

# Course Reminders

## Due Dates:

- A2 due Sunday (11:59 PM)
- A1 grades will be announced on Piazza when posted
- Project proposal feedback should be released monday

# EDA Case Study Review

1. Figure out specific question
2. Determine what information you'd want/need
3. Identify datasets
4. Describe them
  - a. Size of dataset
  - b. missingness
  - c. Quantitative & categorical variables
    - i. Center + variability
    - ii. distribution/shape
5. Explore!
  - a. Relationship between variables
  - b. Lots of plots

## Communication:

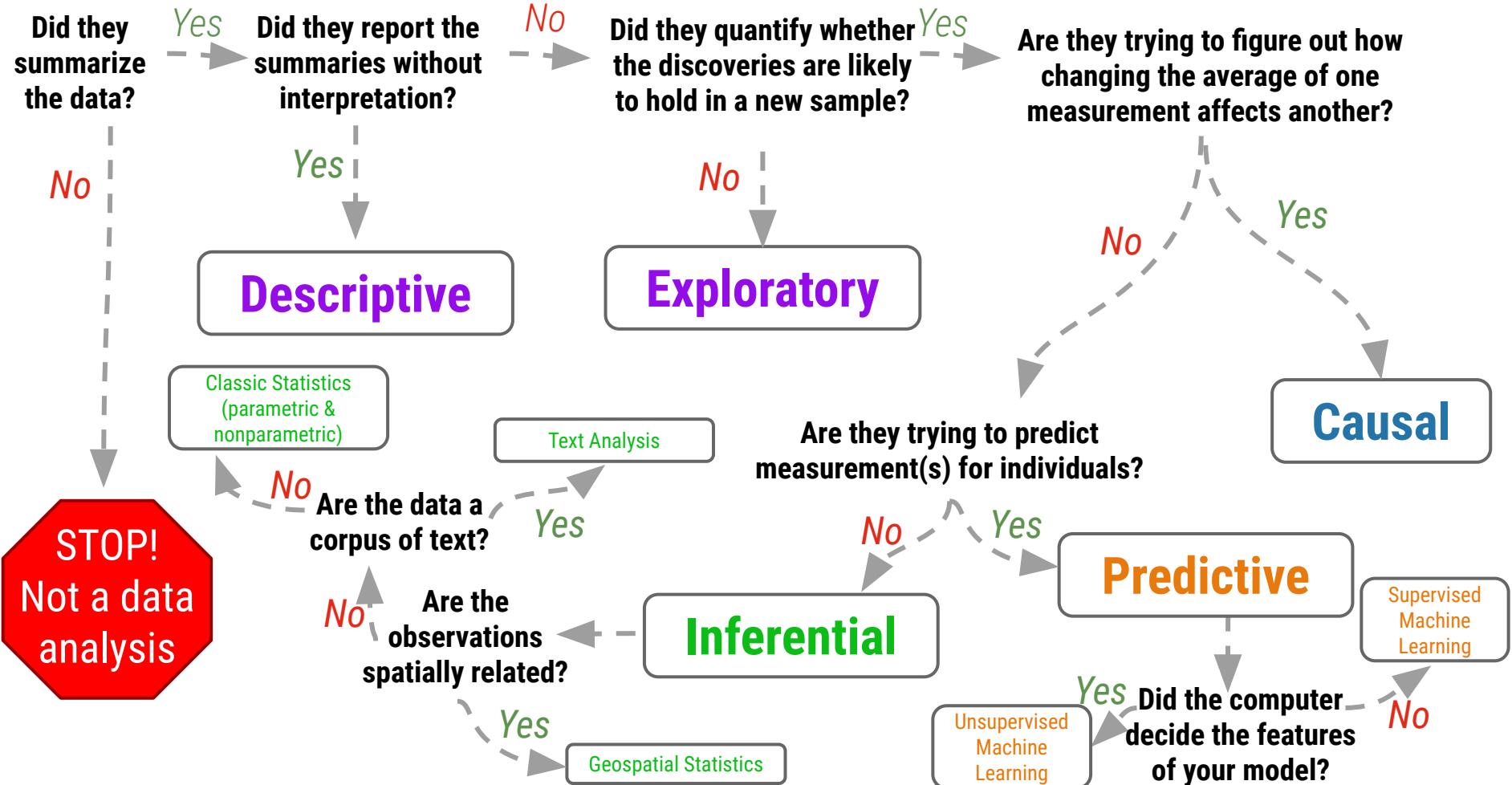
- Consider your audience
- Tell a story
  - Define question
  - Explain analysis
  - Answer question
- Spend time on visualizations here
  - EDA can be less-than-ideal visualizations
  - Final communication should be where colors, font size, etc. of visualization are considered

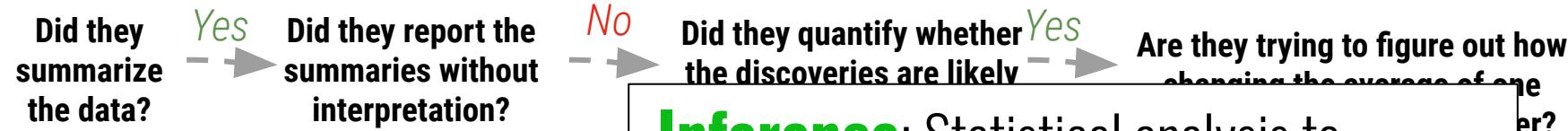
# Inferential Analysis

Shannon E. Ellis, Ph.D  
UC San Diego

• • •

Department of Cognitive Science  
[sellis@ucsd.edu](mailto:sellis@ucsd.edu)





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

No

Yes

**Descriptive**

Classic Statistics  
(parametric & nonparametric)

Text Analysis

**Causal**

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

**STOP!**  
Not a data analysis

No → Are the data a corpus of text?

Yes

No → Are the observations spatially related?

Yes

**Inferential**

Geospatial Statistics

**Predictive**

Supervised Machine Learning

Unsupervised Machine Learning

Yes → Did the computer decide the features of your model?

No

- **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



Sesame Street  
viewership

??

Test scores

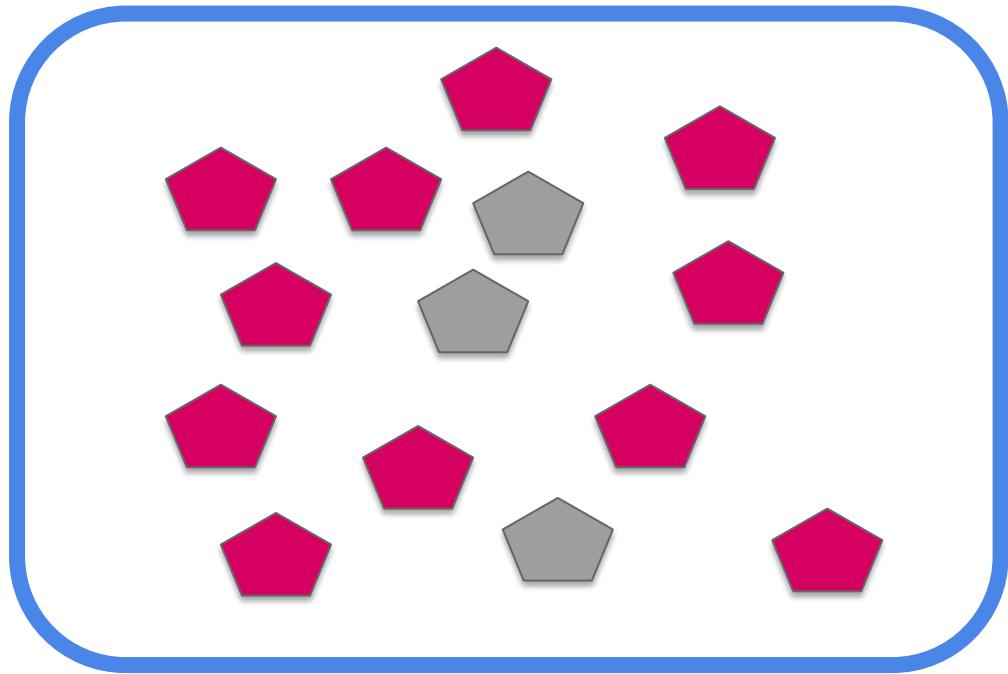
# Establishing & Stating Your Null and Alternative Hypotheses Helps Guide Your Analysis

## Null Hypothesis:

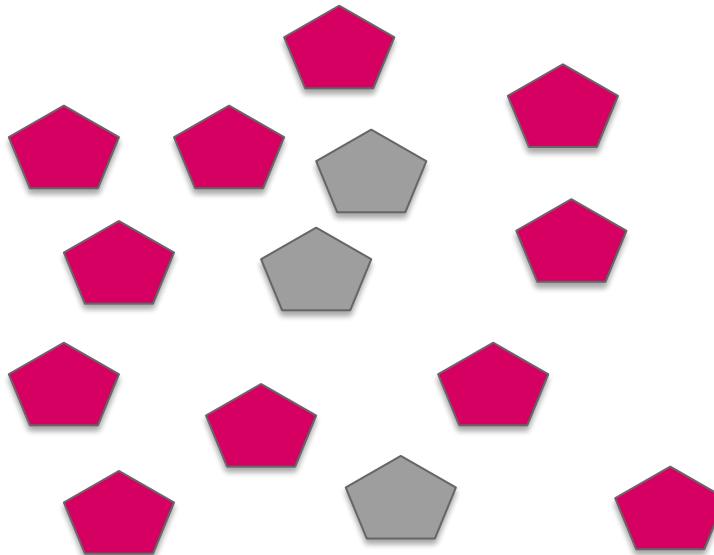
$H_0$ : Sesame Street has *no effect* on kids brain development

## Alternative Hypothesis:

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



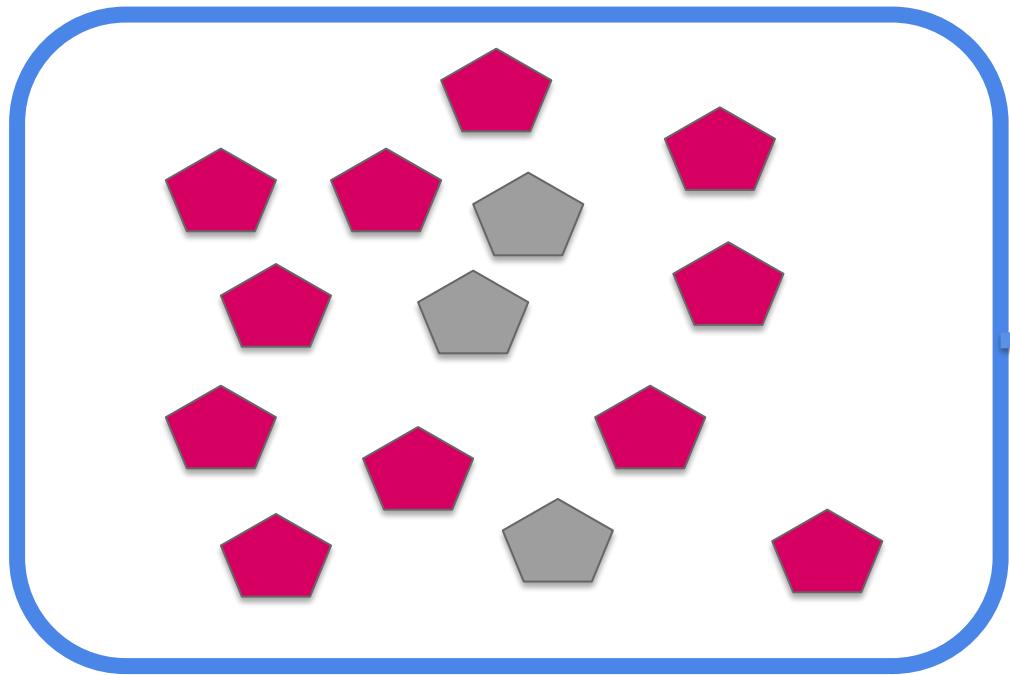
Population



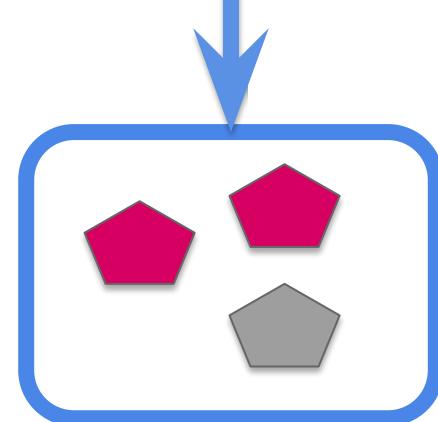
# Population



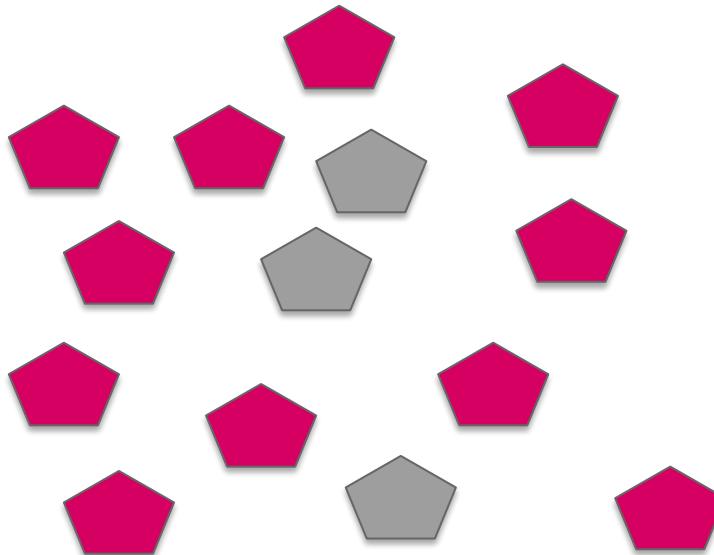
In our Sesame street example, the population would be all children



Population



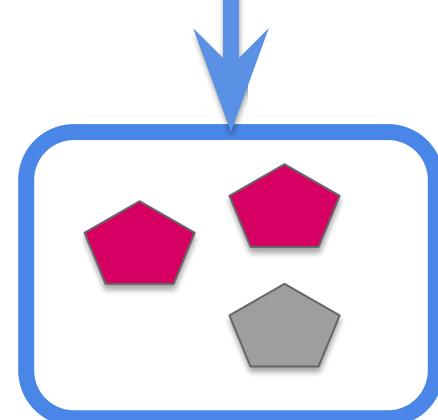
Sample



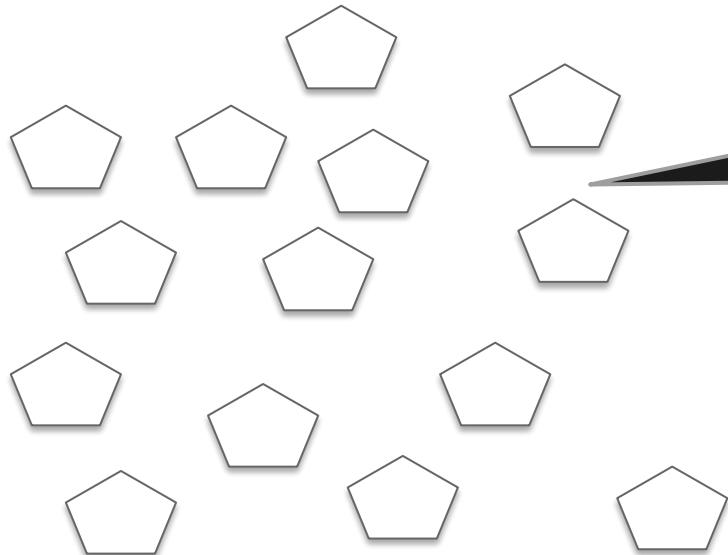
## Population



In our Sesame street example,  
the sample would be the  
children included in the study

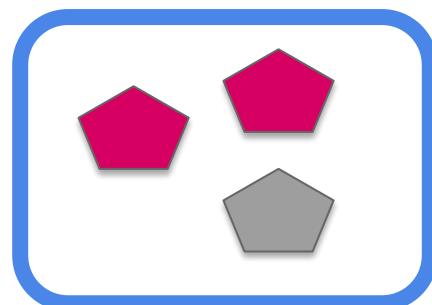


## Sample



Population

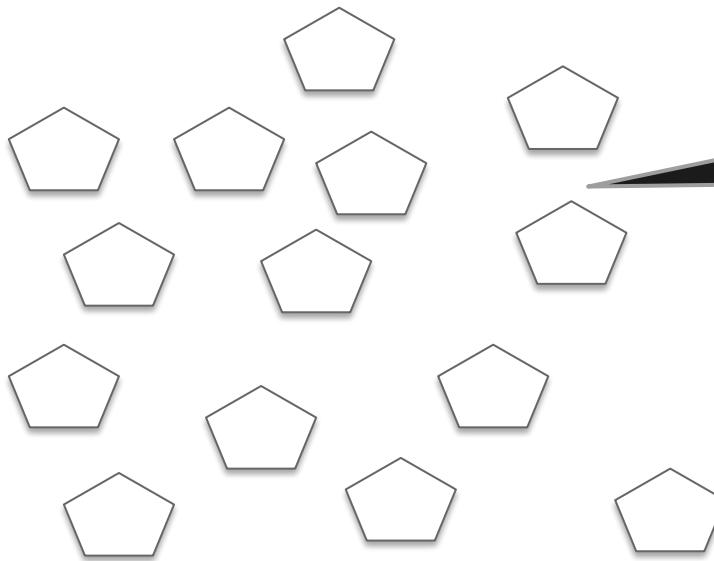
-↖(ツ)↗-



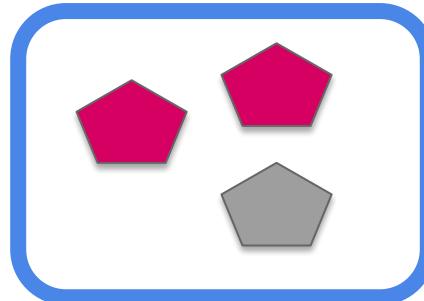
Sample



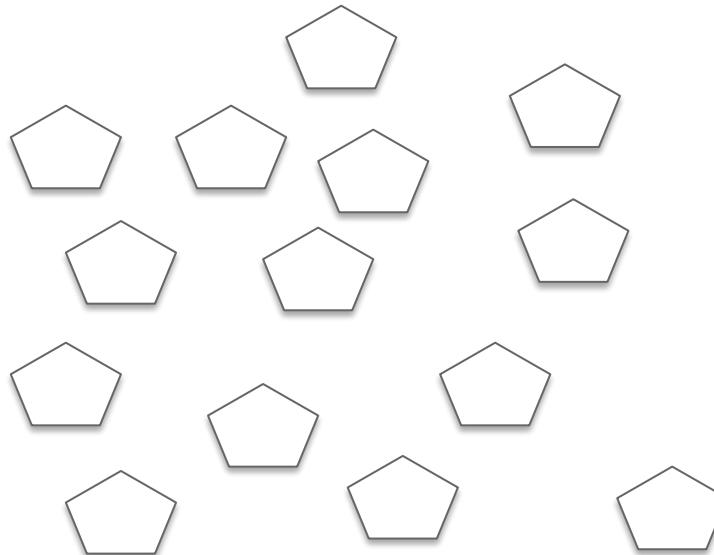
We don't know how much Sesame street was watched by or the tests scores of all kids



Population



Sample

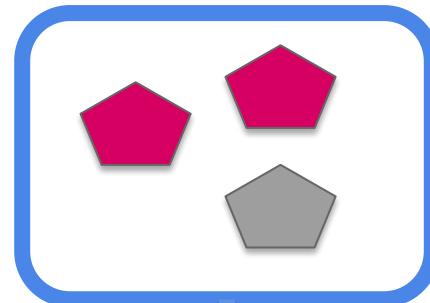


Population

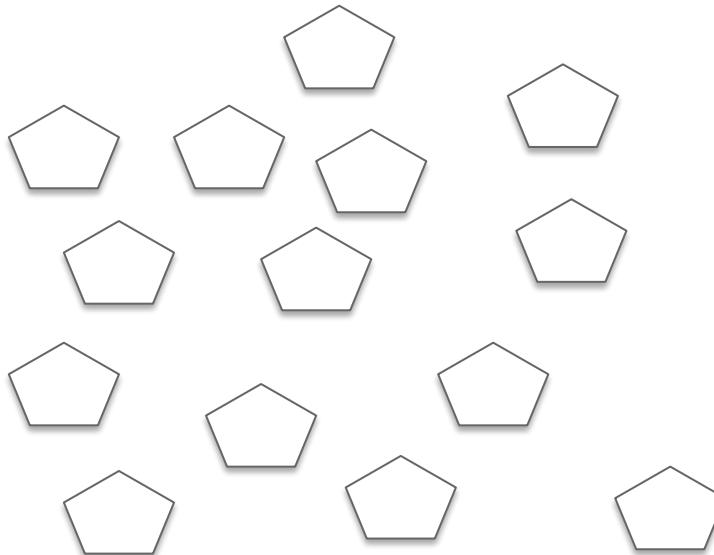


Inference!

