

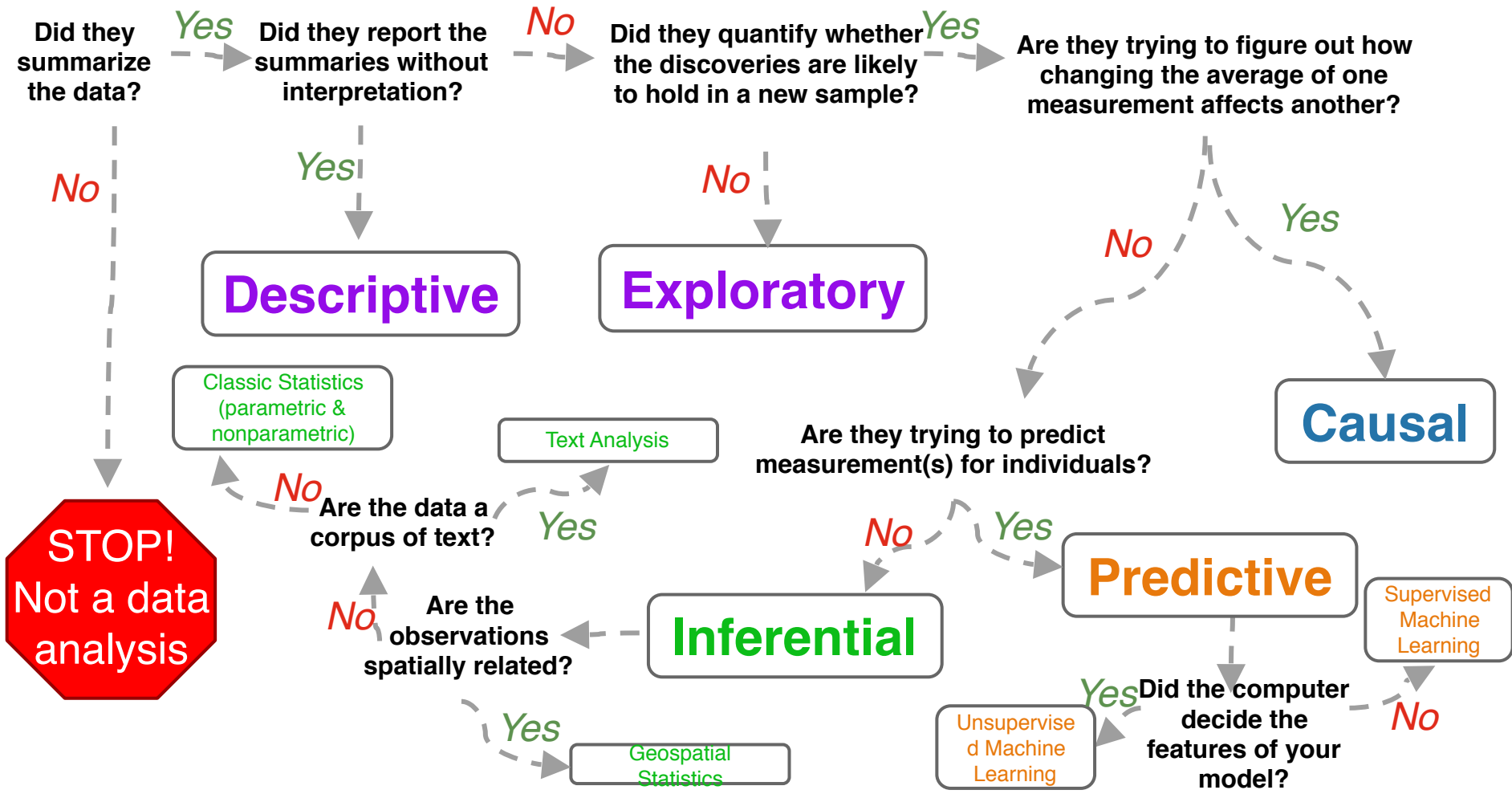
Inferential Analysis

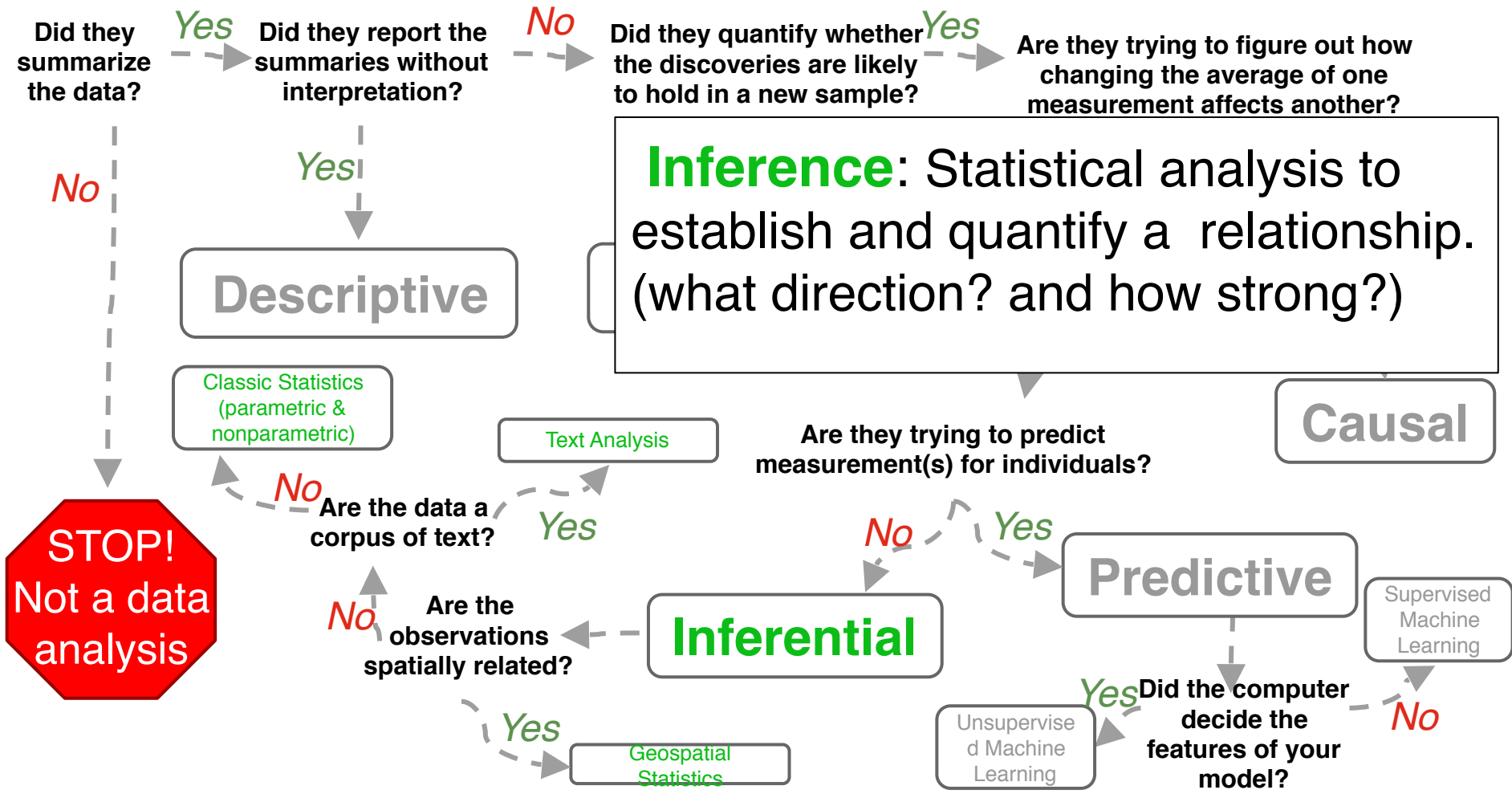
C. Alex Simpkins Jr., Ph.D
UC San Diego, RDPRobotics LLC



Department of Cognitive Science
rdrobotics@gmail.com
csimpkinsjr@ucsd.edu

Lectures : http://casimpkinsjr.radiantdolphinpress.com/pages/cogs108_ss1_23/





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

STOP!
Not a data analysis

Descriptive

Classic Statistics
(parametric & nonparametric)

Text Analysis

Inferential

Geospatial Statistics

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

Predictive

Unsupervised Machine Learning

Supervised Machine Learning

Causal

Did the computer decide the features of your model?

- **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

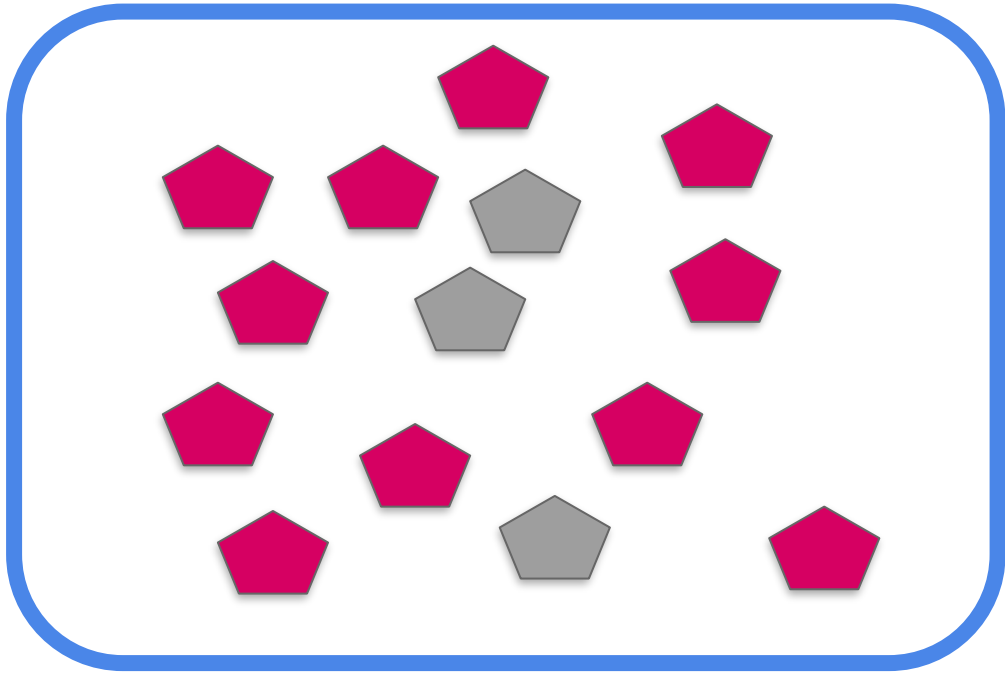
Establishing & Stating Your Null and Alternative Hypotheses Helps Guide Your Analysis

Null Hypothesis:

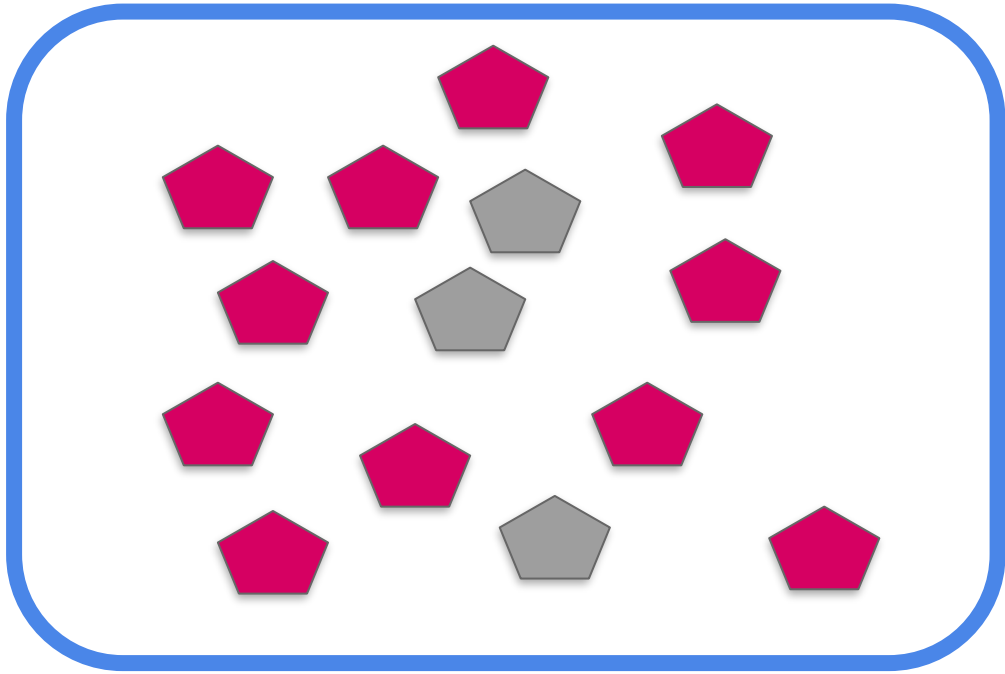
H_0 : Sesame Street has *no effect* on kids brain development

Alternative Hypothesis:

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



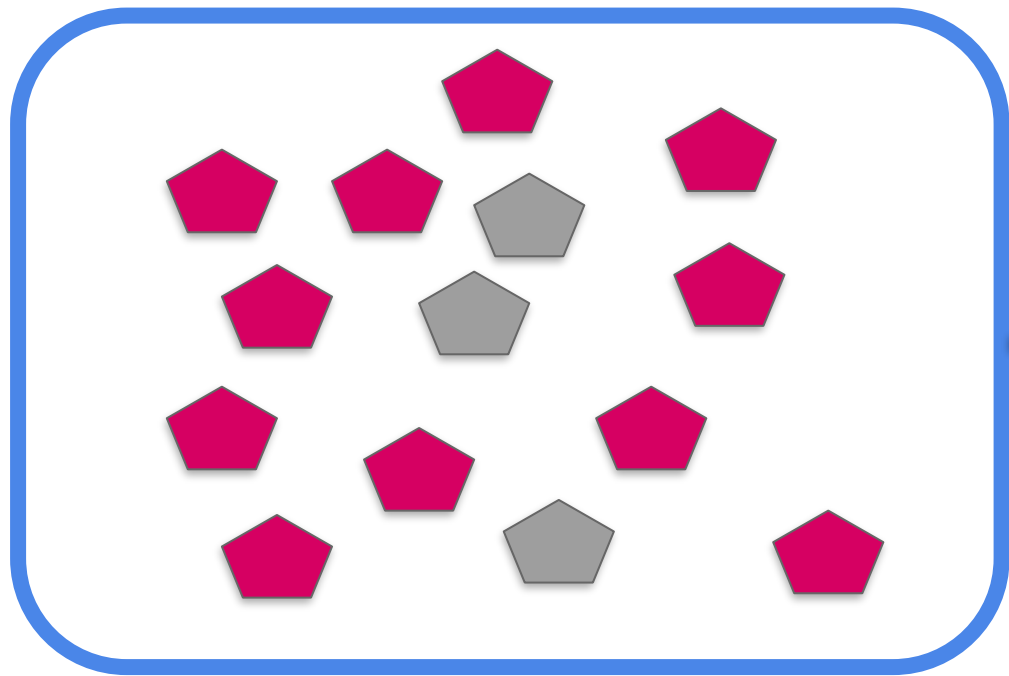
Population



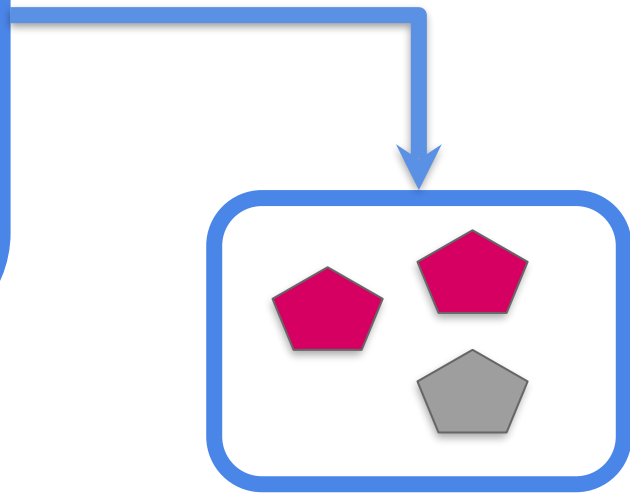
Population



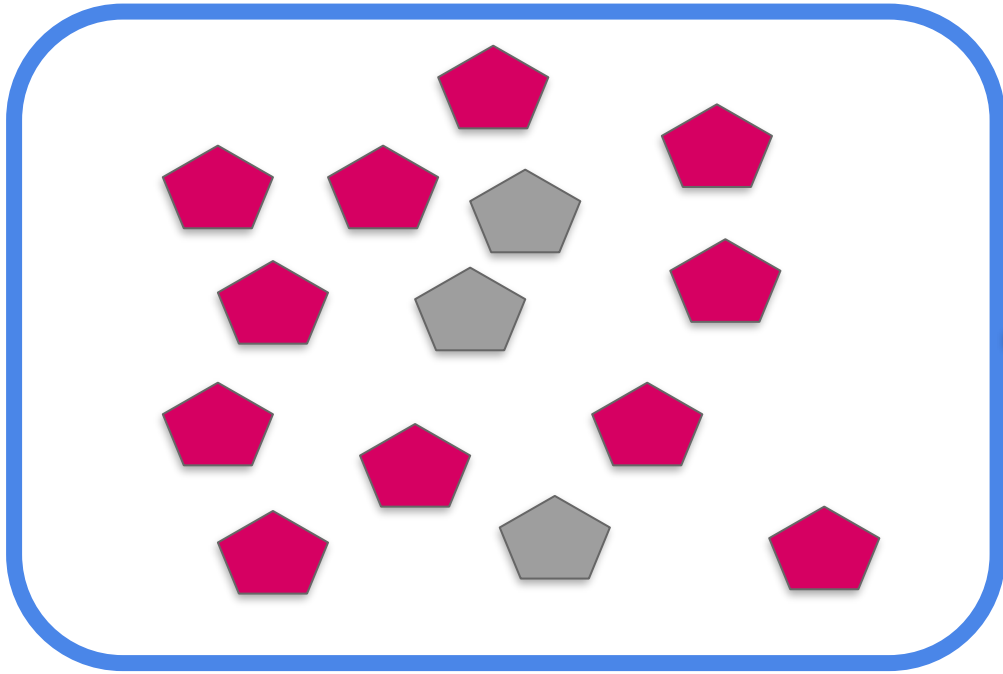
In our Sesame street example, the population would be all children



Population



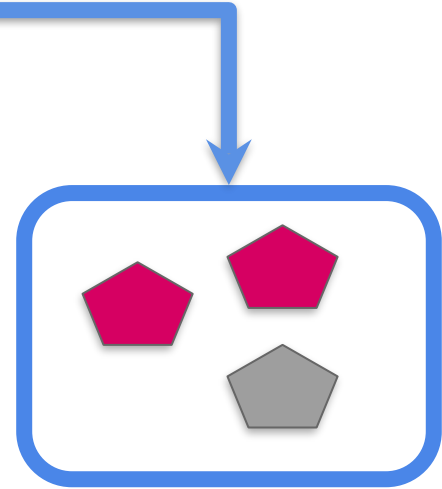
Sample



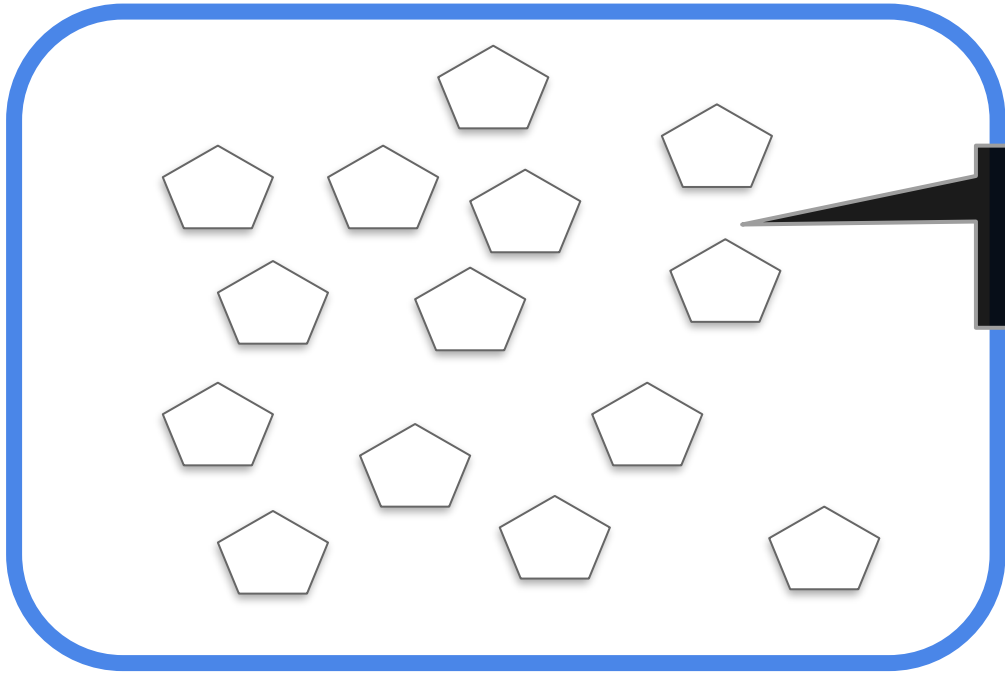
Population



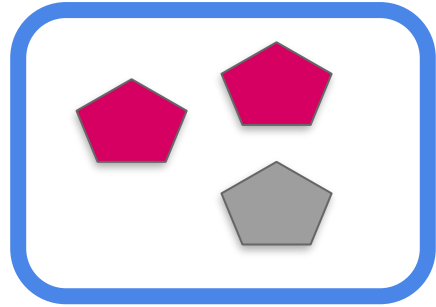
In our Sesame street example, the sample would be the children included in the study



