# Course Announcements

- Due Friday (11:59 PM)
  - D4
  - Q4
  - A2
- PLEASE be precise in autograded notebooks!!!

Grading underway: Project Proposals

# Inferential analysis

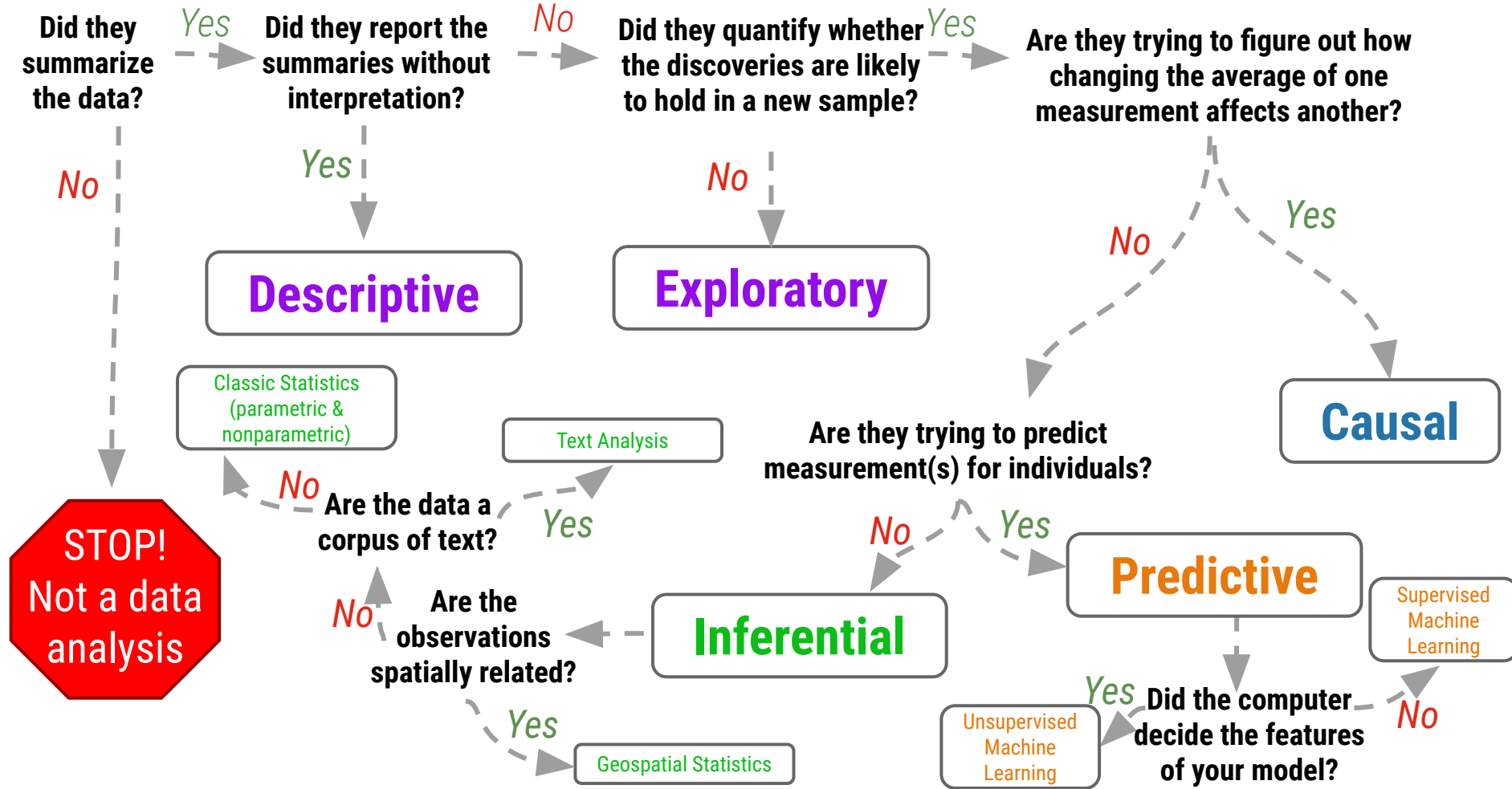**Jason G. Fleischer, Ph.D.**
**Asst. Teaching Professor**
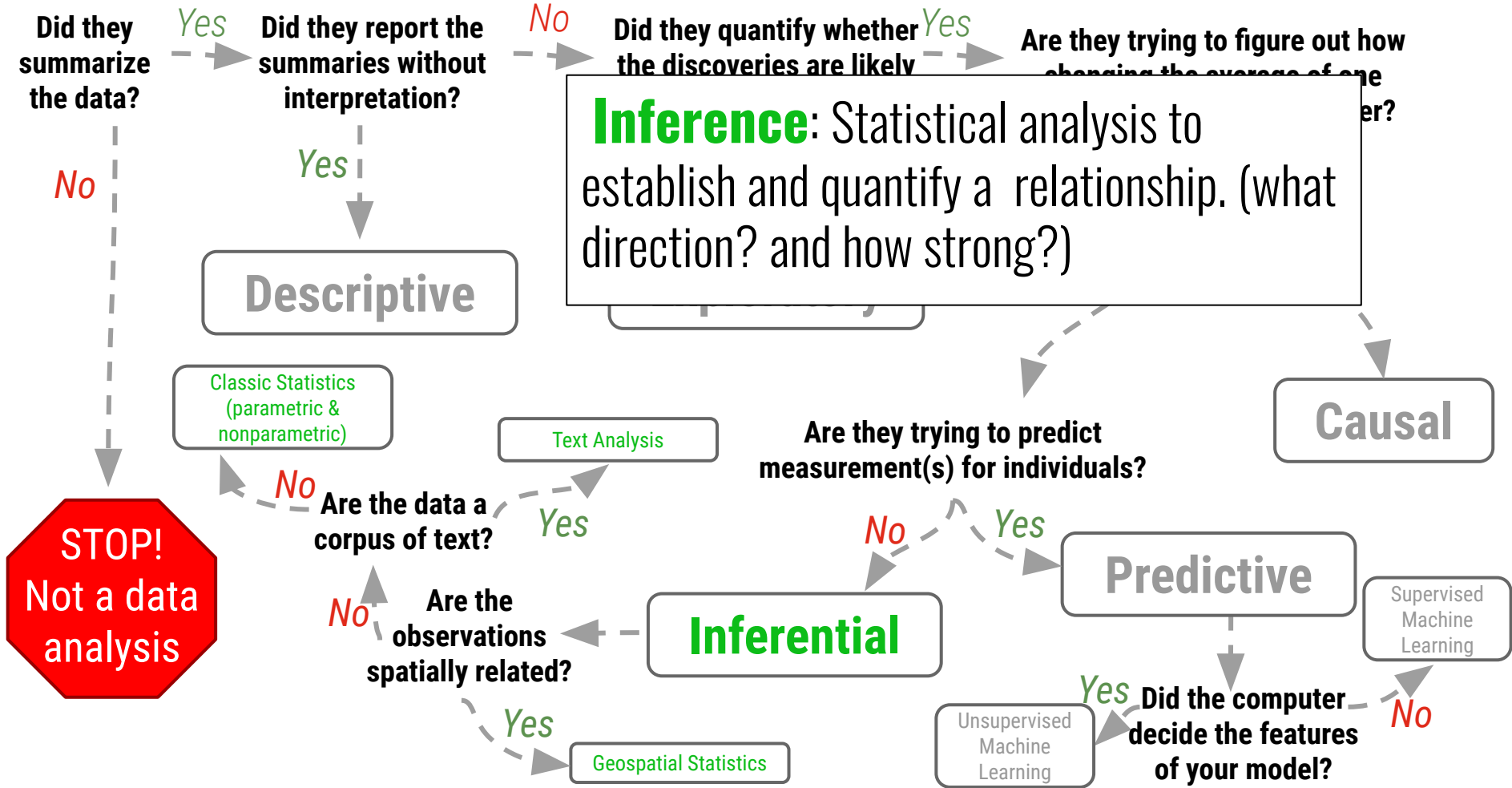**Department of Cognitive Science, UC San Diego**
jfleischer@ucsd.edu

@jasongfleischer

https://jgfleischer.com

**Did they summarize the data?**

*Yes* → **Did they report the summaries without interpretation?**

*No* → **Did they quantify whether the discoveries are likely to hold in a new sample?**

*Yes* → **Are they trying to figure out how changing the average of one measurement affects another?**

*No* (from "Did they summarize the data?") → **STOP! Not a data analysis**

*Yes* (from "Did they report the summaries without interpretation?") → **Descriptive**

*No* (from "Did they quantify whether the discoveries are likely to hold in a new sample?") → **Exploratory**

*No* (from "Are they trying to figure out how changing the average of one measurement affects another?") → **Are they trying to predict measurement(s) for individuals?**

*Yes* (from "Are they trying to figure out how changing the average...") → **Causal**

Classic Statistics (parametric & nonparametric)

Text Analysis

*No* → **Are the data a corpus of text?**

*Yes* → Text Analysis

*No* (from "Are they trying to predict measurement(s) for individuals?") → **Inferential**

*Yes* → **Predictive**

*No* → **Are the observations spatially related?**

*Yes* → Geospatial Statistics

**Are the observations spatially related?**

Supervised Machine Learning

*Yes* → Unsupervised Machine Learning

**Did the computer decide the features of your model?**

*No* → Supervised Machine Learning

**Did they summarize the data?**

*Yes* → **Did they report the summaries without interpretation?**

*No* → **Did they quantify whether the discoveries are likely**

*Yes* → **Are they trying to figure out how changing the average of one ~~changing the average of one~~ ...er?**

*No* (from first question) → **STOP! Not a data analysis**

*Yes* (from second question) → **Descriptive**

**Inference**: Statistical analysis to establish and quantify a relationship. (what direction? and how strong?)

**Causal**

Classic Statistics (parametric & nonparametric)

Text Analysis

*No* **Are the data a corpus of text?** *Yes*

**Are they trying to predict measurement(s) for individuals?**

*No* / *Yes*

**Predictive**

Supervised Machine Learning

*No* **Are the observations spatially related?**

**Inferential**

Geospatial Statistics

Unsupervised Machine Learning

*Yes* **Did the computer decide the features of your model?** *No*

*Yes* (from spatially related)

- **Problem:** Does Sesame Street affect kids brain development?
- **Data science question:** What is the relationship between watching Sesame Street and test scores among children?
- **Type of analysis:** Inferential analysis

# Establishing & Stating Your Null and Alternative Hypotheses Helps Guide Your Analysis

Null Hypothesis:

$H_0$: Sesame Street has *no effect* on kids brain development

Alternative Hypothesis:

$H_a$: Watching Sesame Street *has an effect* on kids' brain development

Population

Population

In our Sesame street example, the population would be all children

Population

Sample

Population

Sample

In our Sesame street example, the <u>sample</u> would be the children included in the study

Population

`¯\_(ツ)_/¯`

Sample

Population

Sample

We don't know how much Sesame street was watched by or the tests scores of all kids

¯\_(ツ)_/¯

Based on the relationship we see in our sample, we can <u>infer</u> the answer to our question in our population

Population

Sample

Inference!

Population

So we look at Sesame street viewing and test scores in a <u>representative sample</u> of kids

Sample

Inference!

