# CSSS/SOC/STAT 321 · Case-Based Statistics I

## **BIVARIATE REGRESSION**

#### Christopher Adolph

**Department of Political Science** 

and

Center for Statistics and the Social Sciences

University of Washington, Seattle

We have cross-national data from several sources:

Fertility The average number of children born per adult female, in 2000 (United Nations)

Education Ratio The ratio of girls to boys in primary and secondary education, in 2000 (Word Bank Development Indicators)

GDP per capita Economic activity in thousands of dollars, purchasing power parity in 2000 (Penn World Tables)

What are the levels of measurement of these variables?

Our question: how are these variables related to each other?

Specifically, we ask:

Specifically, we ask:

• If the level of female education changed by a certain amount, how much would we expect Fertility to change?

Specifically, we ask:

- If the level of female education changed by a certain amount, how much would we expect Fertility to change?
- If the level of GDP per capita changed by a certain amount, how much would we expect Fertility to change?

Specifically, we ask:

- If the level of female education changed by a certain amount, how much would we expect Fertility to change?
- If the level of GDP per capita changed by a certain amount, how much would we expect Fertility to change?
- How much would we expect our predictions to be off because of other random factors (noise)?

Specifically, we ask:

- If the level of female education changed by a certain amount, how much would we expect Fertility to change?
- If the level of GDP per capita changed by a certain amount, how much would we expect Fertility to change?
- How much would we expect our predictions to be off because of other random factors (noise)?
- How much would we expect our predictions to be off because of sampling variability (poor estimation)?

Answering these questions will go far towards towards answering hypotheses about relationships between variables

Review the Univariate Summary Statistics for our Example Explore the Bivariate Relationship between Fertility & Education Ratio Explore the Bivariate Relationship between Fertility & GDP per capita

Throughout, develop an understanding of linear regression

## Summary of Univariate Distribution: Fertility



## Summary of Univariate Distribution: Fertility





## Summary of Univariate Distribution: Education Ratio



## Summary of Univariate Distribution: Education Ratio



## Summary of Univariate Distribution: GDP per capita



## Summary of Univariate Distribution: GDP per capita





How would you describe the relationship between Fertility & Education Ratio?



How would you describe the relationship between Fertility & Education Ratio?

If I asked you to predict Fertility for a country not sampled, how accurate do you expect your prediction to be?







Labelling cases sometimes helps, especially for identifying outliers

What makes a point an outlier?



The best fit line is the line that passes closest to the majority of the points



The best fit line is the line that passes closest to the majority of the points If we take this line to be our model of Fertility, how

interpret it?

do we

## **Best fit lines**

Customarily, in statistics, we write the equation of a line as:

$$\gamma = \beta_0 + \beta_1 x$$

where:

- $\gamma_i$  is the dependent variable
- x is the independent variable,
- β<sub>1</sub> is the slope of the line, or the change in γ for a 1 unit change in x,
- and  $\beta_0$  is the intercept, or value of y when x = 0

# Best fit for fertility against education ratio

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$
  
$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \times \text{EduRatio}$$

The above equation is the best fit line given by linear regression

The  $\hat{\beta}$ 's are the estimated linear regression coefficients

Fertility is the fitted value, or model prediction, of the level of Fertility given the EduRatio

# Intrepreting regression coefficients

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$
  
$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \times \text{EduRatio}$$

Interpreting 
$$\hat{\beta}_1 = -0.10$$
:

Increasing EduRatio by 1 unit lowers Fertility by 0.10 units.

Because EduRatio is measured in percentage points, this means a 10% increase in female education (relative to males) will lower the number of children a woman has over her lifetime by 1 on average.

## Intrepreting regression intercepts

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$

$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \text{EduRatio}$$

Interpreting  $\hat{\beta}_0 = 12.59$ :

If EduRatio is 0, Fertility will be 12.59.

If there are no girls in primary or secondary education, then women are expected to have 12.59 children on average over their lifetimes.

Can we trust this prediction?

## Intrepreting regression intercepts

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$

$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \text{EduRatio}$$

Interpreting  $\hat{\beta}_0 = 12.59$ :

If EduRatio is 0, Fertility will be 12.59.