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



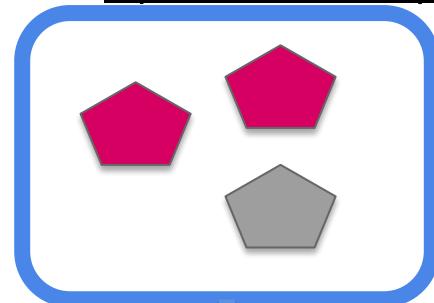
Sample



Population



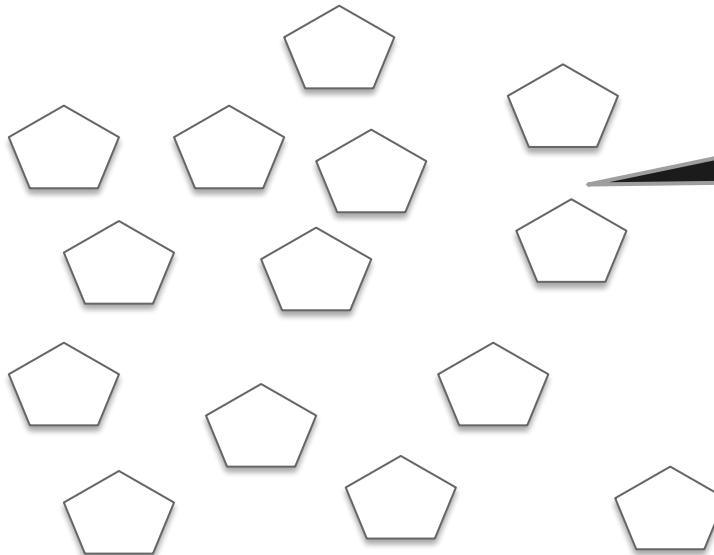
So we look at Sesame street viewing and test scores in a representative sample of kids



Inference!

Sample

Population



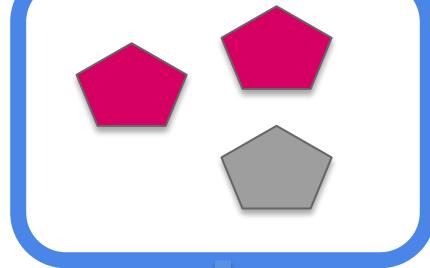
Best guess

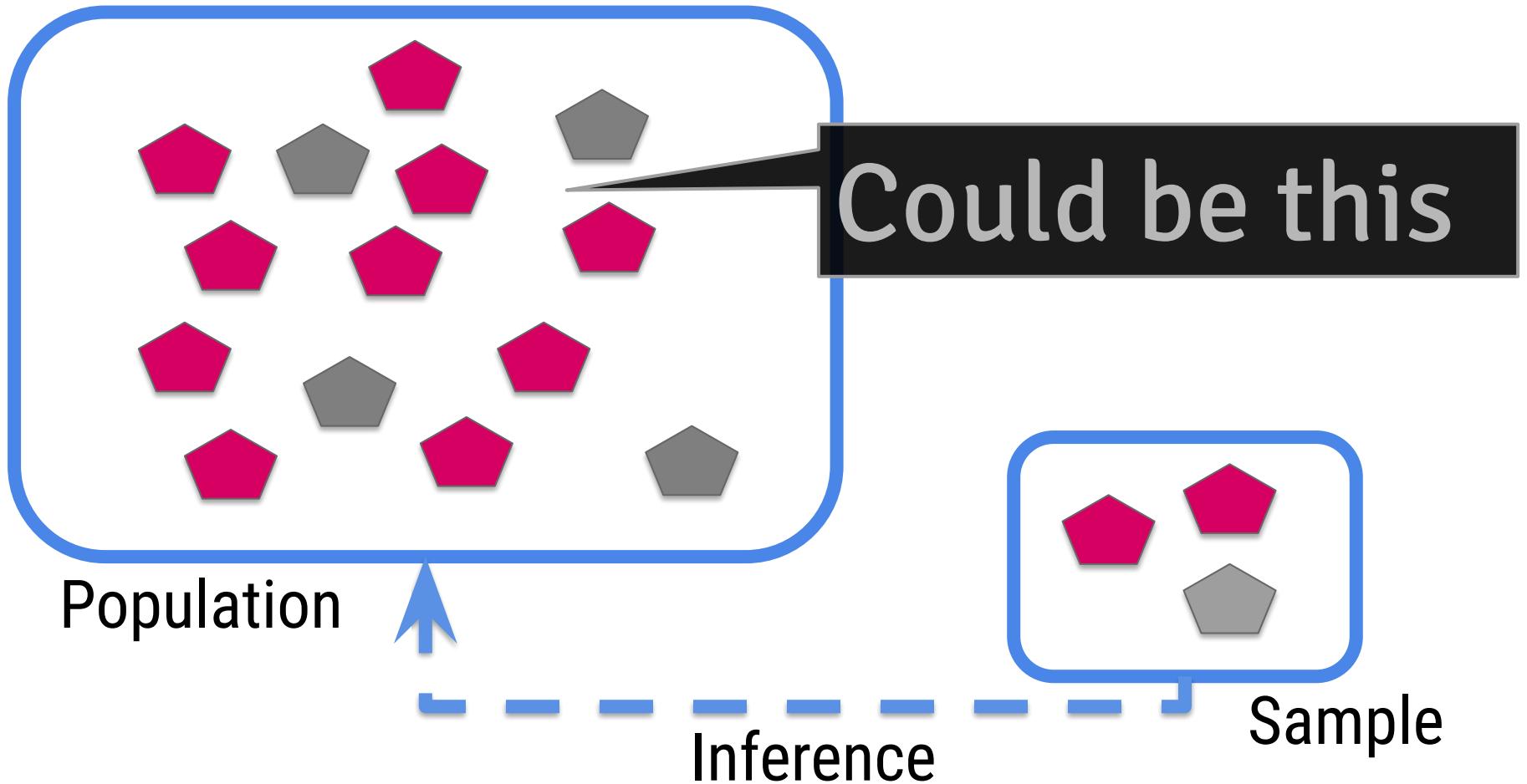


So we look at Sesame street viewing and test scores in a representative sample of kids

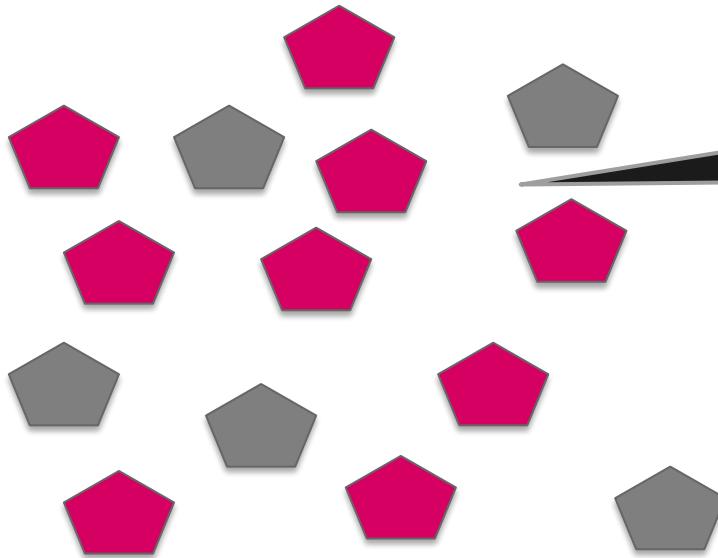
Inference!

Sample



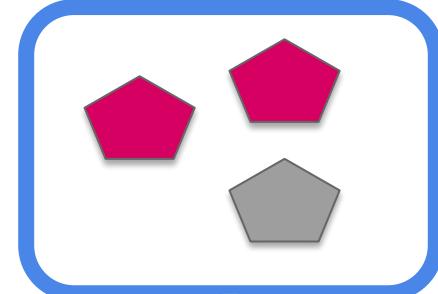


Population

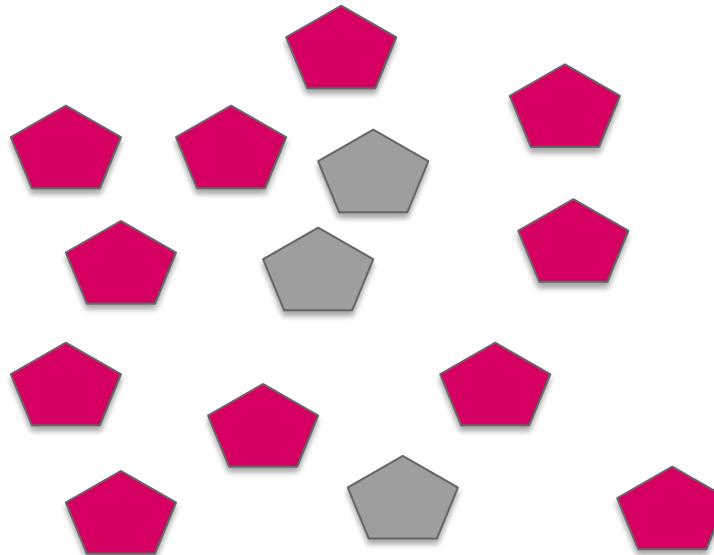


Inference

...or this

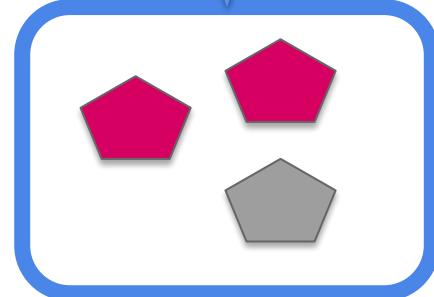


Sample

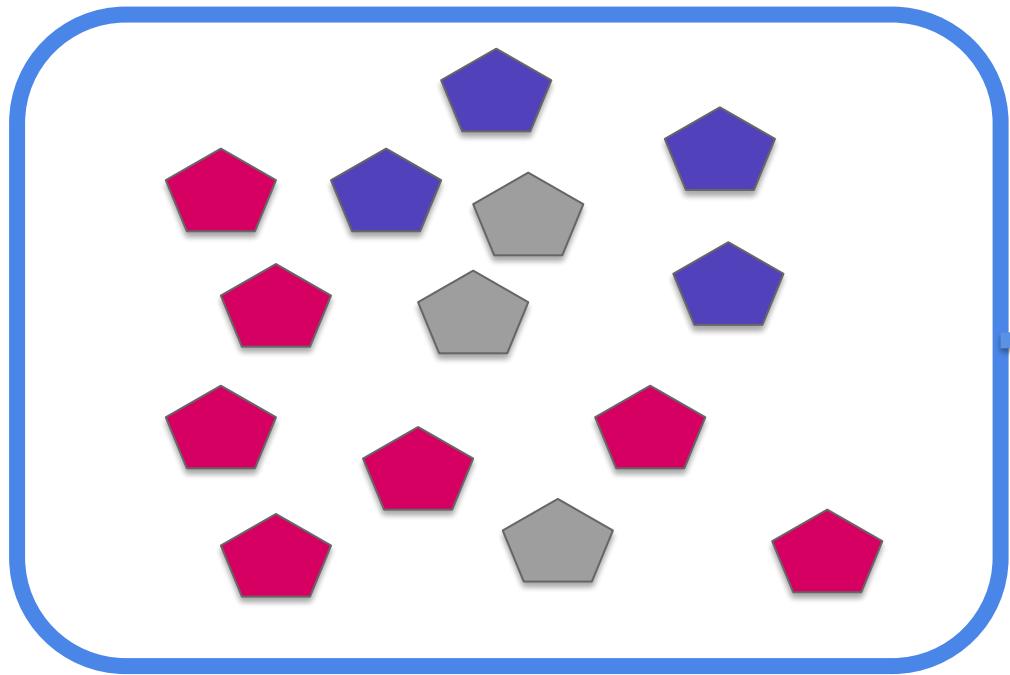


Population

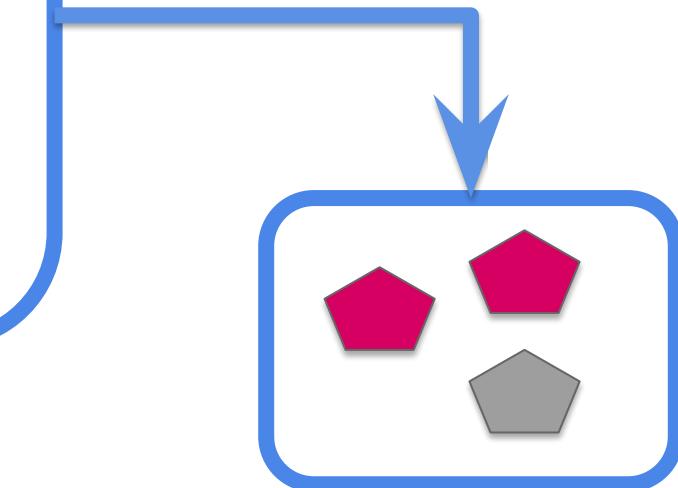
Probability



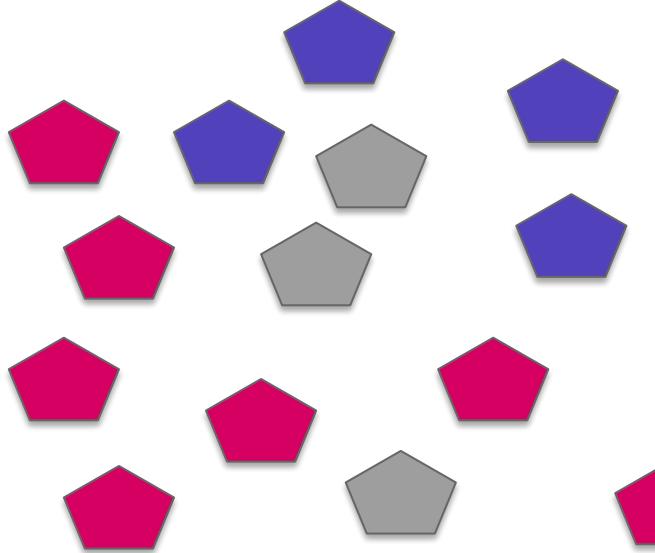
Sample



Population



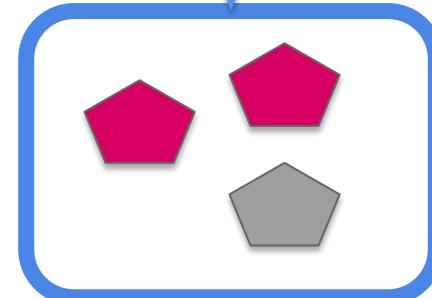
Sample



Population

~~Inference~~

If your sample is *not* representative of your population,  
you can not do inferential analysis.



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**

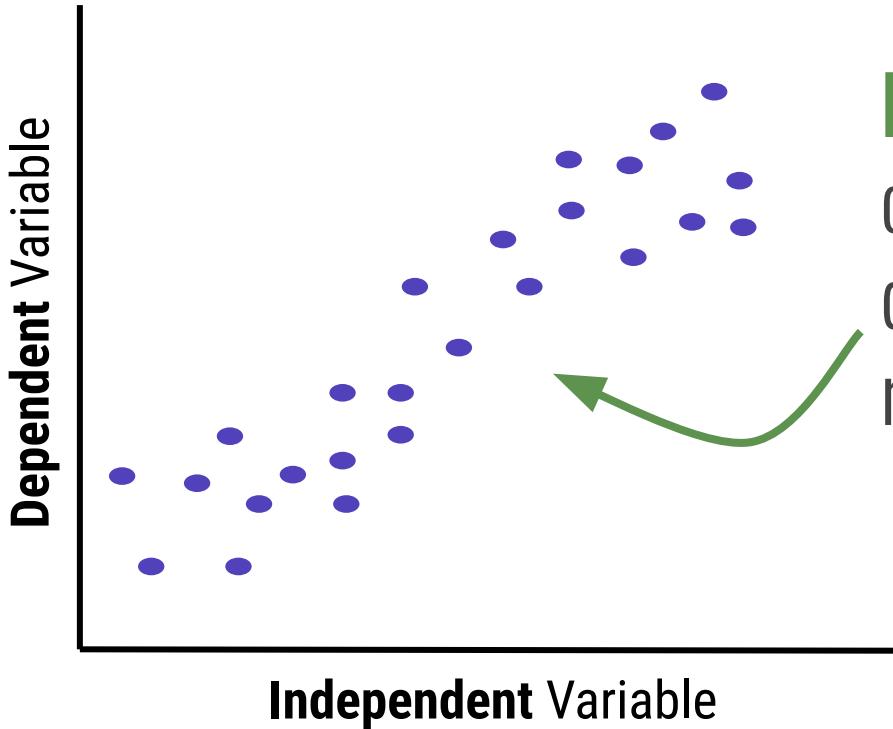
### **DOES CHANGE IN ONE VARIABLE MEAN CHANGE IN ANOTHER?**

i.e. simple regression,  
multiple regression

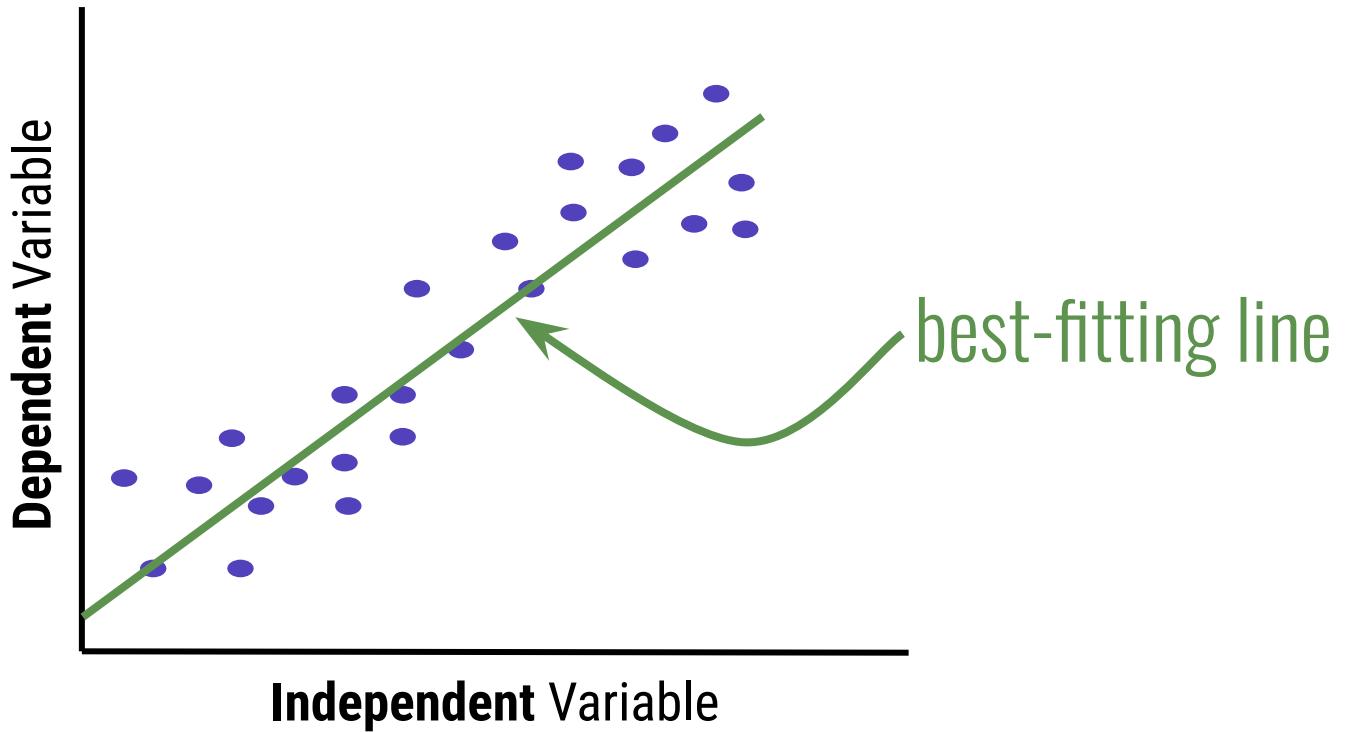
## **NON-PARAMETRIC TESTS**

### **FOR WHEN ASSUMPTIONS IN THESE OTHER 3 CATEGORIES ARE NOT MET**

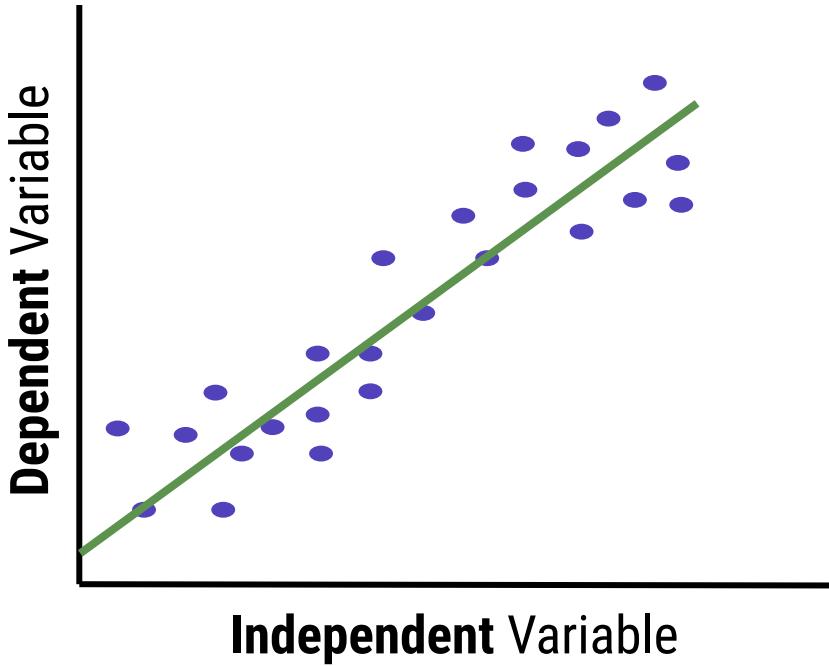
i.e. Wilcoxon rank-sum  
test, Wilcoxon sign-rank  
test, sign test



