

11. Multivariate analyses: introduction

LPO.8800: Statistical Methods in Education Research

Sean P. Corcoran

Last time

- Bivariate regression
- Prediction equation, predicted values, residuals (prediction errors)
- Ordinary least squares (OLS) - the “line of best fit”
- Interpreting regression intercept and slope
- Assessing goodness of fit (R^2)
- Conditional mean interpretation of regression
- Inference about the population slope: confidence intervals and hypothesis tests
- Regression diagnostics with residuals

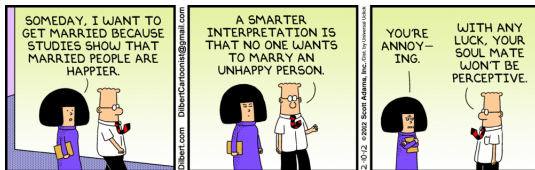
Correlation vs. causality, revisited

Generally speaking, regression slopes (and correlations) *cannot* be interpreted as *causal*. Examples:

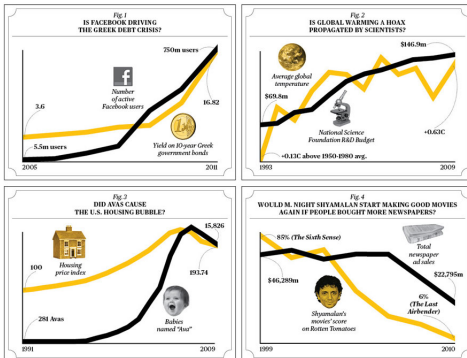
- Russian cholera epidemic: peasants observed that in communities with lots of doctors, there were lots of cholera cases; doctors were murdered.
- SAT prep courses: in 1988 Harvard interviewed its freshmen and found that those who took SAT coaching courses scored 63 points lower than those who did not.
 - ▶ A dean concluded that the SAT courses were unhelpful and that “the coaching industry is playing on parental anxiety.”

Causal questions imply “all else is held equal.”

Correlation vs. causality, revisited

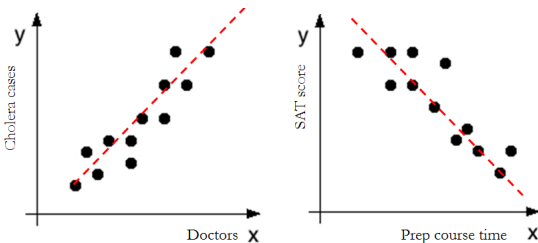


Correlation vs. causality, revisited



Correlation vs. causality, revisited

Imagine collecting data and conducting a simple regression analysis for each case:



In the lefthand figure, each data point is a community. In the righthand figure, each data point is a college applicant.

Correlation vs. causality, revisited

There is clearly an *association* between these pairs of variables, but can we say that changes in X are *causing* changes in Y ($X \rightarrow Y$)?

Not without ruling out alternative explanations. One alternative explanation is **reverse causality** ($Y \rightarrow X$). Another is the presence of one or more other factors (**confounders**) that are affecting both Y and X .

Correlation vs. causality, revisited

Considering the above two examples:

- Russian cholera epidemic: it is unlikely doctors (X) caused the cholera cases (Y), since the presence of cholera preceded the arrival of the doctors ($Y \rightarrow X$).
- SAT prep courses: it is *possible* that the prep course worsened SAT performance (the time ordering is appropriate). But it is more likely a third factor explains *both* enrollment in the prep course *and* low SAT scores (e.g., test anxiety, poor prior academic preparation). The association may be **spurious**.

Correlation vs. causality, revisited



Experimental vs. observational data

Ruling out alternative explanations can be very difficult to do in social science and education research. The researcher is typically working with *observational* data, and has no control over assignment to “treatment” conditions of interest. Consider again these questions:

- Does smoking cause lung cancer?
- Would a smaller class size improve learning?
- Does education increase labor market productivity and earnings?
- Is parental divorce detrimental to childrens' outcomes?
- Do mask mandates reduce the transmission of infectious diseases?

Experimental vs. observational data

This is in contrast to the medical researcher who can randomly assign subjects to receive a new drug or a placebo. With this study design, she can confidently attribute any systematic differences in the subjects' outcomes to the drug (and not due to some third factor).

Ruling out alternative explanations

In the absence of random assignment, the elimination of alternative explanations is difficult to do, and depends on sound research design, data availability, and a good theoretical understanding of factors that affect variation in the outcome Y .

