

# Lecture 2: Fixed Effects

January 21, 2026

# Course Administration

- ① Any problems with summary assignments?
  - I aspire to grade these weekly
- ② Any problems accessing recorded lecture?
- ③ Proposal due next week

# Course Administration

- ① Any problems with summary assignments?
  - I aspire to grade these weekly
- ② Any problems accessing recorded lecture?
- ③ Proposal due next week
- ④ Lab session here after class tonight
  - attend at your discretion

# Course Administration

- ① Any problems with summary assignments?
  - I aspire to grade these weekly
- ② Any problems accessing recorded lecture?
- ③ Proposal due next week
- ④ Lab session here after class tonight
  - attend at your discretion
- ⑤ Problem set 1 due next week
  - Follow link on handout to submit
  - Natalia will provide feedback
  - Speedy quiz on problem set next week
- ⑥ Anything else?

# Today

- ① General problem of selection
- ② Omitted variable bias in terms of regression coefficients
- ③ Indicator variables
- ④ Discussion of Black et al

# 1. General Problem of Selection Bias

## The General Problem

If we assume a homogeneous treatment effect,  $\kappa = Y_{1i} - Y_{0i}$ ,

## The General Problem

If we assume a homogeneous treatment effect,  $\kappa = Y_{1i} - Y_{0i}$ , then

$$\text{Avg}_n[Y_{1i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] =$$

## The General Problem

If we assume a homogeneous treatment effect,  $\kappa = Y_{1i} - Y_{0i}$ , then

$$\begin{aligned} \text{Avg}_n[Y_{1i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] &= \\ \text{Avg}_n[\kappa + Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] &= \end{aligned}$$

## The General Problem

If we assume a homogeneous treatment effect,  $\kappa = Y_{1i} - Y_{0i}$ , then

$$\begin{aligned} & \text{Avg}_n[Y_{1i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] = \\ & \text{Avg}_n[\kappa + Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] = \\ & \kappa + \text{Avg}_n[Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] \end{aligned}$$

## The General Problem

If we assume a homogeneous treatment effect,  $\kappa = Y_{1i} - Y_{0i}$ , then

$$\begin{aligned} & \text{Avg}_n[Y_{1i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] = \\ & \text{Avg}_n[\kappa + Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] = \\ & \kappa + \text{Avg}_n[Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] \end{aligned}$$

Red term is difference in outcome  $Y$  for treated relative to untreated in the absence of treatment. What is it called?

## The General Problem

If we assume a homogeneous treatment effect,  $\kappa = Y_{1i} - Y_{0i}$ , then

$$\begin{aligned} \text{Avg}_n[Y_{1i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] &= \\ \text{Avg}_n[\kappa + Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] &= \\ \kappa + \text{Avg}_n[Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0] \end{aligned}$$

Red term is difference in outcome  $Y$  for treated relative to untreated in the absence of treatment. What is it called? **selection bias**.

## Let's Think of Some Examples of Selection Bias

$$\text{Avg}_n[Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0]$$

## Let's Think of Some Examples of Selection Bias

$$\text{Avg}_n[Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0]$$

A fix: control for covariates  $X_i$  to make selection bias disappear.

## Let's Think of Some Examples of Selection Bias

$$\text{Avg}_n[Y_{0i}|D_i = 1] - \text{Avg}_n[Y_{0i}|D_i = 0]$$

A fix: control for covariates  $X_i$  to make selection bias disappear.

Strong evidence that “controlling for observables” rarely gets rid of selection.

## 2. Omitted Variable Bias Formula

## Long (True) vs. Short (False) Regression

Suppose that the “true” (long) regression is

$$Y = \alpha + \beta'X_1 + \gamma X_2 + \epsilon'$$

## Long (True) vs. Short (False) Regression

Suppose that the “true” (long) regression is

$$Y = \alpha + \beta'X_1 + \gamma X_2 + \epsilon'$$

Unfortunately, you don't observe  $X_2$  – examples?

## Long (True) vs. Short (False) Regression

Suppose that the “true” (long) regression is

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l$$

Unfortunately, you don't observe  $X_2$  – examples?

So instead you estimate the “false” (short) regression

$$Y = \alpha + \beta^s X_1 + \epsilon^s$$

Should you trust  $\beta^s$ ?

## Evaluating Whether to Trust $\beta^s$

Recall

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l \quad (1)$$

$$Y = \alpha + \beta^s X_1 + \epsilon^s \quad (2)$$

## Evaluating Whether to Trust $\beta^s$

Recall

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l \quad (1)$$

$$Y = \alpha + \beta^s X_1 + \epsilon^s \quad (2)$$

Estimate the relationship between the treatment  $X_1$  and the omitted variable  $X_2$ :

$$X_2 = \pi_0 + \pi_1 X_1 + \epsilon^c$$

## Evaluating Whether to Trust $\beta^s$

Recall

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l \quad (1)$$

$$Y = \alpha + \beta^s X_1 + \epsilon^s \quad (2)$$

Estimate the relationship between the treatment  $X_1$  and the omitted variable  $X_2$ :

$$X_2 = \pi_0 + \pi_1 X_1 + \epsilon^c$$

Then (proof in book)

OVB =

## Evaluating Whether to Trust $\beta^s$

Recall

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l \quad (1)$$

$$Y = \alpha + \beta^s X_1 + \epsilon^s \quad (2)$$

Estimate the relationship between the treatment  $X_1$  and the omitted variable  $X_2$ :

$$X_2 = \pi_0 + \pi_1 X_1 + \epsilon^c$$

Then (proof in book)

$$\text{OVB} = \beta^s - \beta^l$$

## Evaluating Whether to Trust $\beta^s$

Recall

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l \quad (1)$$

$$Y = \alpha + \beta^s X_1 + \epsilon^s \quad (2)$$

Estimate the relationship between the treatment  $X_1$  and the omitted variable  $X_2$ :

$$X_2 = \pi_0 + \pi_1 X_1 + \epsilon^c$$

Then (proof in book)

$$\text{OVB} = \beta^s - \beta^l = \pi_1 \gamma$$

## Evaluating Whether to Trust $\beta^s$

Recall

$$Y = \alpha + \beta^l X_1 + \gamma X_2 + \epsilon^l \quad (1)$$

$$Y = \alpha + \beta^s X_1 + \epsilon^s \quad (2)$$

Estimate the relationship between the treatment  $X_1$  and the omitted variable  $X_2$ :

$$X_2 = \pi_0 + \pi_1 X_1 + \epsilon^c$$

Then (proof in book)

$$\text{OVB} = \beta^s - \beta^l = \pi_1 \gamma$$

OVB is one type of selection bias.

## Let's think about this equation

$\pi_1 \equiv$  relationship between  $X_2$  and  $X_1$

$\gamma \equiv$  relationship between  $X_2$  and  $Y$  in long regression

$$\text{OVB} = \beta^s - \beta^l = \pi_1 \gamma$$

- What if the treatment and the omitted variable are not correlated?