Sample



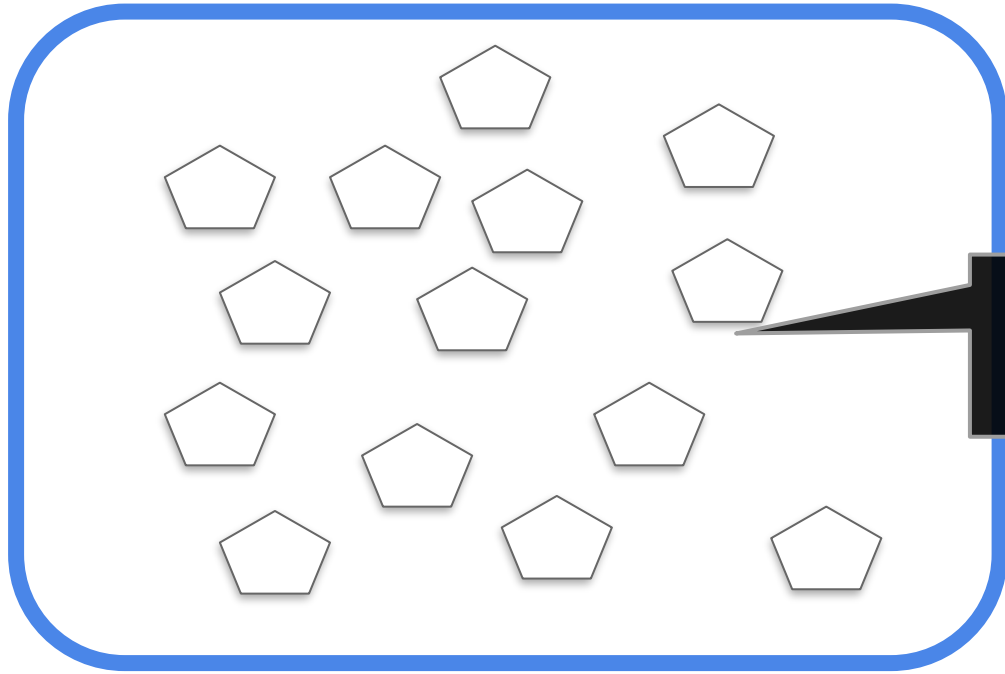
Population



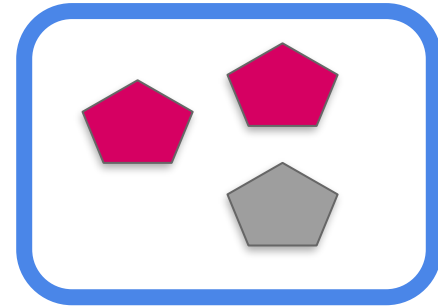
Sample



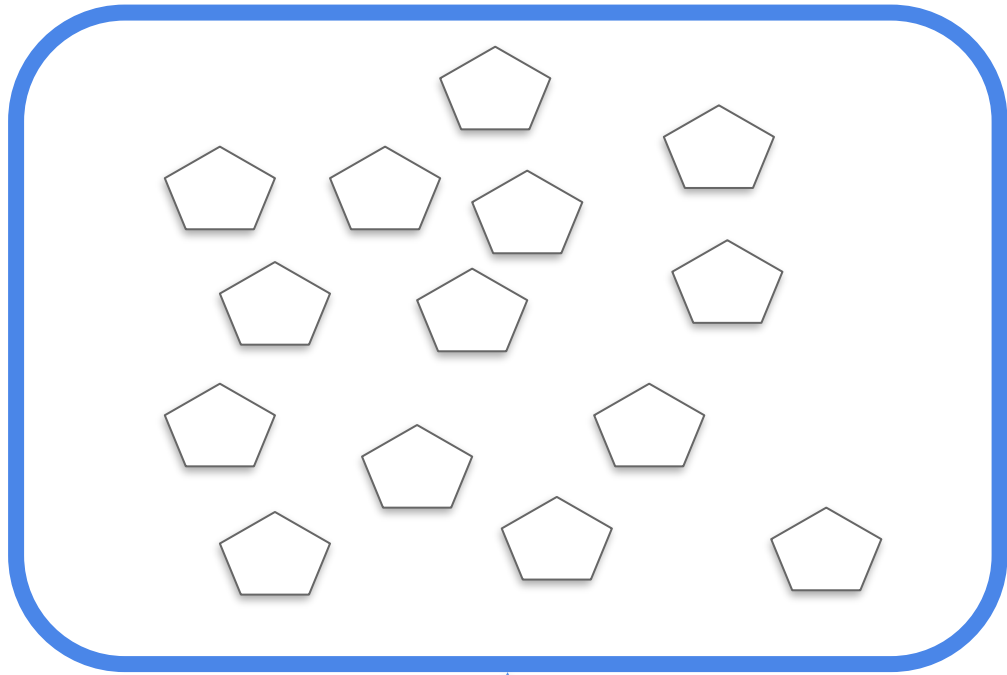
We don't know how much Sesame street was watched by or the tests scores of all kids



Population

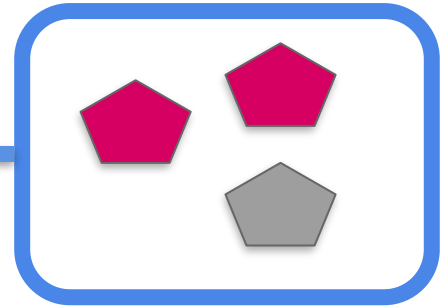


Sample

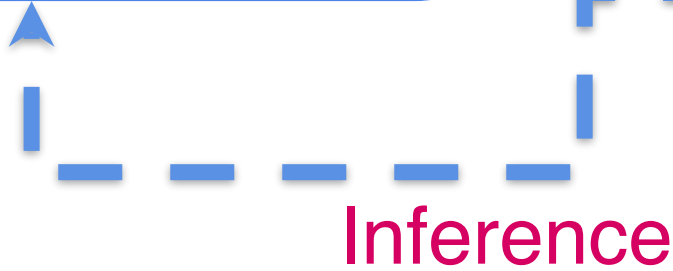


Population

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

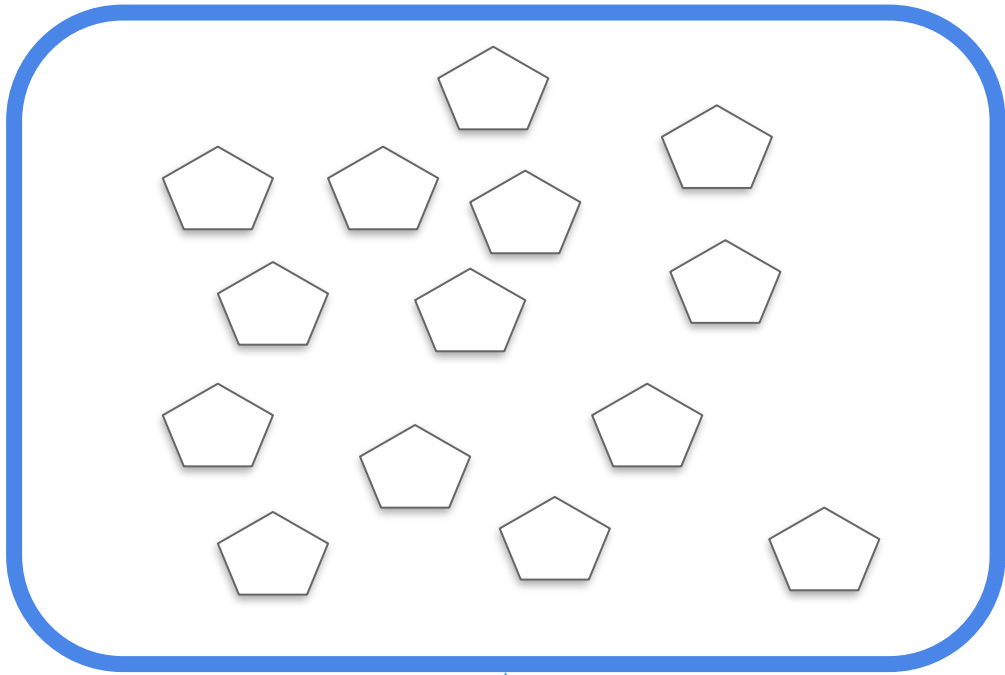


Sample



Inference

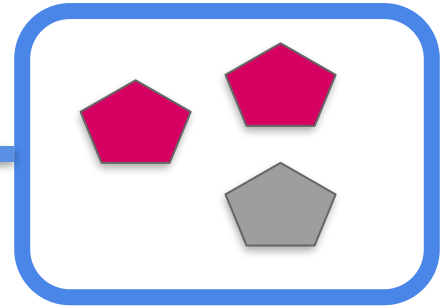




Population



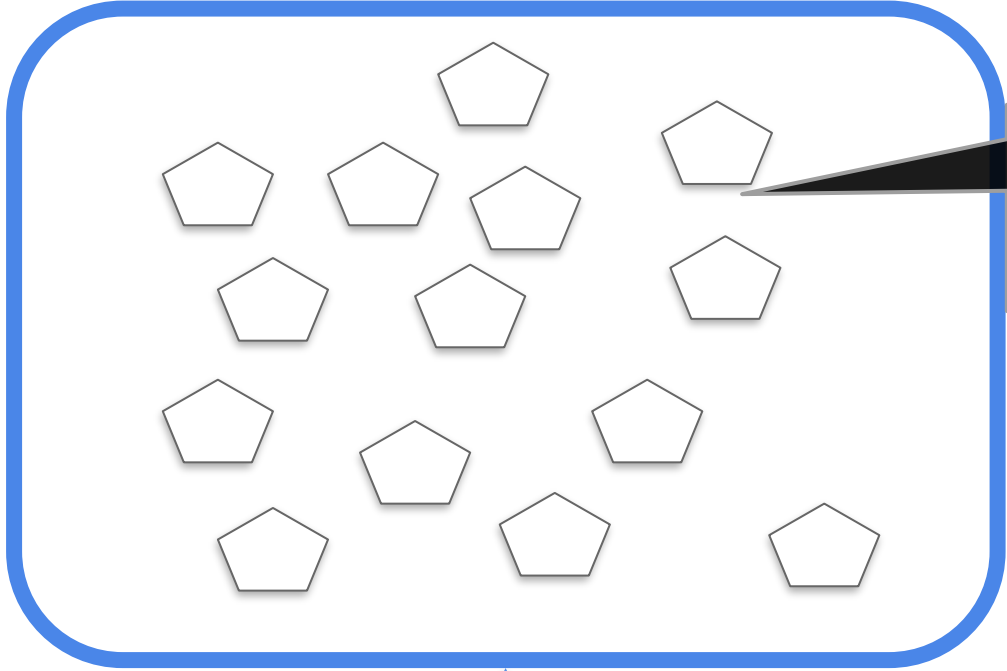
So we look at Sesame street viewing and test scores in a representative sample of kids



Sample

Inference

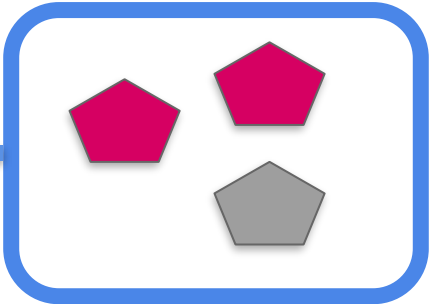




Best guess



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

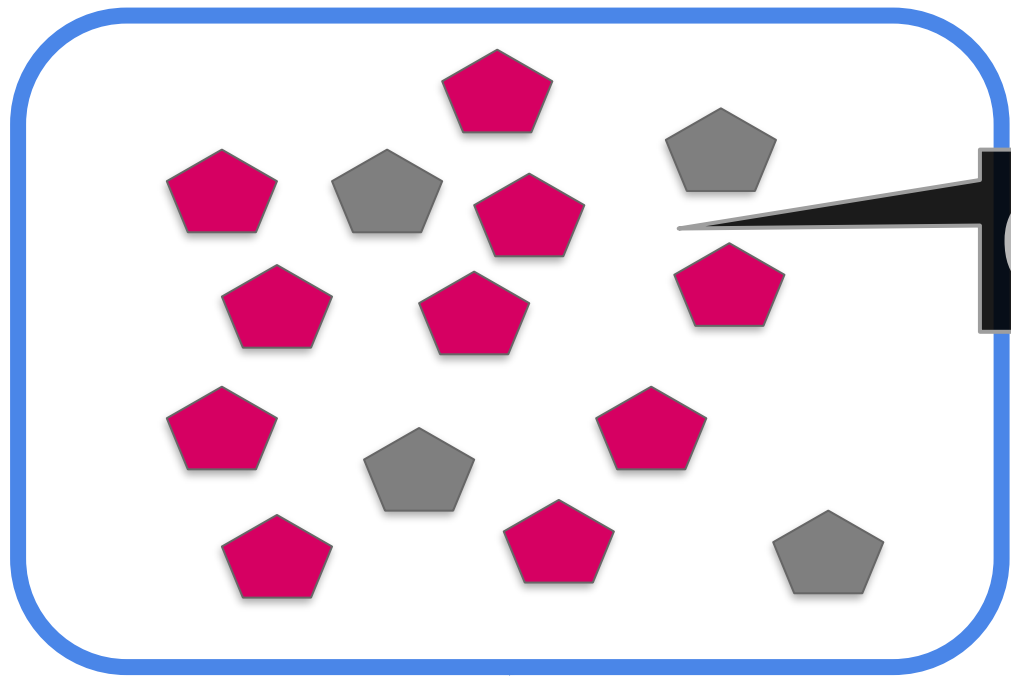


Population

Sample

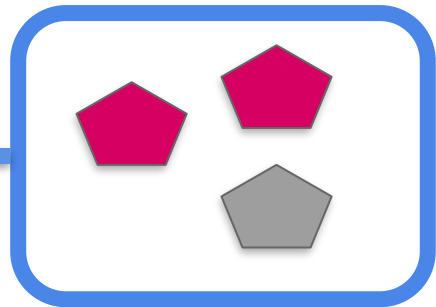


Inference!



Could be this

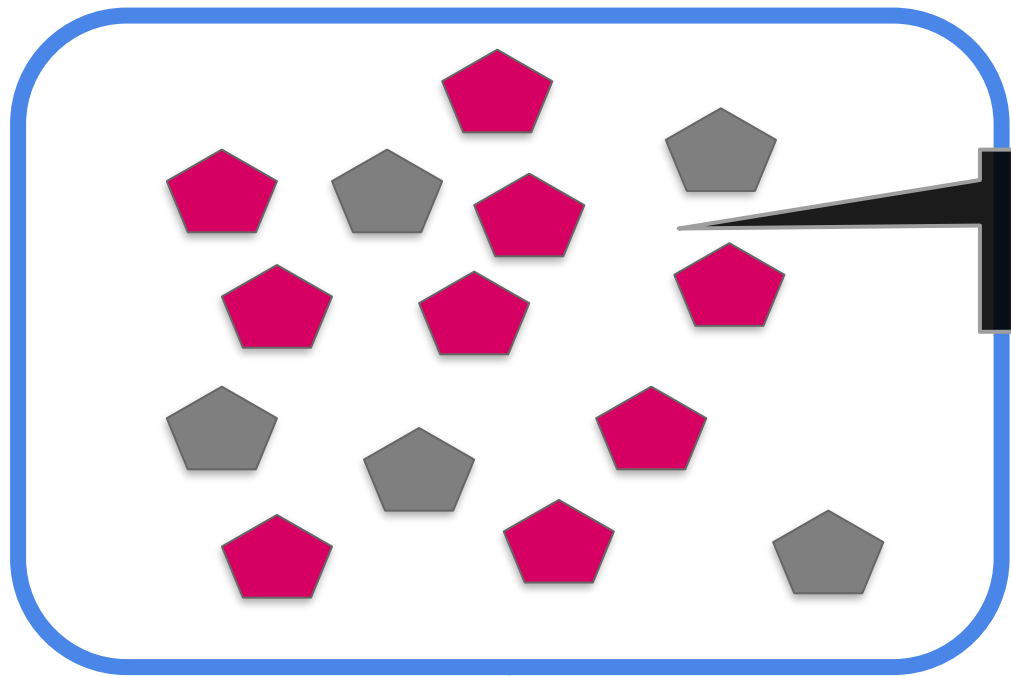
Population



Sample

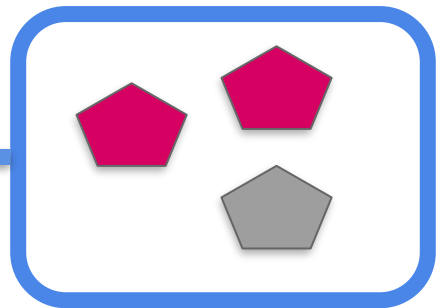


Inference



Population

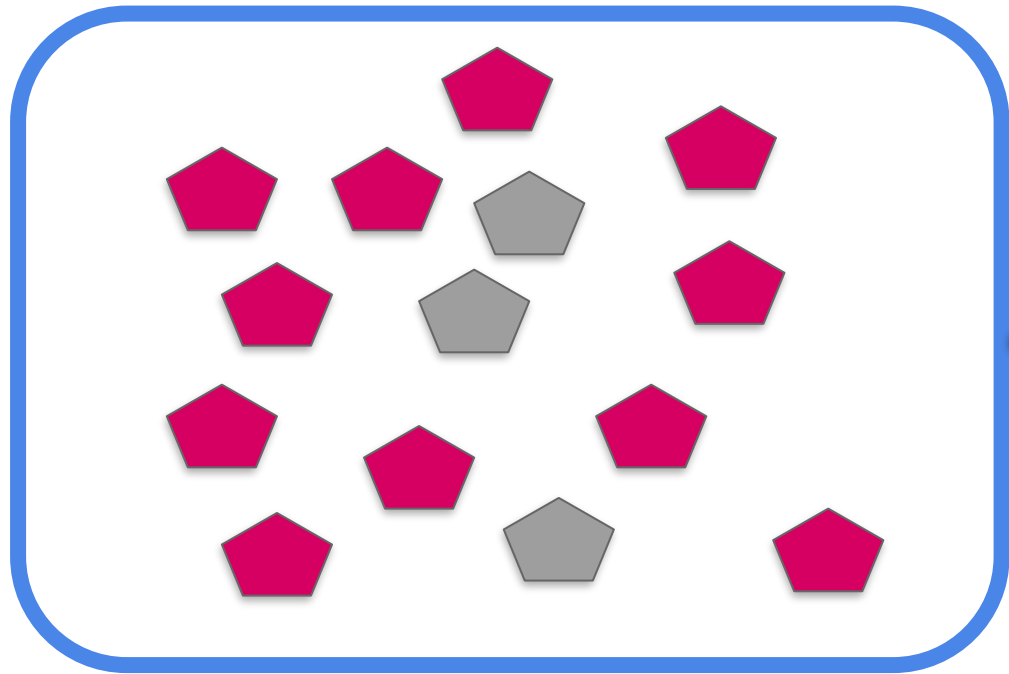
...or this



Sample

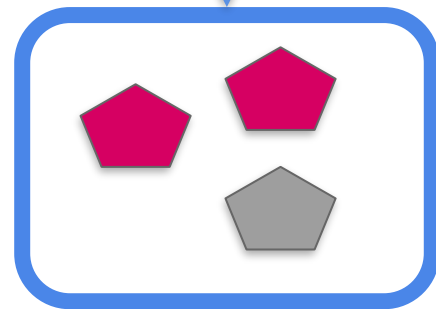


Inference

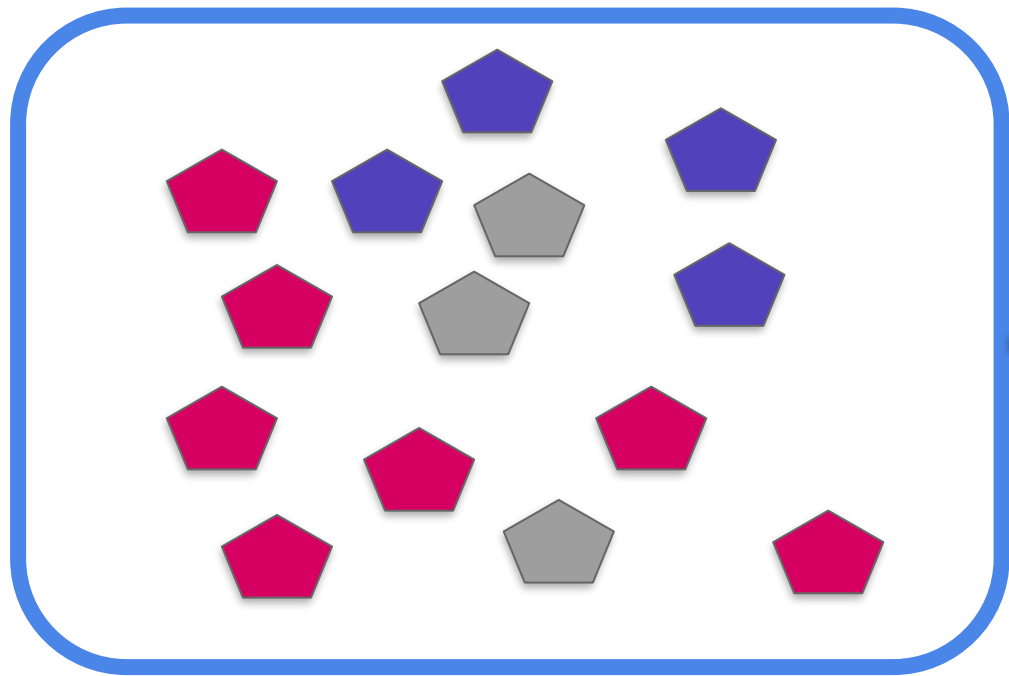


Population

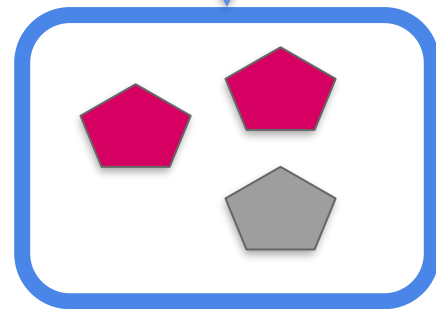
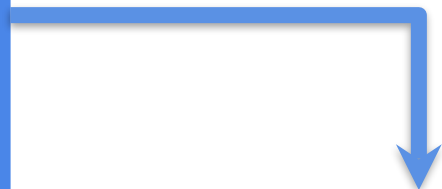
Probability