Note: outliers and anecdotal examples of contradictory cases are **not** sufficient for ruling out causal relationships! Causal effects are a description of how X affects Y *on average*, not in a deterministic sense.

- A high-poverty school that is “beating the odds” does not demonstrate that poverty has no effect on academic achievement.
- A smoker that lives to 102 is not proof that smoking does not cause lung cancer.

Controlling for other variables

In practice how does one rule out alternative explanations for the association between X and Y ? One way is through statistical **controls**. Controlling involves using statistical techniques to find the correlation between two variables, holding the value of other variables constant.

The variables we wish to remove the effects of—i.e., control for—are called **control variables** or **covariates** (e.g., X_2, X_3, \dots, X_k).

- We statistically control for a third variable X_2 by examining the relationship between X_1 and Y *conditional on* X_2 (i.e., for fixed values of X_2).

Example 1

Does computer ownership improve 8th grade math achievement?

```
. tabstat achmat08, by(computer) stat(mean n)
```

Summary for variables: achmat08
by categories of: computer (computer owned by family in eighth grade?)

computer	mean	N
no	54.94897	263
yes	58.41321	237
Total	56.59102	500

```
. reg achmat08 computer
```

Source	SS	df	MS
Model	1496.05776	1	1496.05776
Residual	42030.8486	498	84.3992944
Total	43526.9064	499	87.2282693

Number of obs = 500
F(1, 498) = 17.73
Prob > F = 0.0000
R-squared = 0.0344
Adj R-squared = 0.0324
Root MSE = 9.1869

achmat08	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
computer	3.464233	.8228153	4.21	0.000	1.847616 5.080851
_cons	54.94897	.5664891	97.00	0.000	53.83597 56.06198

Example 1

The association between computer ownership and math achievement could be spurious, explained by a third factor correlated with computer *and* math achievement. Let's try "controlling" for SES using 2 groups (low or high):

```
. egen ses2=cut(ses), group(2)
. * the above command creates a new variable 'ses2' that splits 'ses' into two equal-sized groups (low and high)
. table computer ses2, contents(mean achmat08 n achmat08)
```

computer owned by family in eighth grade?	ses2	
	0	1
no	53.5129 169	57.53085 94
yes	55.46333 78	59.86031 159

The 3.46 point "effect" of computer ownership on math achievement is smaller after conditioning on SES. For high SES students, the "effect" is 2.32 points; for low SES students, 1.95 points.

Example 2

Do AP courses improve high school math achievement?

```
. reg achmat12 approg
```

Source	SS	df	MS	Number of obs	=	493
Model	4688.50685	1	4688.50685	F(1, 491)	=	88.17
Residual	26110.315	491	53.177831	Prob > F	=	0.0000
				R-squared	=	0.1522
				Adj R-squared	=	0.1505
Total	30798.8219	492	62.5992314	Root MSE	=	7.2923

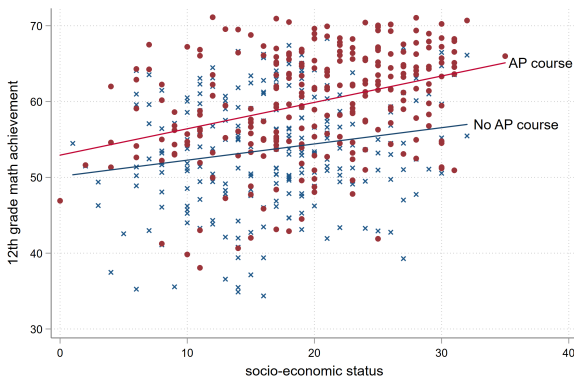
achmat12	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
approg	6.175654	.6577047	9.39	0.000	4.883391 7.467917
_cons	53.69983	.4767134	112.65	0.000	52.76318 54.63648

Example 2

How does AP course taking vary with SES (quintiles)? How does mean 12th grade math achievement vary with AP *conditional* on SES quintile?

SESquint	took AP	Mean:			Count:	
		no AP	AP	Diff	no	yes
1	0.409	52.2	55.7	3.5	55	38
2	0.385	52.3	57.1	4.9	56	35
3	0.465	55.4	59.4	4.1	53	46
4	0.610	54.3	60.4	6.1	39	61
5	0.718	55.3	62.9	7.6	31	79

Example 2



Controlling for other variables

In each of these cases we would like a single estimate that represents the average difference in 12th grade math achievement *conditional on* SES. This implies some kind of weighted average.

