

# COGS 109: Lecture 9



Hypothesis testing

July 20, 2023

***Modeling and Data Analysis***

Summer Session 1, 2023

C. Alex Simpkins Jr., Ph.D.

RDPRobotics LLC | Dept. of CogSci, UCSD

# Plan for today

- Announcements
- Inference using central tendency and variability concepts
- Hypothesis testing
- Confounds

# Plan for today II

- Hypothesis testing
- Introduction to models and the modeling process
- Colormaps - custom

# Announcements

- Q2 posted
- D4 released
- Github repos
- In process on feedback and A1

# Update: the big picture

- Where we are
  - 5 parts of the course
    - We discussed data
      - What is it, how do we manipulate it, import it, python and some matlab implementation
      - Filtering
      - Computing basic statistics
    - We discussed basic visualization
      - Plotting data (2d, histograms, scatterplots, etc)

# Update: the big picture (II)

- Where we're going
  - We will now cover
    - Modeling
      - what is modeling?
      - interpolation, approximation, extrapolation
    - Error analysis
      - How good is your model?

# Update: the big picture (III)

- Where we're going (continued)
  - **What we're going to cover**
    - Basic models
      - **Linear fits, nonlinear fits**
      - **Regression**
      - **Relationship to machine learning**
      - **Interpolation/extrapolation (also data analysis methods)**
    - Advanced models and modeling methods
      - **Fitting models with optimization methods**
      - **Artificial neural networks**
      - **AI**
  - Communicating results
    - **This has been integrated and will continue to be integrated**
    - **Proper forms of inserting figures and tables in scientific communications**
    - **Format in homeworks is designed to teach proper communication methodology**

Extending central tendency and  
variability to inference and  
hypothesis testing



## 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

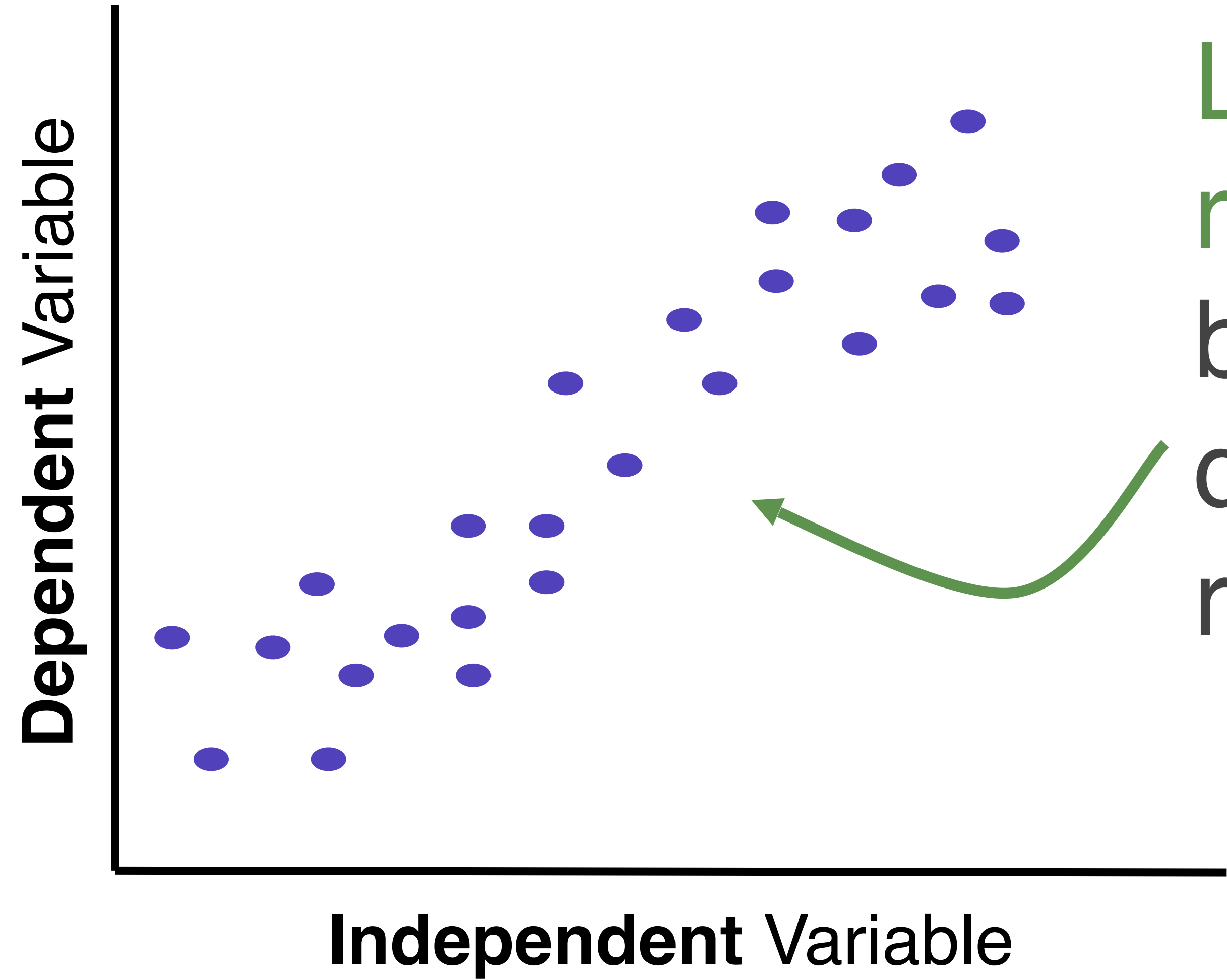
## REGRESSION

DOES CHANGE IN ONE  
VARIABLE MEAN  
CHANGE IN ANOTHER?

I.e. simple  
regression, multiple  
regression

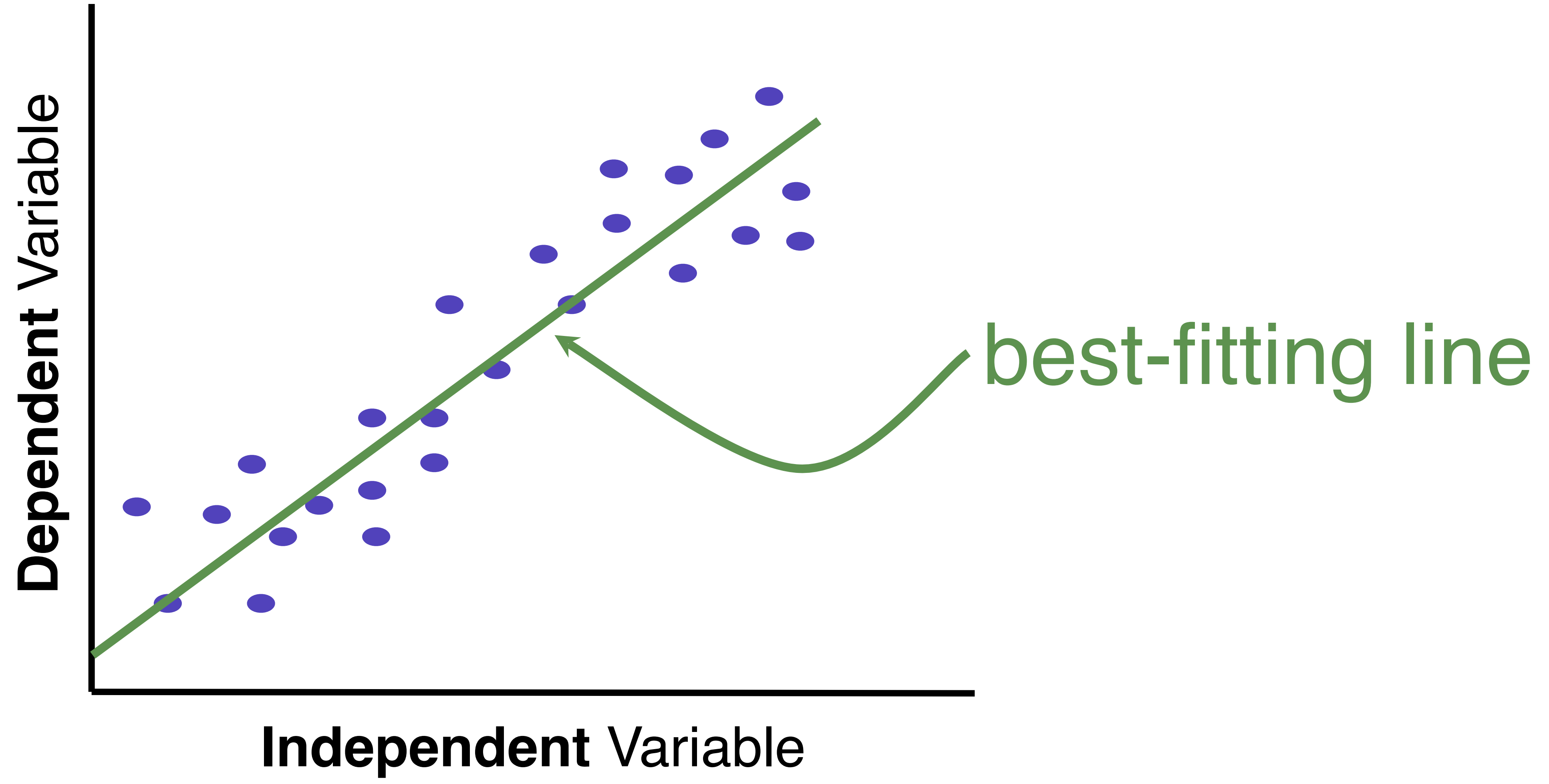
## NON-PARAMETRIC TESTS

FOR WHEN  
ASSUMPTIONS IN  
THESE OTHER 3  
CATEGORIES ARE NOT  
MET  
Levene, Mann-Whitney  
sum test, Wilcoxon  
sign-rank test, sign  
test