If there are no girls in primary or secondary education, then women are expected to have 12.59 children on average over their lifetimes.

Can we trust this prediction? No.

No country has 0 female education, so this is an extrapolation from the model.

# Using regression coefficients to predict specific cases

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$

$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \text{EduRatio}$$

How many children do we expect women to bear if as girls they receive half the education boys do?

If EduRatio is 50, Fertility will be  $12.59 - 0.10 \times 50 = 7.59$ .

How many children do we expect women to bear if as girls they receive the same education boys do?

If EduRatio is 100, Fertility will be  $12.59 - 0.10 \times 100 = 2.59$ .

# Using regression coefficients to predict specific cases

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$

$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \text{EduRatio}$$

If EduRatio is 100, Fertility will be  $12.59 - 0.10 \times 100 = 2.59$ .

Does this hold exactly for any country with education parity?

# Using regression coefficients to predict specific cases

$$\widehat{\text{Fertility}} = \hat{\beta}_0 + \hat{\beta}_1 \text{EduRatio}$$

$$\widehat{\text{Fertility}} = 12.59 - 0.10 \times \text{EduRatio}$$

If EduRatio is 100, Fertility will be  $12.59 - 0.10 \times 100 = 2.59$ .

Does this hold exactly for any country with education parity?

No. It holds on average. In any specific case *i*, there is some error between the expected and actual levels of Fertility

$$\gamma_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$

To account for the random deviation of each case from the underlying trend, we add an error term,  $\varepsilon_i$ .

We will assume our  $y_i$ 's follow the above model

That is, we will assume there is some "true"  $\beta_0$  and  $\beta_1$  which generated the  $\gamma_i$  we observe, and some "true" error from this tendency

# The linear regression model

$$\hat{\gamma}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\varepsilon}_i$$

When we estimate this model,

we designate the estimated parameters by adding "hats" to them

The estimates  $(\hat{\beta}_0, \hat{\beta}_1, \hat{\varepsilon}_i)$  probably differ from the (usually unknown) true values  $(\beta_0, \beta_1, \varepsilon_i)$ 

To emphasize this, we will call  $\hat{arepsilon}_i$  the residual, because it is not the true error, but only an estimate

# Estimating linear regression coefficients

$$\hat{\gamma}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i + \hat{\varepsilon}_i$$

How do we obtain our estimates of the  $\beta$ 's?

The full details are beyond the scope of 321

Key assumption is that  $\varepsilon_i$  is Normally distributed:

 $\varepsilon_i \sim \operatorname{Normal}(0, \sigma^2)$ 



(Source: Larry Gonick & Wollcott Smith, The Cartoon Guide to Statistics)

The distribution of  $\varepsilon_i$  determines how closely or widely the  $\gamma_i$ 's are spaced around the best fit line

Our key simplifying assumption is that everywhere around the line, the  $\gamma_i$ 's are spread with the same Normal distribution



### Perhaps the line that minimizes the total residuals?

Perhaps the line that minimizes the total residuals?

But some residuals are positive, and others negative - their sum is always 0

Perhaps the line that minimizes the total residuals?

But some residuals are positive, and others negative - their sum is always 0

So lets minimize the sum of squared error!

Linear regression is fitted using the least squares procedure


It's an aggregate measure of how much the line's "predicted  $y_i,$  " or  $\hat{y}_i,$  differ from the actual data values  $y_i.$ 

(Source: Larry Gonick & Wollcott Smith, The Cartoon Guide to Statistics)

The least squares estimates are the  $\hat{eta}'$ s that minimize the total area of the above

#### squares



It's an aggregate measure of how much the line's "predicted  $y_i$ ," or  $\hat{y}_i$  , differ from the actual data values  $y_i$ .

(Source: Larry Gonick & Wollcott Smith, The Cartoon Guide to Statistics)

Stata can find these  $\hat{\beta}$ 's easily using the reg command, and R using the function lm()

#### Residuals

Notice the distinction between what we explain and what is left unexplained

LET'S QUANTIFY THIS BY APPORTIONING THE VARIABILITY IN Y. REFER TO THE PICTURE AT RIGHT FOR GUIDANCE. WE LET

$$\widehat{y}_i = a + b x_i$$

Thus,  $\hat{y}_i$  are the predicted weights determined by the regression line.