Multiple regression is one way of obtaining such an average.

Multiple regression

Multiple regression is one way of statistically controlling for other explanatory variables that are ignored in simple regression. Hopefully, these controls will get us closer to a causal interpretation. With 2 explanatory variables the best fit “line” is:

$$\hat{y} = a + b_1x_1 + b_2x_2$$

With k explanatory variables:

$$\hat{y} = a + b_1x_1 + b_2x_2 + \dots + b_kx_k$$

There is now an intercept and k slope coefficients to find.

Multiple regression

As before, the best fit “line” is the one where the intercept a and slope coefficients b_1, b_2, \dots, b_k minimize the sum of the squared deviations between the actual data points and the predicted values:

$$\min_{a, b} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\min_{a, b} \sum_{i=1}^n (y_i - a - b_1x_1 - b_2x_2 - \dots - b_kx_k)^2$$

Multiple regression: example 1

To implement multiple regression in Stata, continue to use the `regress` command, and include the additional explanatory variables in your variable list:

```
. regress achmat08 computer i.ses2
```

Source	SS	df	MS	Number of obs	=	500
Model	3478.87468	2	1739.43734	F(2, 497)	=	21.59
Residual	40048.0317	497	80.5795407	Prob > F	=	0.0000
Total	43526.9064	499	87.2282693	R-squared	=	0.0799
				Adj R-squared	=	0.0762
				Root MSE	=	8.9766

achmat08	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
computer	2.149567	.8465355	2.54	0.011	.4863372	3.812796
i.ses2	4.193894	.8454511	4.96	0.000	2.532795	5.854992
_cons	53.45002	.630632	84.76	0.000	52.21098	54.68905

Multiple regression: example 2

Using SES quintiles as control variables:

```
. reg achmat12 approg i.sesquint
```

Source	SS	df	MS	Number of obs	=	493
Model	6608.74334	5	1321.74867	F(5, 487)	=	26.61
Residual	24190.0785	487	49.6716191	Prob > F	=	0.0000
				R-squared	=	0.2146
				Adj R-squared	=	0.2065
Total	30798.8219	492	62.5992314	Root MSE	=	7.0478

achmat12	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
approg	5.218614	.658095	7.93	0.000	3.925558	6.51167
sesquint						
2	.6226858	1.039324	0.60	0.549	-1.419427	2.664799
3	3.341469	1.018429	3.28	0.001	1.340412	5.342526
4	3.37997	1.023907	3.30	0.001	1.368148	5.391791
5	5.52072	1.013494	5.45	0.000	3.529358	7.512081
_cons	51.49928	.7787234	66.13	0.000	49.9692	53.02935

Multiple regression: example 2

Using continuous SES as control variable:

```
. reg achmat12 approg ses
```

Source	SS	df	MS	Number of obs	=	493
Model	6558.71774	2	3279.35887	F(2, 490)	=	66.29
Residual	24240.1041	490	49.4696003	Prob > F	=	0.0000
				R-squared	=	0.2130
				Adj R-squared	=	0.2097
Total	30798.8219	492	62.5992314	Root MSE	=	7.0335

achmat12	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
approg	5.227734	.6528237	8.01	0.000	3.945055	6.510413
ses	.2895499	.047092	6.15	0.000	.1970227	.3820771
_cons	48.86905	.9103236	53.68	0.000	47.08043	50.65767

Note: this regression constrains the slope on ses to be the same for AP and non-AP students. That is, it finds the best fit regression equation where the slope is the same for these two groups.

Multiple regression: example 2

Using continuous SES as control variable:



LPO.8800 (Corcoran)

Lecture 11

Last update: November 27, 2023 25 / 51

Multiple regression: example 3

Regression of monthly wages on experience, controlling for education.

. regress wage exper						
Source	SS	df	MS	Number of obs = 935		
Model	732.242855	1	732.242855	F(1, 933) = 0.00		
Residual	152715436	933	163682.139	Prob > F = 0.9467		
Total	152716168	934	163507.675	R-squared = 0.0000		
				Adj R-squared = -0.0011		
				Root MSE = 404.58		
wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
exper	.2024031	3.026148	0.07	0.947	-5.786443 6.141249	
_cons	955.6049	37.4111	25.54	0.000	882.1853 1029.025	
. regress wage exper educ						
Source	SS	df	MS	Number of obs = 935		
Model	20747023.1	2	10373511.5	F(2, 932) = 73.26		
Residual	131969145	932	141597.795	Prob > F = 0.0000		
Total	152716168	934	163507.675	R-squared = 0.1359		
				Adj R-squared = 0.1340		
				Root MSE = 376.29		
wage	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
exper	17.63777	3.161775	5.58	0.000	11.43275 23.84279	
educ	76.21639	6.296604	12.10	0.000	63.85922 88.57355	
_cons	-272.5279	107.2627	-2.54	0.011	-483.0323 -62.02344	