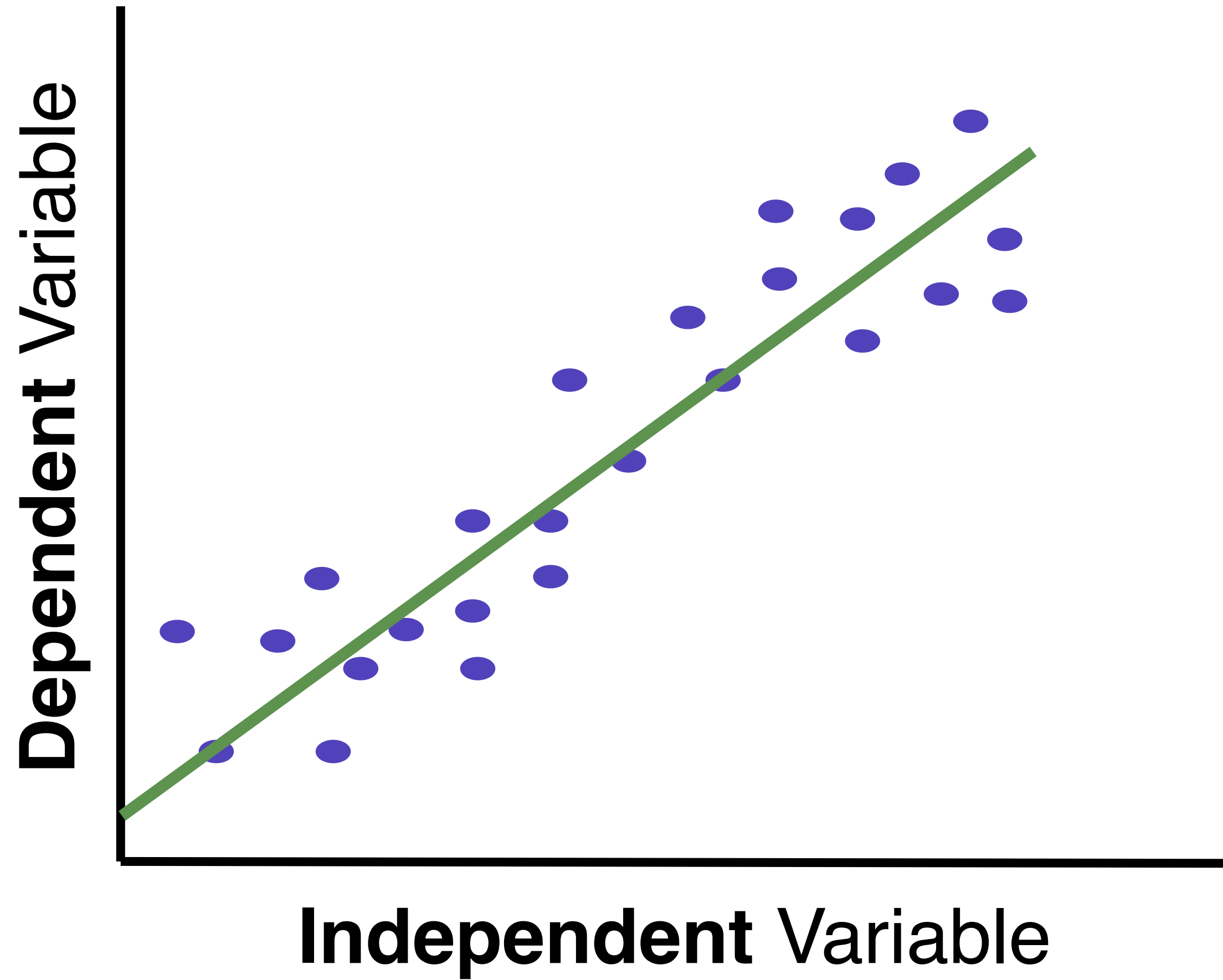


Linear regression can be used to describe this relationship

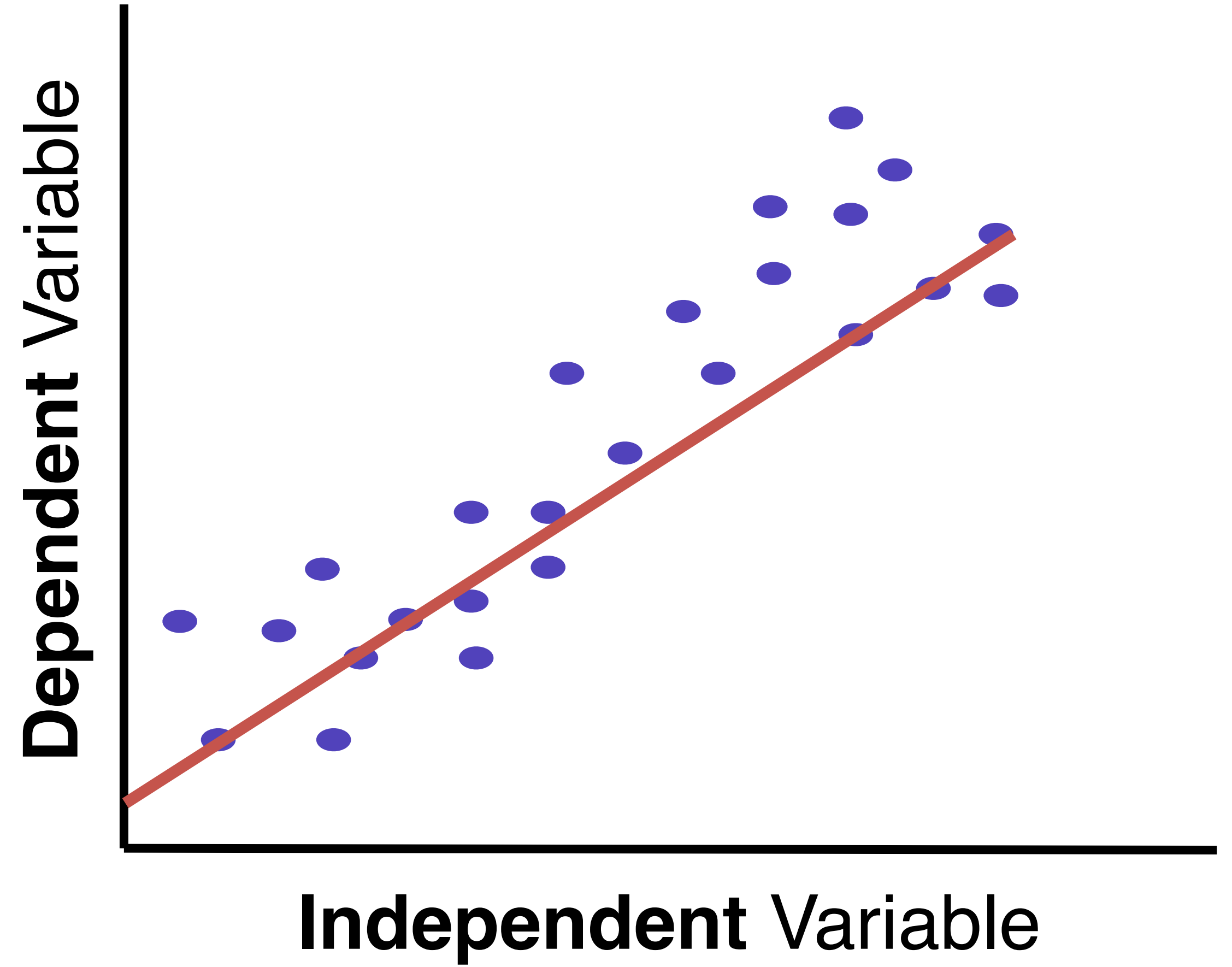


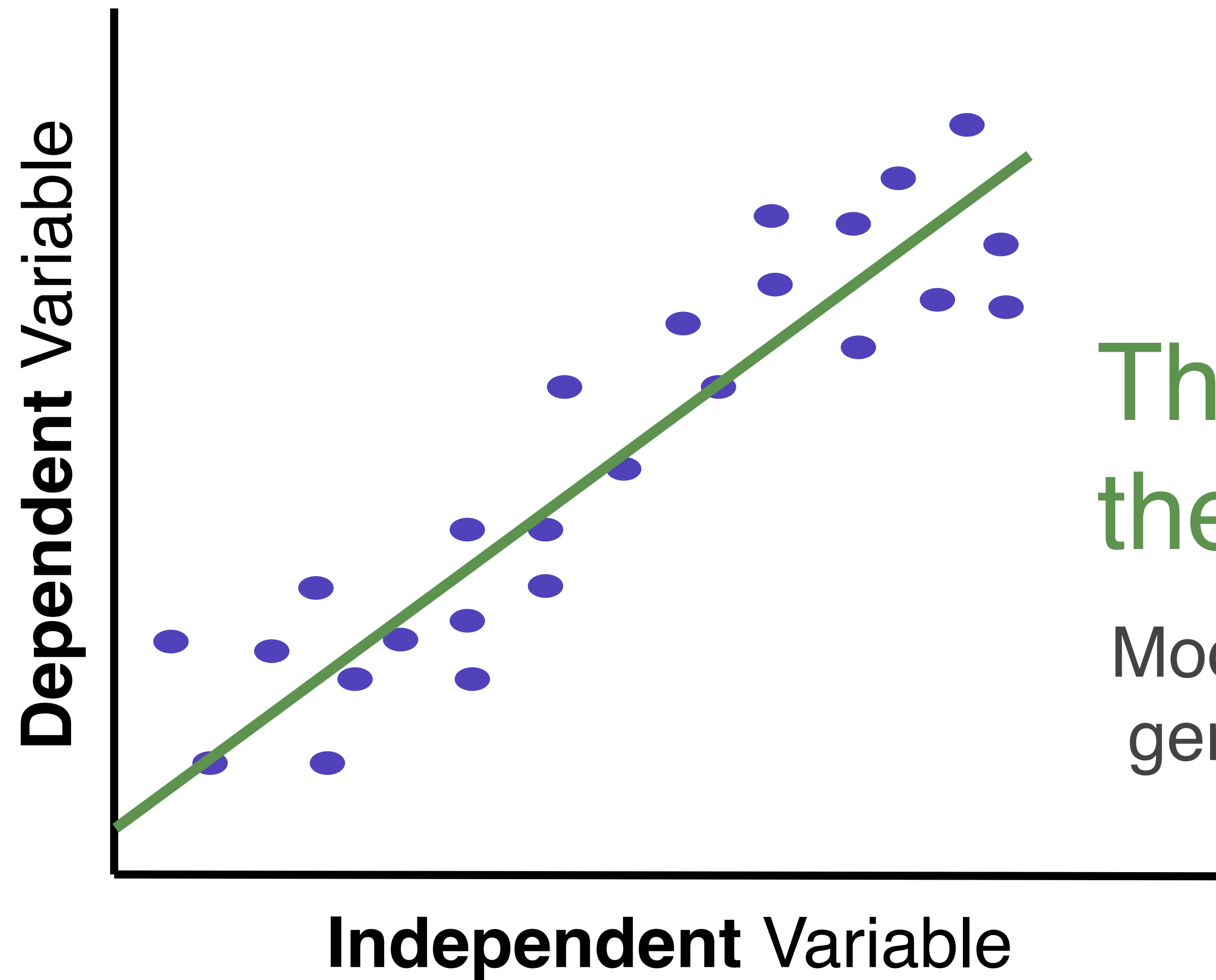


Best-fitting line



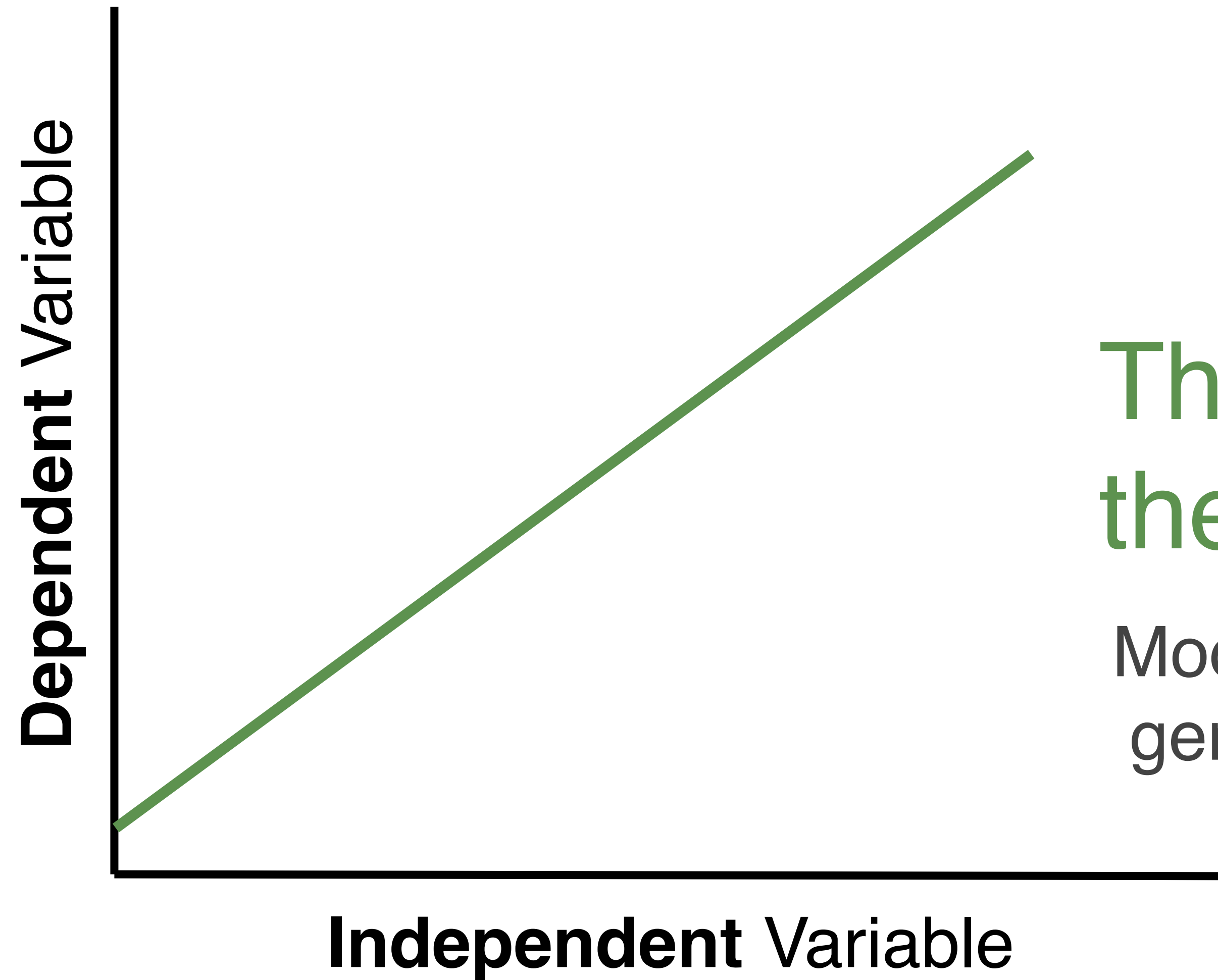
NOT a best-fitting line





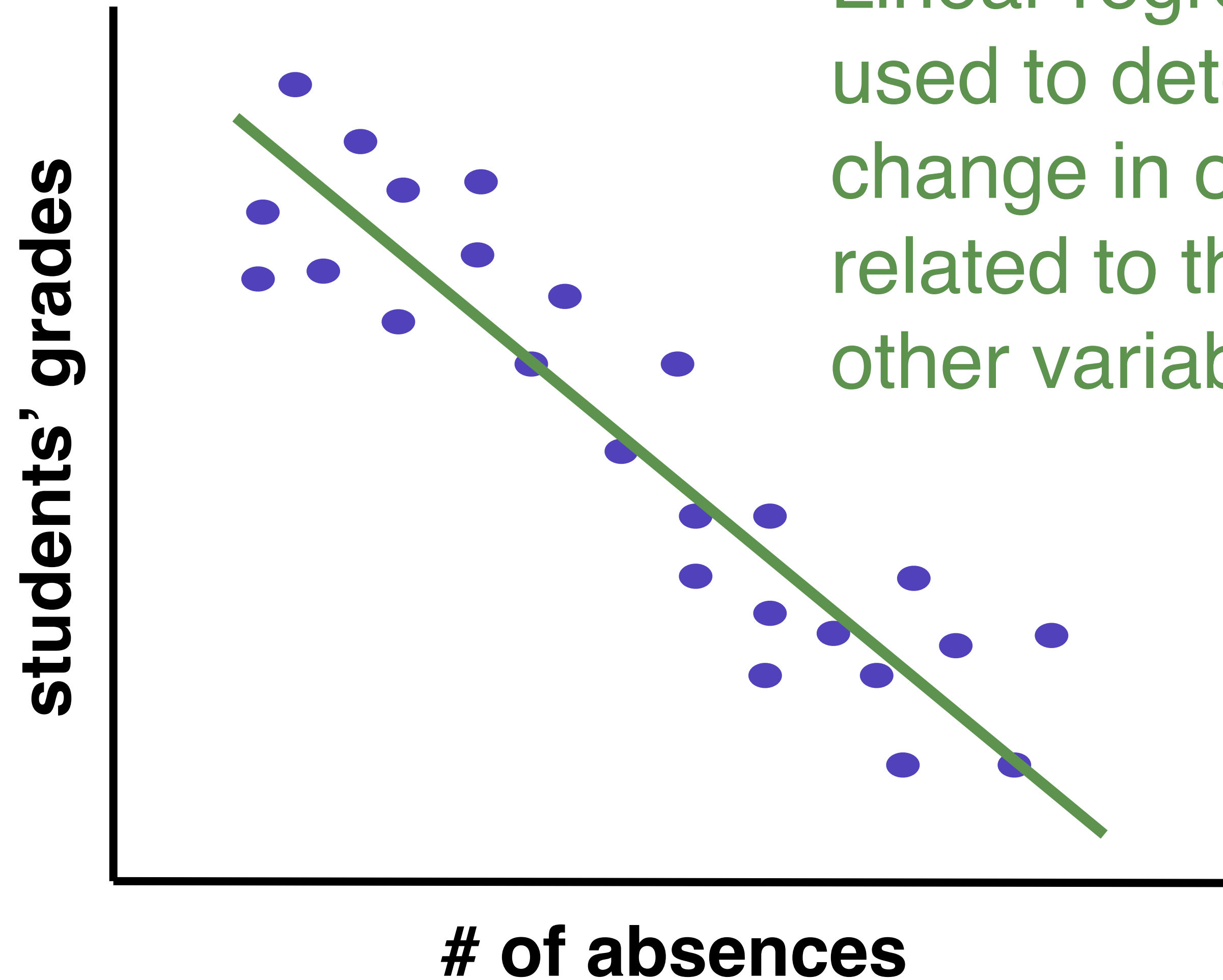
This line is a model of the data

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



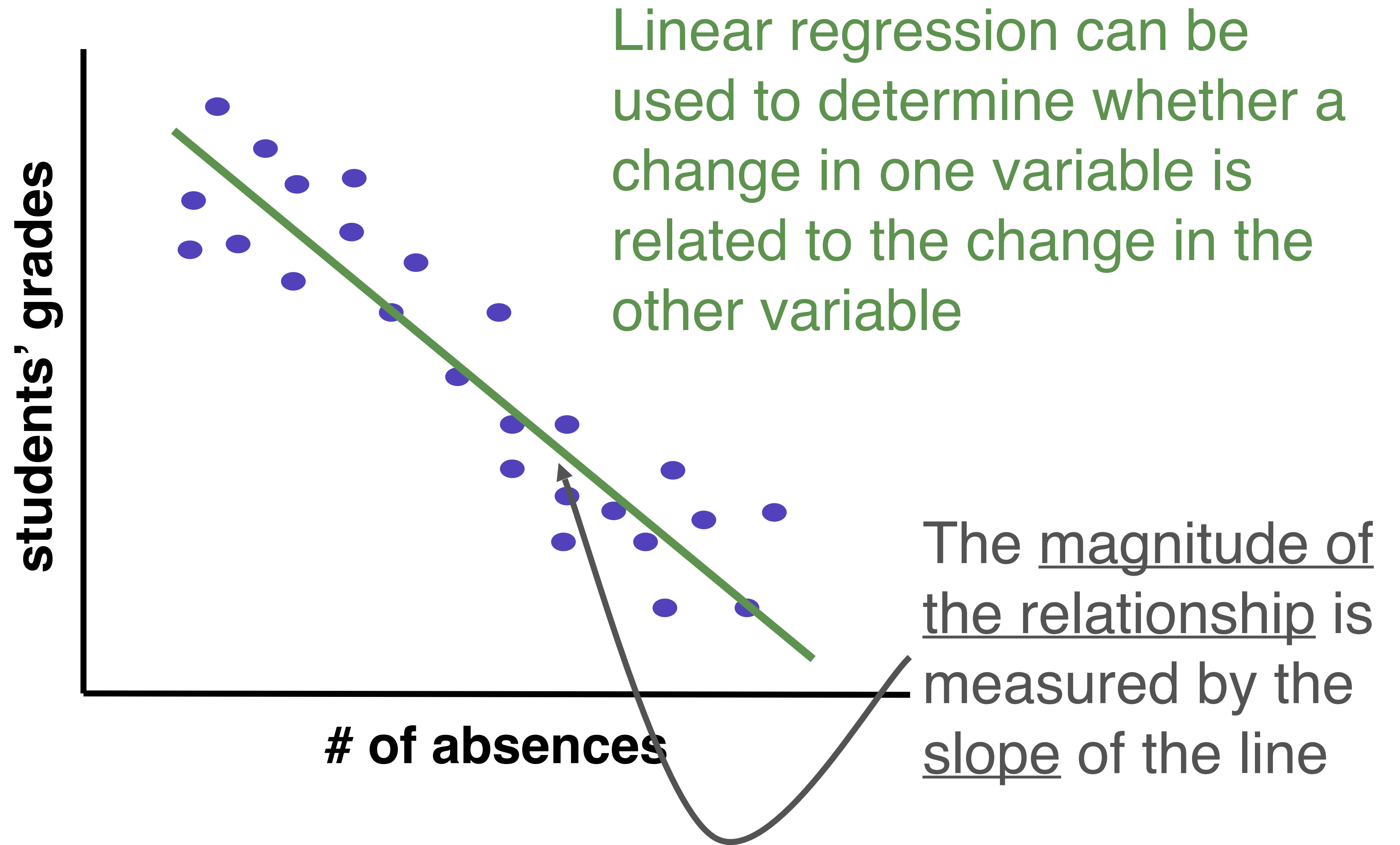