(Source: Larry Gonick & Wollcott Smith, The Cartoon Guide to Statistics)

The total variation in  $\gamma_i$  is its total variance from the mean  $\bar{\gamma}$ , or  $\sum_{i=1}^n (\gamma_i - \bar{\gamma})^2$ 

Using least squares, we can break down the variance in  $y_i$  into two components:

Sum of square errors (SSE)  $\sum_{i=1}^{n} (y_i - \hat{y}_i)^2$ 

The total variation in  $\gamma_i$  is its total variance from the mean  $\bar{\gamma}$ , or  $\sum_{i=1}^n (\gamma_i - \bar{\gamma})^2$ 

Using least squares, we can break down the variance in  $y_i$  into two components:

 $\begin{array}{lll} \text{Sum of square errors (SSE)} & \sum_{i=1}^{n}(\gamma_i - \hat{\gamma}_i)^2 \\ \text{Regression sum of squares (RSS)} & \sum_{i=1}^{n}(\hat{\gamma}_i - \overline{\gamma})^2 \end{array}$ 

The total variation in  $\gamma_i$  is its total variance from the mean  $\bar{\gamma}$ , or  $\sum_{i=1}^n (\gamma_i - \bar{\gamma})^2$ 

Using least squares, we can break down the variance in  $y_i$  into two components:

Sum of square errors (SSE)	$\sum_{i=1}^{n} (\gamma_i - \hat{\gamma}_i)^2$
Regression sum of squares (RSS)	$\sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2$
Total sum of squares (TSS)	$\sum_{i=1}^{n} (\gamma_i - \bar{\gamma})^2$

The Regression sum of squares (RSS) is what we have explained

The Sum of squared errors (SSE) is what is left unexplained

## Analysis of variance

The Sum of squared errors is what is left unexplained:

$$\sum_{i=1}^{n} (\gamma_i - \hat{\gamma}_i)^2$$

#### Analysis of variance

The Sum of squared errors is what is left unexplained:

$$\sum_{i=1}^{n} (\gamma_i - \hat{\gamma}_i)^2 = \sum_{i=1}^{n} \hat{\varepsilon}_i^2$$

A very useful summary of this is the square root of the mean squared error:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}$$

## Analysis of variance

The Sum of squared errors is what is left unexplained:

$$\sum_{i=1}^{n} (\gamma_i - \hat{\gamma}_i)^2 = \sum_{i=1}^{n} \hat{\varepsilon}_i^2$$

A very useful summary of this is the square root of the mean squared error:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (\gamma_i - \hat{\gamma}_i)^2}$$

This is how much a prediction from this regression will differ from the true  $\gamma_i$  on average

Also known as the standard error of the regression







for the regression of Fertility on Education Ratio This line minimizes the squared deviations on the dependent variable

The residuals



The smaller the sum of squared residuals, the better the model fits the data.



The smaller the sum of squared residuals, the better the model fits the data.

The quality of model fit is a separate issue from the substantive strength of the relationship, which is given by  $\beta$ , or the change in  $\gamma$  for a one unit change in x Our model is captured in the  $\beta ' {\rm s},$  or regression coefficients. In contrast to...

The correlation coefficient *r*, a goodness of fit measure; larger values imply better fit of the model to the data

In our example, r between Fertility and Education Ratio is -0.75

Substantively, this number is hard to interpret

```
(What's a "big" r? A "small" r? Arbitrary)
```

# The coefficient of determination, $R^2$

One easy to interpret goodness of fit measure is  $R^2$ , known as the coefficient of determination

In general,  $R^2$  is the ratio of the variance the model explains to the total variance:

$$R^2 = rac{\mathrm{RSS}}{\mathrm{TSS}} = 1 - rac{\mathrm{SSE}}{\mathrm{TSS}}$$

In bivariate regression only,  $R^2$  is also the square of  $r_{X,Y}$ 

In our example,  $R^2 = 0.56$ , which says that Education Ratio "explains" 56% of the variation in Fertility, and vice versa

 $R^2$  is a proportional reduction in error (PRE) statistic



I prefer a more tangible measure of goodness of fit, the root mean squared error (RMSE).