LPO.8800 (Corcoran)

Lecture 11

Last update: November 27, 2023 26 / 51

Example: multiple regression

The slope coefficients are now interpreted as **marginal** or **partial** effects: the linear relationship between Y and X_1 , *conditional on* (or “holding constant”) X_2 and any other included control variables.

- Conditional on years of education (holding constant years of education), we predict that an additional year of work experience is associated with \$17.64 additional monthly earnings.
- Conditional on years of work experience (holding constant work experience), we predict that an additional year of education is associated with \$76.22 additional monthly earnings.

The *prediction equation* can be used to find the “best prediction” of Y given values of X_2, \dots, X_K .

Example: multiple regression

For example, let years of experience be $X_1 = 10$ and years of education completed be $X_2 = 14$. Our best prediction of monthly earnings is:

$$\hat{y} = -272.53 + 17.64 * 10 + 76.22 * 14 = 970.95$$

Causality, revisited

Can multiple regression coefficients can be interpreted as causal? In most cases, unfortunately not. While the regression controls for some confounders, there are likely others. This is where careful research design comes in (see later courses!)

When x_1 and x_2 are uncorrelated

When x_1 and x_2 are uncorrelated, the OLS estimators of b_1 and b_2 are:

$$\hat{b}_1 = r_{y1} \frac{s_y}{s_1}$$

$$\hat{b}_2 = r_{y2} \frac{s_y}{s_2}$$

where r_{y1} is the correlation between y and x_1 , and r_{y2} is the correlation between y and x_2 . (s_1 is the standard deviation of x_1 , and s_2 is the standard deviation of x_2). Notice these are equivalent to the formula for \hat{b} in the simple regression case. The r_{y1} and r_{y2} are sometimes called **zero-order correlations**.

When x_1 and x_2 are uncorrelated

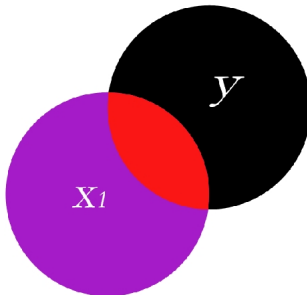
In multiple regression, R^2 can still be used as a measure of fit, interpreted in the same way: the fraction of overall variation in y that is explained by the prediction equation. When x_1 and x_2 are uncorrelated, R^2 is simply:

$$R^2 = r_{y1}^2 + r_{y2}^2$$

(the sum of the two squared zero-order correlations)

- R^2 is the *coefficient of determination*
- R —the square root of R^2 —is the *multiple correlation*

Venn diagram with one explanatory variable

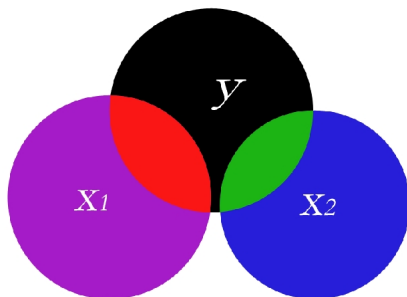


Venn diagram with one explanatory variable

The circle labeled y represents variation in y , and the circle labeled x_1 represents variation in x_1 .

- Think of the overlap (red) as variation in y “explained” by variation in x_1
- The red area represents correlation between x_1 and y : information used by the regression to estimate \hat{b}_1
- The black area is variation in y *unexplained* by variation in x_1 (“residual” variation)
- The proportion of y covered by x_1 represents the R^2

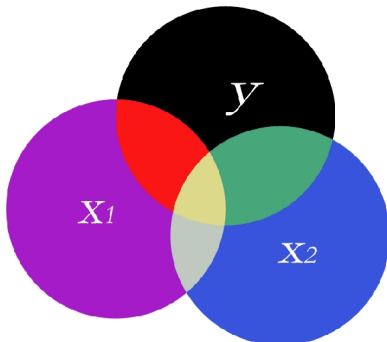
Venn diagram with two explanatory variables - 1