This line is a model of the data

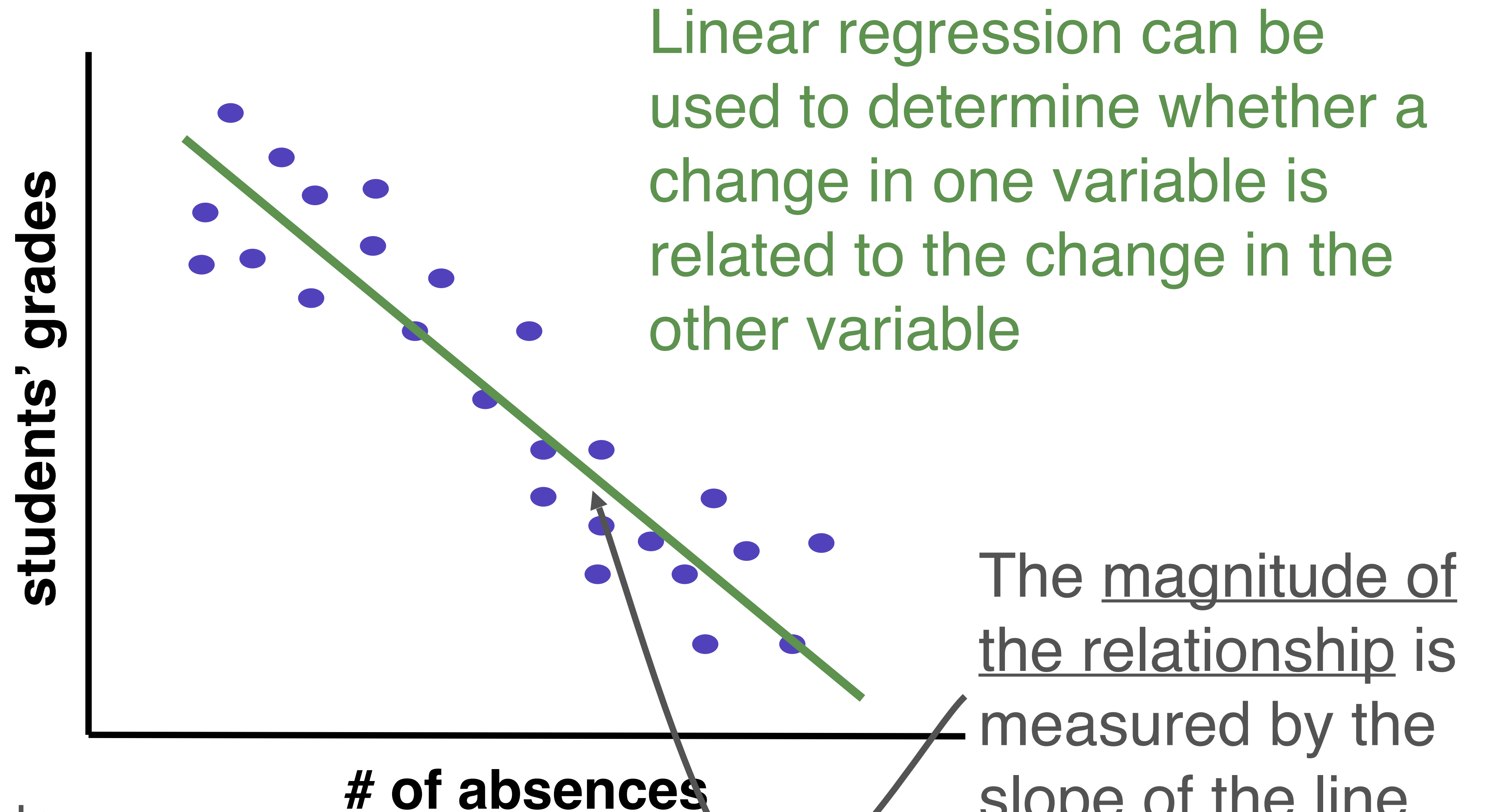
Models are mathematical equations generated to *represent* the real life situation



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

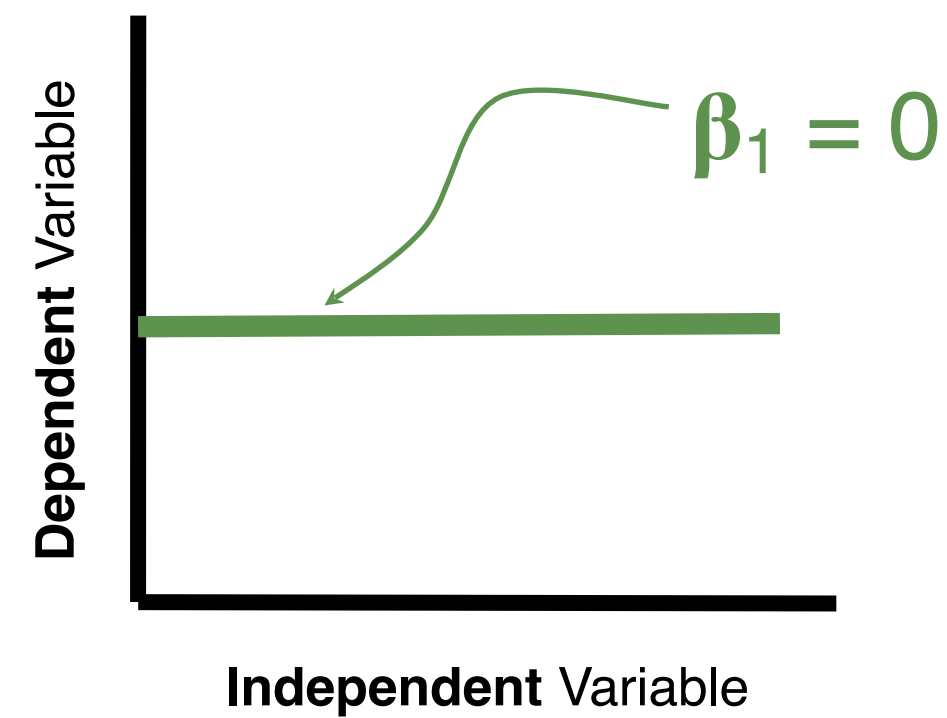




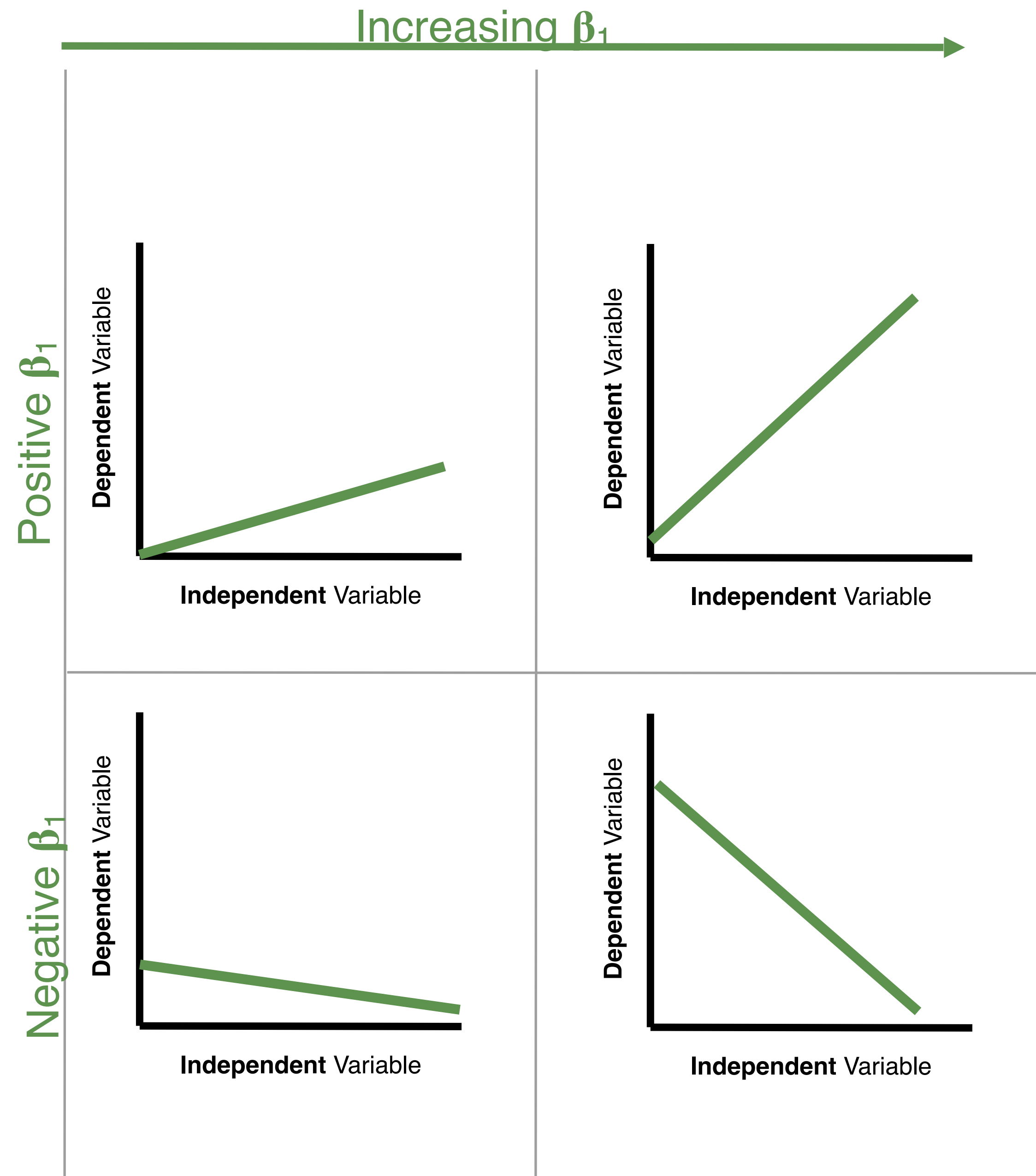
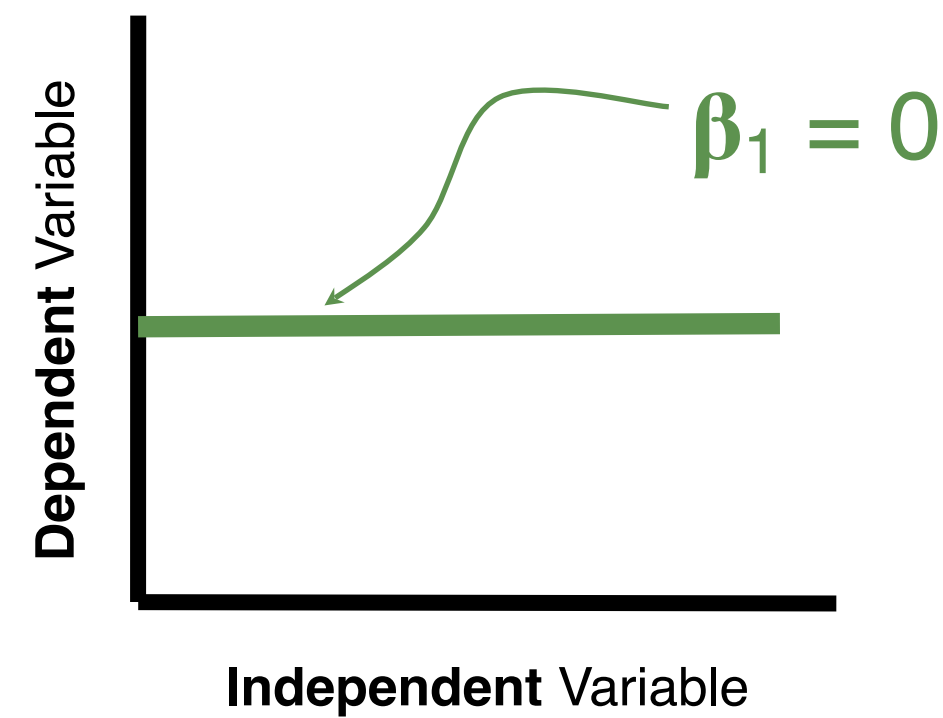


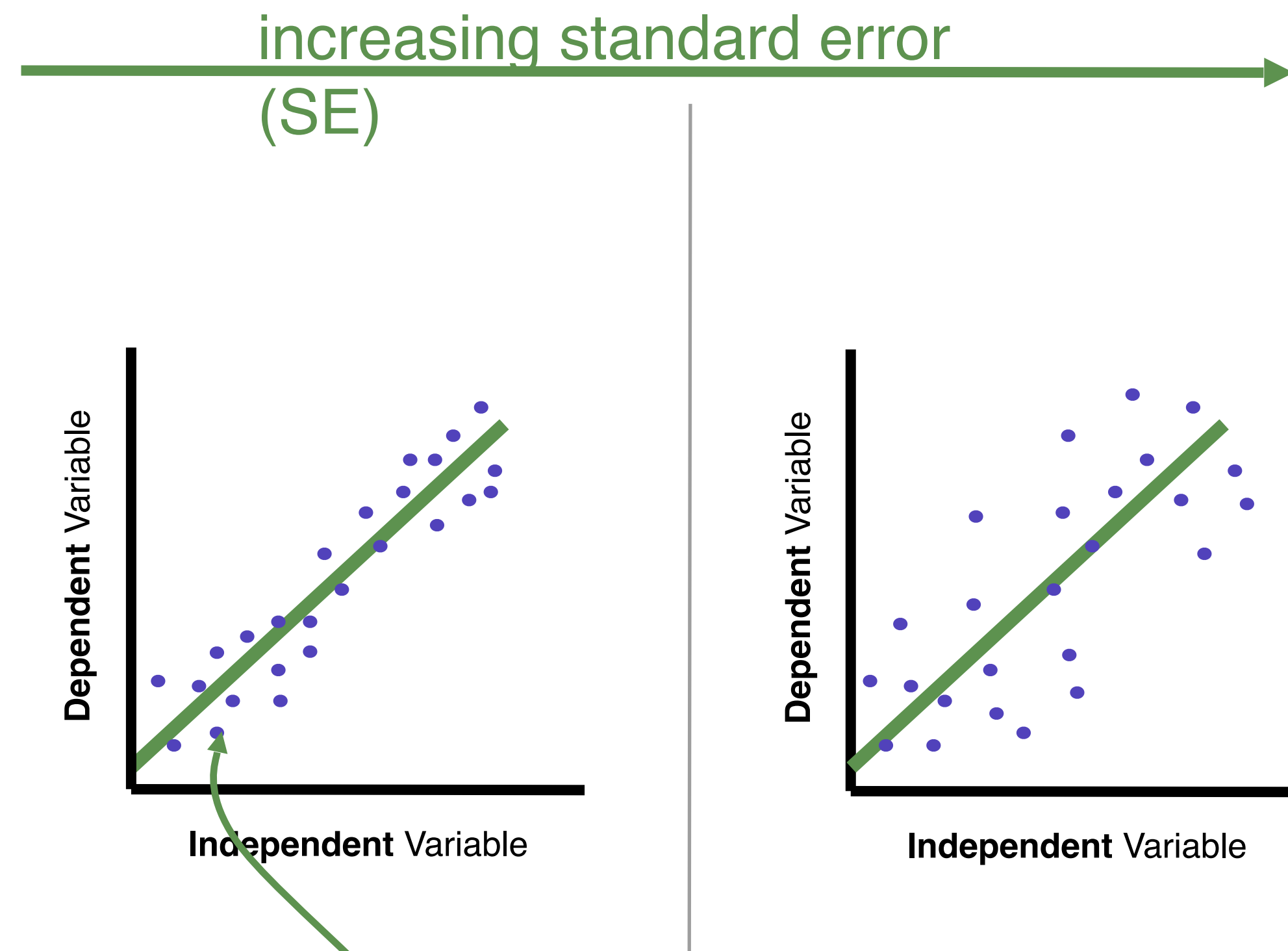
This is also referred to as the model's effect size ( $\beta_1$ )