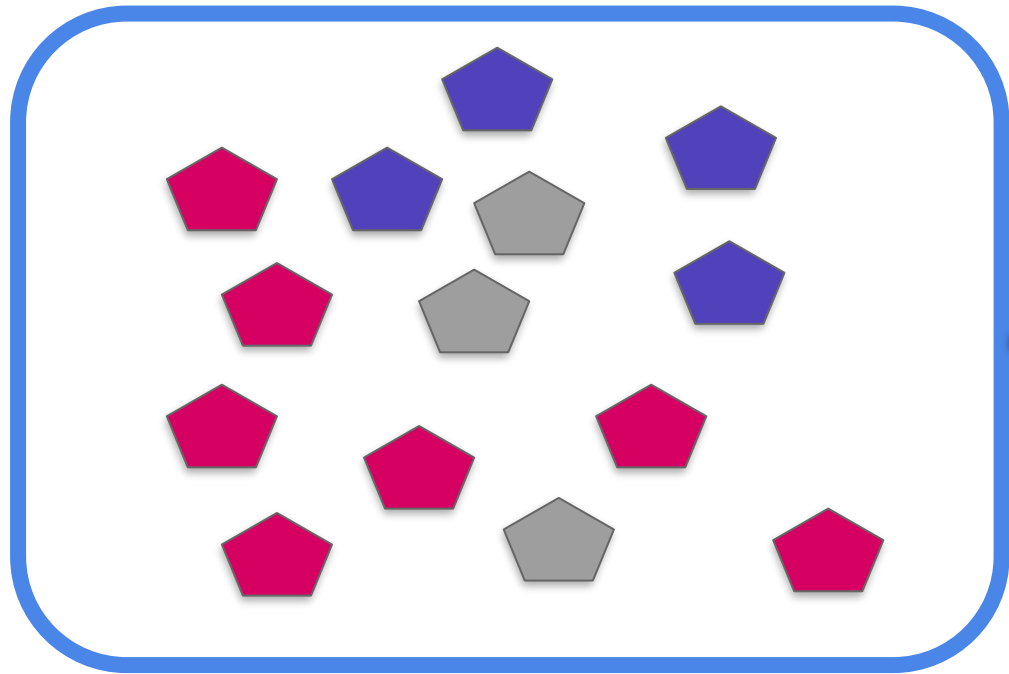
Sample



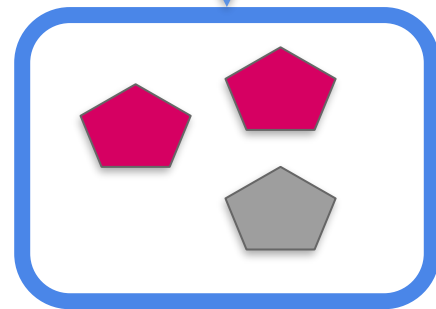
Population



Sample



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



Population

~~Inference~~

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

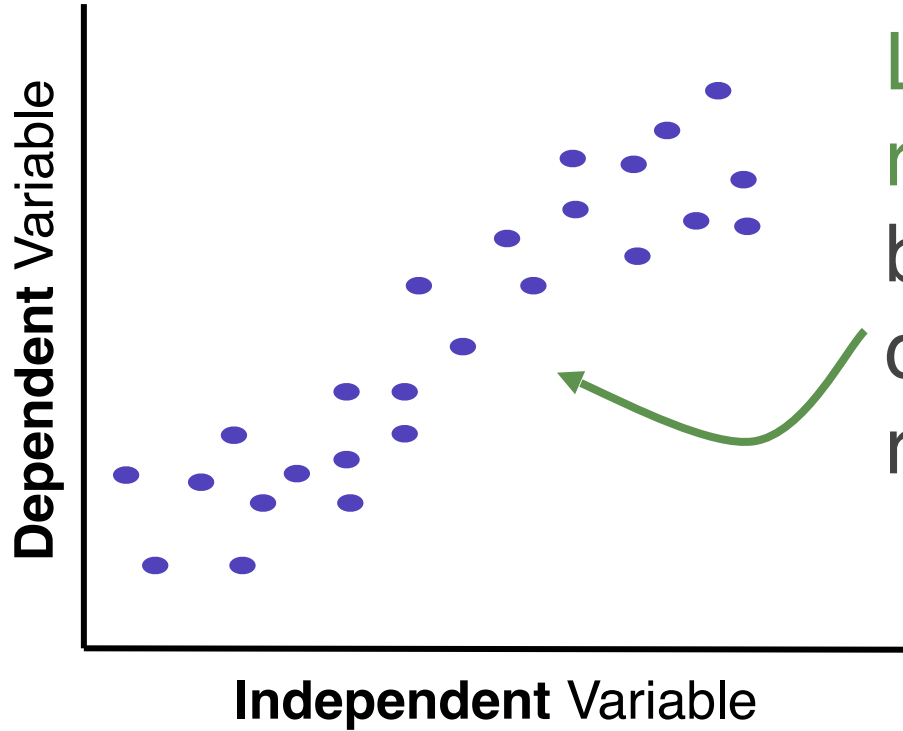
i.e. t-test, ANOVA

REGRESSION
DOES CHANGE IN ONE
VARIABLE MEAN
CHANGE IN ANOTHER?

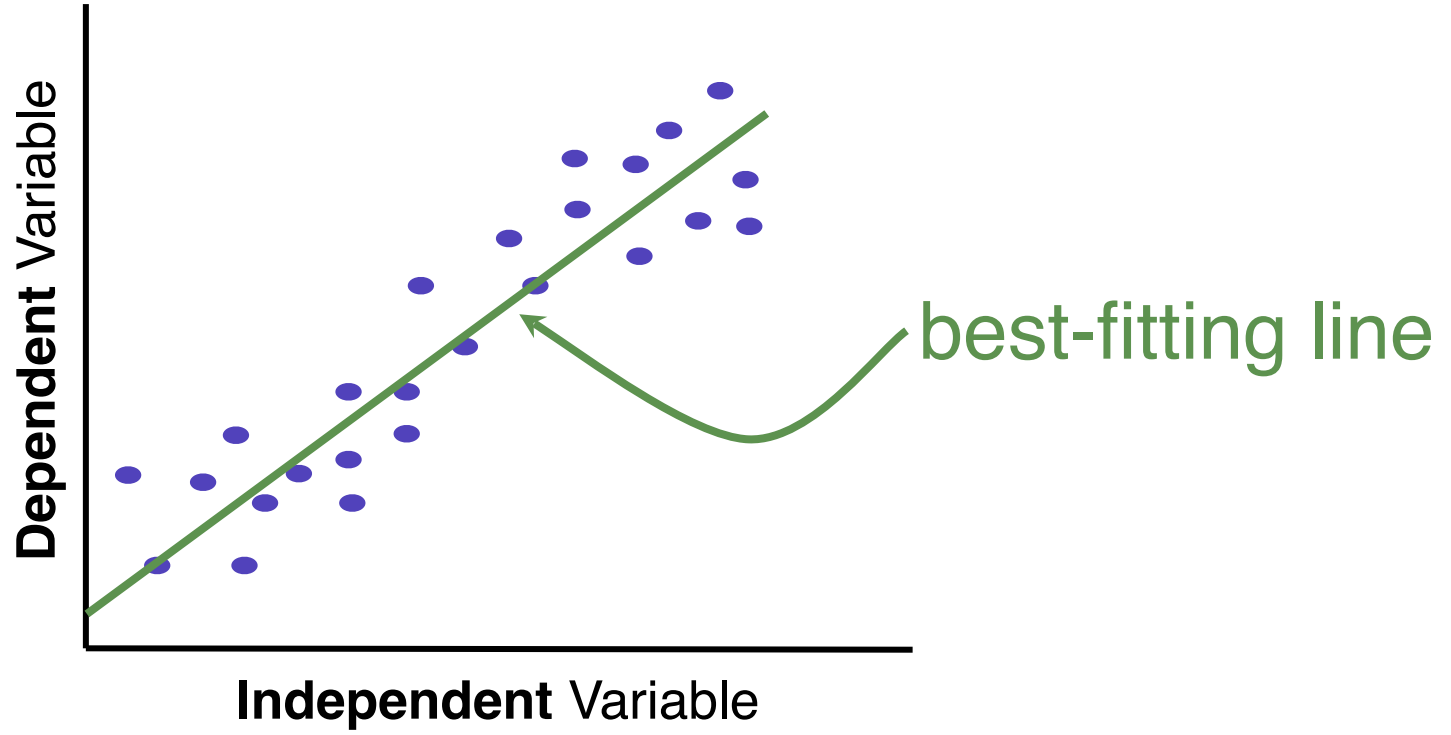
i.e. simple
regression, multiple
regression

**NON-PARAMETRIC
TESTS**
FOR WHEN
ASSUMPTIONS IN
THESE OTHER 3
CATEGORIES ARE NOT
MET

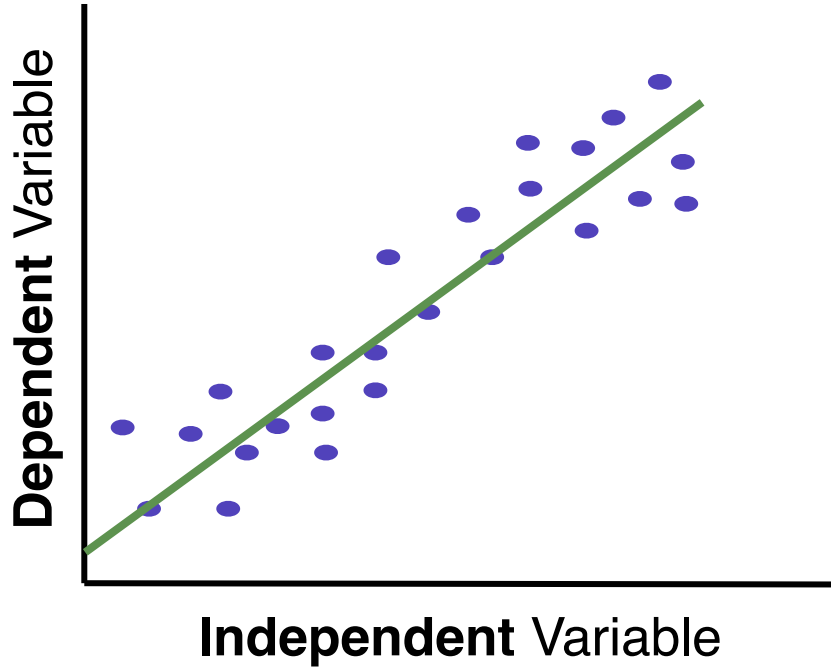
i.e. Wilcoxon rank-
sum test, Wilcoxon
sign-rank test, sign
test



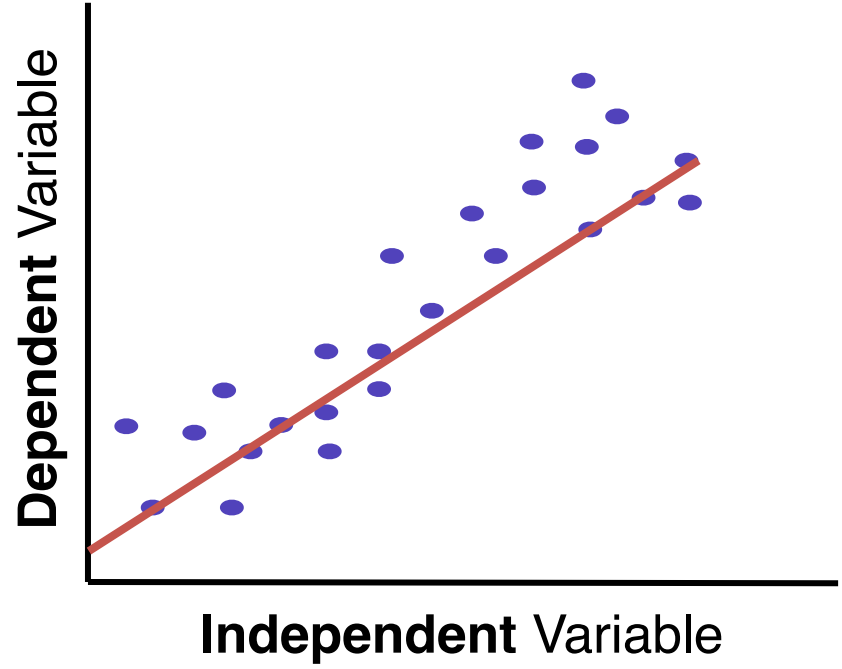
Linear regression can be used to describe this relationship

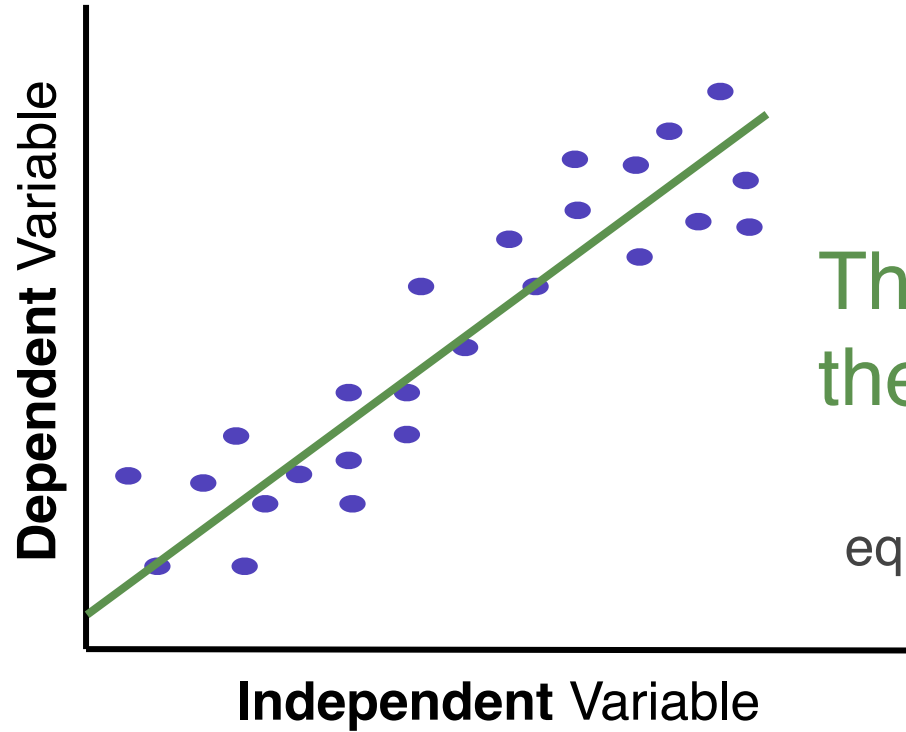


Best-fitting line



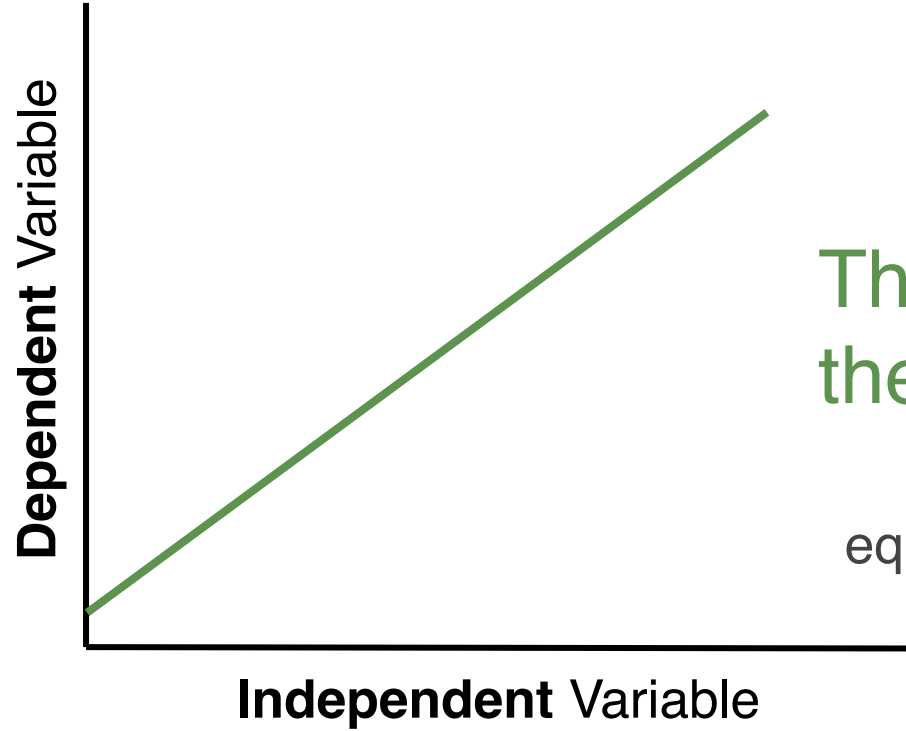
NOT a best-fitting line





This line is a model of the data

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

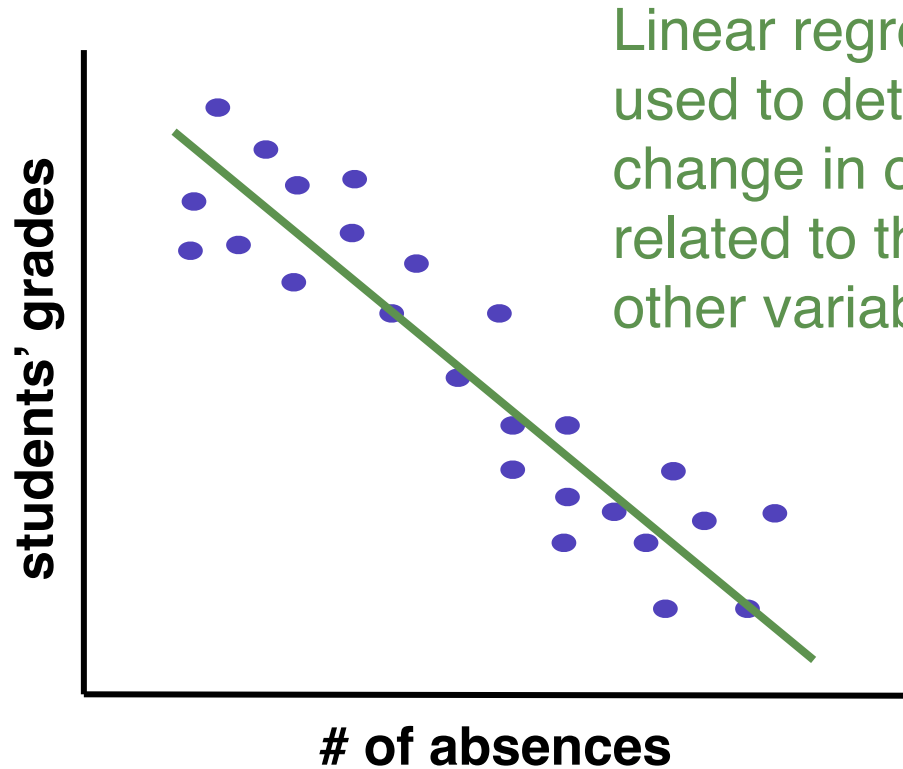


This line is a model of the data

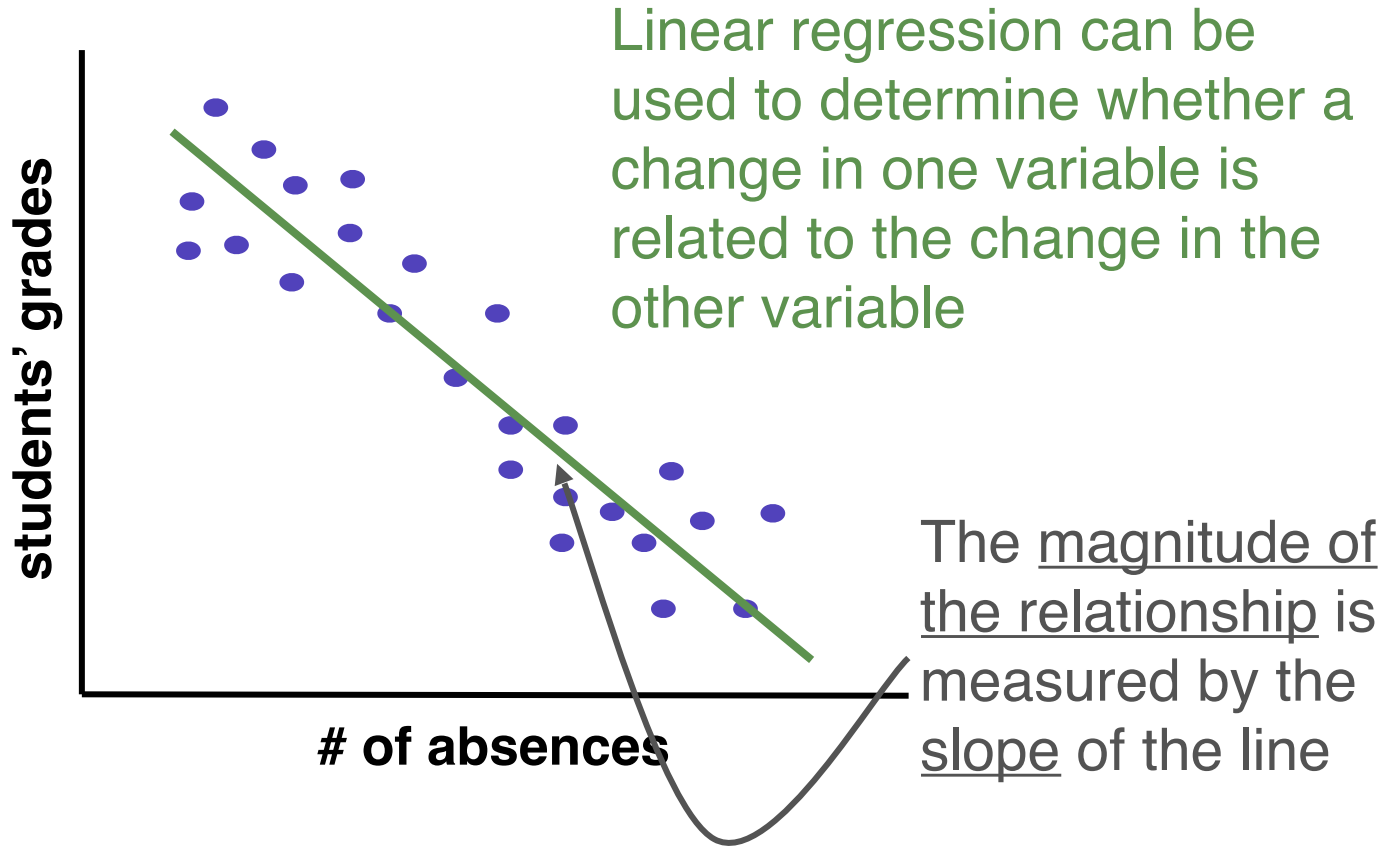
Models are mathematical equations generated to *represent* the real life situation

