

Course Announcements

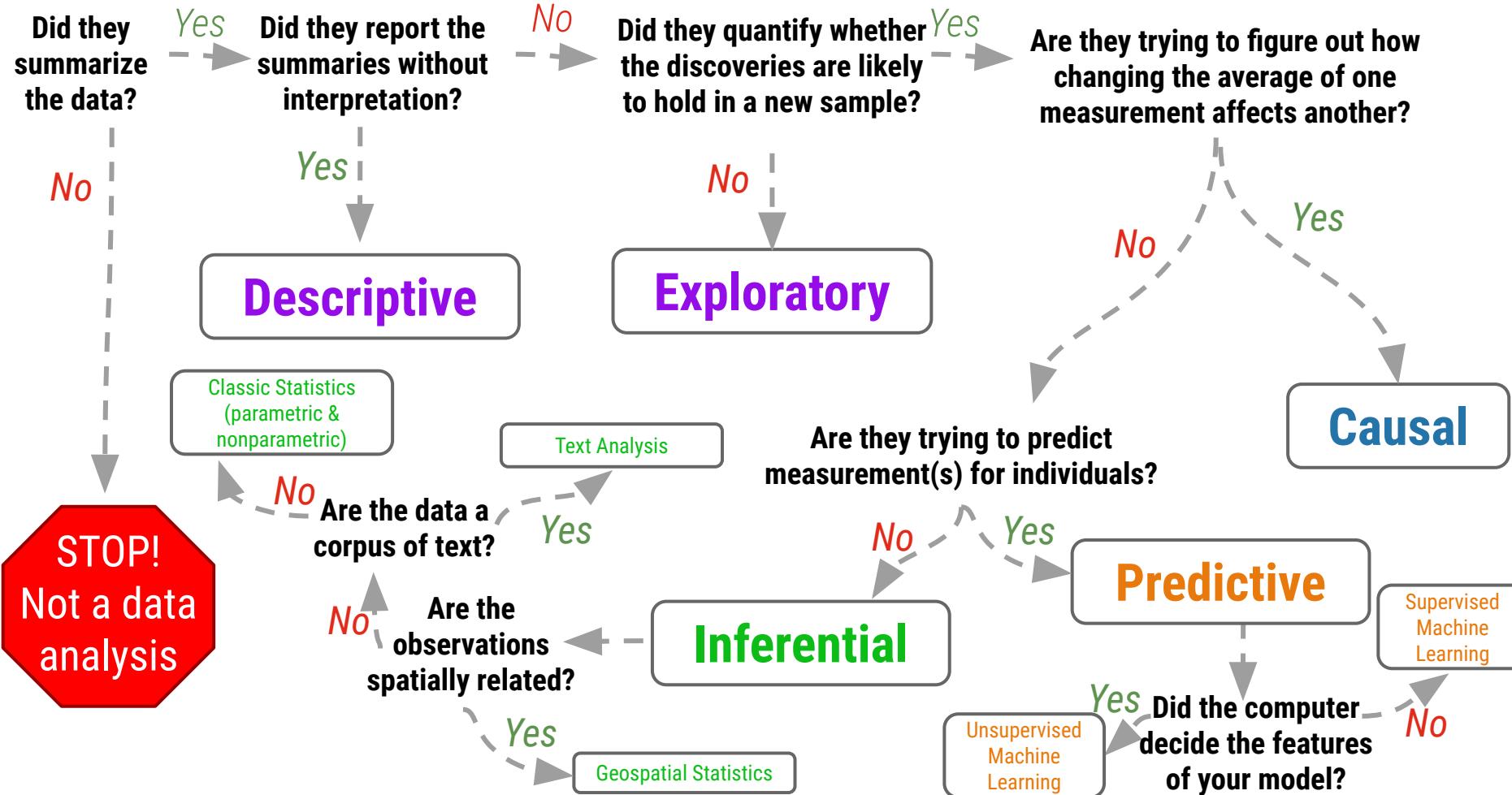
- Guest Lecture Friday!
- Due Friday (11:59 PM)
 - D5
 - Q5
 - Checkpoint #1: Data
 - Demonstrate you have reasonable dataset(s) to answer your question...and that they're in a reasonable format
 - Can be changed/updated/modified later
- A3 will be posted later today
- Project Proposal Feedback will be posted this afternoon
 - Sent as an issue on your GitHub repo
 - Regrades - Campuswire : all will be handled by Prof Ellis
 - Private post to Instructors/TAs
 - Explain where you believe you lost points but should not have
 - Will regrade entire proposal: grade can go up, down, or stay the same
- Grades posted: D4, Q4, mid-course survey EC

Inferential Analysis

Shannon E. Ellis, Ph.D
UC San Diego

• • •

Department of Cognitive Science
sellis@ucsd.edu





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

No

Descriptive

Classic Statistics
(parametric & nonparametric)

Text Analysis

STOP!
Not a data analysis

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

Causal

No → Are the data a corpus of text?

Yes

No → Are the observations spatially related?

Yes

Inferential

Geospatial Statistics

No

Yes

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

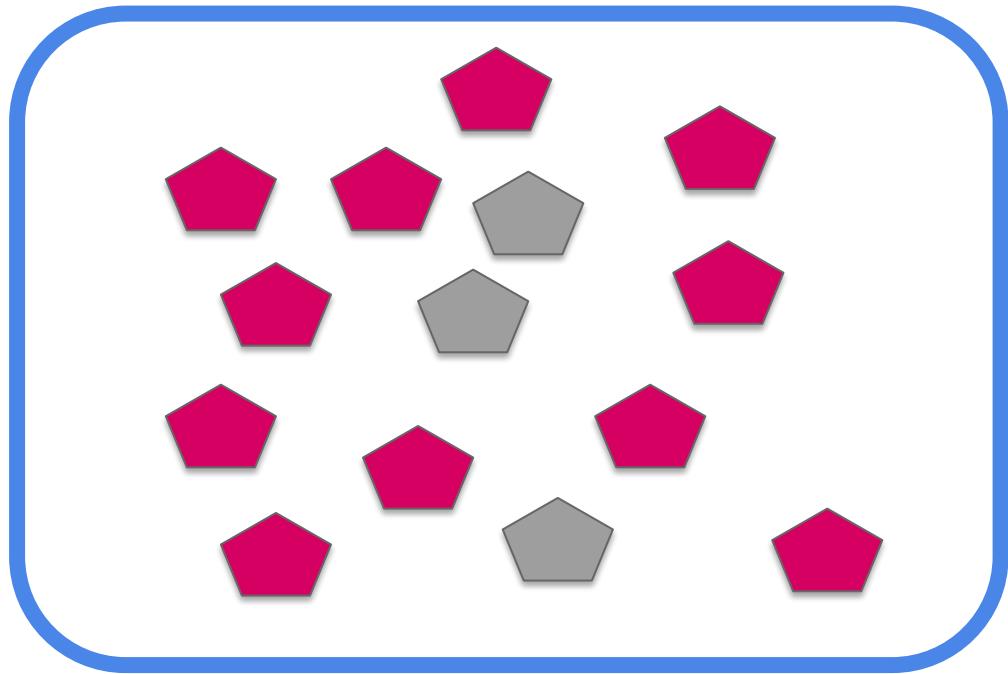
Establishing & Stating Your Null and Alternative Hypotheses Helps Guide Your Analysis

Null Hypothesis:

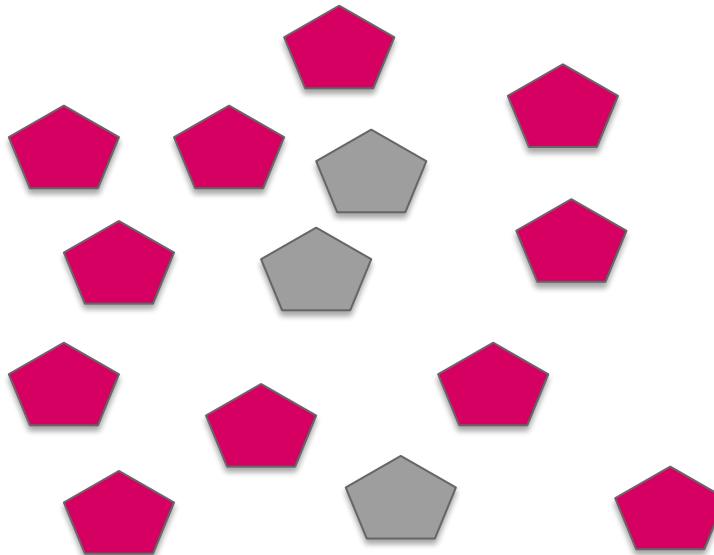
H_0 : Sesame Street has *no effect* on kids brain development

Alternative Hypothesis:

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



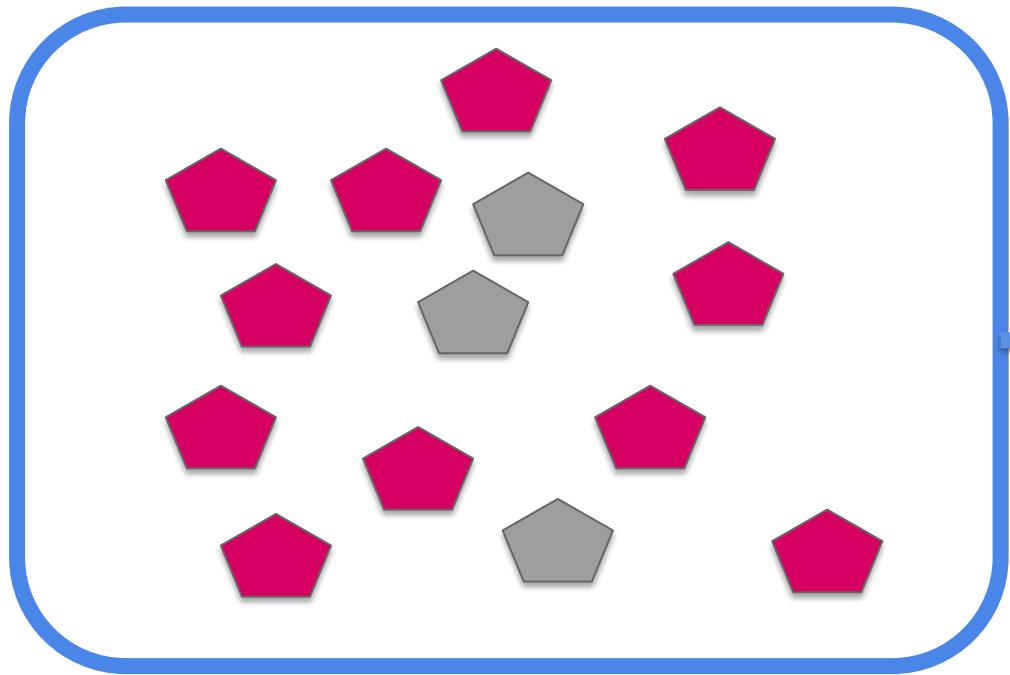
Population



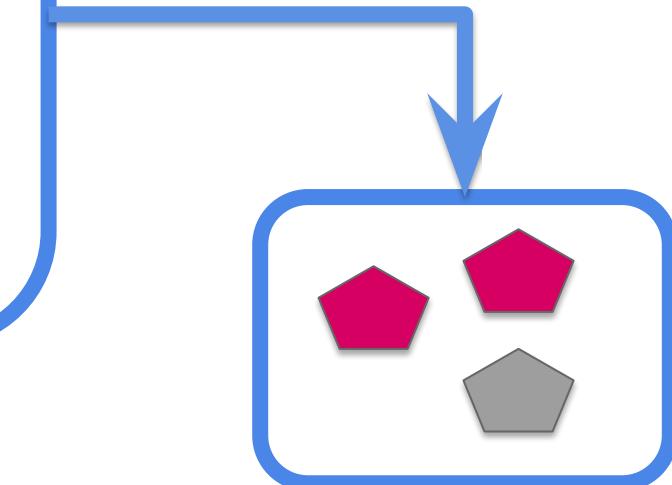
Population



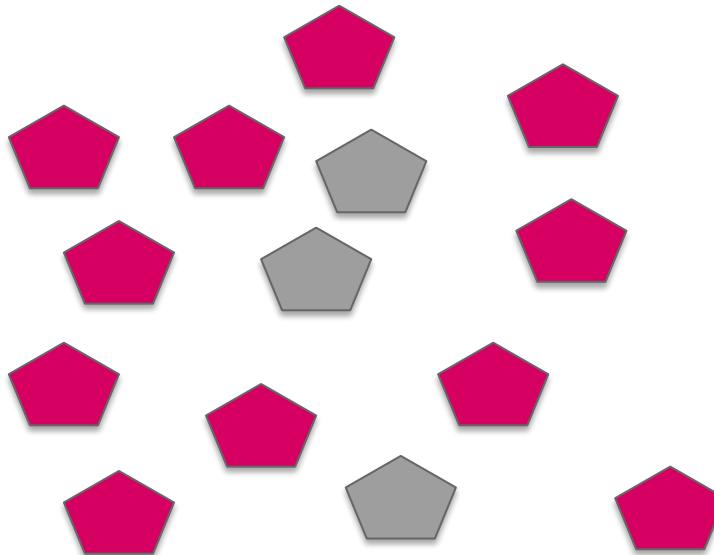
In our Sesame street example, the population would be all children



Population



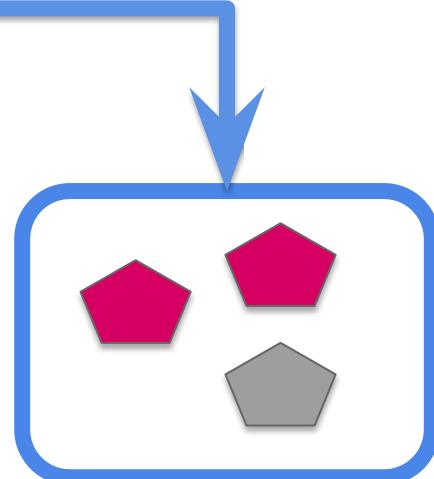
Sample



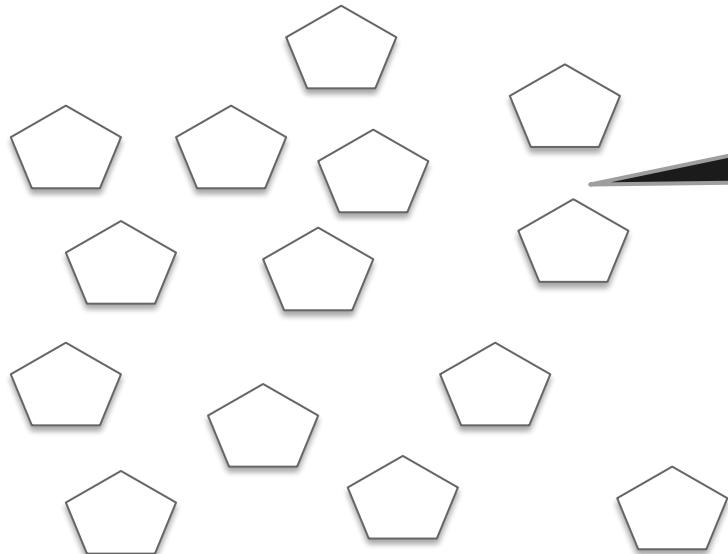
Population



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

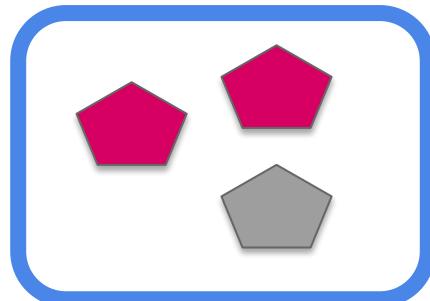


Sample



Population

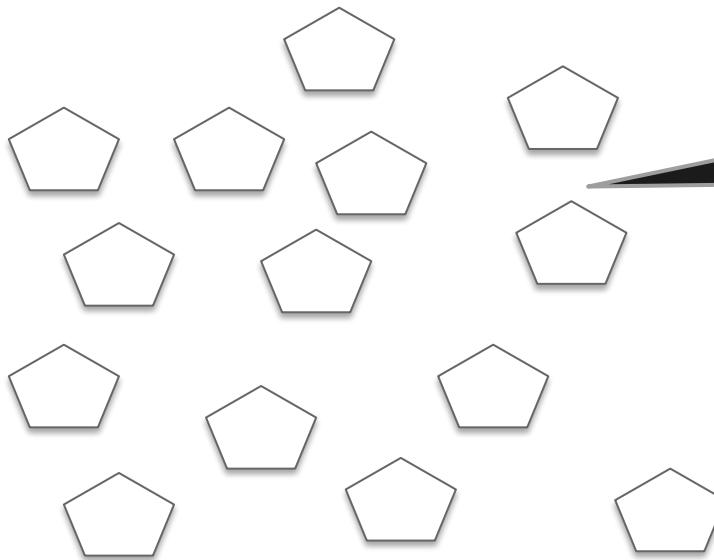
（ツ）



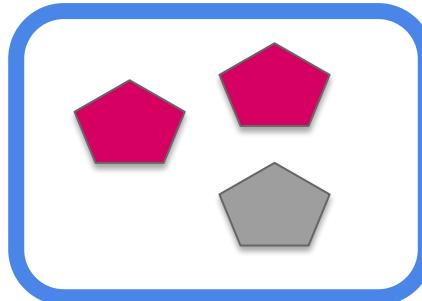
Sample



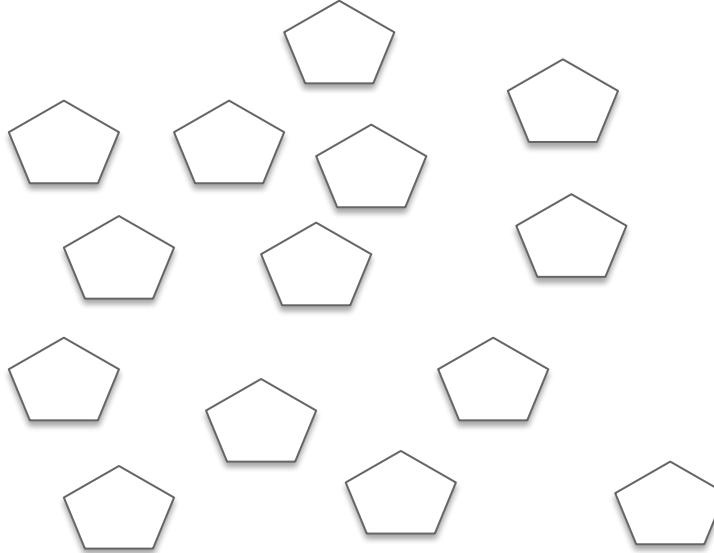
We don't know how much Sesame street was watched by or the tests scores of all kids



Population



Sample

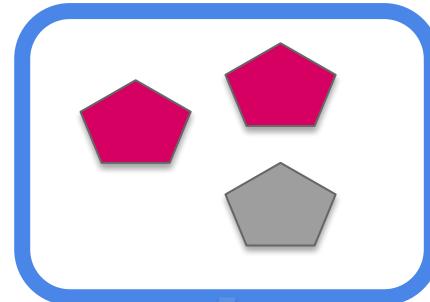


Population

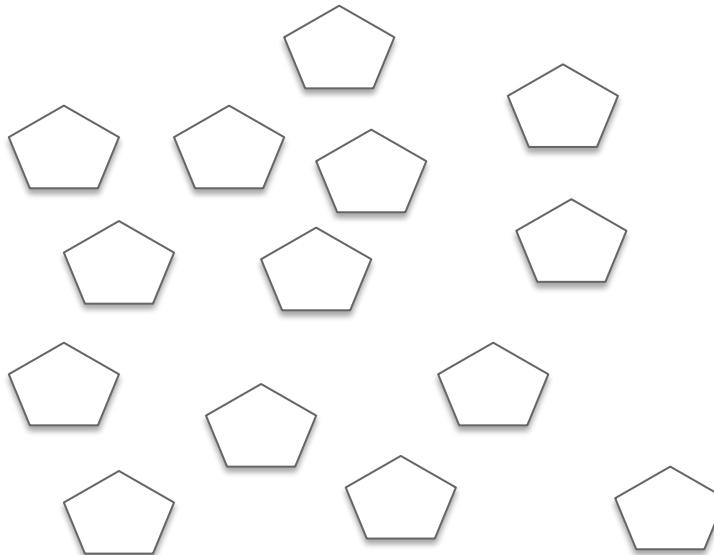


Inference!