## Let's think about this equation

$\pi_1 \equiv$  relationship between  $X_2$  and  $X_1$

$\gamma \equiv$  relationship between  $X_2$  and  $Y$  in long regression

$$\text{OVB} = \beta^s - \beta^l = \pi_1 \gamma$$

- What if the treatment and the omitted variable are not correlated?
- What if the omitted variable is not correlated with the outcome  $Y$ ?

## Let's think about this equation

$\pi_1 \equiv$  relationship between  $X_2$  and  $X_1$

$\gamma \equiv$  relationship between  $X_2$  and  $Y$  in long regression

$$\text{OVB} = \beta^s - \beta^l = \pi_1 \gamma$$

- What if the treatment and the omitted variable are not correlated?
- What if the omitted variable is not correlated with the outcome  $Y$ ?
- Any story about omitted variable bias needs to include **both** parts

## Let's think about this equation

$\pi_1 \equiv$  relationship between  $X_2$  and  $X_1$

$\gamma \equiv$  relationship between  $X_2$  and  $Y$  in long regression

$$\text{OVB} = \beta^s - \beta^l = \pi_1 \gamma$$

- What if the treatment and the omitted variable are not correlated?
- What if the omitted variable is not correlated with the outcome  $Y$ ?
- Any story about omitted variable bias needs to include **both** parts
- Resolving the problem of omitted variable bias in order to generate causal estimates is the key concern of this course

### 3. Indicator Variables

# What is an indicator variable?

All these things are the same

- dummy variable
- indicator variable
- fixed effect
- $1\{\text{condition}\}$

## What is an indicator variable?

All these things are the same

- dummy variable
- indicator variable
- fixed effect
- $1\{\text{condition}\}$

All are coded 1 if true and 0 otherwise

## Interpreting Indicator Variables

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

- $\text{female} \in \{0, 1\}$
- how do we interpret  $\beta_1$ ?

## Interpreting Indicator Variables

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

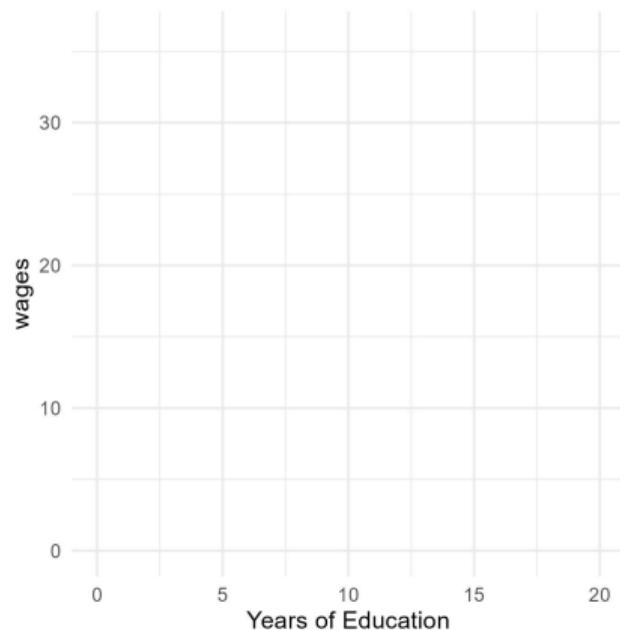
- $\text{female} \in \{0, 1\}$
- how do we interpret  $\beta_1$ ?
- let's draw in a figure

## Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?

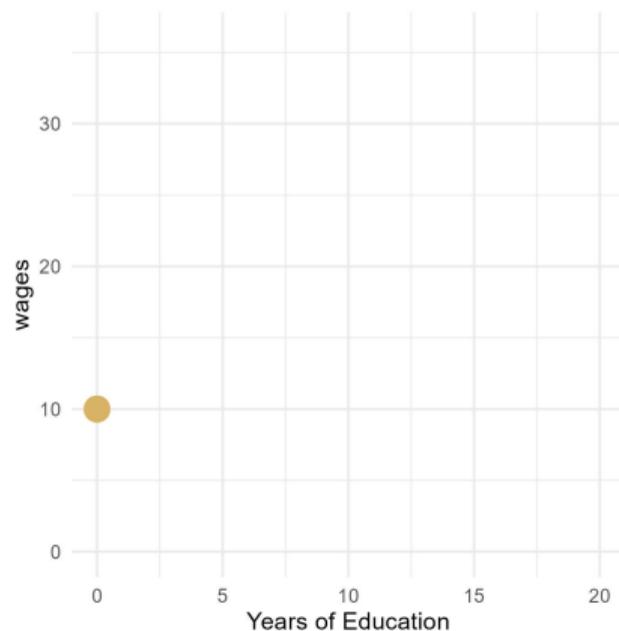


## Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?

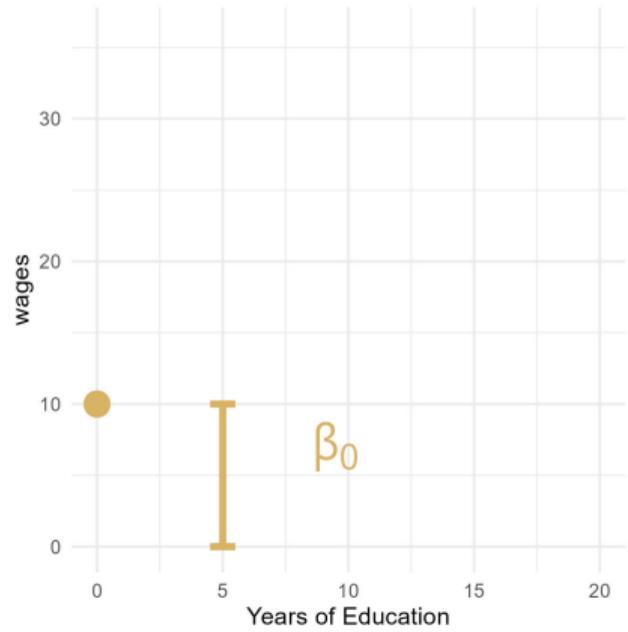


# Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?

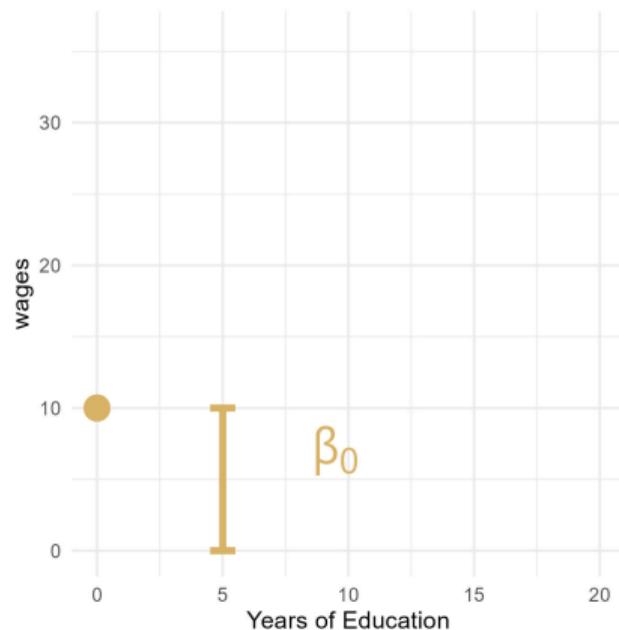


# Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?
- where is  $\beta_2$ ?