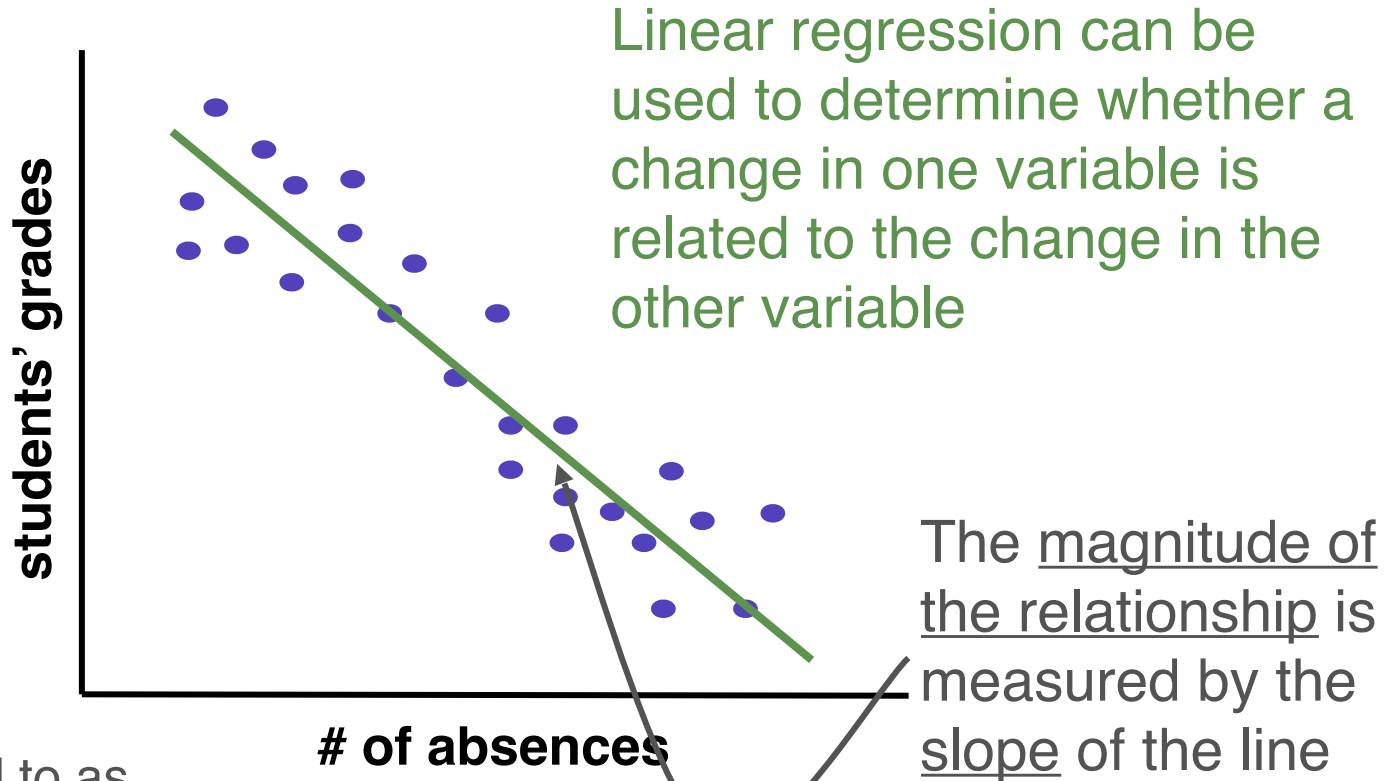
“All models are wrong, but some are useful”

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



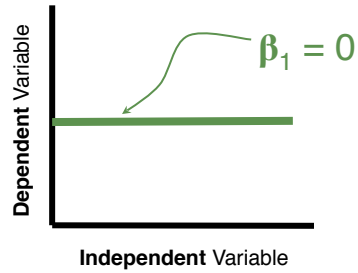
Linear regression can be used to determine whether a change in one variable is related to the change in the other variable



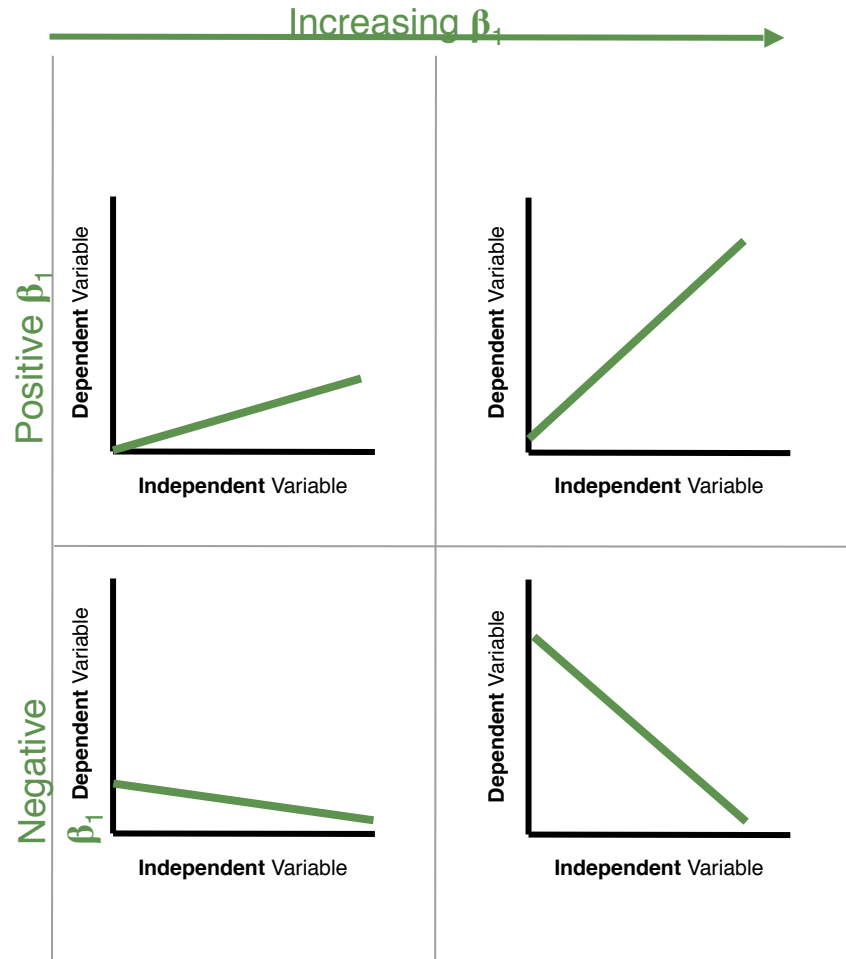
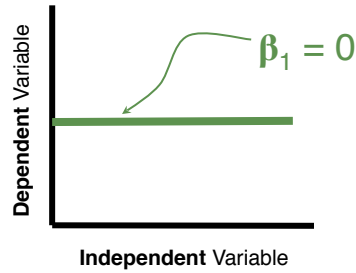


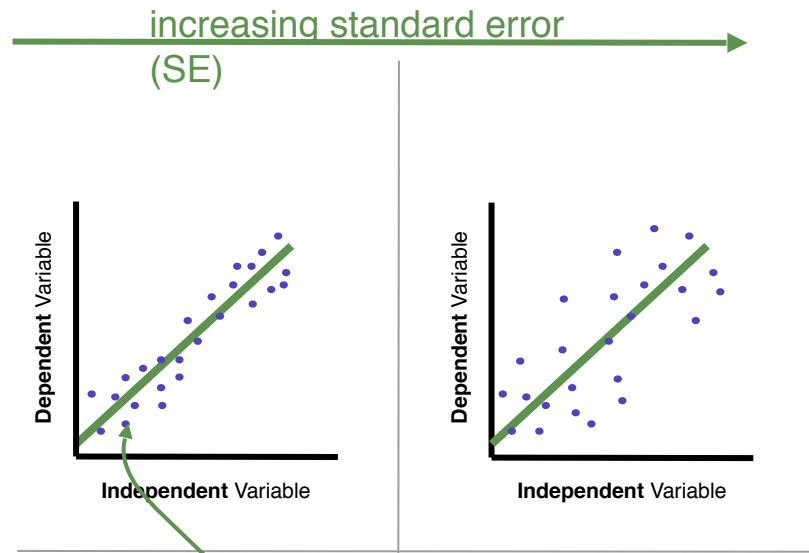
This is also referred to as the model's effect size (β_1)

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



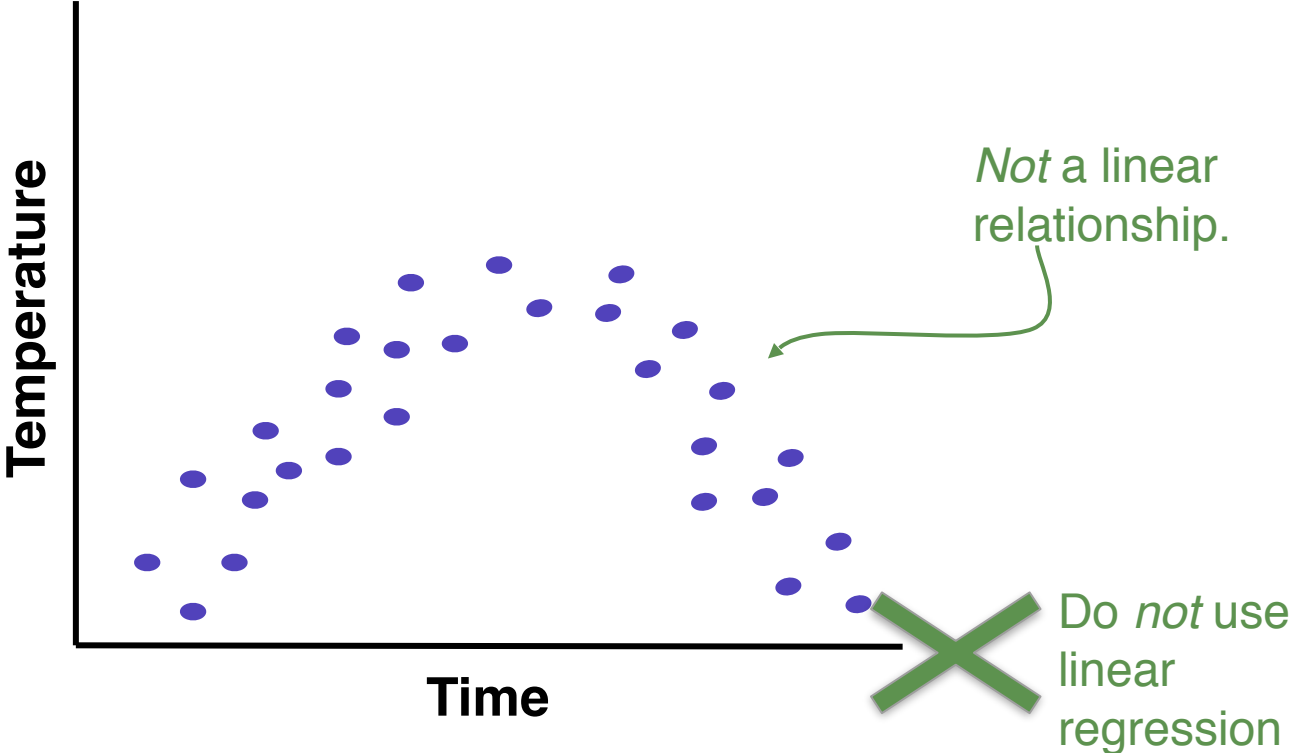


The *closer* the points are to the regression line, the *less uncertain* we are in our estimate

Assumptions of linear regression

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

Linearity



Multicollinearity

- Linear regression assumes no multicollinearity. **Multicollinearity** occurs when the independent variables (in multiple linear regression) are too highly correlated with each other.
- 2 variables are perfectly correlated if they have a correlation coefficient of 1.0

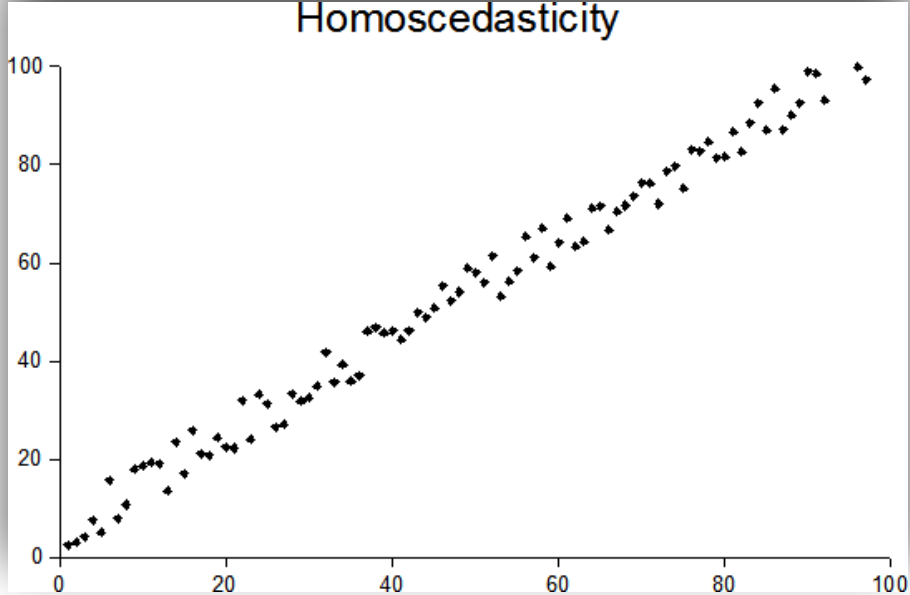
Autocorrelation

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

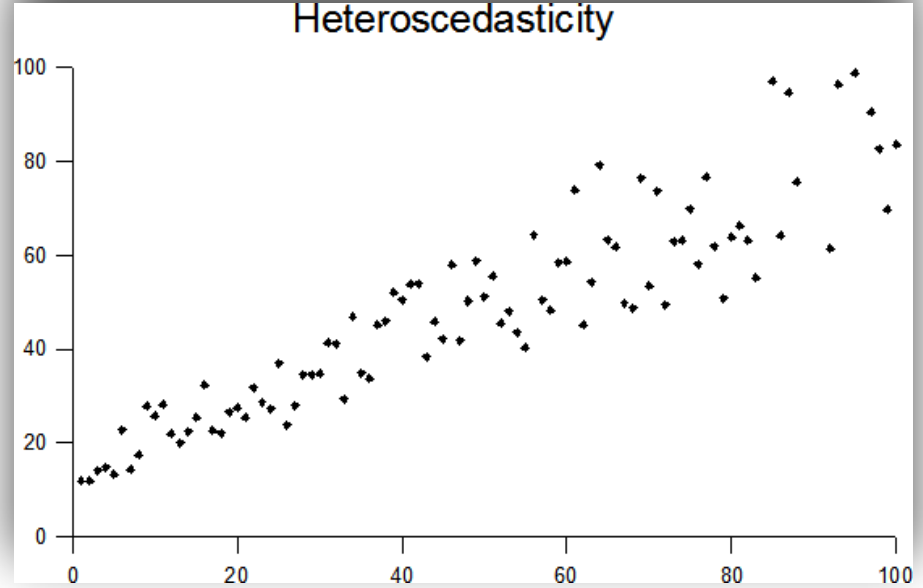


Homoscedasticity - a reminder of what that is

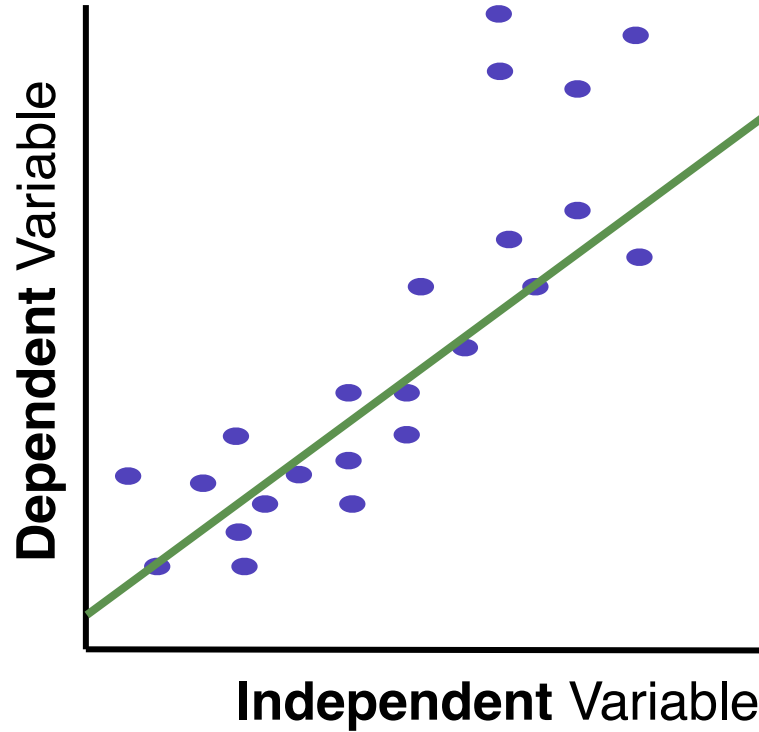
Homoscedasticity



Heteroscedasticity



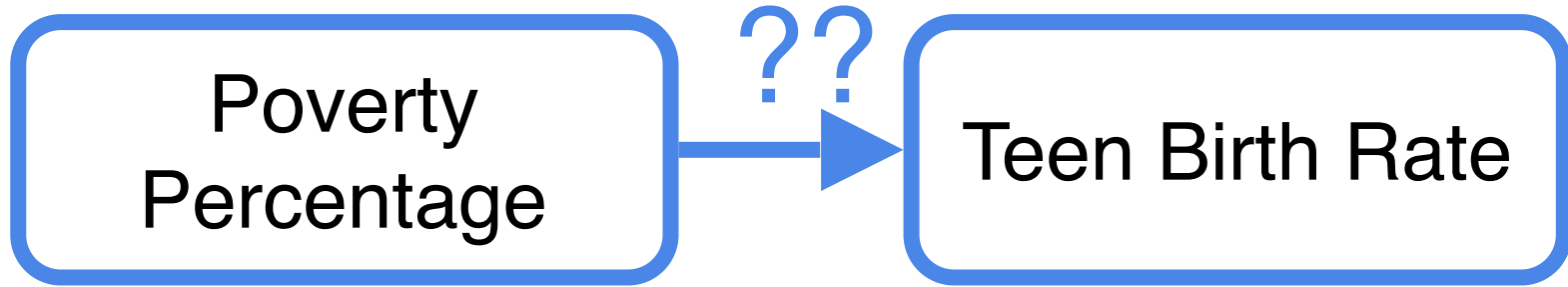
Homoscedasticity



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

X Do *not* use
linear
regression

Does Poverty
Percentage affect Teen
Birth Rate?



Null Hypothesis:

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

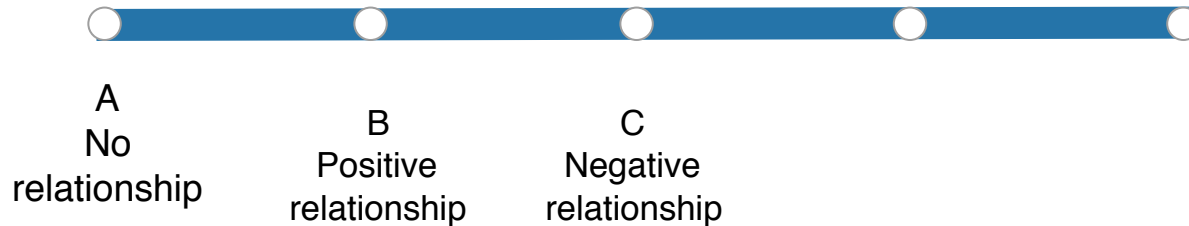
Alternative Hypothesis:

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



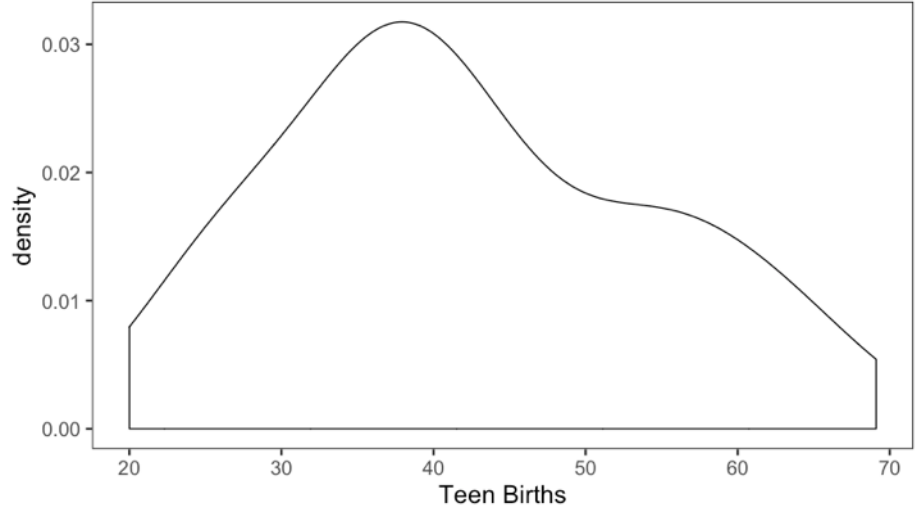
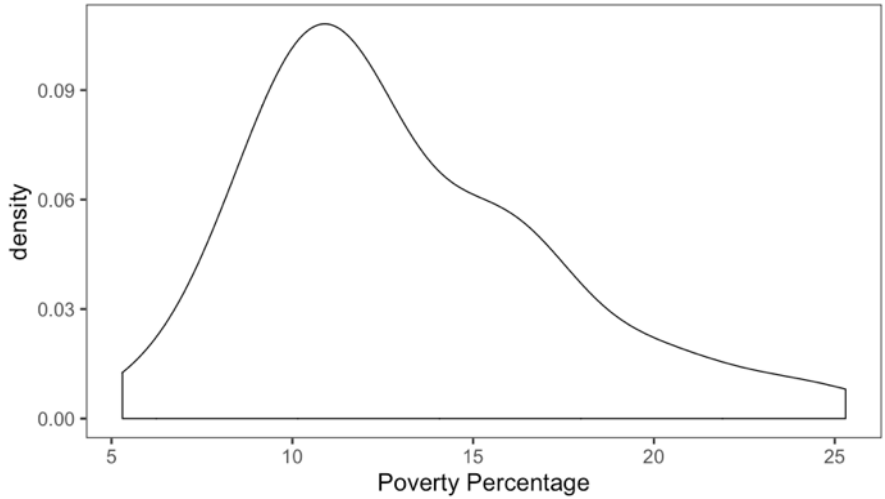
What is the relationship between Poverty Percentage & Teen Birth Rate?