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



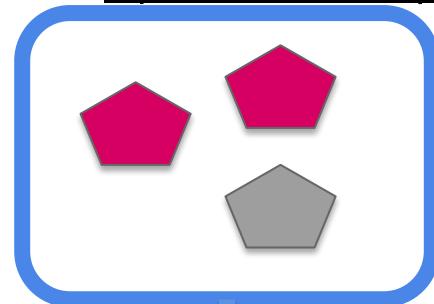
Sample



Population



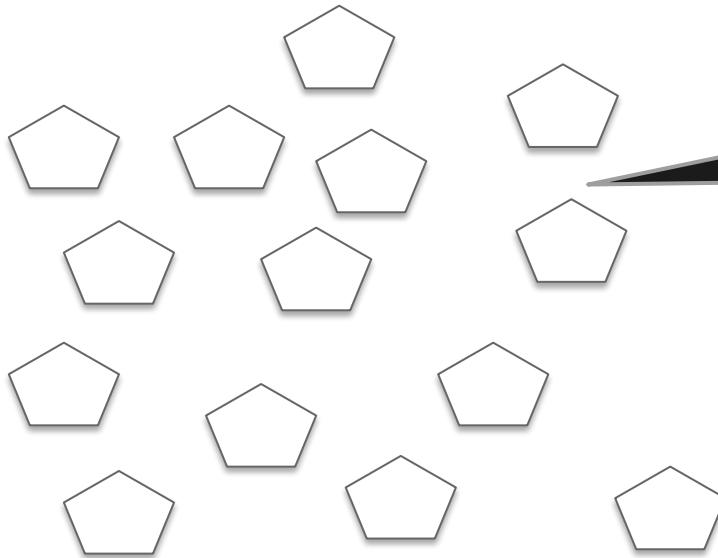
So we look at Sesame street viewing and test scores in a representative sample of kids



Inference!

Sample

Population



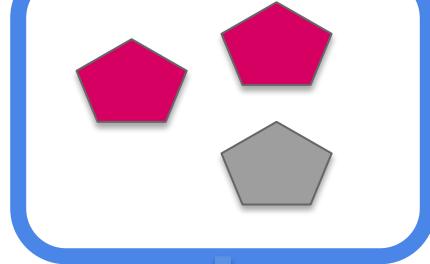
Best guess

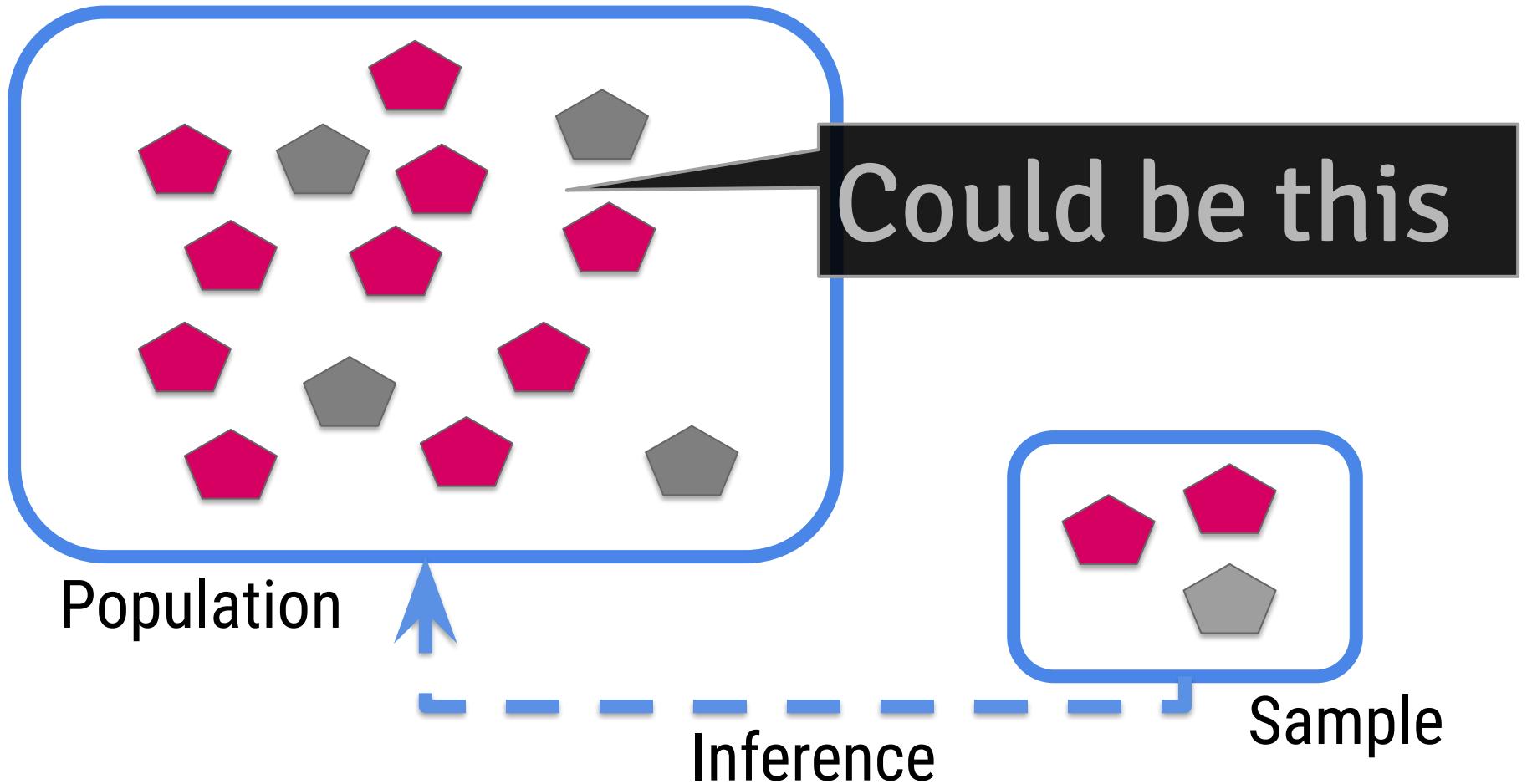


So we look at Sesame street viewing and test scores in a representative sample of kids

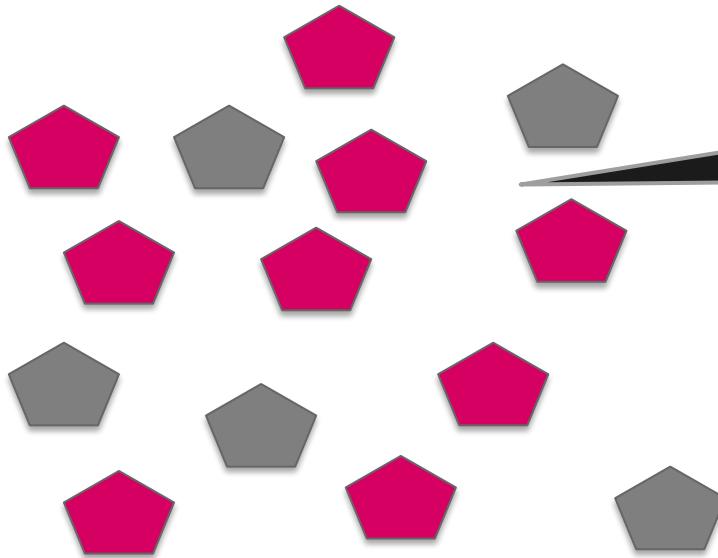
Inference!

Sample



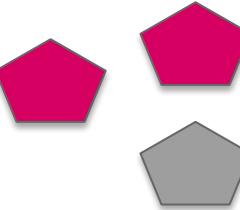


Population

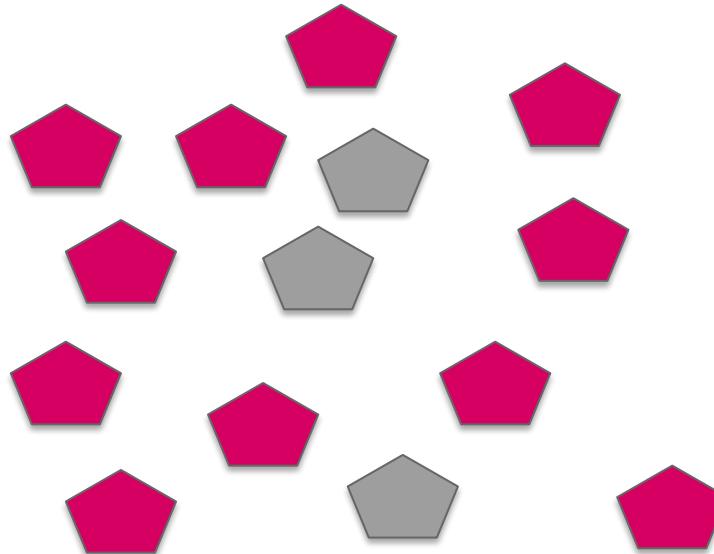


Inference

...or this

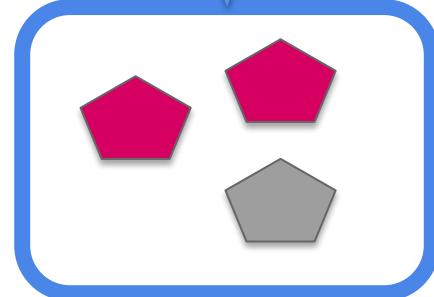


Sample

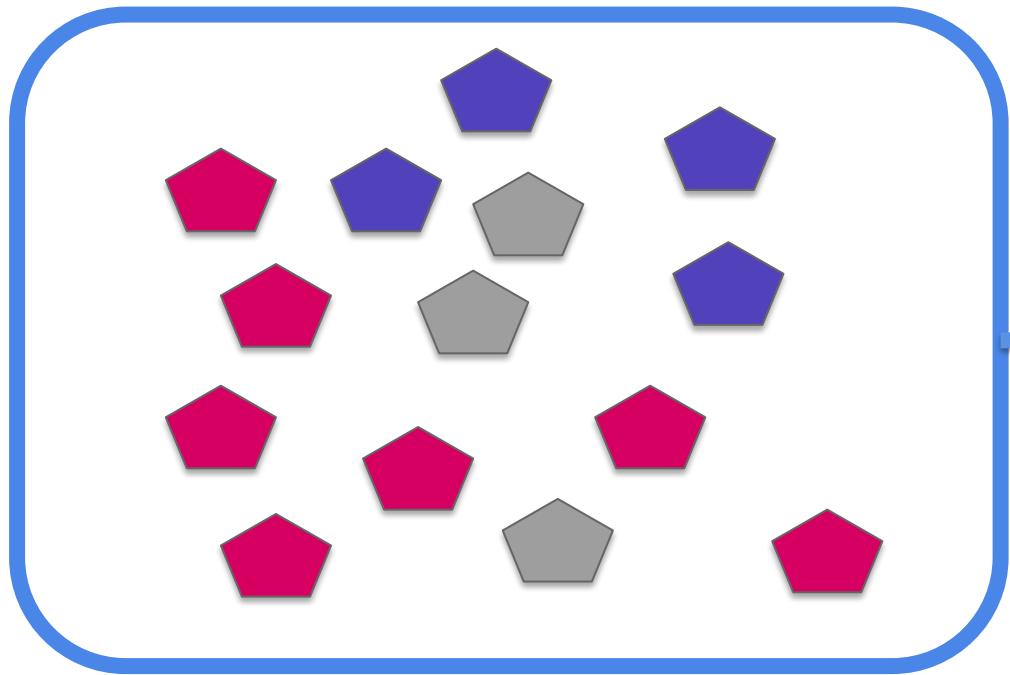


Population

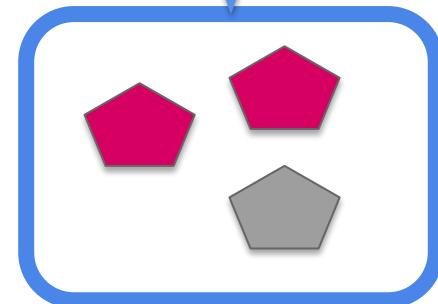
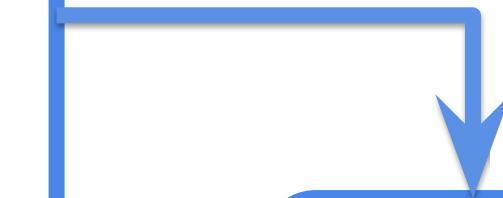
Probability



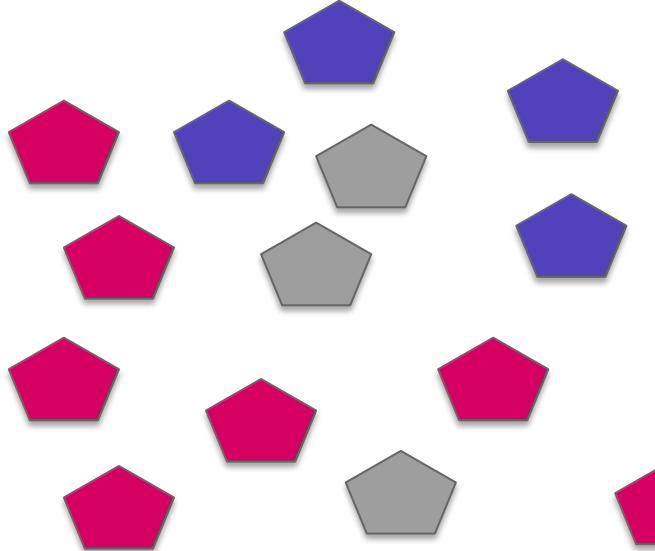
Sample



Population



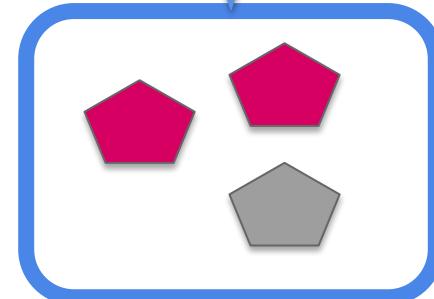
Sample



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

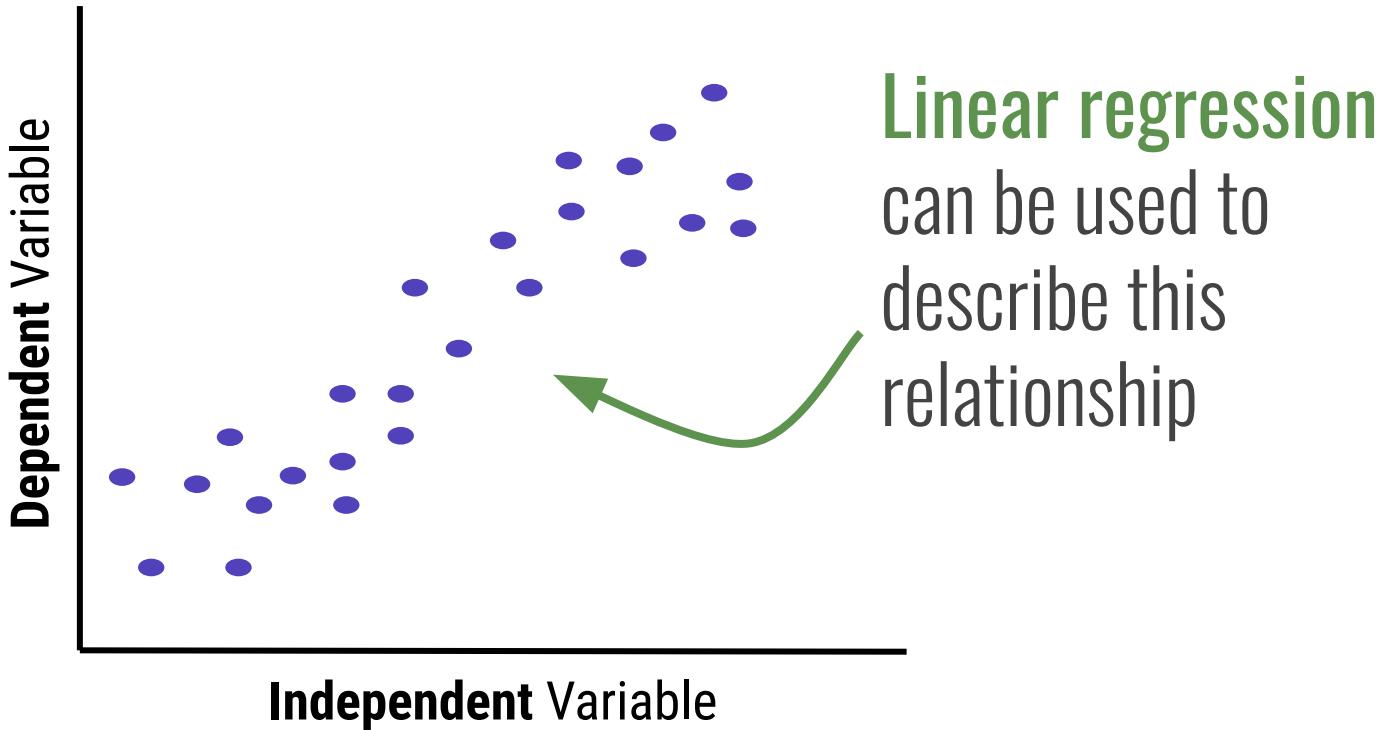
DOES CHANGE IN ONE VARIABLE MEAN CHANGE IN ANOTHER?