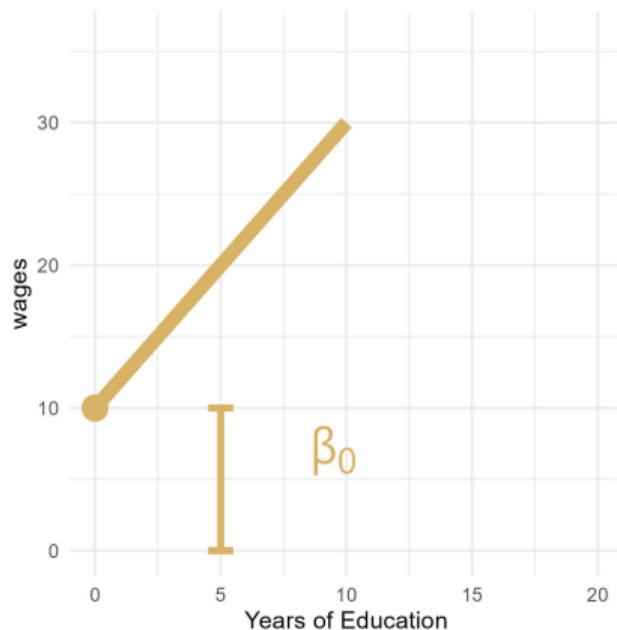


# Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?
- where is  $\beta_2$ ?

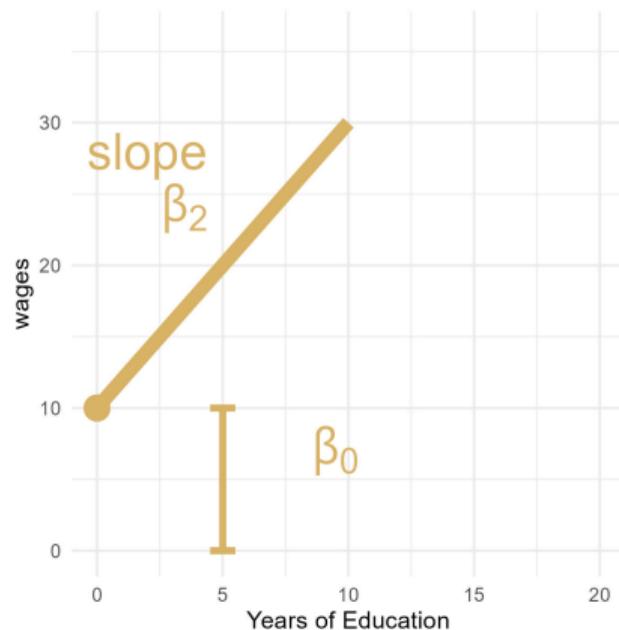


## Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?
- where is  $\beta_2$ ?

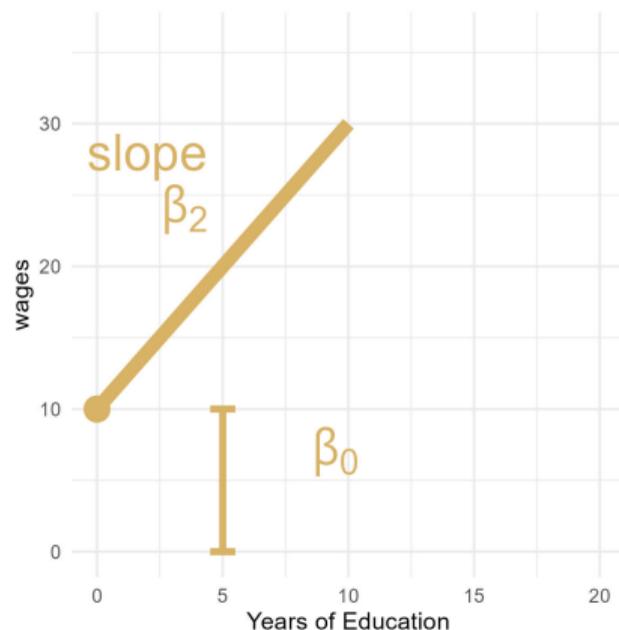


## Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?
- where is  $\beta_2$ ?
- how do we draw wages for women as a function of education?

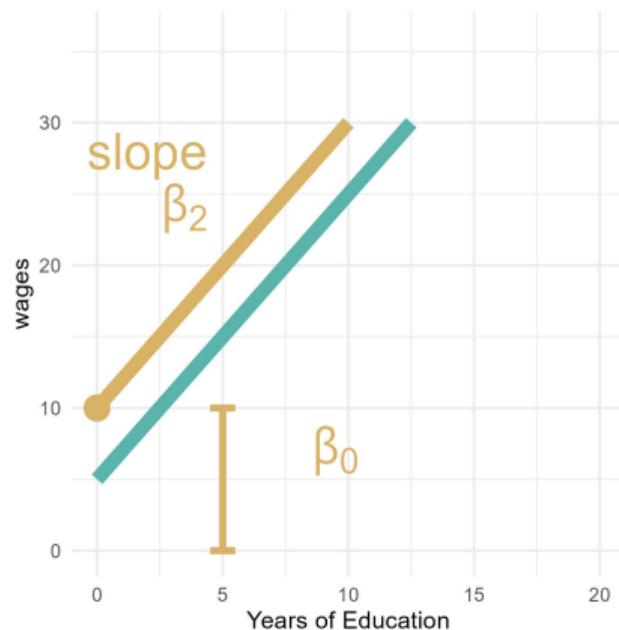


## Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?
- where is  $\beta_2$ ?
- how do we draw wages for women as a function of education?  
 $\beta_2 * \text{education} + \beta_1$

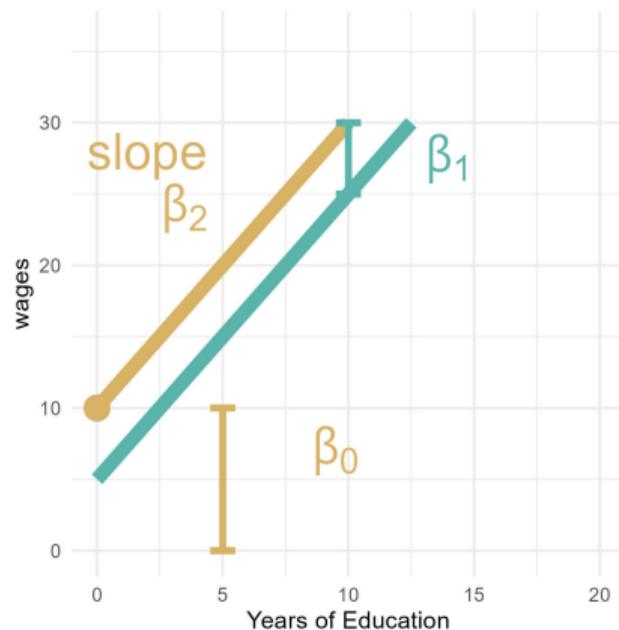


## Interpreting Coefficients

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

Draw the relationship

- x axis is education
- y axis is wage
- where is  $\beta_0$ ?
- where is  $\beta_2$ ?
- how do we draw wages for women as a function of education?  
 $\beta_2 * \text{education} + \beta_1$



## Coding Variables

- Suppose we want to look at the relationship between gender and wages:

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

- Data are coded 1 for men, 2 for women
- Why don't we just use this coding? Why do we make a dummy variable?