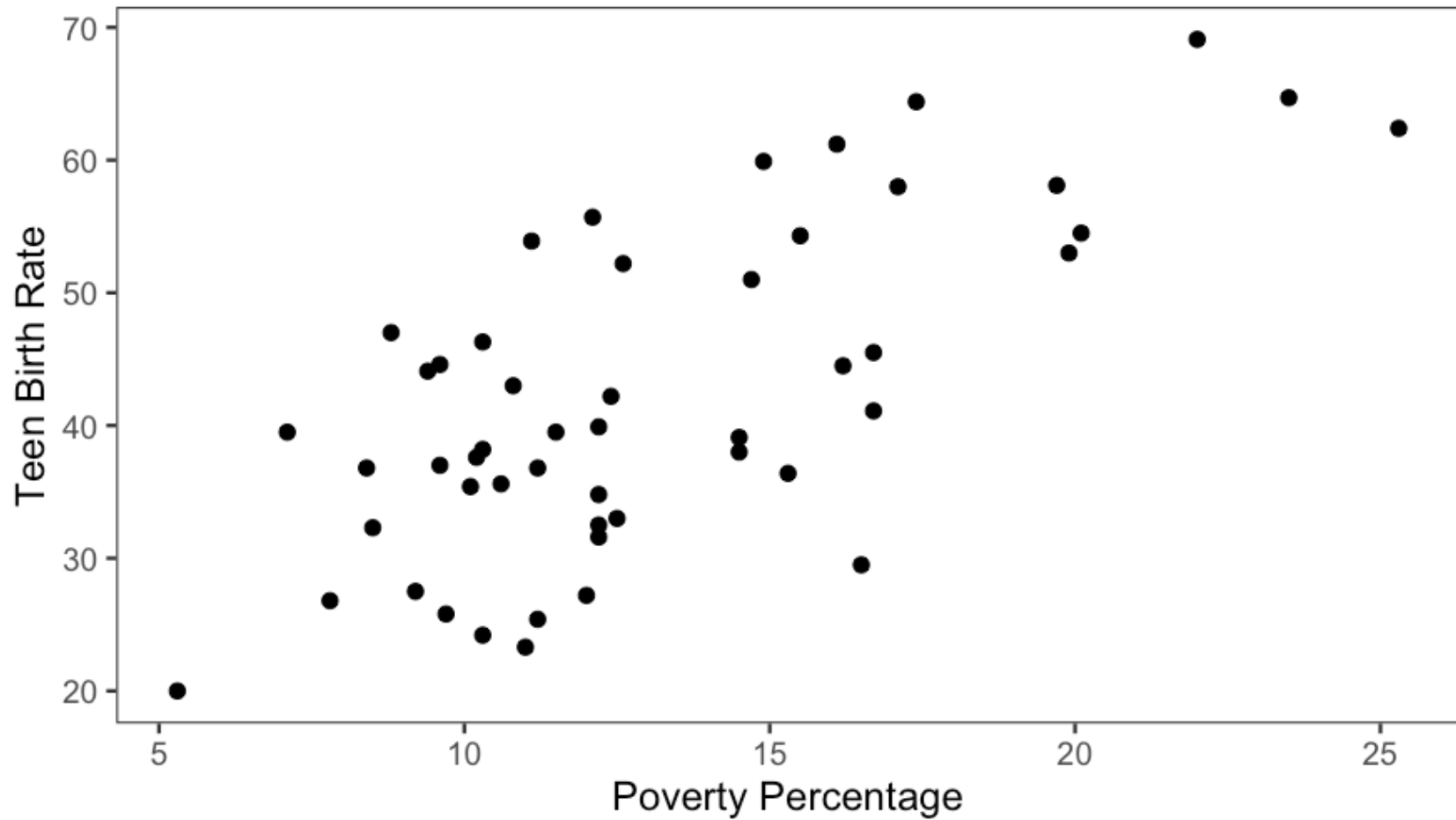
What's your hypothesis?

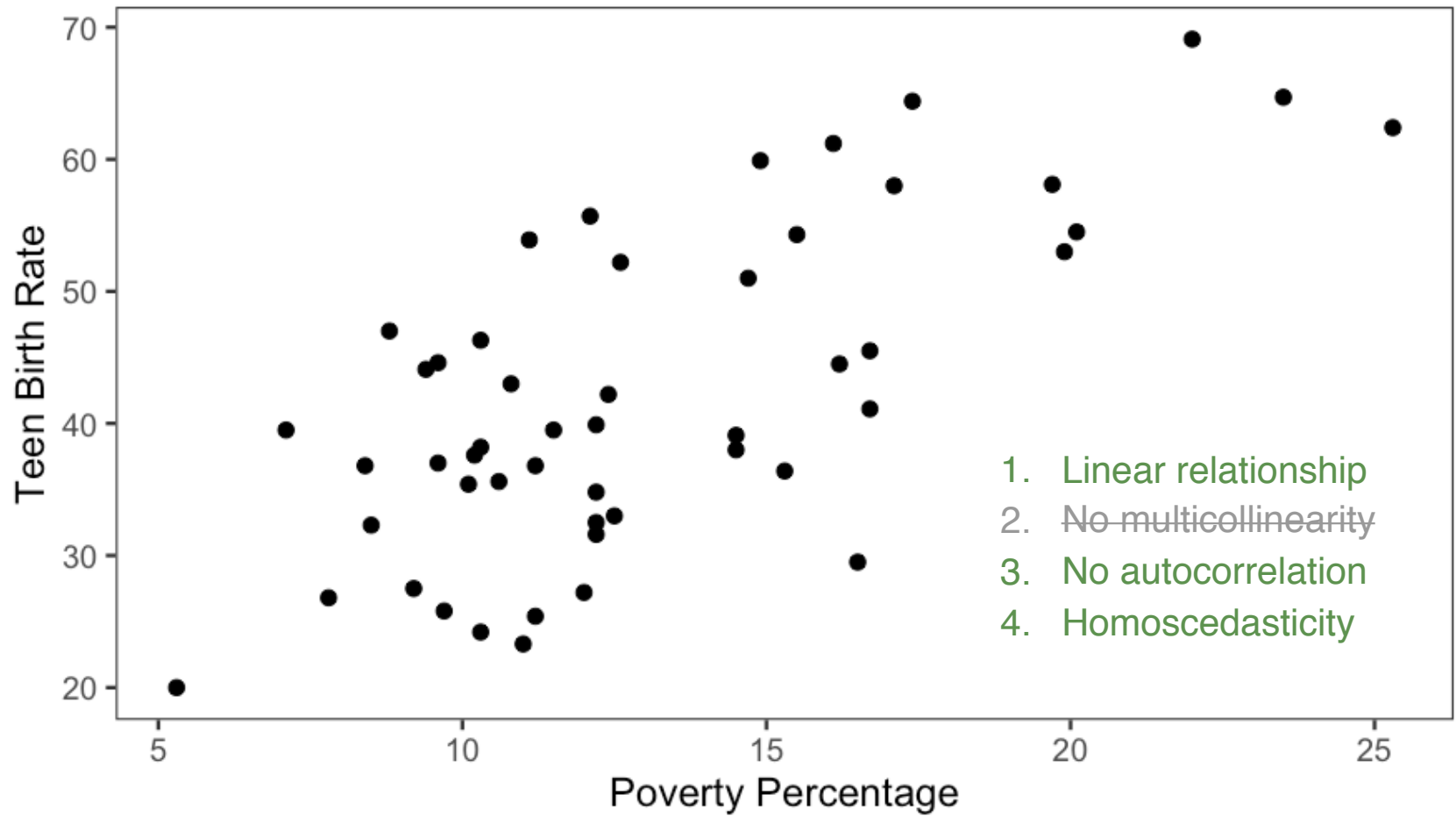


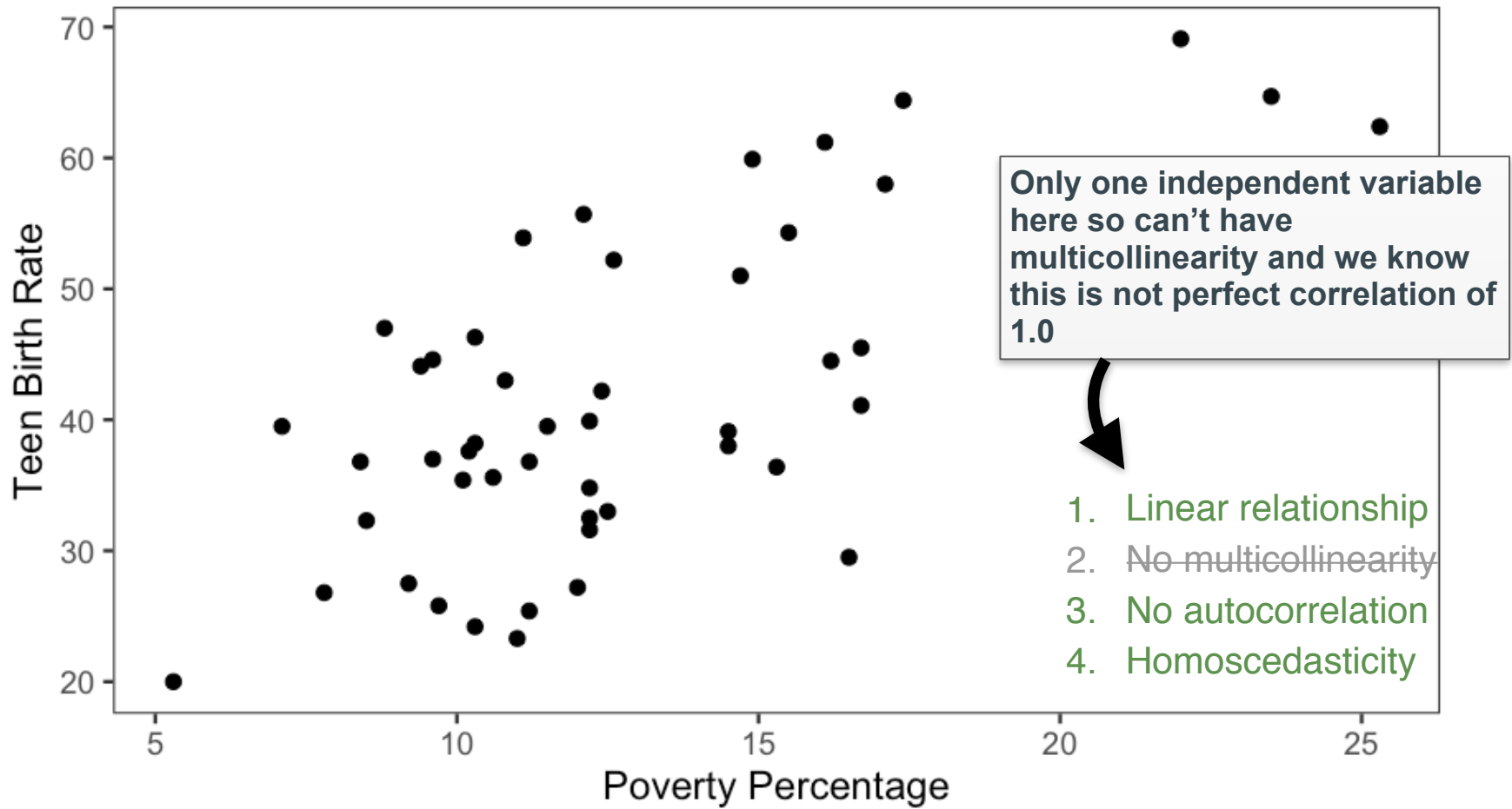
	Location	PovPct	Brth15to17	Brth18to19	ViolCrime	TeenBrth
1	Alabama	20.1	31.5	88.7	11.2	54.5
2	Alaska	7.1	18.9	73.7	9.1	39.5
3	Arizona	16.1	35.0	102.5	10.4	61.2
4	Arkansas	14.9	31.6	101.7	10.4	59.9
5	California	16.7	22.6	69.1	11.2	41.1
6	Colorado	8.8	26.2	79.1	5.8	47.0
7	Connecticut	9.7	14.1	45.1	4.6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