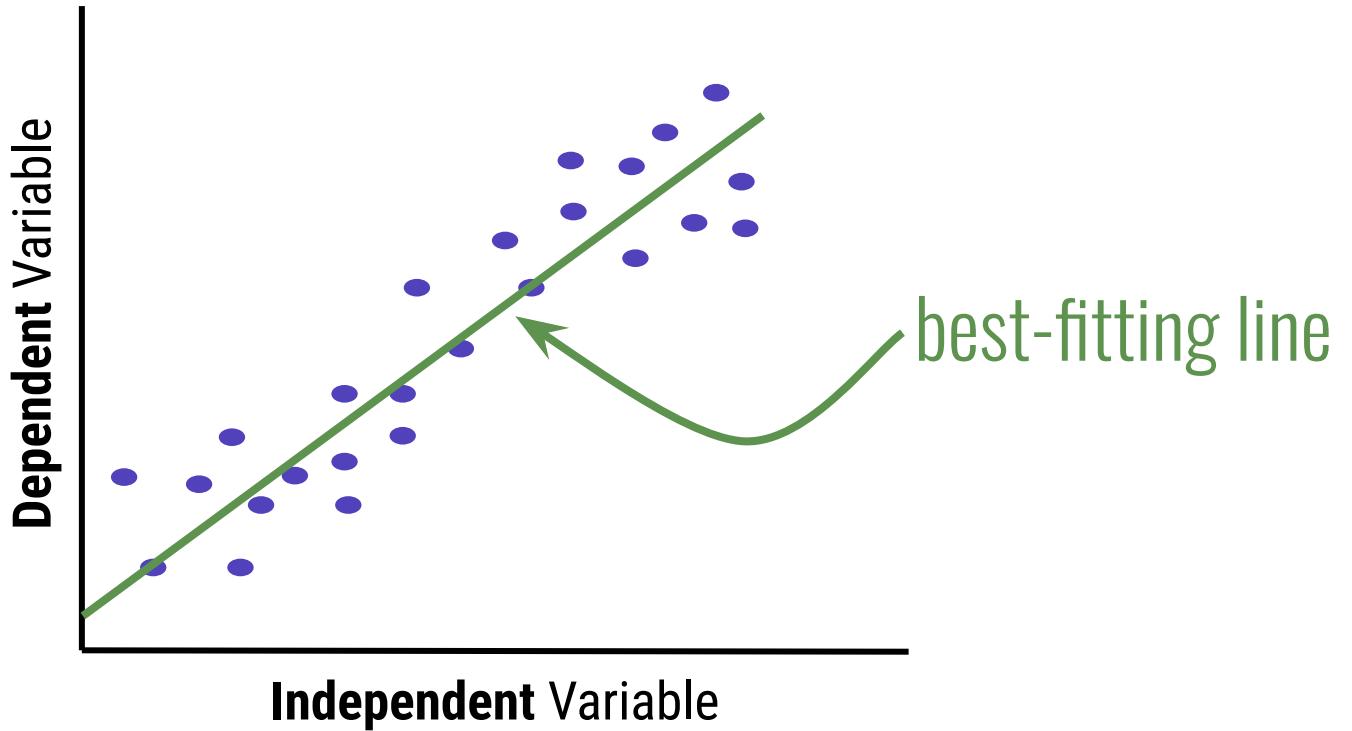
i.e. simple regression,
multiple regression

NON-PARAMETRIC TESTS

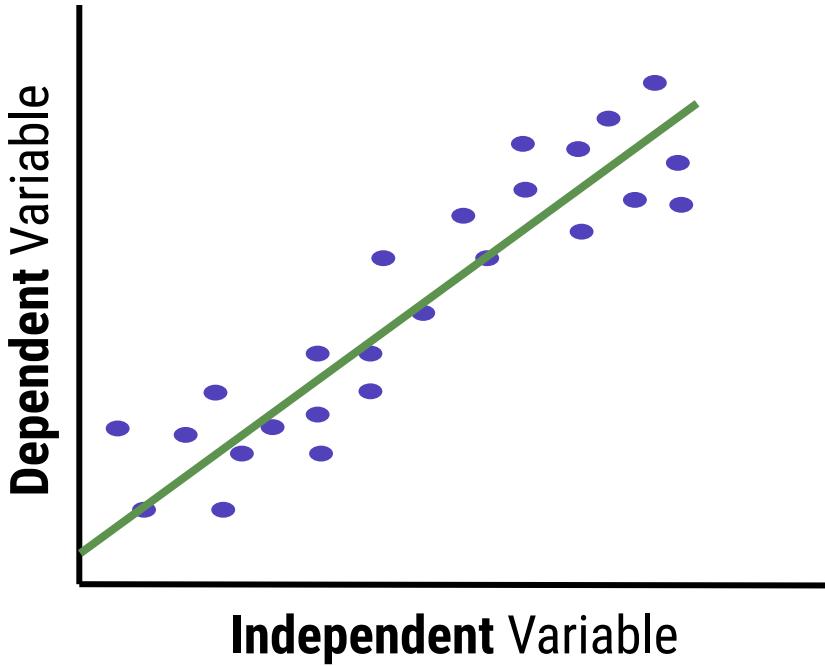
FOR WHEN ASSUMPTIONS IN THESE OTHER 3 CATEGORIES ARE NOT MET

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

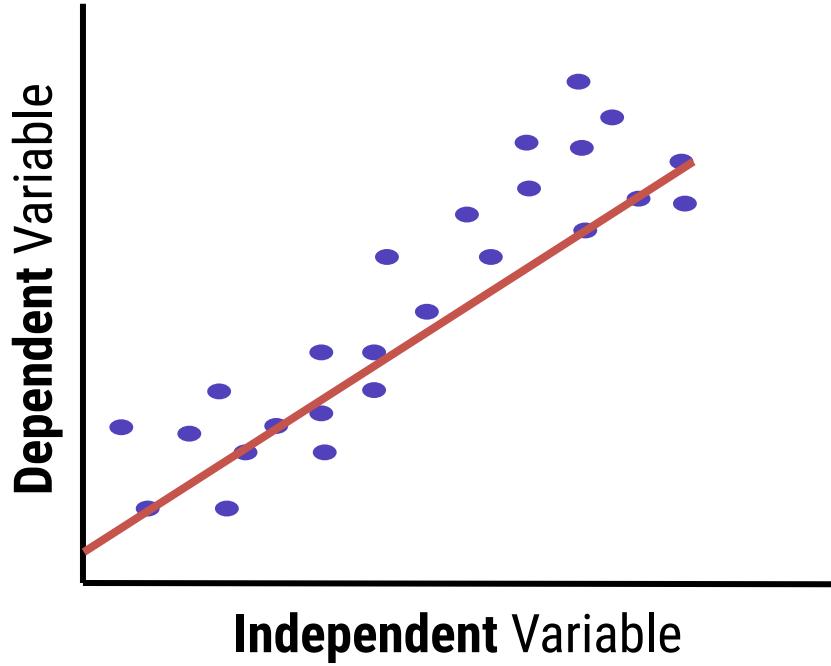


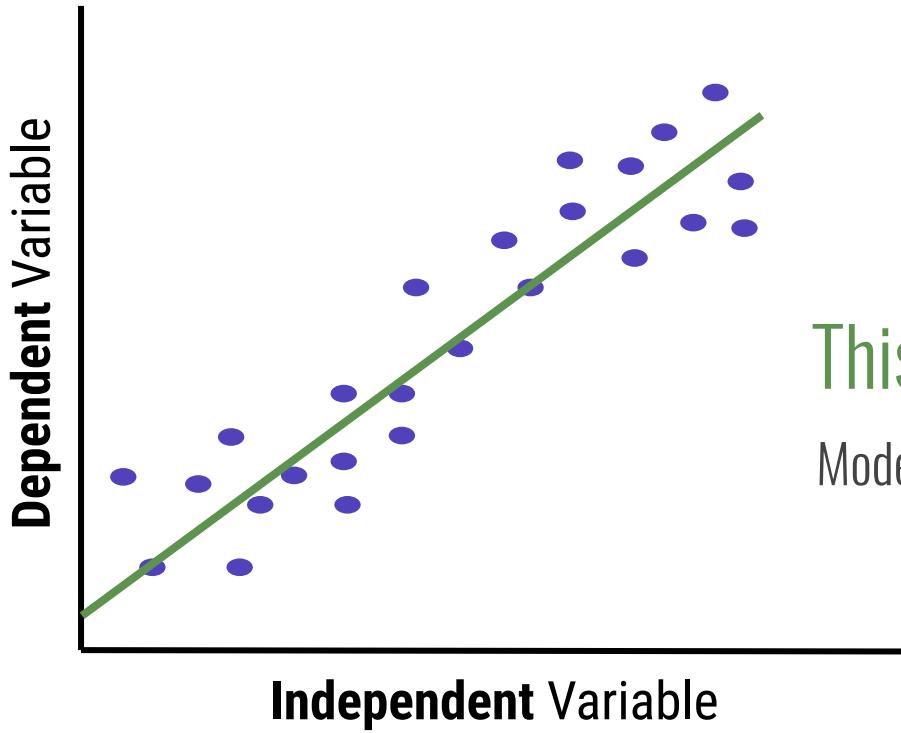


Best-fitting line



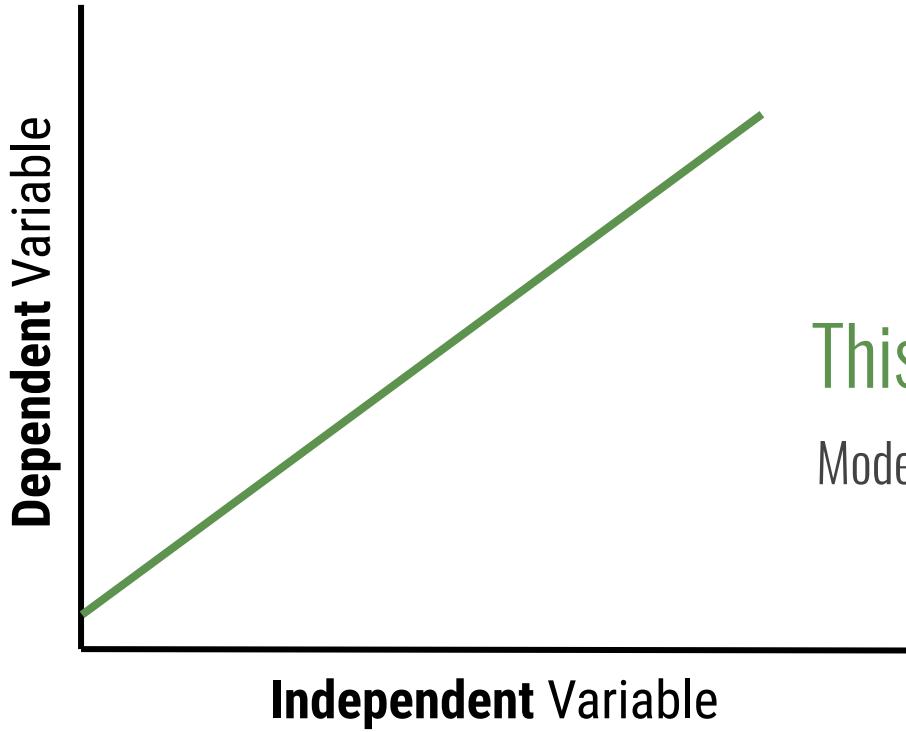
NOT a best-fitting line





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

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

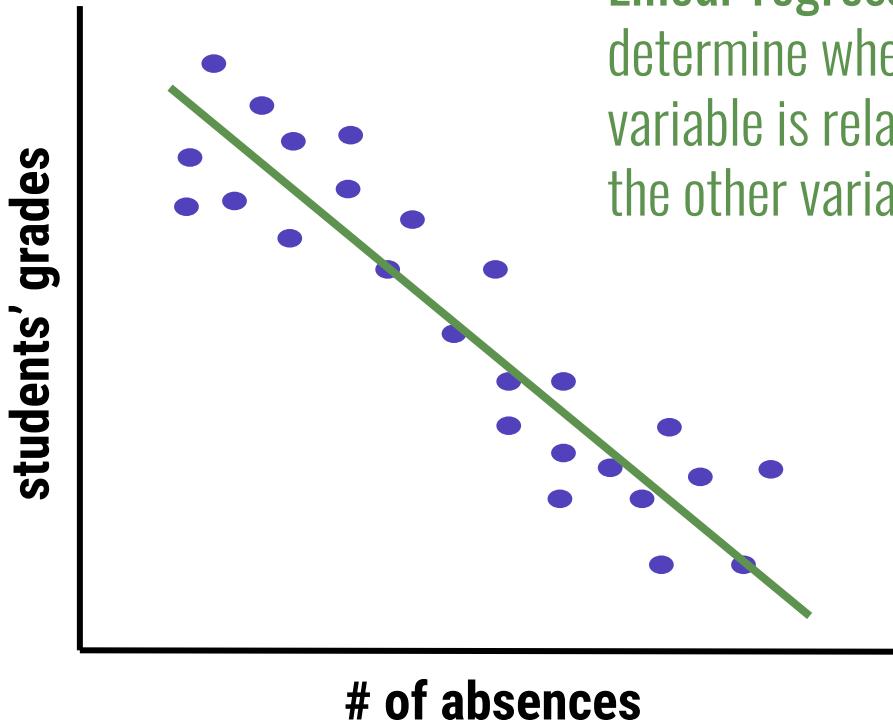


This line is a **model** of the data

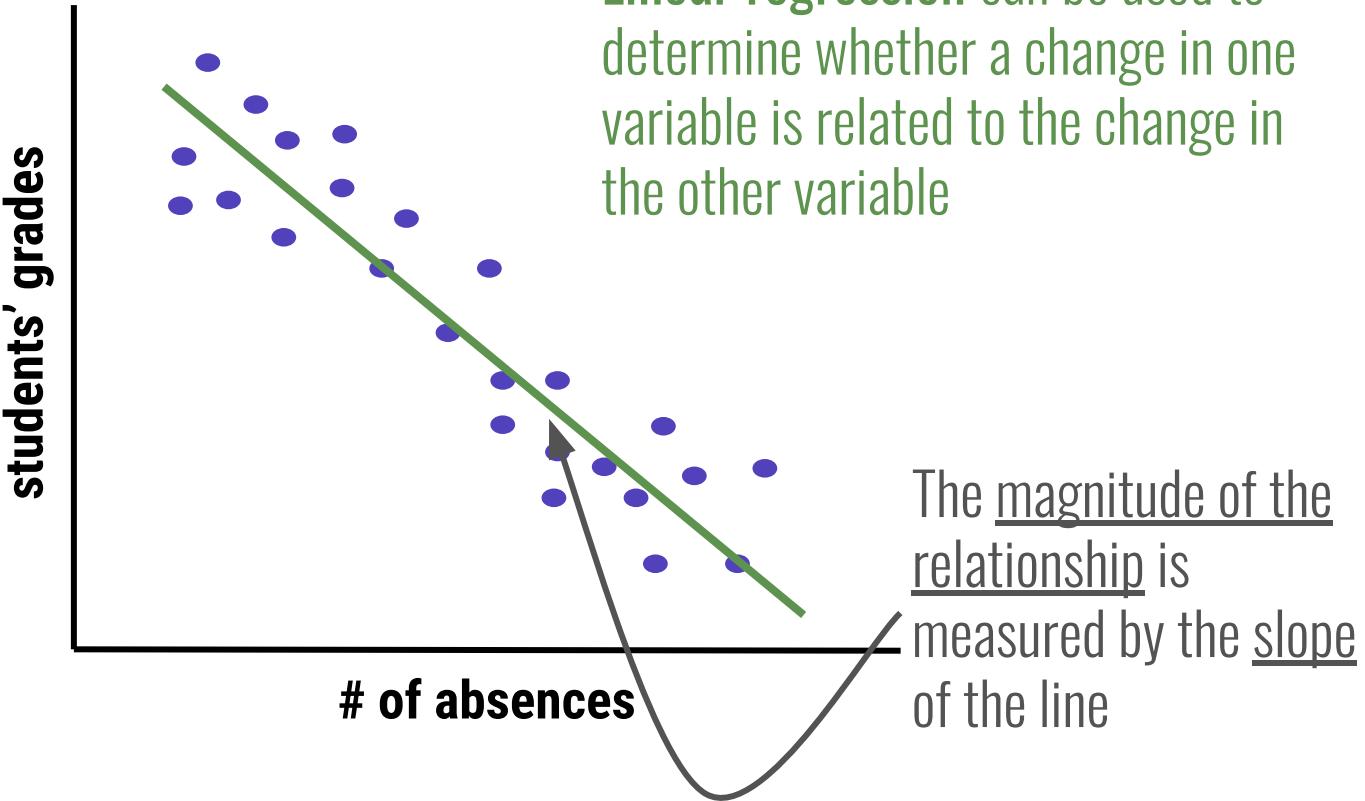
Models are mathematical equations generated
to *represent* the real life situation

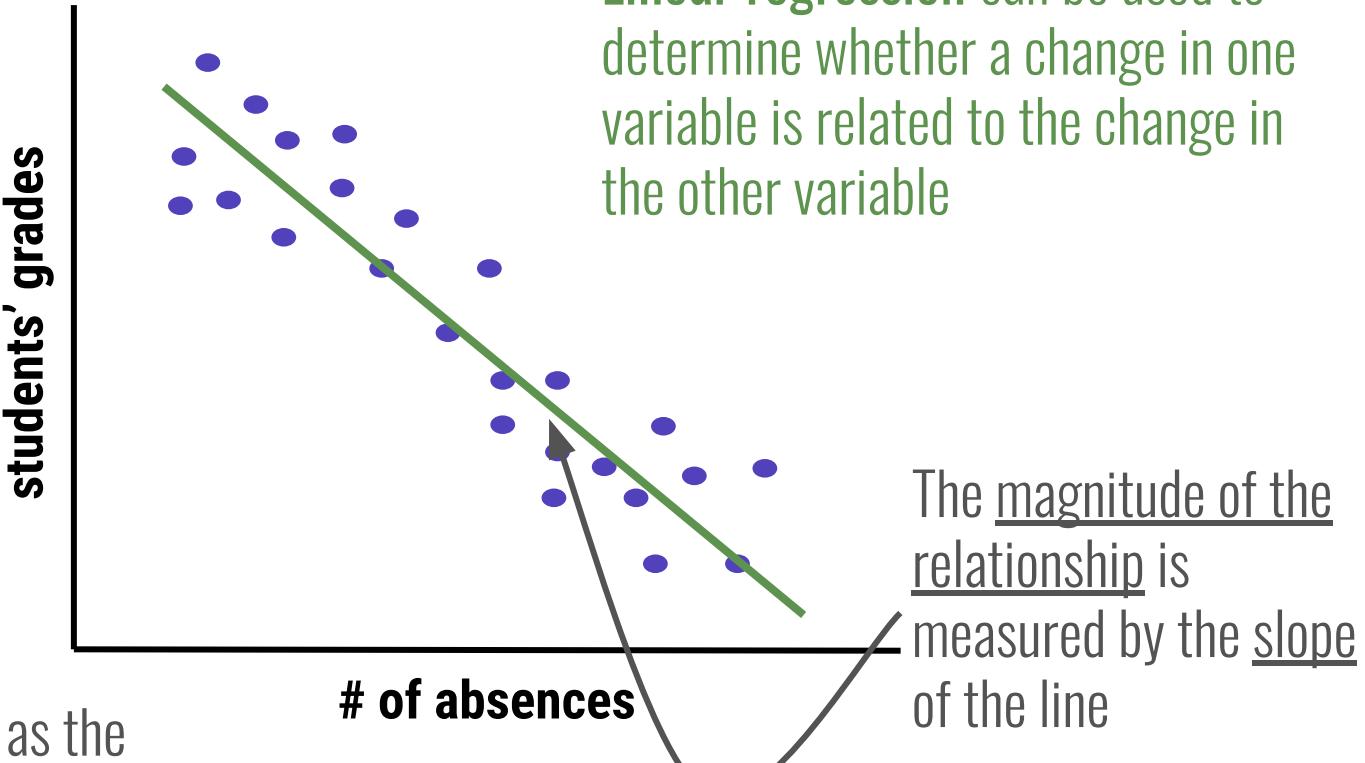
“All models are wrong, but some are useful”

-George Box (British Statistician, JASA 1976)

