

# Gov 2000 - 8. Regression with Two Independent Variables

Matthew Blackwell

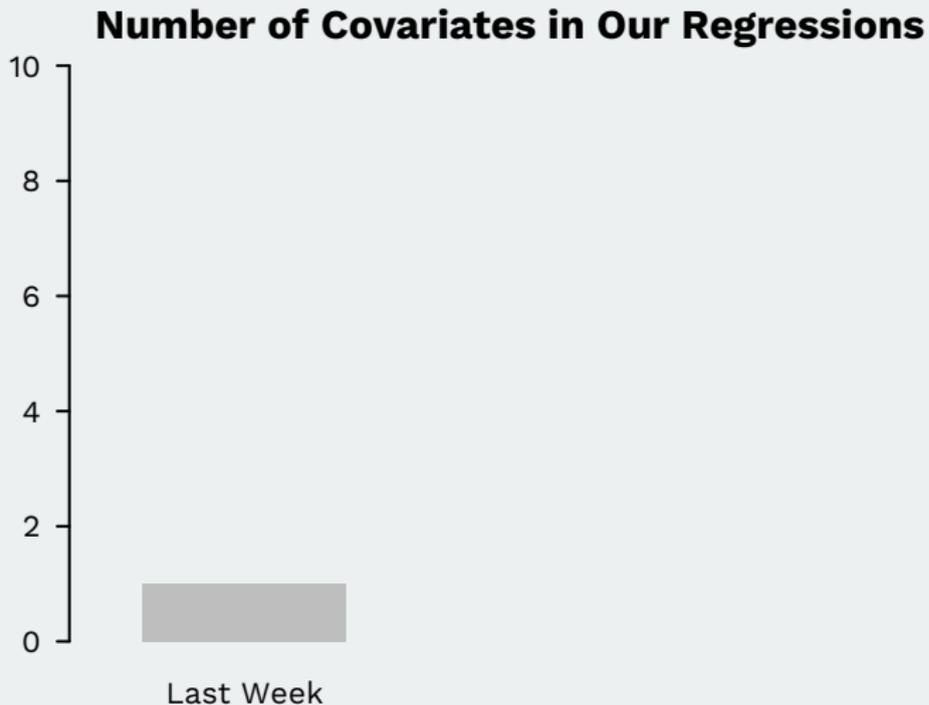
November 3, 2015

1. Why add variables to the regression?
2. Adding a binary variable
3. Adding a continuous variable
4. OLS mechanics with 2 variables
5. OLS assumptions & inference with 2 variables
6. Omitted Variable Bias
7. Multicollinearity

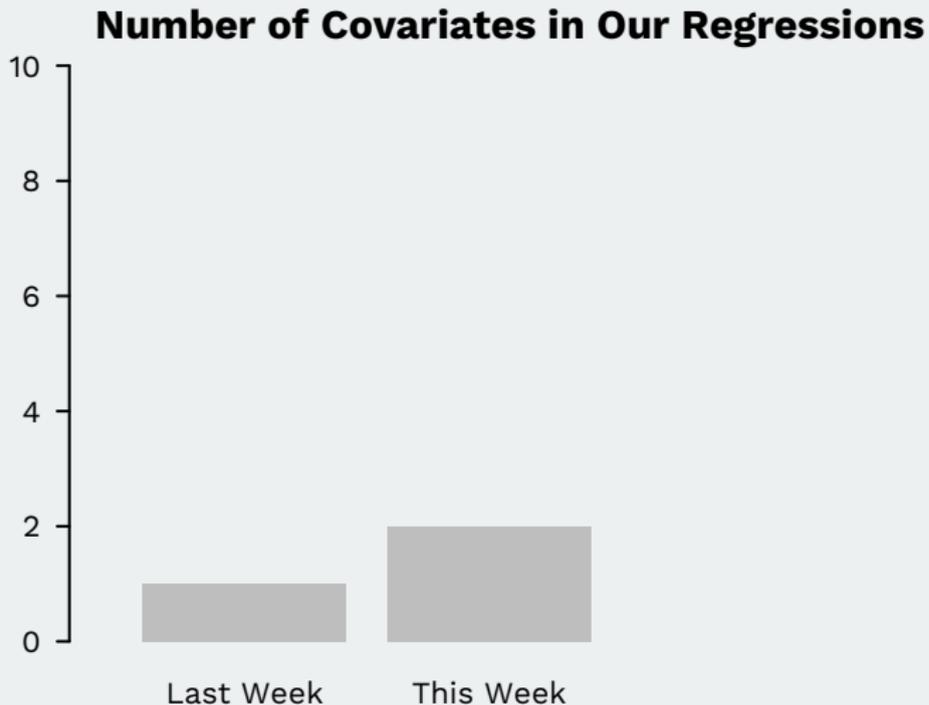
# Announcements

- Midterm:
  - ▶ Mean: 23 out of 30 (rescaled for 1000/2000)
  - ▶ SD: 5.84
  - ▶ Excellent job! We're really happy with scores!
  - ▶ Grades coming this week
- Matrix algebra (Wooldridge, Appendix D)