Effect size ( $\beta_1$ ) can be estimated using the slope of the line



Effect size ( $\beta_1$ ) can be estimated using the slope of the line



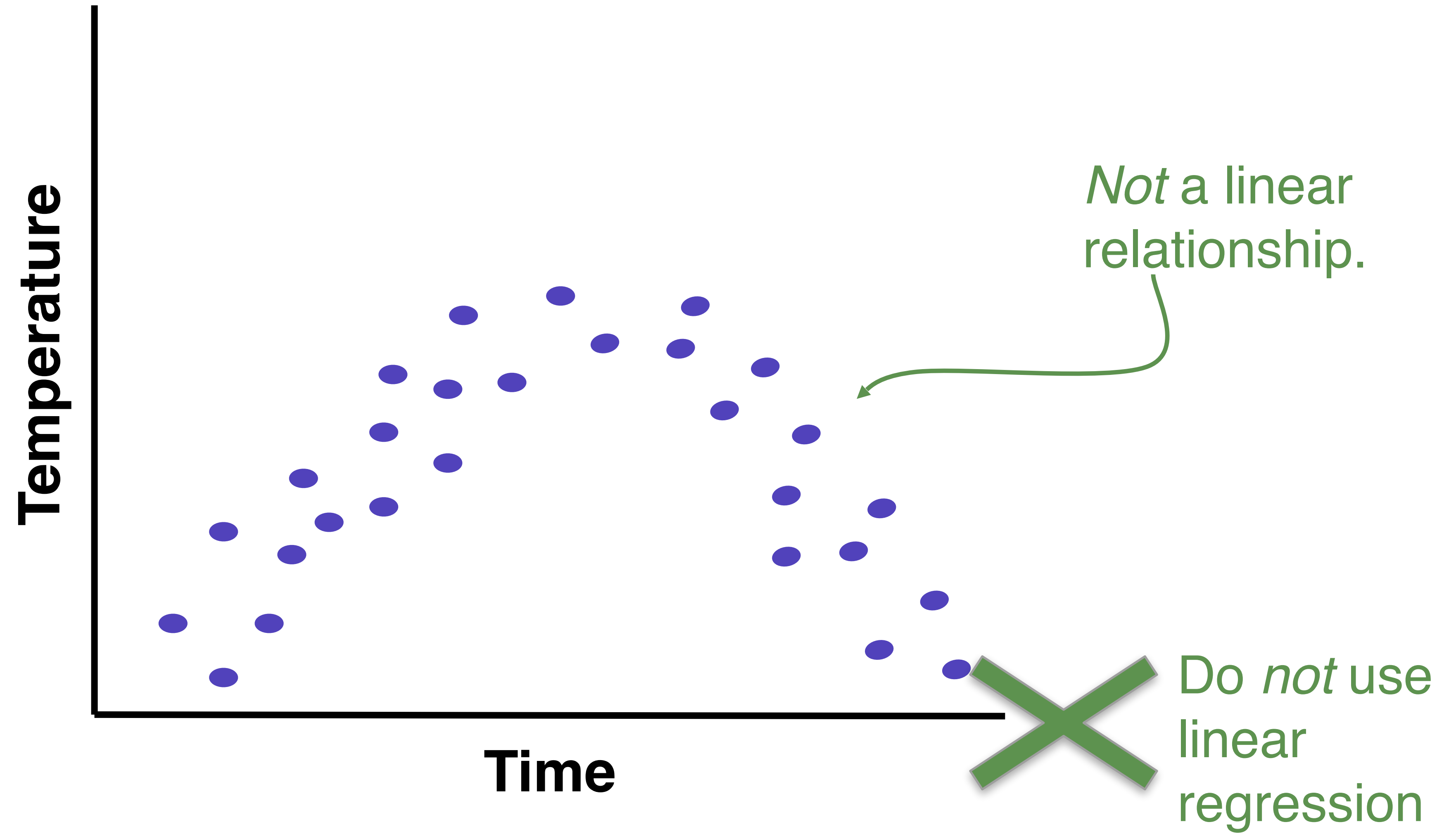


The *closer* the points are to the regression line, the *less uncertain* we are in our estimate

# 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



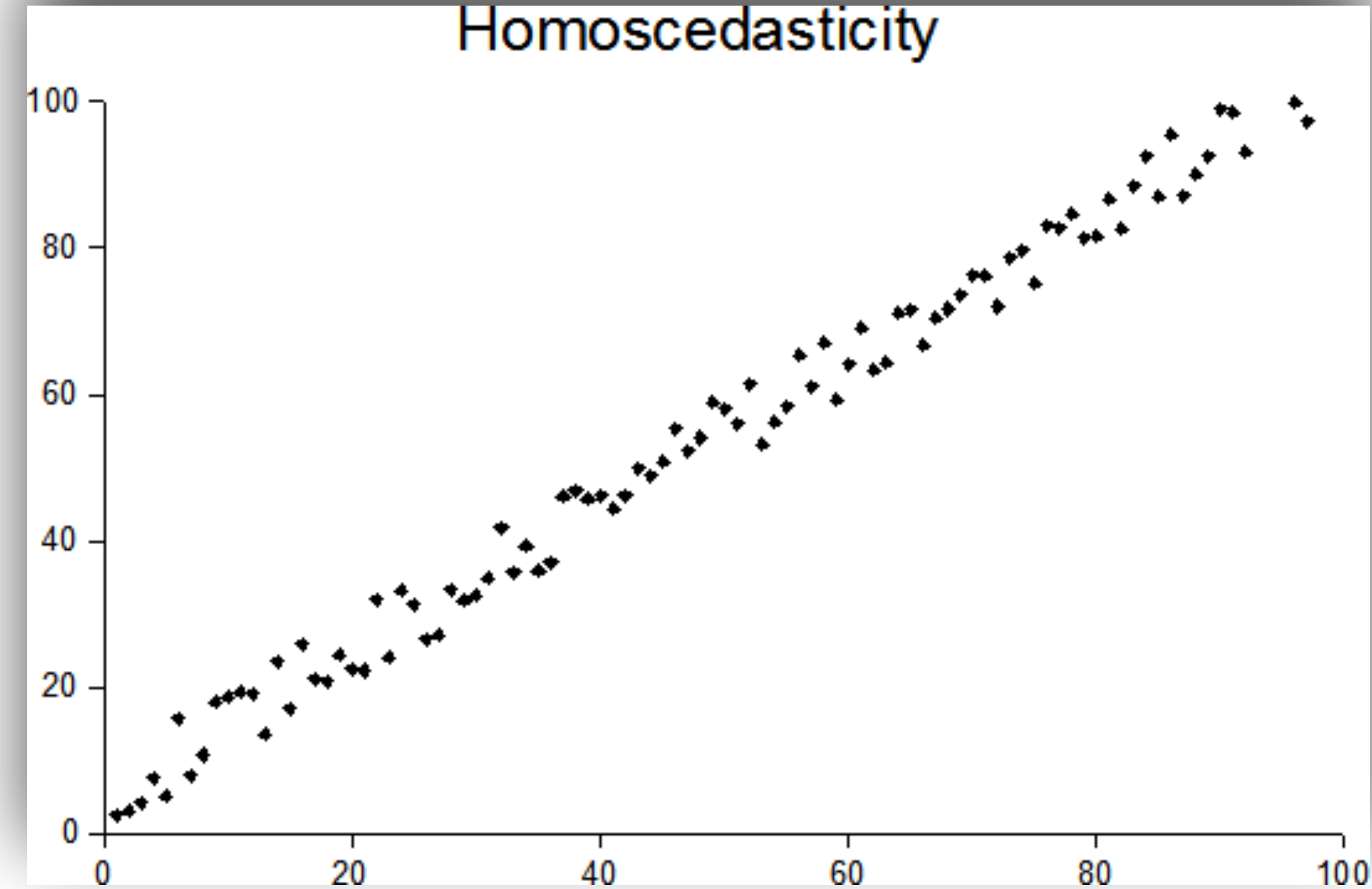
# Autocorrelation

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

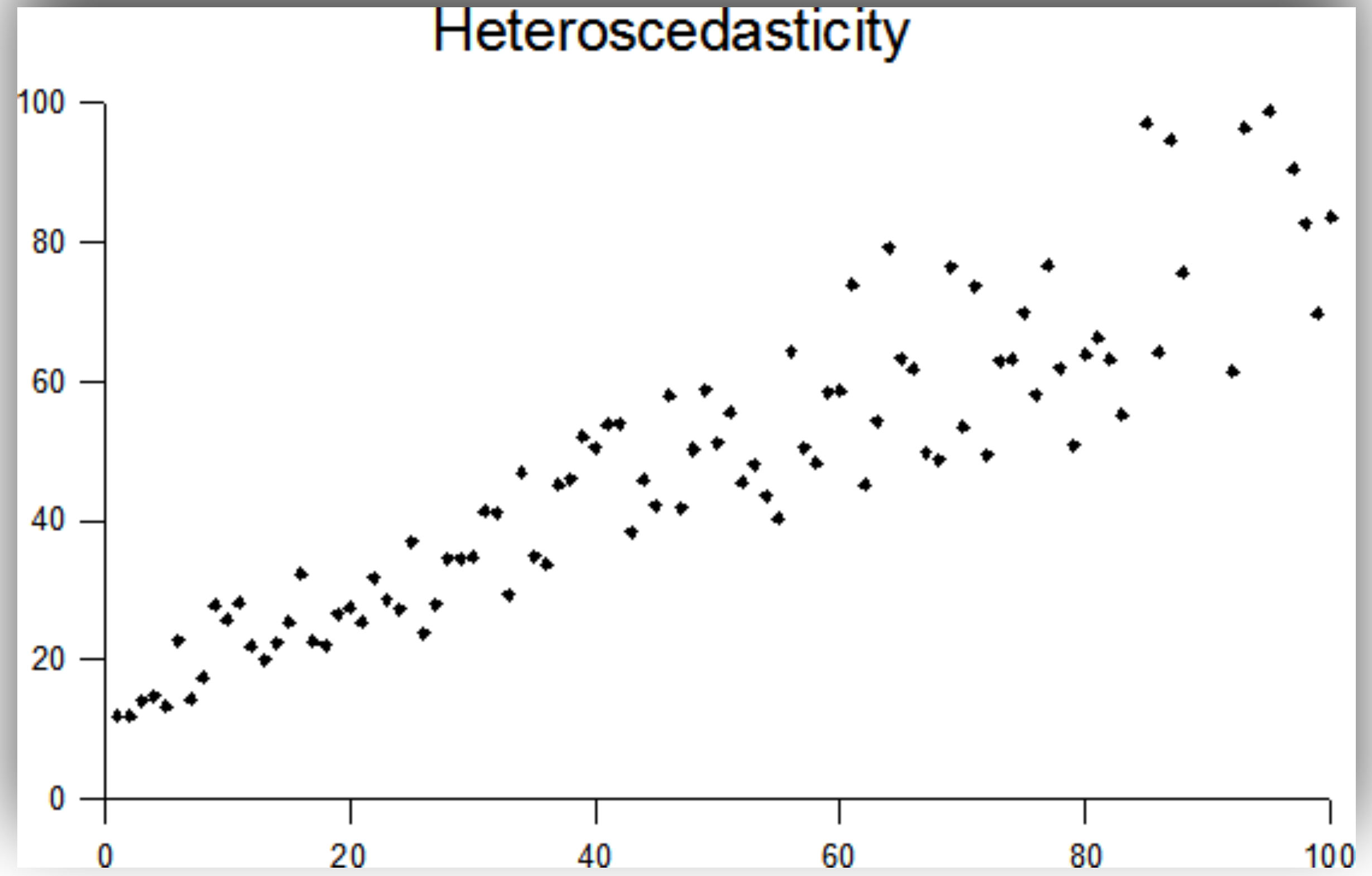


# Homoscedasticity - a reminder of what that is

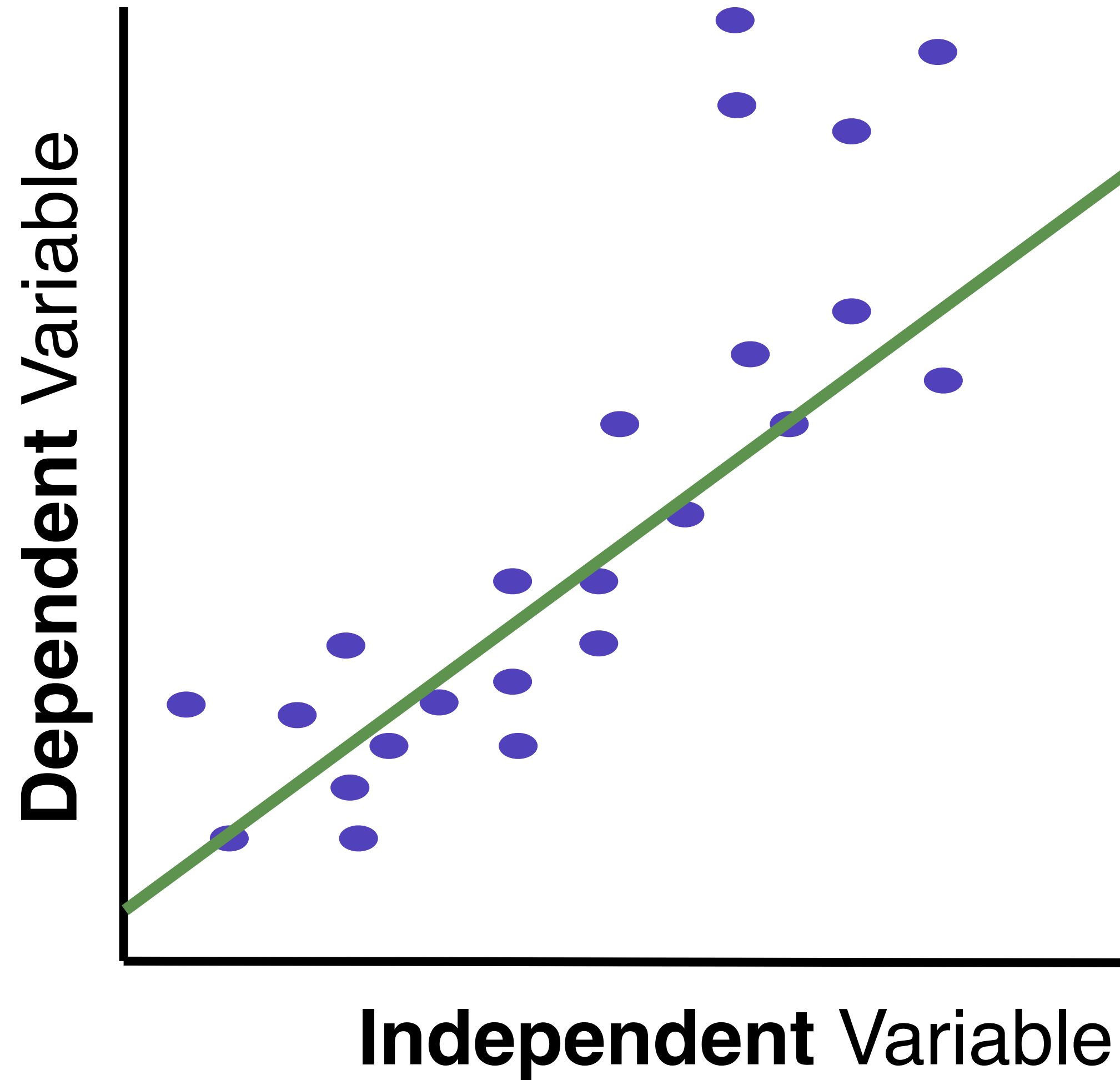
Homoscedasticity



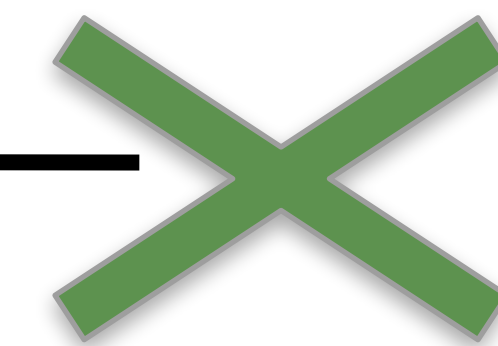
Heteroscedasticity



# Homoscedasticity

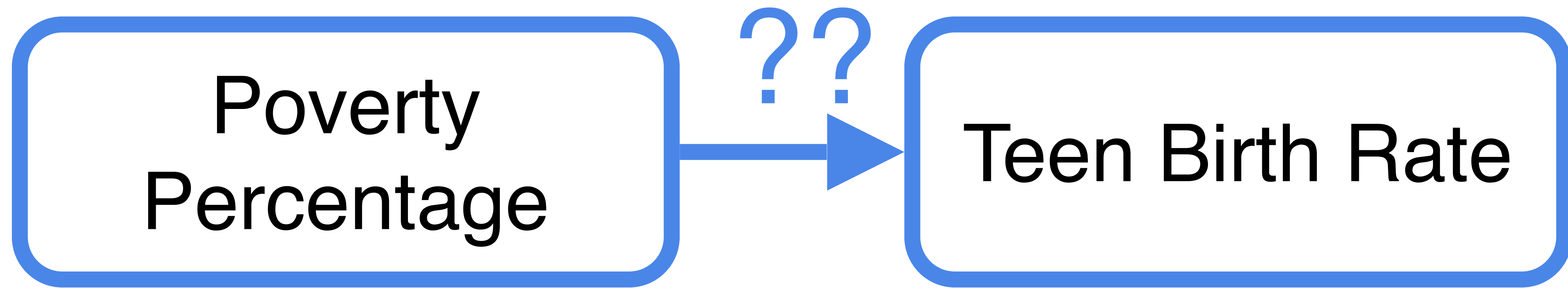


*Not homoscedastic:*  
points at this end are much  
further from the line than at  
the other end



Do *not* use  
linear  
regression

Does Poverty  
Percentage affect Teen  
Birth Rate?