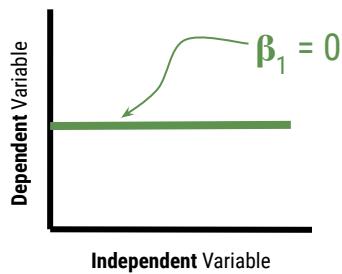


Linear regression can be used to determine whether a change in one variable is related to the change in the other variable

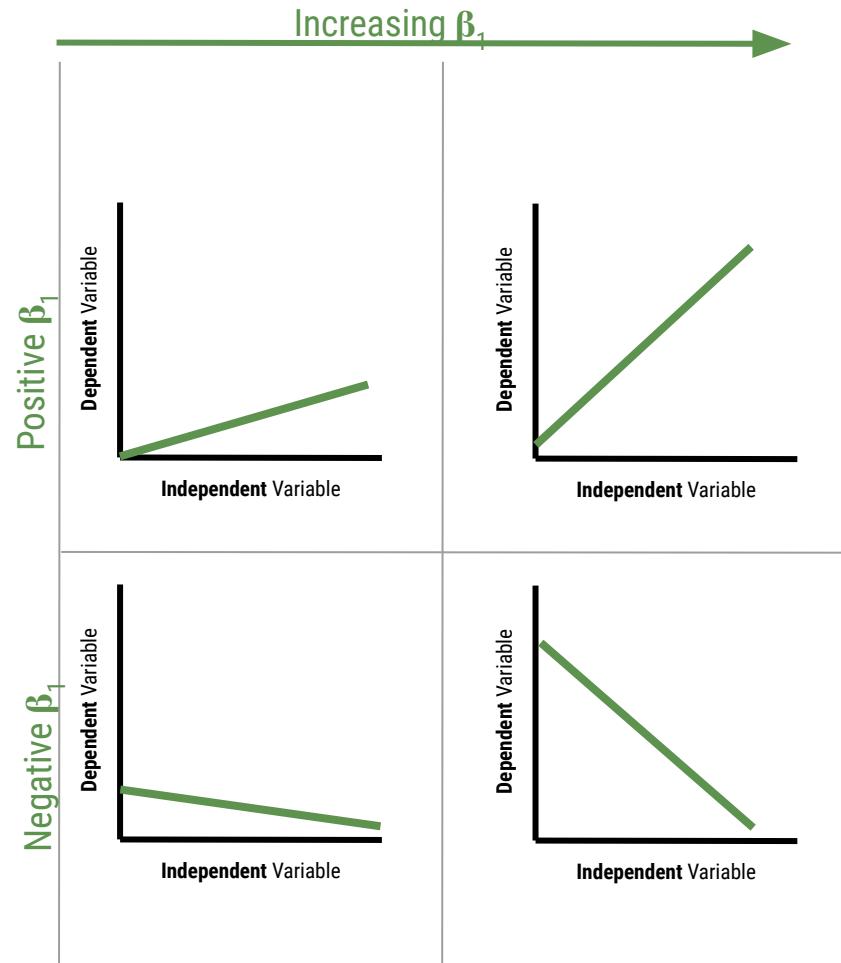
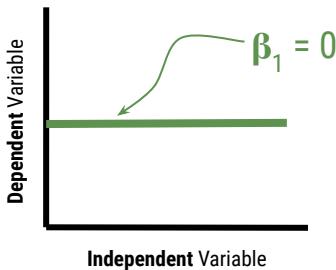




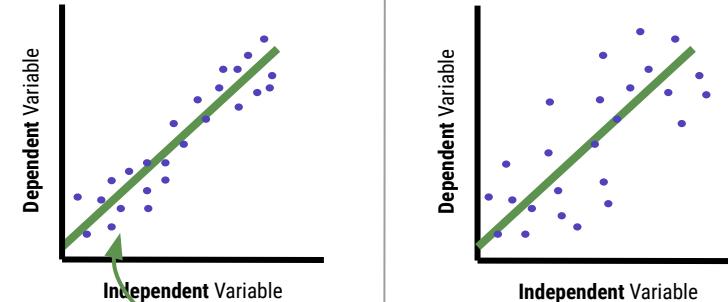
Effect size (β_1) can
be estimated using
the slope of the line



Effect size (β_1) can
be estimated using
the slope of the line



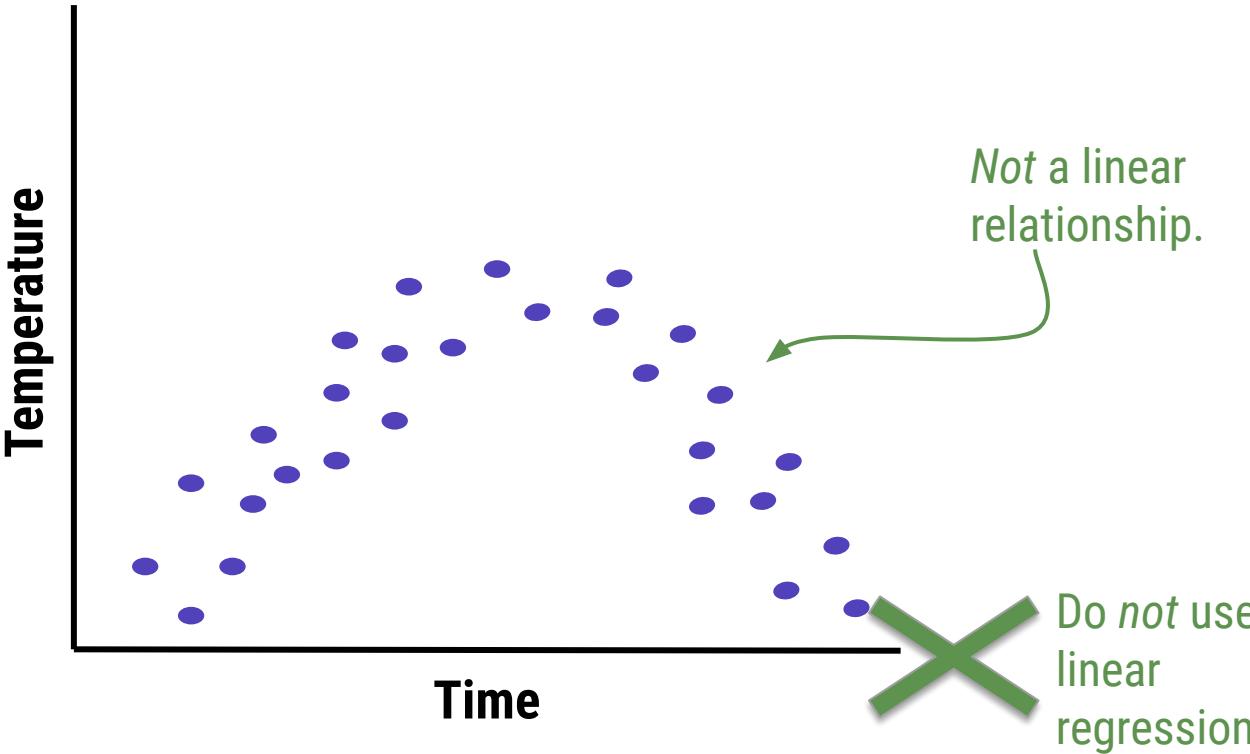
increasing standard error (SE) →



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

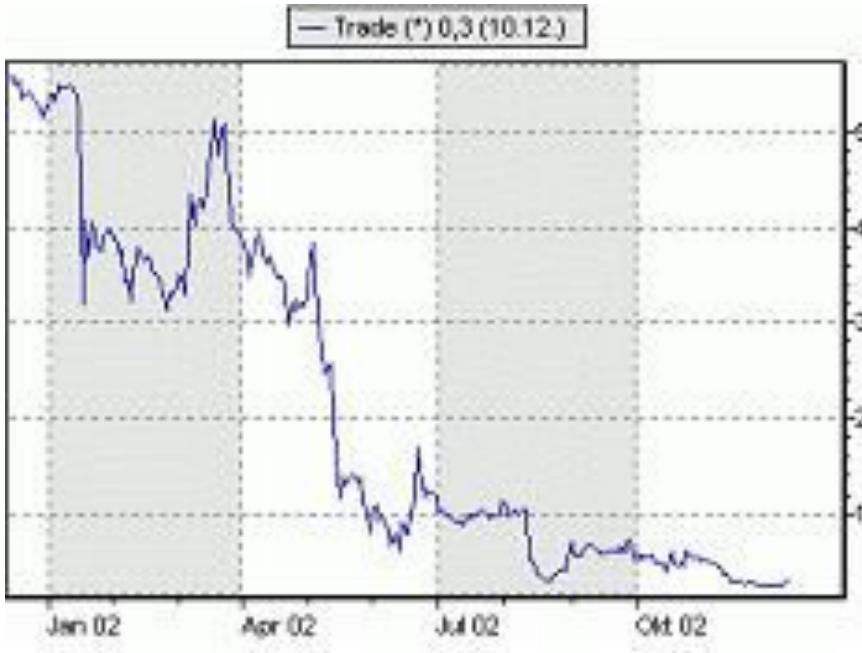
Assumptions of linear regression

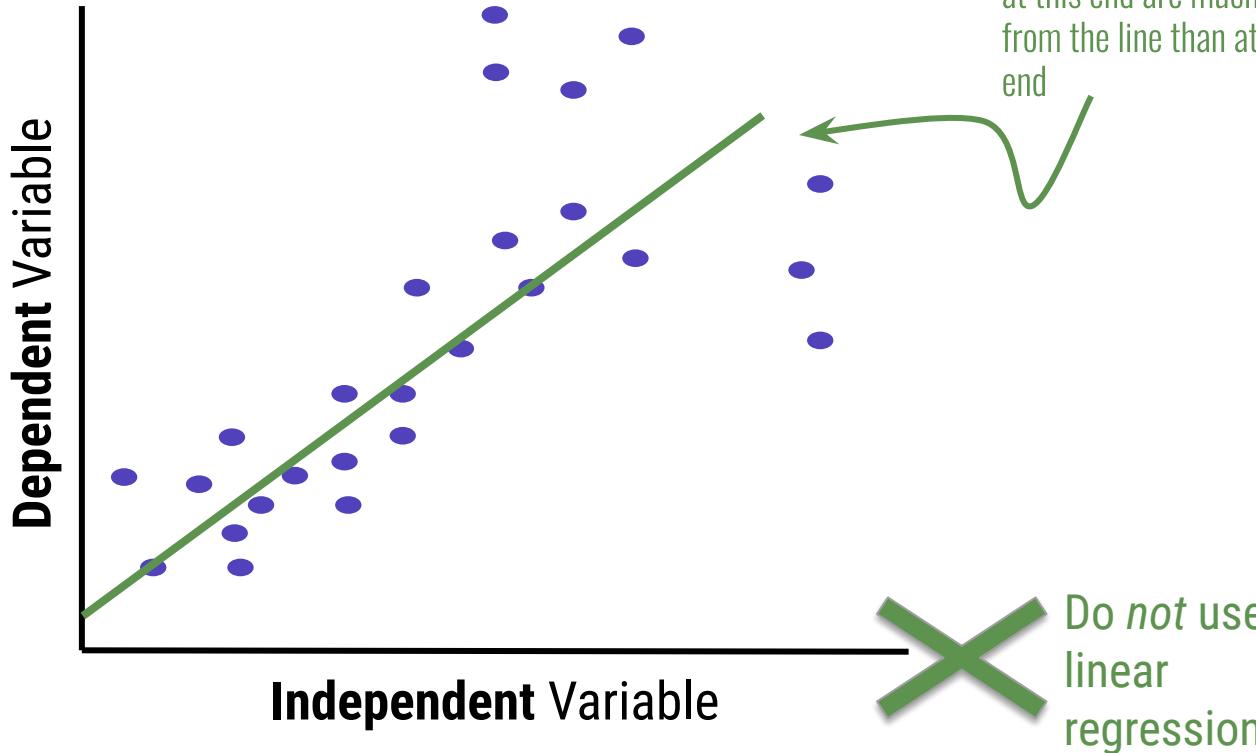
1. Linear relationship
2. No multicollinearity
3. No auto-correlation
4. Homoscedasticity



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)

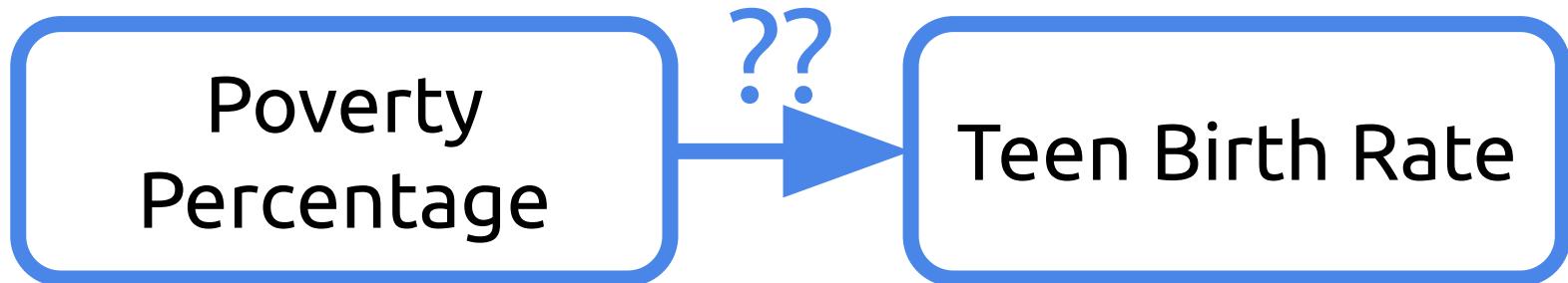




Not homoscedastic: points at this end are much further from the line than at the other end

Do not use linear regression

Does Poverty Percentage
affect Teen Birth Rate?



Null Hypothesis:

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

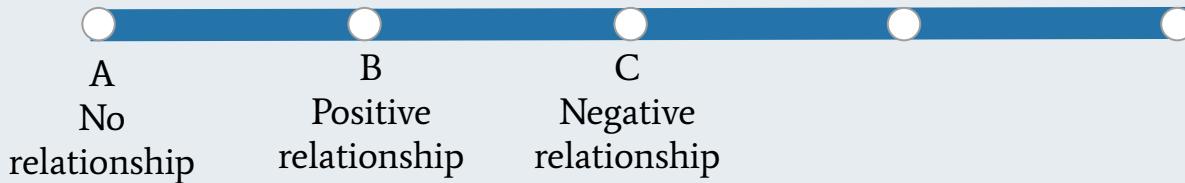
Alternative Hypothesis:

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



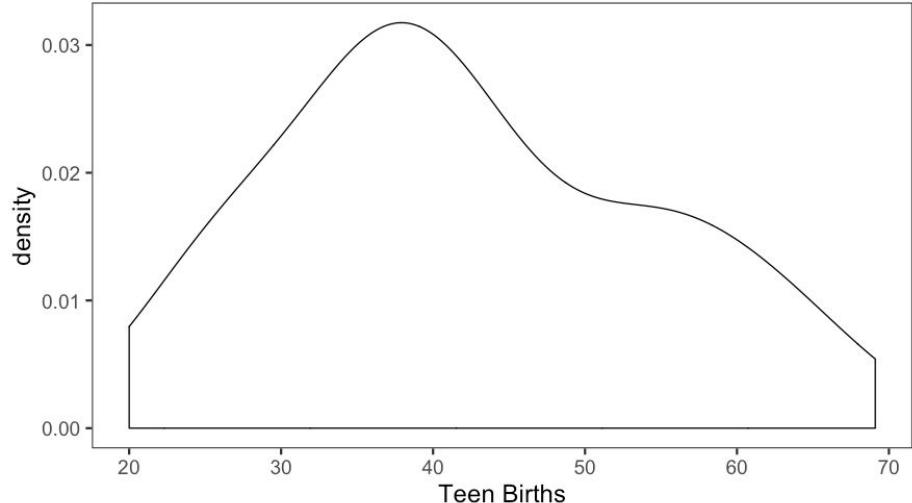
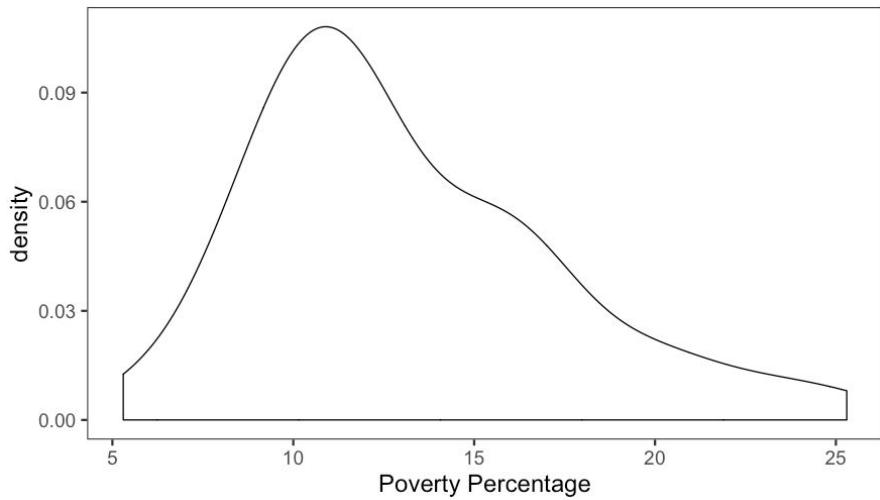
What is the relationship between Poverty Percentage & Teen Birth Rate?

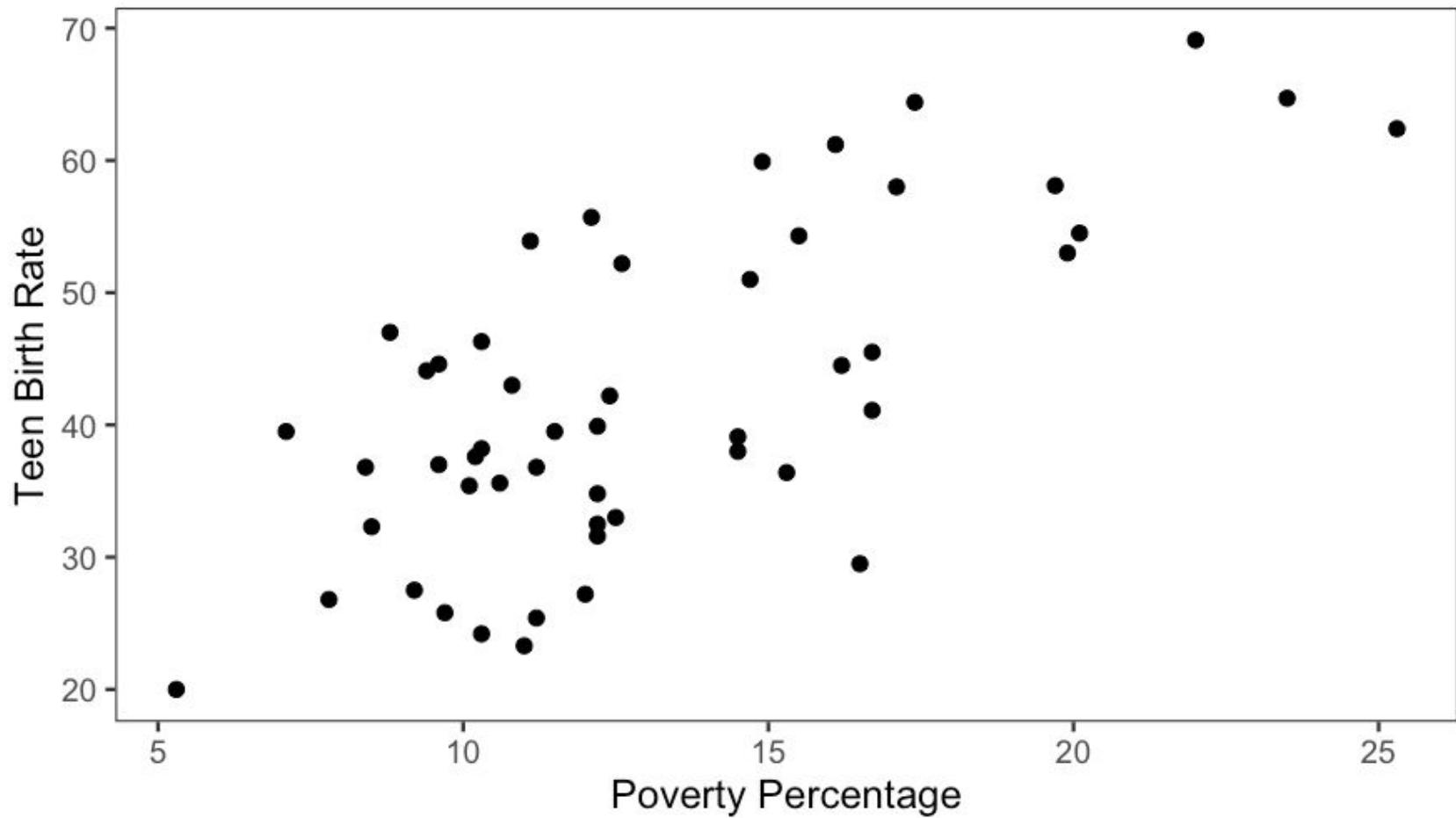
What's your hypothesis?