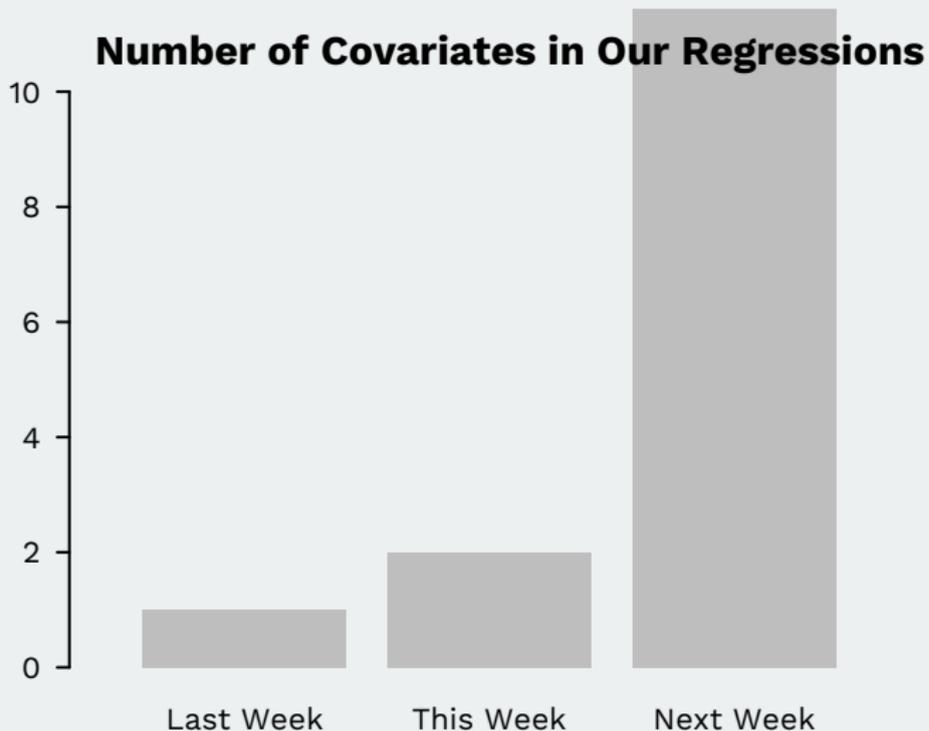
# Where are we? Where are we going?



# Where are we? Where are we going?



# Where are we? Where are we going?



**1/** Why add variables to the regression?



# Berkeley gender bias

- Graduate admissions data from Berkeley, 1973
- Acceptance rates:
  - ▶ Men: 8442 applicants, 44% admission rate
  - ▶ Women: 4321 applicants, 35% admission rate
- Evidence of discrimination toward women in admissions?
- This is a **marginal relationship**
- What about the **conditional relationship** within departments?

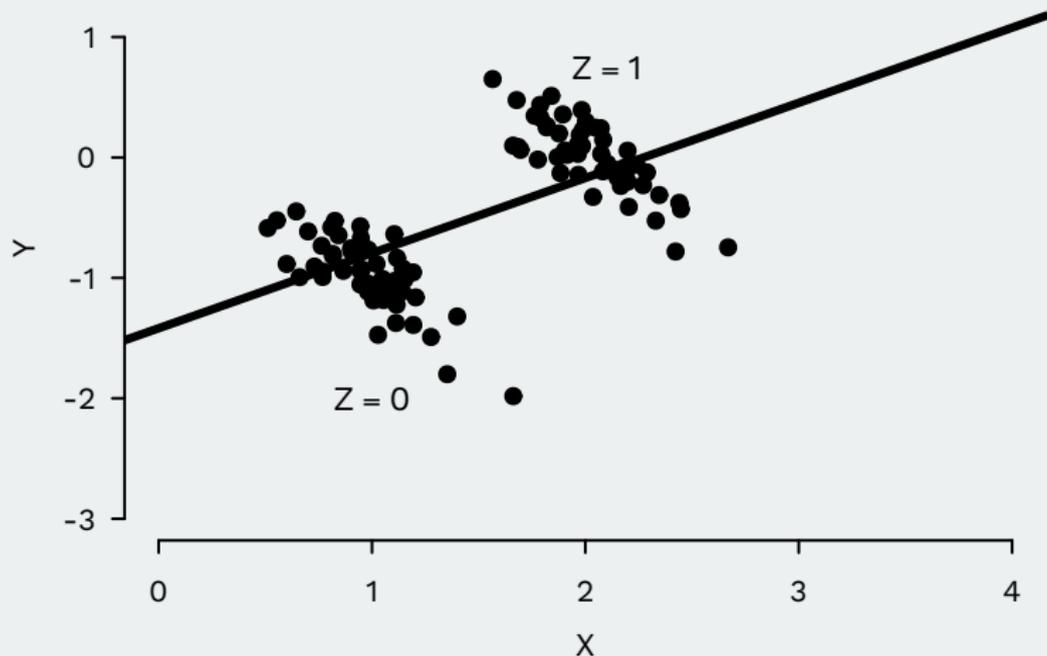
# Berkeley gender bias, II

- Within departments:

Dept	Men		Women	
	Applied	Admitted	Applied	Admitted
A	825	62%	108	82%
B	560	63%	25	68%
C	325	37%	593	34%
D	417	33%	375	35%
E	191	28%	393	24%
F	373	6%	341	7%

