to word**

```
outreg2 using results, word
```

```
. outreg2 using results, word  
"results.rtf"  
seeout
```

Click here to see the file



In excel

	A	B
1	v1	v2
2		(3)
3	COEFFICIENT	csat
4		
5	percent	-6.520***
6		(0.49)
7	percent2	0.0537***
8		(0.0056)
9	high	2.987***
10		(0.55)
11	Constant	844.8***
12		(38.8)
13	Observations	51
14	R-squared	0.93
15	Robust standard errors in parentheses	
16	*** p<0.01, ** p<0.05, * p<0.10	

In word

COEFFICIENT	csat
percent	-6.520***
	(0.49)
percent2	0.0537***
	(0.0056)
high	2.987***
	(0.55)
Constant	844.8***
	(38.8)
Observations	51
R-squared	0.93

Robust standard errors in parentheses  
\*\*\* p<0.01, \*\* p<0.05, \* p<0.10

# Regression: publishing regression output (outreg2)

You can add more models to compare. Lets say you want to add another model without percent2:

```
regress csat percent high, robust
```

Now type to export the results to word (**notice** we add the append option)

```
outreg2 using results, word append
```

In excel

	COEFFICIENT	csat	csat
3			
4			
5	percent	-6.520***	-2.315***
6		(0.49)	(0.17)
7	percent2	0.0537***	
8		(0.0056)	
9	high	2.987***	2.561***
10		(0.55)	(0.72)
11	Constant	844.8***	831.6***
12		(38.8)	(53.5)
13	Observations	51	51
14	R-squared	0.93	0.81
15	Robust standard errors in parentheses		
16	*** p<0.01, ** p<0.05, * p<0.1		

In word

COEFFICIENT	csat	csat
percent	-6.520***	-2.315***
	(0.49)	(0.17)
percent2	0.0537***	
	(0.0056)	
high	2.987***	2.561***
	(0.55)	(0.72)
Constant	844.8***	831.6***
	(38.8)	(53.5)
Observations	51	51
R-squared	0.93	0.81
Robust standard errors in parentheses		
*** p<0.01, ** p<0.05, * p<0.1		

**NOTE:** If you run logit/probit regression with odds ratios you need to add the option `eform` to export the odd ratios

Type `help outreg2` for more details. If you do not see `outreg2`, you may have to install it by typing `ssc install outreg2`. If this does not work type `findit outreg2`, select from the list and click "install".

Note: If you get the following error message (when you use the option `append` or `replace` it means that you need to close the excel/word window.

**file results.rtf is read-only; cannot be modified or erased**

# Regression: publishing regression output (outreg2) continue

For a customized look, here are some options:

\*\*\* Excel

```
outreg2 using results, bdec(2) tdec(2) rdec(2) addec(2) alpha (0.01, 0.05, 0.10)
addstat(Adj. R-squared, e(r2_a)) excel
```

\*\*\* Word

```
outreg2 using results, bdec(2) tdec(2) rdec(2) addec(2) alpha (0.01, 0.05, 0.10)
addstat(Adj. R-squared, e(r2_a)) word
```

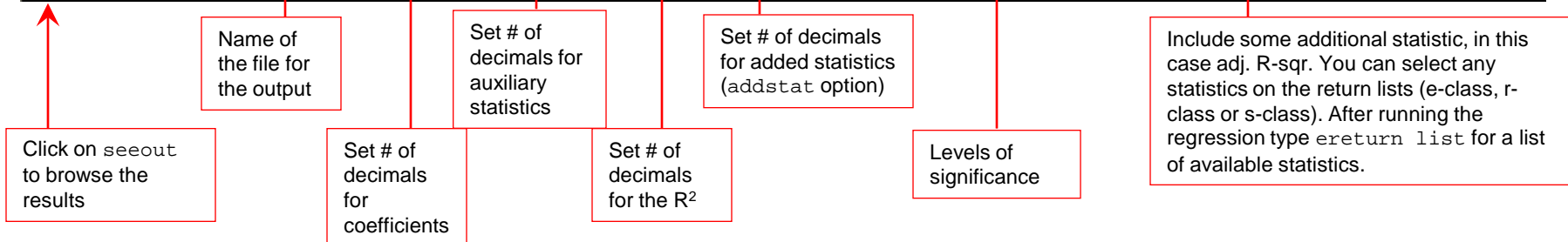
For excel

```
. outreg2 using results, bdec(2) tdec(2) rdec(2) addec(2) alpha(0.001, 0.01, 0.05) addstat(Adj. R-squared, e
> (r2_a)) excel
"results.xml"
seeout
```