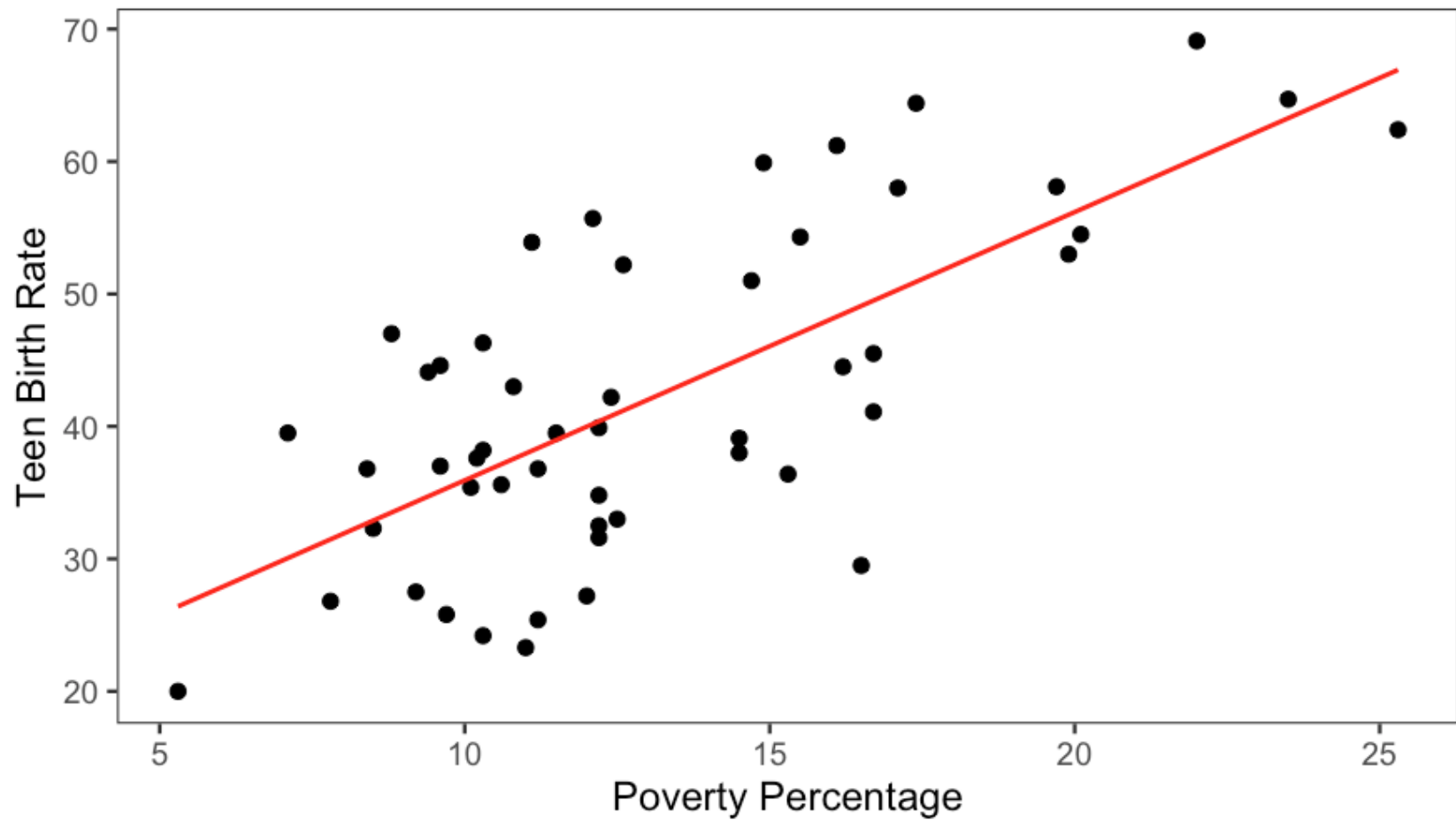
EDA: distributions

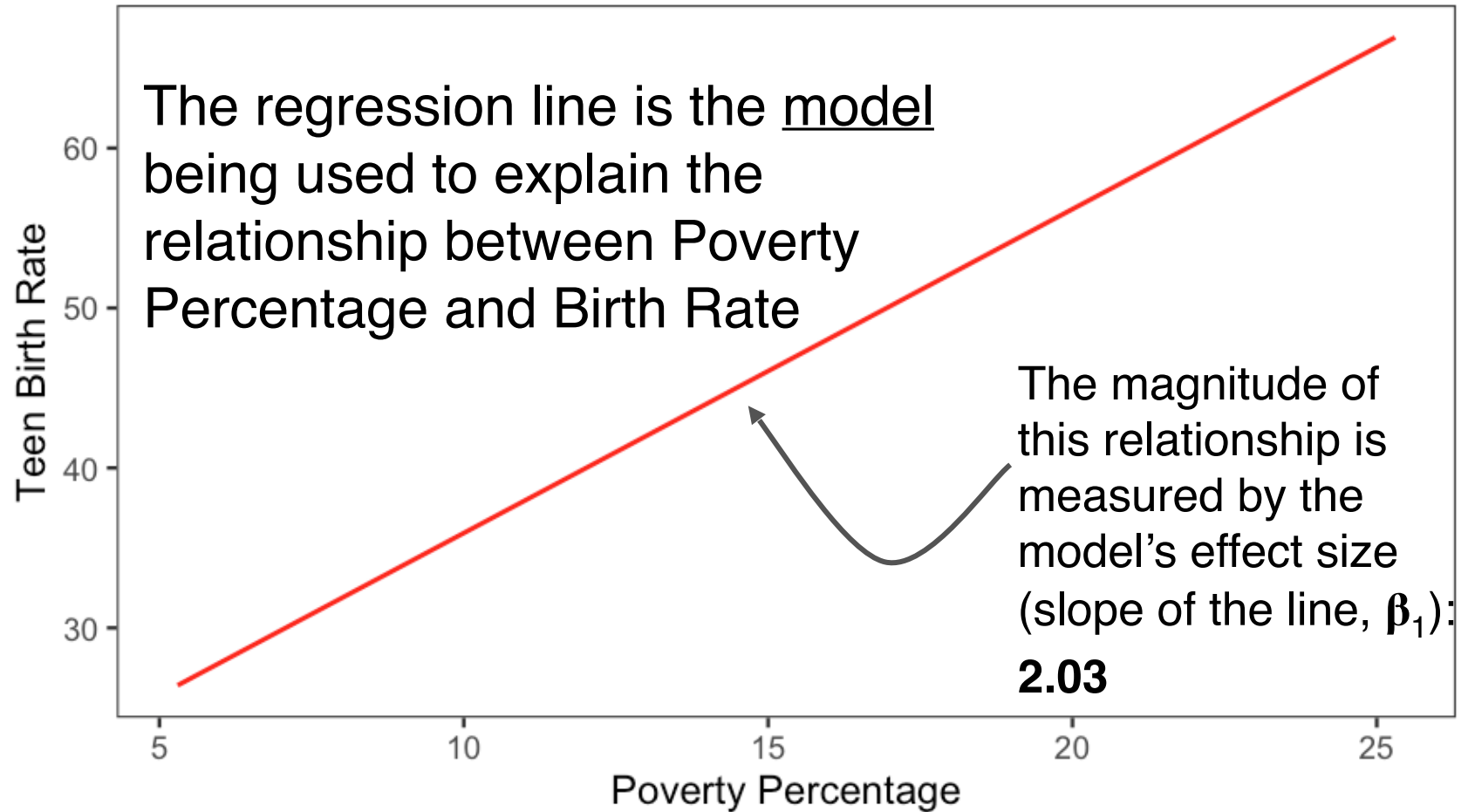


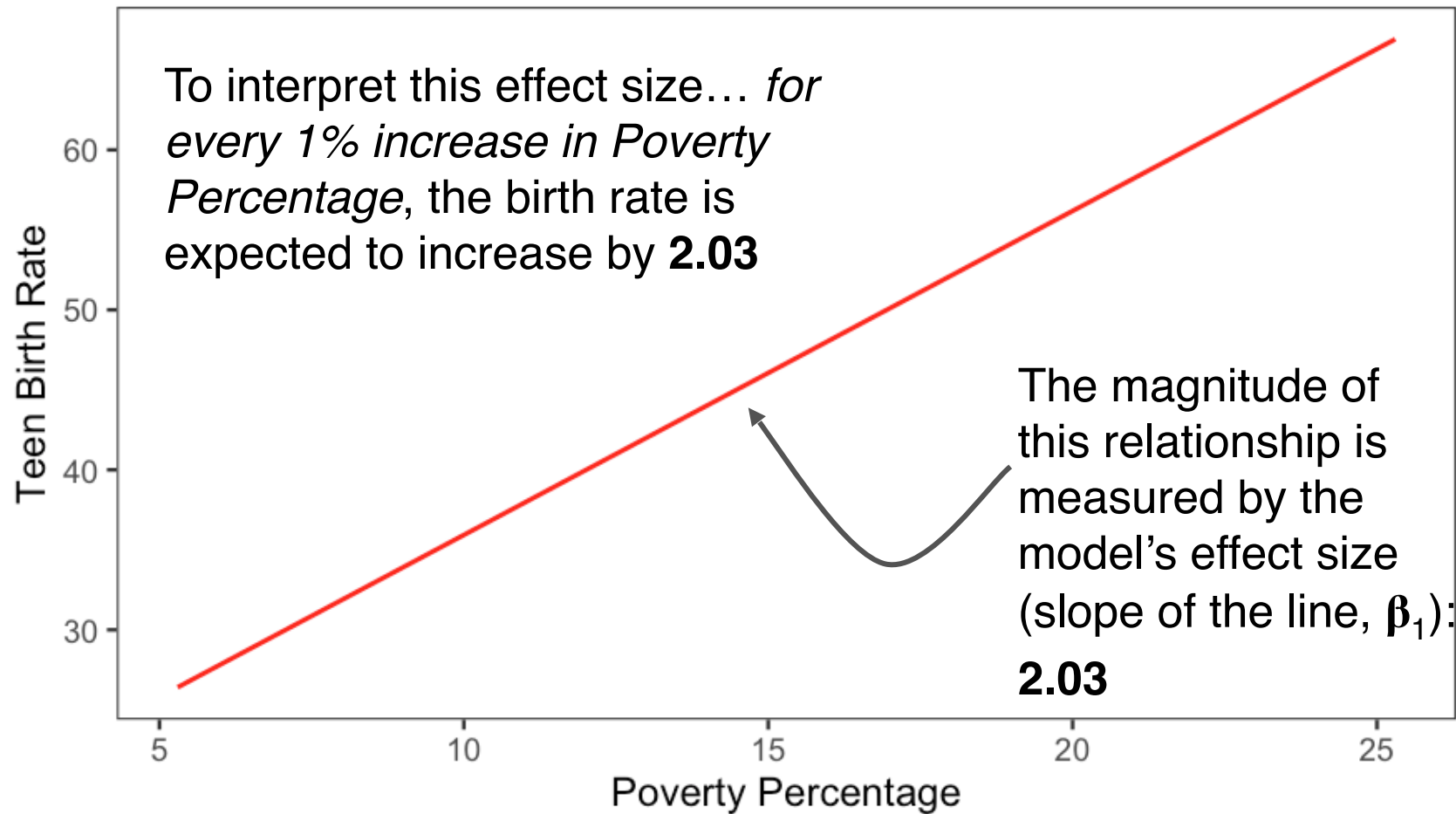






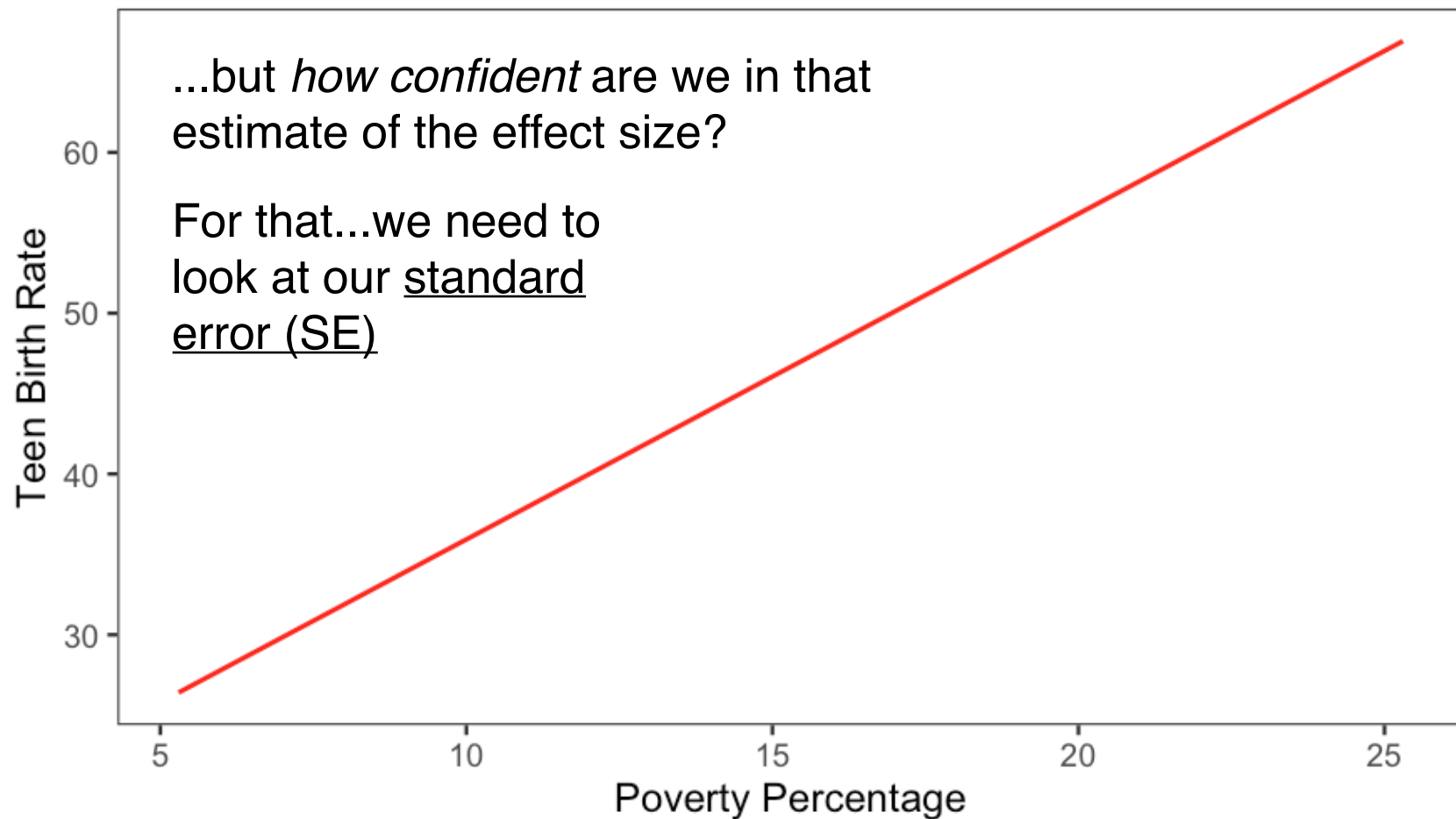


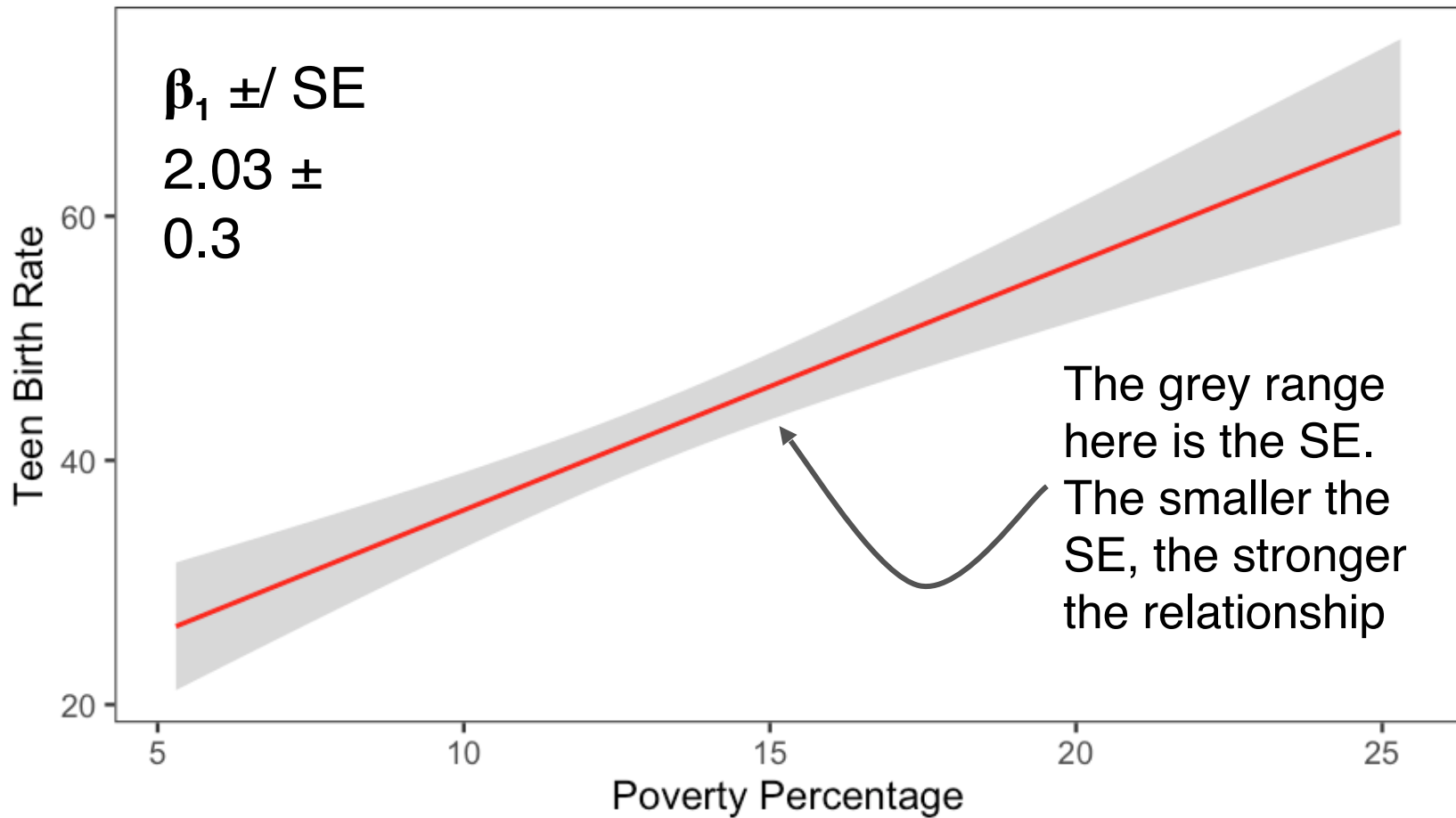


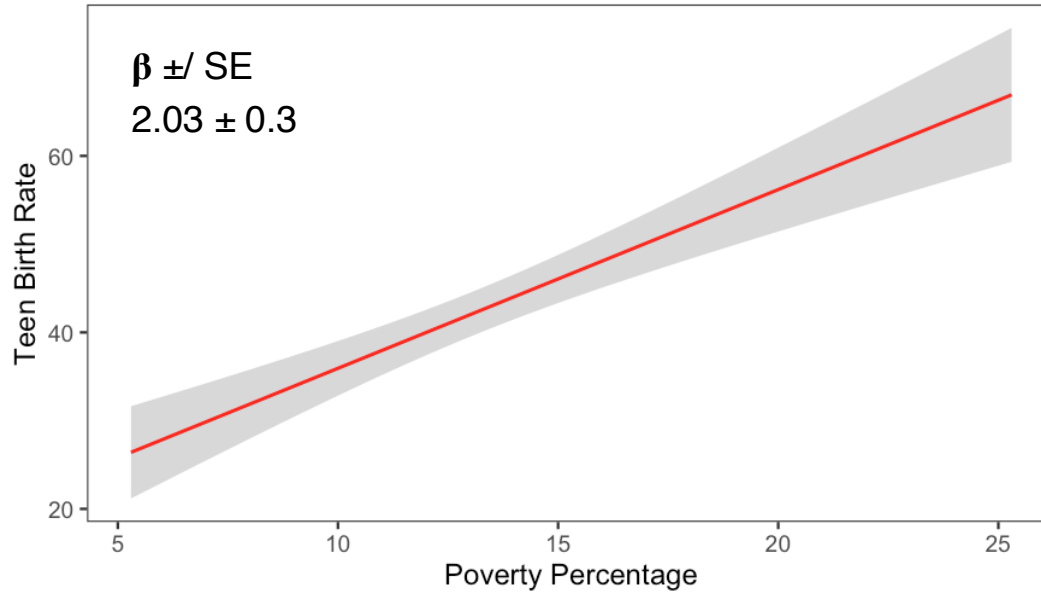


...but *how confident* are we in that estimate of the effect size?

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

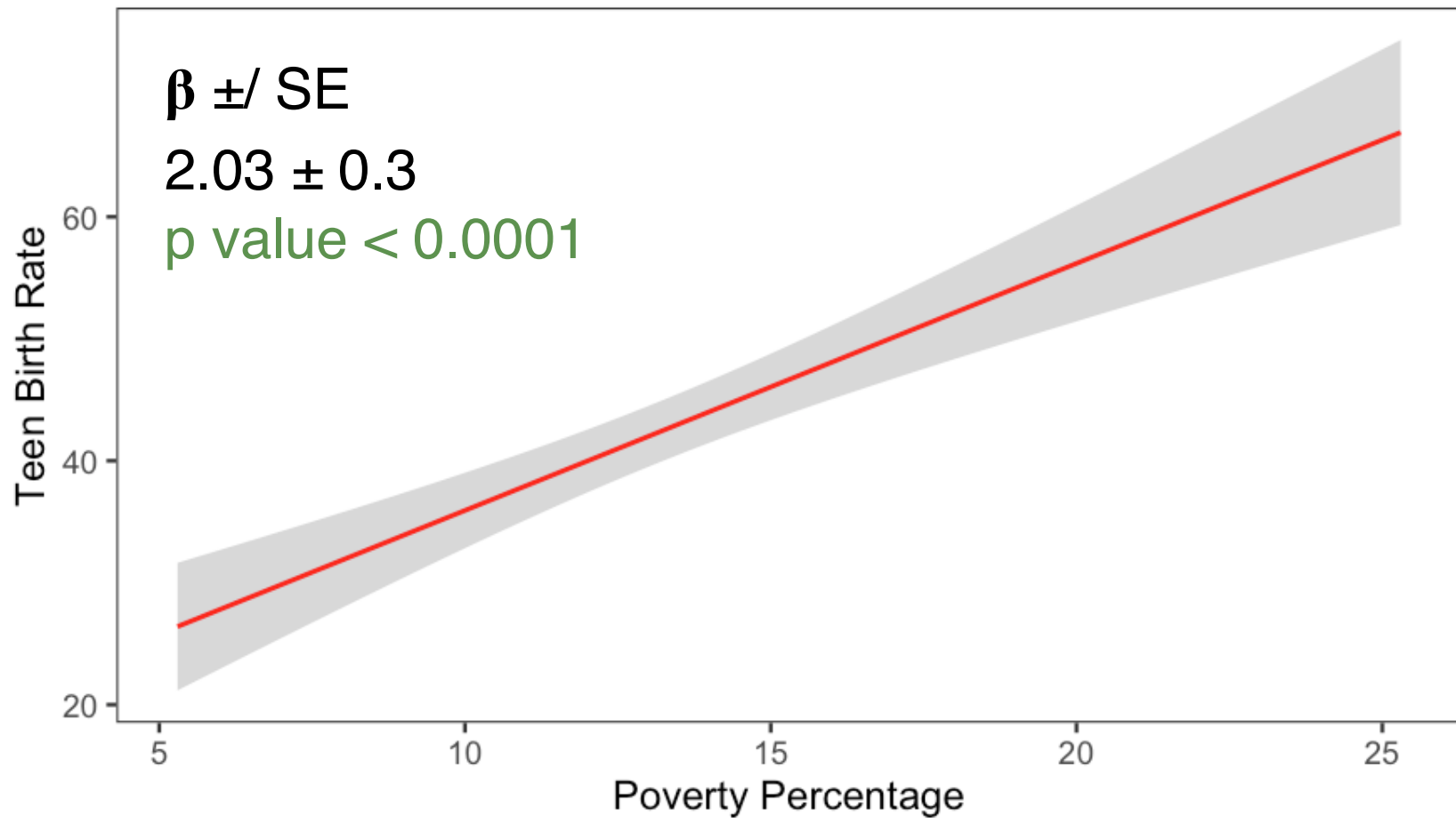




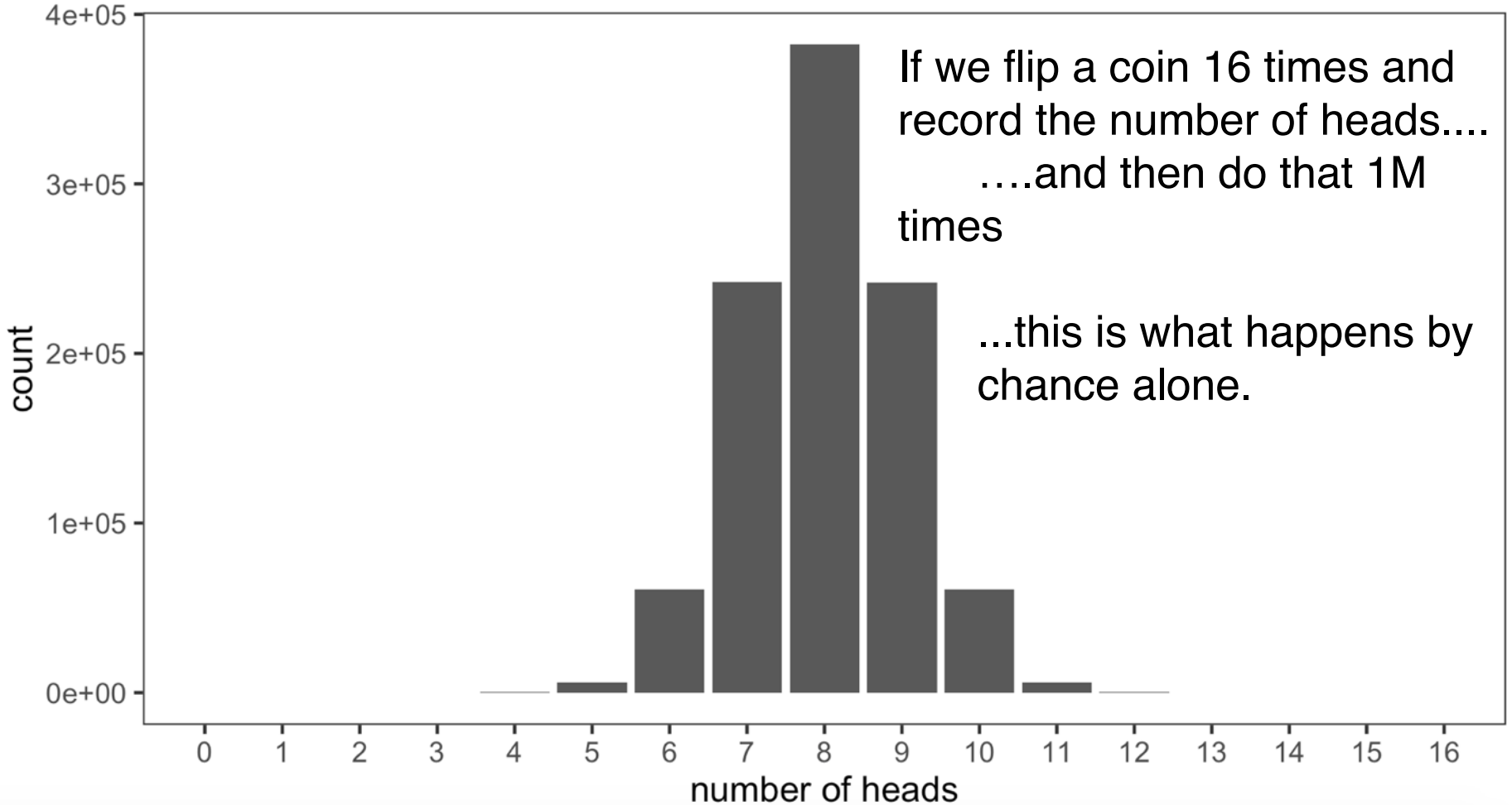


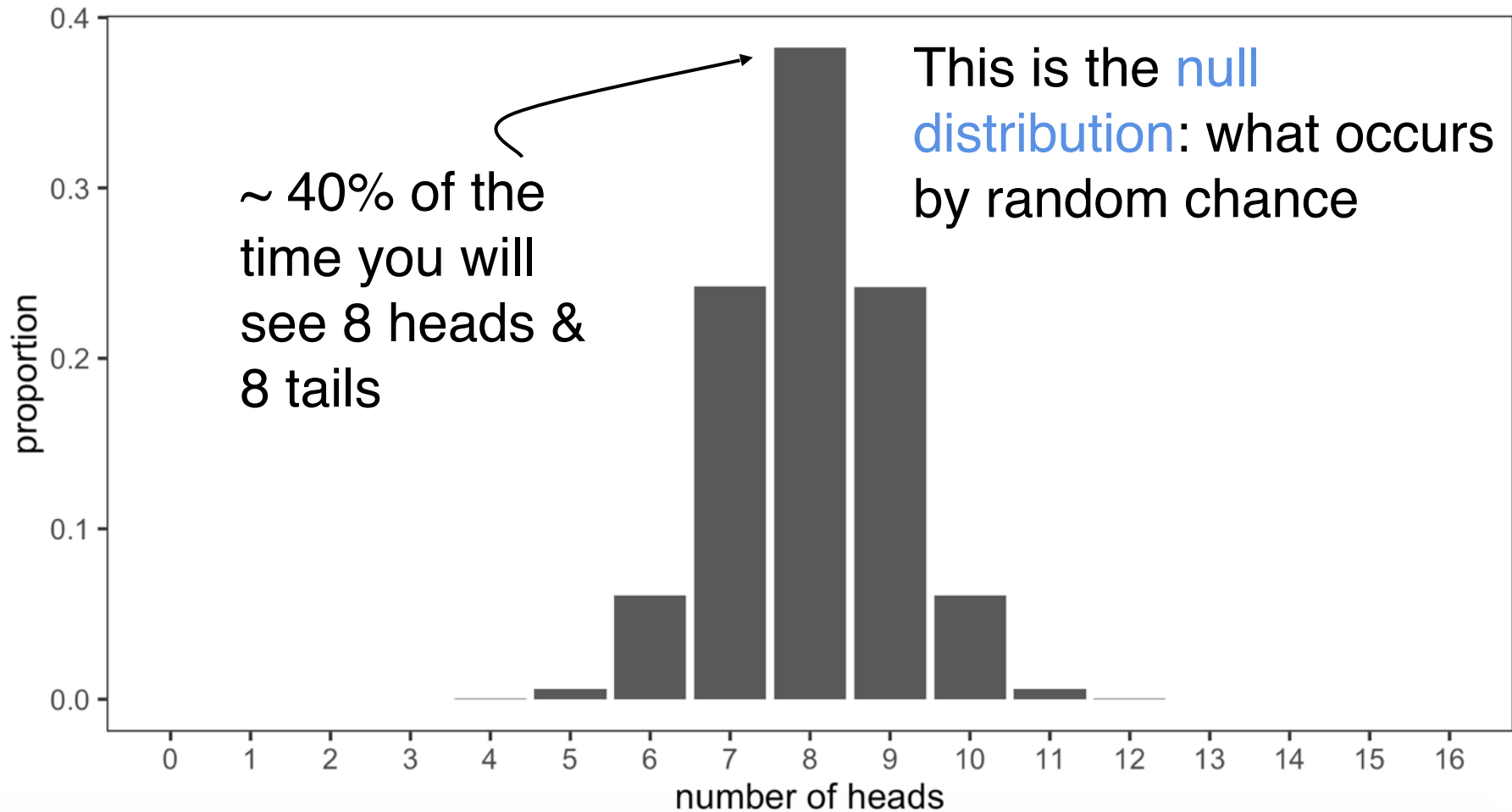
If there were a stronger effect of Poverty on Birth rate, what would β_1 be?