Venn diagram with two explanatory variables - 1

- Think of the overlap between y and x_1 (red) as variation in y “explained” by variation in x_1
- Think of the overlap between y and x_2 (green) as variation in y “explained” by variation in x_2
- x_1 and x_2 do not overlap (they are uncorrelated), so it is easy to attribute variation in y separately to x_1 and x_2
- The red area represents information used by the regression to estimate \hat{b}_1
- The green area represents information used by the regression to estimate \hat{b}_2
- The black area is variation in y unexplained by x_1 or x_2
- The proportion of y covered by x_1 and x_2 represents R^2

Venn diagram with two explanatory variables - 2



Venn diagram with two explanatory variables - 2

- In this case x_1 and x_2 overlap—they are *correlated* (represented by the yellow area), thus it is not as clear how to attribute variation in y separately to x_1 and x_2
- The red area represents the unique information used by the regression to estimate \hat{b}_1
- The green area represents the unique information used by the regression to estimate \hat{b}_2
- Both the red and green areas are *smaller* than those in example 1—we have less certainty about how much of y can be attributed to each explanatory variable

When x_1 and x_2 are correlated

When x_1 and x_2 are correlated, the OLS estimators of b_1 and b_2 can be written:

$$\hat{b}_1 = \left(\frac{r_{y1} - r_{y2}r_{12}}{1 - r_{12}^2} \right) \left(\frac{s_y}{s_1} \right)$$
$$\hat{b}_2 = \left(\frac{r_{y2} - r_{y1}r_{12}}{1 - r_{12}^2} \right) \left(\frac{s_y}{s_2} \right)$$

where r_{12} is the correlation between x_1 and x_2 (and other terms were defined previously). Notice what happens if $r_{12} = 0$ (i.e. if there is no correlation between x_1 and x_2).

Example: private school attendance and math achievement

```
. reg achmat12 private
```

Source	SS	df	MS	Number of obs	=	500
Model	665.476286	1	665.476286	F(1, 498)	=	10.92
Residual	30351.3075	498	60.9464005	Prob > F	=	0.0010
				R-squared	=	0.0215
				Adj R-squared	=	0.0195
Total	31016.7838	499	62.1578833	Root MSE	=	7.8068

achmat12	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
private	2.630139	.7959513	3.30	0.001	1.066303 4.193976
_cons	56.22278	.4058571	138.53	0.000	55.42538 57.02019

```
. reg achmat12 private ses
```

Source	SS	df	MS	Number of obs	=	500
Model	3264.62388	2	1632.31194	F(2, 497)	=	29.23
Residual	27752.1599	497	55.8393559	Prob > F	=	0.0000
				R-squared	=	0.1053
				Adj R-squared	=	0.1017
Total	31016.7838	499	62.1578833	Root MSE	=	7.4726

achmat12	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
private	.5417838	.821064	0.66	0.510	-1.071401 2.154968
ses	.3552102	.0520643	6.82	0.000	.2529169 .4575035
_cons	50.21781	.9620882	52.20	0.000	48.32755 52.10807

Example: private school attendance and math achievement

```
. corr achmat12 private ses
(obs=500)
```

	achmat12	private	ses
achmat12	1.0000		
private	0.1465	1.0000	
ses	0.3232	0.3728	1.0000

```
. summ achmat12 private ses
```

Variable	Obs	Mean	Std. Dev.
achmat12	500	56.90662	7.884027
private	500	.26	.4390735
ses	500	18.434	6.924271

$$\hat{b}_1 = \left(\frac{r_{y1} - r_{y2}r_{12}}{1 - r_{12}^2} \right) \left(\frac{s_y}{s_1} \right)$$

$$\hat{b}_1 = \left(\frac{0.1465 - 0.3232 * 0.3728}{1 - 0.3728^2} \right) \left(\frac{7.884}{0.439} \right) = 0.542$$

Example: private school attendance and math achievement

```
. corr achmat12 private ses
(obs=500)
```

	achmat12	private	ses
achmat12	1.0000		
private	0.1465	1.0000	
ses	0.3232	0.3728	1.0000


```
. summ achmat12 private ses
```

Variable	Obs	Mean	Std. Dev.
achmat12	500	56.90662	7.884027
private	500	.26	.4390735
ses	500	18.434	6.924271

$$\hat{b}_2 = \left(\frac{r_{y2} - r_{y1}r_{12}}{1 - r_{12}^2} \right) \left(\frac{s_y}{s_2} \right)$$
$$\hat{b}_2 = \left(\frac{0.3232 - 0.1465 * 0.3728}{1 - 0.3728^2} \right) \left(\frac{7.884}{6.924} \right) = 0.3552$$