For word

```
. outreg2 using results, bdec(2) tdec(2) rdec(2) addec(2) alpha(0.001, 0.01, 0.05) addstat(Adj. R-squared, e
> (r2_a)) word
"results.rtf"
seeout
```

Click here to see the output, a excel/word window will open



# 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 testscr hi_el hi_str str_el, robust
```

Linear regression

Number of obs =	420
F( 3, 416) =	60.20
Prob > F =	0.0000
R-squared =	0.2956
Root MSE =	16.049

testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
hi_el	-18.16295	2.345952	-7.74	0.000	-22.77435 -13.55155
hi_str	-1.907842	1.932215	-0.99	0.324	-5.705964 1.890279
str_el	-3.494335	3.121226	-1.12	0.264	-9.629677 2.641006
_cons	664.1433	1.388089	478.46	0.000	661.4147 666.8718

The equation is  $\text{testscr}_{\text{hat}} = 664.1 - 18.1 \cdot \text{hi\_el} - 1.9 \cdot \text{hi\_str} - 3.5 \cdot \text{str\_el}$

The effect of `hi_str` on the test scores is -1.9 but given the interaction term (and assuming all coefficients are significant), the net effect is  $-1.9 - 3.5 \cdot \text{hi\_el}$ . 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)
(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).

```
. summarize yhat1 yhat2 yhat3 yhat4
```

variable	Obs	Mean	Std. Dev.	Min	Max
yhat1	149	664.1433	0	664.1433	664.1433
yhat2	79	662.2355	0	662.2355	662.2355
yhat3	89	645.9803	0	645.9803	645.9803
yhat4	103	640.5782	0	640.5782	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 regression

Number of obs =	420
F( 3, 416) =	63.67
Prob > F =	0.0000
R-squared =	0.3103
Root MSE =	15.88

testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]
str	-.9684601	.5891016	-1.64	0.101	-2.126447 .1895268
hi_el	5.639141	19.51456	0.29	0.773	-32.72029 43.99857
str_el2	-1.276613	.9669194	-1.32	0.187	-3.17727 .6240436
_cons	682.2458	11.86781	57.49	0.000	658.9175 705.5742

The equation is  $\text{testscr}_{\text{hat}} = 682.2 - 0.97 \cdot \text{str} + 5.6 \cdot \text{hi\_el} - 1.28 \cdot \text{str\_el2}$

The effect of `str` on `testscr` will be mediated by `hi_el`.

- If `hi_el` is 0 (low) then the effect of `str` is  $682.2 - 0.97 \cdot \text{str}$ .
- If `hi_el` is 1 (high) then the effect of `str` is  $682.2 - 0.97 \cdot \text{str} + 5.6 - 1.28 \cdot \text{str} = 687.8 - 2.25 \cdot \text{str}$

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).

\*The data used in this section is the "California Test Score" data set (`caschool.dta`) from chapter 6 of the book *Introduction to Econometrics* from Stock and Watson, 2003. Data can be downloaded from [http://wps.aw.com/aw\\_stock\\_ie\\_2/50/13016/3332253.cw/index.html](http://wps.aw.com/aw_stock_ie_2/50/13016/3332253.cw/index.html). For a detailed discussion please refer to the respective section in the book.

# 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 testscr str el_pct str_el3, robust
```