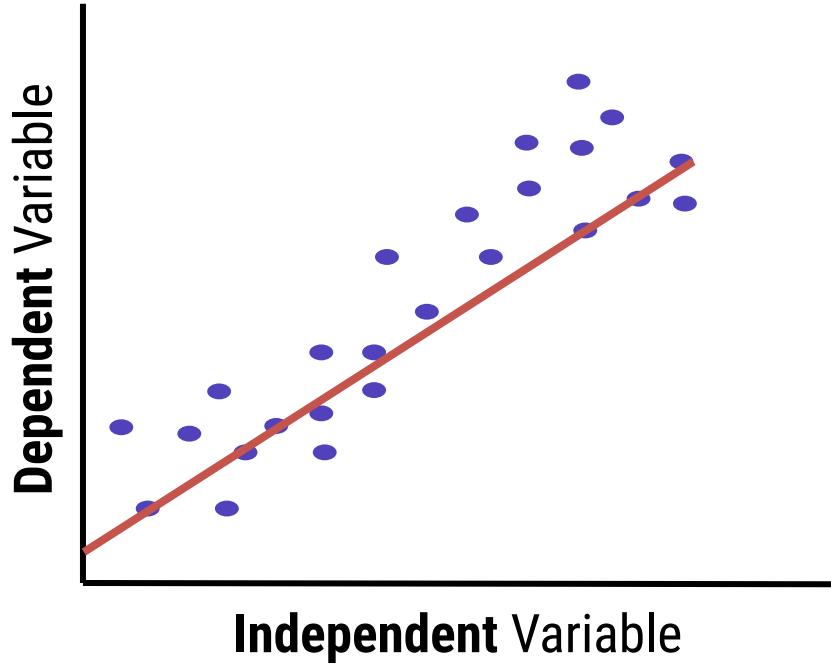
**Linear regression**  
can be used to  
describe this  
relationship

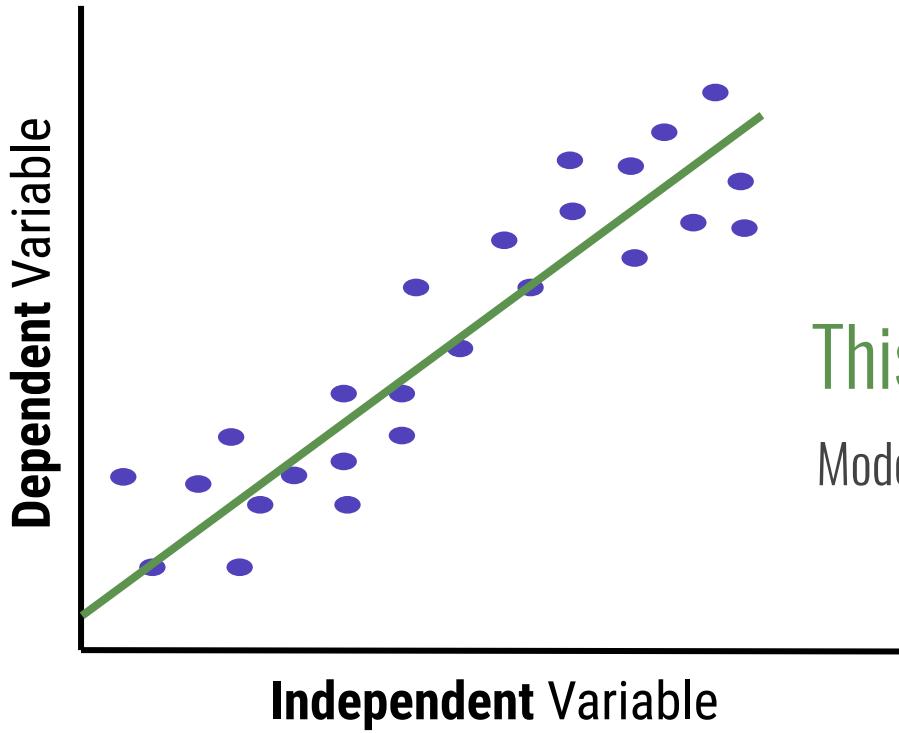


Best-fitting line

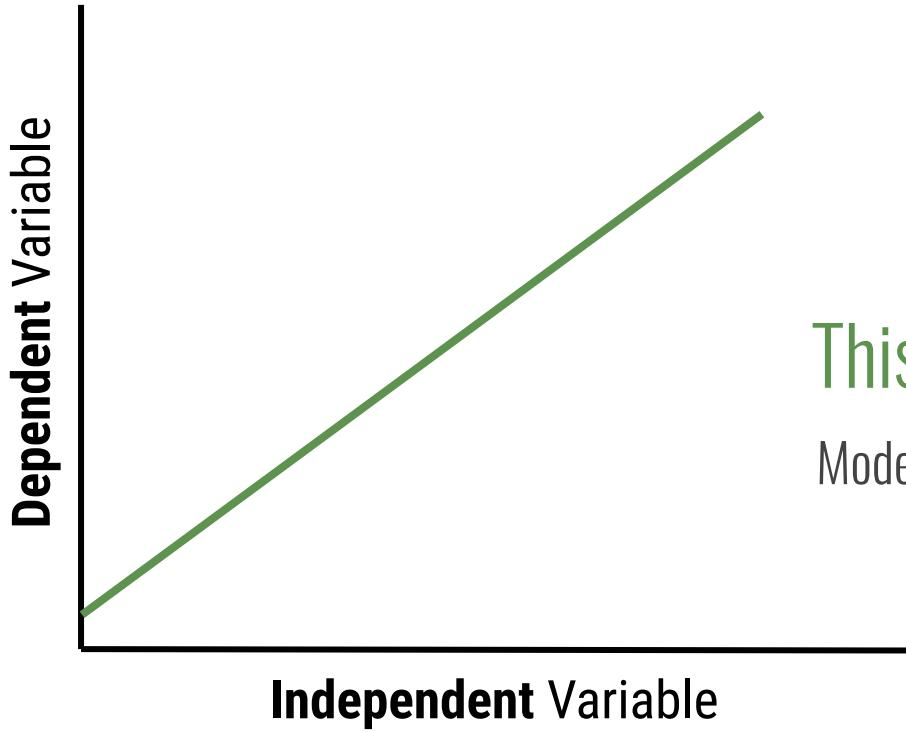


NOT a best-fitting line





This line is a **model** of the data  
Models are mathematical equations generated  
to *represent* the real life situation



This line is a **model** of the data

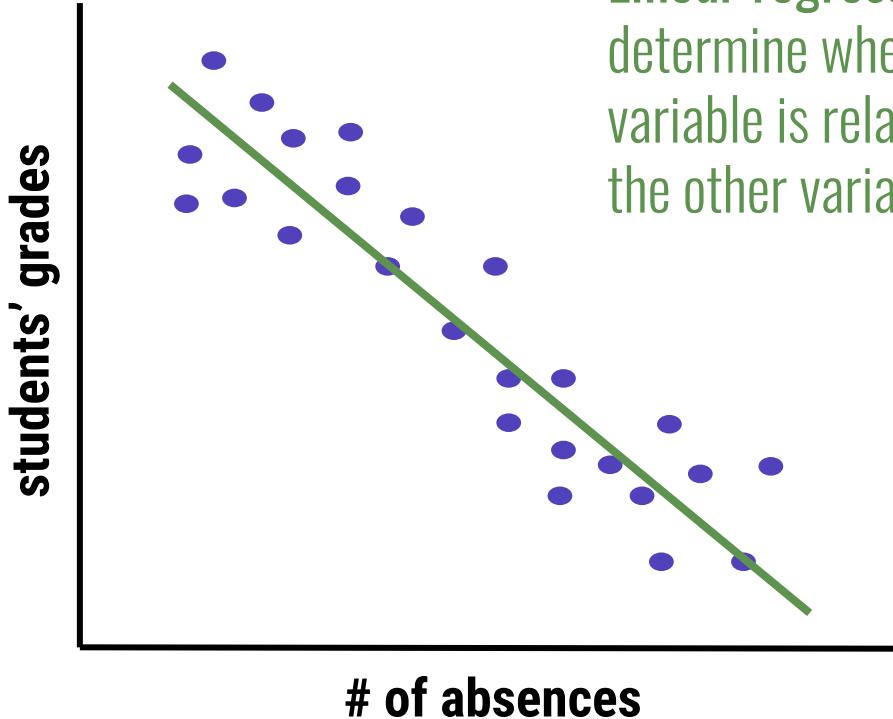
Models are mathematical equations generated  
to *represent* the real life situation

## **2.3 Parsimony**

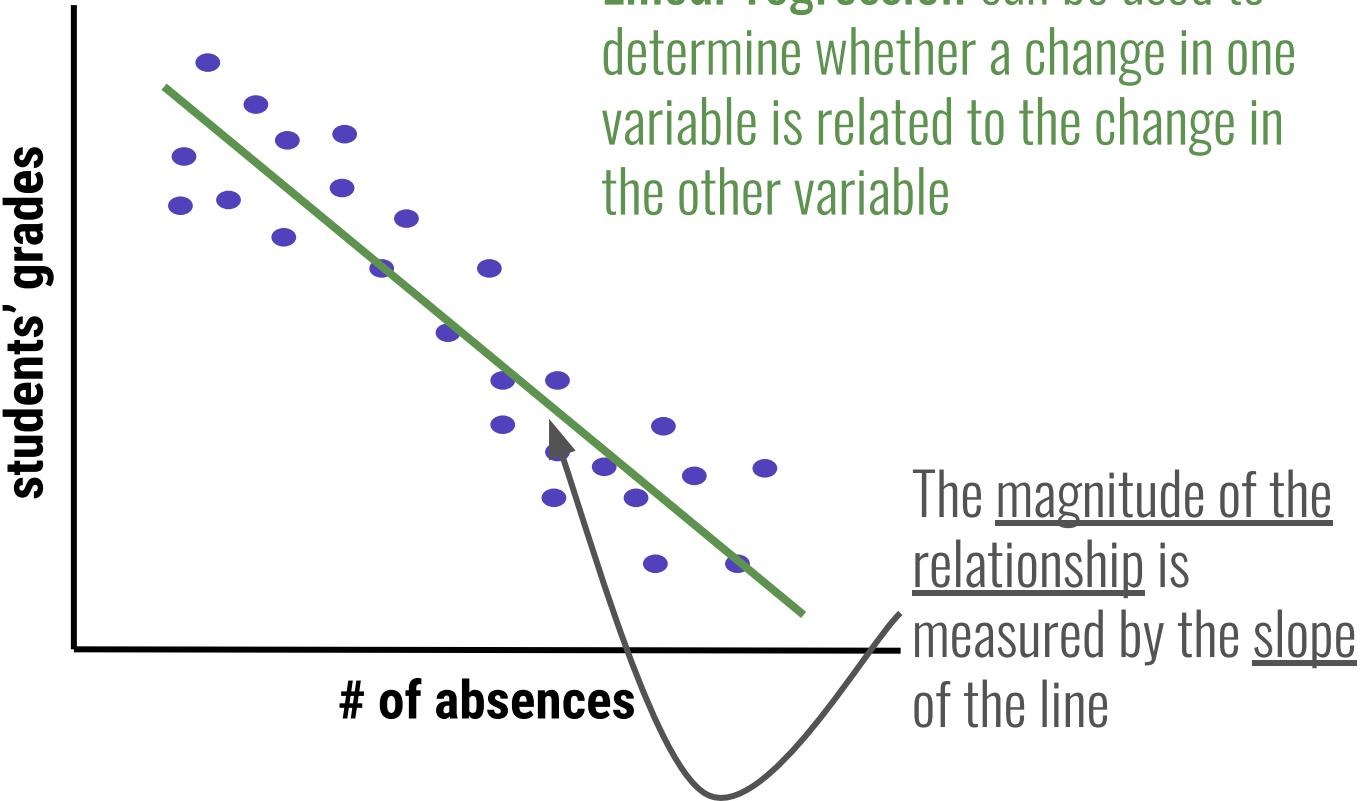
Since all models are wrong the scientist cannot obtain a “correct” one by excessive elaboration. On the contrary following William of Occam he should seek an economical description of natural phenomena. Just as the ability to devise simple but evocative models is the signature of the great scientist so overelaboration and overparameterization is often the mark of mediocrity.

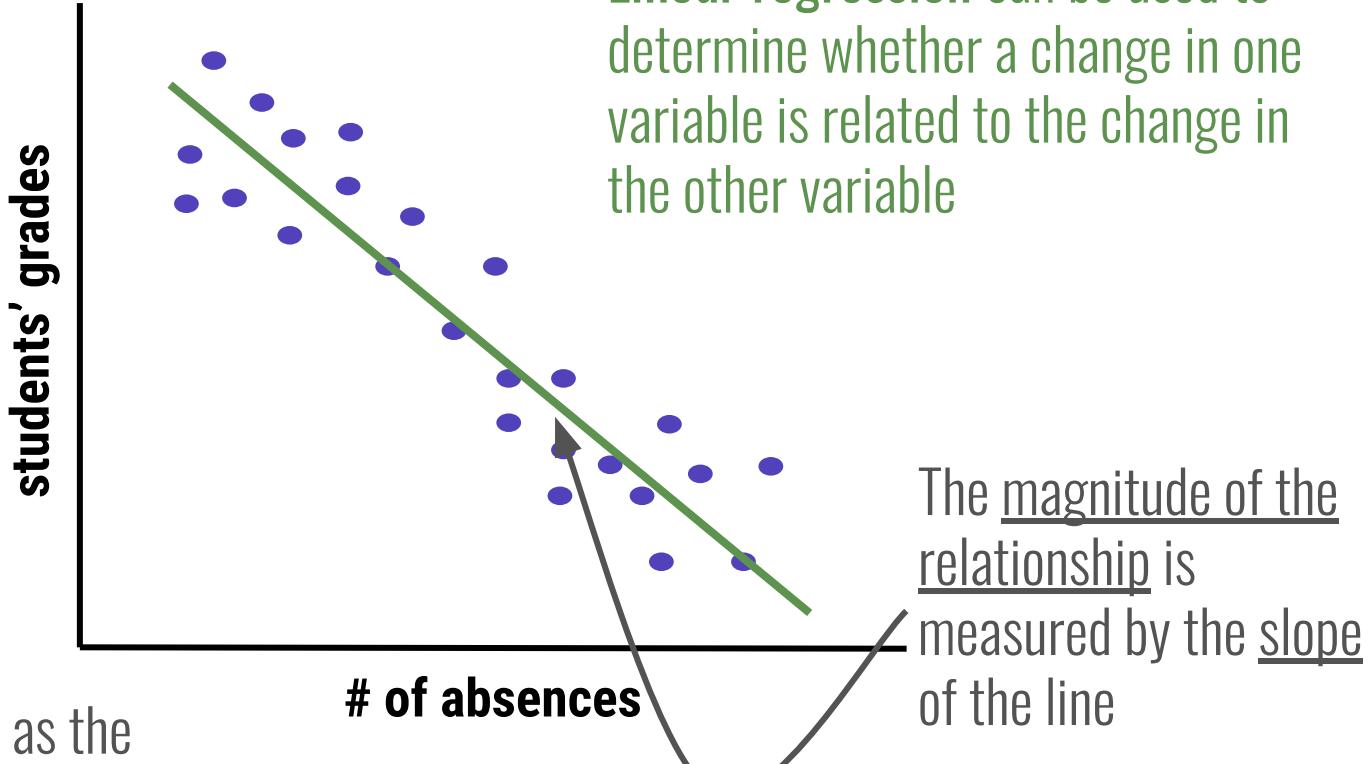
## **2.4 Worrying Selectively**

Since all models are wrong the scientist must be alert to what is importantly wrong. It is inappropriate to be concerned about mice when there are tigers abroad.

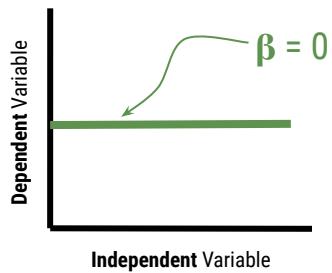


Linear regression can be used to determine whether a change in one variable is related to the change in the other variable

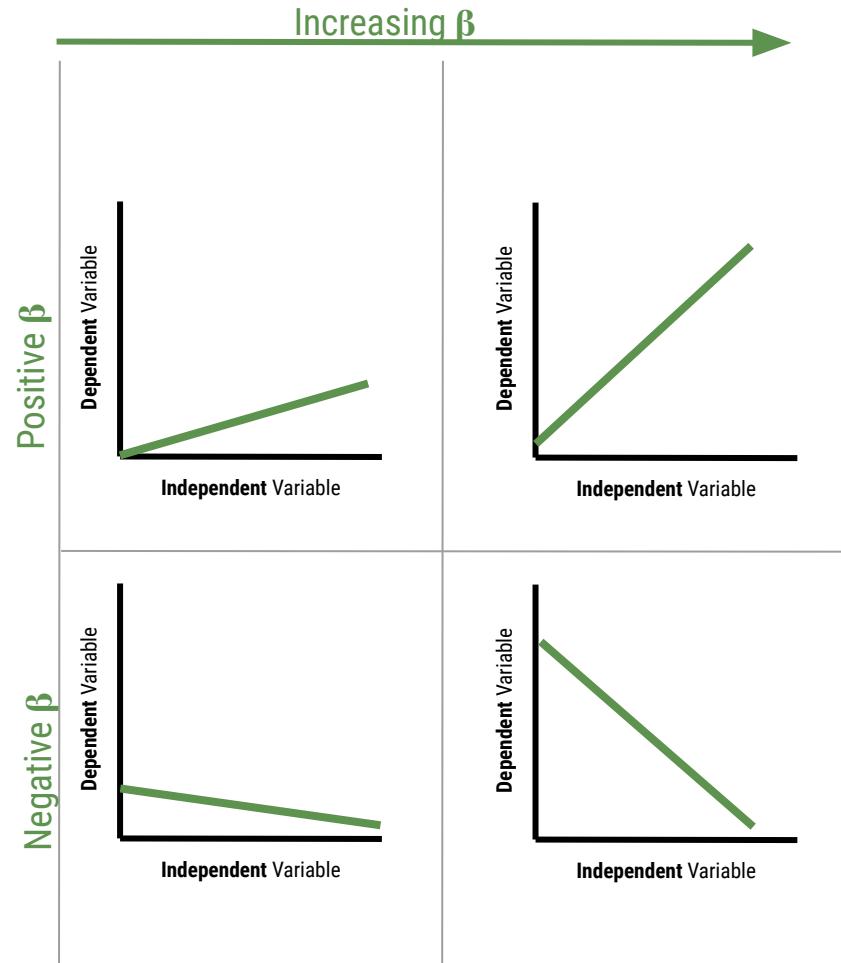
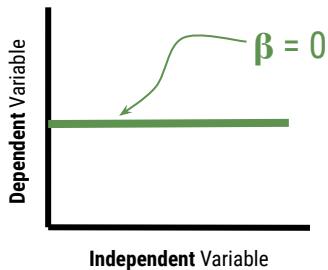