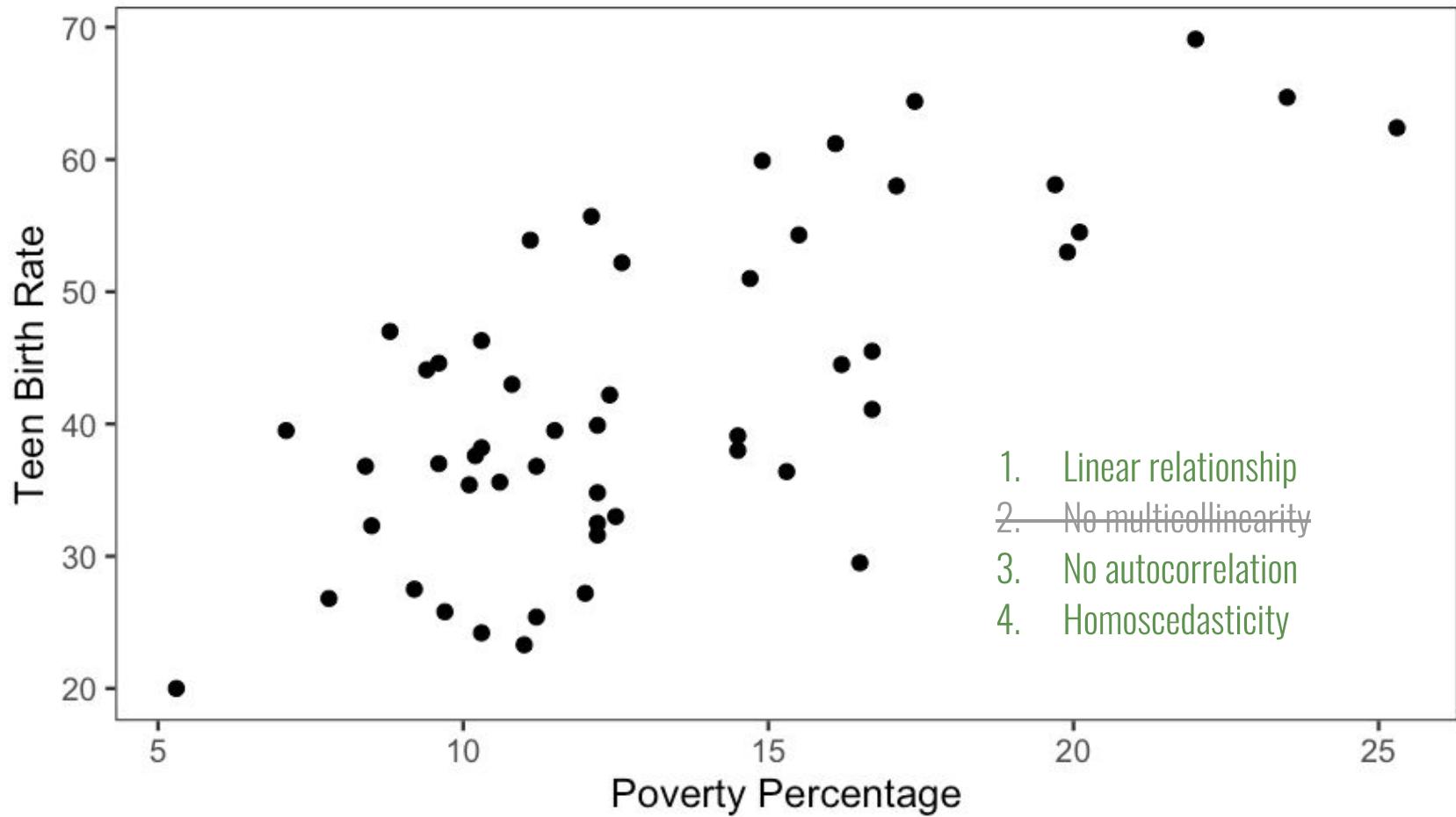
	Location	PovPct	Brth15to17	Brth18to19	ViolCrime	TeenBrth
1	Alabama	20.1	31.5	88.7	11.2	54.5
2	Alaska	7.1	18.9	73.7	9.1	39.5
3	Arizona	16.1	35.0	102.5	10.4	61.2
4	Arkansas	14.9	31.6	101.7	10.4	59.9
5	California	16.7	22.6	69.1	11.2	41.1
6	Colorado	8.8	26.2	79.1	5.8	47.0
7	Connecticut	9.7	14.1	45.1	4.6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

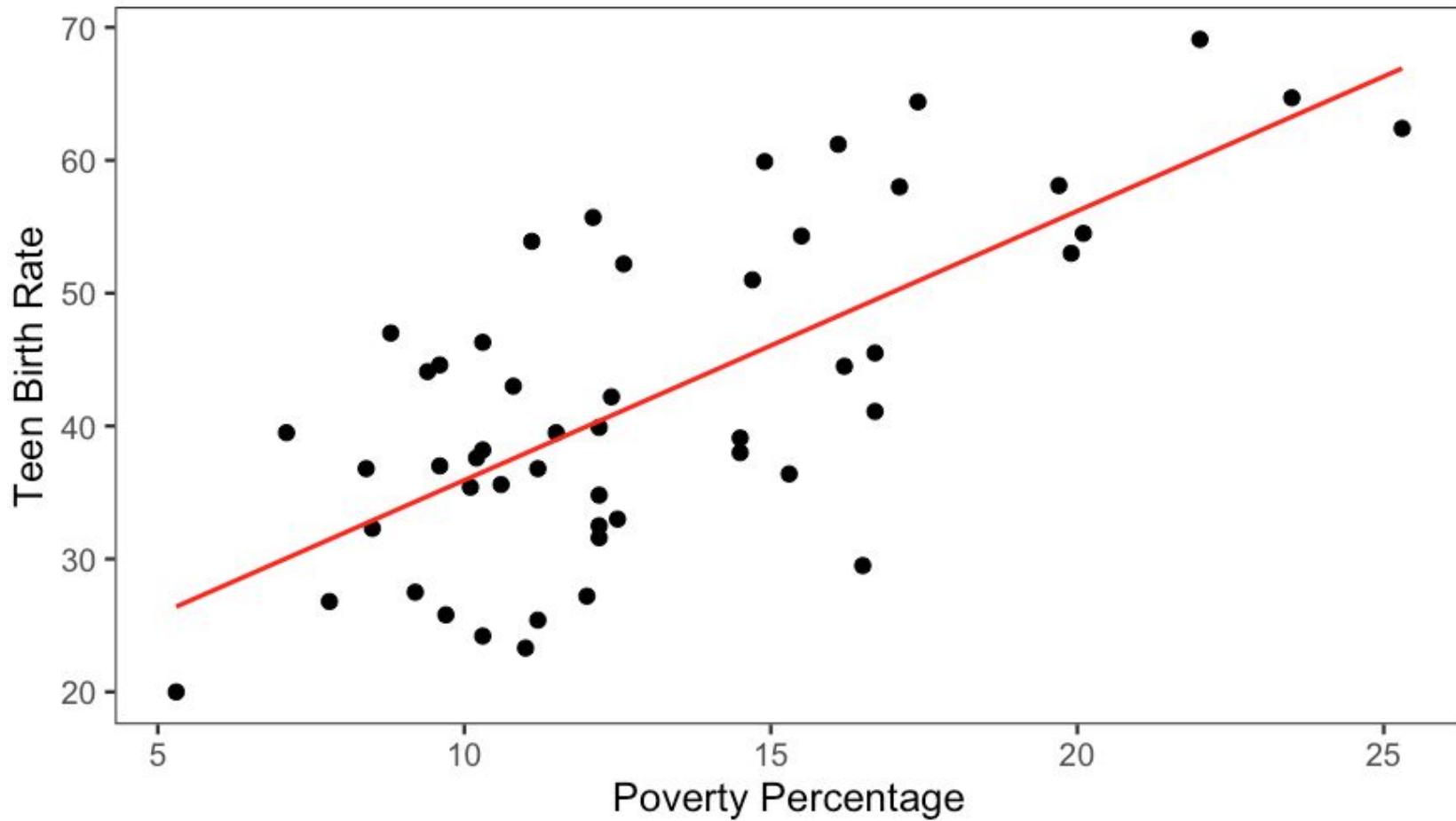
EDA: distributions





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





Teen Birth Rate

60

50

40

30

5

10

15

20

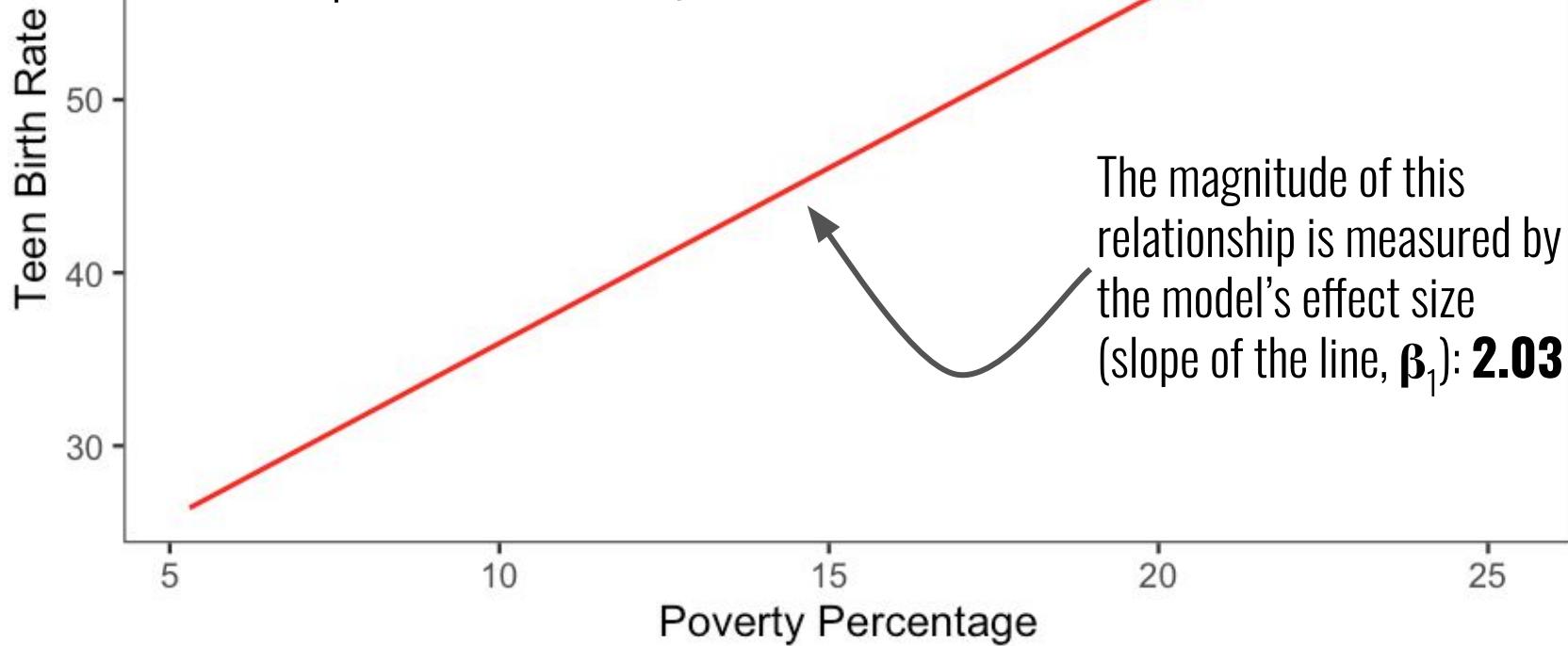
25

Poverty Percentage

The regression line is the model being used to explain the relationship between Poverty Percentage and Birth Rate

The magnitude of this relationship is measured by the model's effect size (slope of the line, β_1): **2.03**





To interpret this effect size... *for every 1% increase in Poverty Percentage*, the birth rate is expected to increase by **2.03**

The magnitude of this relationship is measured by the model's effect size (slope of the line, β_1): **2.03**