## Coding Variables

- Suppose we want to look at the relationship between gender and wages:

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

- Data are coded 1 for men, 2 for women
- Why don't we just use this coding? Why do we make a dummy variable?
- Why do we not make one dummy variable for each gender?

## Coding Variables

- Suppose we want to look at the relationship between gender and wages:

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

- Data are coded 1 for men, 2 for women
- Why don't we just use this coding? Why do we make a dummy variable?
- Why do we not make one dummy variable for each gender?
- How can you modify the specification so that the impact of education may differ by gender?

## Coding Variables

- Suppose we want to look at the relationship between gender and wages:

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \epsilon$$

- Data are coded 1 for men, 2 for women
- Why don't we just use this coding? Why do we make a dummy variable?
- Why do we not make one dummy variable for each gender?
- How can you modify the specification so that the impact of education may differ by gender?

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

## Interpreting Indicator Variables in Interaction

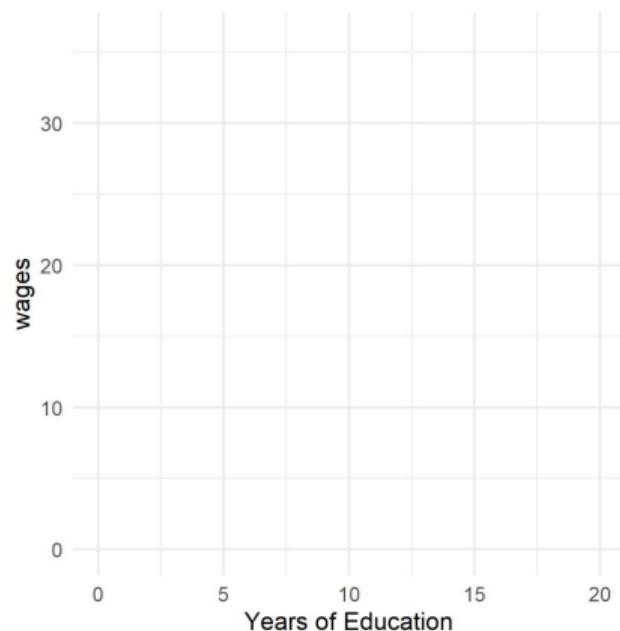
$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

- $\text{female} \in \{0, 1\}$
- what is this specification doing differently?

## Interpreting Coefficients in Interacted Specification

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

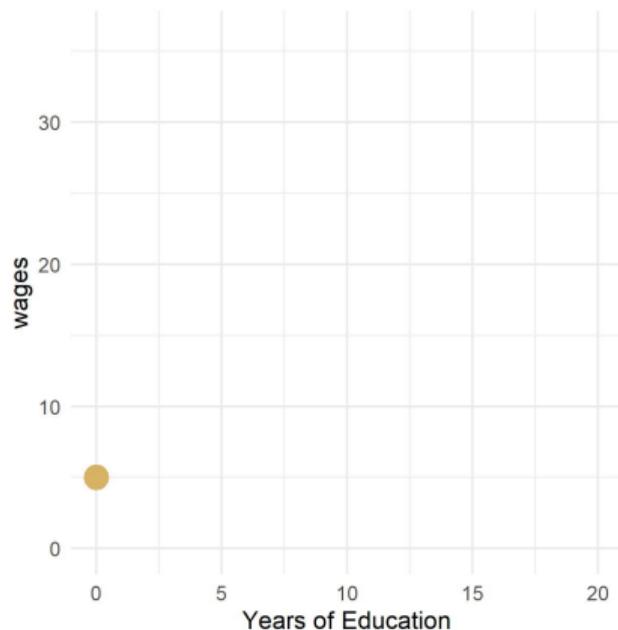
- what are men's wages with no education?



## Interpreting Coefficients in Interacted Specification

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

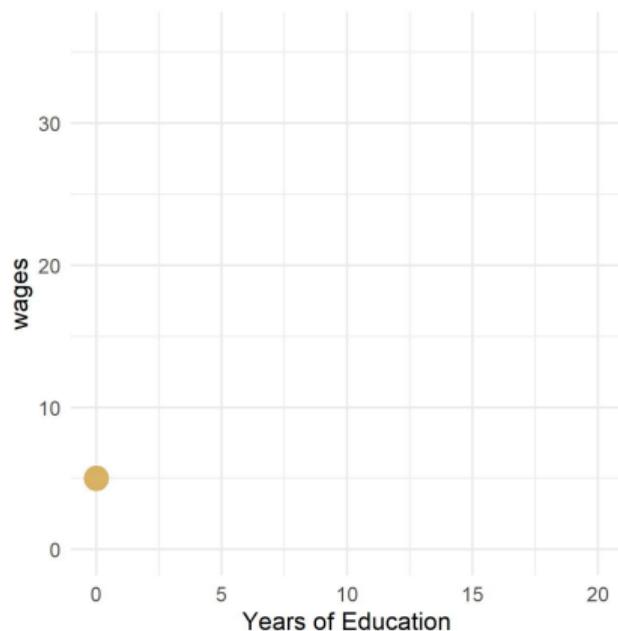
- what are men's wages with no education?  $\beta_0$



## Interpreting Coefficients in Interacted Specification

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

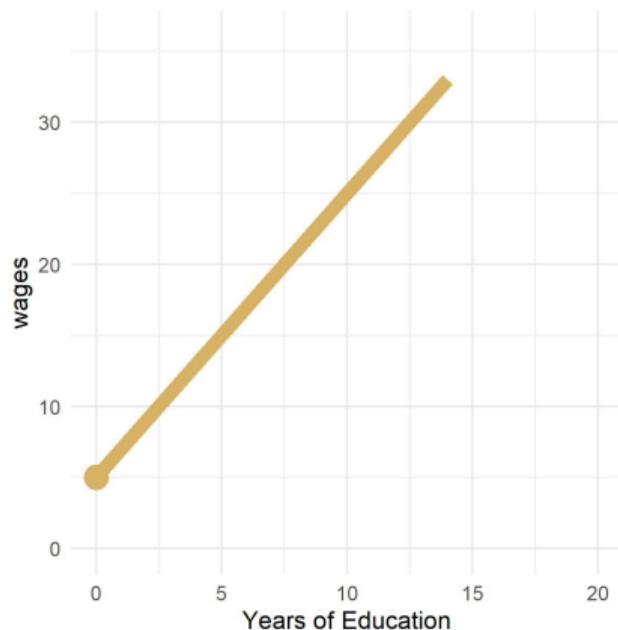
- what are men's wages with no education?  $\beta_0$
- how do men's wages change with education?