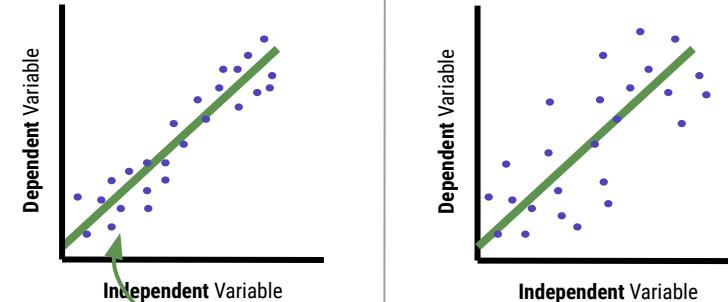
Effect size ( $\beta$ ) can  
be estimated using  
the slope of the line



Effect size ( $\beta$ ) can  
be estimated using  
the slope of the line



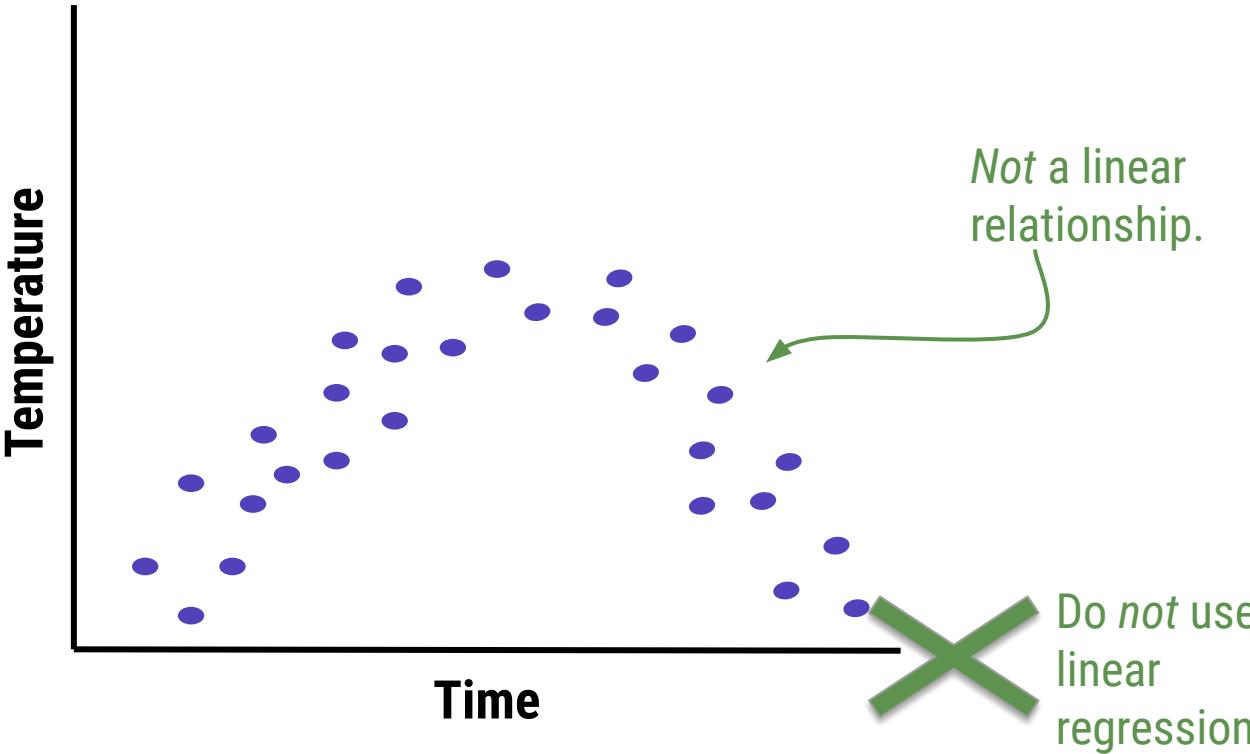
increasing standard error (SE) →



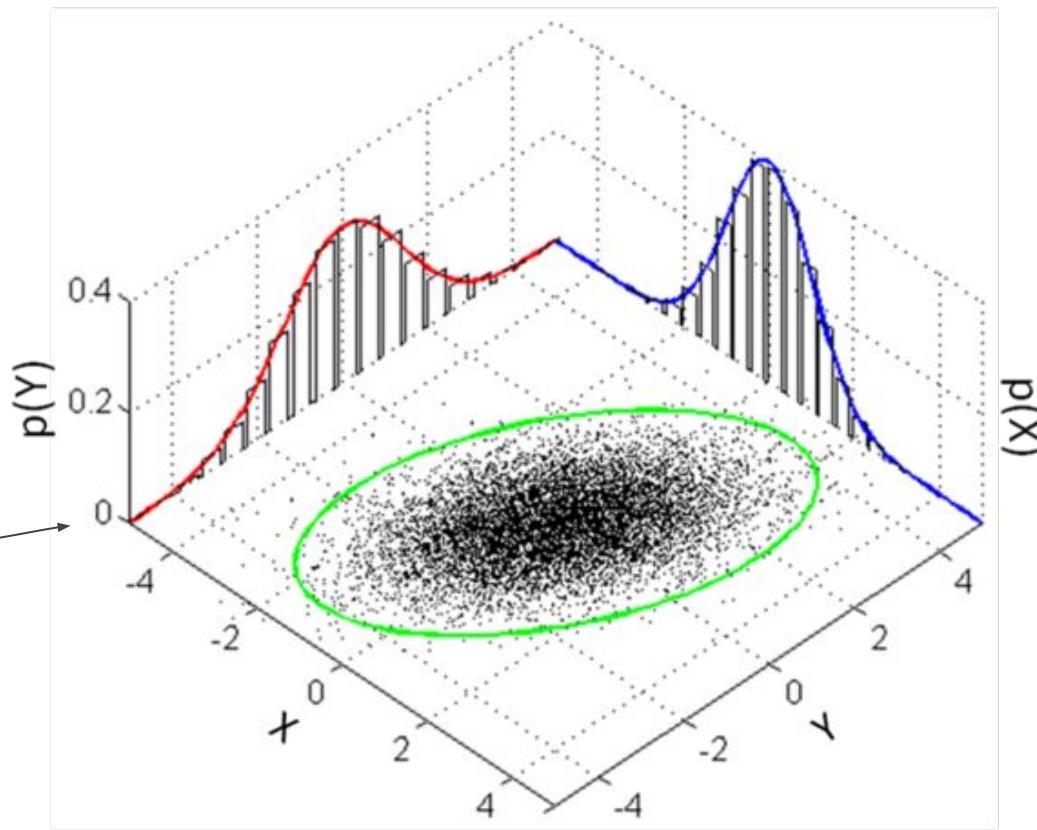
The *closer* the points  
are to the regression  
line, the *less uncertain*  
we are in our estimate

# Assumptions of linear regression

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

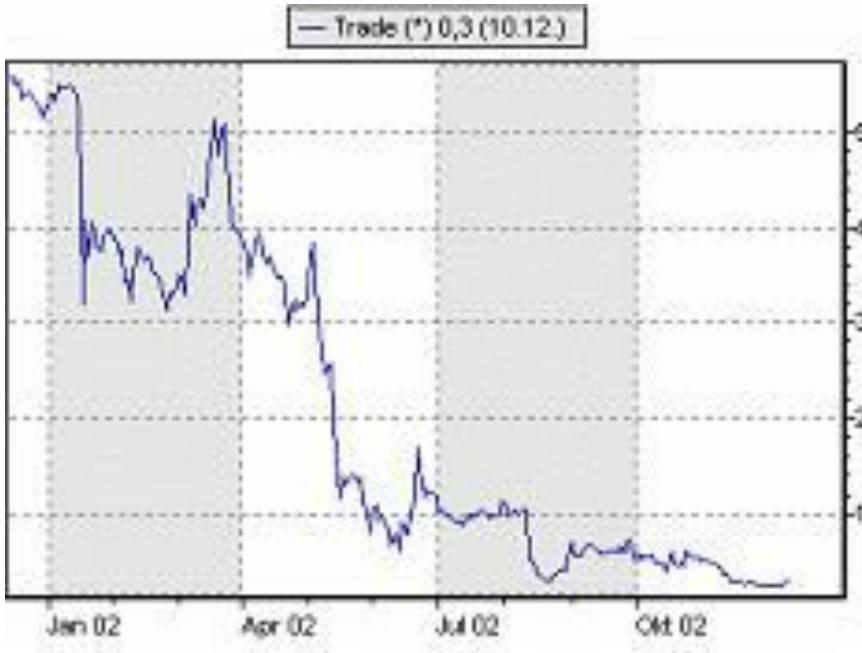


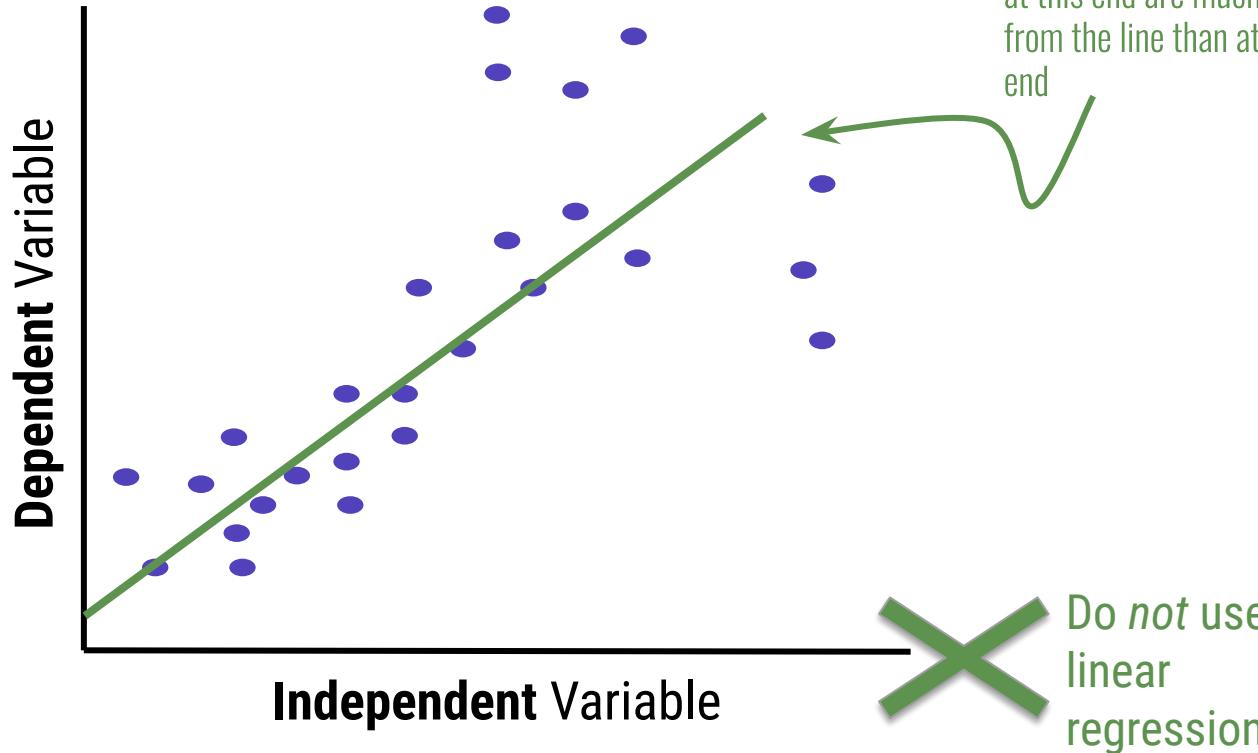
A multivariate normal  
probability distribution  
(joint normal)



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)

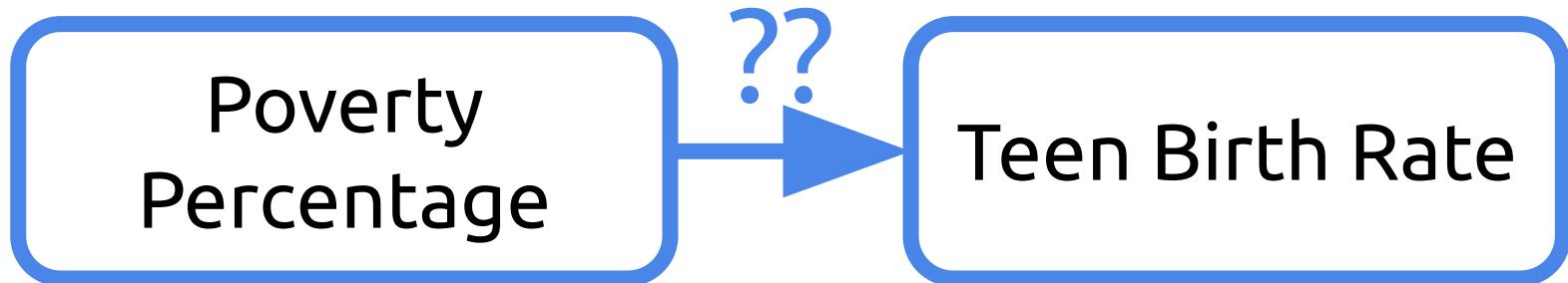




*Not homoscedastic:* points at this end are much further from the line than at the other end

***Do not use linear regression***

Does Poverty Percentage  
affect Teen Birth Rate?



Null Hypothesis:

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

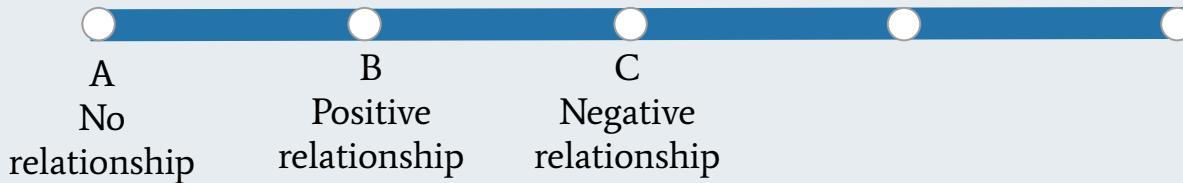
Alternative Hypothesis:

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



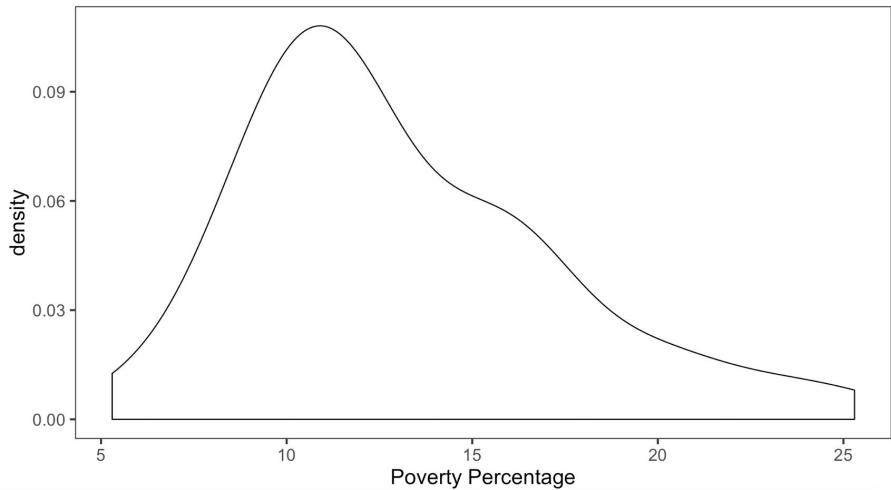
# What is the relationship between Poverty Percentage & Teen Birth Rate?

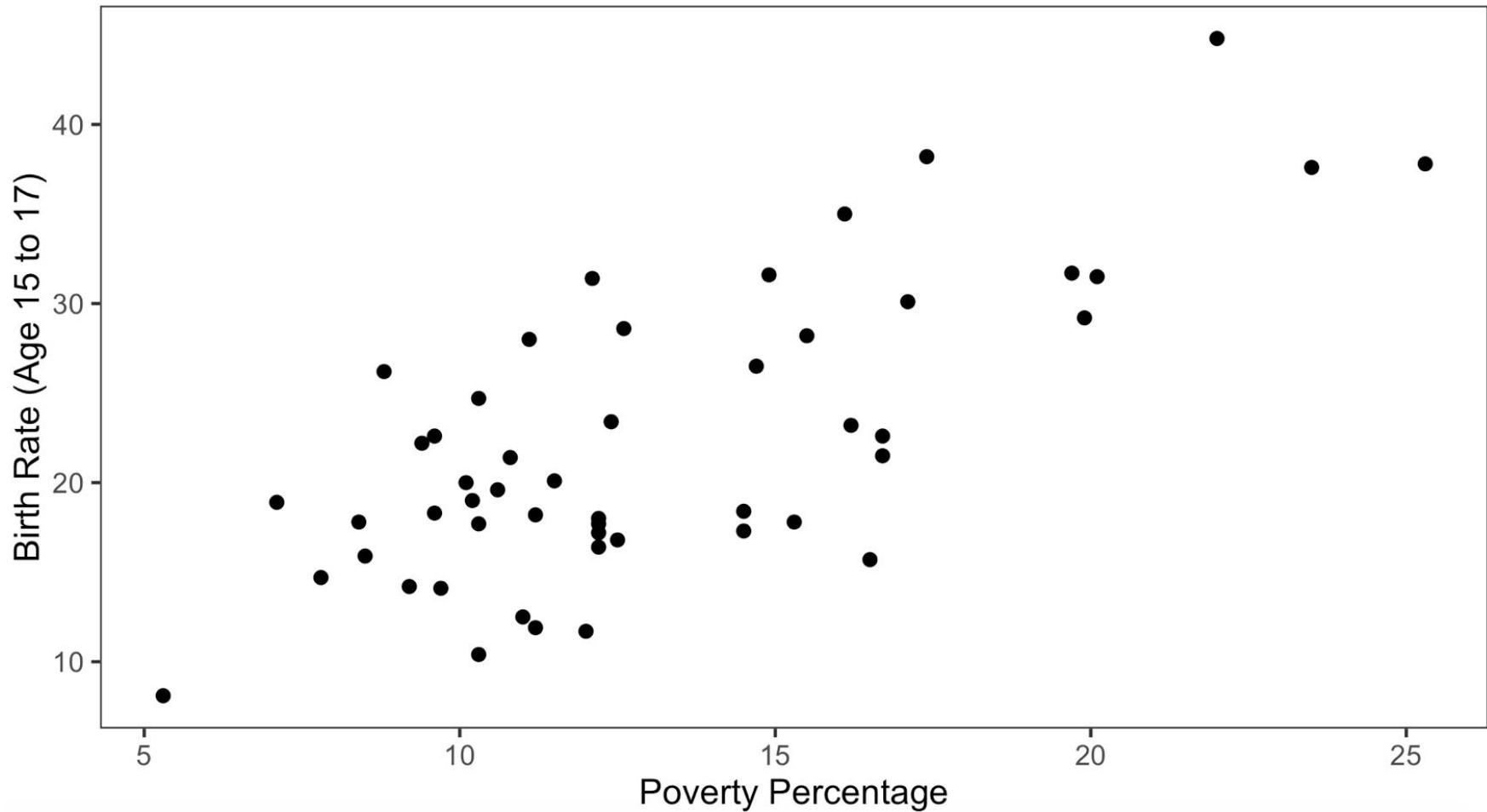
What's your hypothesis?