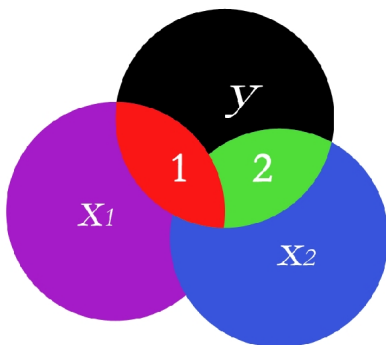
When x_1 and x_2 are correlated

With two explanatory variables the R^2 can be written as:

$$R^2 = r_{y1}^2 + r_{y2|1}^2$$

- First part: the proportion of variation in y explained by x_1
- Second part: the proportion of variation in y explained by x_2 *beyond that explained by x_1* (a *semi-partial* correlation)

When x_1 and x_2 are correlated



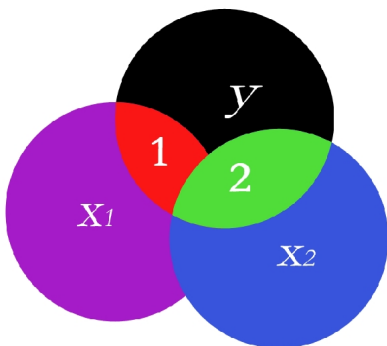
When x_1 and x_2 are correlated

Equivalently, the R^2 can be written as:

$$R^2 = r_{y2}^2 + r_{y1|2}^2$$

- First part: the proportion of variation in y explained by x_2
- Second part: the proportion of variation in y explained by x_1 *beyond that explained by x_2* (a *semi-partial correlation*)

When x_1 and x_2 are correlated



Semi-partial correlations

The correlations $r_{y2|1}^2$ and $r_{y1|2}^2$ are called *semi-partial* or *part* correlations. They represent the correlation observed between y and that part of x_1 (or x_2) that is uncorrelated with x_2 (or x_1).

Multiple regression and R^2

Some facts about multiple regression and R^2 :

- R^2 still ranges between 0 and 1
- R^2 will be high when the x 's are highly correlated with y
- R^2 will not fall below the highest R^2 with an individual x
- R^2 cannot *decrease* when additional x s are added to the regression equation
- R^2 will be larger when the explanatory variables are not *redundant*—i.e. their intercorrelation is low
- There is usually diminishing returns to additional explanatory variables (a greater chance of redundancy)

Adjusted R^2

The calculated R^2 tends to overestimate the population R^2 (it is upwardly biased). The smaller is N relative to the number of explanatory variables K , the more R^2 will be inflated. An **adjusted R^2** is often used instead:

$$R_{ADJ}^2 = 1 - (1 - R^2) \frac{(N - 1)}{N - K - 1}$$

Holding N constant, the adjusted R^2 “penalizes” you for including additional explanatory variables K .

Multicollinearity

Multicollinearity is the condition when explanatory variables in a regression are highly correlated. The consequence of this is that it becomes more difficult to discern how much of the variation in y is “due to” each individual x . This is a bigger problem the smaller the sample size.

Semi-partial (part) correlations

The **semi-partial (or part) correlation** between y and x_1 is the correlation observed between y and that part of x_1 that is uncorrelated with the other x variables.

- The *square* of the semi-partial correlation is the amount by which R^2 decreases when that explanatory variable is excluded.
- It is also the proportion of the variation in y that is explained by x_1 only
- This can be used to assess the relative importance of the explanatory variables (in terms of independent predictive power).
- Could be used to guide model specification.

Can obtain semi-partial correlations in Stata using `pcorr y x1 x2`

Semi-partial (part) correlations

Try this using the math achievement and private school example above.

- Regress *achmat08* on *ses* and *private*, note R^2 (0.1053)
- Regress *achmat08* on *ses* alone, note R^2 (0.1045)
- Regress *achmat08* on *private* alone, note R^2 (0.0215)
- Get squared semi-partial correlations `pcorr achmat08 private ses`

```
. pcorr achmat12 private ses  
(obs=500)
```

Partial and semipartial correlations of achmat12 with

Variable	Partial Corr.	Semipartial Corr.	Partial Corr.^2	Semipartial Corr.^2	Significance Value
private	0.0296	0.0280	0.0009	0.0008	0.5097
ses	0.2926	0.2895	0.0856	0.0838	0.0000