Linear regression

	Number of obs = 420
	F( 3, 416) = 155.05
	Prob > F = 0.0000
	R-squared = 0.4264
	Root MSE = 14.482

testscr	Coef.	Robust Std. Err.	t	P> t	[95% Conf. Interval]	
str	-1.117018	.5875135	-1.90	0.058	-2.271884	.0378468
el_pct	-.6729116	.3741231	-1.80	0.073	-1.408319	.0624958
str_el3	.0011618	.0185357	0.06	0.950	-.0352736	.0375971
_cons	686.3385	11.75935	58.37	0.000	663.2234	709.4537

The equation is  $\text{testscr}_{\text{hat}} = 686.3 - 1.12 \cdot \text{str} - 0.67 \cdot \text{el\_pct} + 0.0012 \cdot \text{str\_el3}$

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

$-1.12 + 0.0012 \cdot \text{el\_pct}$  (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.

\*The data used in this section is the "California Test Score" data set (`caschool.dta`) from chapter 6 of the book *Introduction to Econometrics* from Stock and Watson, 2003. Data can be downloaded from [http://wps.aw.com/aw\\_stock\\_ie\\_2/50/13016/3332253.cw/index.html](http://wps.aw.com/aw_stock_ie_2/50/13016/3332253.cw/index.html). For a detailed discussion please refer to the respective section in the book.

# Creating dummies

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

```
tab major, generate(major_dum)
```

```
. tab major, generate(major_dum)
```

Major	Freq.	Percent	Cum.
Econ	10	33.33	33.33
Math	10	33.33	66.67
Politics	10	33.33	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

Name	Label
city	City
state	State
gender	Gender
status	Status: grad or undergrad
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

```
. tab1 major_dum1 major_dum2 major_dum3
```

-> tabulation of major\_dum1

major==Econ	Freq.	Percent	Cum.
0	20	66.67	66.67
1	10	33.33	100.00
Total	30	100.00	

-> tabulation of major\_dum2

major==Math	Freq.	Percent	Cum.
0	20	66.67	66.67
1	10	33.33	100.00
Total	30	100.00	

-> tabulation of major\_dum3

major==Politics	Freq.	Percent	Cum.
0	20	66.67	66.67
1	10	33.33	100.00
Total	30	100.00	40



Here is another example:

```
tab agegrups, generate(agegrups_dum)
```

```
. tab agegrups, generate(agegrups_dum)
```

Age by groups	Freq.	Percent	Cum.
18 to 19	10	33.33	33.33
20 to 29	9	30.00	63.33
30 to 39	11	36.67	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

Name	Label
status	Status: grad or undergrad
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
agegrups	Age by groups
sex	Gender
major_dum1	major==Econ
major_dum2	major==Math
major_dum3	major==Politics
agegrups_dum1	agegrups==18 to 19
agegrups_dum2	agegrups==20 to 29
agegrups_dum3	agegrups==30 to 39

```
. tab1 agegrups_dum1 agegrups_dum2 agegrups_dum3
```

-> tabulation of agegrups\_dum1

agegrups==	Freq.	Percent	Cum.
18 to 19			
0	20	66.67	66.67
1	10	33.33	100.00
Total	30	100.00	

-> tabulation of agegrups\_dum2

agegrups==	Freq.	Percent	Cum.
20 to 29			
0	21	70.00	70.00
1	9	30.00	100.00
Total	30	100.00	

-> tabulation of agegrups\_dum3

agegrups==	Freq.	Percent	Cum.
30 to 39			
0	19	63.33	63.33
1	11	36.67	100.00
Total	30	100.00	

# Frequently used Stata commands

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

Type `help [command name]` in the windows command for details

Source: <http://www.ats.ucla.edu/stat/stata/notes2/commands.htm>

Basic data reporting	<b>describe</b> <b>codebook</b> <b>inspect</b> <b>list</b> <b>browse</b> <b>count</b> <b>assert</b> <b>summarize</b> <b>Table (tab)</b> <b>tabulate</b>
Data manipulation	<b>generate</b> <b>replace</b> <b>egen</b> <b>recode</b> <b>rename</b> <b>drop</b> <b>keep</b> <b>sort</b> <b>encode</b> <b>decode</b> <b>order</b> <b>by</b> <b>reshape</b>
Formatting	<b>format</b> <b>label</b>
Keeping track of your work	<b>log</b> <b>notes</b>
Convenience	<b>display</b>