Teen Birth Rate

60

50

40

30

5

10

15

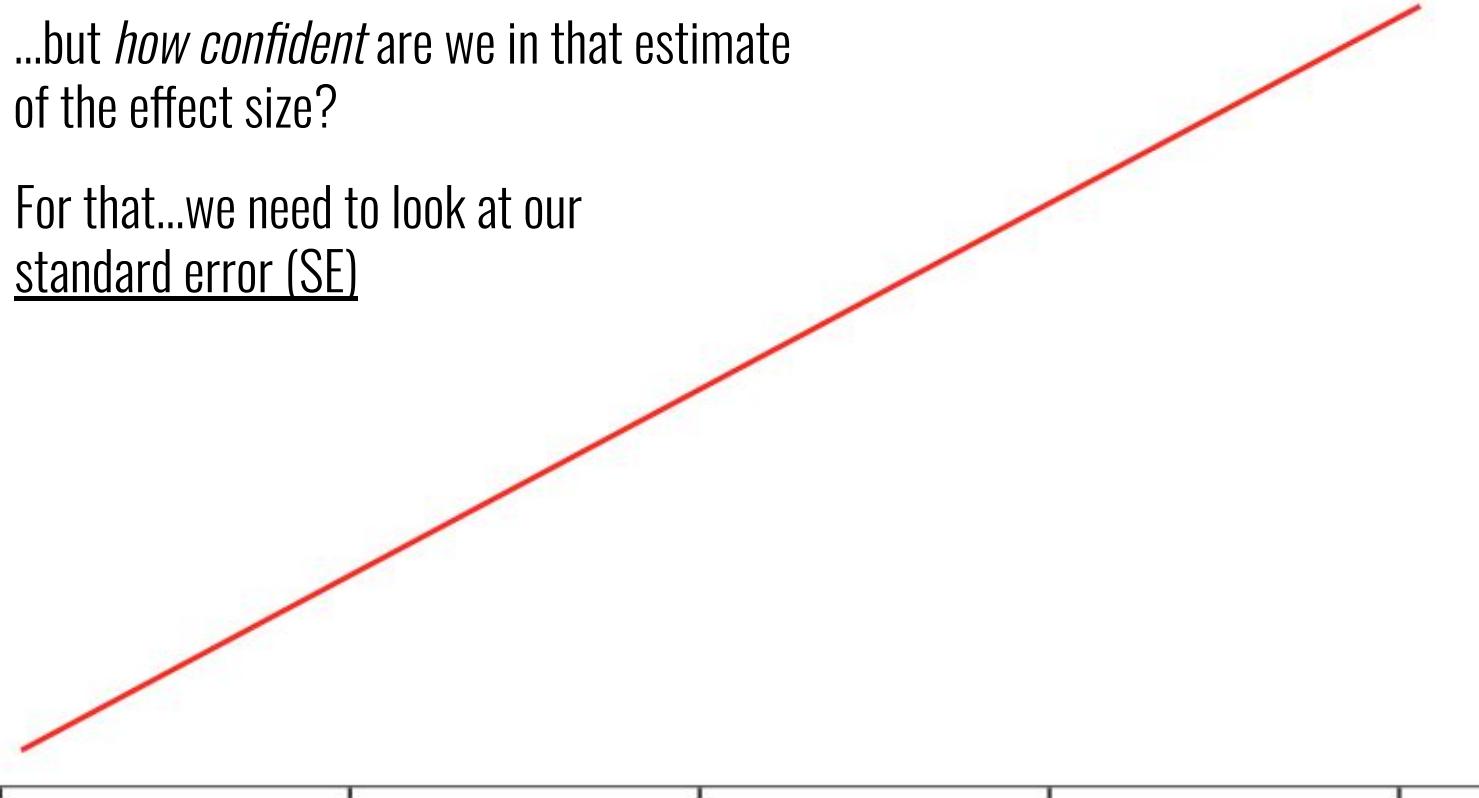
20

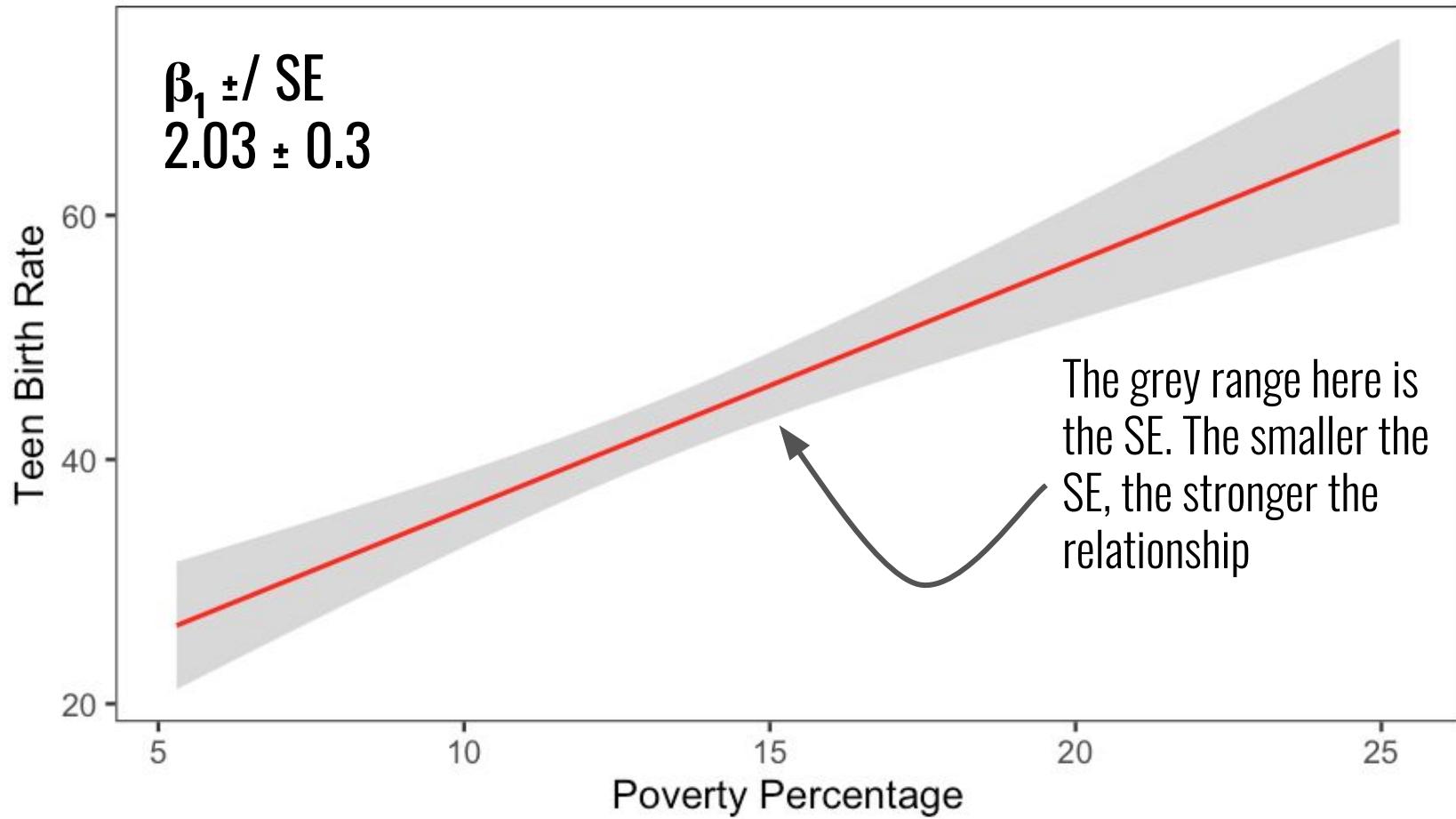
25

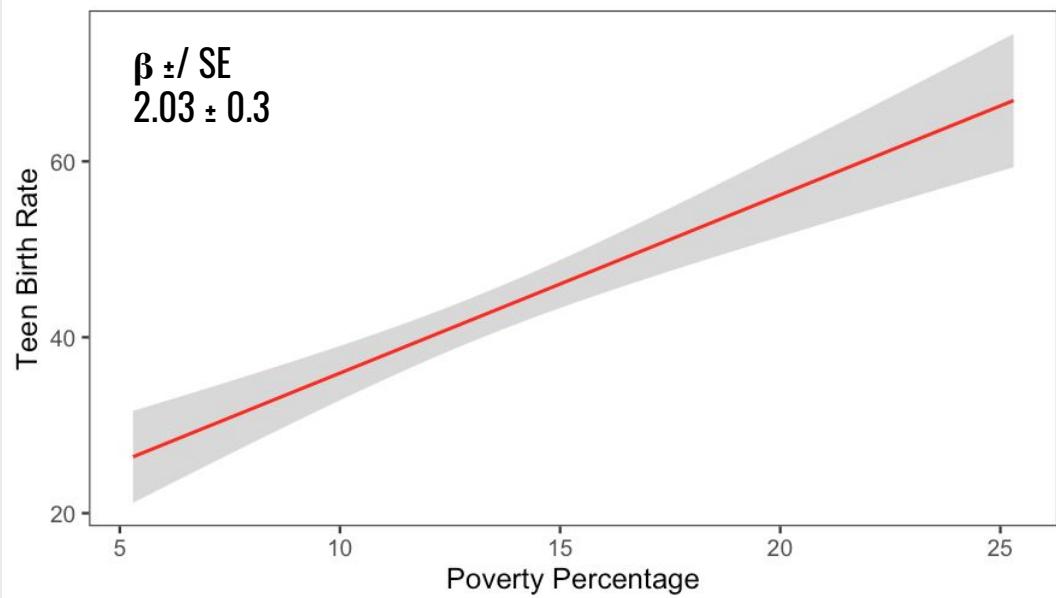
Poverty Percentage

...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 β be?