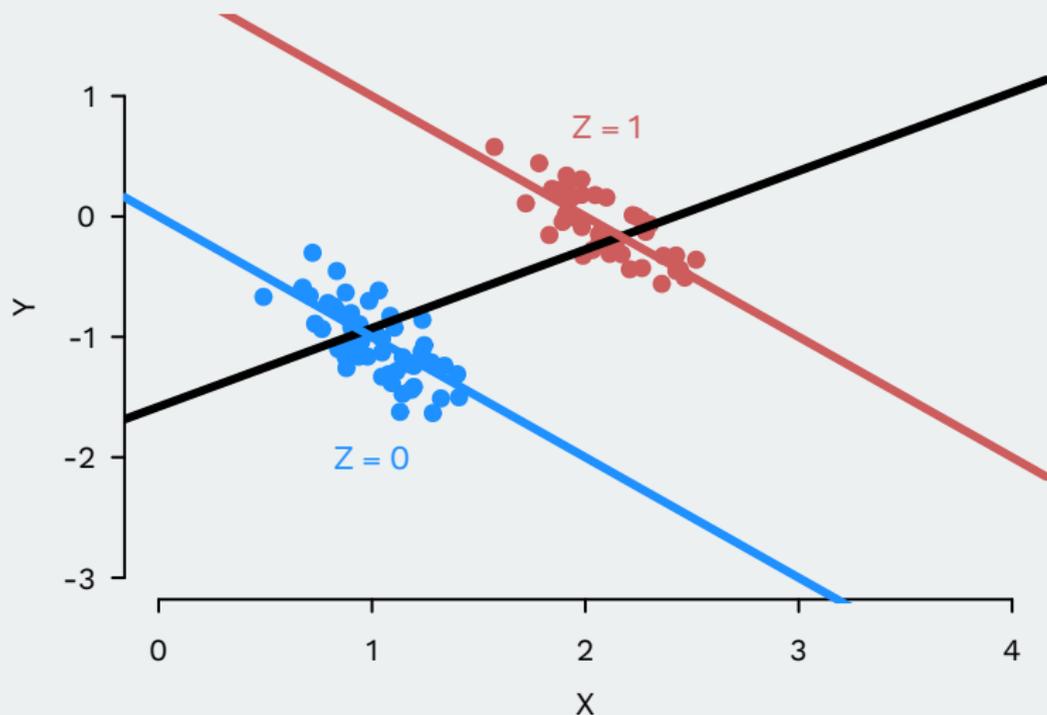
- Within departments, women do somewhat better than men!
- Women apply to more challenging departments.
- Marginal relationships (admissions and gender)  $\neq$  conditional relationship given third variable (department)

# Simpson's paradox



- Overall a positive relationship between  $Y_i$  and  $X_i$  here

# Simpson's paradox



- Overall a positive relationship between  $Y_i$  and  $X_i$  here
- But within levels of  $Z_i$ , the opposite

# Basic idea

- **Old goal:** estimate the mean of  $Y$  as a function of some independent variable,  $X$ :

$$\mathbb{E}[Y_i|X_i]$$

- For continuous  $X$ 's, we modeled the CEF/regression function with a line:

$$Y_i = \beta_0 + \beta_1 X_i + u_i$$

- **New goal:** estimate the relationship of two variables,  $Y_i$  and  $X_i$ , conditional on a third variable,  $Z_i$ :

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i + u_i$$

- $\beta$ 's are the population parameters we want to estimate

# Why control for another variable

- Descriptive
  - ▶ Get a sense for the relationships in the data.
  - ▶ Conditional on the number of steps I've taken, does higher activity levels correlate with less weight?
- Predictive
  - ▶ We can usually make better predictions about the dependent variable with more information on independent variables.
- Causal
  - ▶ Block potential **confounding**, which is when  $X$  doesn't cause  $Y$ , but only appears to because a third variable  $Z$  causally affects both of them.
  - ▶  $X_i$ : ice cream sales on day  $i$
  - ▶  $Y_i$ : drowning deaths on day  $i$
  - ▶  $Z_i$ : ??

# Plan of attack

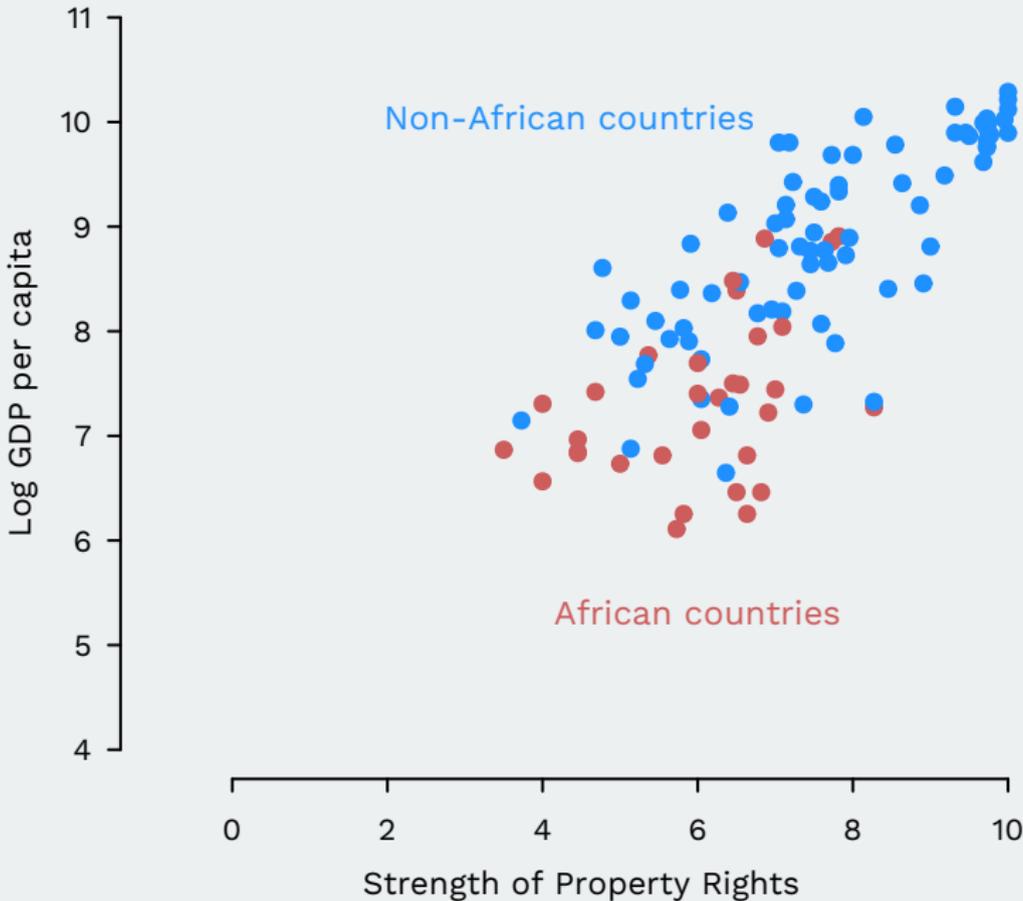
1. Adding a binary  $Z_i$
2. Adding a continuous  $Z_i$
3. Mechanics of OLS with 2 covariates
4. OLS assumptions with 2 covariates:
  - ▶ Omitted variable bias
  - ▶ Multicollinearity