Population

Sample

# Approaches to Inference

**CORRELATION**

ASSOCIATION
BETWEEN VARIABLES

i.e. Pearson Correlation,
Spearman Correlation,
chi-square test

**COMPARISON OF MEANS**

DIFFERENCE IN MEANS
BETWEEN VARIABLES
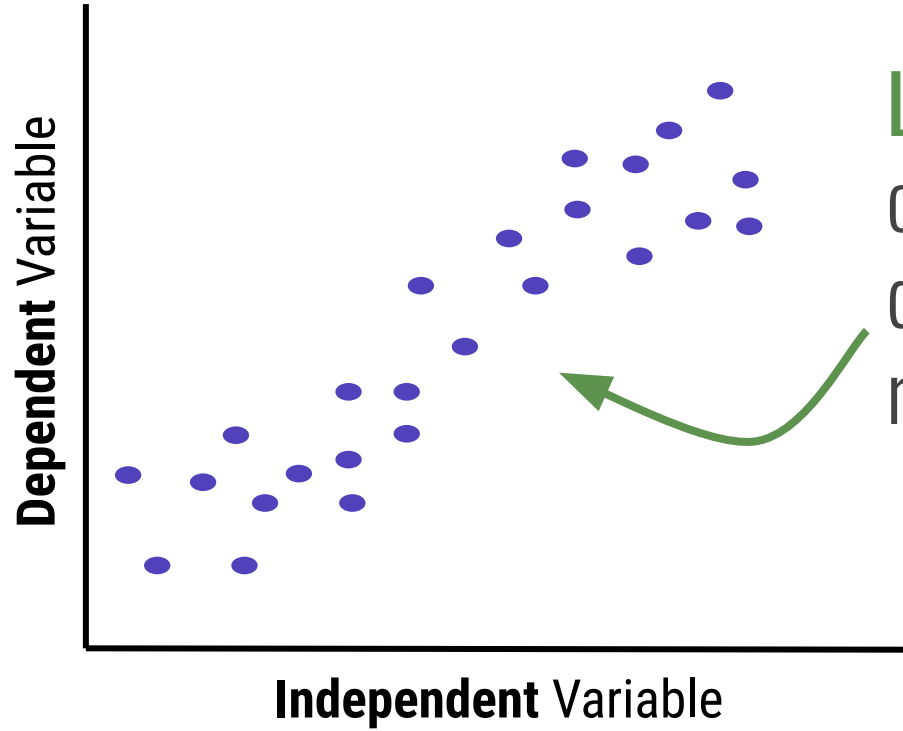
i.e. t-test, ANOVA

**REGRESSION**

DOES CHANGE IN ONE
VARIABLE MEAN CHANGE IN
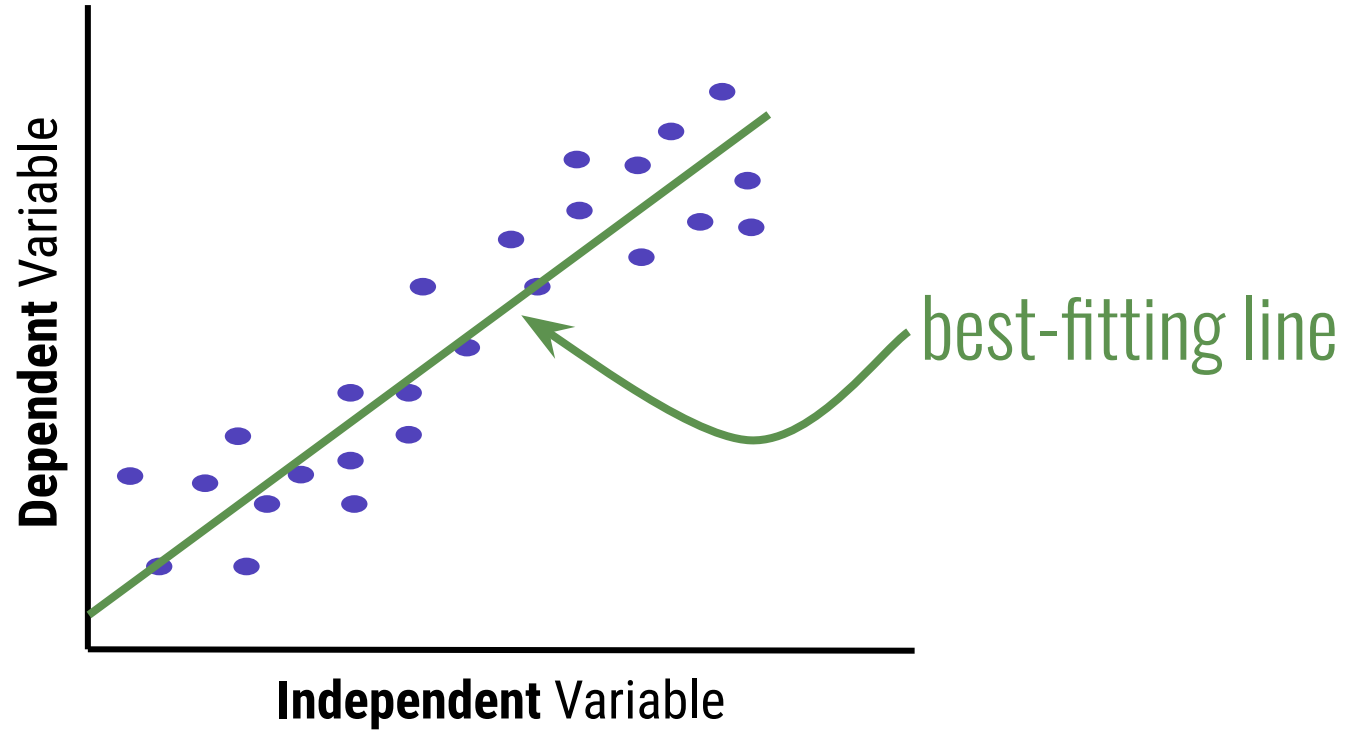ANOTHER?

I.e. simple regression,
multiple regression

**NON-PARAMETRIC TESTS**

FOR WHEN ASSUMPTIONS IN
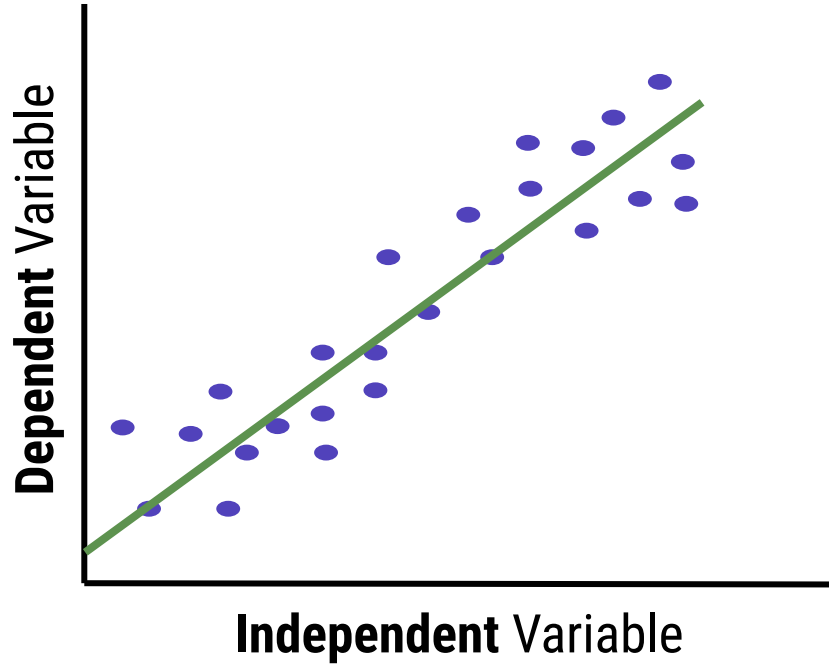THESE OTHER 3 CATEGORIES
ARE NOT MET

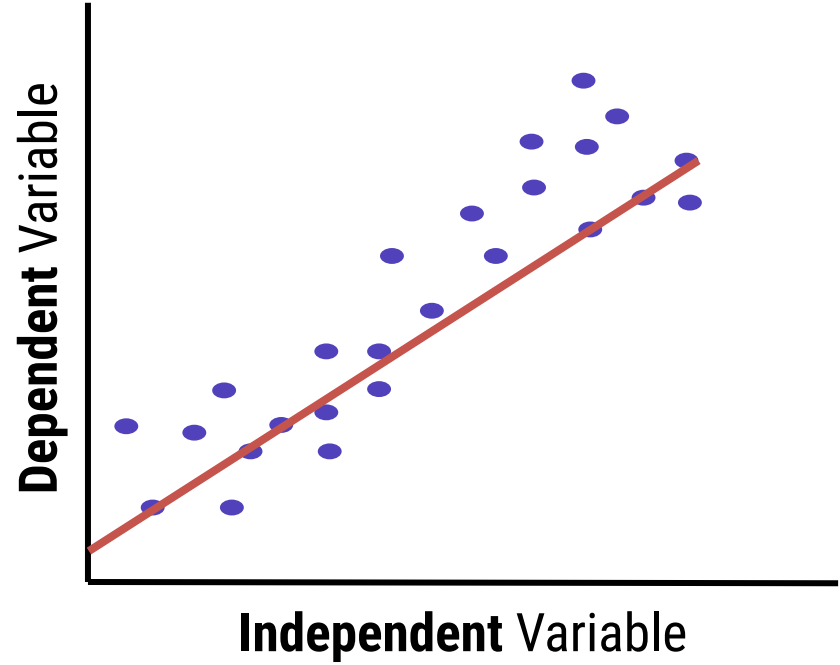i.e. Wilcoxon rank-sum
test, Wilcoxon sign-rank
test, sign test

**CORRELATION**

ASSOCIATION
BETWEEN VARIABLES

i.e. Pearson Correlation,
Spearman Correlation,
chi-square test

**COMPARISON OF MEANS**

DIFFERENCE IN MEANS
BETWEEN VARIABLES

i.e. t-test, ANOVA

**REGRESSION**

DOES CHANGE IN ONE
VARIABLE MEAN CHANGE IN
ANOTHER?