I prefer a more tangible measure of goodness of fit, the root mean squared error (RMSE).

RMSE is "how much your model predictions miss by":



I prefer a more tangible measure of goodness of fit, the root mean squared error (RMSE).

RMSE is "how much your model predictions miss by":

here, 1.12 children per female



RMSE is better than  $R^2$  because it can be compared across models and datasets –  $R^2$  can't.



RMSE is better than  $R^2$  because it can be compared across models and datasets –  $R^2$  can't.

A question: we assumed the errors would be Normal – are they?



Recall that linear regression assumes the  $\varepsilon_i$ 's are Normally distributed.





When estimating a mean or difference of means, we worried that by chance, our sample might not reflect the population

That's a worry in linear regression as well

Does  $\hat{\beta}$  estimated from our sample reflect the true population  $\beta$ ?

Or did we get an unusual result due to sampling variability?

 $\mathrm{se}(\hat{\beta})$  is the amount we expect to miss the population  $\beta$  by on average over regression using repeated samples

 $\mathrm{se}(\hat{\beta})$  is the amount we expect to miss the population  $\beta$  by on average over regression using repeated samples

Remarkably, as  $N \to \infty$ , the  $\hat{\beta}$ 's become Normally distributed, no matter what distribution  $\gamma_i$  follows

 $\mathrm{se}(\hat{\beta})$  is the amount we expect to miss the population  $\beta$  by on average over regression using repeated samples

Remarkably, as  $N o \infty$ , the  $\hat{eta}$ 's become Normally distributed, no matter what distribution  $\gamma_i$  follows

So we can use a *t*-test to see if our  $\hat{\beta}$ 's would differ from the null hypothesis purely by chance sample from the population

 $\mathrm{se}(\hat{\beta})$  is the amount we expect to miss the population  $\beta$  by on average over regression using repeated samples

Remarkably, as  $N o \infty$ , the  $\hat{eta}$ 's become Normally distributed, no matter what distribution  $\gamma_i$  follows

So we can use a *t*-test to see if our  $\hat{\beta}$ 's would differ from the null hypothesis purely by chance sample from the population

Often, we will consider the null hypothesis to be  $\beta^{null} = 0$ , but sometimes we might want a different null

We can also construct confidence intervals around  $\hat{eta}_0$  and  $\hat{eta}_1$ 

These CIs reflect the uncertainty created by randomly sampling our data from the population

In 95% of samples, the true population eta's should lie in the reported 95% confidence intervals

If we have a lot of data, these intervals will be roughly  $\pm 2$  standard errors,

but in presenting our own results, we should lookup the correct critical t as in past examples



The standard errors of  $\hat{\beta}$  reflect the fact that in 95% of randomly sampled datasets, the true best fit line for the population lies within range of the estimated line



The standard errors of  $\hat{\beta}$  reflect the fact that in 95% of randomly sampled datasets, the true best fit line for the population lies within range of the estimated line

We can capture this "wiggle room" graphically



Why don't 95% of the observations lie inside this interval?



Why don't 95% of the observations lie inside this interval?

Because of fundametal uncertainty (statistical "noise"), measured by RMSE



Why don't 95% of the observations lie inside this interval?

Because of fundametal uncertainty (statistical "noise"), measured by RMSE

The CIs just report uncertainty in the best fit line, not in the data itself

Regression	of Fertility on Ed	ucation Rat	io			
	Variable	Estimates	se	t-stat	p-value	
	Intercept	12.59	(0.75)	16.75	< 0.001	
	Education Ratio	-0.10	(0.01)	-12.71	<0.001	
	Ν	130				
	$R^2$	0.56				
	RMSE	1.12				

The most common presentation of a linear regression is the above table

Usually, graphics are more informative and easier to read, but older articles rely heavily on this tabular format

Understanding these tables will be important for the final exam. Let's take this one apart
## A standard regression table