# What we won't cover in lecture

1. The OLS formulas for 2 covariates
2. Proofs
3. The second covariate being a function of the first:  $Z_i = X_i^2$
4. Goodness of fit

## **2/** Adding a binary variable

# Example



# Basics

- Ye olde model:

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i$$

- $Z_i = 1$  to indicate that  $i$  is an African country
- $Z_i = 0$  to indicate that  $i$  is a non-African country
- Concern: AJR might be picking up an “African effect”:
  - ▶ African countries might have low incomes and weak property rights
  - ▶ “Control for” country being in Africa or not to remove this
  - ▶ Effects are now within Africa or within non-Africa, not between

- New model:

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i$$

# AJR model

- Let's see an example with the AJR data:

```
ajr.mod <- lm(logpgp95 ~ avexpr + africa, data = ajr)
summary(ajr.mod)
```

```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5.6556     0.3134  18.04 <2e-16 ***
## avexpr         0.4242     0.0397  10.68 <2e-16 ***
## africa        -0.8784     0.1471  -5.97  3e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.625 on 108 degrees of freedom
## (52 observations deleted due to missingness)
## Multiple R-squared:  0.708, Adjusted R-squared:  0.702
## F-statistic: 131 on 2 and 108 DF, p-value: <2e-16
```

# Two lines, one regression

- How can we interpret this model?
- Plug in two possible values for  $Z_i$  and rearrange
- When  $Z_i = 0$ :

$$\begin{aligned}\widehat{Y}_i &= \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i \\ &= \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 \times 0 \\ &= \widehat{\beta}_0 + \widehat{\beta}_1 X_i\end{aligned}$$

- When  $Z_i = 1$ :

$$\begin{aligned}\widehat{Y}_i &= \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i \\ &= \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 \times 1 \\ &= (\widehat{\beta}_0 + \widehat{\beta}_2) + \widehat{\beta}_1 X_i\end{aligned}$$

- Two different intercepts, same slope

# Interpretation of the coefficients

- Let's review what we've seen so far:

	Intercept for $X_i$	Slope for $X_i$
Non-African country ( $Z_i = 0$ )	$\hat{\beta}_0$	$\hat{\beta}_1$
African country ( $Z_i = 1$ )	$\hat{\beta}_0 + \hat{\beta}_2$	$\hat{\beta}_1$

- In this example, we have:

$$\hat{Y}_i = 5.656 + 0.424 \times X_i - 0.878 \times Z_i$$

- We can read these as:

- $\hat{\beta}_0$ : average log income for non-African country ( $Z_i = 0$ ) with property rights measured at 0 is 5.656
- $\hat{\beta}_1$ : A one-unit increase in property rights is associated with a 0.424 increase in average log incomes for two African countries (or for two non-African countries)
- $\hat{\beta}_2$ : there is a -0.878 average difference in log income per capita between African and non-African counties conditional on property rights

# General interpretation of the coefficients

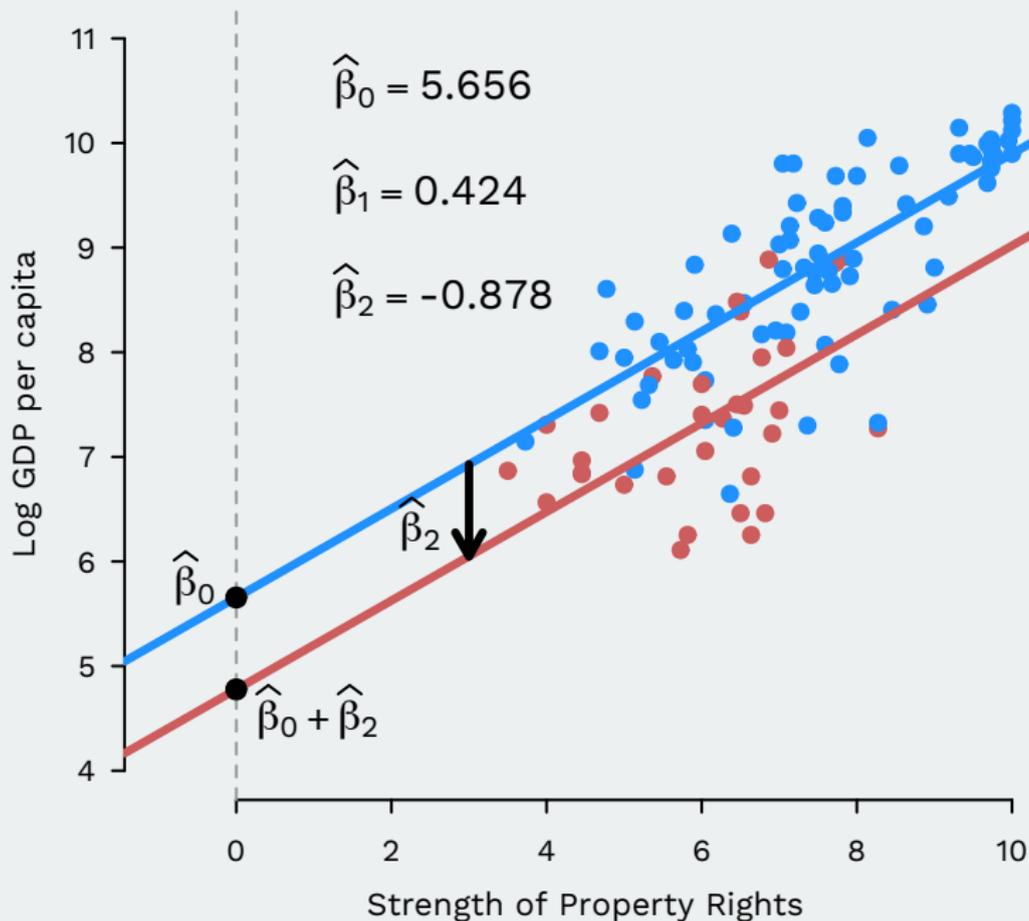
$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i$$