I.e. simple regression,
multiple regression

**NON-PARAMETRIC TESTS**

FOR WHEN ASSUMPTIONS IN
THESE OTHER 3 CATEGORIES
ARE NOT MET

i.e. Wilcoxon rank-sum
test, Wilcoxon sign-rank
test, sign test

Best-fitting line

NOT a best-fitting line

**Dependent** Variable

**Independent** Variable

**Dependent** Variable

**Independent** Variable
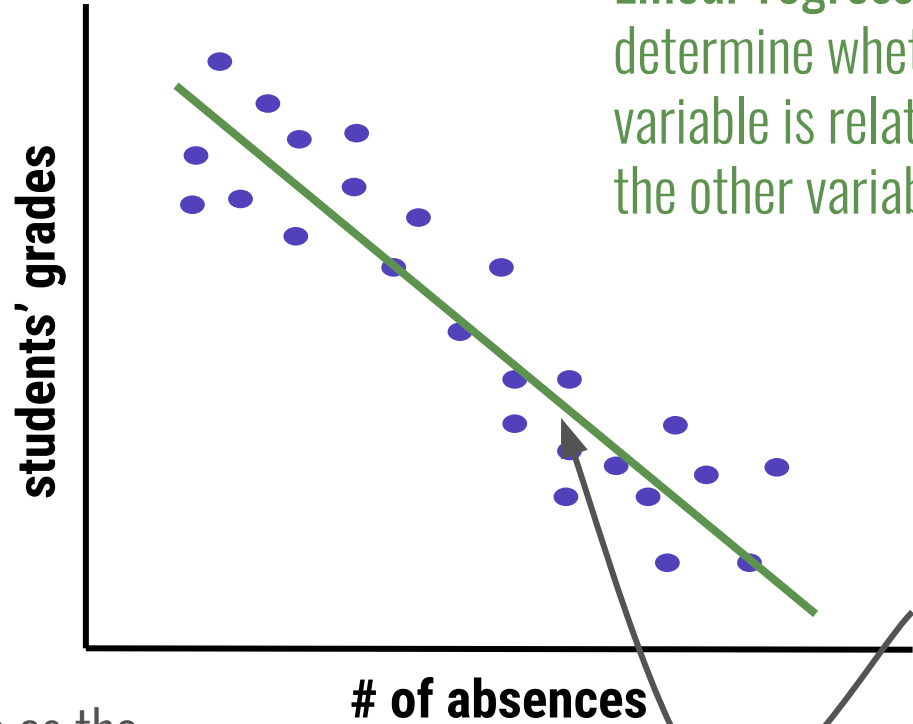
*"All models are wrong, but some are useful"*

-George Box (British Statistician, *JASA* 1976)

**Linear regression** can be used to determine whether a change in one variable is related to the change in the other variable
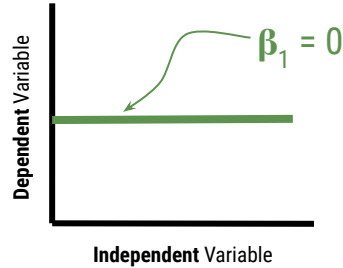
**Linear regression** can be used to determine whether a change in one variable is related to the change in the other variable

**students' grades**

**# of absences**

The <u>magnitude of the relationship</u> is measured by the <u>slope</u> of the line

**Linear regression** can be used to determine whether a change in one variable is related to the change in the other variable
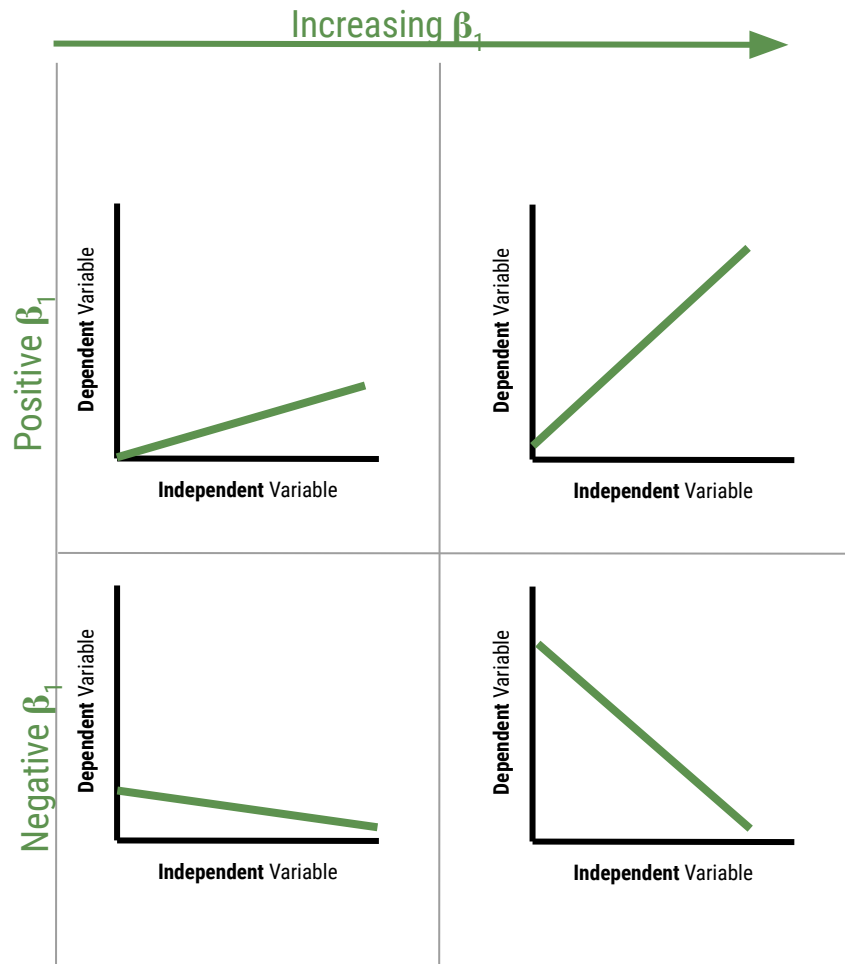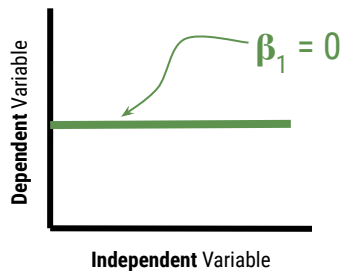
students' grades

# of absences

The <u>magnitude of the relationship</u> is measured by the <u>slope</u> of the line

This is also referred to as the model's <u>effect size</u> ($\beta_1$)
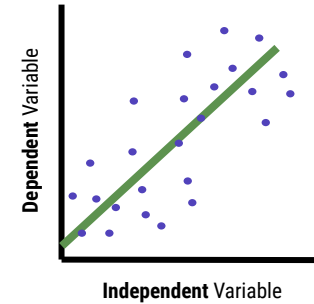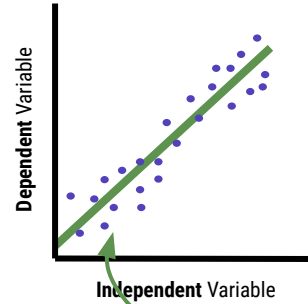
# Effect size (β₁) can be estimated using the slope of the line

Dependent Variable

Independent Variable

Dependent Variable

Independent Variable
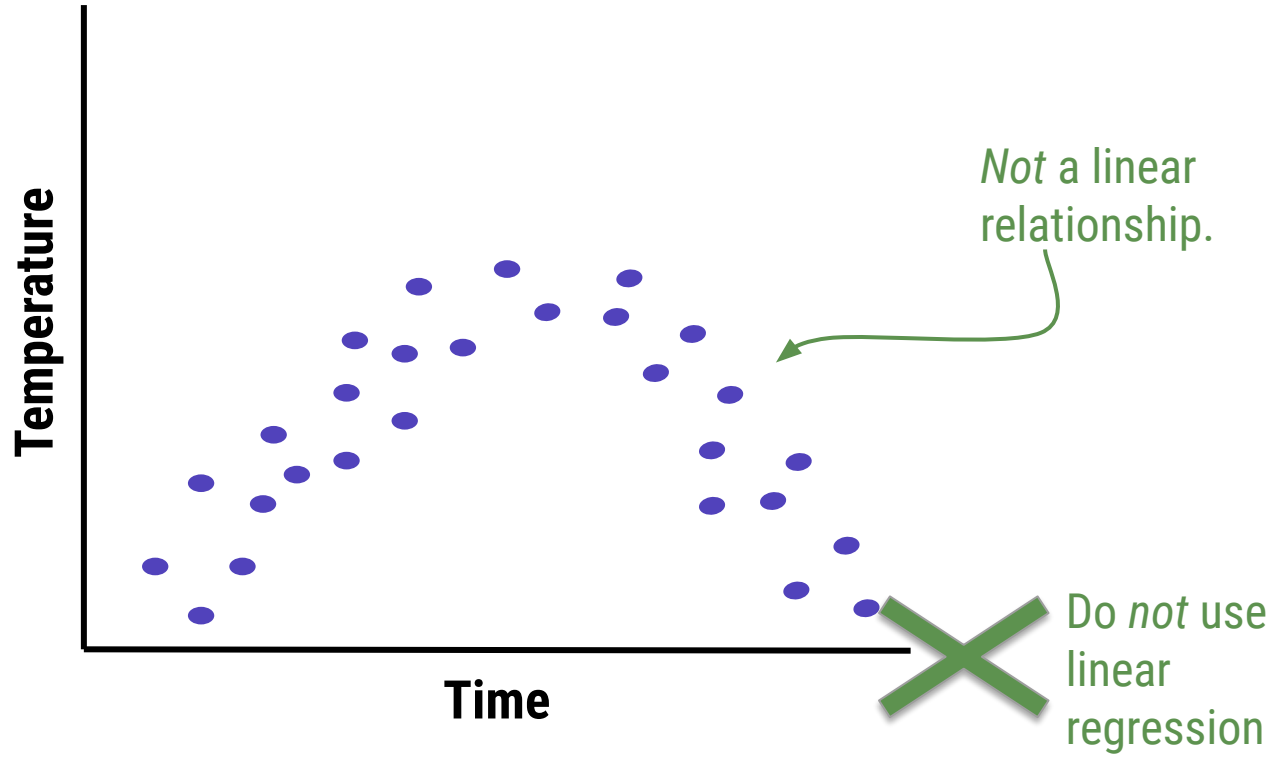
The *closer* the points
are to the regression
line, the *less uncertain*
we are in our estimate
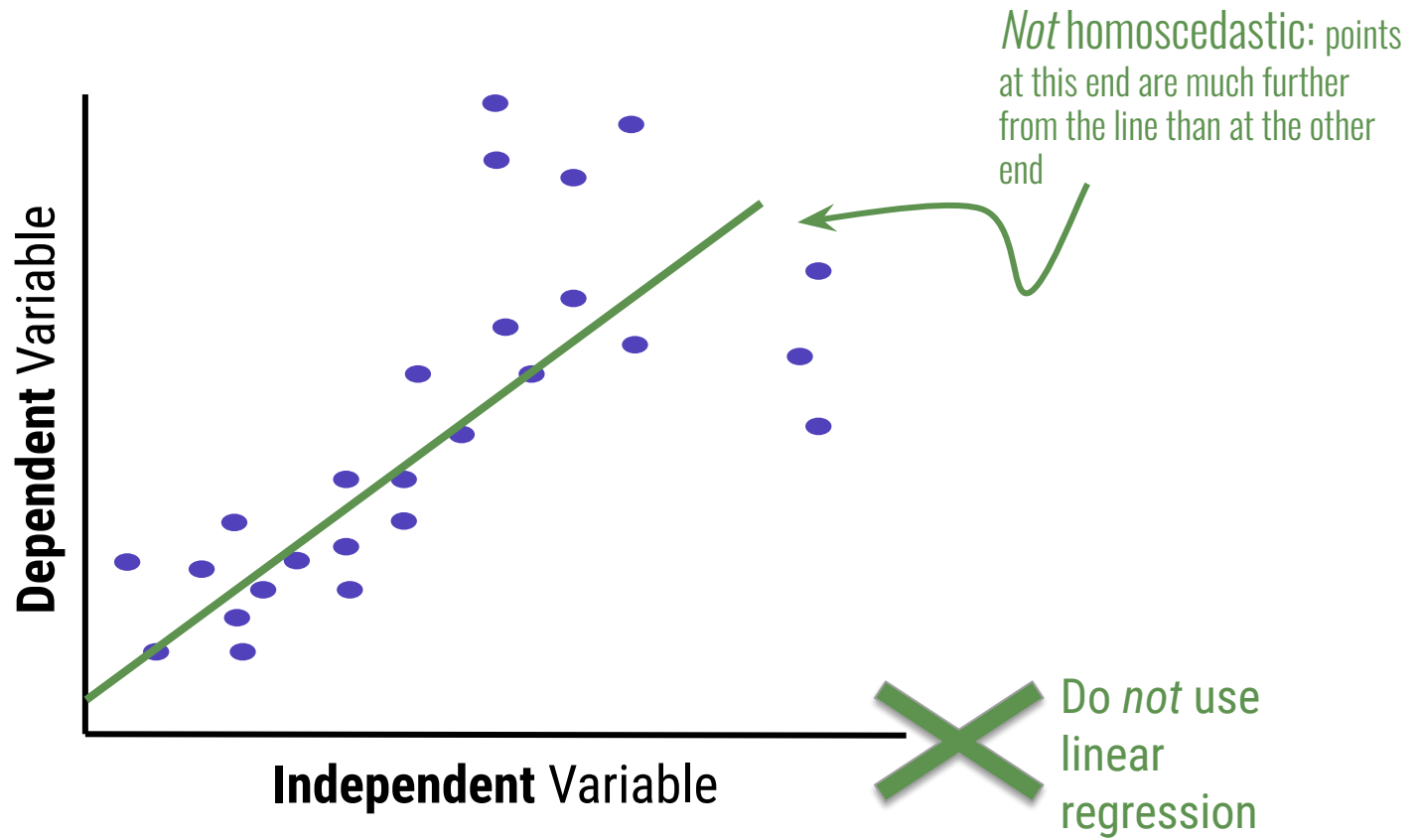
# Assumptions of linear regression