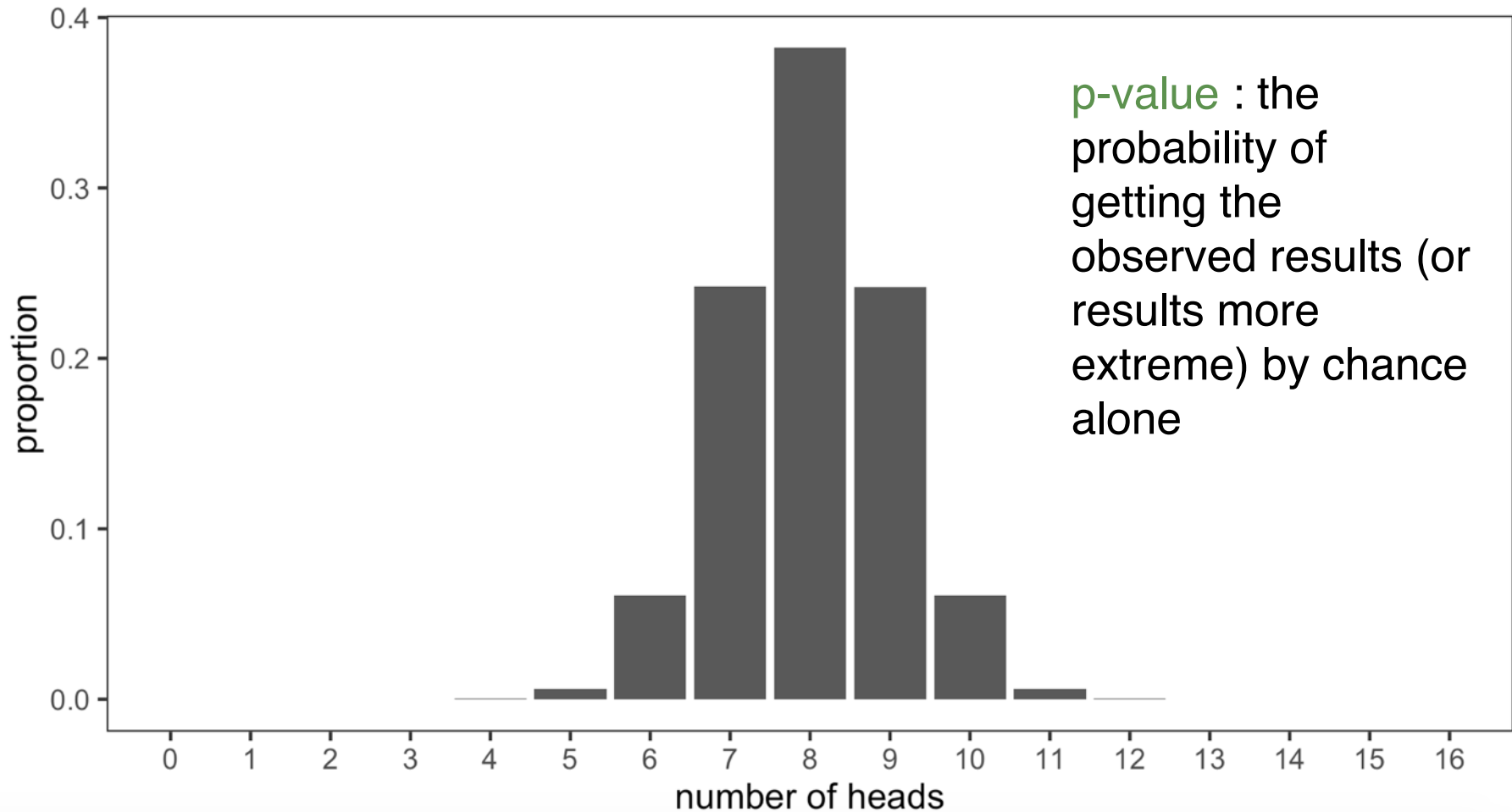


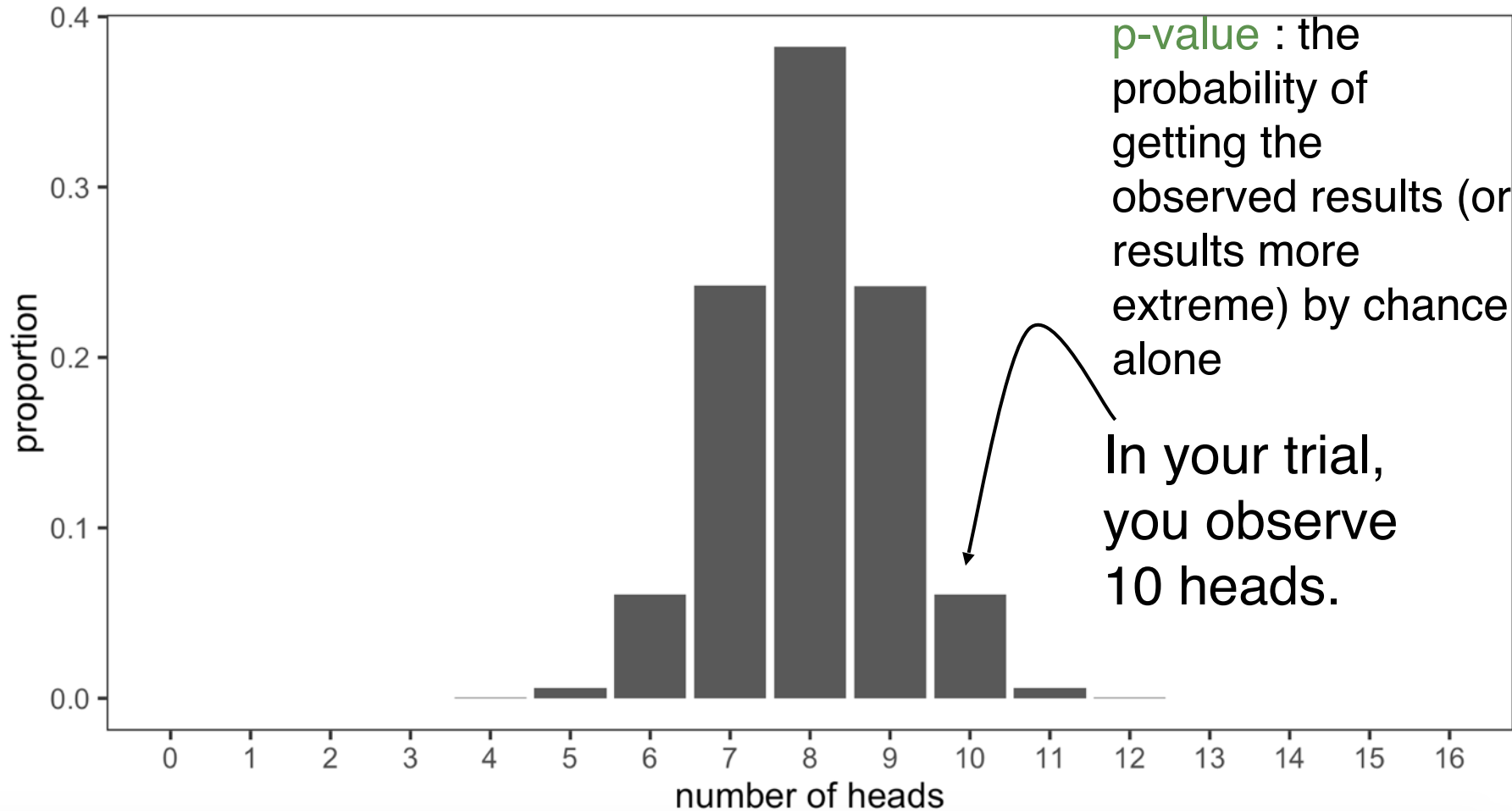


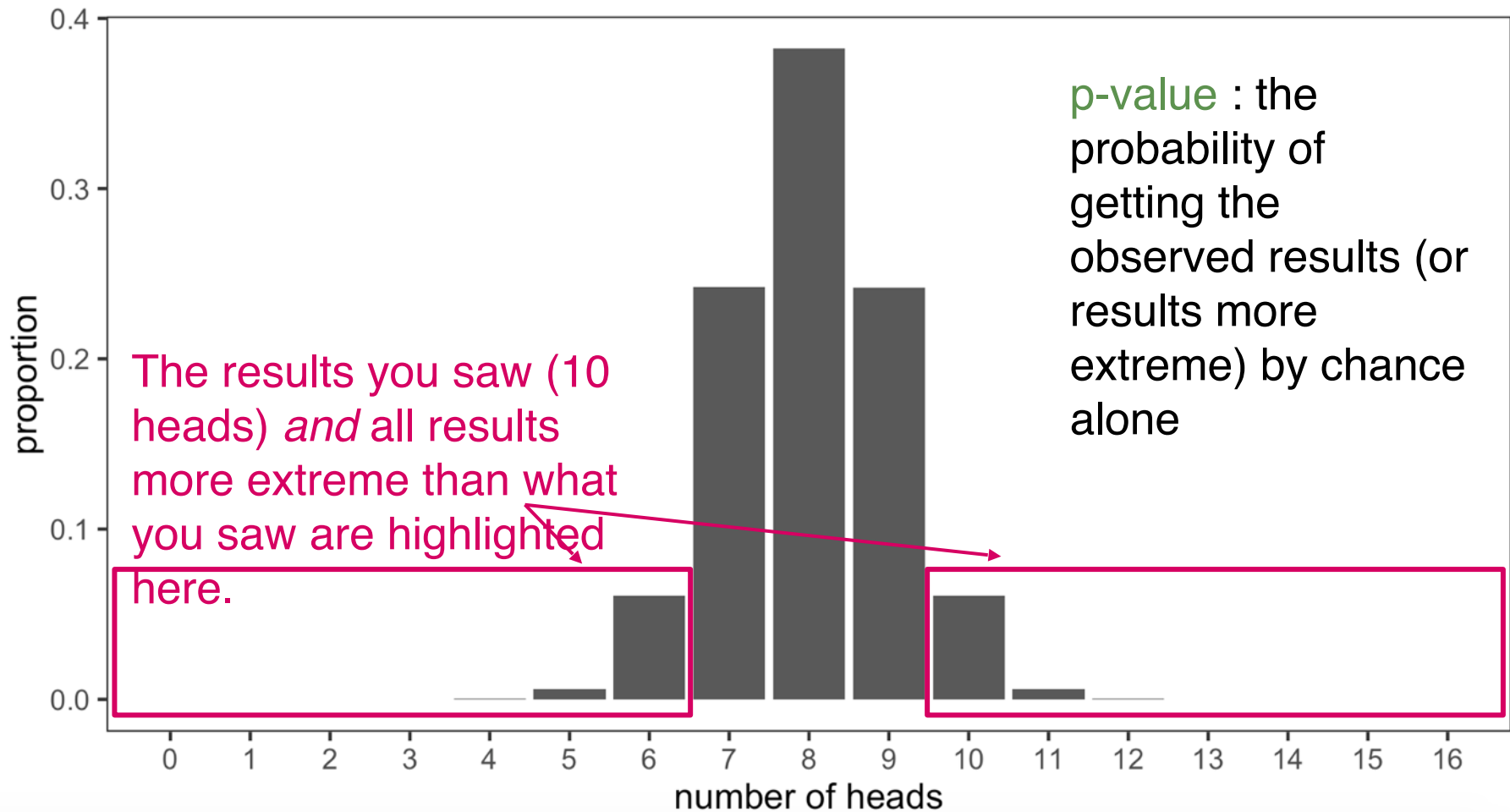
p-value : the probability of getting the observed results (or results more extreme) by chance alone