- $\widehat{\beta}_0$ : average value of  $Y_i$  when both  $X_i$  and  $Z_i$  are equal to 0
- $\widehat{\beta}_1$ : A one-unit increase in  $X_i$  is associated with a  $\widehat{\beta}_1$ -unit change in  $Y_i$  **conditional on**  $Z_i$
- $\widehat{\beta}_2$ : average difference in  $Y_i$  between  $Z_i = 1$  group and  $Z_i = 0$  group **conditional on**  $X_i$

# Adding a binary variable, visually



# Adding a binary variable, visually



# **3/** Adding a continuous variable

# Adding a continuous variable

- Ye olde model:

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i$$

- $Z_i$ : mean temperature in country  $i$  (continuous)
- Concern: geography is confounding the effect
  - ▶ geography might affect political institutions
  - ▶ geography might affect average incomes (through diseases like malaria)
- New model:

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i$$

# AJR model, revisited

```
ajr.mod2 <- lm(logpgp95 ~ avexpr + meantemp, data = ajr)
summary(ajr.mod2)
```

```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   6.8063     0.7518   9.05 1.3e-12 ***
## avexpr        0.4057     0.0640   6.34 3.9e-08 ***
## meantemp     -0.0602     0.0194  -3.11 0.003 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.643 on 57 degrees of freedom
## (103 observations deleted due to missingness)
## Multiple R-squared:  0.615, Adjusted R-squared:  0.602
## F-statistic: 45.6 on 2 and 57 DF, p-value: 1.48e-12
```

# Interpretation with a continuous Z

	Intercept for $X_i$	Slope for $X_i$
$Z_i = 0^\circ\text{C}$	$\widehat{\beta}_0$	$\widehat{\beta}_1$
$Z_i = 21^\circ\text{C}$	$\widehat{\beta}_0 + \widehat{\beta}_2 \times 21$	$\widehat{\beta}_1$
$Z_i = 24^\circ\text{C}$	$\widehat{\beta}_0 + \widehat{\beta}_2 \times 24$	$\widehat{\beta}_1$
$Z_i = 26^\circ\text{C}$	$\widehat{\beta}_0 + \widehat{\beta}_2 \times 26$	$\widehat{\beta}_1$

- In this example we have:

$$\widehat{Y}_i = 6.806 + 0.406 \times X_i - 0.06 \times Z_i$$

- $\widehat{\beta}_0$ : average log income for a country with property rights measured at 0 and a mean temperature of 0 is 6.806
- $\widehat{\beta}_1$ : A one-unit increase in property rights is associated with a 0.406 change in average log incomes conditional on a country's mean temperature
- $\widehat{\beta}_2$ : A one-degree increase in mean temperature is associated with a -0.06 change in average log incomes conditional on strength of property rights

# General interpretation

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i$$

- The coefficient  $\widehat{\beta}_1$  measures how the predicted outcome varies in  $X_i$  for a fixed value of  $Z_i$ .
- The coefficient  $\widehat{\beta}_2$  measures how the predicted outcome varies in  $Z_i$  for a fixed value of  $X_i$ .

# 4/ OLS mechanics with 2 variables

# Fitted values and residuals

- Where do we get our hats?  $\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2$
- To answer this, we first need to redefine some terms from simple linear regression.
- **Fitted values** for  $i = 1, \dots, n$ :

$$\hat{Y}_i = \widehat{\mathbb{E}}[Y_i|X_i, Z_i] = \hat{\beta}_0 + \hat{\beta}_1 X_i + \hat{\beta}_2 Z_i$$

- **Residuals** for  $i = 1, \dots, n$ :

$$\hat{u}_i = Y_i - \hat{Y}_i$$

# Least squares is still least squares

- How do we estimate  $\hat{\beta}_0$ ,  $\hat{\beta}_1$ , and  $\hat{\beta}_2$ ?
- Minimize the sum of the squared residuals, just like before:

$$(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2) = \arg \min_{b_0, b_1, b_2} \sum_{i=1}^n (Y_i - b_0 - b_1 X_i - b_2 Z_i)^2$$

- The calculus is the same as last week, with 3 partial derivatives instead of 2

# OLS estimator recipe using two steps

- No explicit OLS formulas this week, but a recipe instead
- “Partialling out” OLS recipe:
  1. Run regression of  $X_i$  on  $Z_i$ :

$$\widehat{X}_i = \widehat{\mathbb{E}}[X_i|Z_i] = \widehat{\delta}_0 + \widehat{\delta}_1 Z_i$$