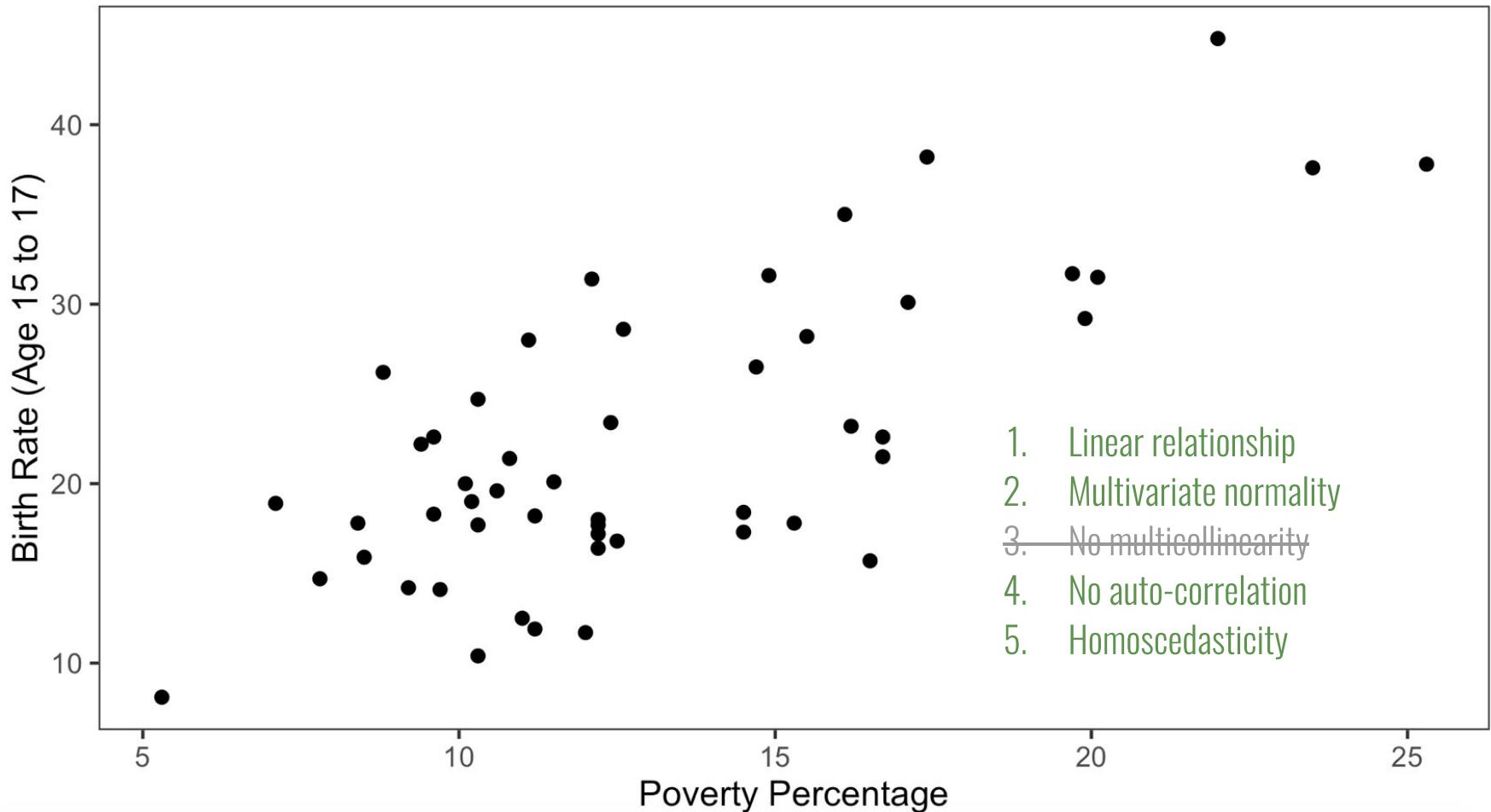
	Location	PovPct	Brth15to17	Brth18to19	ViolCrime	TeenBrth
1	Alabama	20.1	31.5	88.7	11.2	54.5
2	Alaska	7.1	18.9	73.7	9.1	39.5
3	Arizona	16.1	35.0	102.5	10.4	61.2
4	Arkansas	14.9	31.6	101.7	10.4	59.9
5	California	16.7	22.6	69.1	11.2	41.1
6	Colorado	8.8	26.2	79.1	5.8	47.0
7	Connecticut	9.7	14.1	45.1	4.6	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

# Normal(ish) distributions

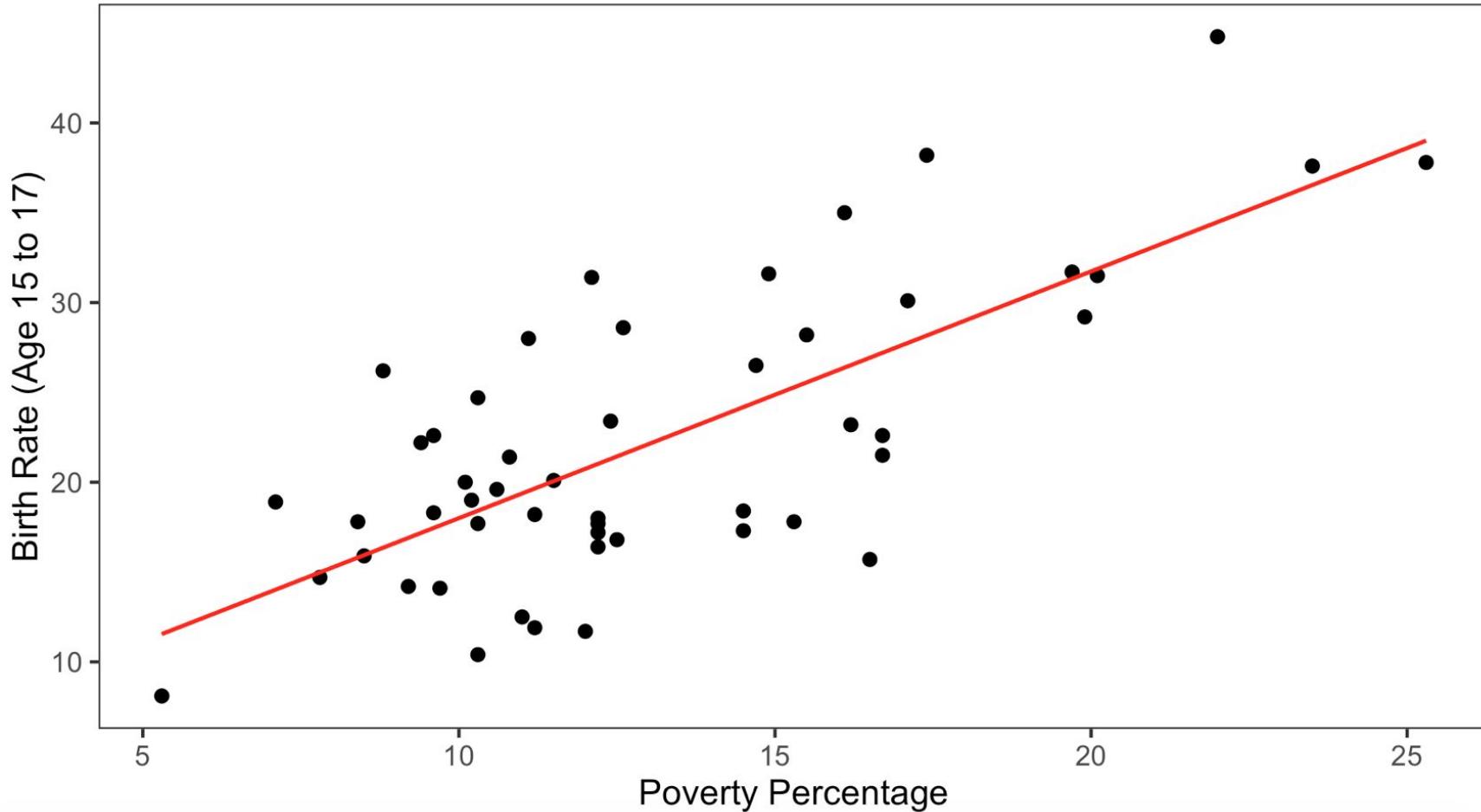


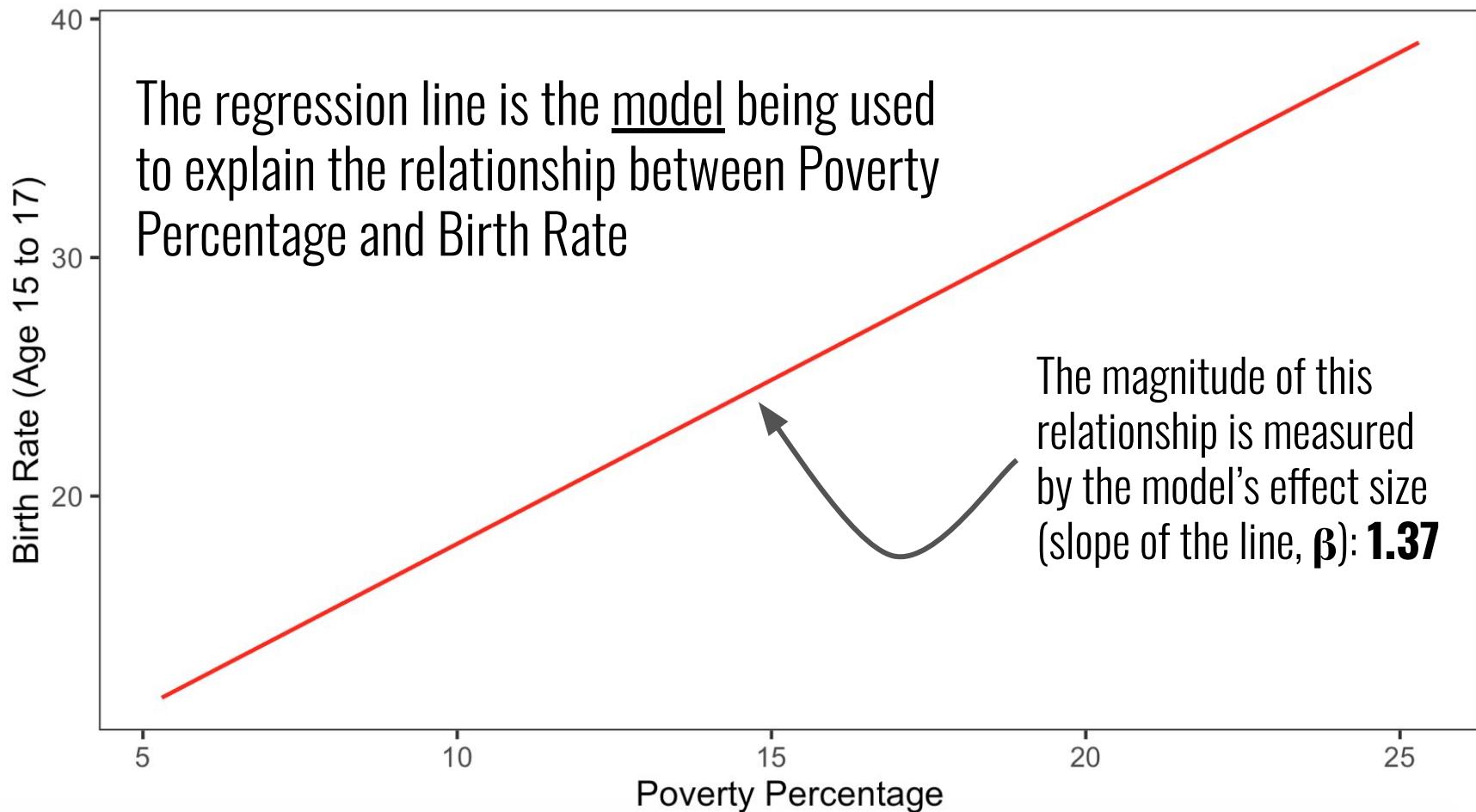


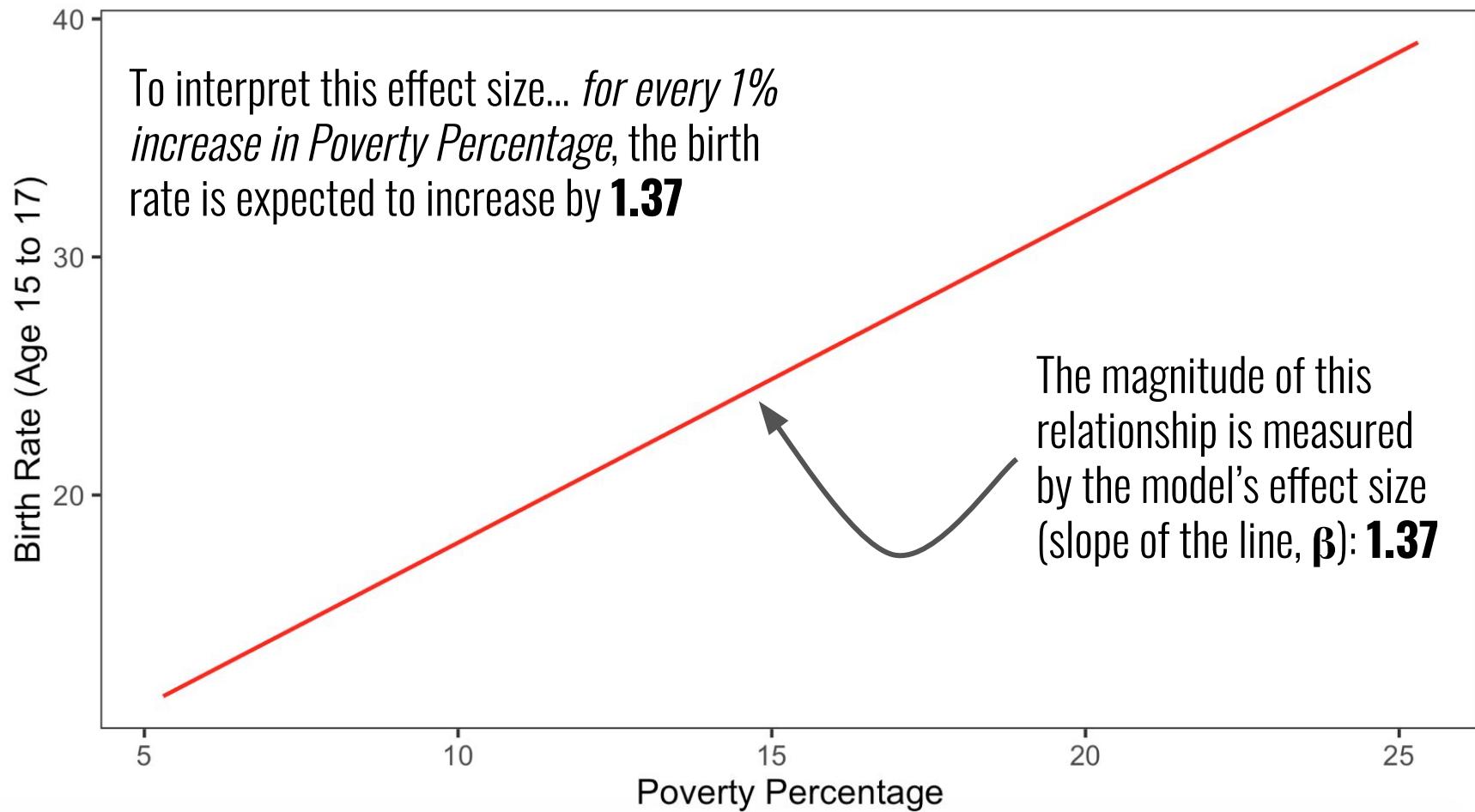
Data source: *Mind On Statistics*, 3rd edition, Utts and Heckard.