1. Linear relationship
2. No multicollinearity
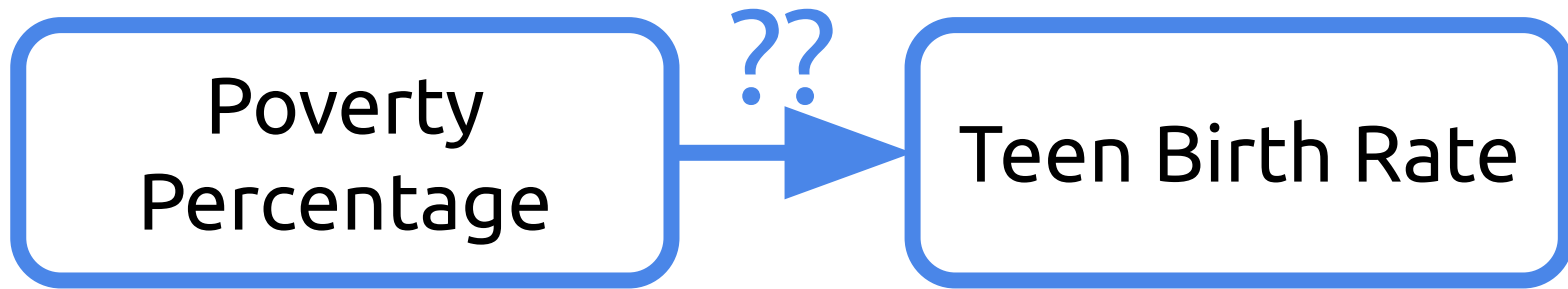3. No auto-correlation
4. Homoscedasticity

Linear regression assumes no multicollinearity. **Multicollinearity** occurs when the independent variables (in multiple linear regression) are too highly correlated with each other.

Autocorrelation occurs when the observations are *not* independent of one another (i.e. stock prices)

*Not* homoscedastic: points at this end are much further from the line than at the other end

**Dependent** Variable

**Independent** Variable

Do *not* use linear regression

# Does Poverty Percentage affect Teen Birth Rate?

Poverty Percentage ?? → Teen Birth Rate

Null Hypothesis:

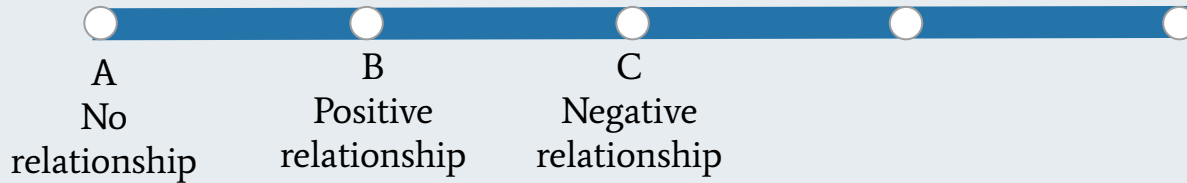$H_0$: Poverty Rate does not affect Teen Birth Rate ($\beta_1 = 0$)

Alternative Hypothesis:

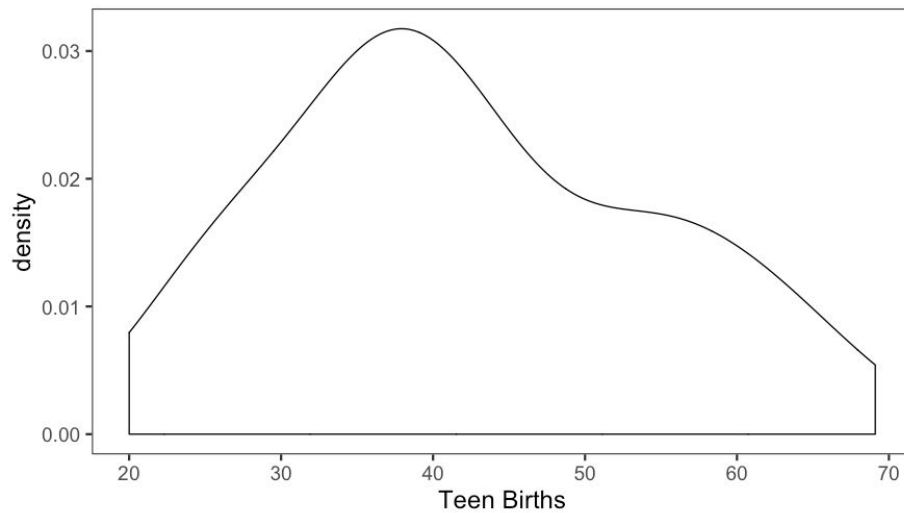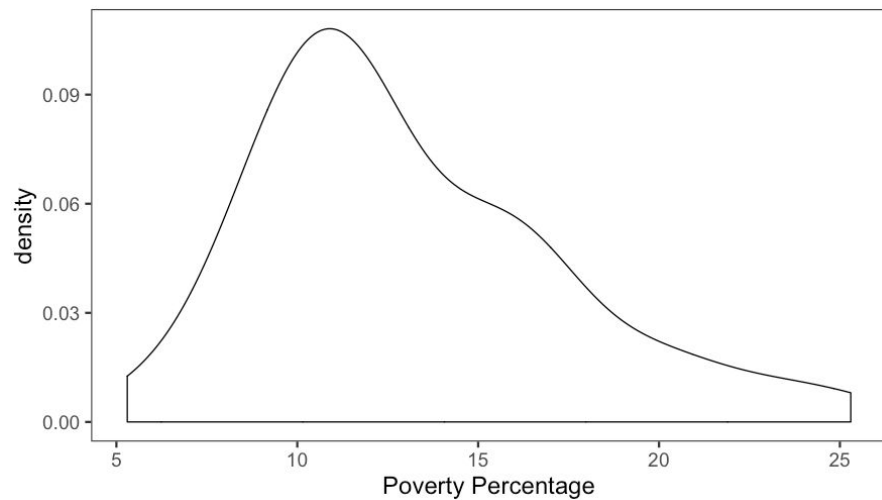$H_a$: Poverty Rate affects Teen Birth Rate ($\beta_1 \neq 0$)

| | Location | PovPct | Brth15to17 | Brth18to19 | ViolCrime | TeenBrth |
|---|---|---|---|---|---|---|
| 1 | Alabama | 20.1 | 31.5 | 88.7 | 11.2 | 54.5 |
| 2 | Alaska | 7.1 | 18.9 | 73.7 | 9.1 | 39.5 |
| 3 | Arizona | 16.1 | 35.0 | 102.5 | 10.4 | 61.2 |
| 4 | Arkansas | 14.9 | 31.6 | 101.7 | 10.4 | 59.9 |
| 5 | California | 16.7 | 22.6 | 69.1 | 11.2 | 41.1 |
| 6 | Colorado | 8.8 | 26.2 | 79.1 | 5.8 | 47.0 |
| 7 | Connecticut | 9.7 | 14.1 | 45.1 | 4.6 | 25.8 |
| 8 | Delaware | 10.3 | 24.7 | 77.8 | 3.5 | 46.3 |
| 9 | District_of_Columbia | 22.0 | 44.8 | 101.5 | 65.0 | 69.1 |
| 10 | Florida | 16.2 | 23.2 | 78.4 | 7.3 | 44.5 |
| 11 | Georgia | 12.1 | 31.4 | 92.8 | 9.5 | 55.7 |
| 12 | Hawaii | 10.3 | 17.7 | 66.4 | 4.7 | 38.2 |
| 13 | Idaho | 14.5 | 18.4 | 69.1 | 4.1 | 39.1 |
| 14 | Illinois | 12.4 | 23.4 | 70.5 | 10.3 | 42.2 |
| 15 | Indiana | 9.6 | 22.6 | 78.5 | 8.0 | 44.6 |
| 16 | Iowa | 12.2 | 16.4 | 55.4 | 1.8 | 32.5 |
| 17 | Kansas | 10.8 | 21.4 | 74.2 | 6.2 | 43.0 |