Null Hypothesis:

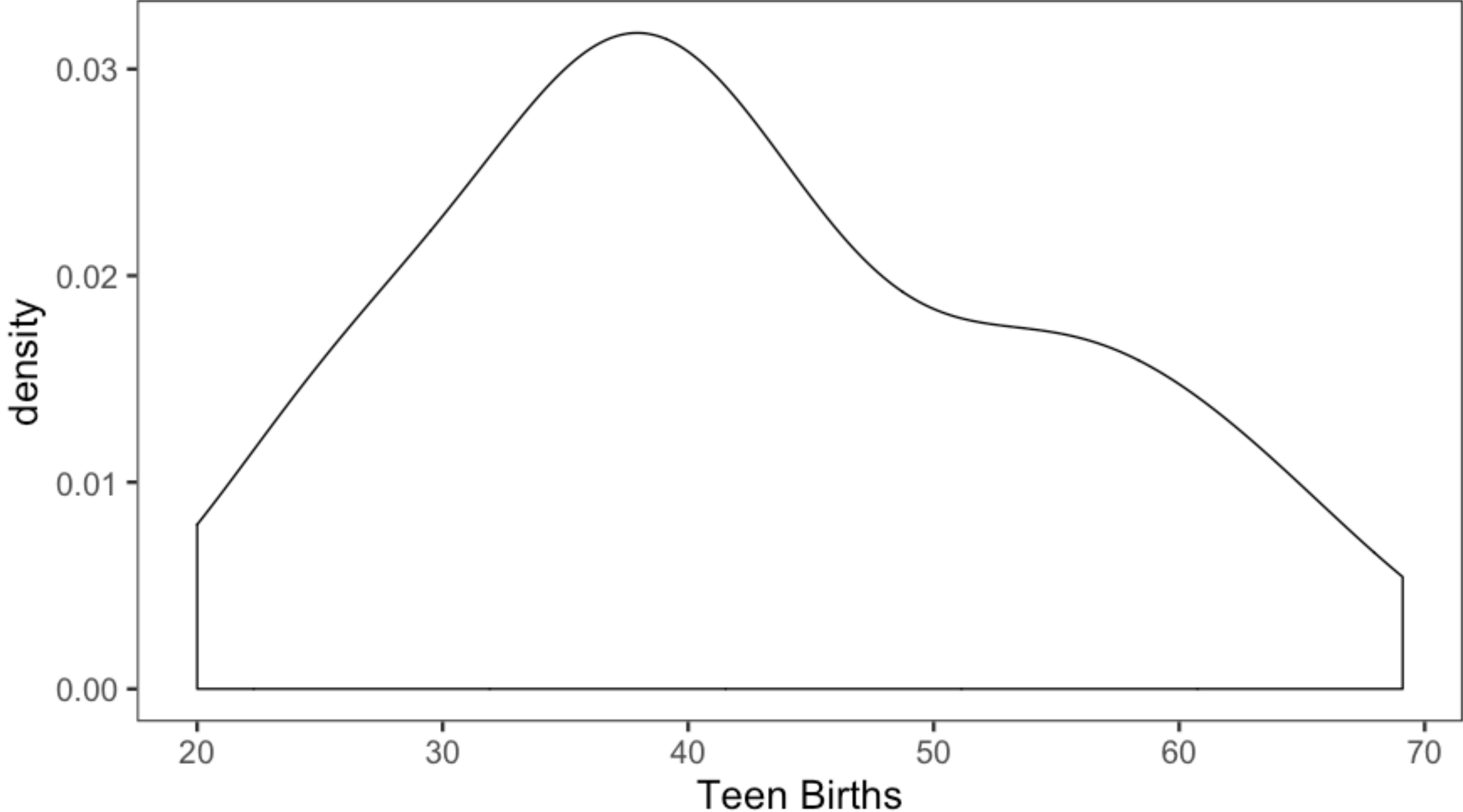
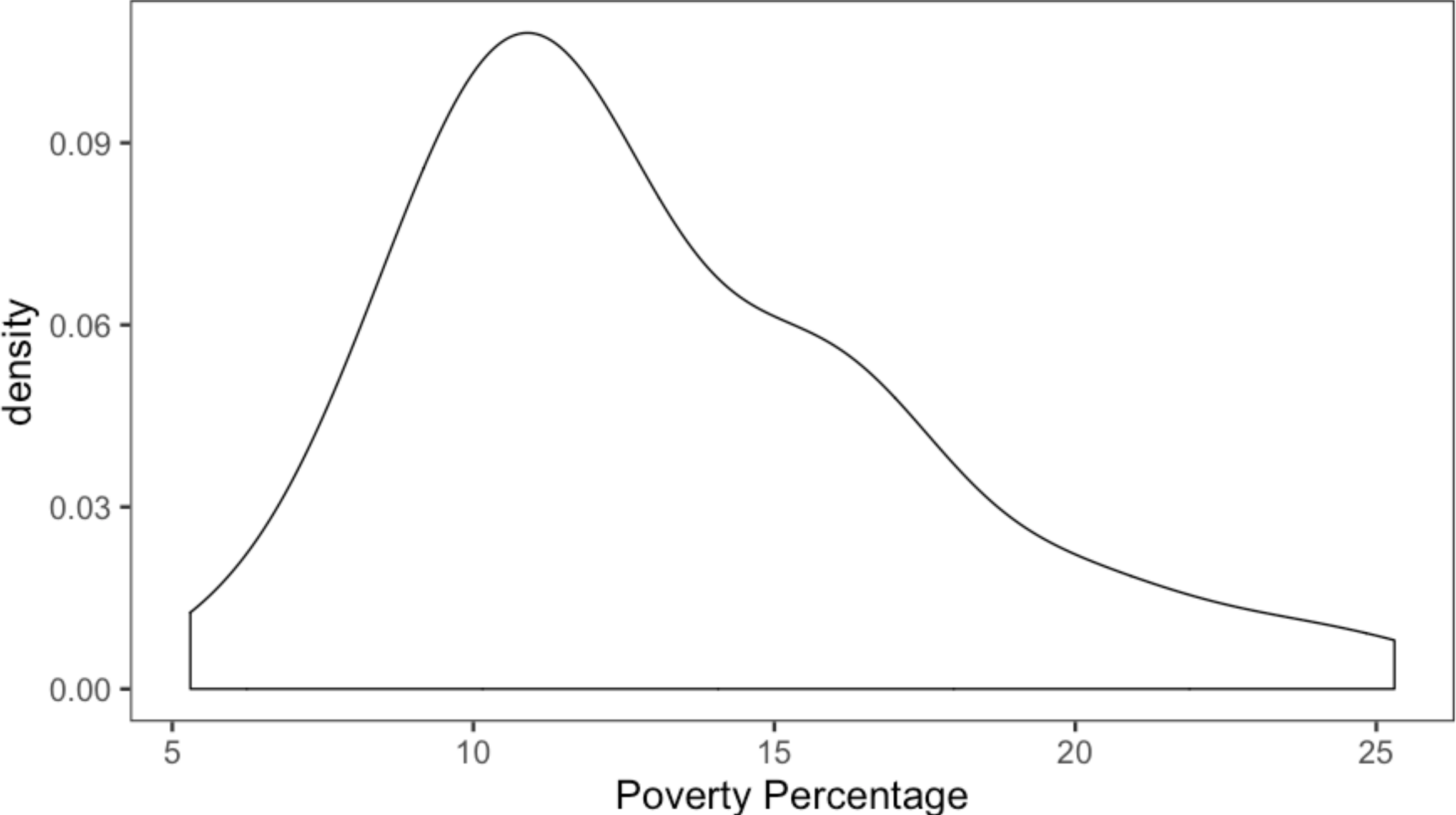
$H_0$ : Poverty Rate does not affect Teen Birth Rate ( $\beta_1=0$ )

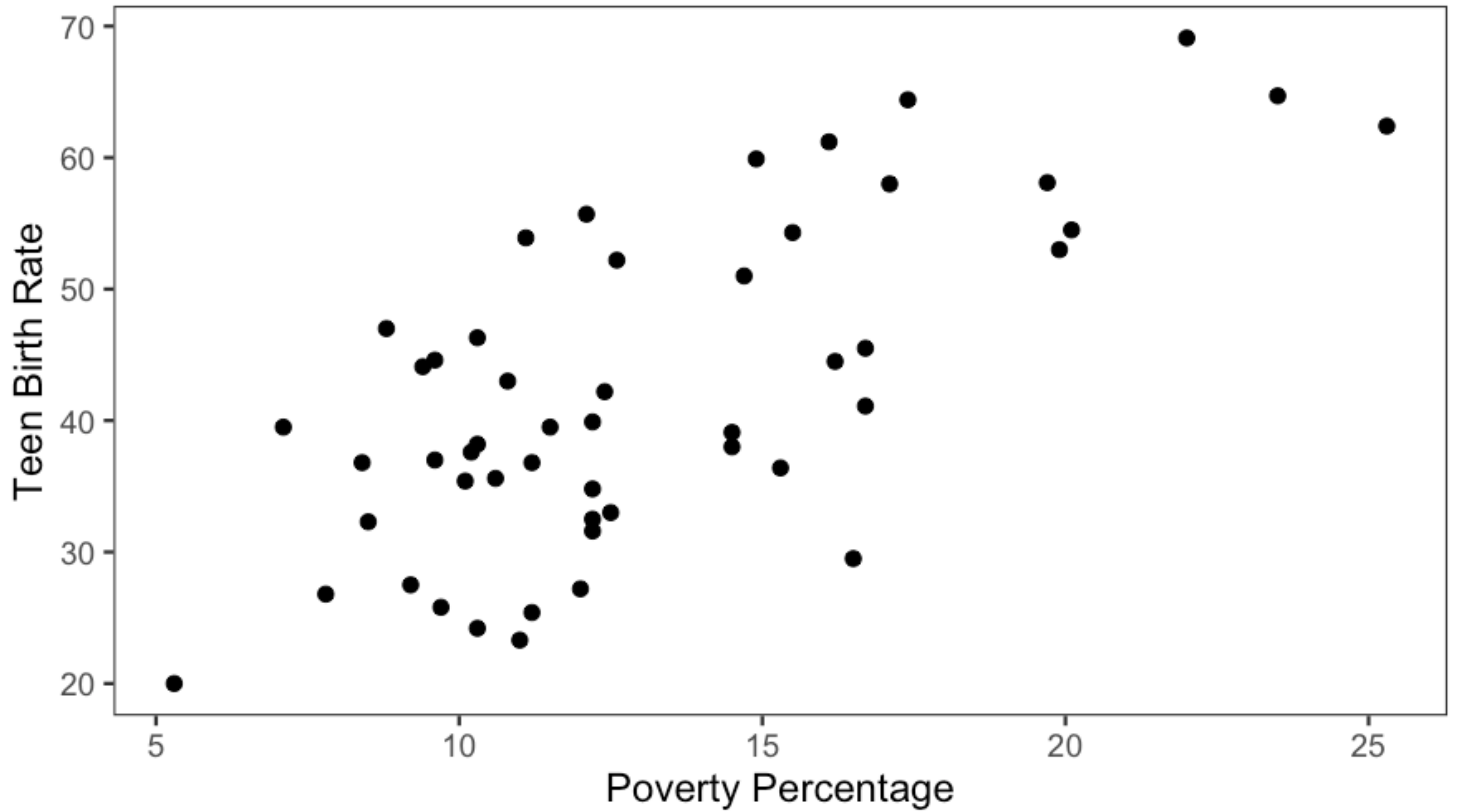
Alternative Hypothesis:

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

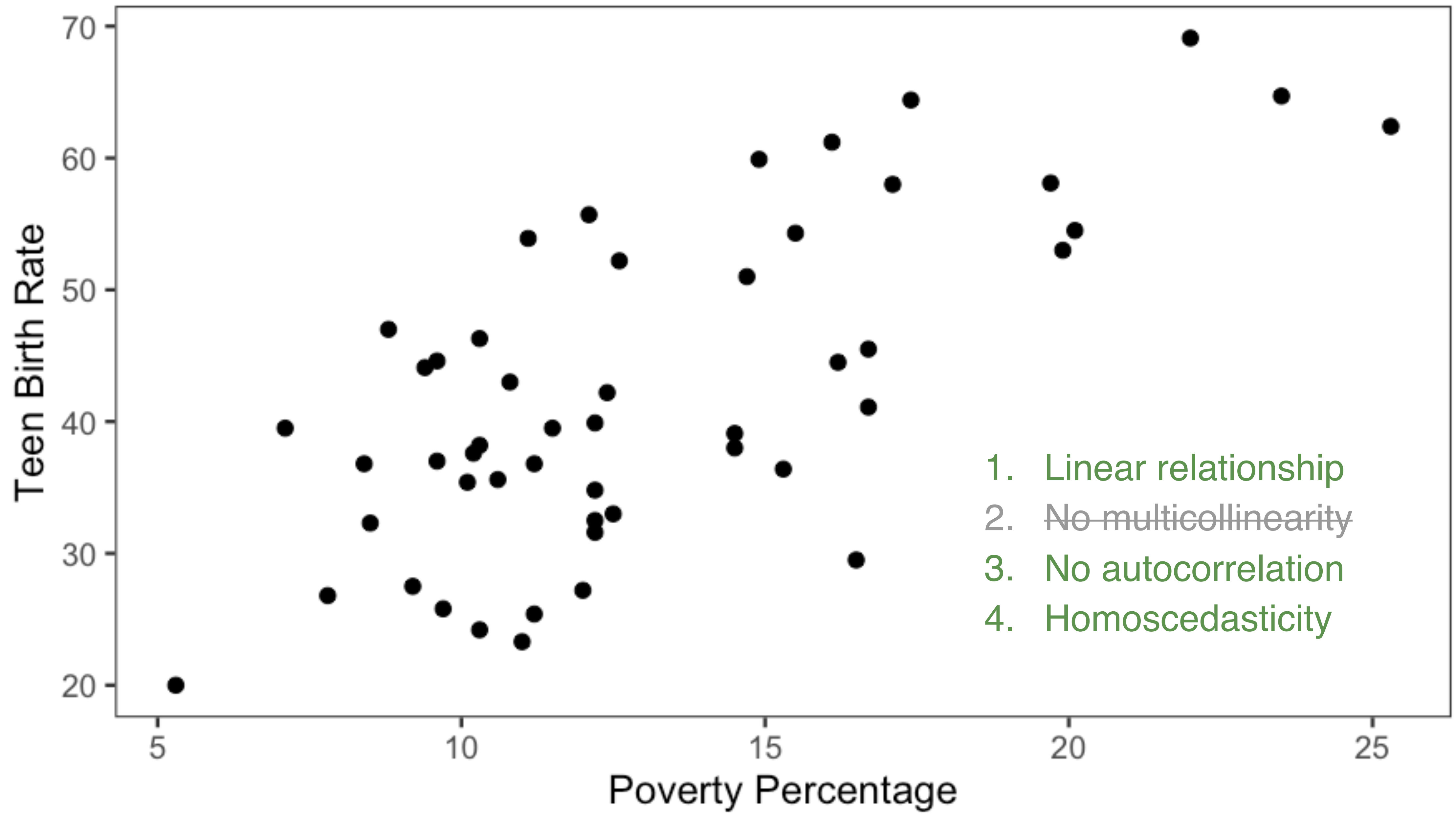
	Location	PovPct	Brth15to17	Brth18to19	ViolCrime	TeenBrth
1	Alabama	20.1	31.5	88.7	11.2	54.5
2	Alaska	7.1	18.9	73.7	9.1	39.5
3	Arizona	16.1	35.0	102.5	10.4	61.2
4	Arkansas	14.9	31.6	101.7	10.4	59.9
5	California	16.7	22.6	69.1	11.2	41.1
6	Colorado	8.8	26.2	79.1	5.8	47.0
7	Connecticut	9.7	14.1	45.1	4.6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