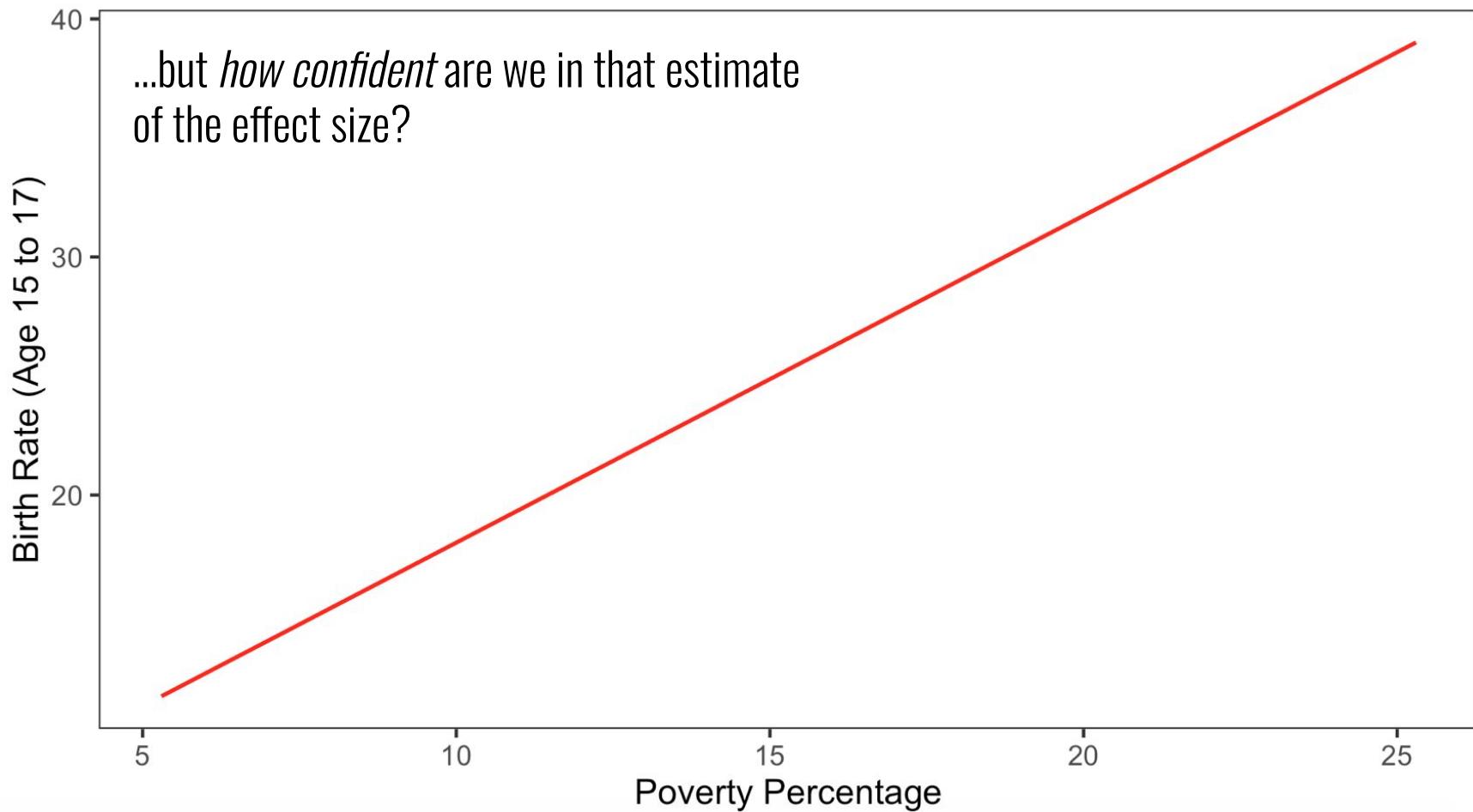


1. Linear relationship
2. Multivariate normality
3. ~~No multicollinearity~~
4. No auto-correlation
5. Homoscedasticity

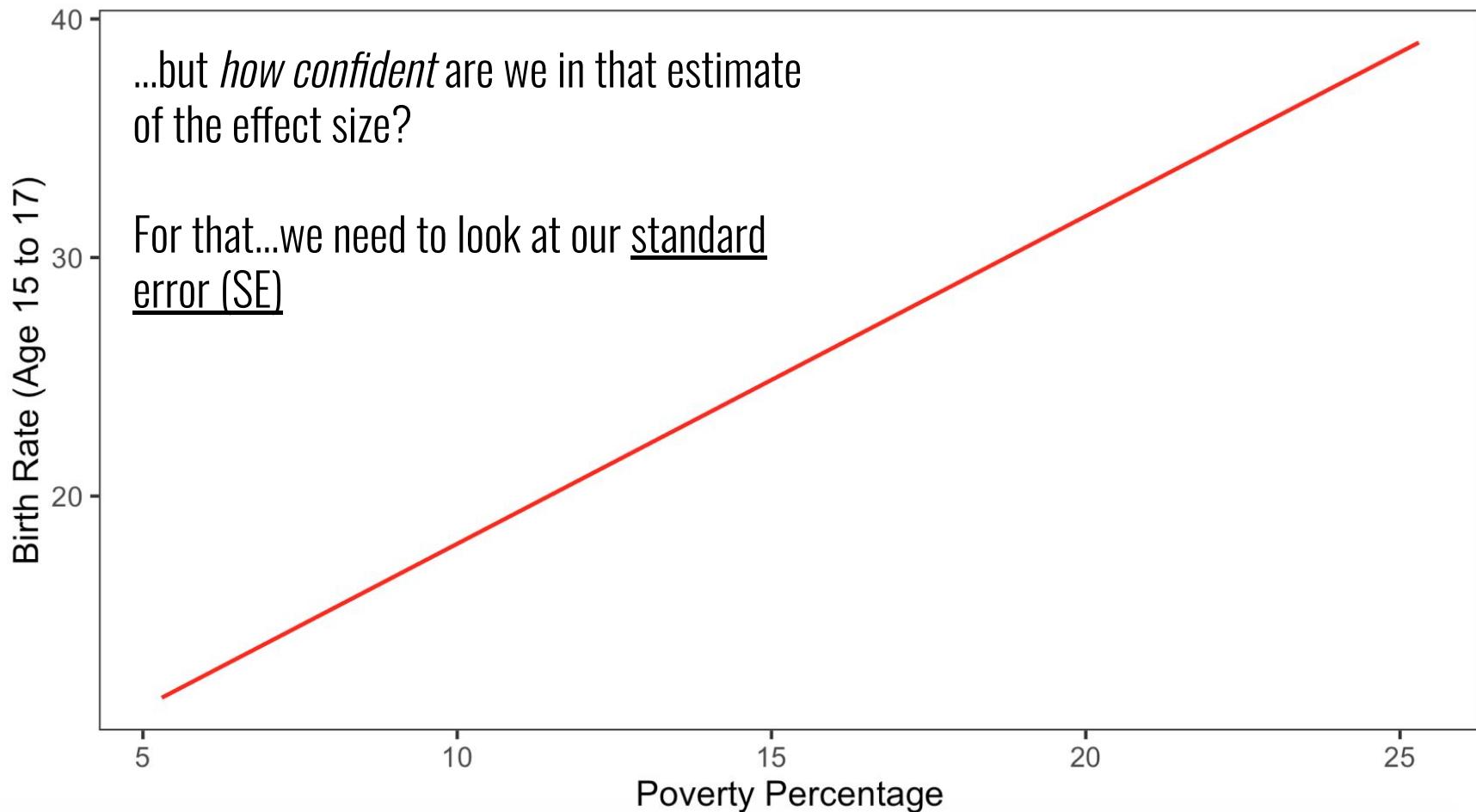






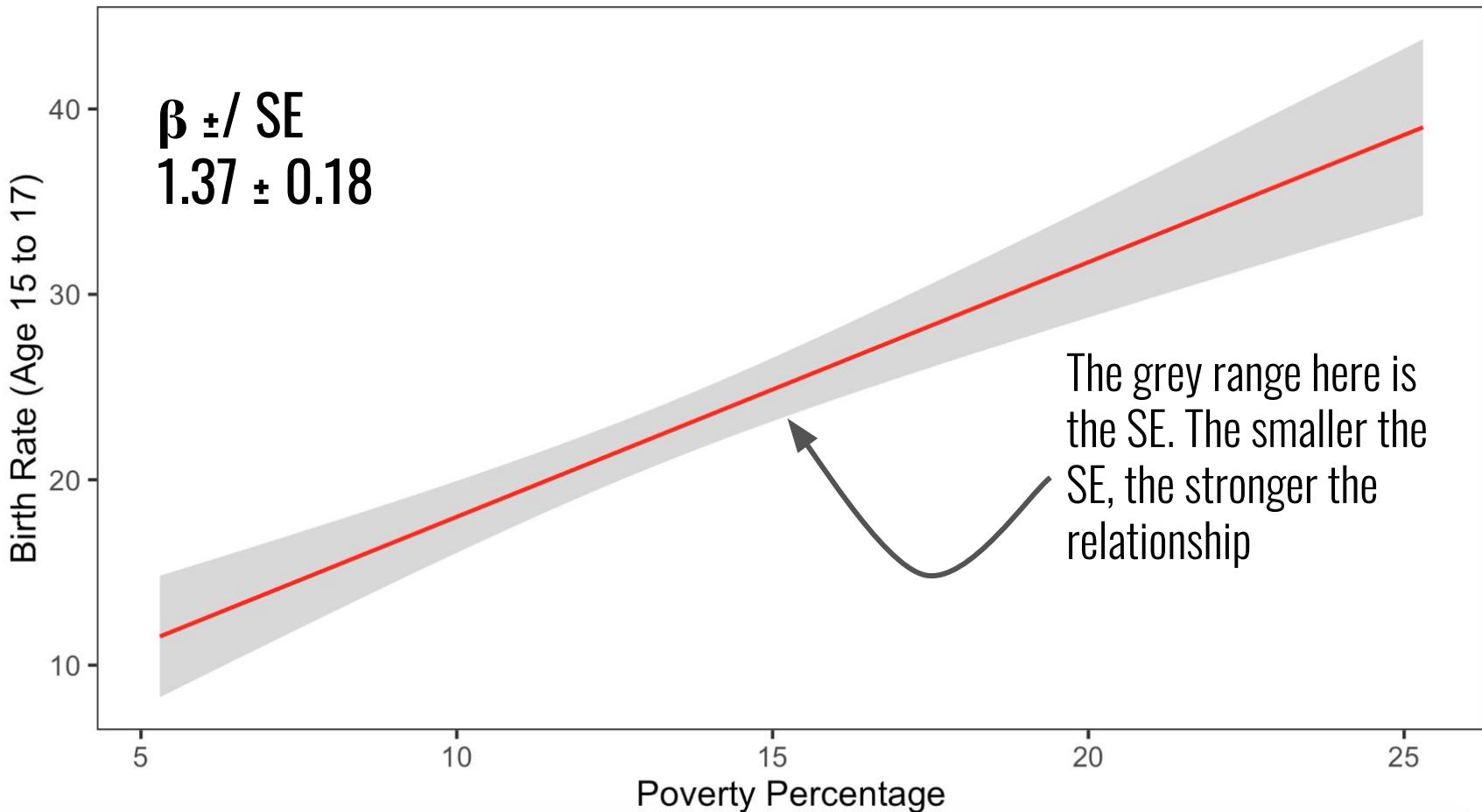


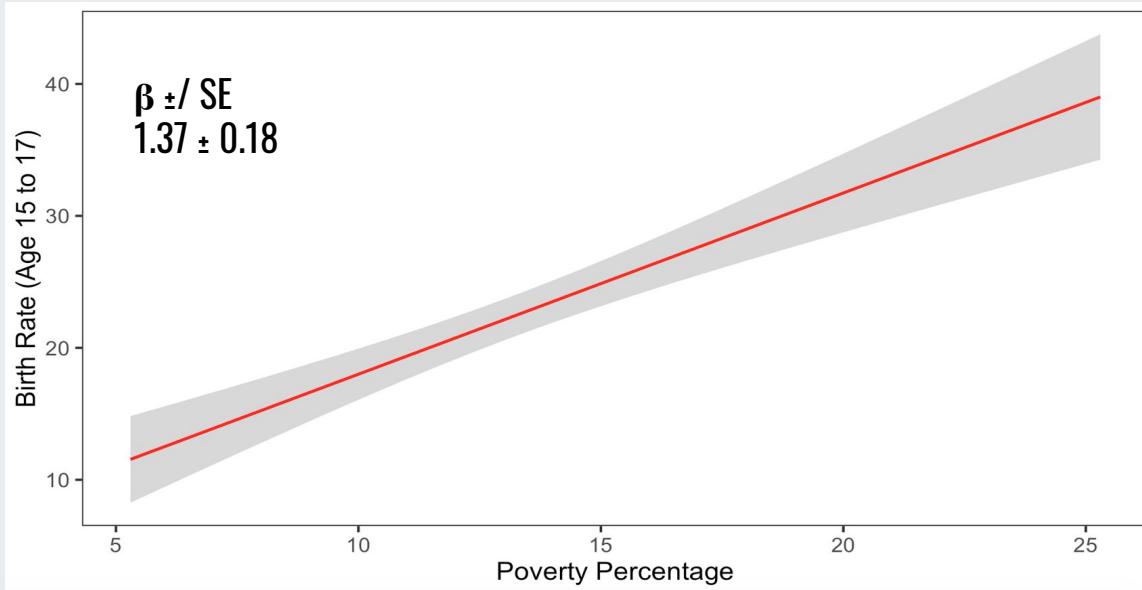
*...but how confident are we in that estimate  
of the effect size?*



...but *how confident* are we in that estimate  
of the effect size?

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





If there were a stronger effect of Poverty on Birth rate, what would  $\beta$  be?