## Interpreting Coefficients in Interacted Specification

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

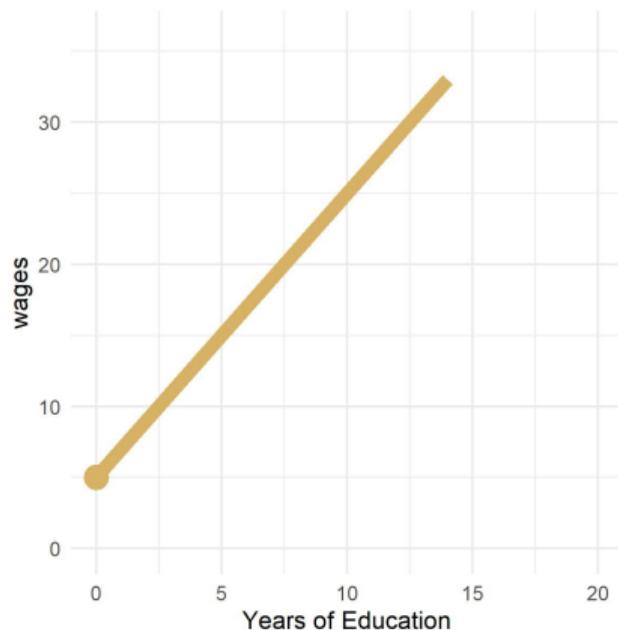
- what are men's wages with no education?  $\beta_0$
- how do men's wages change with education?  $\beta_2 * \text{education}$



## Interpreting Coefficients in Interacted Specification

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

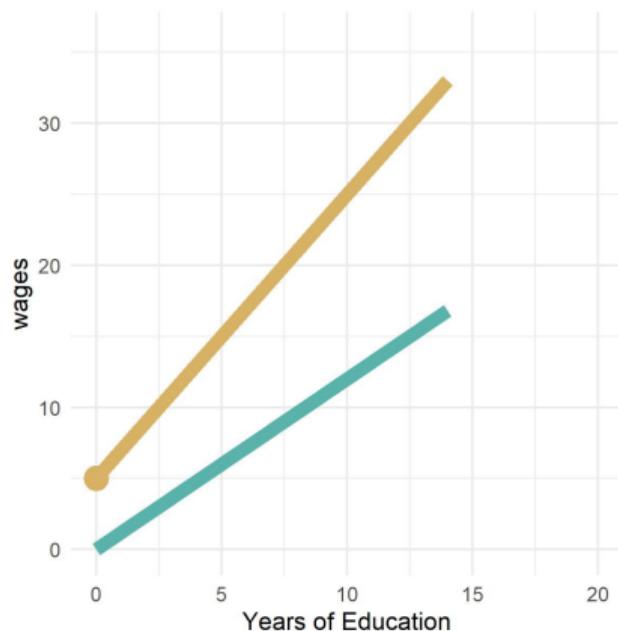
- what are men's wages with no education?  $\beta_0$
- how do men's wages change with education?  $\beta_2 * \text{education}$
- how do women's wages change with education?



## Interpreting Coefficients in Interacted Specification

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

- what are men's wages with no education?  $\beta_0$
- how do men's wages change with education?  $\beta_2 * \text{education}$
- how do women's wages change with education?  
start at  $\beta_0 + \beta_1$   
change by  
 $\beta_2 * \text{education} + \beta_3 * \text{education}$



## Formal Testing

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

- How to test whether the relationship between education and wages differs by gender?

## Formal Testing

$$\text{wage} = \beta_0 + \beta_1 \text{female} + \beta_2 \text{education} + \beta_3 \text{female} * \text{education} + \epsilon$$

- How to test whether the relationship between education and wages differs by gender?
- Test  $\beta_3 = 0$

## 4. Black et al on family size

# Paper Overview

What is this paper about?

- what is the theory that they rebut in this paper?

# Paper Overview

What is this paper about?

- what is the theory that they rebut in this paper? theory about quality vs. quantity in kids
- to whom is it due?

# Paper Overview

What is this paper about?

- what is the theory that they rebut in this paper? theory about quality vs. quantity in kids
- to whom is it due? Nobel laureate Becker and some buddies

# Paper Overview

What are the data?

What is this paper about?

- what is the theory that they rebut in this paper? theory about quality vs. quantity in kids
- to whom is it due? Nobel laureate Becker and some buddies

# Paper Overview

What is this paper about?

- what is the theory that they rebut in this paper? theory about quality vs. quantity in kids
- to whom is it due? Nobel laureate Becker and some buddies

What are the data?

- people aged 16-74 from 1986-2000 (would you be in this sample? )
- parents and kids must both appear in the dataset
- can match parents to kids
- about each person they know year of birth, completed education, earnings
- about each family, they know family size

# Paper Overview

What is this paper about?

- what is the theory that they rebut in this paper? theory about quality vs. quantity in kids
- to whom is it due? Nobel laureate Becker and some buddies

What are the data?

- people aged 16-74 from 1986-2000 (would you be in this sample? )
- parents and kids must both appear in the dataset
- can match parents to kids
- about each person they know year of birth, completed education, earnings
- about each family, they know family size
- what is the unit of observation?

# What Can We Learn from Summary Statistics?

TABLE III  
AVERAGE EDUCATION BY NUMBER OF CHILDREN IN FAMILY AND BIRTH ORDER

	Average education	Average mother's education	Average father's education	Fraction with <12 years	Fraction with 12 years	Fraction with >12 years
Family size						
1	12.0	9.2	10.1	.44	.25	.31
2	12.4	9.9	10.8	.34	.31	.35
3	12.3	9.7	10.6	.37	.30	.33
4	12.0	9.3	10.1	.43	.29	.28
5	11.7	8.8	9.5	.49	.27	.24
6	11.4	8.5	9.1	.54	.25	.20
7	11.2	8.3	8.9	.57	.24	.19
8	11.1	8.2	8.8	.58	.24	.18
9	11.0	8.0	8.6	.59	.25	.16
10+	11.0	7.9	8.8	.59	.26	.15
Birth order						
1	12.2	9.7	10.6	.38	.28	.34
2	12.2	9.6	10.5	.38	.30	.31
3	12.0	9.3	10.2	.40	.31	.29
4	11.9	9.0	9.7	.43	.32	.25
5	11.7	8.6	9.2	.46	.31	.22
6	11.6	8.3	8.9	.49	.31	.20
7	11.5	8.1	8.7	.51	.30	.19
8	11.6	8.0	8.6	.49	.31	.20
9	11.3	7.9	8.4	.53	.32	.15
10+	11.3	7.8	8.7	.52	.32	.15
			All			
	12.2	9.5	10.4	.39	.29	.32

- We ignore instrumental variables and twins
- Focus only on the regular estimations
- But start with summary stats
- What does Table 3 tell us about education as family size increases?