## *Is my model OK? (links)*

### ***Regression diagnostics: A checklist***

<http://www.ats.ucla.edu/stat/stata/webbooks/reg/chapter2/statareg2.htm>

### ***Logistic regression diagnostics: A checklist***

<http://www.ats.ucla.edu/stat/stata/webbooks/logistic/chapter3/statalog3.htm>

### ***Times series diagnostics: A checklist (pdf)***

<http://homepages.nyu.edu/~mrg217/timeseries.pdf>

### ***Times series: dfueller test for unit roots (for R and Stata)***

<http://www.econ.uiuc.edu/~econ472/tutorial9.html>

### ***Panel data tests: heteroskedasticity and autocorrelation***

- <http://www.stata.com/support/faqs/stat/panel.html>
- <http://www.stata.com/support/faqs/stat/xtreg.html>
- <http://www.stata.com/support/faqs/stat/xt.html>
- [http://dss.princeton.edu/online\\_help/analysis/panel.htm](http://dss.princeton.edu/online_help/analysis/panel.htm)

***I can't read the output of my model!!!*** (links)

***Data Analysis: Annotated Output***

<http://www.ats.ucla.edu/stat/AnnotatedOutput/default.htm>

***Data Analysis Examples***

<http://www.ats.ucla.edu/stat/dae/>

***Regression with Stata***

<http://www.ats.ucla.edu/STAT/stata/webbooks/reg/default.htm>

***Regression***

<http://www.ats.ucla.edu/stat/stata/topics/regression.htm>

***How to interpret dummy variables in a regression***

<http://www.ats.ucla.edu/stat/Stata/webbooks/reg/chapter3/statareg3.htm>

***How to create dummies***

<http://www.stata.com/support/faqs/data/dummy.html>

<http://www.ats.ucla.edu/stat/stata/faq/dummy.htm>

***Logit output: what are the odds ratios?***

[http://www.ats.ucla.edu/stat/stata/library/odds\\_ratio\\_logistic.htm](http://www.ats.ucla.edu/stat/stata/library/odds_ratio_logistic.htm)

## ***Topics in Statistics (links)***

***What statistical analysis should I use?***

[http://www.ats.ucla.edu/stat/mult\\_pkg/whatstat/default.htm](http://www.ats.ucla.edu/stat/mult_pkg/whatstat/default.htm)

***Statnotes: Topics in Multivariate Analysis, by G. David Garson***

<http://www2.chass.ncsu.edu/garson/pa765/statnote.htm>

***Elementary Concepts in Statistics***

<http://www.statsoft.com/textbook/stathome.html>

***Introductory Statistics: Concepts, Models, and Applications***

<http://www.psychstat.missouristate.edu/introbook/sbk00.htm>

***Statistical Data Analysis***

<http://math.nicholls.edu/badie/statdataanalysis.html>

***Stata Library. Graph Examples (some may not work with STATA 10)***

<http://www.ats.ucla.edu/STAT/stata/library/GraphExamples/default.htm>

***Comparing Group Means: The T-test and One-way ANOVA Using STATA, SAS, and SPSS***

<http://www.indiana.edu/~statmath/stat/all/ttest/>

## Useful links / Recommended books

- DSS Online Training Section <http://dss.princeton.edu/training/>
- UCLA Resources to learn and use STATA <http://www.ats.ucla.edu/stat/stata/>
- DSS help-sheets for STATA [http://dss/online\\_help/stats\\_packages/stata/stata.htm](http://dss/online_help/stats_packages/stata/stata.htm)
- *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.” <http://fmwww.bc.edu/GStat/docs/StataIntro.pdf>
- STATA FAQ website <http://stata.com/support/faqs/>
- Princeton DSS Libguides <http://libguides.princeton.edu/dss>

### 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