# EDA: distributions

1. Linear relationship
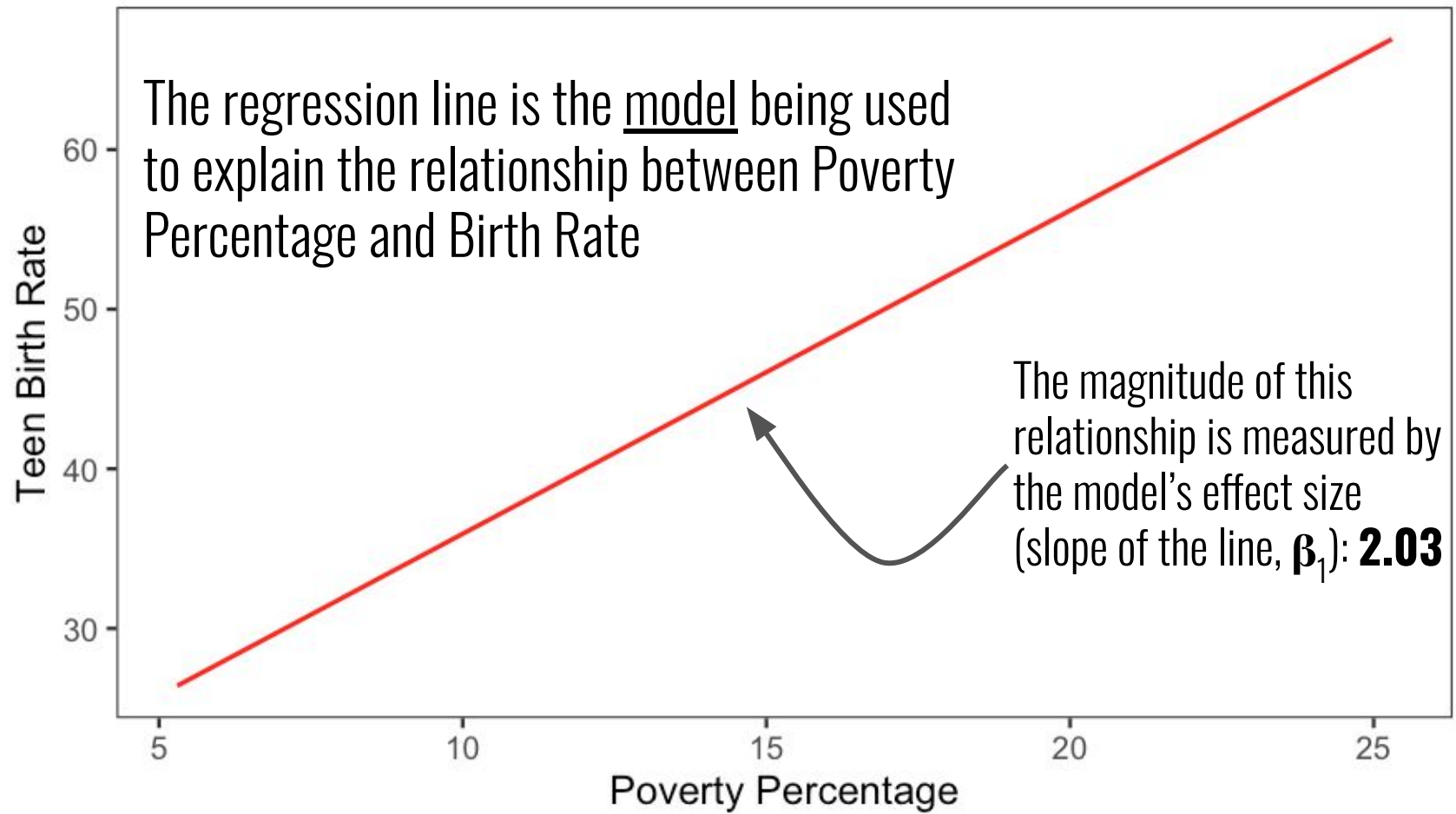2. ~~No multicollinearity~~
3. No autocorrelation
4. Homoscedasticity
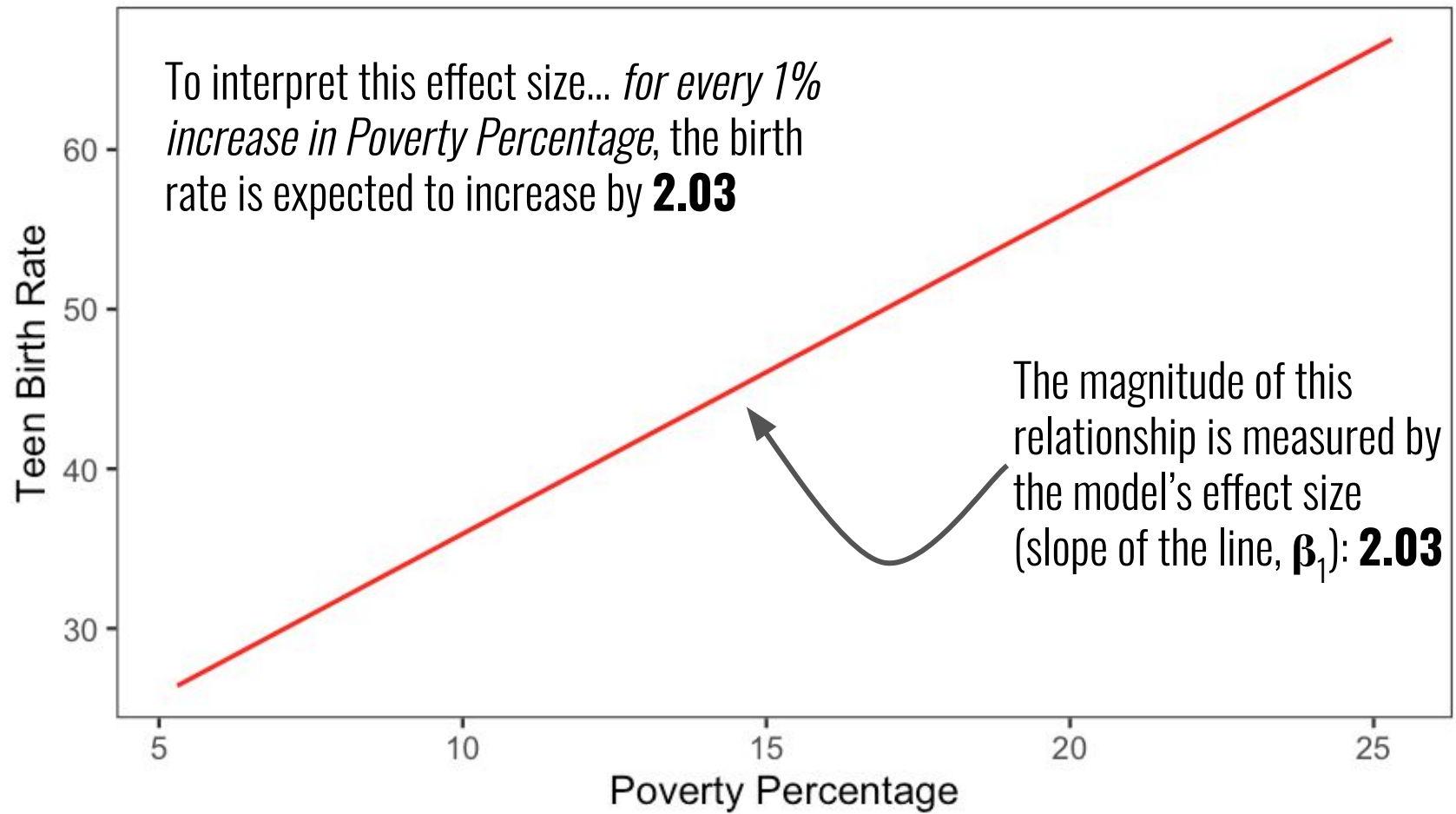
...but *how confident* are we in that estimate of the effect size?