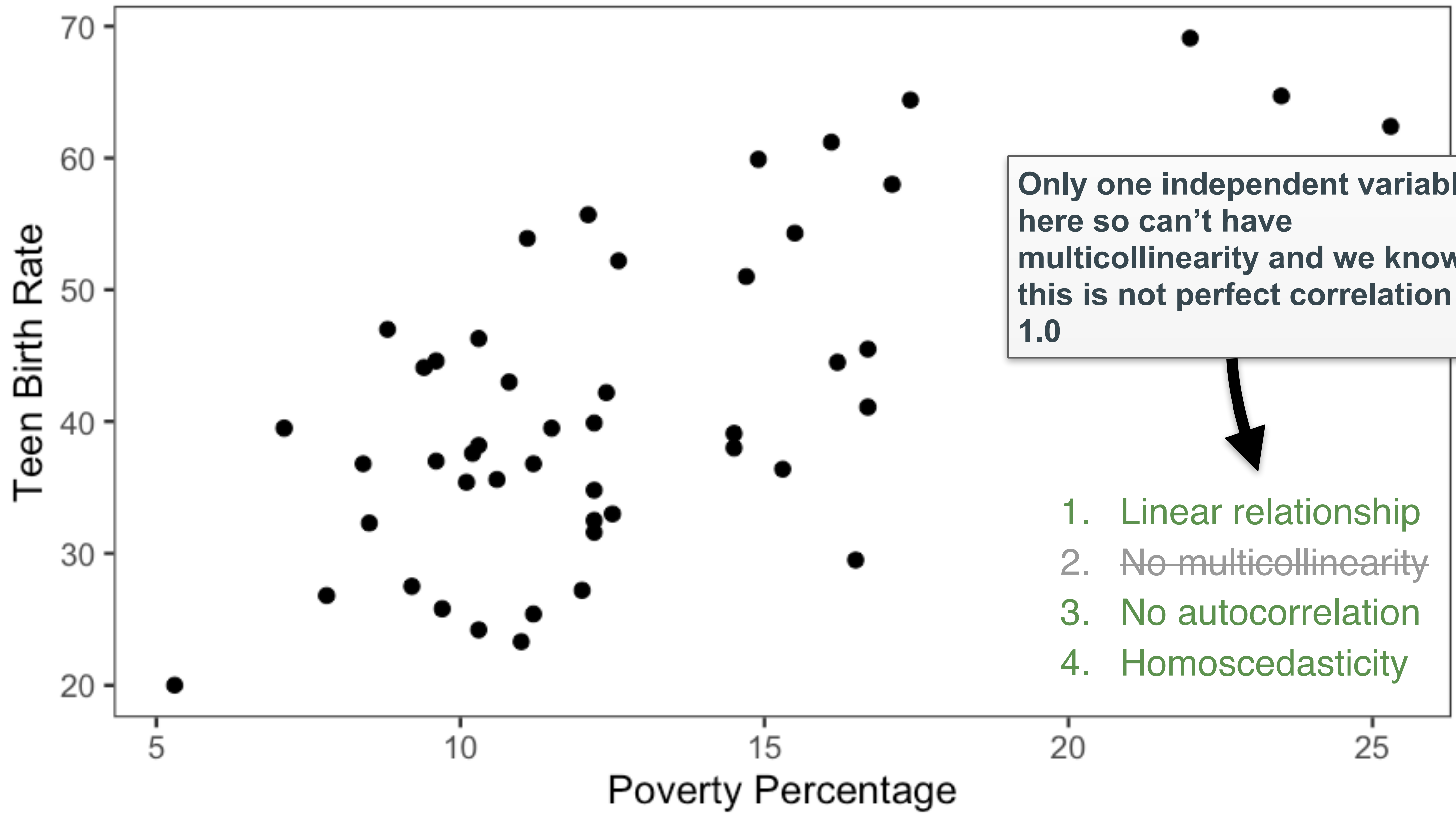
# EDA: distributions





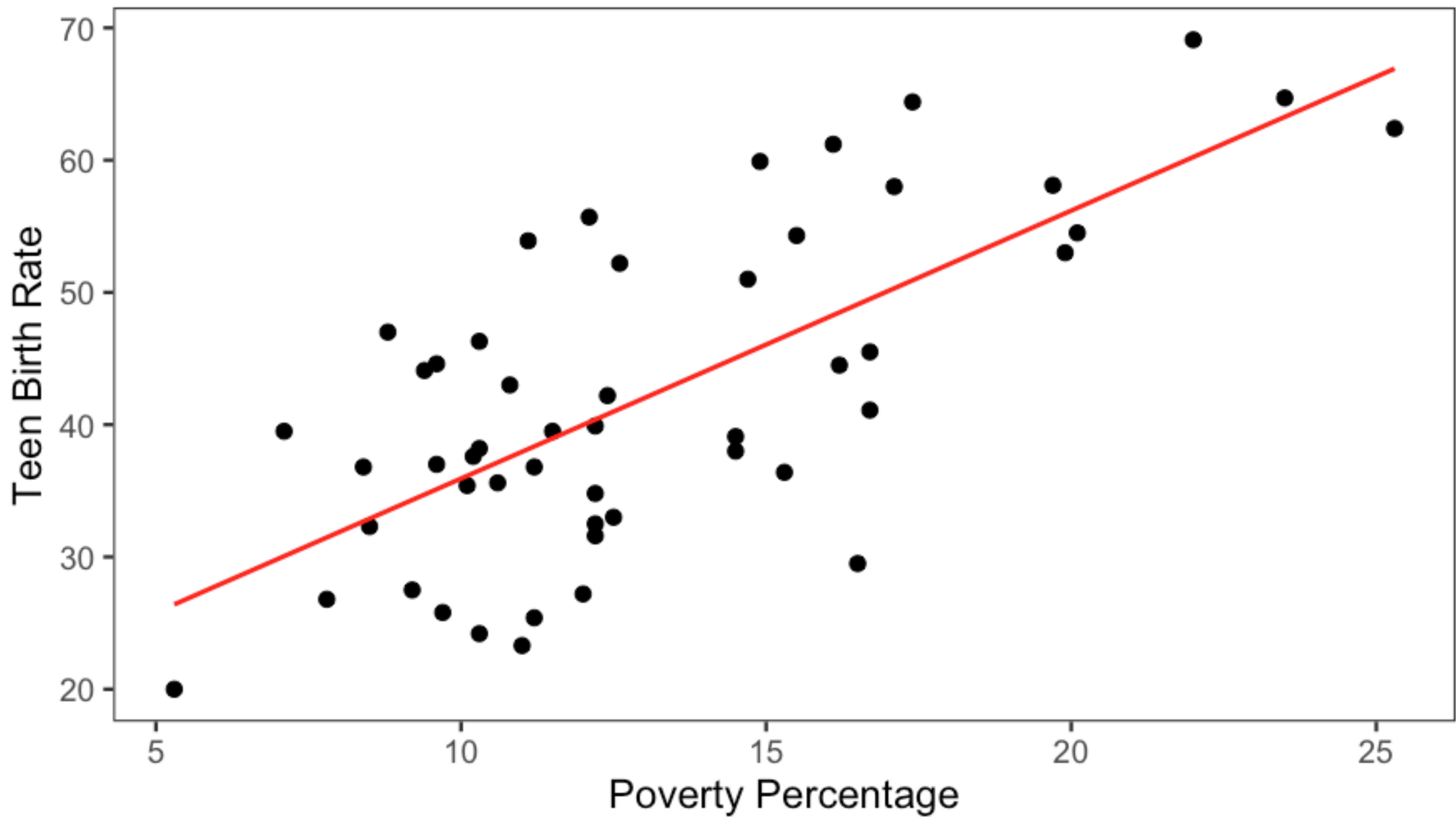


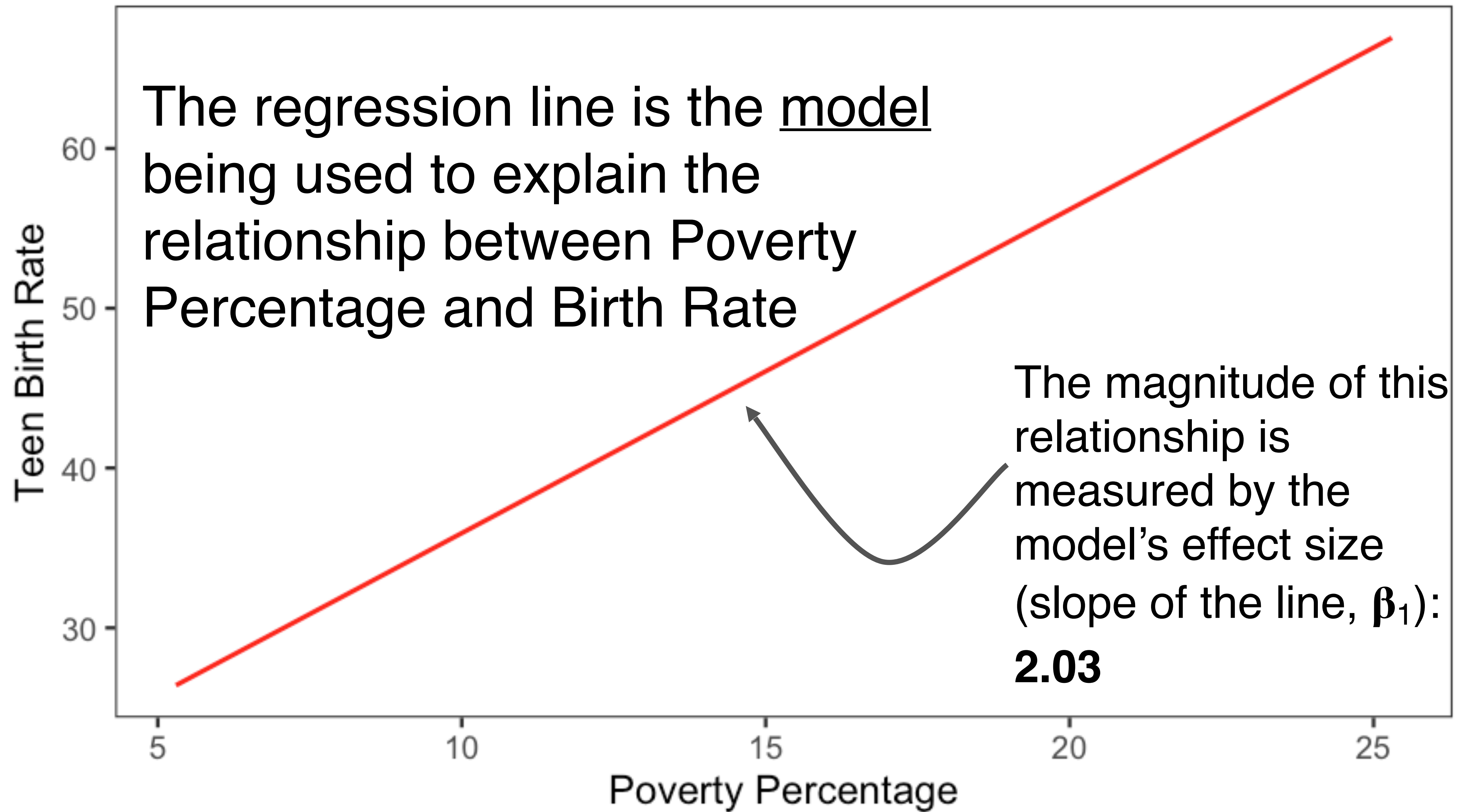


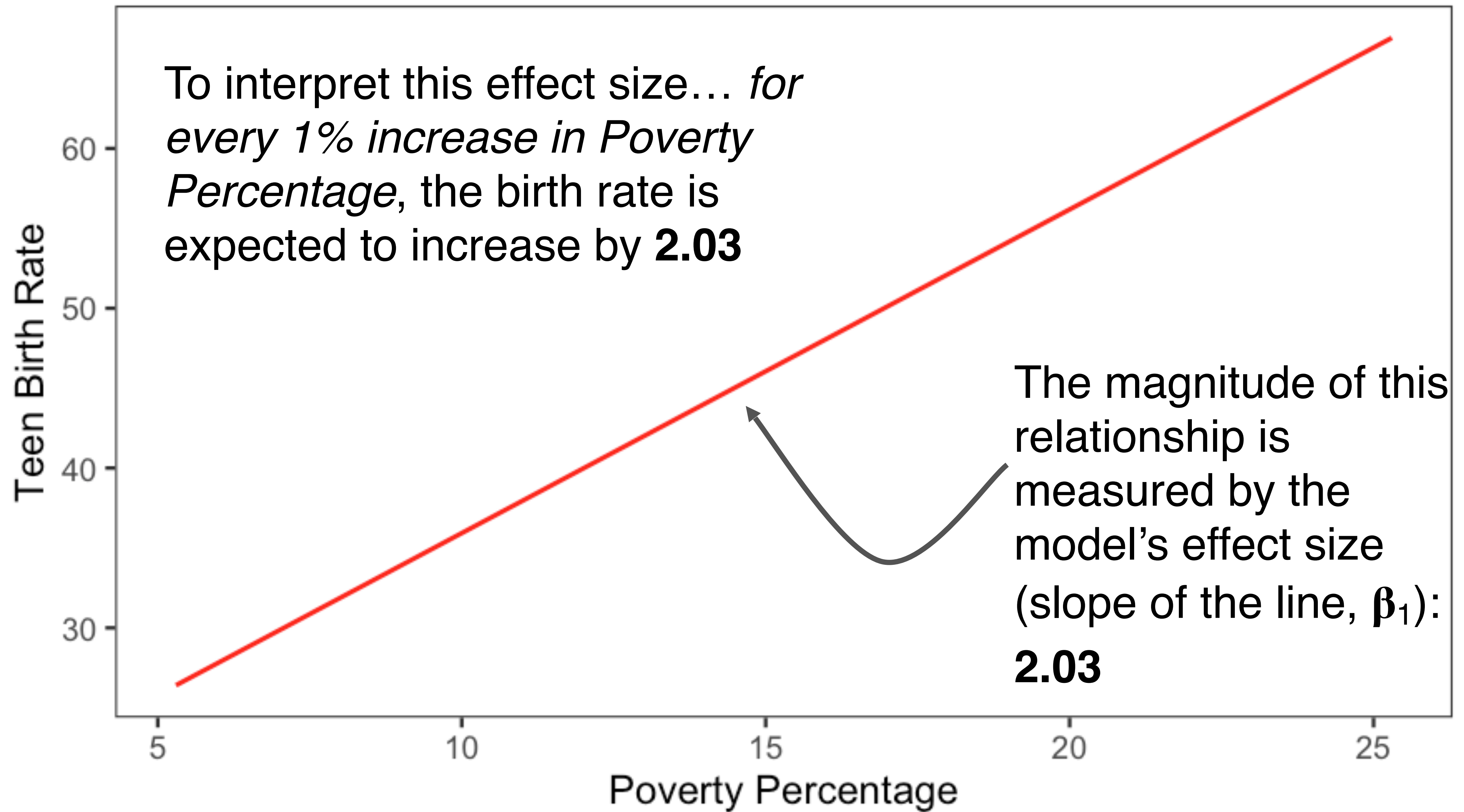


Only one independent variable here so can't have multicollinearity and we know this is not perfect correlation of 1.0