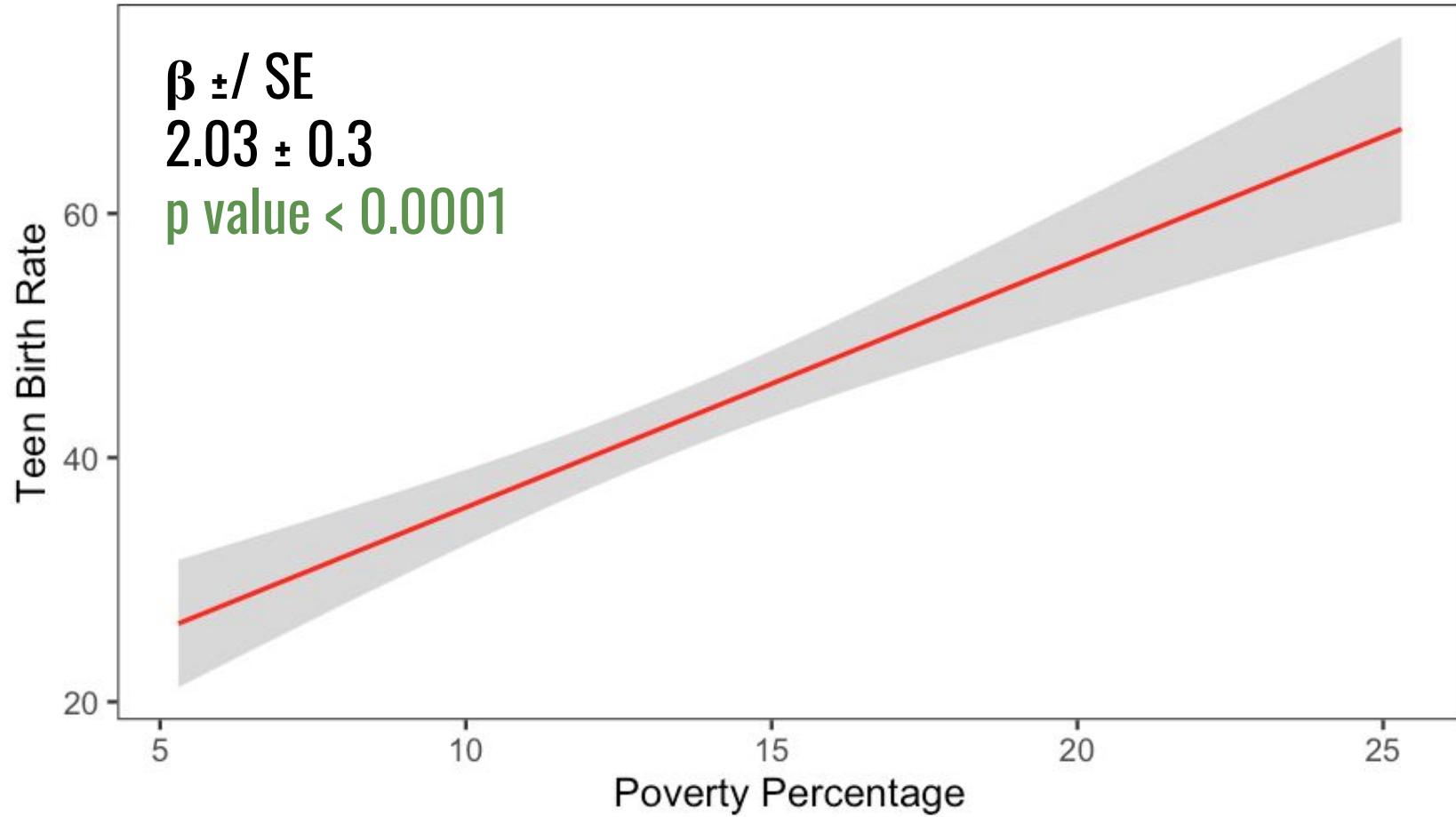


A
 < 2.03

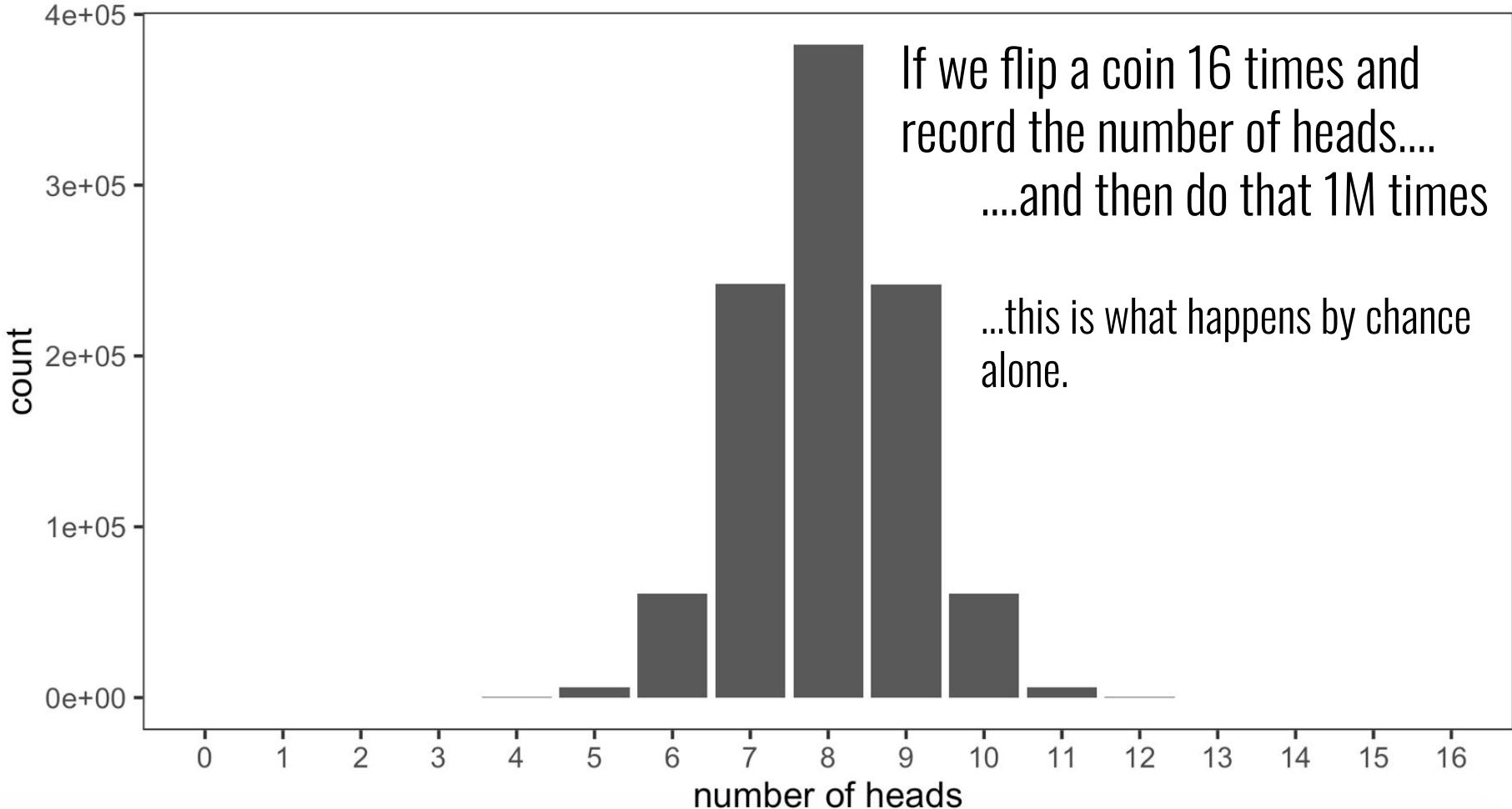


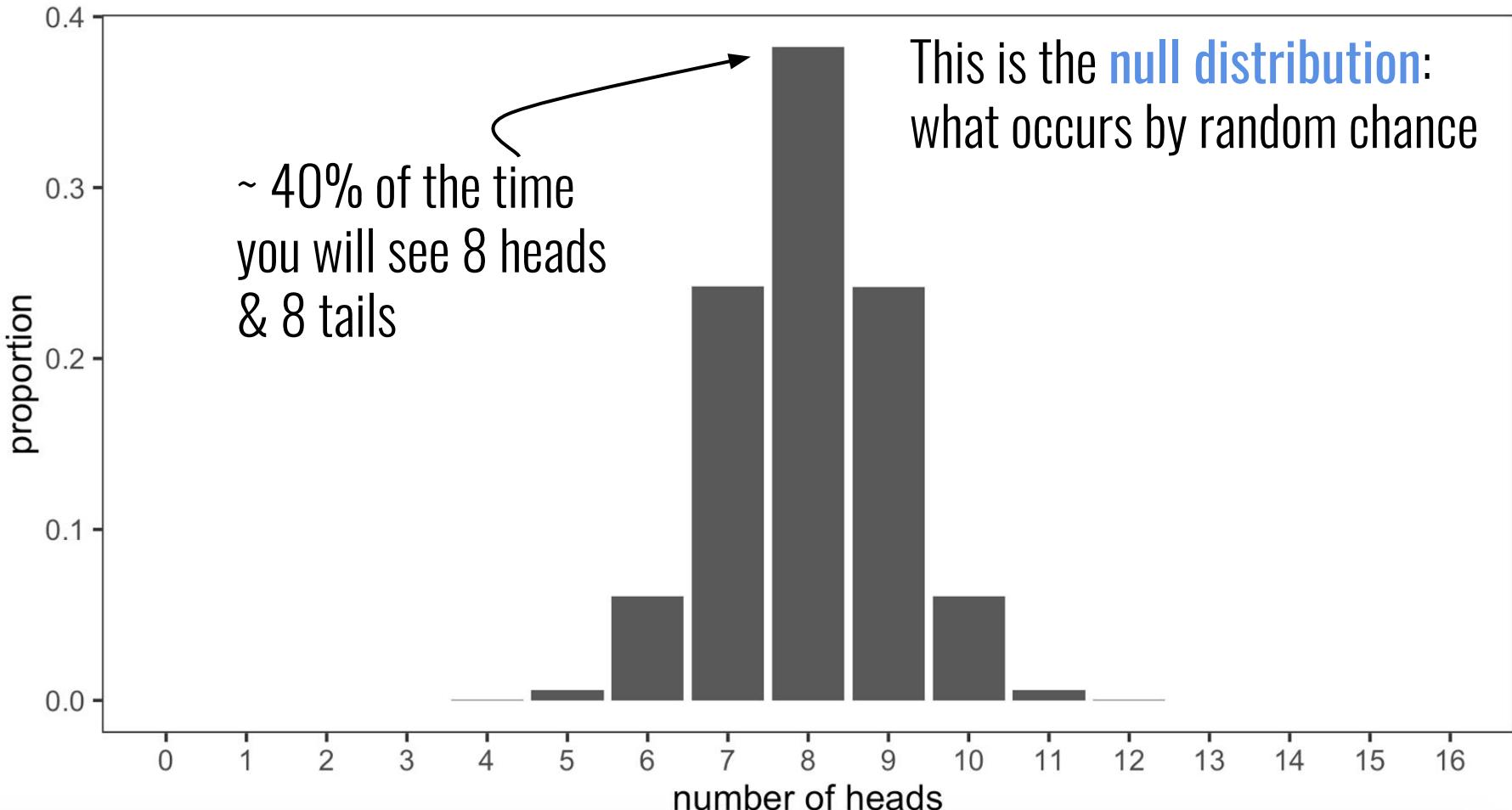
B
 > 2.03

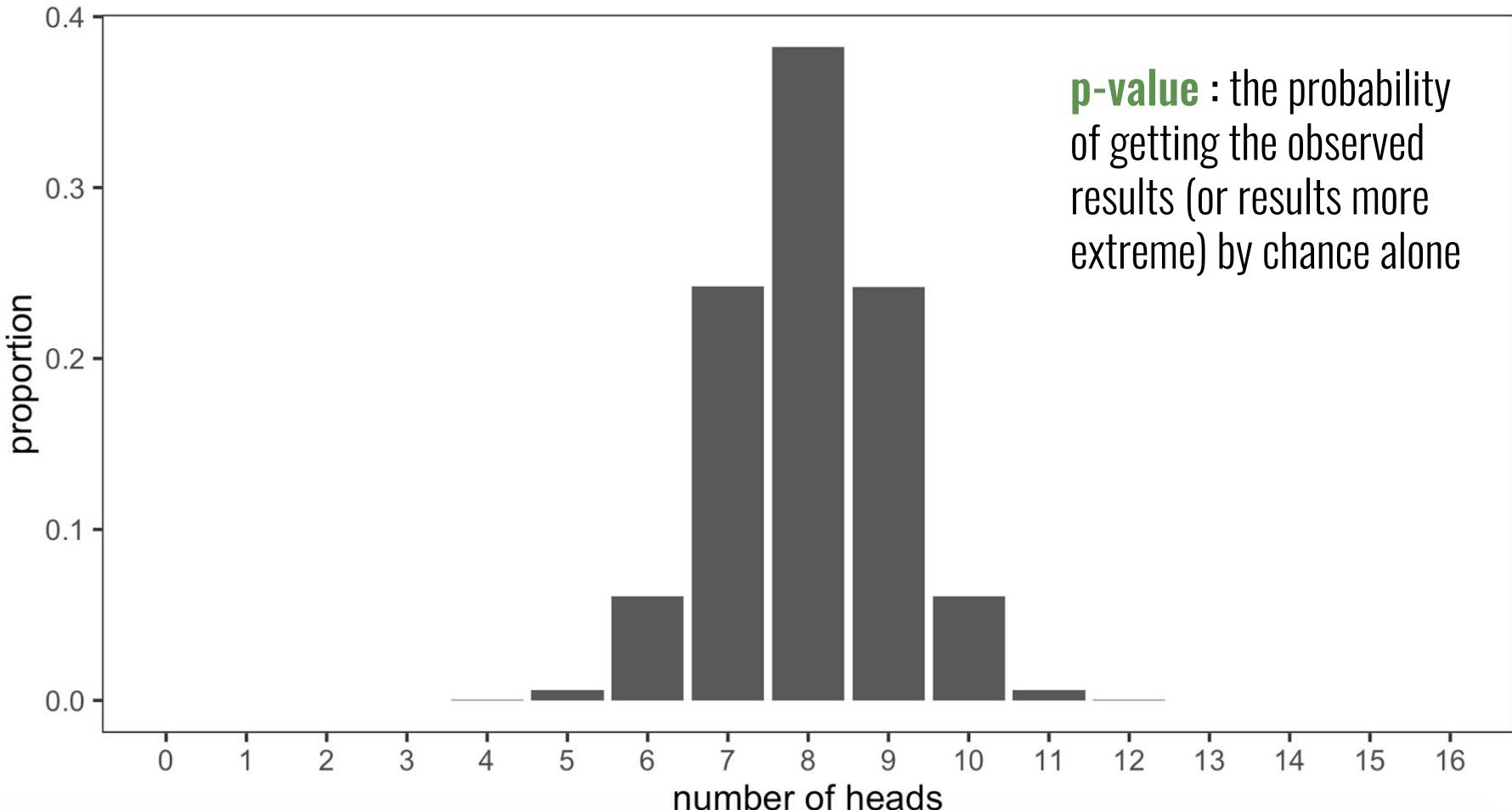




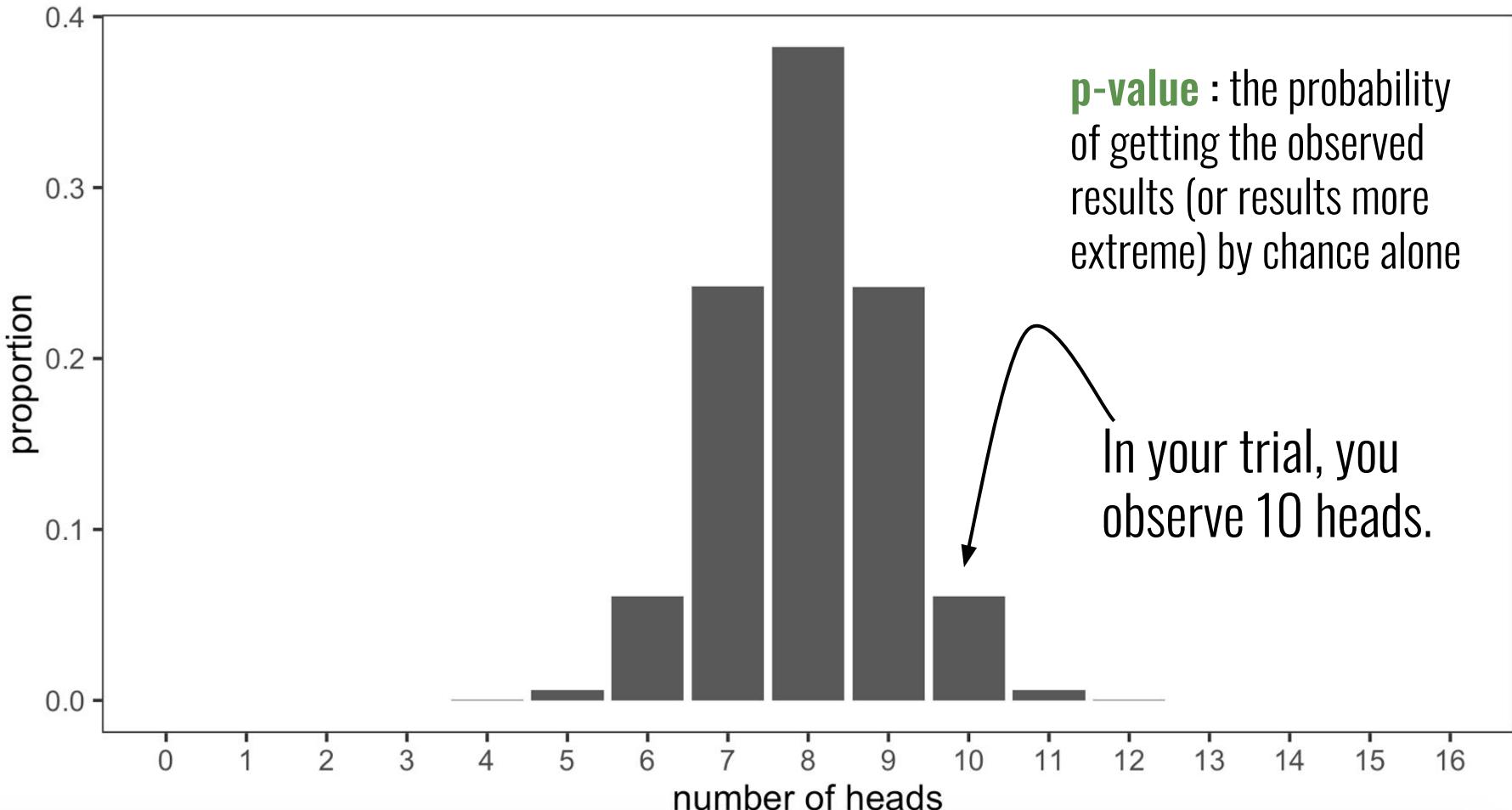
p-value : the probability of getting the observed results (or results more extreme) by chance alone

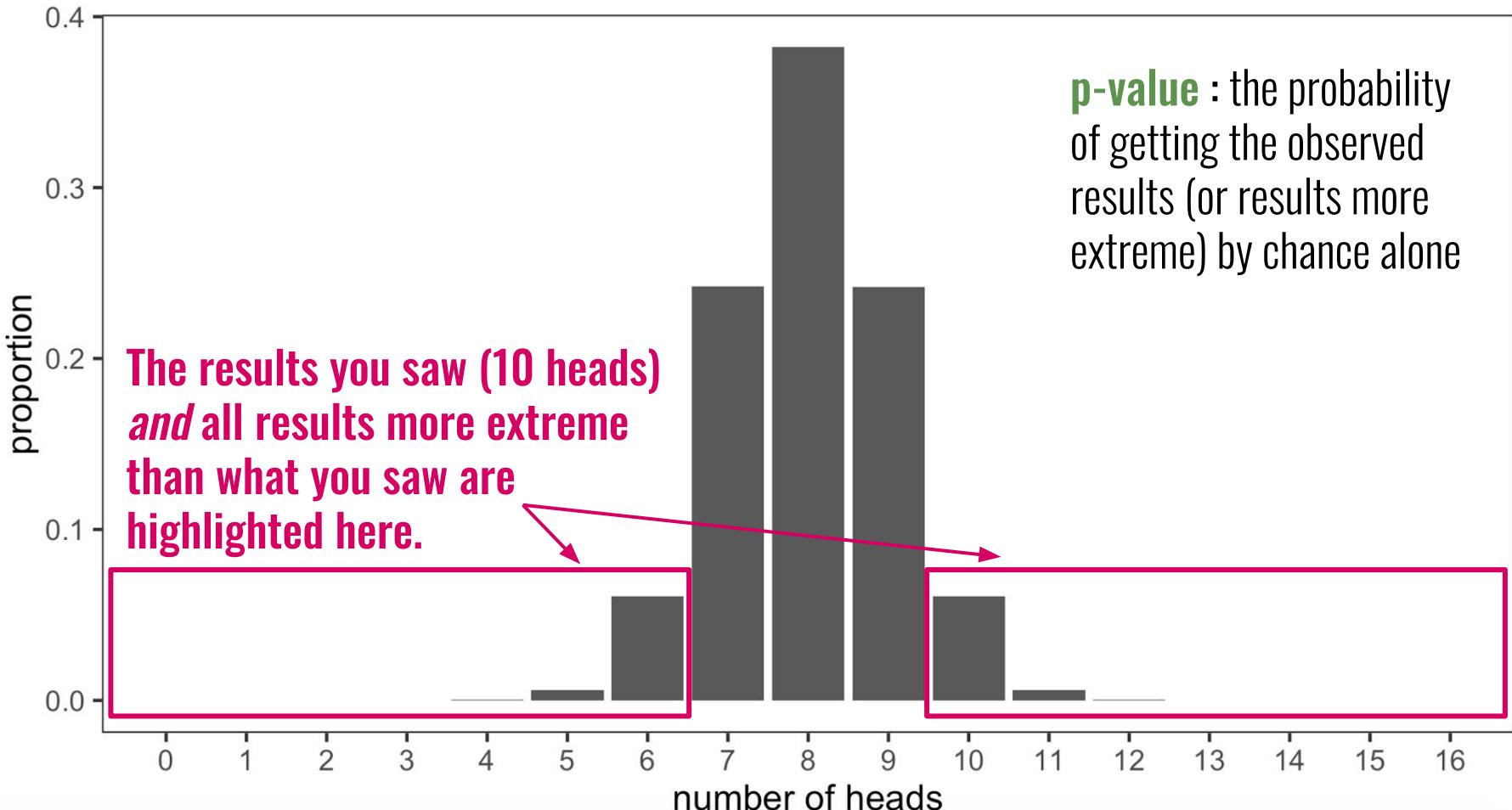


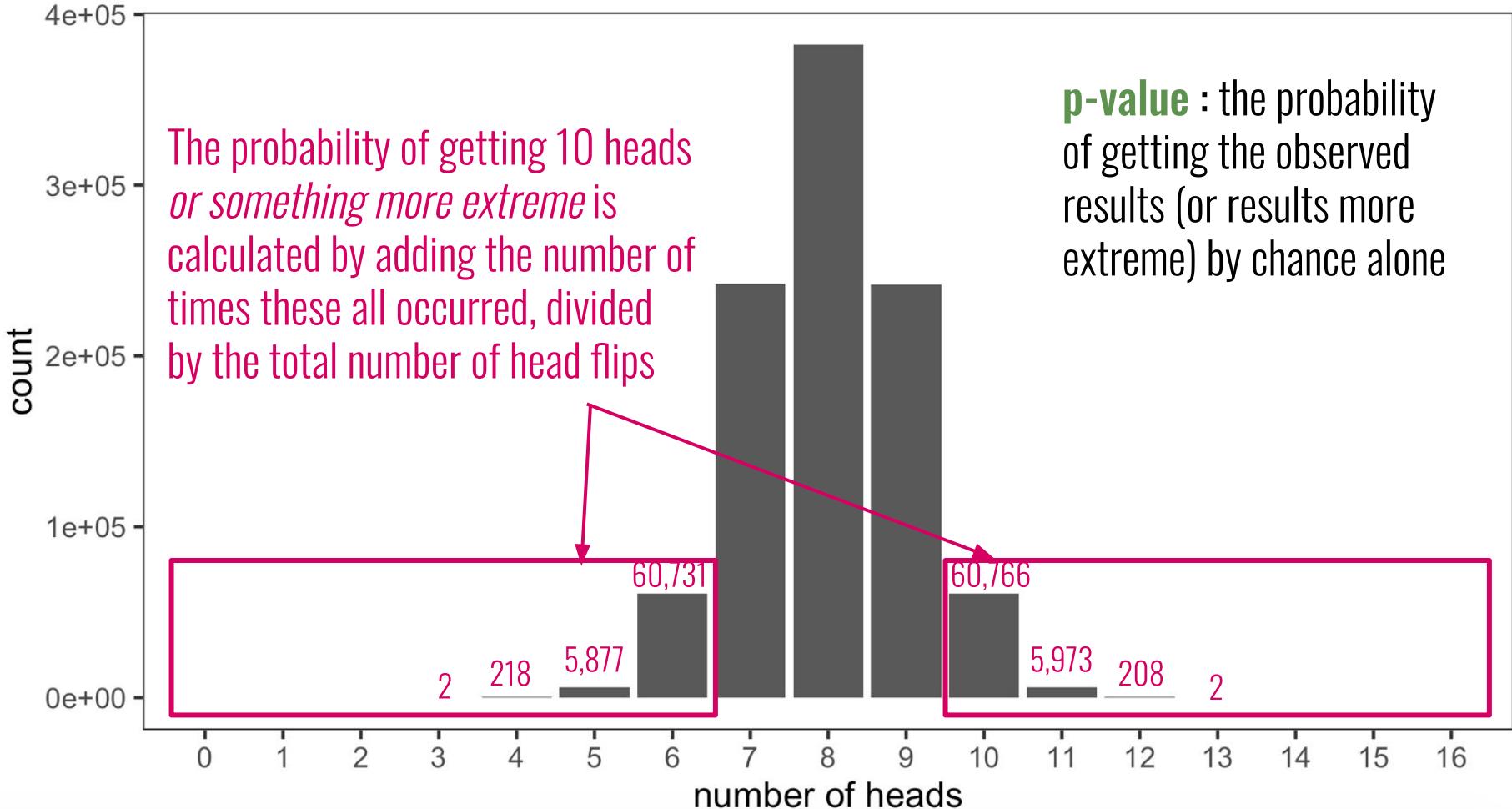


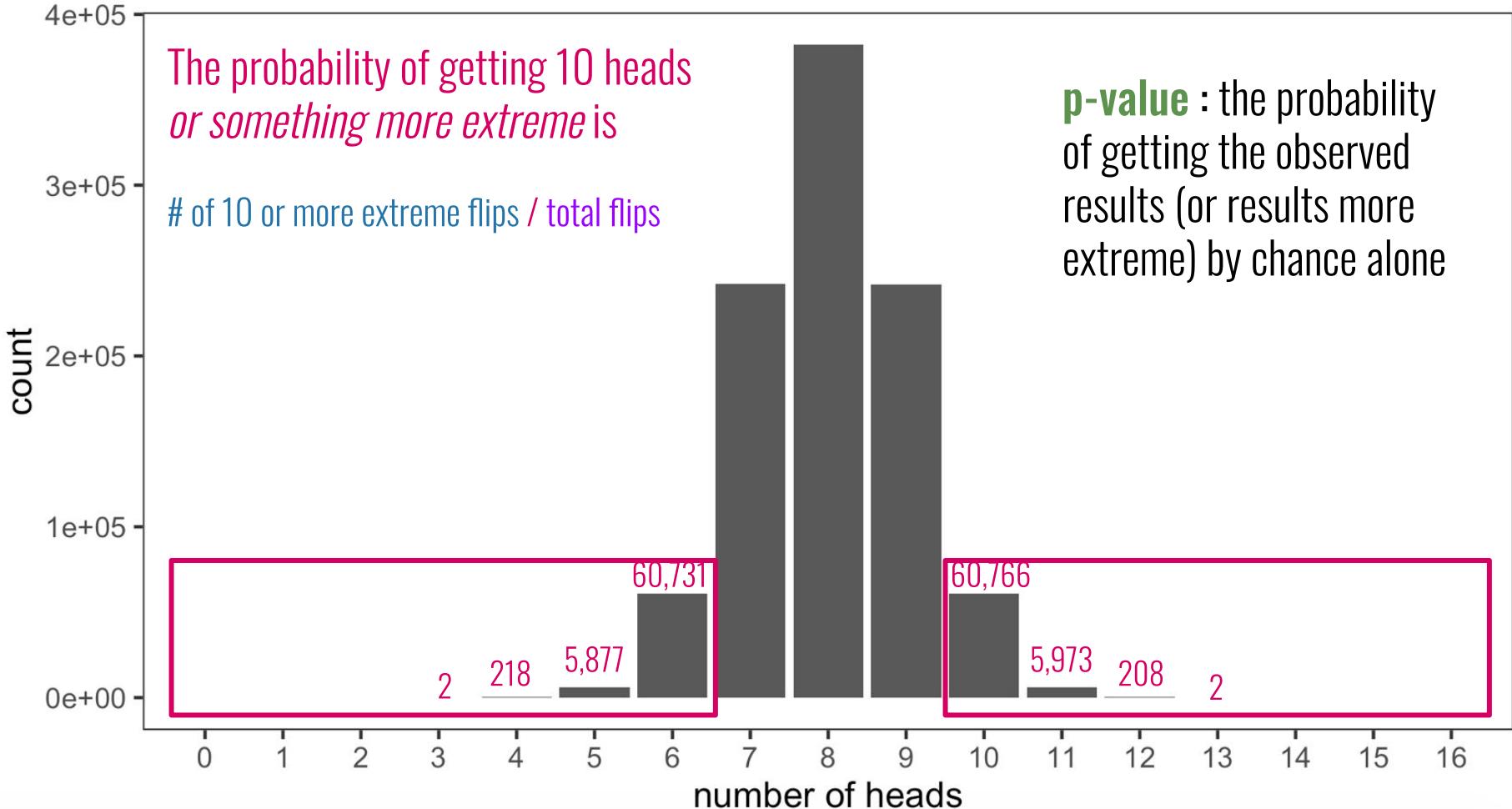


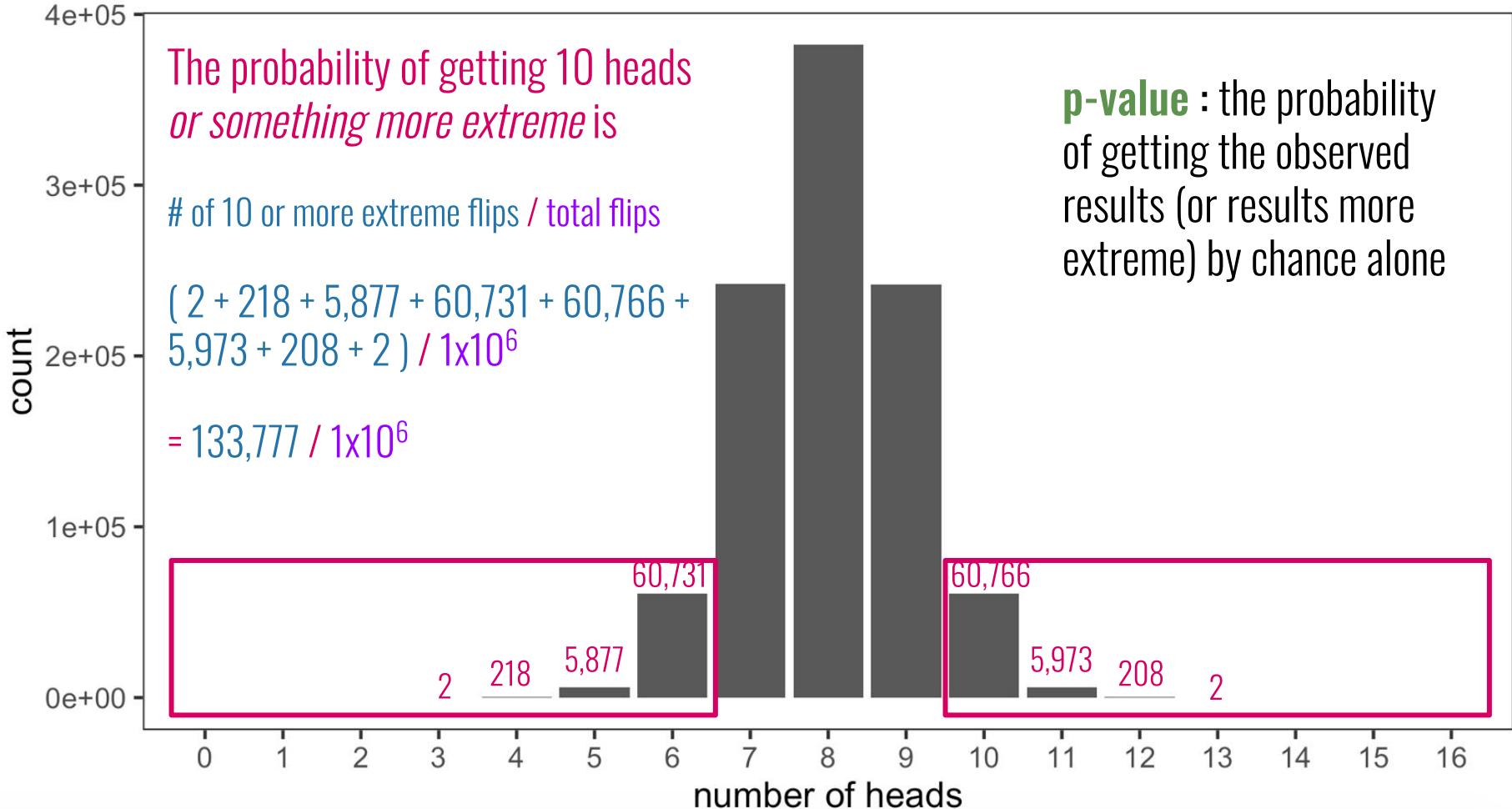
p-value : the probability
of getting the observed
results (or results more
extreme) by chance alone