For that...we need to look at our standard error (SE)

$\beta_1 \pm / \text{SE}$
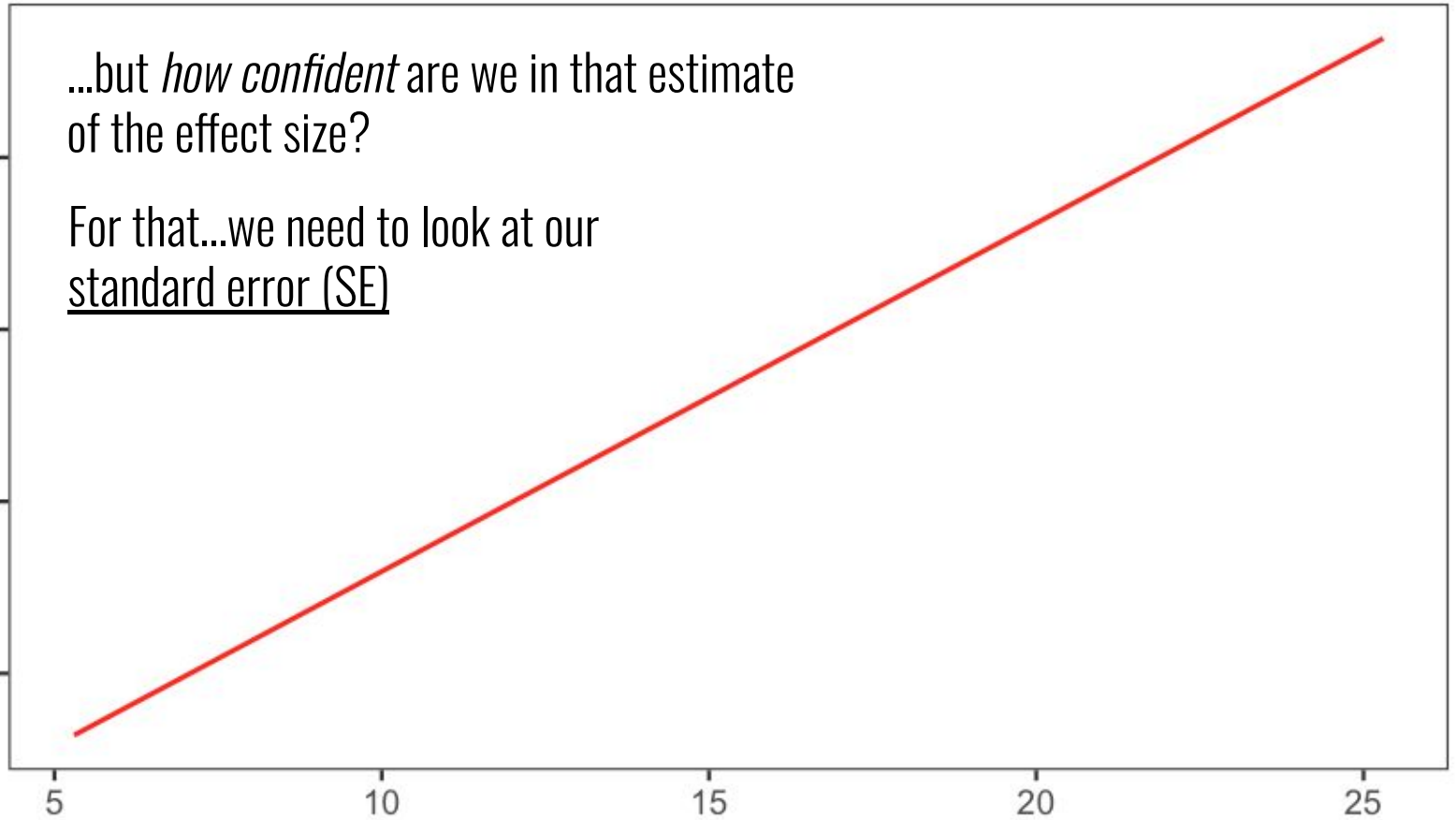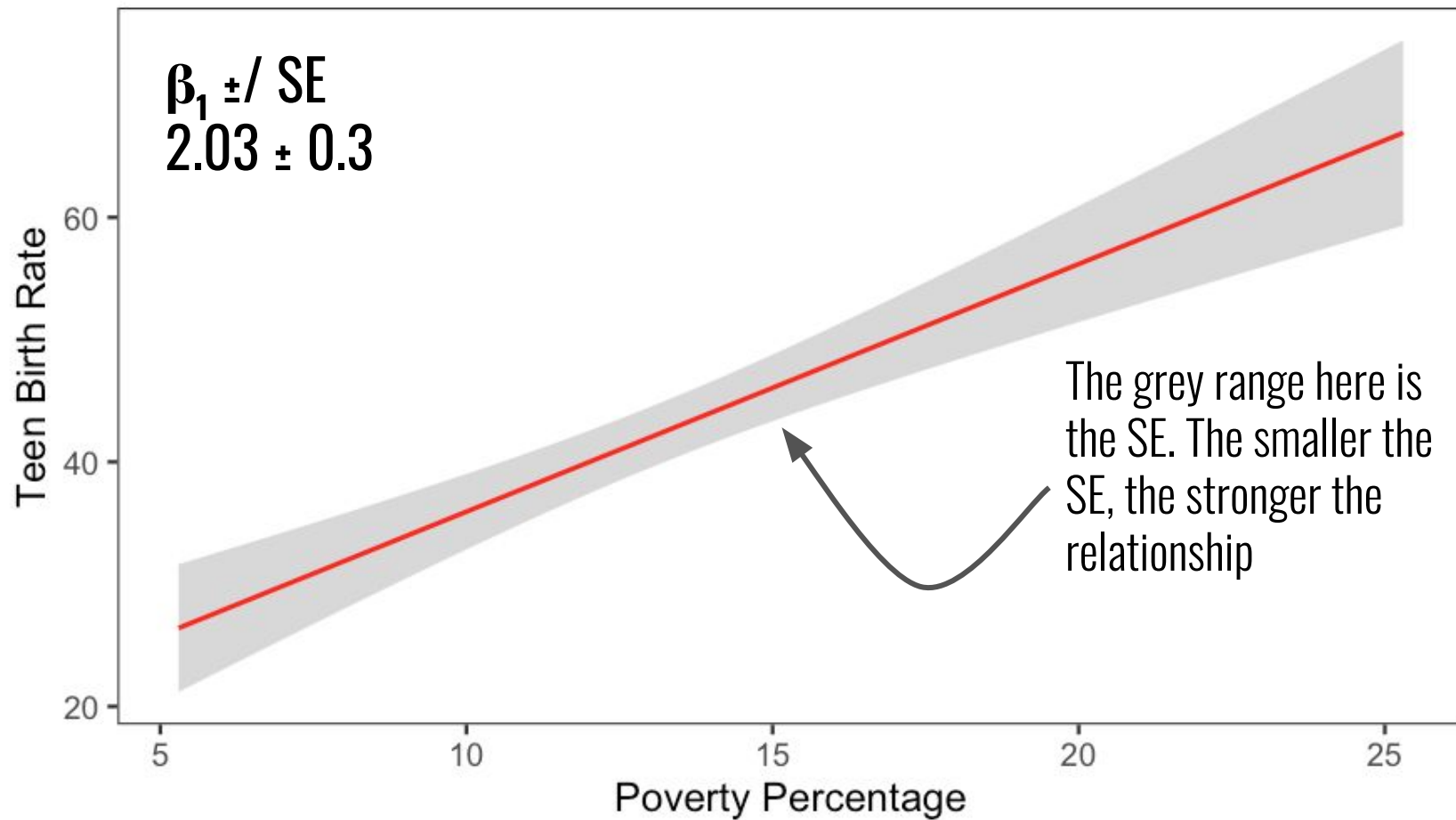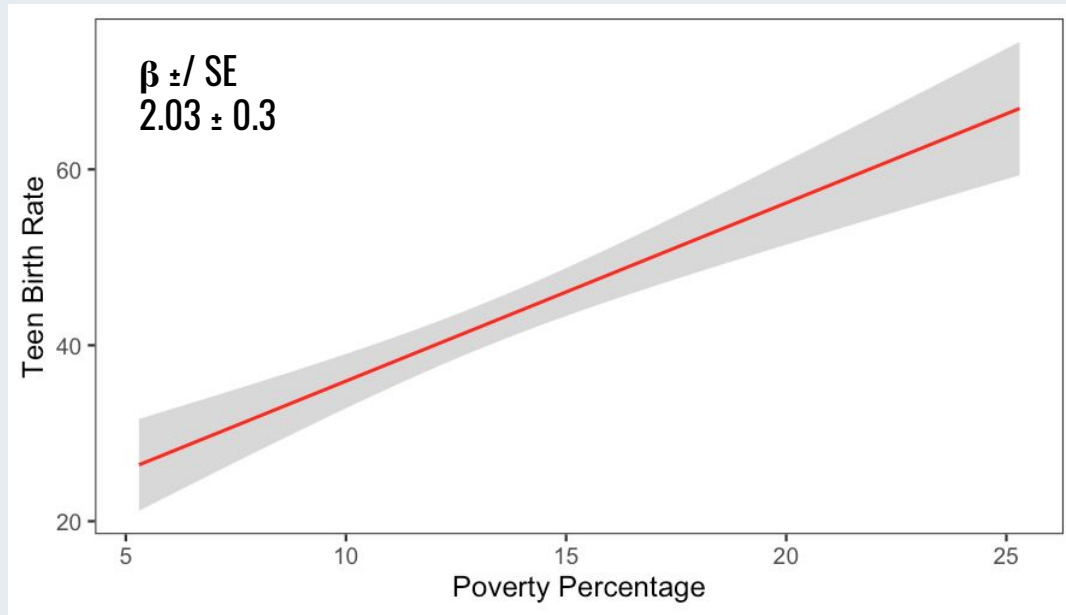$2.03 \pm 0.3$