- 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)

Teen Birth Rate

60

50

40

30

5

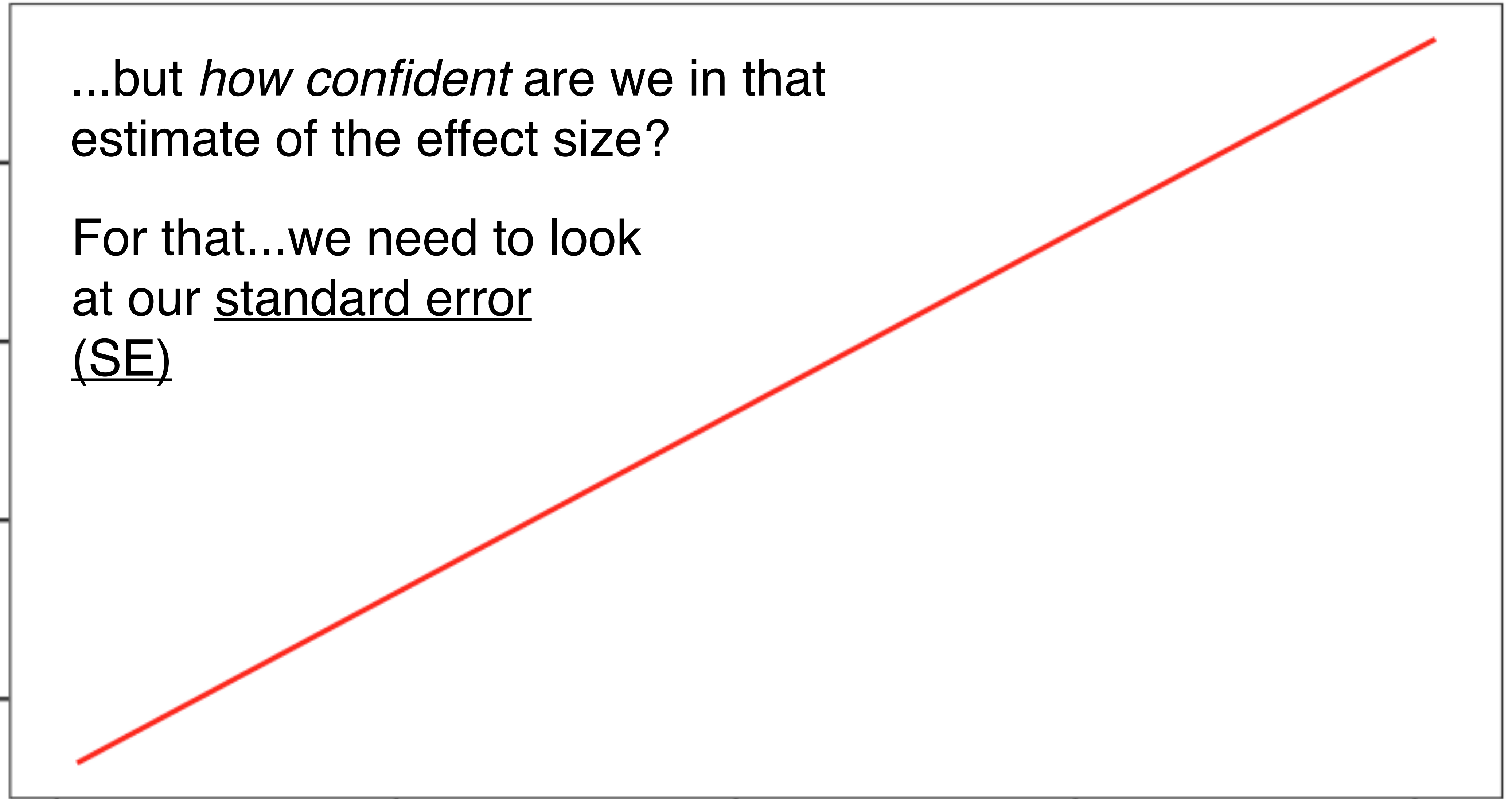
10

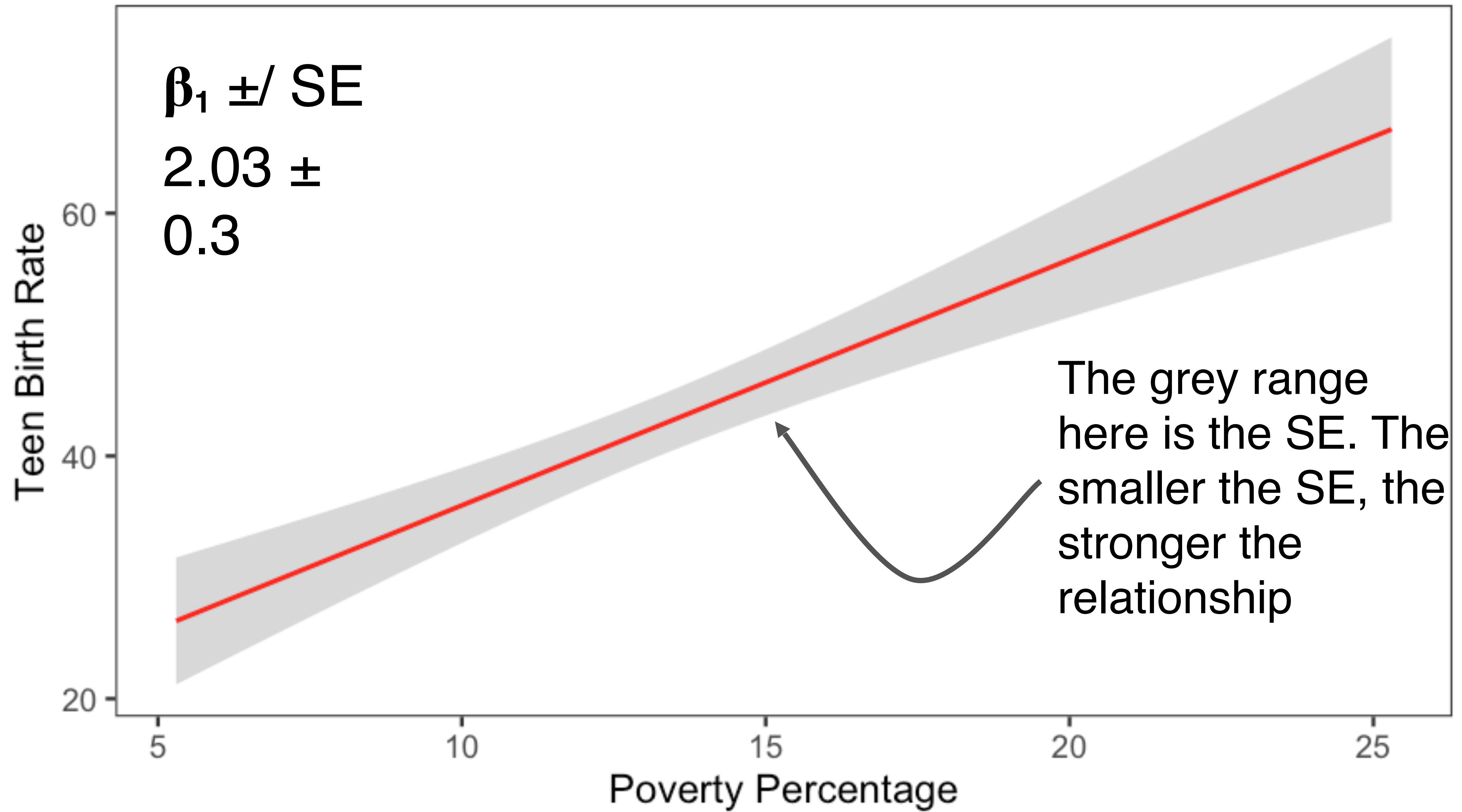
15

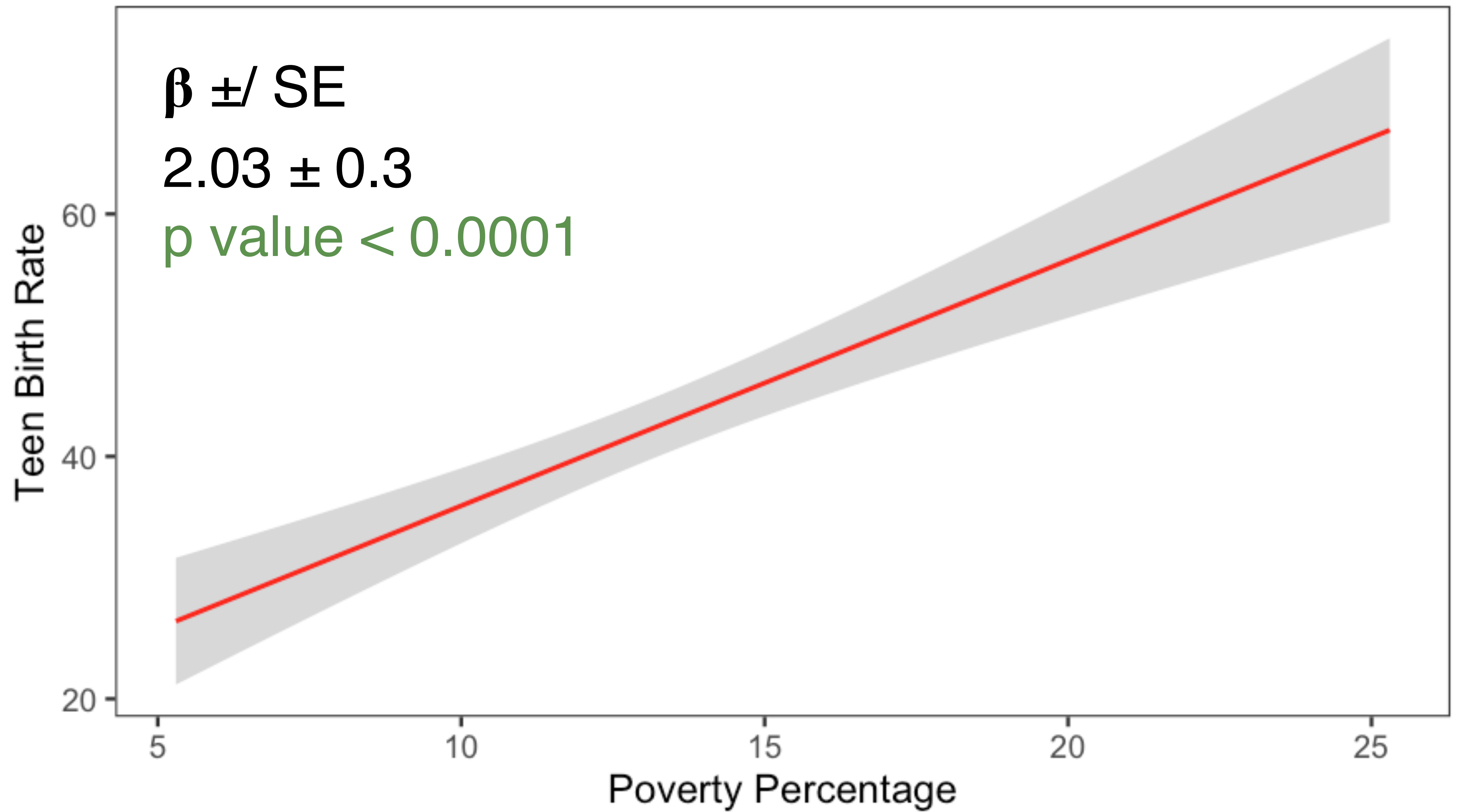
20

25

Poverty Percentage

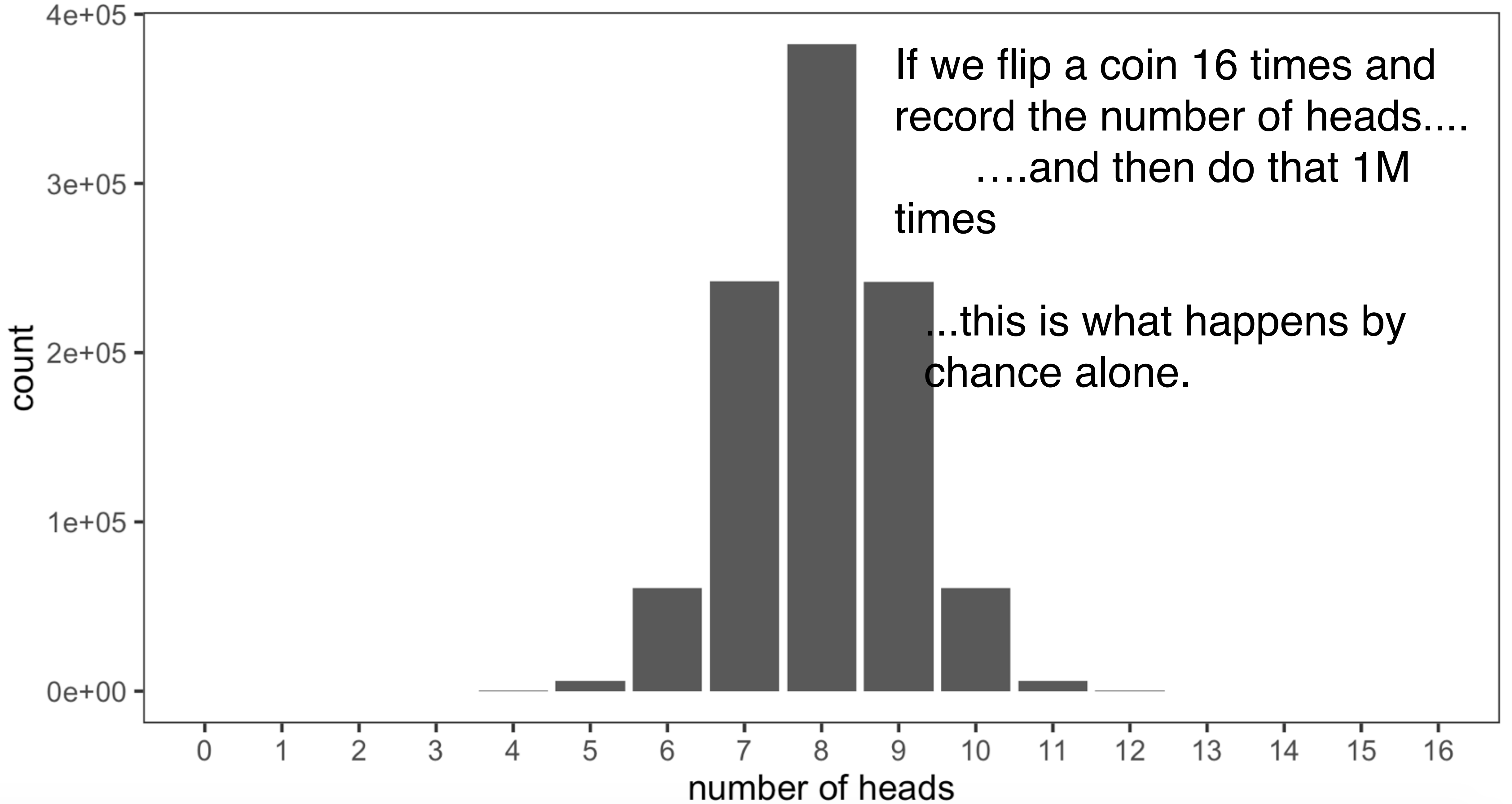