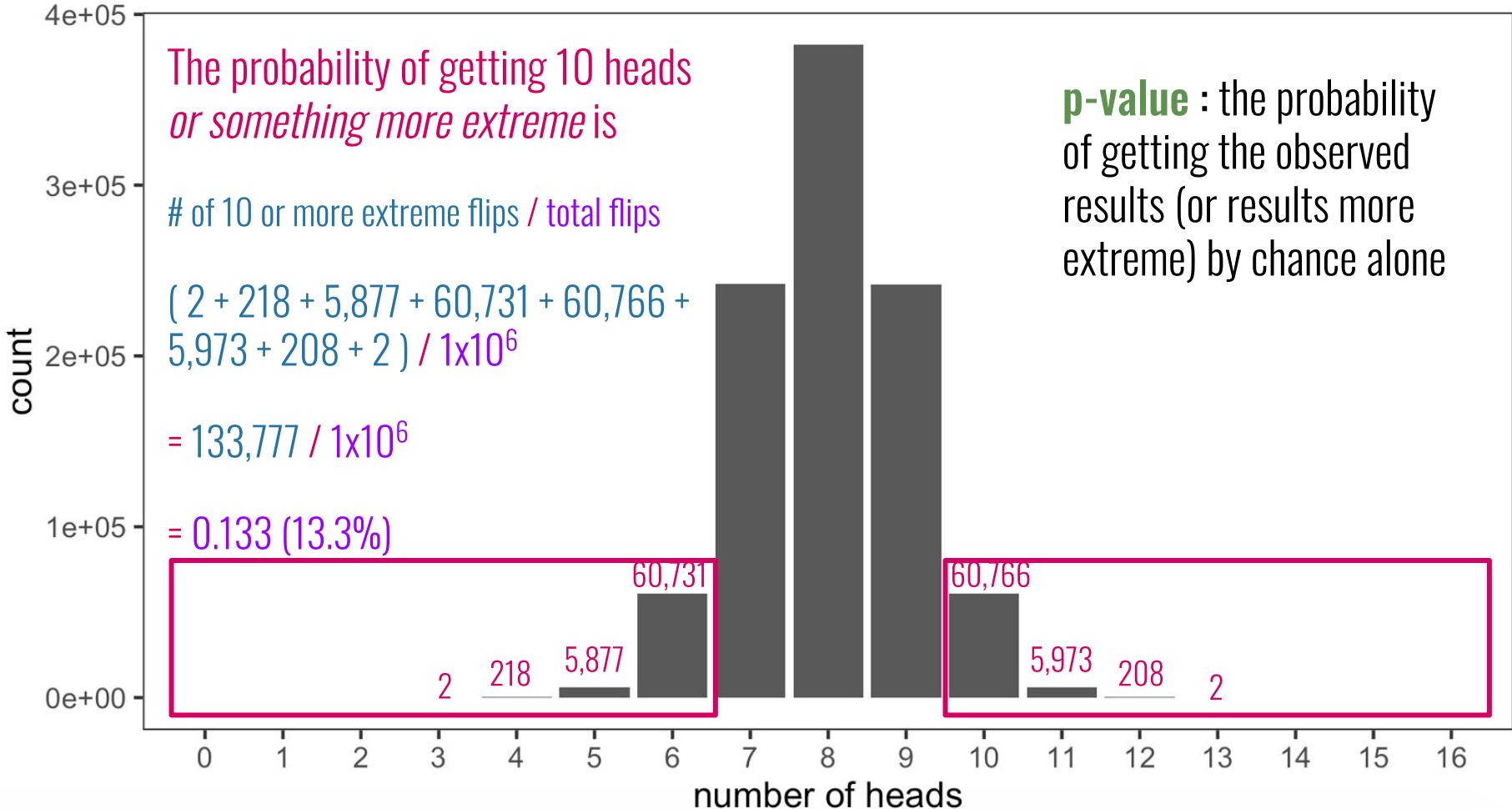


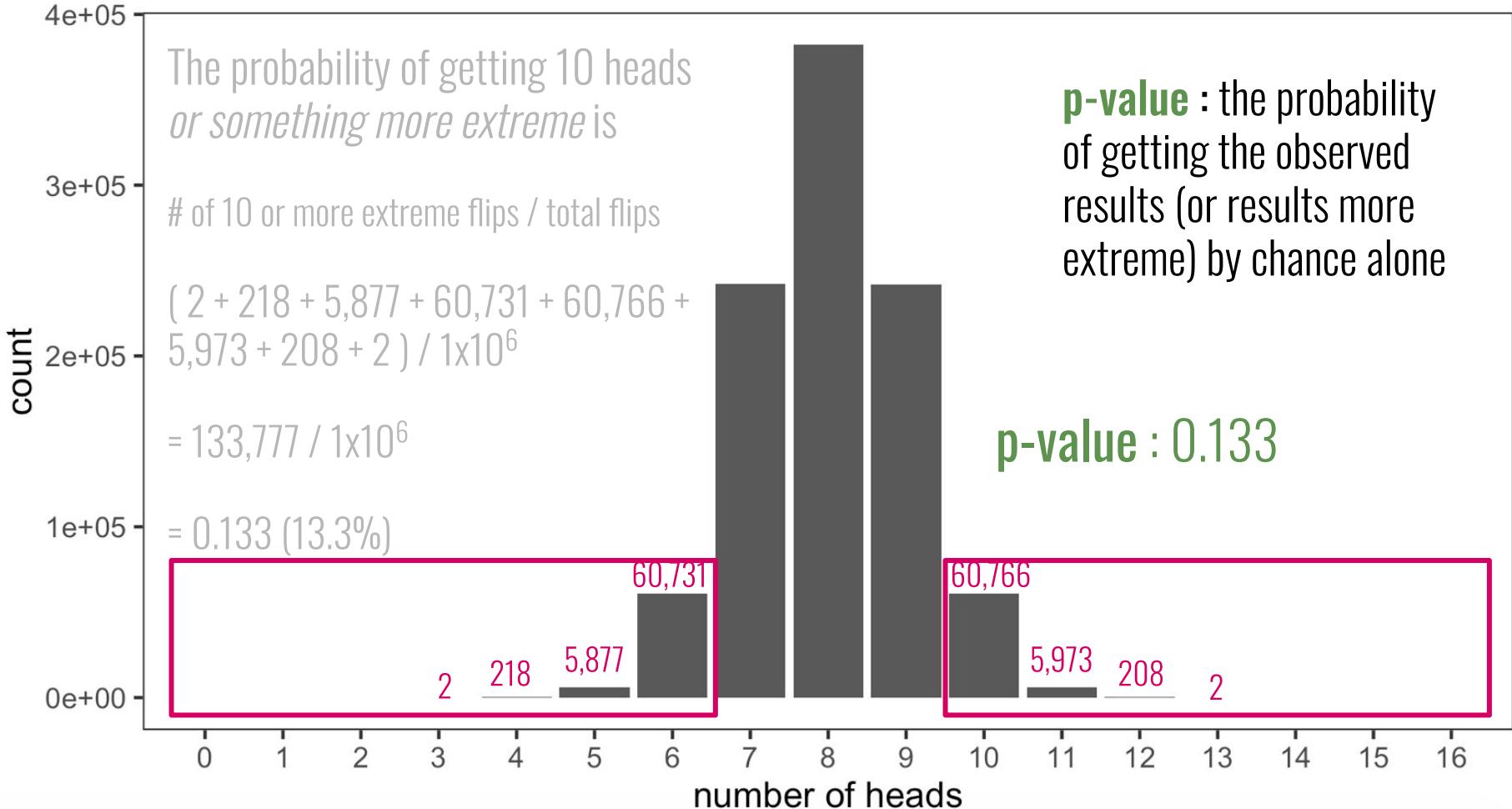


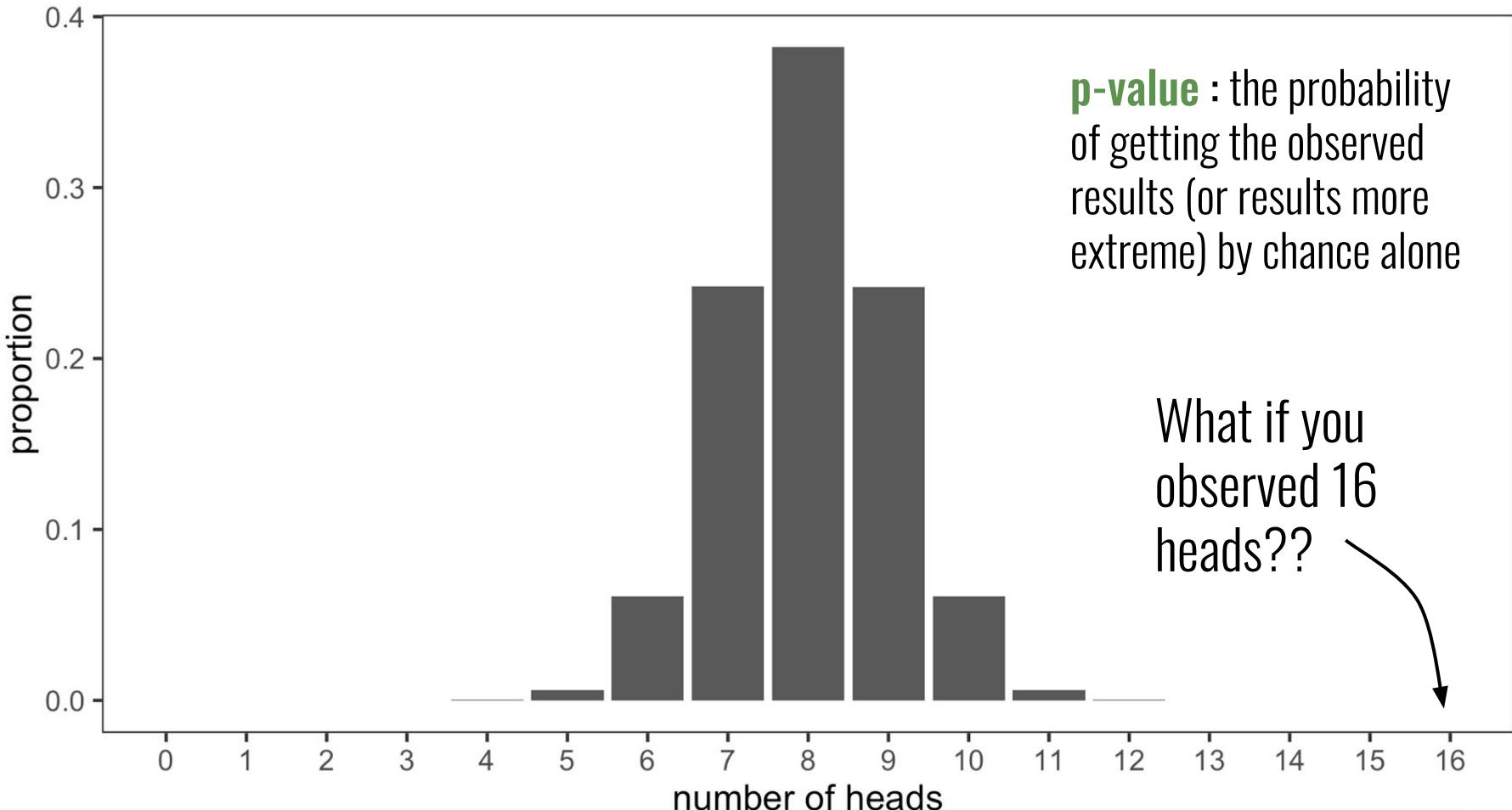


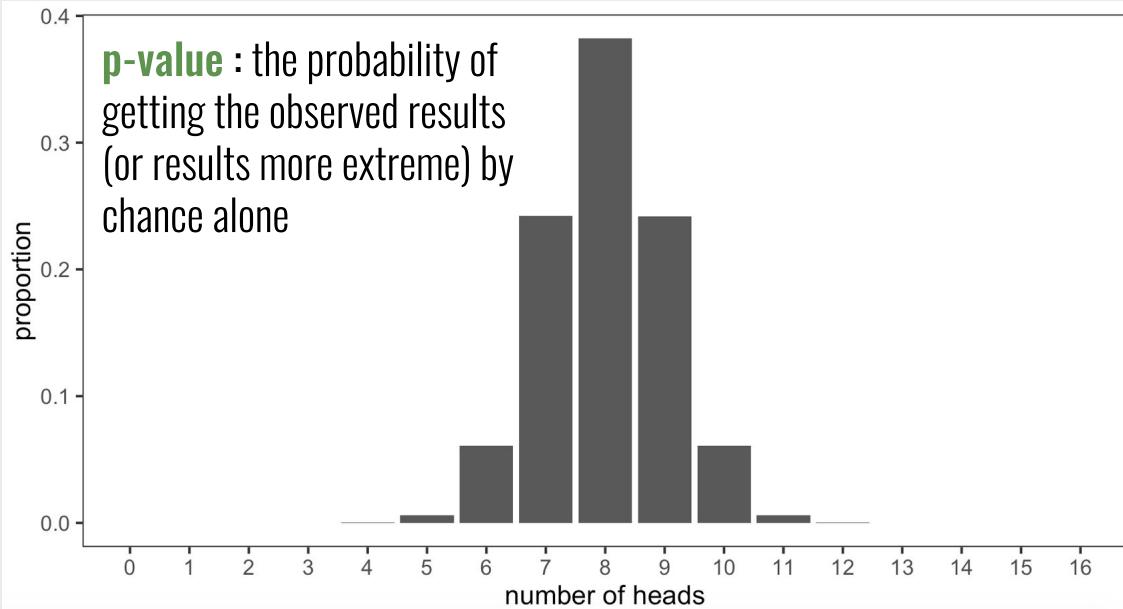












What would be the p-value of you flipping 16 heads?



A

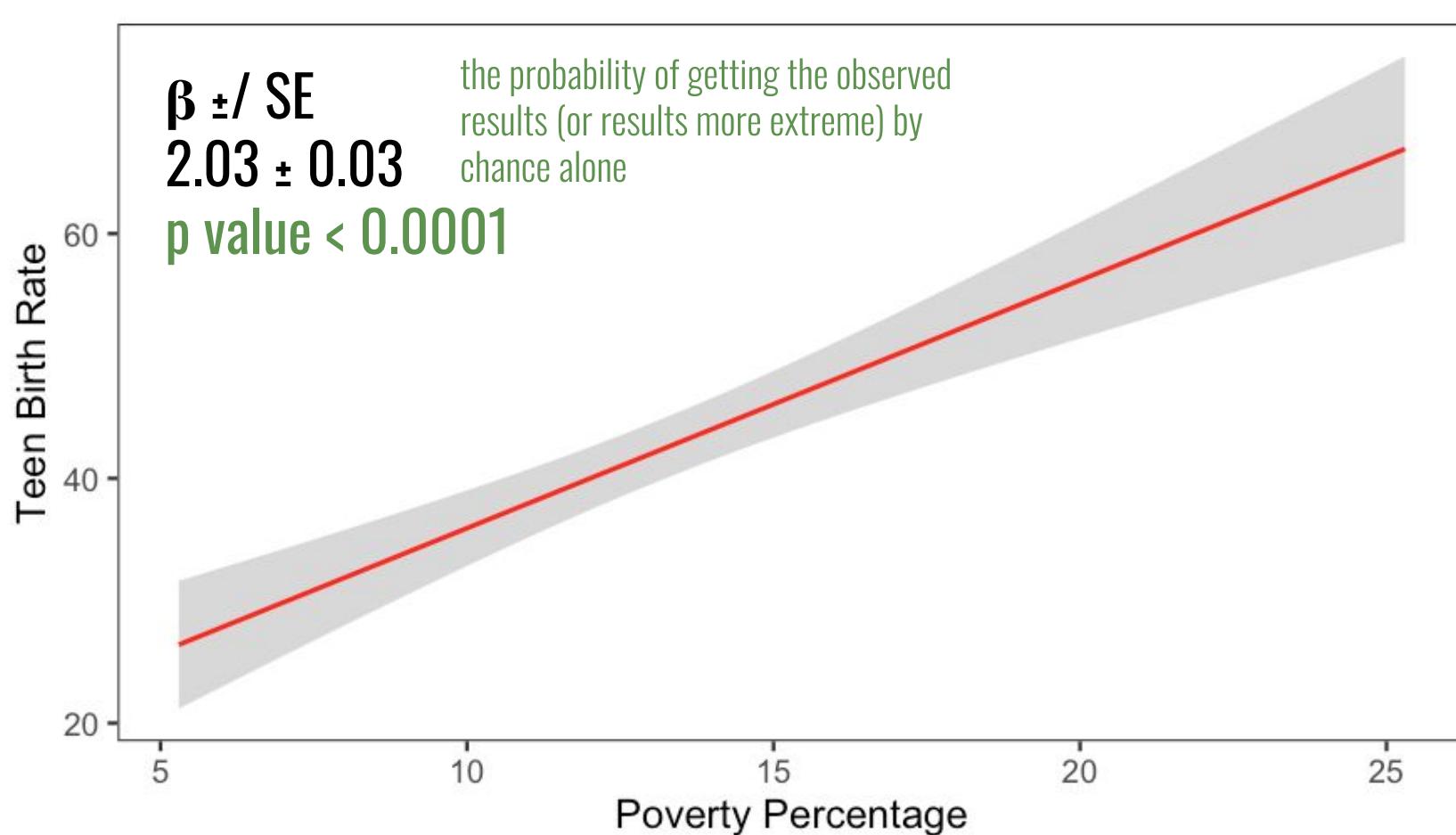
< 0.13



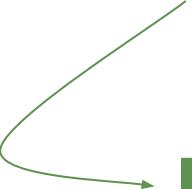
B

> 0.13





Takes into account the effect size (β_1) and the SE



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

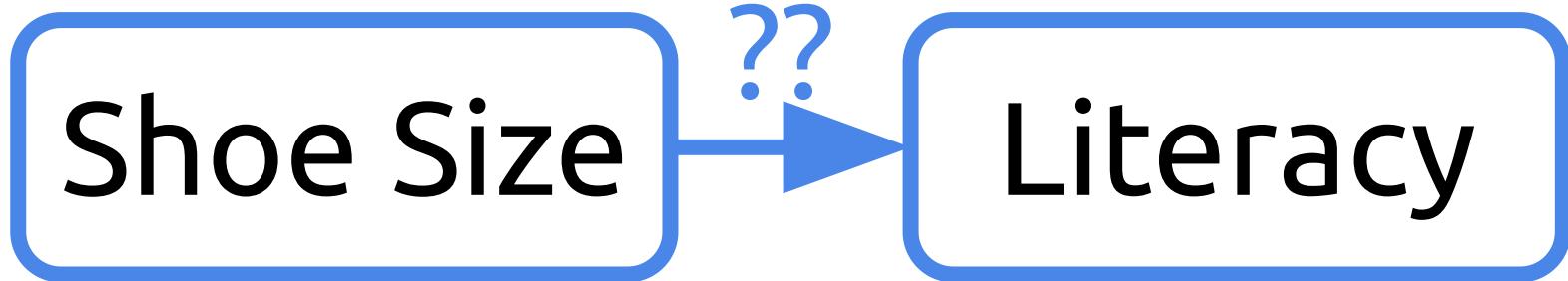
Confounding





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

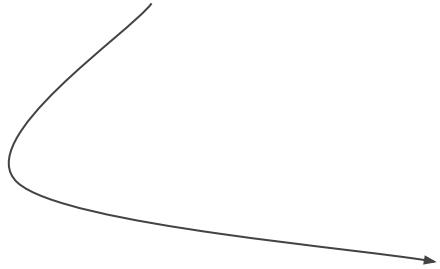
Confounding



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



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...

	Vertebroplasty N = 200	Conservative care N = 200
Age, y, mean \pm SD	78.2 ± 4.1	79.0 ± 5.2
Weight, kg, mean \pm 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