# What Can We Learn from Summary Statistics?

TABLE III  
AVERAGE EDUCATION BY NUMBER OF CHILDREN IN FAMILY AND BIRTH ORDER

	Average education	Average mother's education	Average father's education	Fraction with <12 years	Fraction with 12 years	Fraction with >12 years
Family size						
1	12.0	9.2	10.1	.44	.25	.31
2	12.4	9.9	10.8	.34	.31	.35
3	12.3	9.7	10.6	.37	.30	.33
4	12.0	9.3	10.1	.43	.29	.28
5	11.7	8.8	9.5	.49	.27	.24
6	11.4	8.5	9.1	.54	.25	.20
7	11.2	8.3	8.9	.57	.24	.19
8	11.1	8.2	8.8	.58	.24	.18
9	11.0	8.0	8.6	.59	.25	.16
10+	11.0	7.9	8.8	.59	.26	.15
Birth order						
1	12.2	9.7	10.6	.38	.28	.34
2	12.2	9.6	10.5	.38	.30	.31
3	12.0	9.3	10.2	.40	.31	.29
4	11.9	9.0	9.7	.43	.32	.25
5	11.7	8.6	9.2	.46	.31	.22
6	11.6	8.3	8.9	.49	.31	.20
7	11.5	8.1	8.7	.51	.30	.19
8	11.6	8.0	8.6	.49	.31	.20
9	11.3	7.9	8.4	.53	.32	.15
10+	11.3	7.8	8.7	.52	.32	.15
			All			
	12.2	9.5	10.4	.39	.29	.32

- We ignore instrumental variables and twins
- Focus only on the regular estimations
- But start with summary stats
- What does Table 3 tell us about education as family size increases? increases (for 1 to 2), then declines
- What does Table 3 tell us about education as birth order increases?

# What Can We Learn from Summary Statistics?

TABLE III  
AVERAGE EDUCATION BY NUMBER OF CHILDREN IN FAMILY AND BIRTH ORDER

	Average education	Average mother's education	Average father's education	Fraction with <12 years	Fraction with 12 years	Fraction with >12 years
Family size						
1	12.0	9.2	10.1	.44	.25	.31
2	12.4	9.9	10.8	.34	.31	.35
3	12.3	9.7	10.6	.37	.30	.33
4	12.0	9.3	10.1	.43	.29	.28
5	11.7	8.8	9.5	.49	.27	.24
6	11.4	8.5	9.1	.54	.25	.20
7	11.2	8.3	8.9	.57	.24	.19
8	11.1	8.2	8.8	.58	.24	.18
9	11.0	8.0	8.6	.59	.25	.16
10+	11.0	7.9	8.8	.59	.26	.15
Birth order						
1	12.2	9.7	10.6	.38	.28	.34
2	12.2	9.6	10.5	.38	.30	.31
3	12.0	9.3	10.2	.40	.31	.29
4	11.9	9.0	9.7	.43	.32	.25
5	11.7	8.6	9.2	.46	.31	.22
6	11.6	8.3	8.9	.49	.31	.20
7	11.5	8.1	8.7	.51	.30	.19
8	11.6	8.0	8.6	.49	.31	.20
9	11.3	7.9	8.4	.53	.32	.15
10+	11.3	7.8	8.7	.52	.32	.15
			All			
	12.2	9.5	10.4	.39	.29	.32

- We ignore instrumental variables and twins
- Focus only on the regular estimations
- But start with summary stats
- What does Table 3 tell us about education as family size increases? increases (for 1 to 2), then declines
- What does Table 3 tell us about education as birth order increases? declines
- Give an example of an omitted variable when studying the impact of family size on wages

## Make a Class Dataset

- Get four families as an example to match paper
- What info do we need?

## Make a Class Dataset

- Get four families as an example to match paper
- What info do we need?
  - year of birth of each sibling
  - education of each family member

## Make a Class Dataset

- Get four families as an example to match paper
- What info do we need?
  - year of birth of each sibling
  - education of each family member
- Make this into a dataset you could do the sort of regressions that Black et al did.
- Make a copy of the google sheet I sent and enter data there
- Some hints
  - What's the unit of observation?

## Make a Class Dataset

- Get four families as an example to match paper
- What info do we need?
  - year of birth of each sibling
  - education of each family member
- Make this into a dataset you could do the sort of regressions that Black et al did.
- Make a copy of the google sheet I sent and enter data there
- Some hints
  - What's the unit of observation? person
  - What variables do you need?

## Make a Class Dataset

- Get four families as an example to match paper
- What info do we need?
  - year of birth of each sibling
  - education of each family member
- Make this into a dataset you could do the sort of regressions that Black et al did.
- Make a copy of the google sheet I sent and enter data there
- Some hints
  - What's the unit of observation? person
  - What variables do you need?
    - you need to be able to know who is in the same family
    - you need a variable for birth order
    - you need a variable for family size

## Understanding Main Estimates: Table 4

What's the estimating equation for Table 4 column 1? (read p. 678, pp under 3.A.)

## Understanding Main Estimates: Table 4

What's the estimating equation for Table 4 column 1? (read p. 678, pp under 3.A.)

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \epsilon_{i,f}$$

# Modifying Dataset to Estimate

## Estimating Column 1

- To estimate column 1, what additional variable does your dataset need?

## Modifying Dataset to Estimate

### Estimating Column 1

- To estimate column 1, what additional variable does your dataset need?  
birth year FE
- Why do we include year of birth fe?
- How do we interpret the coeff -0.182?

## Modifying Dataset to Estimate

### Estimating Column 1

- To estimate column 1, what additional variable does your dataset need?  
birth year FE
- Why do we include year of birth fe?
- How do we interpret the coeff -0.182?  
increasing family size by one more child decreases the average child's education by .18 of a year (20% of a year)

### Estimating Column 2 – New regression equation?

## Modifying Dataset to Estimate

### Estimating Column 1

- To estimate column 1, what additional variable does your dataset need?  
birth year FE
- Why do we include year of birth fe?
- How do we interpret the coeff -0.182?  
increasing family size by one more child decreases the average child's education by .18 of a year (20% of a year)

### Estimating Column 2 – New regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{kids in fam FE}_f + \beta_2 \text{year of birth FE}_i + \epsilon_{i,f}$$

## Modifying Dataset to Estimate

### Estimating Column 1

- To estimate column 1, what additional variable does your dataset need?  
birth year FE
- Why do we include year of birth fe?
- How do we interpret the coeff -0.182?  
increasing family size by one more child decreases the average child's education by .18 of a year (20% of a year)

### Estimating Column 2 – New regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{kids in fam FE}_f + \beta_2 \text{year of birth FE}_i + \epsilon_{i,f}$$

- what does our dataset need to estimate it?