Teen Birth Rate

60

40

20

Poverty Percentage

5    10    15    20    25

The grey range here is the SE. The smaller the SE, the stronger the relationship
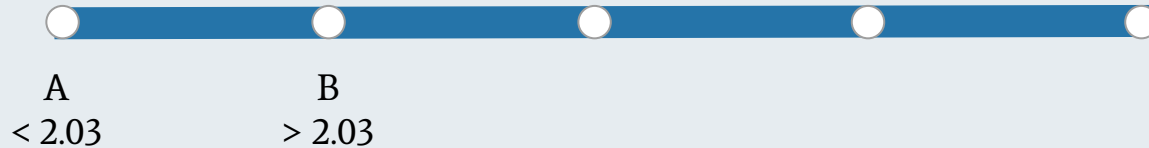
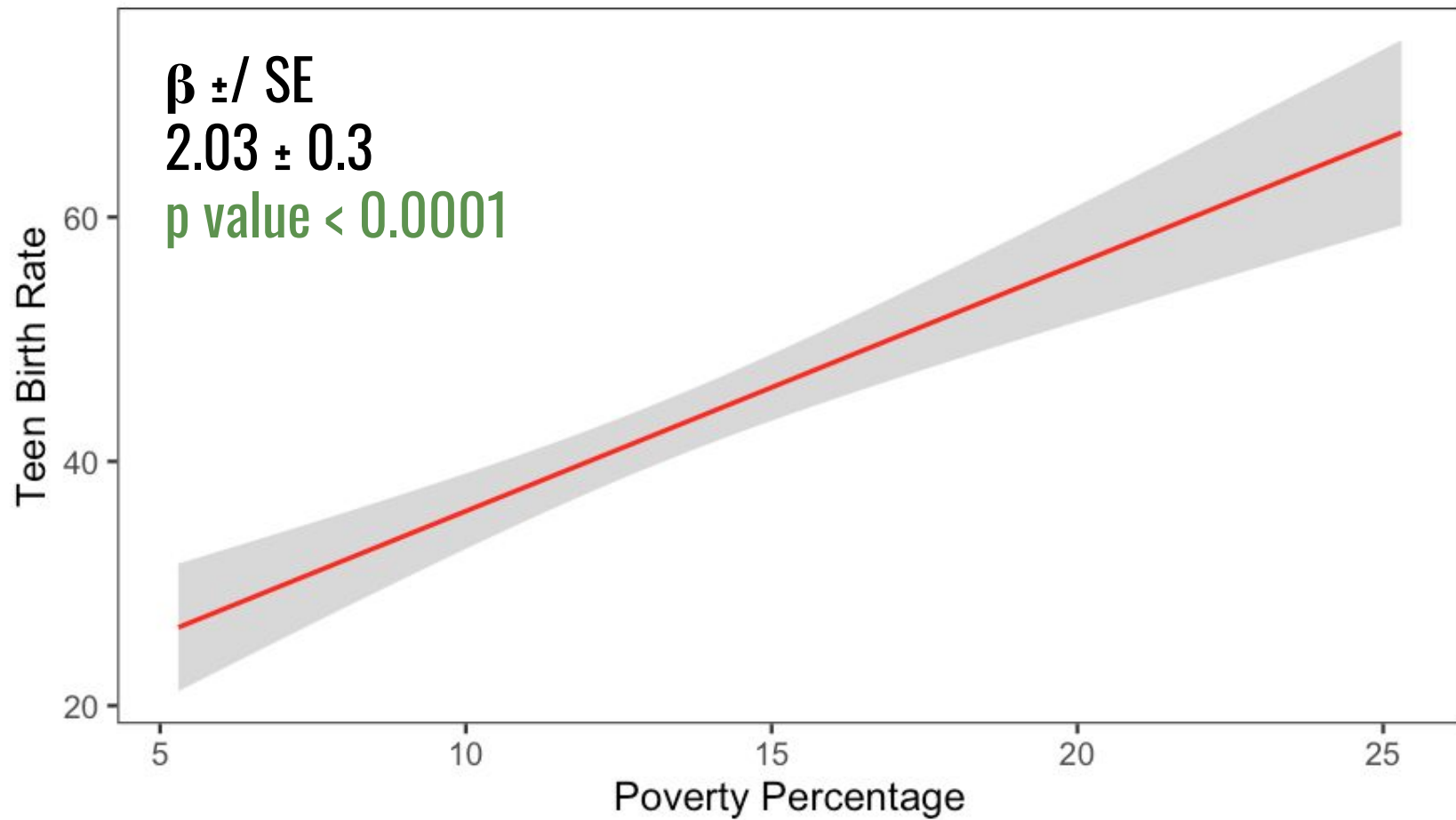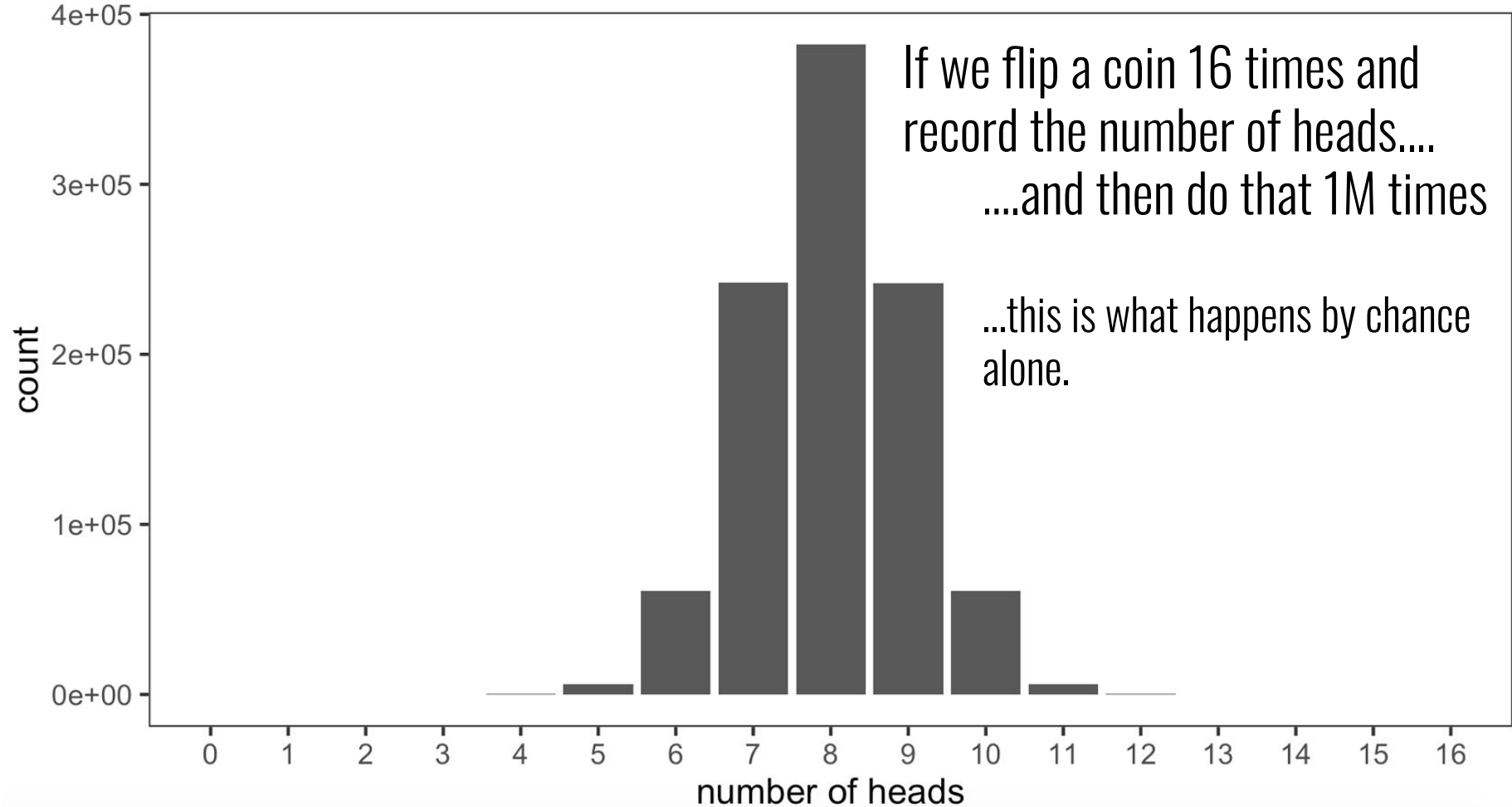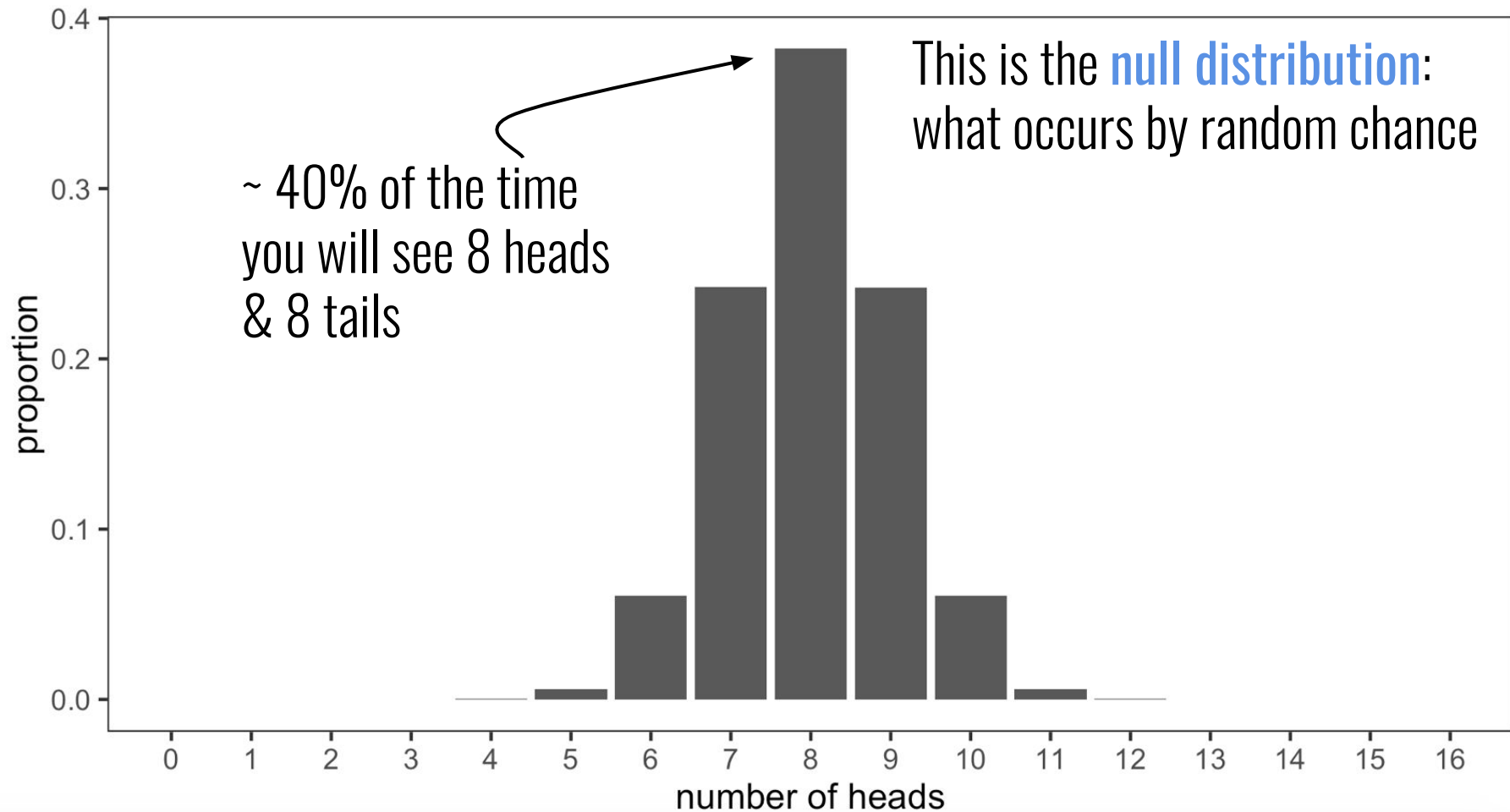If there were a stronger effect of Poverty on Birth rate, what would $\beta_1$ be?
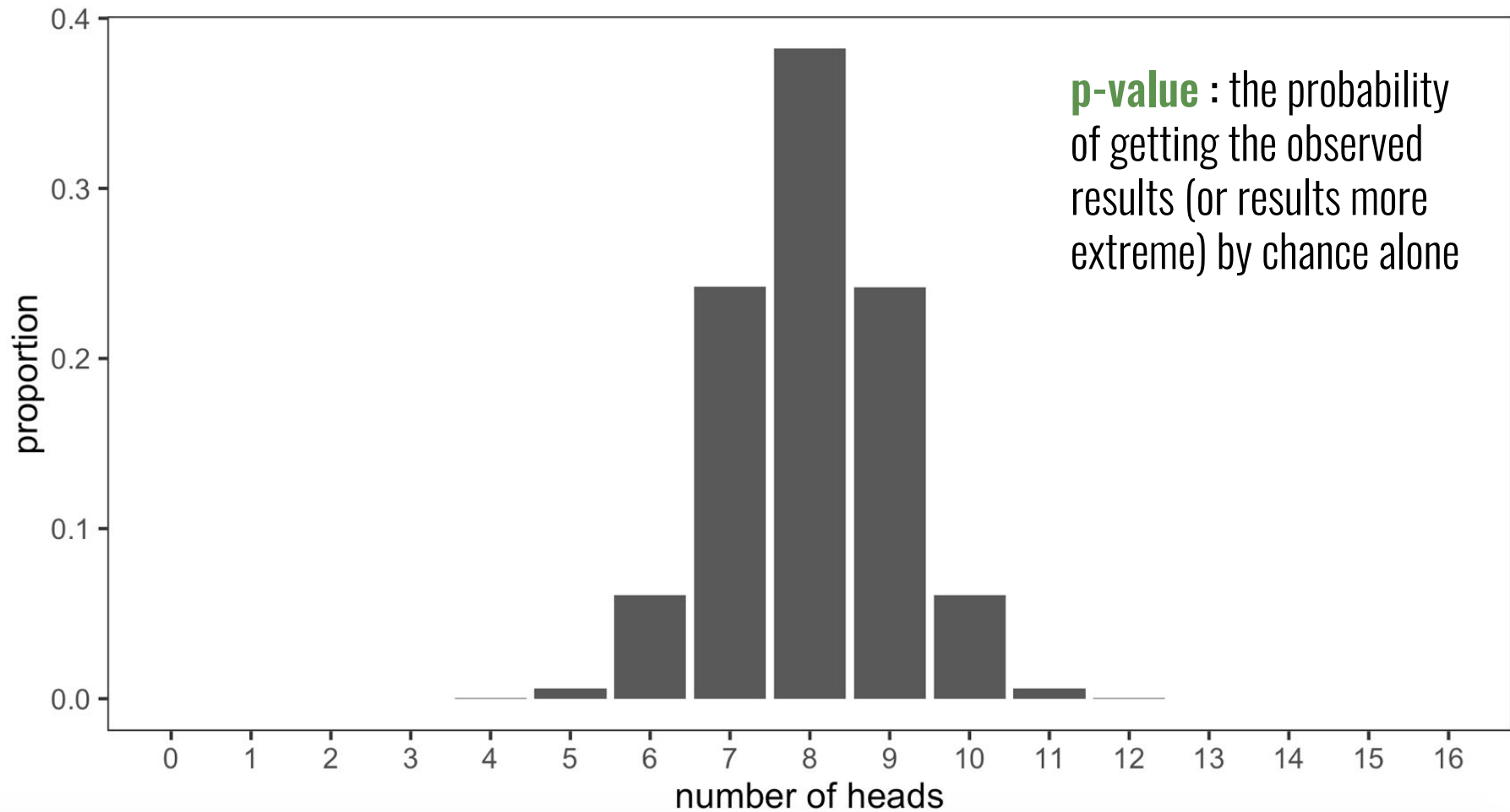
A
< 2.03