2. Calculate residuals from this regression:

$$\widehat{r}_{xz,i} = X_i - \widehat{X}_i$$

3. Run a simple regression of  $Y_i$  on residuals,  $\widehat{r}_{xz,i}$ :

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 \widehat{r}_{xz,i}$$

- Estimate of  $\widehat{\beta}_1$  will be the same as running:

$$\widehat{Y}_i = \widehat{\beta}_0 + \widehat{\beta}_1 X_i + \widehat{\beta}_2 Z_i$$

# First regression

- Regress  $X_i$  on  $Z_i$ :

```
## when missing data exists, need the na.action in order
## to place residuals or fitted values back into the data
ajr.first <- lm(avexpr ~ meantemp, data = ajr,
               na.action = na.exclude)
summary(ajr.first)
```

```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.9568      0.8202   12.1 < 2e-16 ***
## meantemp     -0.1490      0.0347   -4.3 0.000067 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.32 on 58 degrees of freedom
## (103 observations deleted due to missingness)
## Multiple R-squared:  0.241, Adjusted R-squared:  0.228
## F-statistic: 18.4 on 1 and 58 DF, p-value: 0.0000673
```

# Regression of log income on the residuals

- Save residuals:

```
## store the residuals  
ajr$avexpr.res <- residuals(ajr.first)
```

- Now we compare the estimated slopes:

```
coef(lm(logpgp95 ~ avexpr.res, data = ajr))
```

```
## (Intercept)  avexpr.res  
##      8.0543      0.4057
```

```
coef(lm(logpgp95 ~ avexpr + meantemp, data = ajr))
```

```
## (Intercept)      avexpr      meantemp  
##      6.80627      0.40568     -0.06025
```

# Residual/partial regression plot

- Can plot the **conditional relationship** between property rights and income given temperature:



# **5/** OLS assumptions & inference with 2 variables

# OLS assumptions for unbiasedness

- Last week we made Assumptions 1-4 for unbiasedness of OLS:
  1. Linearity
  2. Random/iid sample
  3. Variation in  $X_i$
  4. Zero conditional mean error:  $\mathbb{E}[u_i|X_i] = 0$
- Small modification to these with 2 covariates:

1. Linearity

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i + u_i$$

2. Random/iid sample
3. No perfect collinearity
4. Zero conditional mean error

$$\mathbb{E}[u_i|X_i, Z_i] = 0$$

# New assumption

## Assumption 3: No perfect collinearity

(1) No independent variable is constant in the sample and (2) there are no exactly linear relationships among the independent variables.

- Two components
  1. Both  $X_i$  and  $Z_i$  have to vary.
  2.  $Z_i$  cannot be a deterministic, linear function of  $X_i$ .
- Part 2 rules out anything of the form:

$$Z_i = a + bX_i$$

- Notice how this is linear (equation of a line) and there is no error, so it is deterministic.
- What's the correlation between  $Z_i$  and  $X_i$ ? 1!

# Perfect collinearity example (I)

- Simple example:
  - ▶  $X_i = 1$  if a country is **not** in Africa and 0 otherwise.
  - ▶  $Z_i = 1$  if a country **is** in Africa and 0 otherwise.
- But, clearly we have the following:

$$Z_i = 1 - X_i$$

- These two variables are perfectly collinear.
- What about the following:
  - ▶  $X_i = \text{property rights}$
  - ▶  $Z_i = X_i^2$
- Do we have to worry about collinearity here?
- No! Because while  $Z_i$  is a deterministic function of  $X_i$ , it is a **nonlinear function** of  $X_i$ .

# R and perfect collinearity

- R, Stata, et al will drop one of the variables if there is perfect collinearity:

```
ajr$nonafrica <- 1 - ajr$africa  
summary(lm(logpgp95 ~ africa + nonafrica, data = ajr))
```

```
##  
## Coefficients: (1 not defined because of singularities)  
##           Estimate Std. Error t value Pr(>|t|)  
## (Intercept)  8.7164      0.0899   96.94 < 2e-16 ***  
## africa      -1.3612      0.1631  -8.35  4.9e-14 ***  
## nonafrica           NA           NA      NA      NA  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 0.913 on 146 degrees of freedom  
## (15 observations deleted due to missingness)  
## Multiple R-squared:  0.323, Adjusted R-squared:  0.318  
## F-statistic: 69.7 on 1 and 146 DF, p-value: 4.87e-14
```

# Perfect collinearity example (II)

- Another example:
  - ▶  $X_i$  = mean temperature in Celsius
  - ▶  $Z_i = 1.8X_i + 32$  (mean temperature in Fahrenheit)

```
ajr$meantemp.f <- 1.8 * ajr$meantemp + 32
coef(lm(logpgp95 ~ meantemp + meantemp.f, data = ajr))
```

```
## (Intercept)      meantemp  meantemp.f
##      10.8455      -0.1207           NA
```

# OLS assumptions for large-sample inference

- For large-sample inference and calculating SEs, we need the two-variable version of the Gauss-Markov assumptions:

1. Linearity

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i + u_i$$

2. Random/iid sample
3. No perfect collinearity
4. Zero conditional mean error

$$\mathbb{E}[u_i | X_i, Z_i] = 0$$

5. Homoskedasticity

$$\text{var}[u_i | X_i, Z_i] = \sigma_u^2$$

# Large-sample inference with 2 covariates

- We have our OLS estimate  $\hat{\beta}_1$  and an estimated SE:  $\widehat{SE}[\hat{\beta}_1]$ .
- Under [assumption 1-5](#), in large samples, we'll have the following:

$$\frac{\hat{\beta}_1 - \beta_1}{\widehat{SE}[\hat{\beta}_1]} \sim N(0, 1)$$

- The same holds for the other coefficient:

$$\frac{\hat{\beta}_2 - \beta_2}{\widehat{SE}[\hat{\beta}_2]} \sim N(0, 1)$$

- Inference is exactly the same in large samples!
- Hypothesis tests and CIs are good to go
- The SE's will change, though

# OLS assumptions for small-sample inference

- For **small-sample inference**, we need the Gauss-Markov plus Normal errors:

1. Linearity

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i + u_i$$

2. Random/iid sample
3. No perfect collinearity
4. Zero conditional mean error

$$\mathbb{E}[u_i | X_i, Z_i] = 0$$

5. Homoskedasticity

$$\text{var}[u_i | X_i, Z_i] = \sigma_u^2$$

6. Normal conditional errors

$$u_i \sim N(0, \sigma_u^2)$$

# Small-sample inference with 2 covariates

- Under assumptions 1-6, we have the following small change to our small- $n$  sampling distribution:

$$\frac{\widehat{\beta}_1 - \beta_1}{\widehat{SE}[\widehat{\beta}_1]} \sim t_{n-3}$$

- The same is true for the other coefficient:

$$\frac{\widehat{\beta}_2 - \beta_2}{\widehat{SE}[\widehat{\beta}_2]} \sim t_{n-3}$$

- Why  $n - 3$ ?
  - ▶ We've estimated another parameter, so we need to take off another degree of freedom.
- ↪ small adjustments to the critical values and the t-values for our hypothesis tests and confidence intervals.

# 6/ Omitted Variable Bias

# Unbiasedness revisited

- True model:

$$Y_i = \beta_0 + \beta_1 X_i + \beta_2 Z_i + u_i$$

- Assumptions 1-4  $\Rightarrow$  we get unbiased estimates of the coefficients
- What happens if we ignore the  $Z_i$  and just run the simple linear regression with just  $X_i$ ?
- Misspecified model:

$$Y_i = \beta_0 + \beta_1 X_i + u_i^* \quad u_i^* = \beta_2 Z_i + u_i$$

- OLS estimates from the misspecified model:

$$\widehat{Y}_i = \tilde{\beta}_0 + \tilde{\beta}_1 X_i$$

# Omitted variable bias

$$Y_i = \beta_0 + \beta_1 X_i + u_i^* \quad u_i^* = \beta_2 Z_i + u_i$$

- Which of the four assumptions might we violate?
  - ▶ Zero conditional mean error!
  - ▶  $E[u_i^*|X_i] \neq 0$  because  $E[Z_i|X_i]$  might not be zero (show on the board, Blackwell)
- **Intuition** Correlation between  $X_i$  and  $Z_i \rightsquigarrow$  correlation between  $X_i$  and the misspecified error  $u_i^*$
- Question: will  $\mathbb{E}[\tilde{\beta}_1] = \beta_1$ ? If not, what will be the bias?

# Omitted variable bias, derivation

- Bias for the misspecified estimator (derived in notes):

$$\text{Bias}(\tilde{\beta}_1) = \mathbb{E}[\tilde{\beta}_1] - \beta_1 = \beta_2 \delta_1$$

- Where the  $\delta_1$  is the coefficient on  $Z_i$  from a regression of  $Z_i$  on  $X_i$ :

$$Z_i = \delta_0 + \delta_1 X_i + e_i$$

- In other words:

omitted variable bias = (effect of  $Z_i$  on  $Y_i$ ) $\times$ (effect of  $X_i$  on  $Z_i$ )

# Omitted variable bias, summary

- Remember that by OLS, the effect of  $X_i$  on  $Z_i$  is:

$$\delta_1 = \frac{\text{cov}(Z_i, X_i)}{\text{var}(X_i)}$$

- We can summarize the direction of bias like so:

	$\text{cov}(X_i, Z_i) > 0$	$\text{cov}(X_i, Z_i) < 0$	$\text{cov}(X_i, Z_i) = 0$
$\beta_2 > 0$	Positive bias	Negative Bias	No bias
$\beta_2 < 0$	Negative bias	Positive Bias	No bias
$\beta_2 = 0$	No bias	No bias	No bias

- Very relevant if  $Z_i$  is unobserved for some reason!

# Including irrelevant variables

- What if we do the opposite and **include an irrelevant variable**?
- What would it mean for  $Z_i$  to be an irrelevant variable?  
Basically, that we have

$$Y_i = \beta_0 + \beta_1 X_i + 0 \times Z_i + u_i$$

- So in this case, the true value of  $\beta_2 = 0$ . But under Assumptions 1-4, OLS is unbiased for all the parameters:

$$\mathbb{E}[\widehat{\beta}_0] = \beta_0$$

$$\mathbb{E}[\widehat{\beta}_1] = \beta_1$$

$$\mathbb{E}[\widehat{\beta}_2] = 0$$

- Including an irrelevant variable will increase the standard errors for  $\widehat{\beta}_1$ .