## Modifying Dataset to Estimate

### Estimating Column 1

- To estimate column 1, what additional variable does your dataset need?  
birth year FE
- Why do we include year of birth fe?
- How do we interpret the coeff -0.182?  
increasing family size by one more child decreases the average child's education by .18 of a year (20% of a year)

### Estimating Column 2 – New regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{kids in fam FE}_f + \beta_2 \text{year of birth FE}_i + \epsilon_{i,f}$$

- what does our dataset need to estimate it?
- how do we interpret 0.272?

## Table 4: Columns 3 and 4

Eq for Table 4, Column 3:

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \epsilon_{i,f}$$

- Add controls. Any questions about how they do that?

## Table 4: Columns 3 and 4

Eq for Table 4, Column 3:

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \epsilon_{i,f}$$

- Add controls. Any questions about how they do that?
- What do we learn by comparing columns 3 and 4 to 1 and 2?

## Table 4: Columns 3 and 4

Eq for Table 4, Column 3:

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \epsilon_{i,f}$$

- Add controls. Any questions about how they do that?
- What do we learn by comparing columns 3 and 4 to 1 and 2?
- Controls are important, but they don't account for the entire effect

## Table 4: Columns 5 and 6

- Column 5
  - what is the regression equation?

## Table 4: Columns 5 and 6

- Column 5

- what is the regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \beta_4 \text{birth order FE}_i + \epsilon_{i,f}$$

- fix your dataset to have enough variables to estimate this
- how do we interpret these coefficients?

## Table 4: Columns 5 and 6

- Column 5
  - what is the regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \beta_4 \text{birth order FE}_i + \epsilon_{i,f}$$

- fix your dataset to have enough variables to estimate this
  - how do we interpret these coefficients?
- then column 6
  - what is the regression equation?

## Table 4: Columns 5 and 6

- Column 5

- what is the regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{no. kids in fam}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \beta_4 \text{birth order FE}_i + \epsilon_{i,f}$$

- fix your dataset to have enough variables to estimate this
- how do we interpret these coefficients?

- then column 6

- what is the regression equation?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{kids in fam FE}_f + \beta_2 \text{year of birth FE}_i + \beta_3 X_{i,f} + \beta_4 \text{birth order FE}_i + \epsilon_{i,f}$$

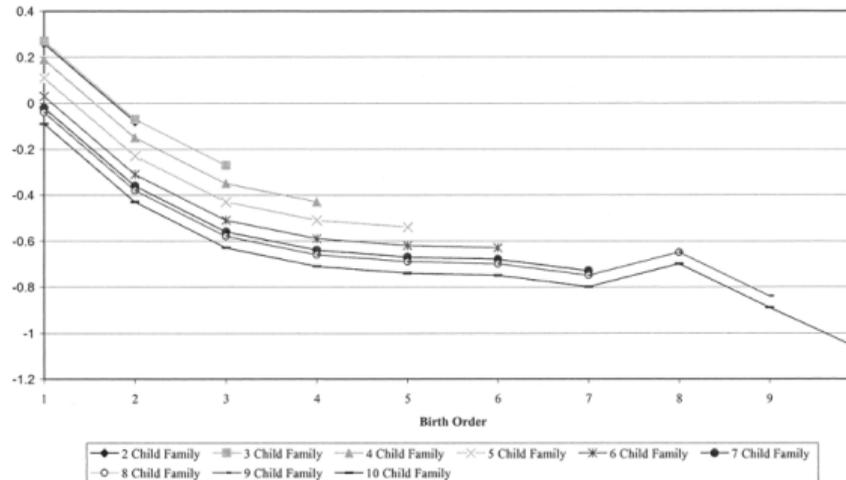
- fix your dataset so that you have enough variables to estimate this

## Visual Representation of Findings

- How does this translate to figure 1 (p. 689)?
- Or, what are they plotting there and what does it mean?
  - warning: the note is not correct – it says predicted values, but these are coefficients

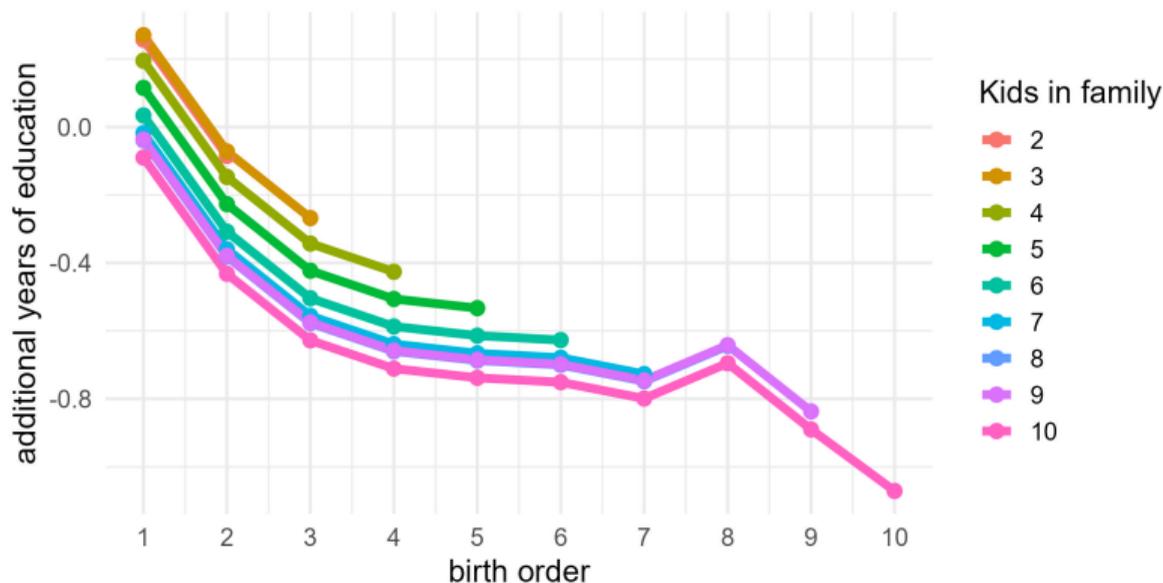
# Visual Representation of Findings

- How does this translate to figure 1 (p. 689)?
- Or, what are they plotting there and what does it mean?
  - warning: the note is not correct – it says predicted values, but these are coefficients



## My Version: Visual Representation of Findings

- How does this translate to figure 1 (p. 689)?
- Or, what are they plotting there and what does it mean?
  - warning: the note is not correct – it says predicted values, but these are coefficients



## Making the Figure, Family Size = 2

- no info for family size = 1

## Making the Figure, Family Size = 2

- no info for family size = 1
- family size of 2
  - first child?