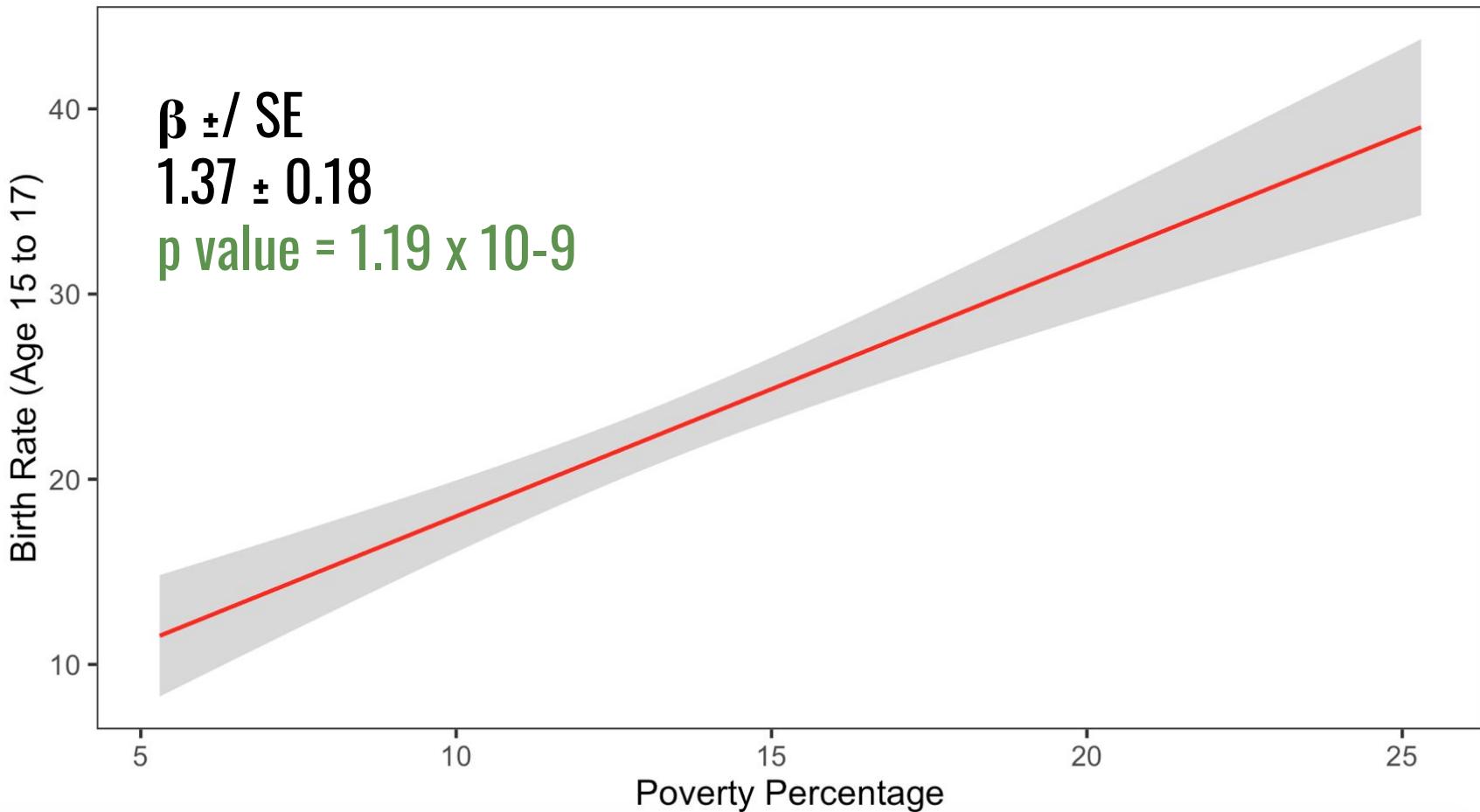


A  
 $< 1.37$

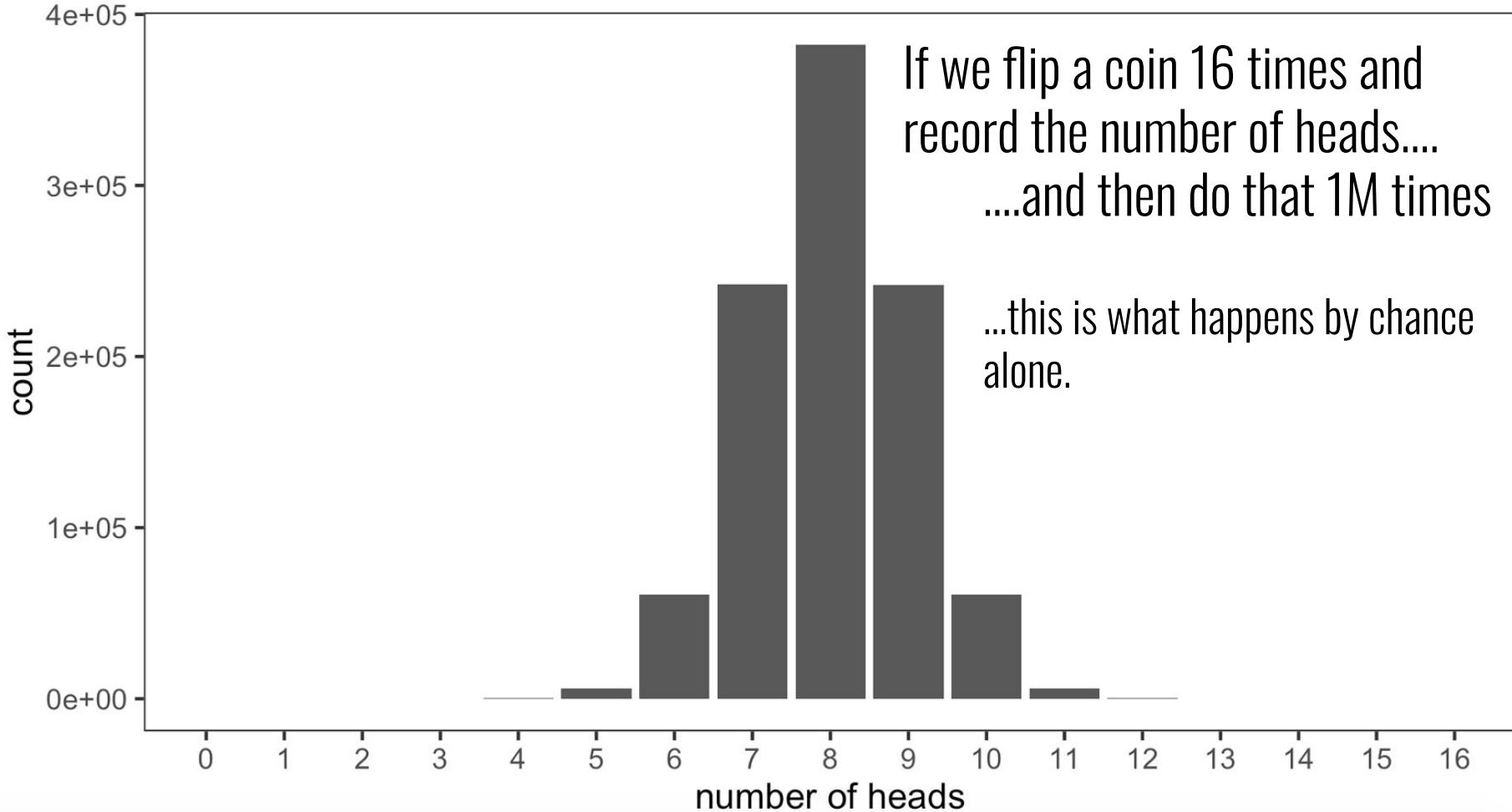


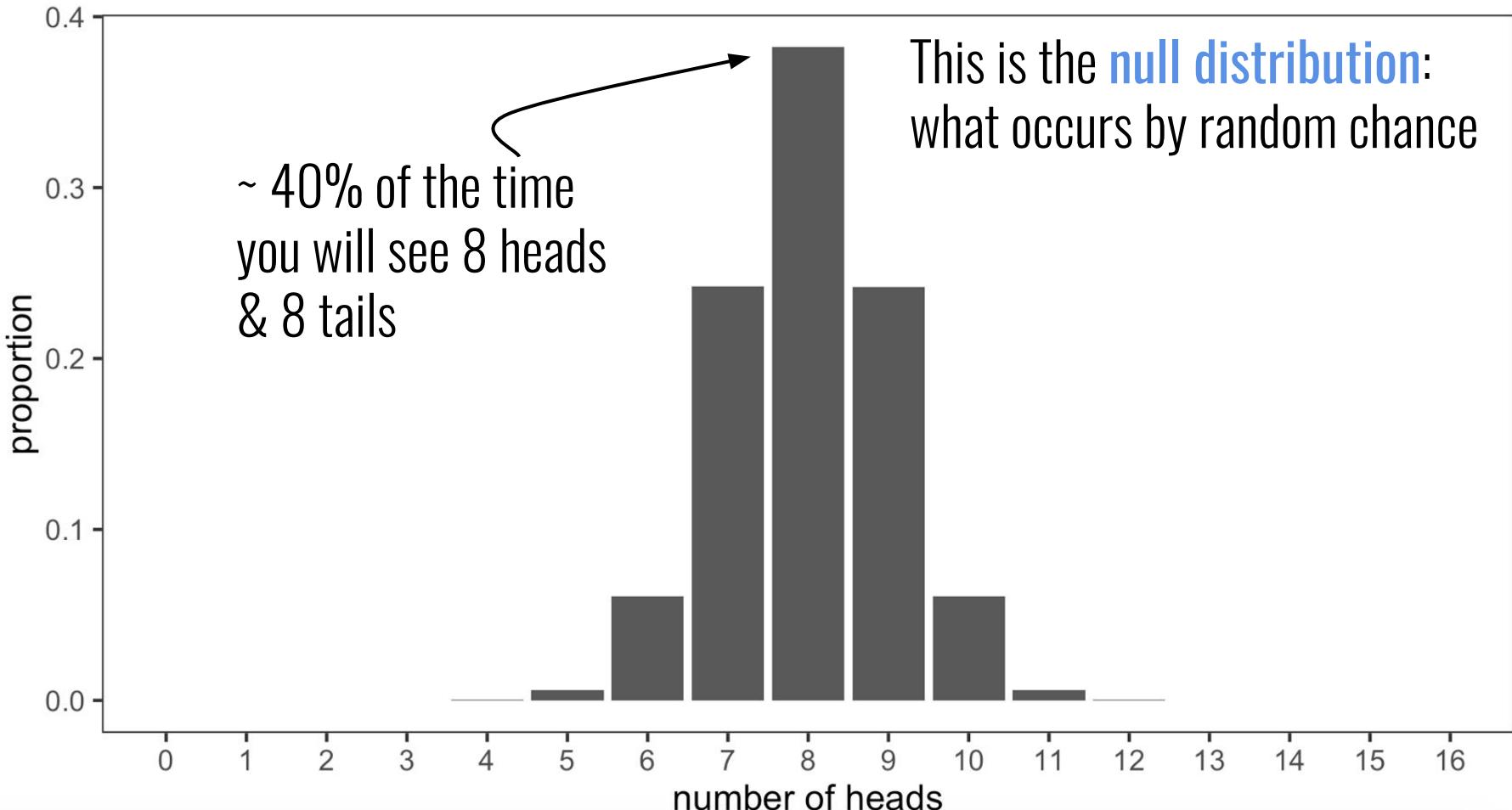
B  
 $> 1.37$

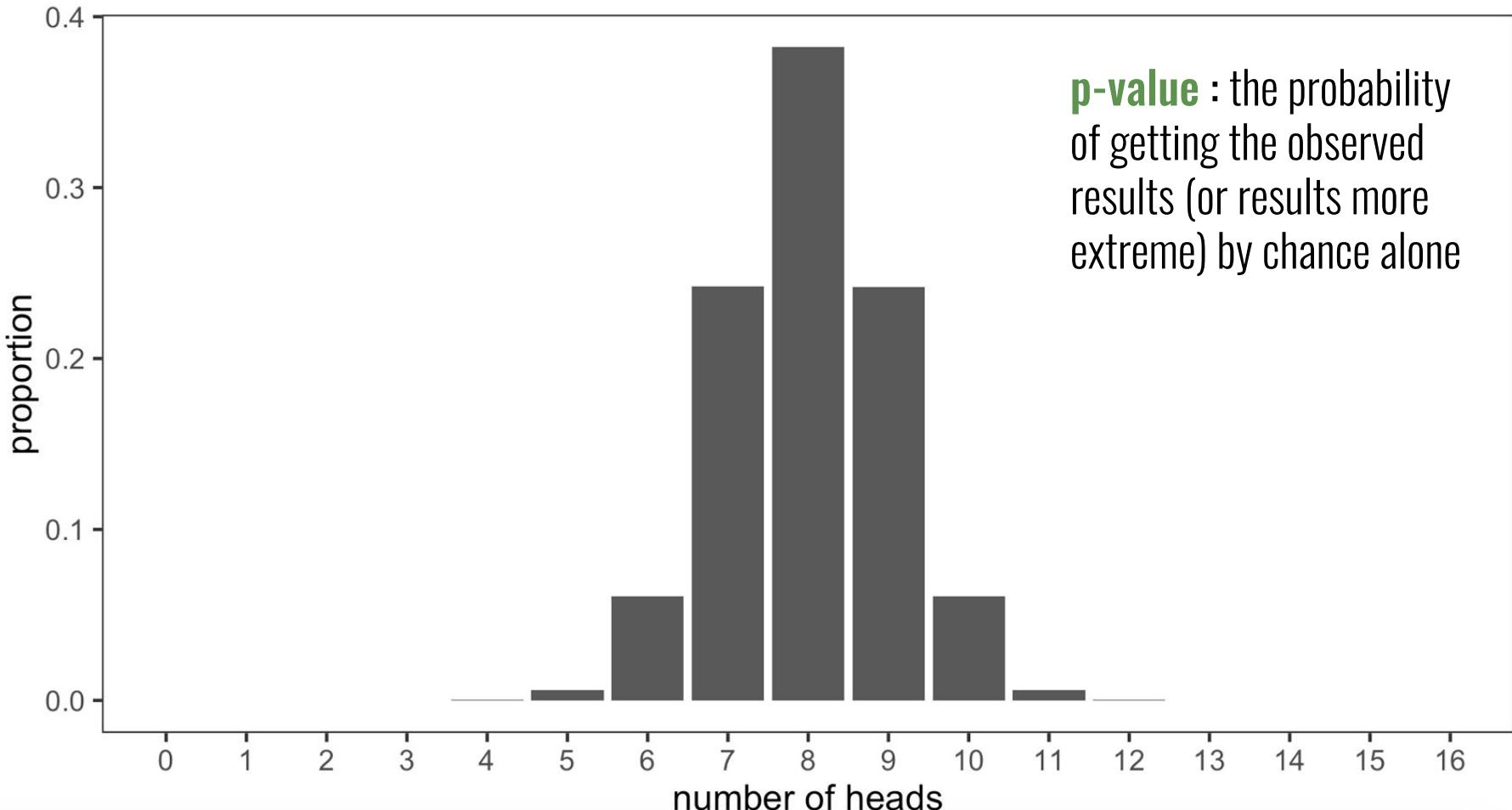




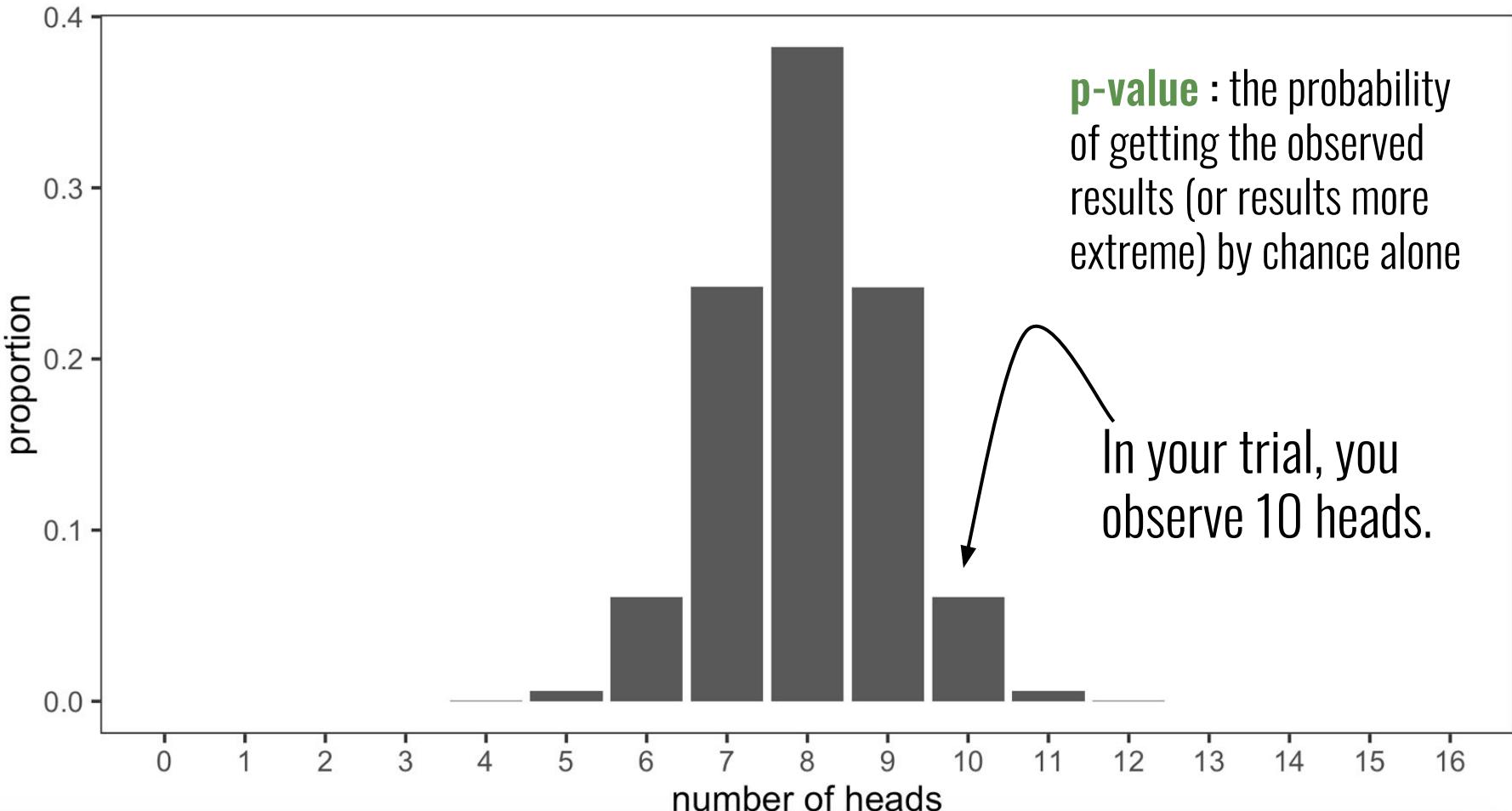
**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

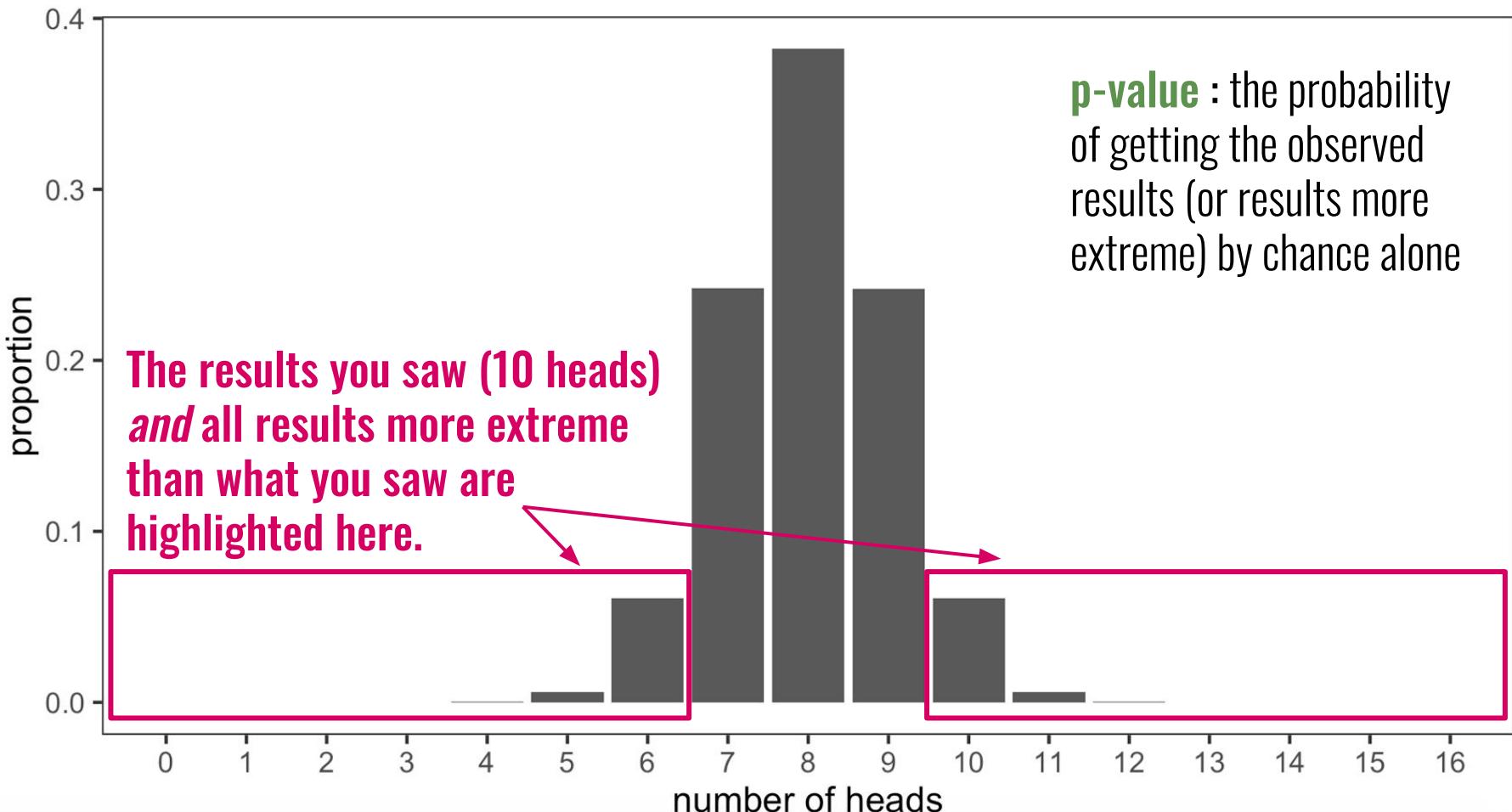


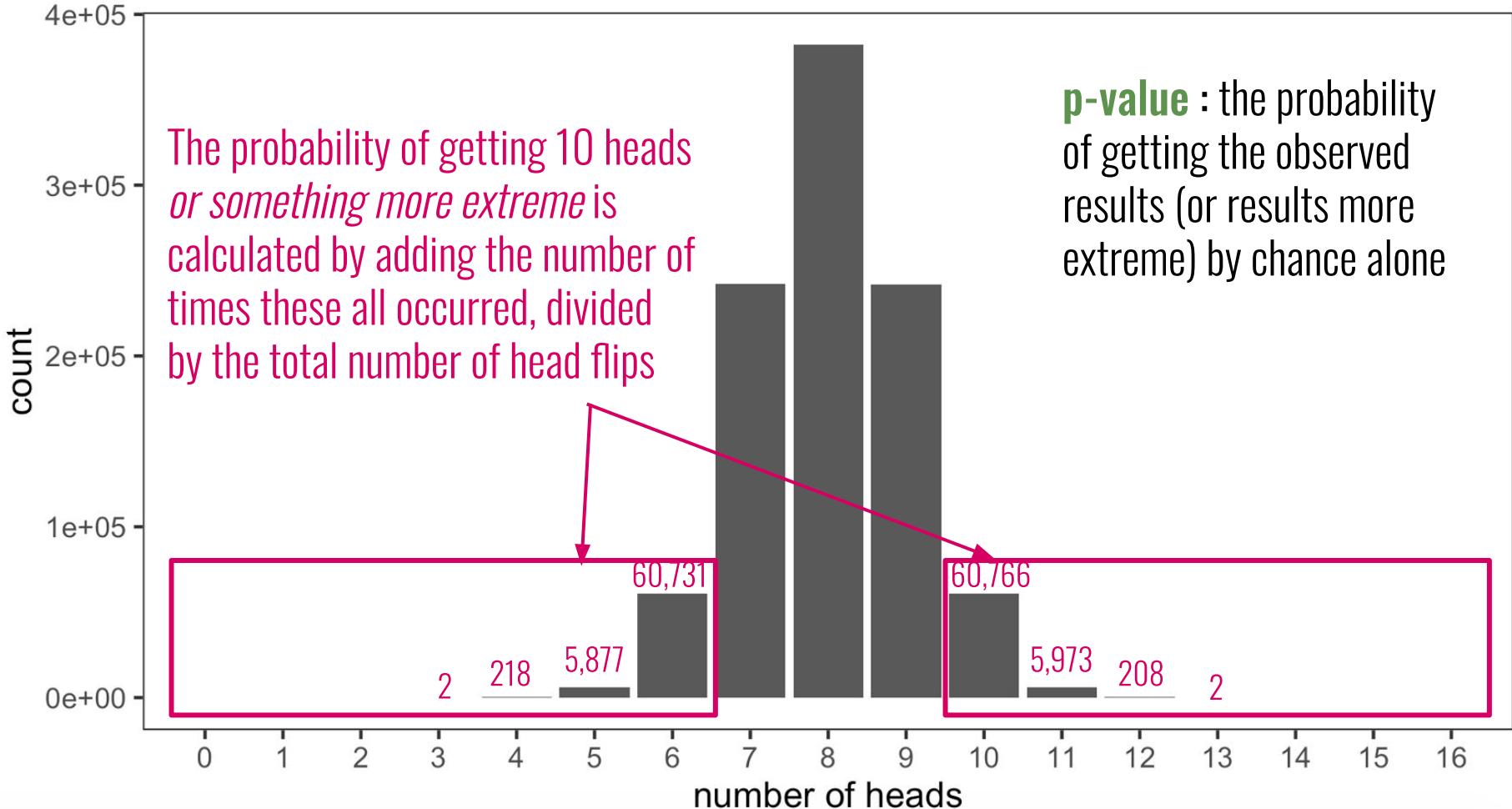


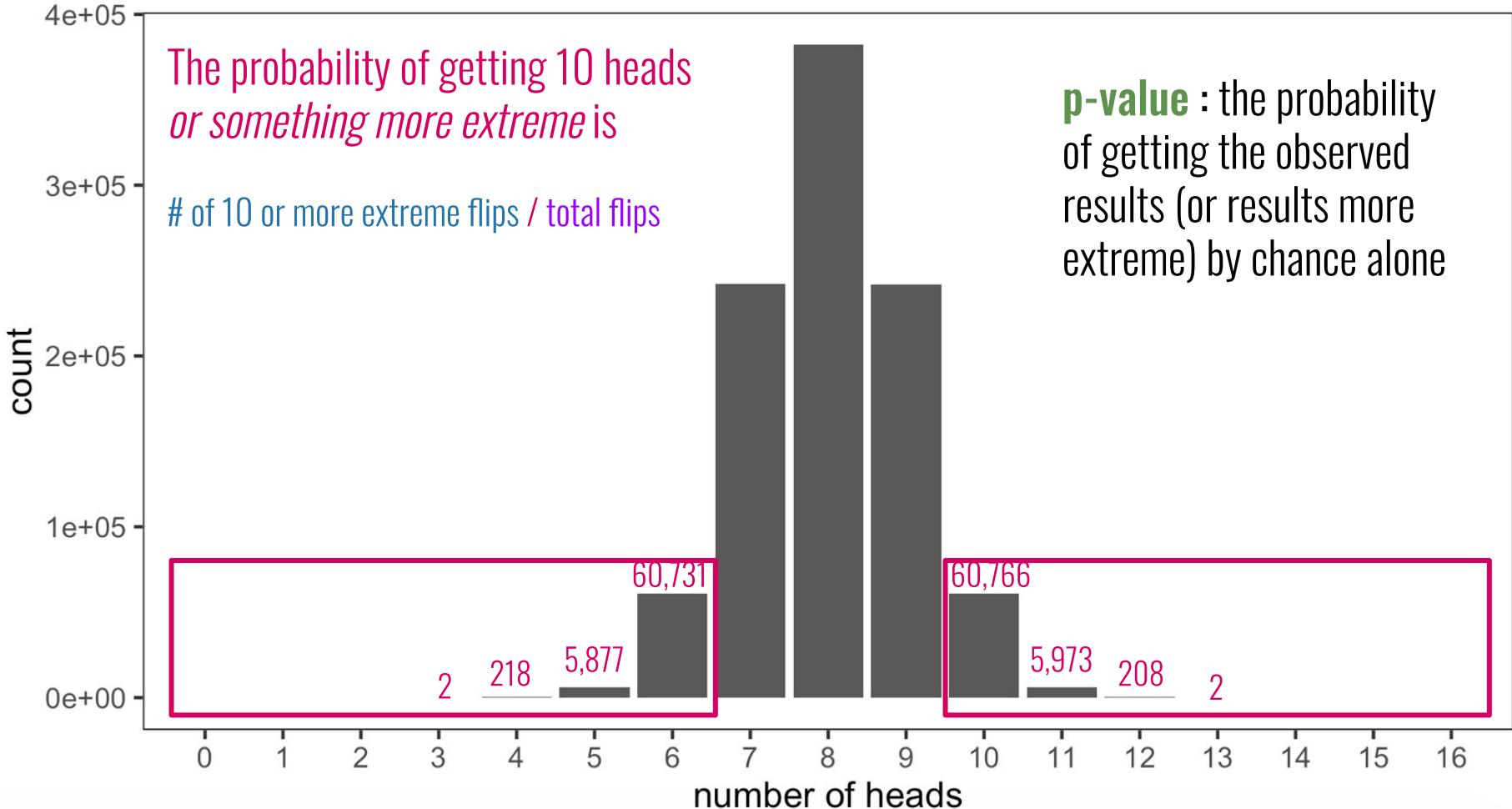


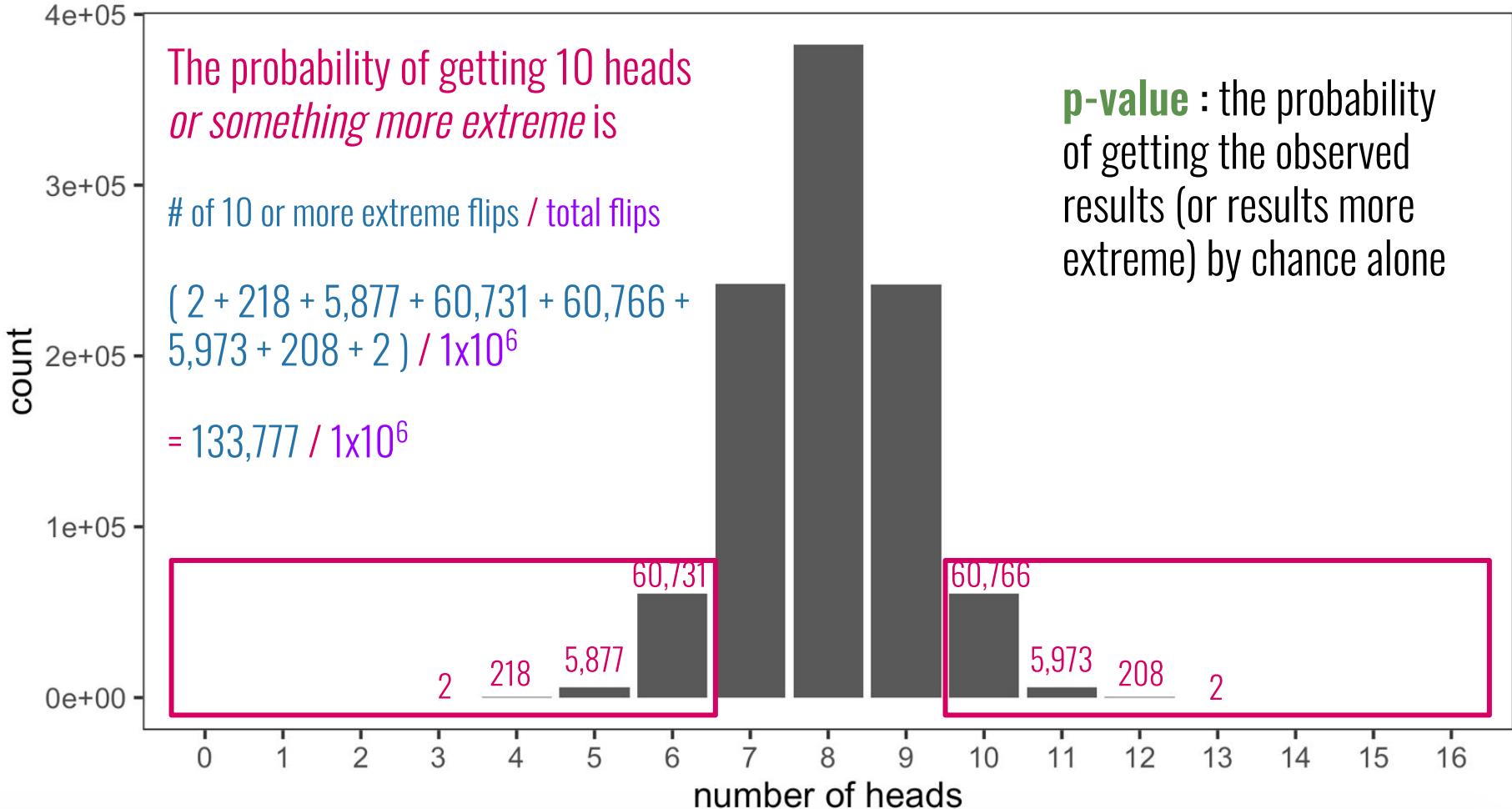
**p-value** : the probability  
of getting the observed  
results (or results more  
extreme) by chance alone