Regression	of Fertility on Ec	lucation Rat	io		
	Variable	Estimates	se	t-stat	p-value
	Intercept	12.59	(0.75)	16.75	<0.001
	Education Ratio	-0.10	(0.01)	-12.71	<0.001
	Ν	130			
	$R^2$	0.56			
	RMSE	1.12			

The top of the table contains important quantities regarding our independent variable(s):

- Estimates: the  $\hat{eta}$ 's, or regression coefficients
- se: the standard errors of  $\hat{eta}$
- t-stat: the t-statistic for the regression coefficient, or  $\hat{eta}/{
  m se}(\hat{eta})$
- p-value: the probability of seeing such a large t-stat by chance

Regression o	f Fertility on GDI	P per capita	I	
			95% Confic	lence Interval
	Variable	Estimates	Lower	Upper
	Intercept	12.59	[11.11 ,	14.08]
	Education Ratio	-0.10	[-0.12 ,	-0.08]
	Ν	130		
	$R^2$	0.36		
	RMSE	1.35		

Just as will our other estimates, we can construct confidence intervals around our eta's

Our results show 95% confidence that a 1 unit (1%) increase in education of girls relative to boys lowers fertility by between 0.08 and 0.12 children per woman

We would only expect the truth to lie outside the reported confidence interval in 1 of 20 random samples When we considered the relationship of female education and fertility, we also hypothesized an effect of GDP per capita

We suspected this might be an indirect effect, flowing through female education

Can we use regression to check for an effect of GDP?



















This is the least squares fit (What does that mean?)



This is the least squares fit (What does that mean?) How good

does this fit look?



Can you imagine an alternative model that would reduce the sum of squared residuals further?



Can you imagine an alternative model that would reduce the sum of squared residuals further? Perhaps a concave curve?







How do we interpret this 95% confidence interval?



How do we interpret this 95% confidence interval?

Why don't 95% of the points lie inside it?

Regressi	on of Fertility on GDP	per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	Ν	130				
	$R^2$	0.36				
	RMSE	1.35				

How do we interpret this table?

Regressi	on of Fertility on GDP	) per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	Ν	130				
	$R^2$	0.36				
	RMSE	1.35				

O How much do we expect Fertility to change when we increase GDP by \$1000?

Regressi	on of Fertility on GDF	Per capita			
	Variable	Estimates	se	t-stat	p-value
	Intercept	4.13	(0.17)	24.57	< 0.001
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001
	N	130			
	$R^2$	0.36			
	RMSE	1.35			

- How much do we expect Fertility to change when we increase GDP by \$1000? decrease by 0.1 children
- What would Fertility be if GDP were \$1000? \$10,000? \$30,000?

Regressi	on of Fertility on GDF	Per capita			
	Variable	Estimates	se	t-stat	p-value
	Intercept	4.13	(0.17)	24.57	< 0.001
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001
	N	130			
	$R^2$	0.36			
	RMSE	1.35			

- How much do we expect Fertility to change when we increase GDP by \$1000? decrease by 0.1 children
- What would Fertility be if GDP were \$1000? \$10,000? \$30,000? 4.03, 3.13, and 1.13, respectively.
- What would Fertility be if GDP were 0? Do you trust this estimate?

Regressi	on of Fertility on GDF	Per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	N	130				
	$R^2$	0.36				
	RMSE	1.35				

- How much do we expect Fertility to change when we increase GDP by \$1000? decrease by 0.1 children
- What would Fertility be if GDP were \$1000? \$10,000? \$30,000?
   4.03, 3.13, and 1.13, respectively.
- What would Fertility be if GDP were 0? Do you trust this estimate?
   4.13. No this is an extrapolation.

Regressi	on of Fertility on GDF	Per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	N	130				
	$R^2$	0.36				
	RMSE	1.35				

Suppose we drew another sample of countries. Would we expect to see a GDP different from zero in that case?

Regressi	on of Fertility on GDF	Per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	N	130				
	$R^2$	0.36				
	RMSE	1.35				

- Suppose we drew another sample of countries. Would we expect to see a GDP different from zero in that case? Yes.
- Why?