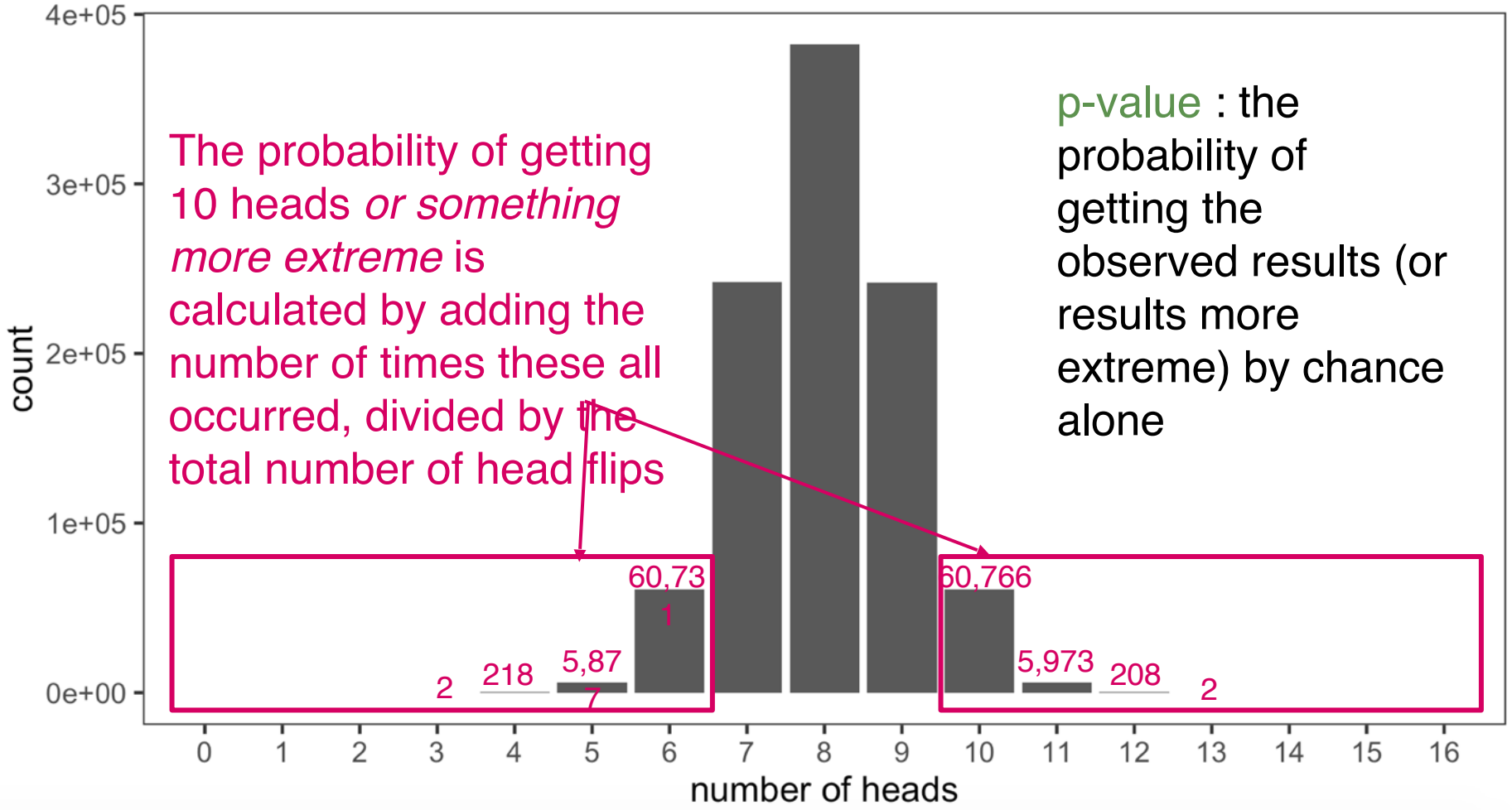






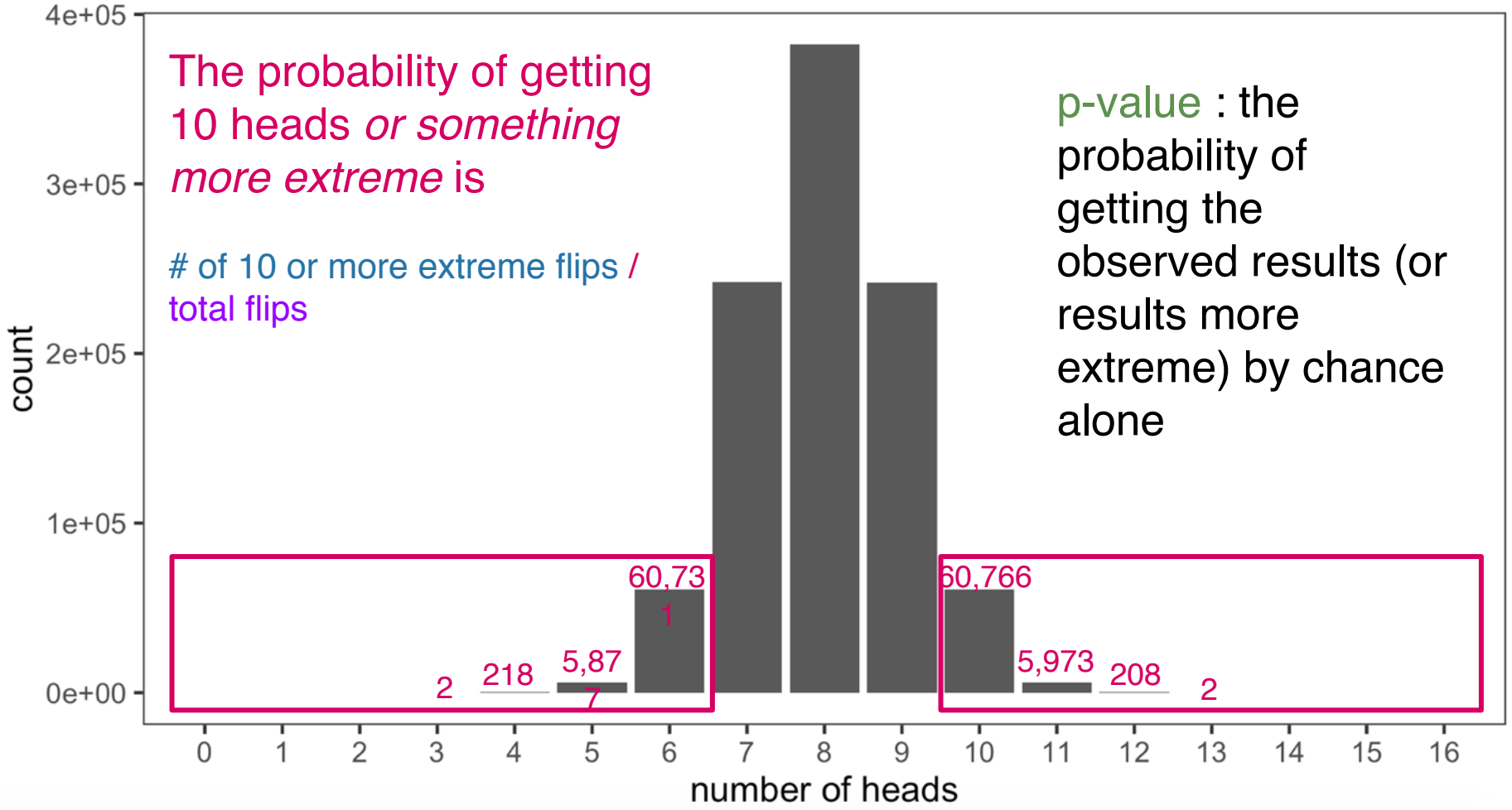


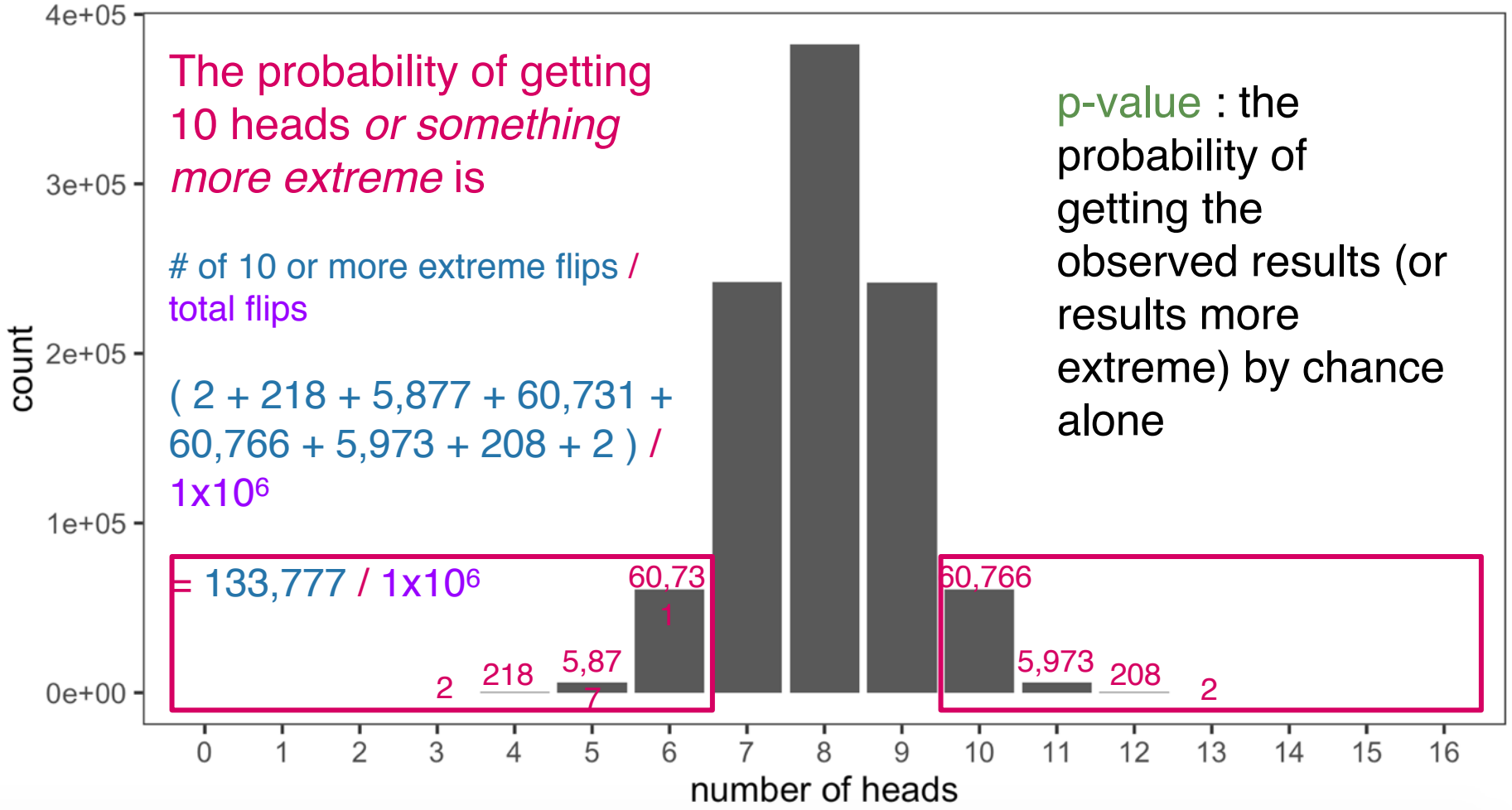


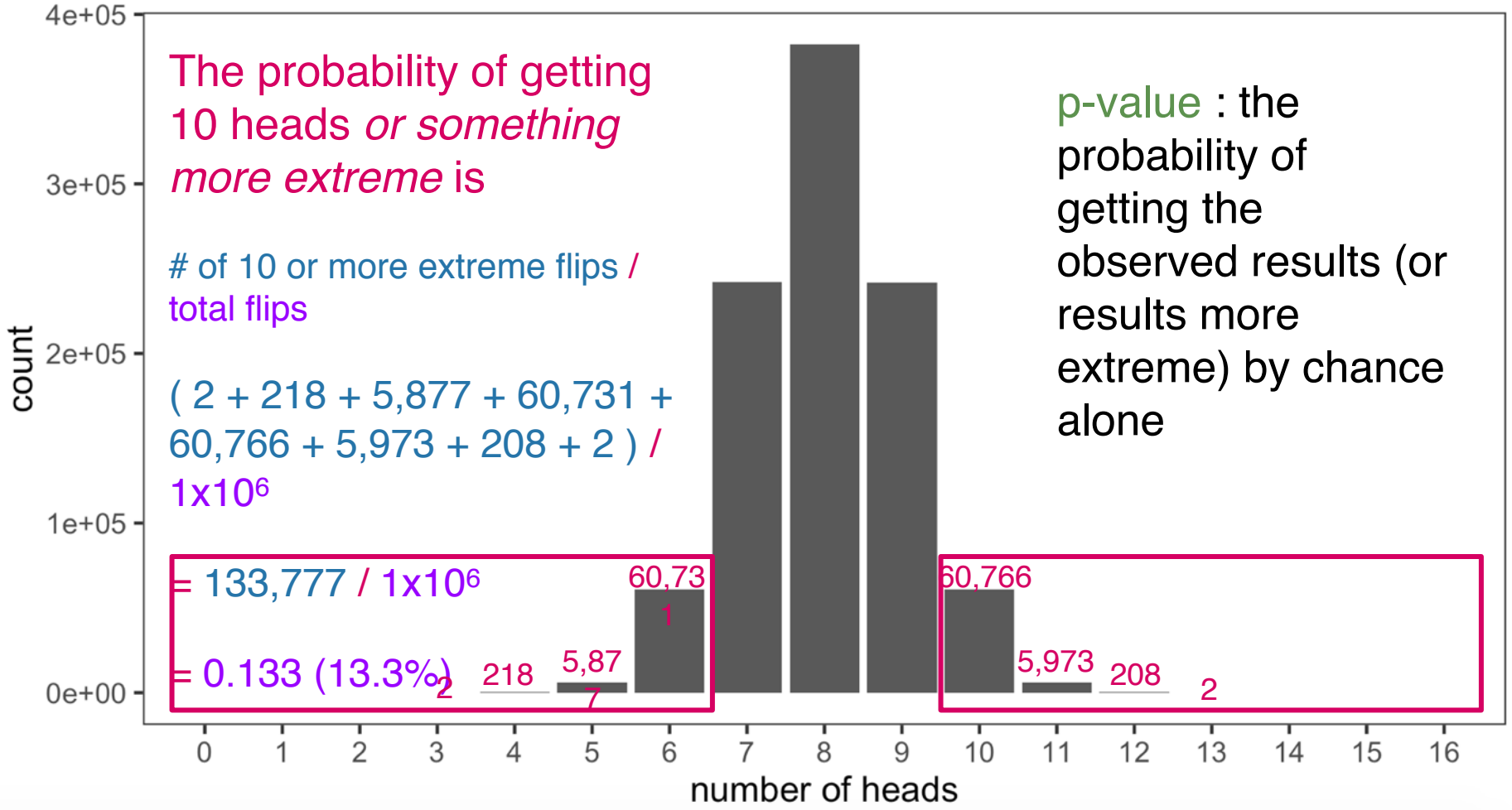


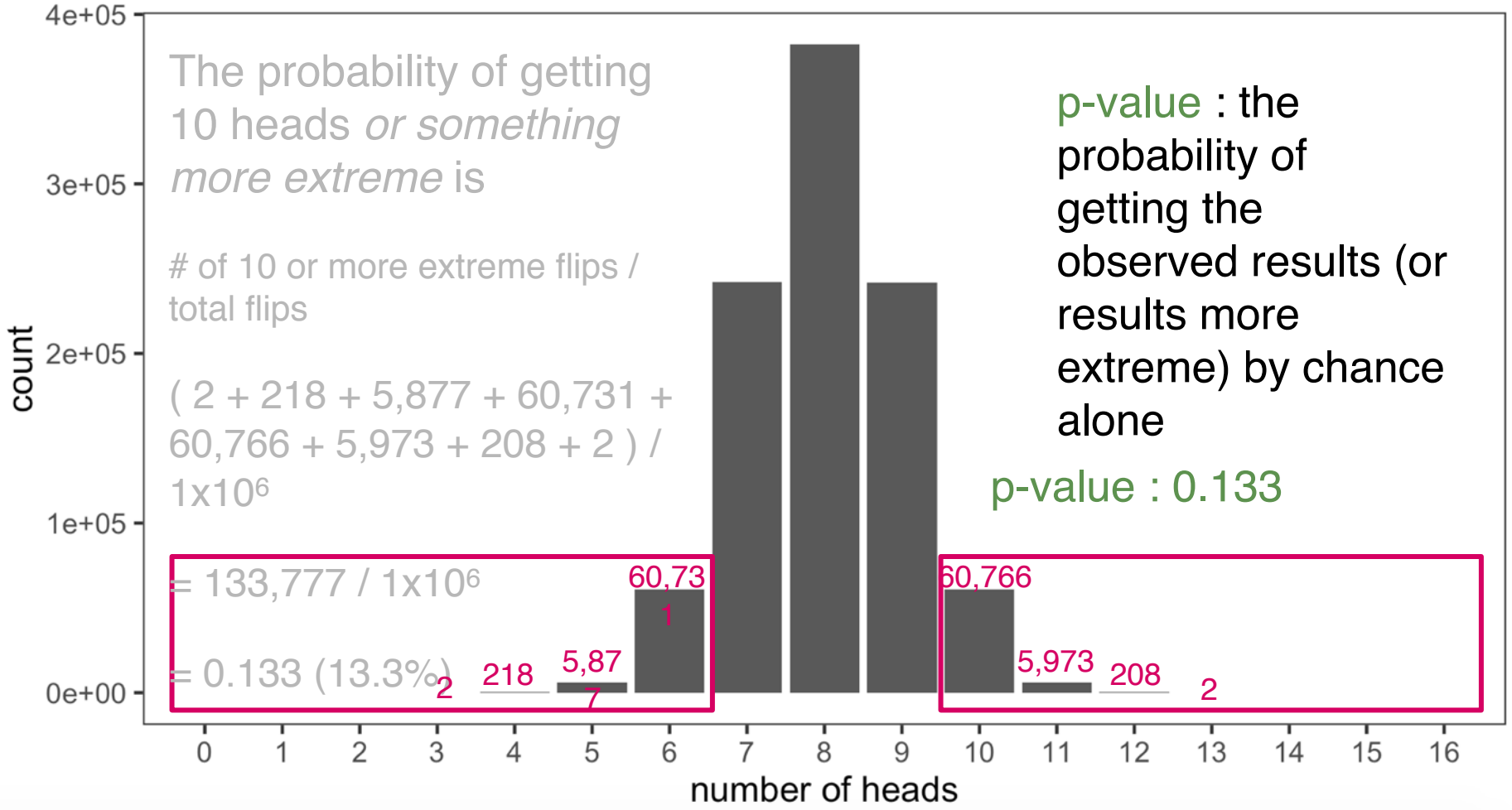
The probability of getting 10 heads *or something more extreme* is calculated by adding the number of times these all occurred, divided by the total number of head flips

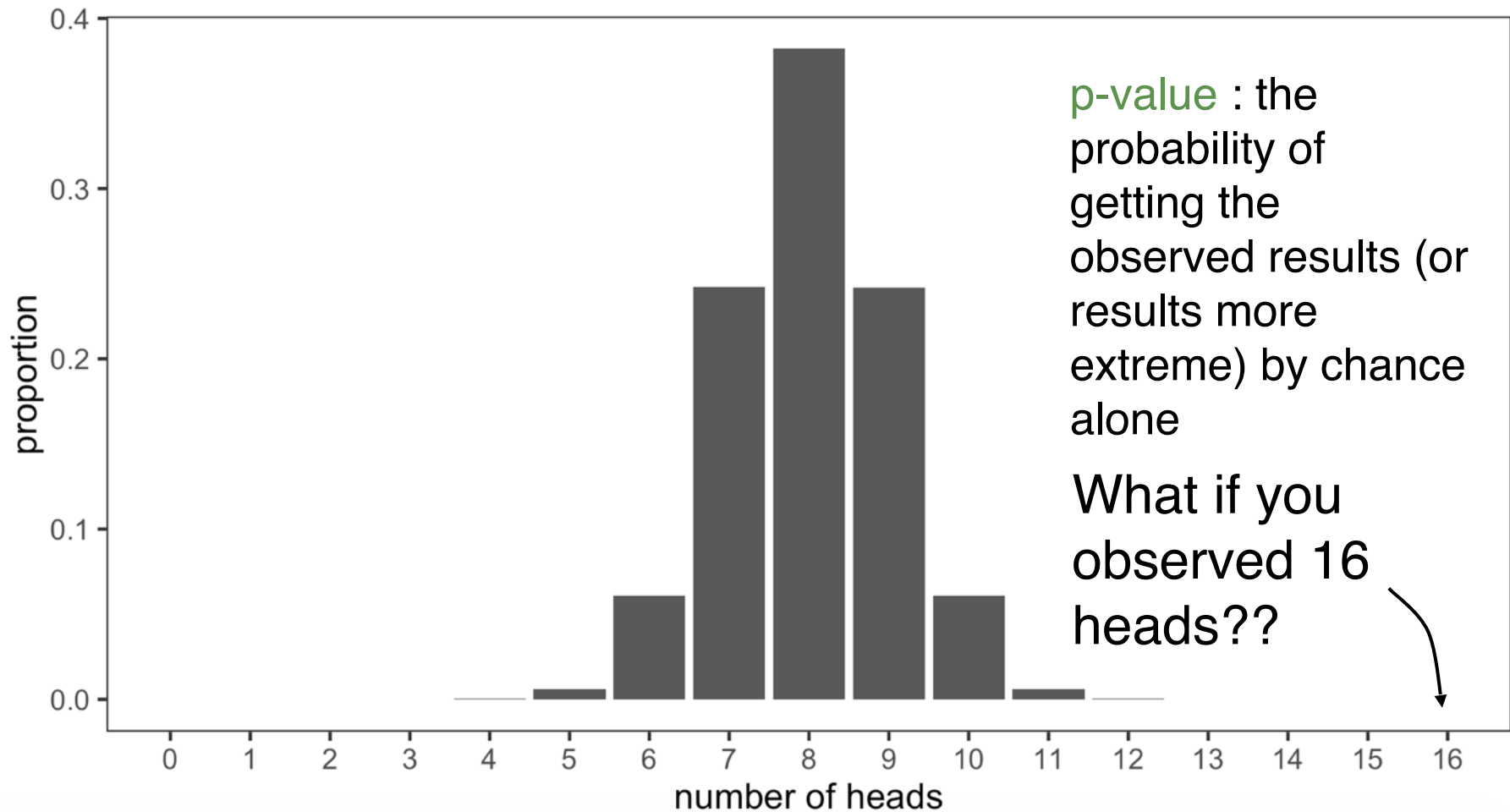
p-value : the probability of getting the observed results (or results more extreme) by chance alone

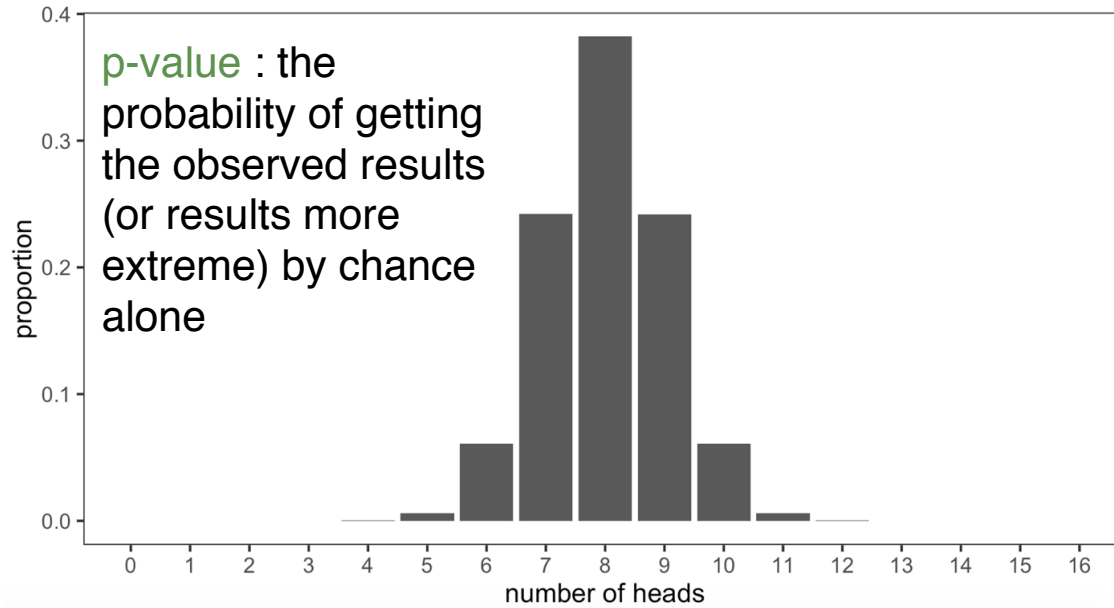






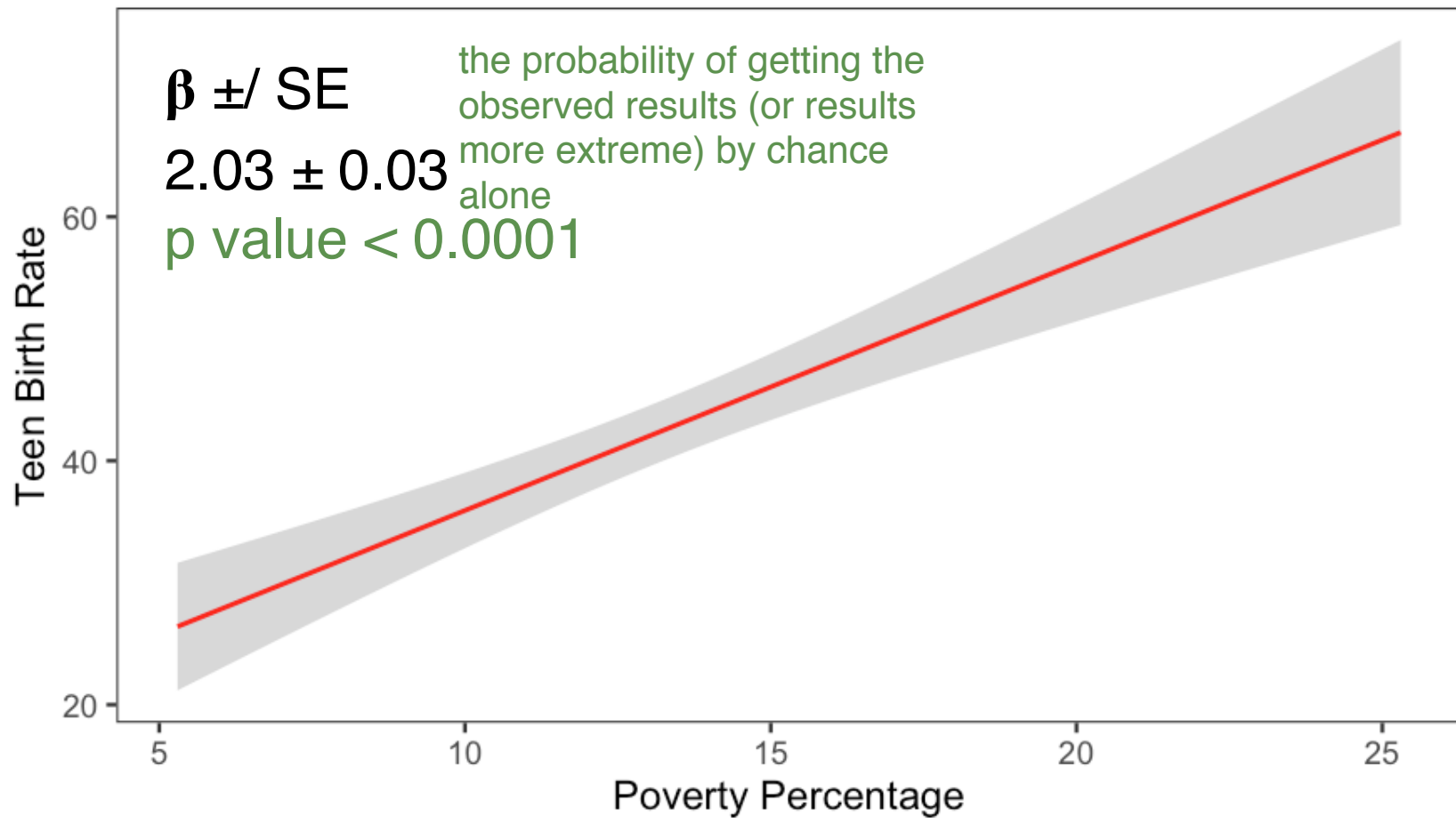







What would be the p-value of you flipping 16 heads?






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

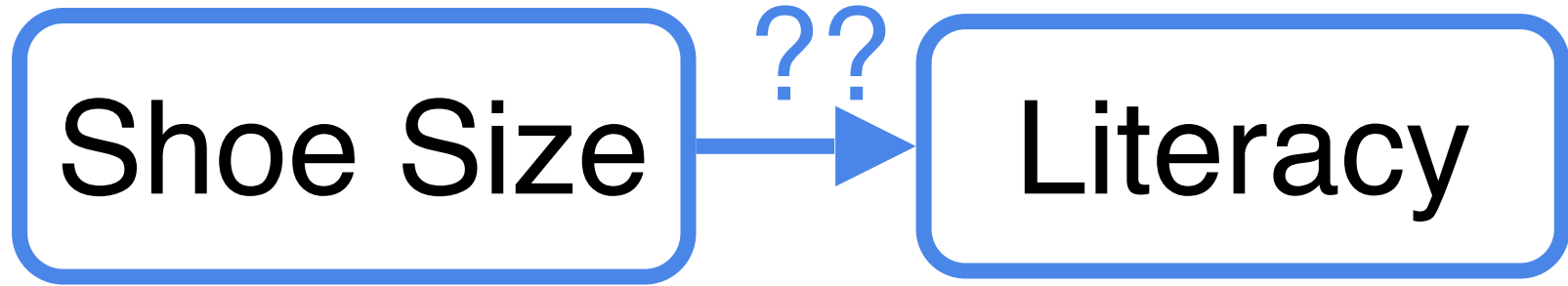
Confounding





Small shoes
Not literate

Big shoes
Literate





Small shoes
Not literate
Child

Big shoes
Literate
Adult

Shoe
Size

Literacy

Age

Variable1

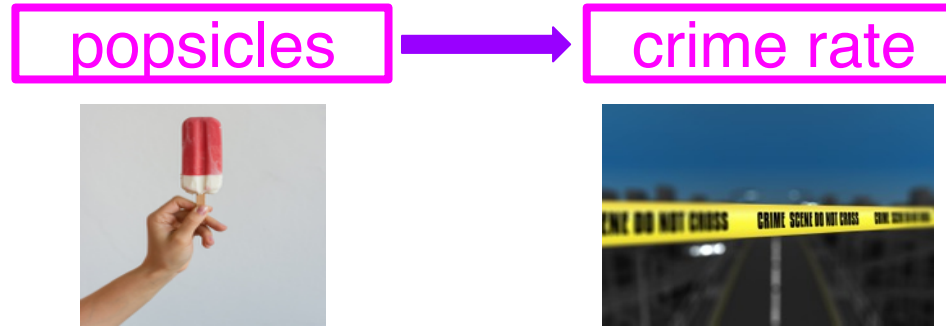
Variable2

Confounder

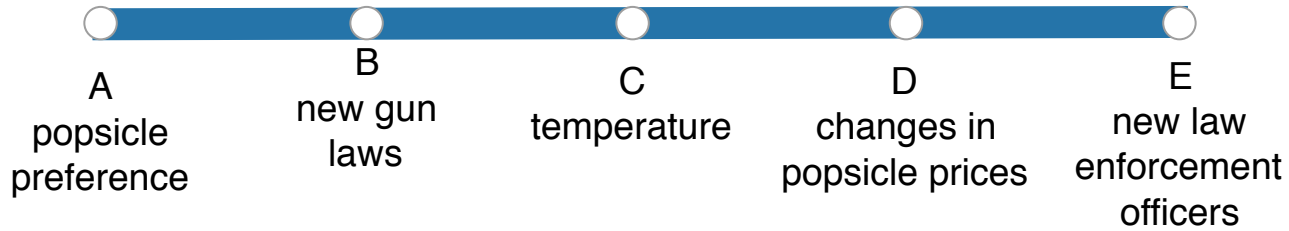
```
graph BT; C[Confounder] -.-> V1[Variable1]; C -.-> V2[Variable2];
```

The diagram illustrates a causal relationship where a confounder, represented by a dashed blue box at the bottom, influences two variables, Variable1 and Variable2, which are shown in solid blue boxes above. Two blue arrows point from the top of the dashed box to the bottom of each variable box, indicating that the confounder is a common cause for both variables.

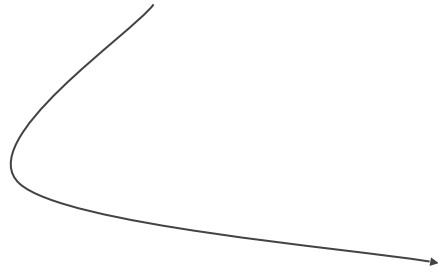
Confounding



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



We'll discuss additional approaches of how to account for confounding in your analysis in another 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	Conservative care
	N = 200	N = 200
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 quickly

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

Confounding



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