## Making the Figure, Family Size = 2

- no info for family size = 1
- family size of 2
  - first child? 0.257

## Making the Figure, Family Size = 2

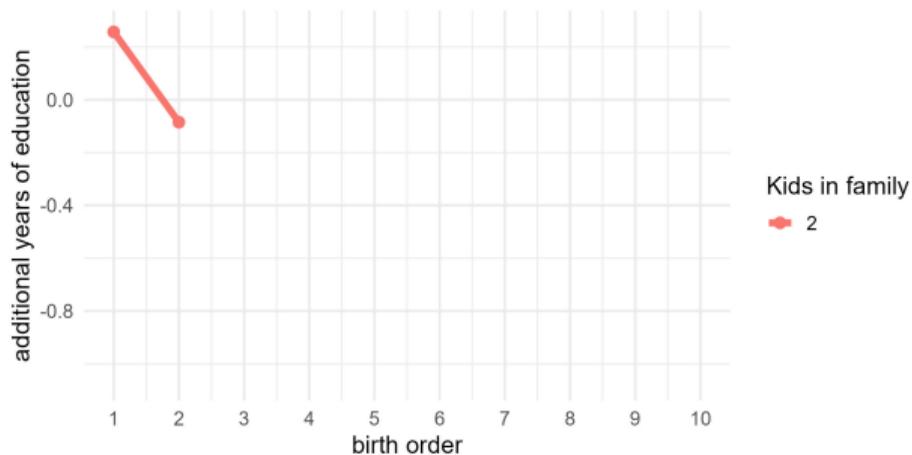
- no info for family size = 1
- family size of 2
  - first child? 0.257
  - second child?

## Making the Figure, Family Size = 2

- no info for family size = 1
- family size of 2
  - first child? 0.257
  - second child? 0.257-0.342

## Making the Figure, Family Size = 2

- no info for family size = 1
- family size of 2
  - first child? 0.257
  - second child? 0.257-0.342



## Making the Figure, Family Size = 3

- family size of 3
  - first child?

## Making the Figure, Family Size = 3

- family size of 3
  - first child? 0.270

## Making the Figure, Family Size = 3

- family size of 3
  - first child? 0.270
  - second child?

## Making the Figure, Family Size = 3

- family size of 3
  - first child? 0.270
  - second child?  
0.270-0.342

## Making the Figure, Family Size = 3

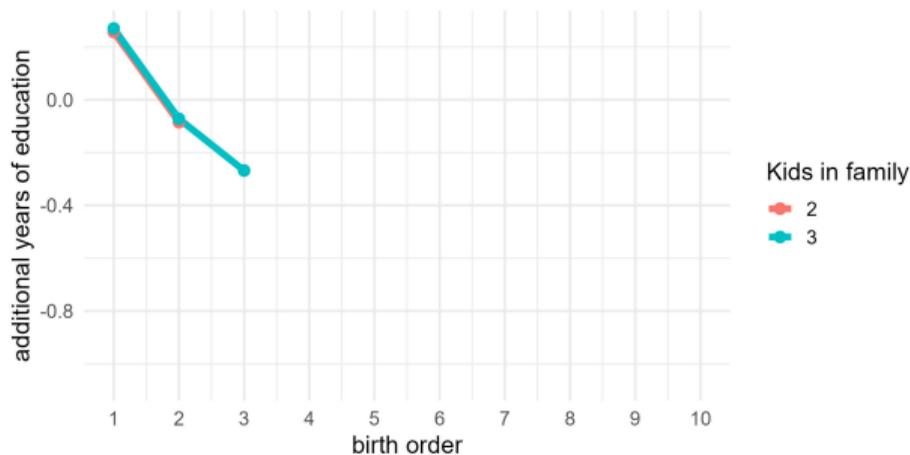
- family size of 3
  - first child? 0.270
  - second child?  
0.270-0.342
  - third child?

## Making the Figure, Family Size = 3

- family size of 3
  - first child? 0.270
  - second child?  
0.270-0.342
  - third child?  
0.270-0.538

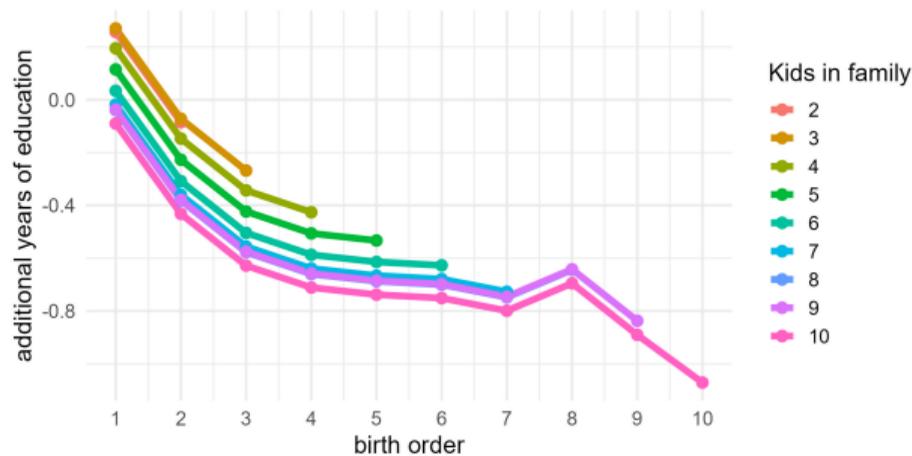
## Making the Figure, Family Size = 3

- family size of 3
  - first child? 0.270
  - second child? 0.270-0.342
  - third child? 0.270-0.538
- why are the lines in the figure parallel?



## Revisiting the Final Figure

- Which estimate would allow us to plot non-parallel lines?



## Understanding Table 6

- what's the estimating eqn for table 6, col 1 (p 687)?

## Understanding Table 6

- what's the estimating eqn for table 6, col 1 (p 687)?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{year of birth dum}_i + \beta_2 X_i + \beta_3 \{1 \text{ if child 2}\}_i + \epsilon_{i,f}$$

## Understanding Table 6

- what's the estimating eqn for table 6, col 1 (p 687)?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{year of birth dum}_i + \beta_2 X_i + \beta_3 \{1 \text{ if child 2}\}_i + \epsilon_{i,f}$$

- do you have the data for these?
- why are these different than the last column of Table 3?

## Understanding Table 6

- what's the estimating eqn for table 6, col 1 (p 687)?

$$\text{educ}_{i,f} = \beta_0 + \beta_1 \text{year of birth dum}_i + \beta_2 X_i + \beta_3 \{1 \text{ if child } 2\}_i + \epsilon_{i,f}$$

- do you have the data for these?
- why are these different than the last column of Table 3?
- Because they allow the effect of birth order to vary by family size

## Next Lecture

- Read *Causal Mixtape*, Chapter 9.1 and 9.2
- Read linked Milligan article, section 5 optional
- Due next week
  - One page proposal
  - Problem set 1
- Next week handout – Problem Set 2, with two week work period