Regressi	on of Fertility on GDF	Per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	Ν	130				
	$R^2$	0.36				
	RMSE	1.35				

- Suppose we drew another sample of countries. Would we expect to see a GDP different from zero in that case? Yes.
- Ø Why?

The se is small relative to  $\hat{\beta}$ , so the true  $\beta$  is probably far from 0.

• How likely is it that we would see a t statistic this large if  $\beta = 0$ ?

Regressi	on of Fertility on GDF	Per capita				
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001	
	Ν	130				
	$R^2$	0.36				
	RMSE	1.35				

- Suppose we drew another sample of countries. Would we expect to see a GDP different from zero in that case? Yes.
- Ø Why?

The se is small relative to  $\hat{\beta}$ , so the true  $\beta$  is probably far from 0.

How likely is it that we would see a t statistic this large if β = 0?
 Very unlikely - less than 1 in 1000 samples.

Regression	n of Fertility on GDP p	oer capita			
			95% Confid	ence Interval	
	Variable	Estimates	Lower	Upper	
	Intercept	4.13	[3.80,	4.46]	
	GDP per capita (\$k)	-0.10	[-0.12 ,	-0.08]	
	N	130			
	$R^2$	0.36			
	RMSE	1.35			



What do these confidence intervals mean?

Regression of Fertility on GDP per capita							
		95% Confidence Interval					
	Variable	Estimates	Lower	Upper			
	Intercept	4.13	[3.80 ,	4.46]			
	GDP per capita (\$k)	-0.10	[-0.12 ,	-0.08]			
	Ν	130					
	$R^2$	0.36					
	RMSE	1.35					

What do these confidence intervals mean?
 In 95% of random samples, the true β's will lie inside the reported confidence intervals

Regression of Fertility on GDP per capita						
	Variable	Estimates	se	t-stat	p-value	
	Intercept	4.13	(0.17)	24.57	< 0.001	
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	< 0.001	
	Ν	130				
	$R^2$	0.36				
	RMSE	1.35				

• How much of the variance in Fertility does this model explain?

Regressi	ion of Fertility on GDI	P per capita			
	Variable	Estimates	se	t-stat	p-value
	Intercept	4.13	(0.17)	24.57	< 0.001
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001
	N	130			
	$R^2$	0.36			
	RMSE	1.35			

- How much of the variance in Fertility does this model explain?
   36 percent
- When using the model to predict fertility for a specific country, how much does it miss by on average?

Regressi	on of Fertility on GDF	Per capita			
	Variable	Estimates	se	t-stat	p-value
	Intercept	4.13	(0.17)	24.57	< 0.001
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001
	N	130			
	$R^2$	0.36			
	RMSE	1.35			

- How much of the variance in Fertility does this model explain?
   36 percent
- When using the model to predict fertility for a specific country, how much does it miss by on average? 1.35
- O How many cases were used in this analysis?

Regressi	on of Fertility on GDF	Per capita			
	Variable	Estimates	se	t-stat	p-value
	Intercept	4.13	(0.17)	24.57	< 0.001
	GDP per capita (\$k)	-0.10	(0.01)	-8.44	<0.001
	N	130			
	$R^2$	0.36			
	RMSE	1.35			

- How much of the variance in Fertility does this model explain?
   36 percent
- When using the model to predict fertility for a specific country, how much does it miss by on average? 1.35
- O How many cases were used in this analysis? 130

How do we reconcile our two sets of results?

Which model, if any, is right?

To solve this conundrum, we need *multiple regression*: A method for regressing a dependent variable on several independent variables at once

Then, at last, we can say something about confounders

Fortunately, all of today's concepts will carry over to multiple regression

## Important linear regression concepts

Regression coefficient	β
Estimate of regression coefficient	$\hat{oldsymbol{eta}}$
Standard error of est. of reg. coef.	$\operatorname{se}(\hat{eta})$
Fitted values	$\hat{\gamma}_i$
Regression errors	$\varepsilon_i$
Residuals	$arepsilon_i$
Coefficient of determination	$R^2$
Sum of squared errors (SSE)	$\sum_{i=1}^{n} \varepsilon_i$
Regression sum of squares (SSR)	$\sum_{i=1}^{n} \hat{y}_i - \bar{y}$