# 7/ Multicollinearity

# Sampling variance for bivariate regression

- Under simple linear regression, we found that the sampling variance of the slope was the following:

$$\text{var}(\hat{\beta}_1) = \frac{\sigma_u^2}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

- Question** What do we call the square root of the sampling variance? Standard error!
- Factors affecting the standard errors:
  - The error variance  $\sigma_u^2$  (higher conditional variance of  $Y_i$  leads to bigger SEs)
  - The total variation in  $X_i$ :  $\sum_{i=1}^n (X_i - \bar{X})^2$  (lower variation in  $X_i$  leads to bigger SEs)

# Sampling variation with 2 covariates

- Regression with an additional independent variable:

$$\text{var}(\widehat{\beta}_1) = \frac{\sigma_u^2}{(1 - R_1^2) \sum_{i=1}^n (X_i - \bar{X})^2}$$

- Here,  $R_1^2$  is the  $R^2$  from the regression of  $X_i$  on  $Z_i$ :

$$\widehat{X}_i = \widehat{\delta}_0 + \widehat{\delta}_1 Z_i$$

- Factors now affecting the standard errors:
  - ▶ The error variance (higher conditional variance of  $Y_i$  leads to bigger SEs)
  - ▶ The total variation of  $X_i$  (lower variation in  $X_i$  leads to bigger SEs)
  - ▶ The strength of the relationship between  $X_i$  and  $Z_i$  (stronger relationships mean higher  $R_1^2$  and thus bigger SEs)
- What happens with perfect collinearity?

# Multicollinearity

## Definition

Multicollinearity is defined to be high, but not perfect, correlation between two independent variables in a regression.

- With multicollinearity, we'll have  $R_1^2 \approx 1$ , but not exactly.
- The stronger the relationship between  $X_i$  and  $Z_i$ , the closer the  $R_1^2$  will be to 1, and the higher the SEs will be:

$$\text{var}(\hat{\beta}_1) = \frac{\sigma_u^2}{(1 - R_1^2) \sum_{i=1}^n (X_i - \bar{X})^2}$$

- Given the symmetry, it will also increase  $\text{var}(\hat{\beta}_2)$  as well.

# Intuition for multicollinearity

- Remember the OLS recipe:
  - ▶  $\hat{\beta}_1$  from regression of  $Y_i$  on  $\hat{r}_{xz,i}$
  - ▶  $\hat{r}_{xz,i}$  are the residuals from the regression of  $X_i$  on  $Z_i$
- Estimated coefficient:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n \hat{r}_{xz,i} Y_i}{\sum_{i=1}^n \hat{r}_{xz,i}^2}$$

- When  $Z_i$  and  $X_i$  have a strong relationship, then the residuals will have low variation (draw this)
- We explain away a lot of the variation in  $X_i$  through  $Z_i$ .
- Low variation in an independent variable (here,  $\hat{r}_{xz,i}$ )  $\rightsquigarrow$  high SEs
- Basically, there is less residual variation left in  $X_i$  after “partialling out” the effect of  $Z_i$

# Effects of multicollinearity

- No effect on the bias of OLS.
- Only increases the standard errors.
- Really just a sample size problem:
  - ▶ If  $X_i$  and  $Z_i$  are extremely highly correlated, you're going to need a much bigger sample to accurately differentiate between their effects.



# Conclusion

- In this brave new world with 2 independent variables:
  1.  $\beta$ 's have slightly different interpretations
  2. OLS still minimizing the sum of the squared residuals
  3. Small adjustments to OLS assumptions and inference
  4. Adding or omitting variables in a regression can affect the bias and the variance of OLS
- Remainder of class:
  1. Regression in most general glory (matrices)
  2. How to diagnose and fix violations of the OLS assumptions

# Deriving of the misspecified coefficient

- Here we'll use  $\widehat{\text{cov}}$  to mean the sample covariance, and  $\widehat{\text{var}}$  to be the sample variance.

$$\begin{aligned}\tilde{\beta}_1 &= \frac{\widehat{\text{cov}}(Y_i, X_i)}{\widehat{\text{var}}(X_i)} \\ &= \frac{\widehat{\text{cov}}(\beta_0 + \beta_1 X_i + \beta_2 Z_i + u_i, X_i)}{\widehat{\text{var}}(X_i)} \\ &= \frac{\widehat{\text{cov}}(\beta_0, X_i)}{\widehat{\text{var}}(X_i)} + \frac{\widehat{\text{cov}}(\beta_1 X_i, X_i)}{\widehat{\text{var}}(X_i)} + \frac{\widehat{\text{cov}}(\beta_2 Z_i, X_i)}{\widehat{\text{var}}(X_i)} + \frac{\widehat{\text{cov}}(u_i, X_i)}{\widehat{\text{var}}(X_i)} \\ &= 0 + \frac{\widehat{\text{cov}}(\beta_1 X_i, X_i)}{\widehat{\text{var}}(X_i)} + \frac{\widehat{\text{cov}}(\beta_2 Z_i, X_i)}{\widehat{\text{var}}(X_i)} + 0 \\ &= \beta_1 \frac{\widehat{\text{var}}(X_i)}{\widehat{\text{var}}(X_i)} + \beta_2 \frac{\widehat{\text{cov}}(Z_i, X_i)}{\widehat{\text{var}}(X_i)} \\ &= \beta_1 + \beta_2 \frac{\widehat{\text{cov}}(Z_i, X_i)}{\widehat{\text{var}}(X_i)} \\ &= \beta_1 + \beta_2 \widehat{\delta}_1\end{aligned}$$

# Next step

- Let's take expectations:

$$\begin{aligned}\mathbb{E}[\tilde{\beta}_1] &= \mathbb{E}[\beta_1 + \beta_2 \hat{\delta}_1] \\ &= \beta_1 + \beta_2 \mathbb{E}[\hat{\delta}_1] \\ &= \beta_1 + \beta_2 \delta_1\end{aligned}$$

- Thus, we can calculate the bias here:

$$\text{Bias}(\tilde{\beta}_1) = \mathbb{E}[\tilde{\beta}_1] - \beta_1 = \beta_2 \delta_1$$