B
> 2.03

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

If we flip a coin 16 times and record the number of heads....
....and then do that 1M times
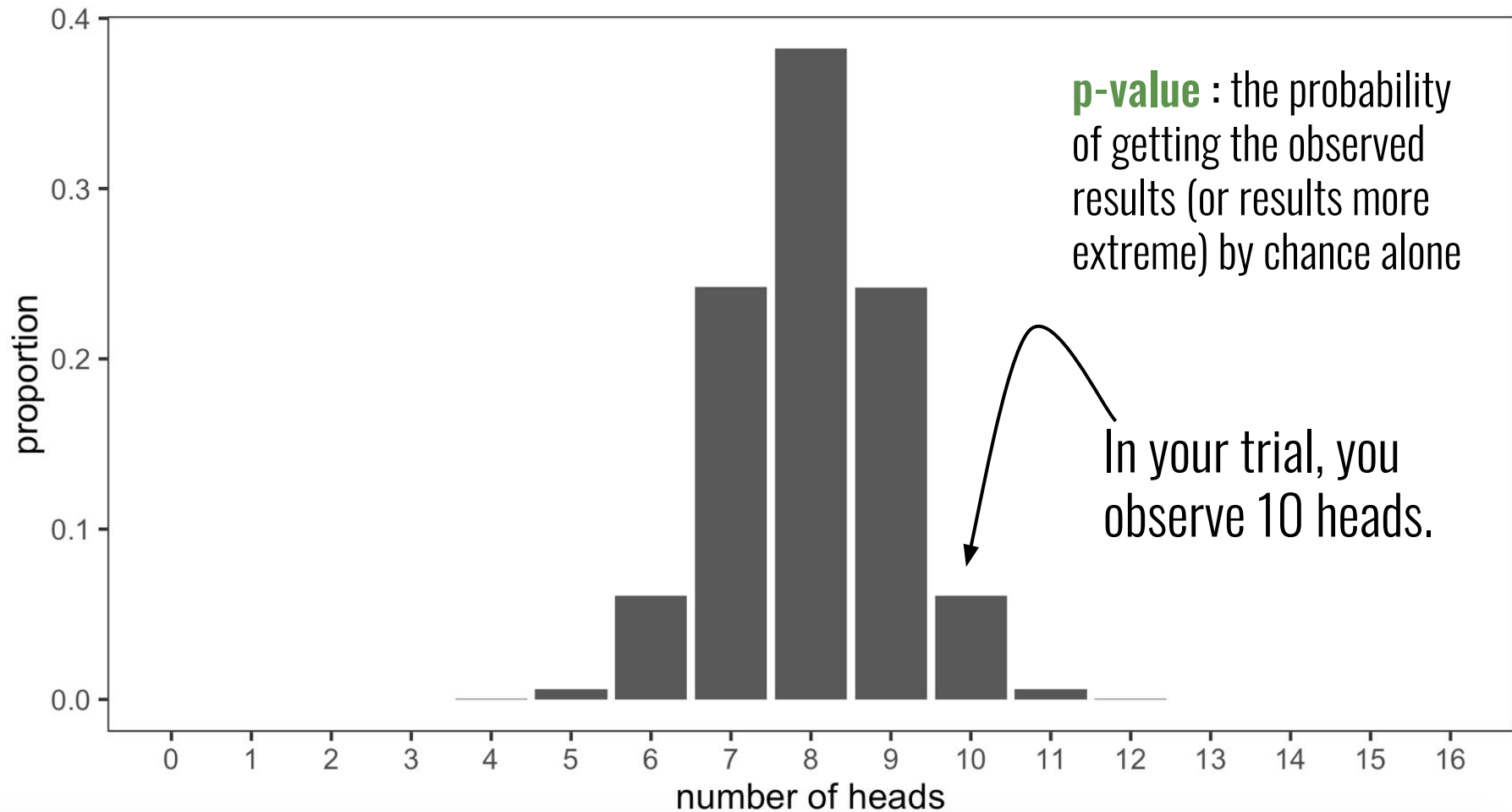
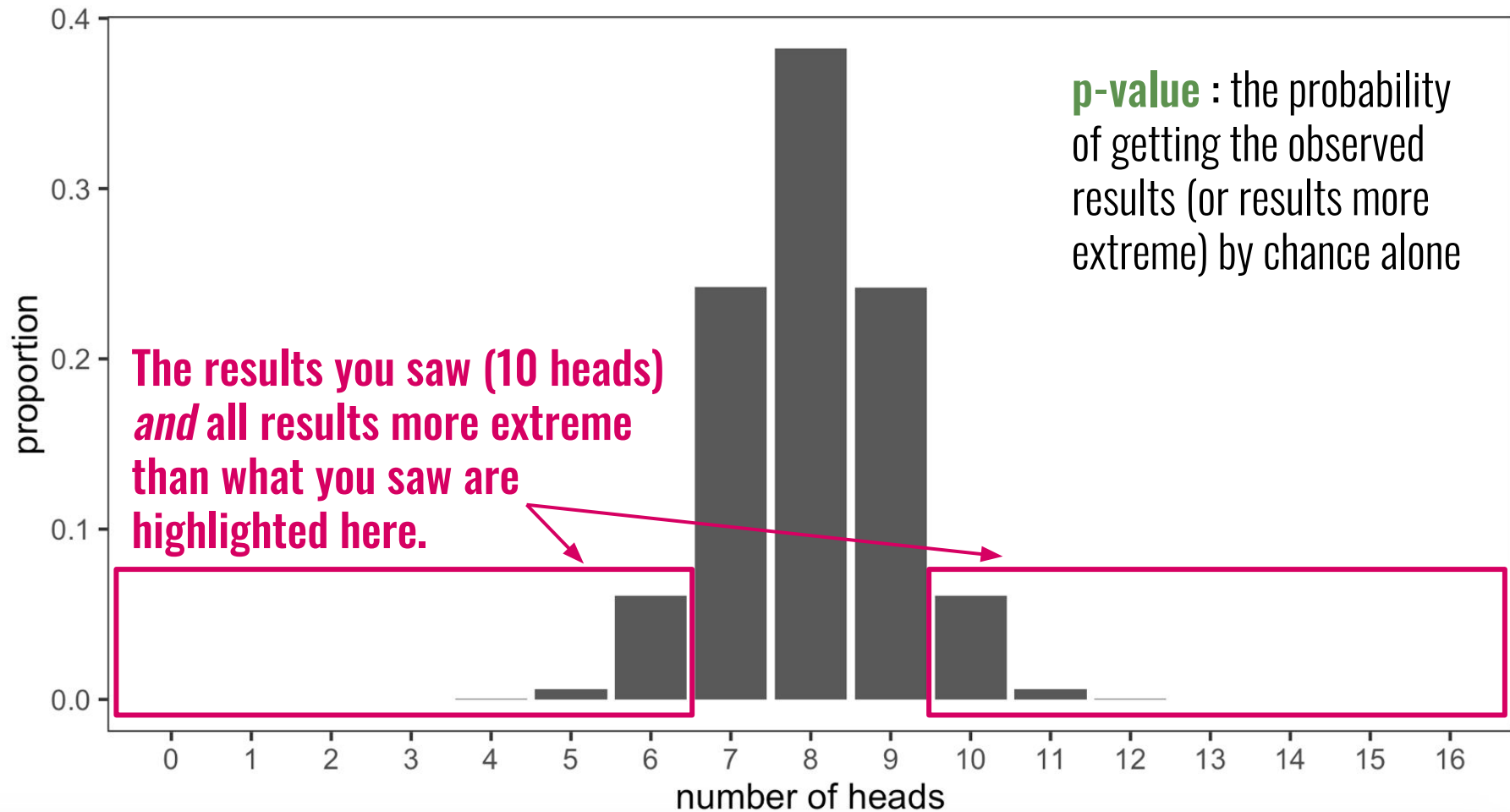...this is what happens by chance alone.

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

In your trial, you observe 10 heads.

The probability of getting 10 heads *or something more extreme* is

# of 10 or more extreme flips / total flips

( 2 + 218 + 5,877 + 60,731 + 60,766 + 5,973 + 208 + 2 ) / $1 \times 10^6$

= 133,777 / $1 \times 10^6$

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone
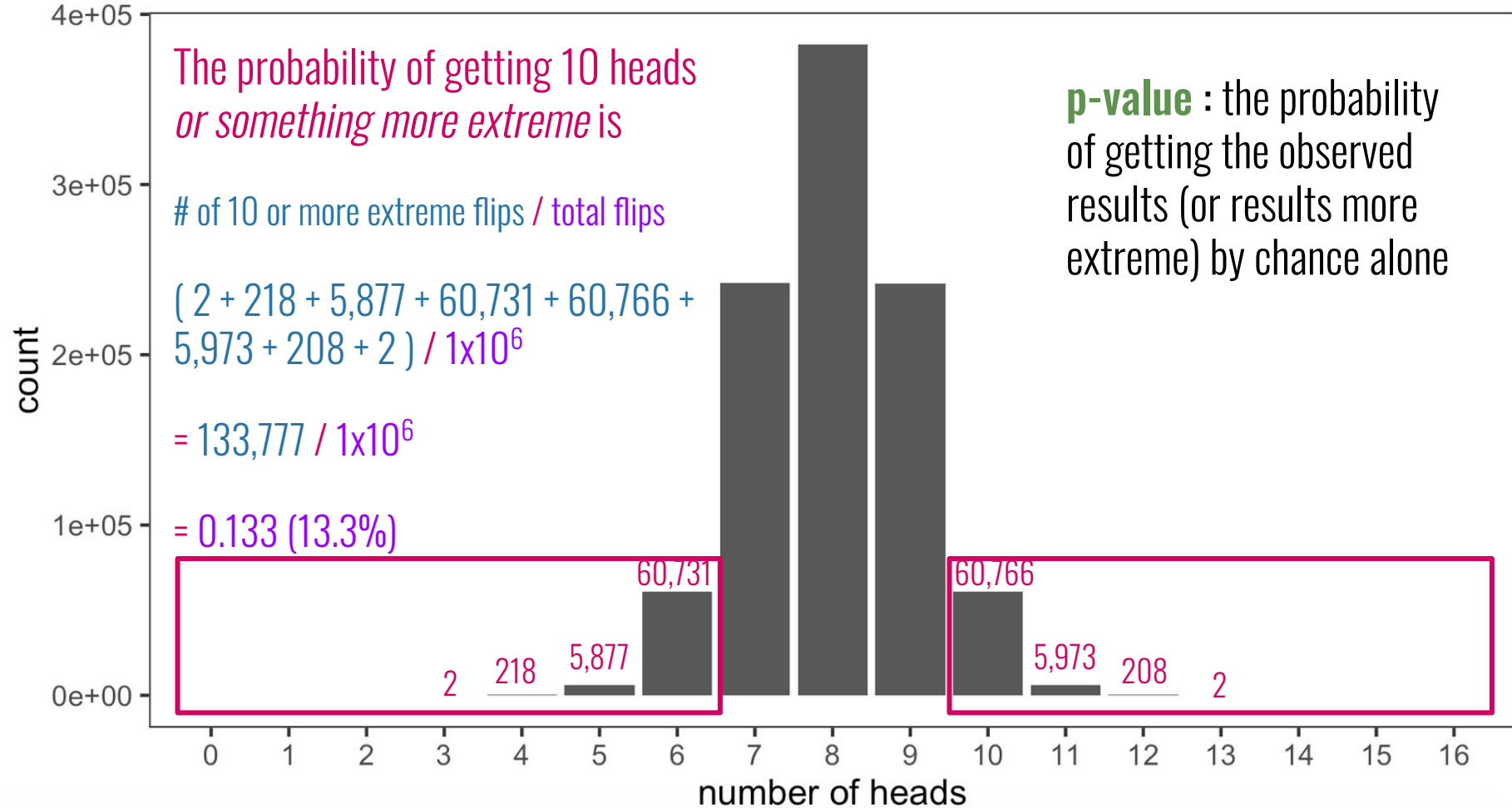
The probability of getting 10 heads *or something more extreme* is

# of 10 or more extreme flips / total flips

( 2 + 218 + 5,877 + 60,731 + 60,766 + 5,973 + 208 + 2 ) / $1 \times 10^6$

= 133,777 / $1 \times 10^6$

= 0.133 (13.3%)

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

The probability of getting 10 heads *or something more extreme* is

\# of 10 or more extreme flips / total flips

( 2 + 218 + 5,877 + 60,731 + 60,766 + 5,973 + 208 + 2 ) / $1 \times 10^6$

= 133,777 / $1 \times 10^6$

= 0.133 (13.3%)

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

**p-value** : 0.133

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

What if you observed 16 heads??

**What would be the p-value of you flipping 16 heads?**

A
< 0.13

B
> 0.13

β ±/ SE
2.03 ± 0.03
p value < 0.0001

the probability of getting the observed results (or results more extreme) by chance alone

Takes into account the effect size ($\beta_1$) and the SE

**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

# Confounding

We'll discuss additional
approaches of how to account
for confounding in your analysis
in the next lecture.

Ignoring confounders will lead
you to draw incorrect conclusions
from your analyses

# Spine Surgery Results

**Sample:** 400 patients with index vertebral fractures

| **Vertebroplasty** | **Conservative care** | **Relative risk (95% confidence interval)** |
| --- | --- | --- |
| 30/200 (15%) | 15/200 (7.5%) | 2.0 (1.1–3.6) |

subsequent fractures

Eek....looks like vertebroplasty was *way* worse for patients!

# But wait...at time of initial fracture...

| | Vertebroplasty N = 200 | Conservative care N = 200 |
|---|---|---|
| Age, y, mean ± SD | 78.2 ± 4.1 | 79.0 ± 5.2 |
| Weight, kg, mean ± SD | 54.4 ± 2.3 | 53.9 ± 2.1 |
| Smoking status, No. (%) | 110 (55) | 16 (8) |

Age and weight are similar between groups. **Smoking Status** differs vastly.

# So...let's stratify those results real quick

| Smoke | | | | No smoke | | |
|---|---|---|---|---|---|---|
| Vertebroplasty | Conservative | RR (95% confidence interval) | | Vertebroplasty | Conservative | RR (95% confidence interval) |
| 23/110 (21%) | 3/16 (19%) | 1.1 (0.4, 3.3) | | 7/90 (8%) | 12/184(7%) | 1.2 (0.5, 2.9) |

Risk of re-fracture is now
similar within group

# What are possible confounders for our analysis of the effect of poverty on teen birth rate?

A
I have some ideas

B
Not sure