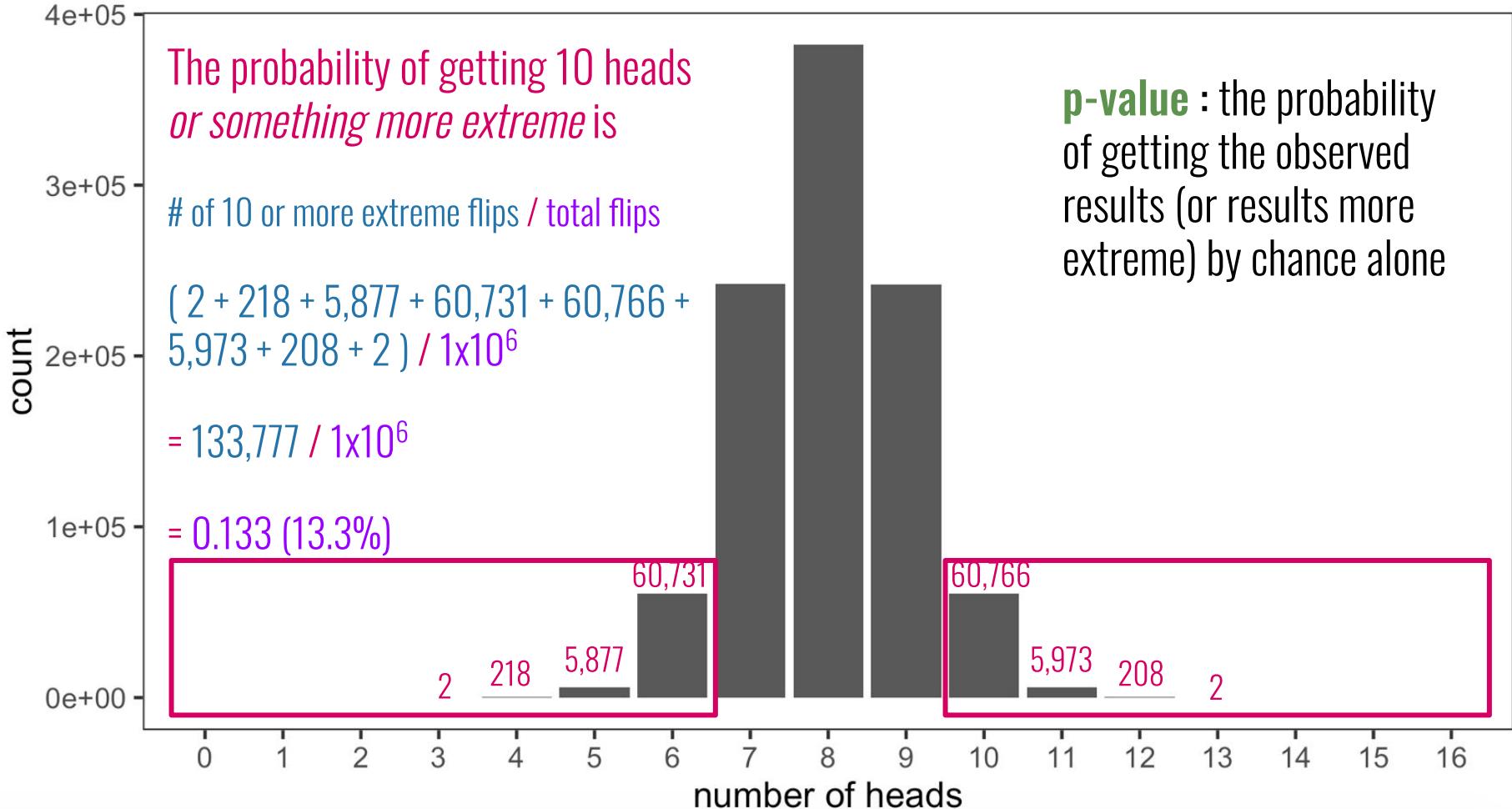


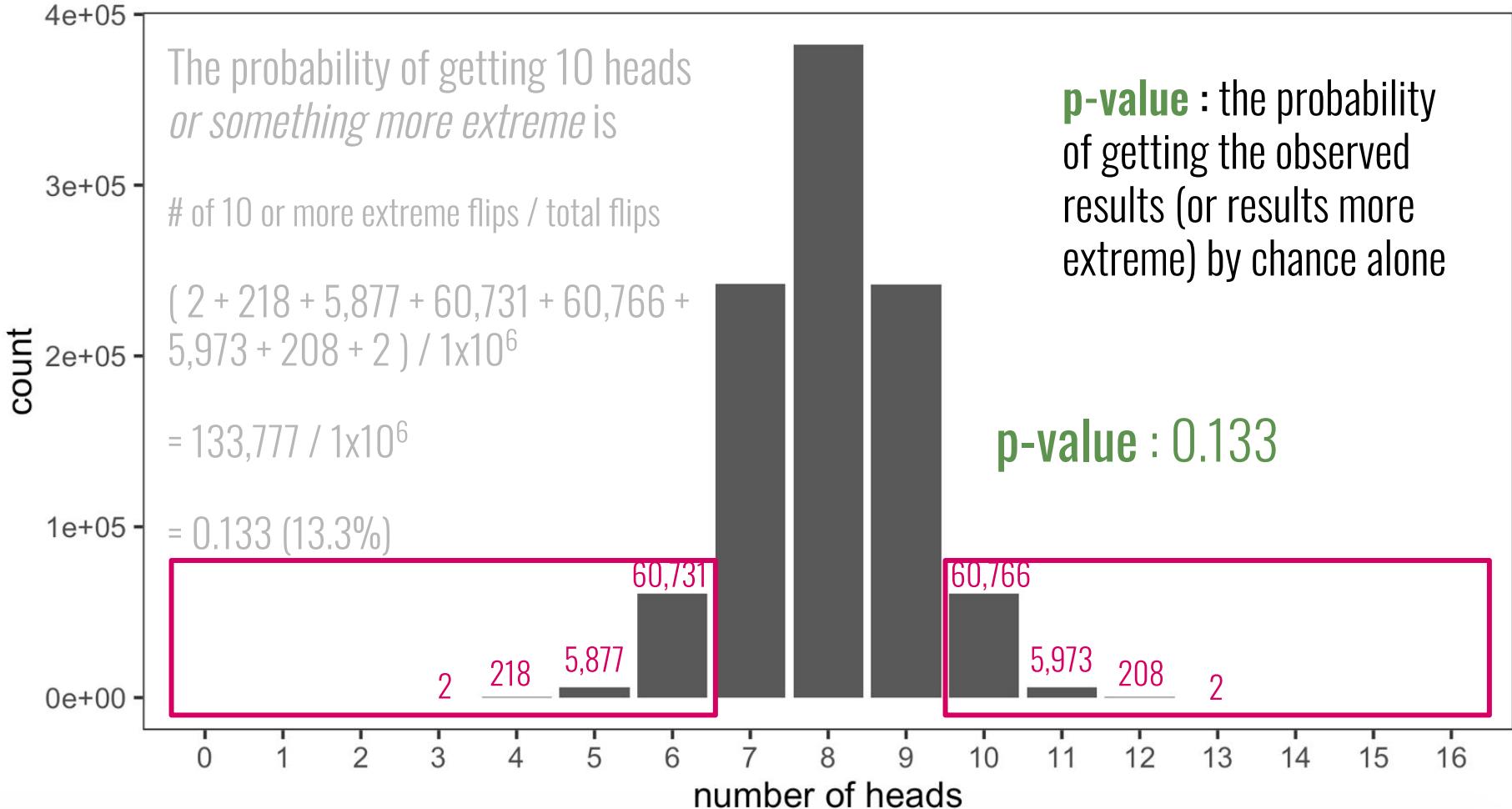


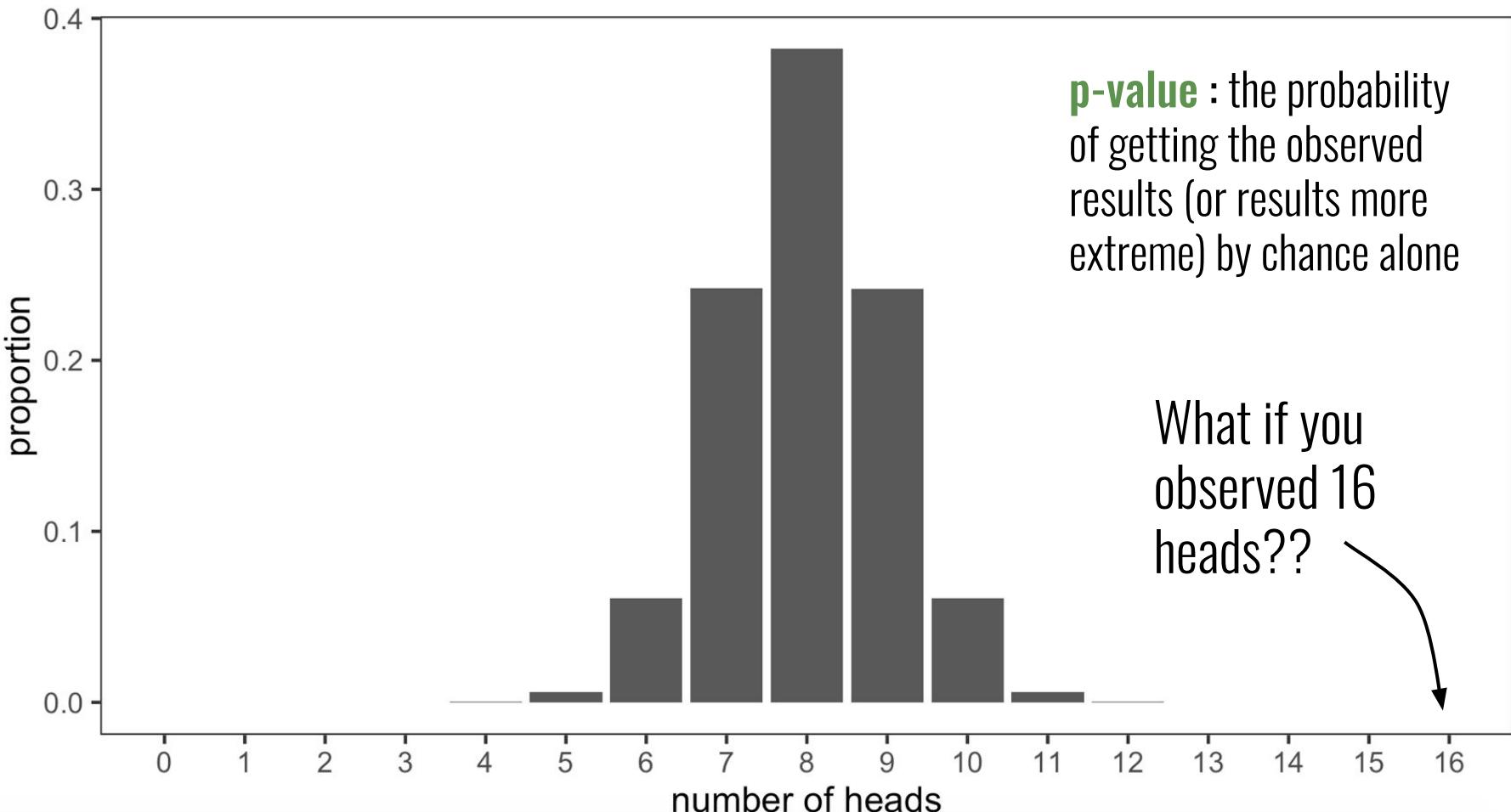


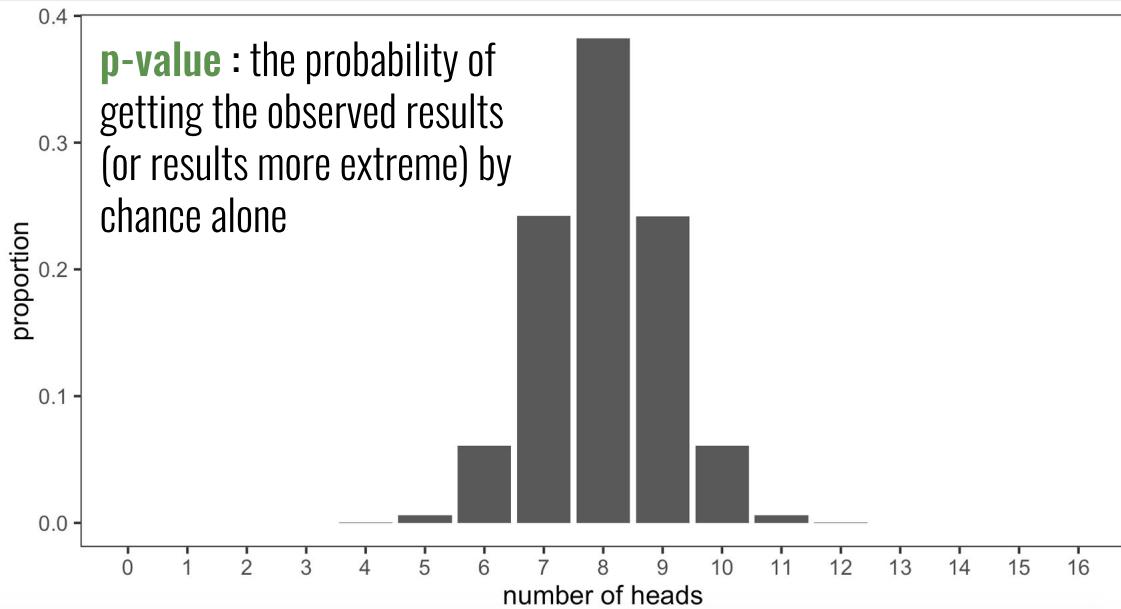












What would be the p-value of you flipping 16 heads?



A

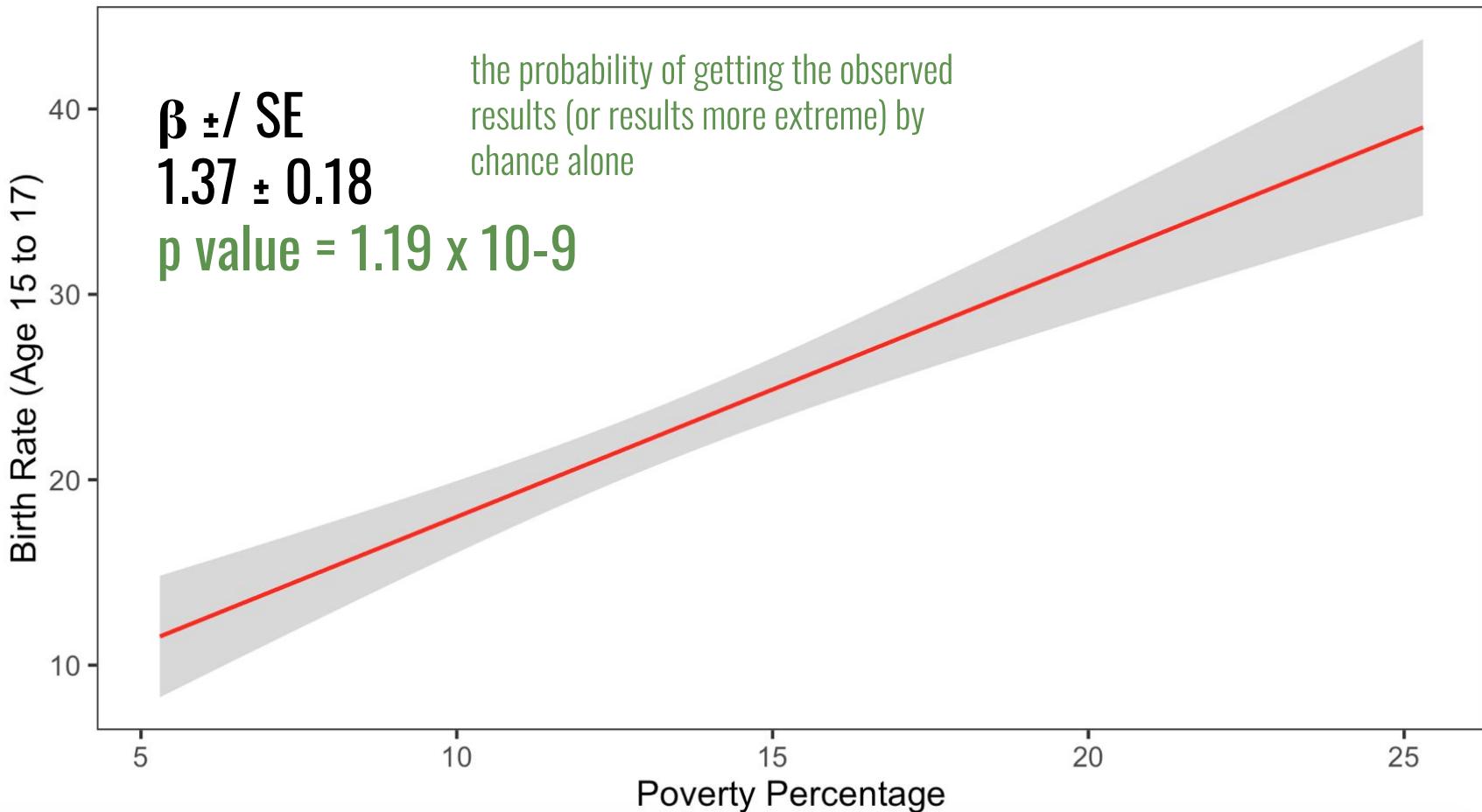
< 0.13



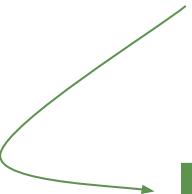
B

> 0.13





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



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

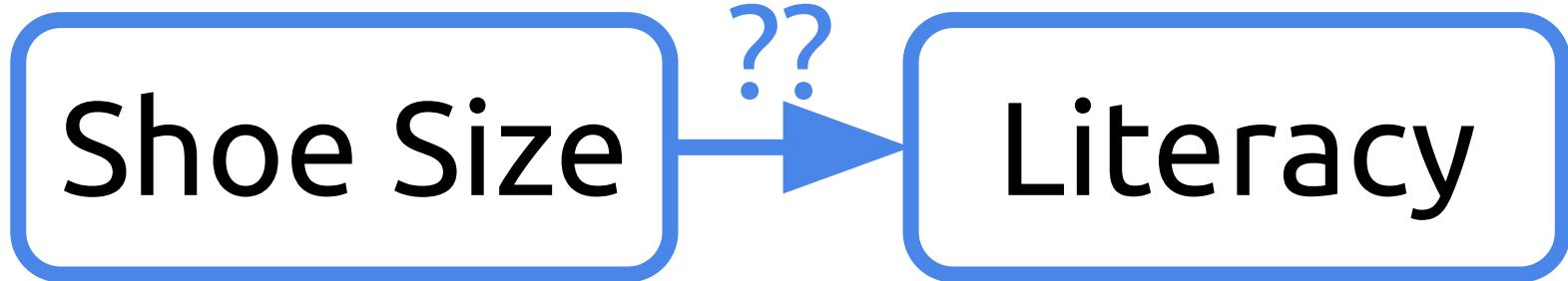
# Confounding





Small shoes  
Not literate

Big shoes  
Literate





Small shoes  
Not literate  
Child

Big shoes  
Literate  
Adult

**Shoe Size**

**Literacy**

**Age**

Variable1

Variable2

Confounder



# Confounding

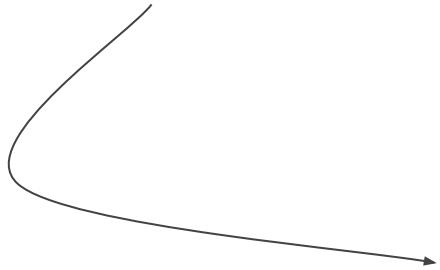
popsicles → crime rate



Your analysis sees an increase in crime rate whenever popsicle sales increase. What could confound this analysis?

- A  
popsicle preference
- B  
new gun laws
- C  
temperature
- D  
changes in  
popsicle prices
- E  
new law  
enforcement  
officers

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)

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

subsequent fractures

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

	<b>Vertebroplasty</b> <b>N = 200</b>	<b>Conservative care</b> <b>N = 200</b>
Age, y, mean $\pm$ SD	$78.2 \pm 4.1$	$79.0 \pm 5.2$
Weight, kg, mean $\pm$ SD	$54.4 \pm 2.3$	$53.9 \pm 2.1$
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