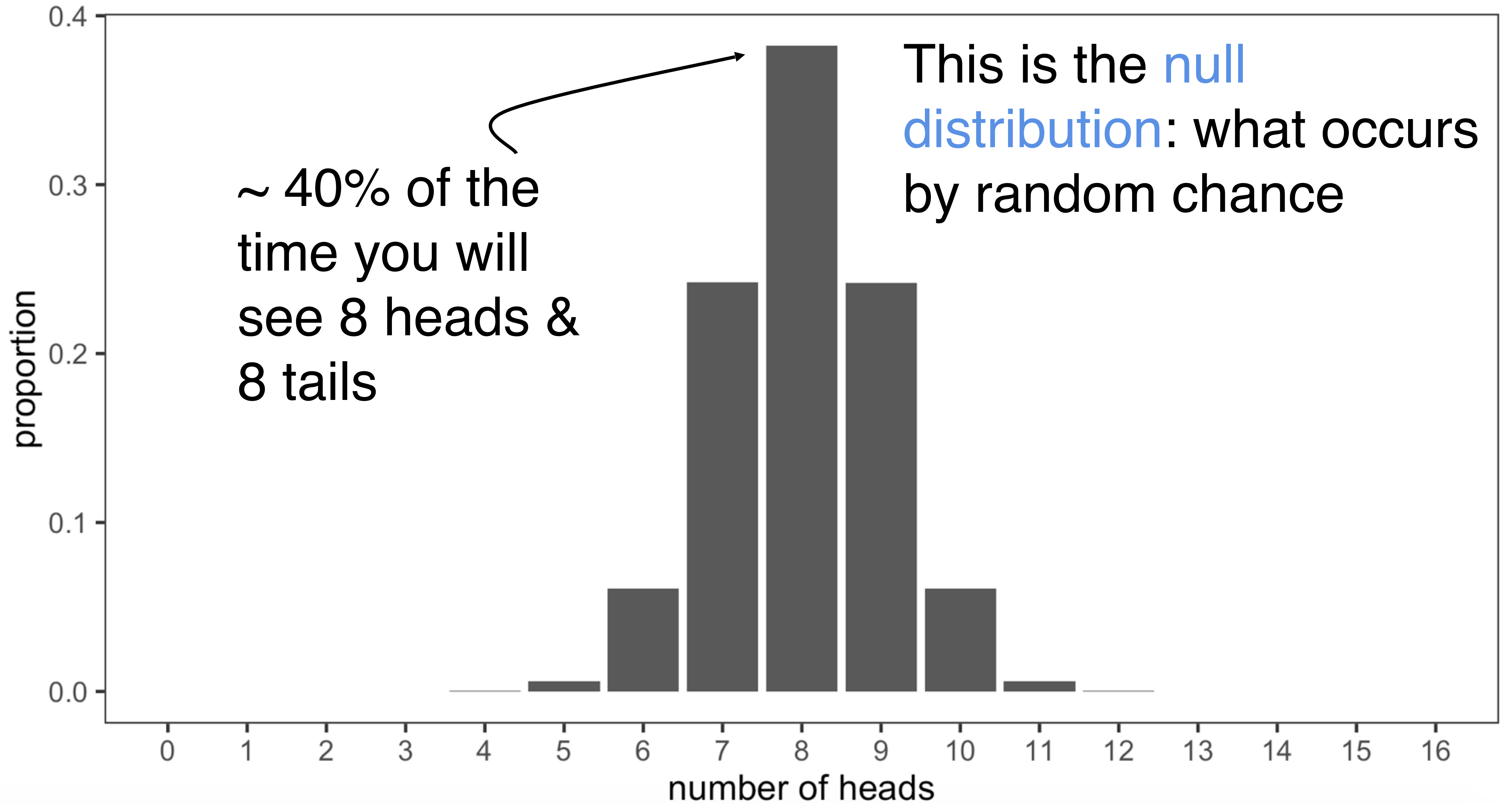


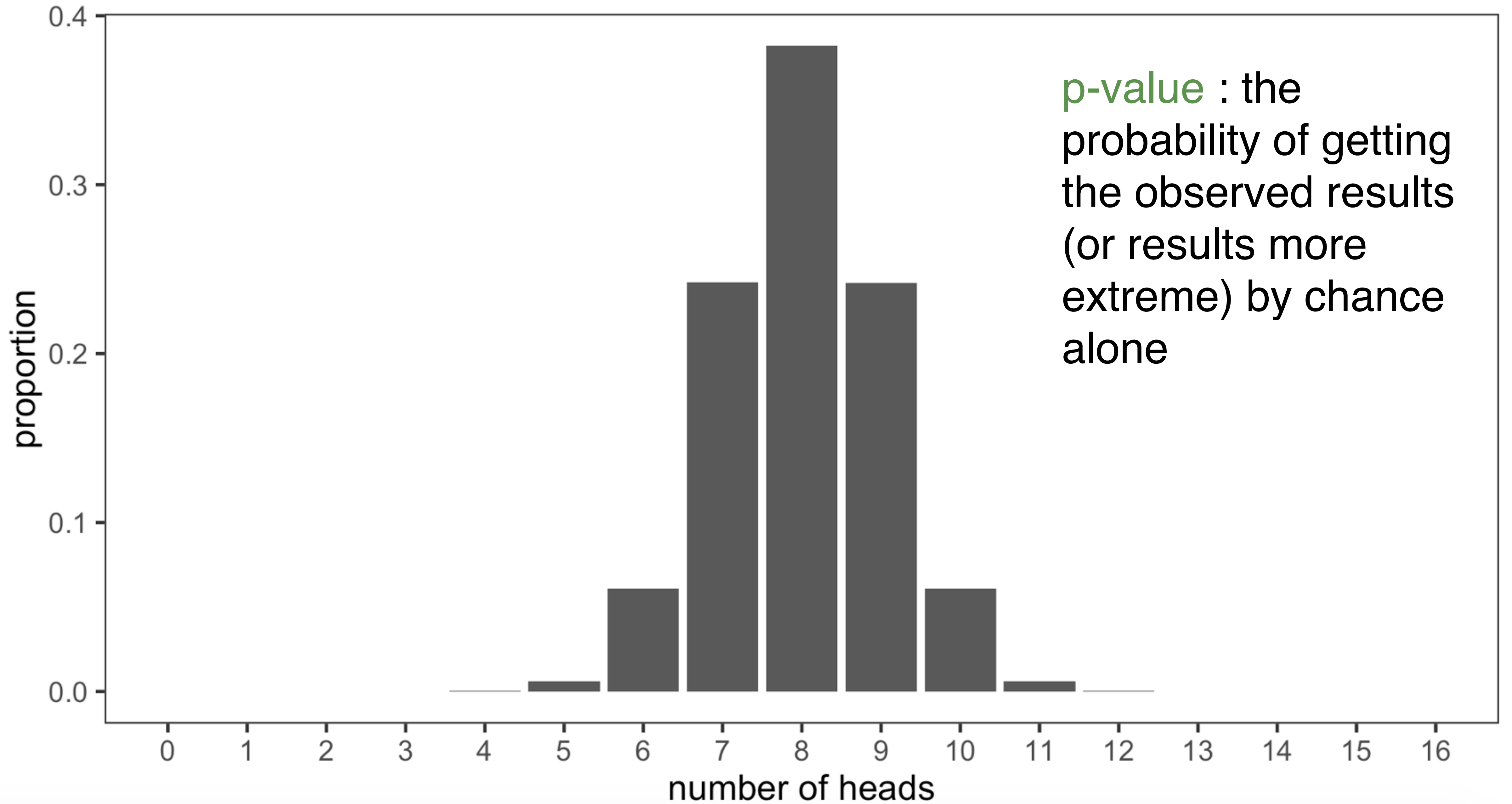


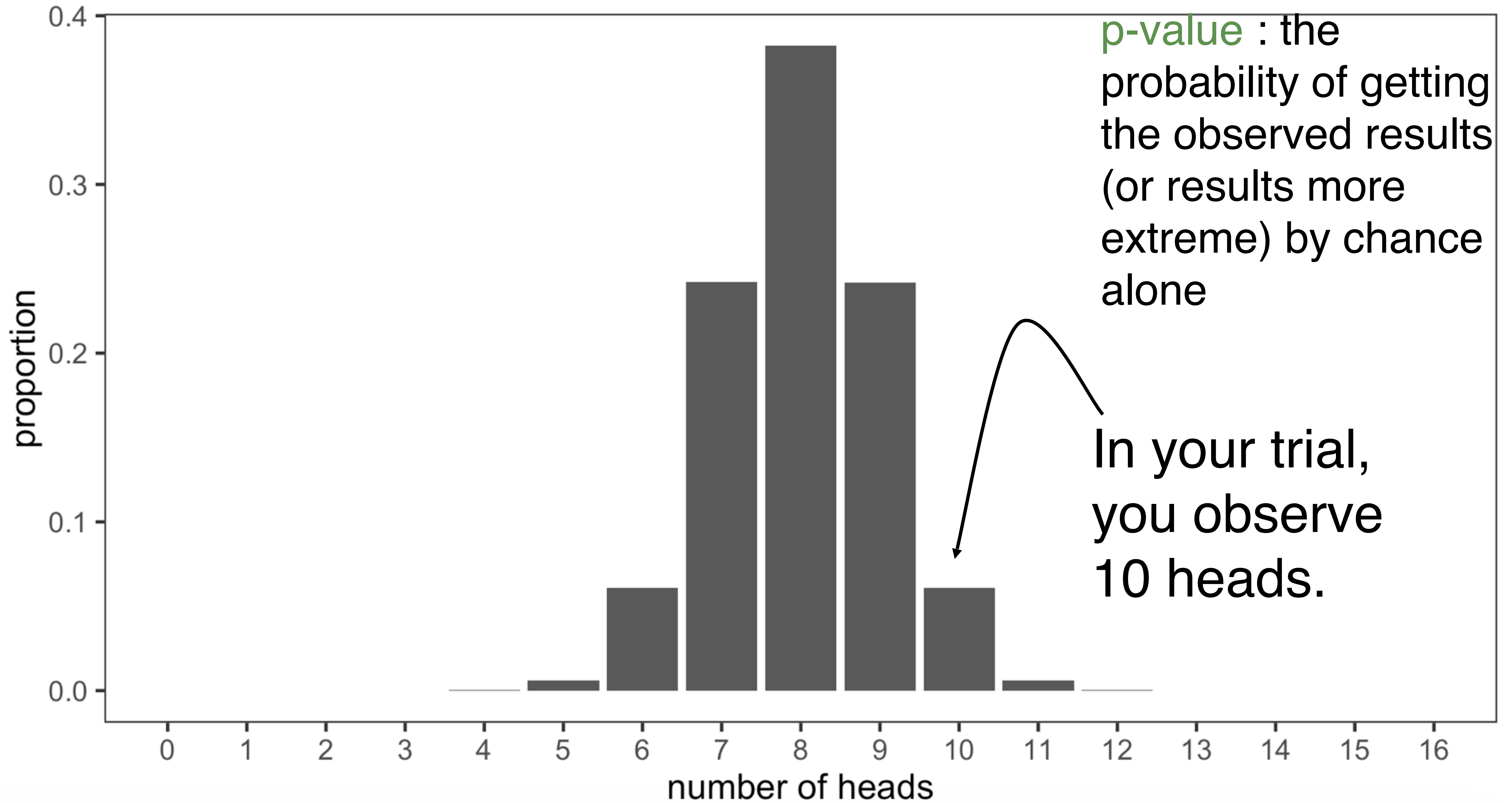


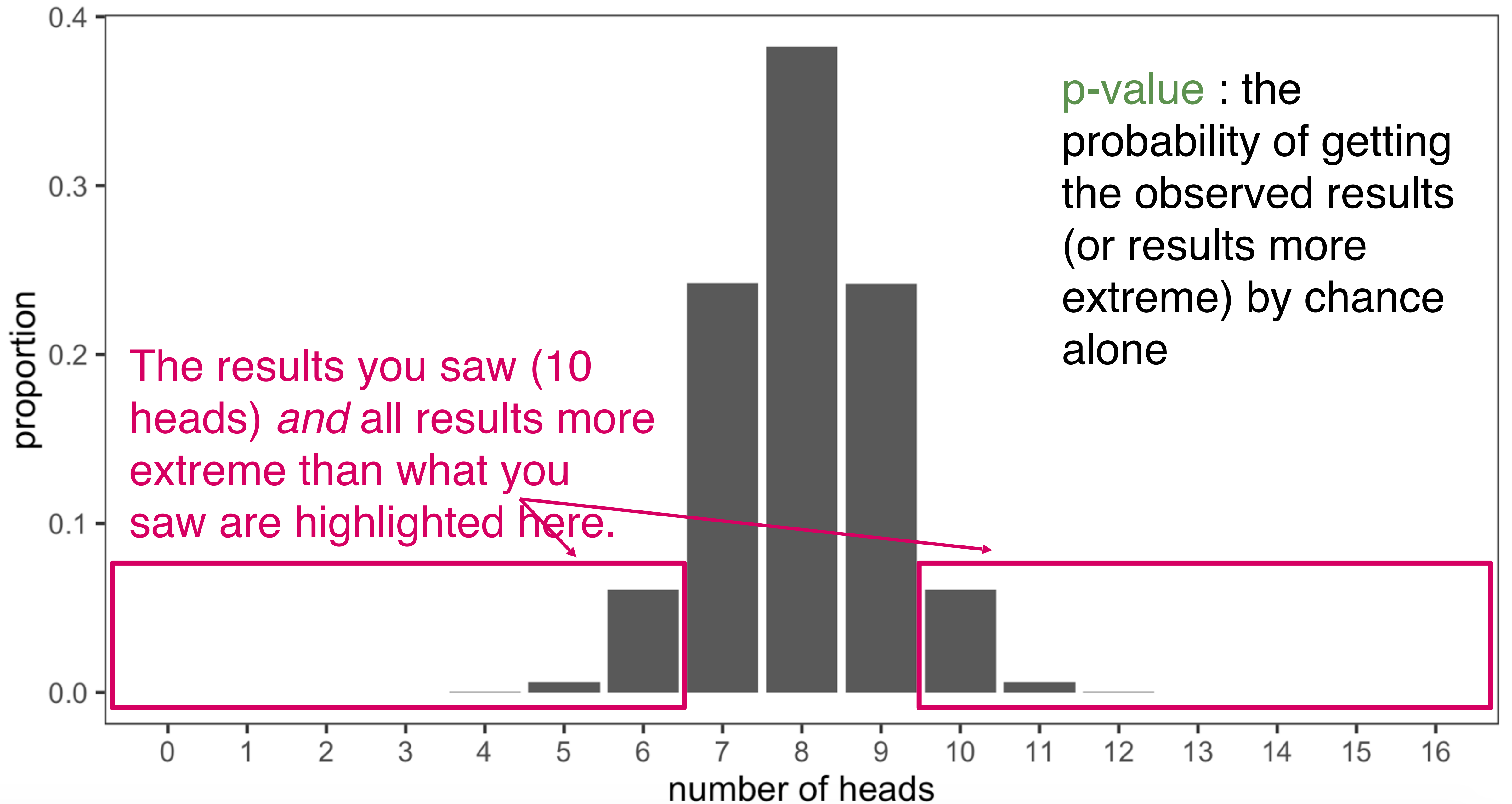
**p-value** : the probability of getting the observed results (or results more extreme) by chance alone

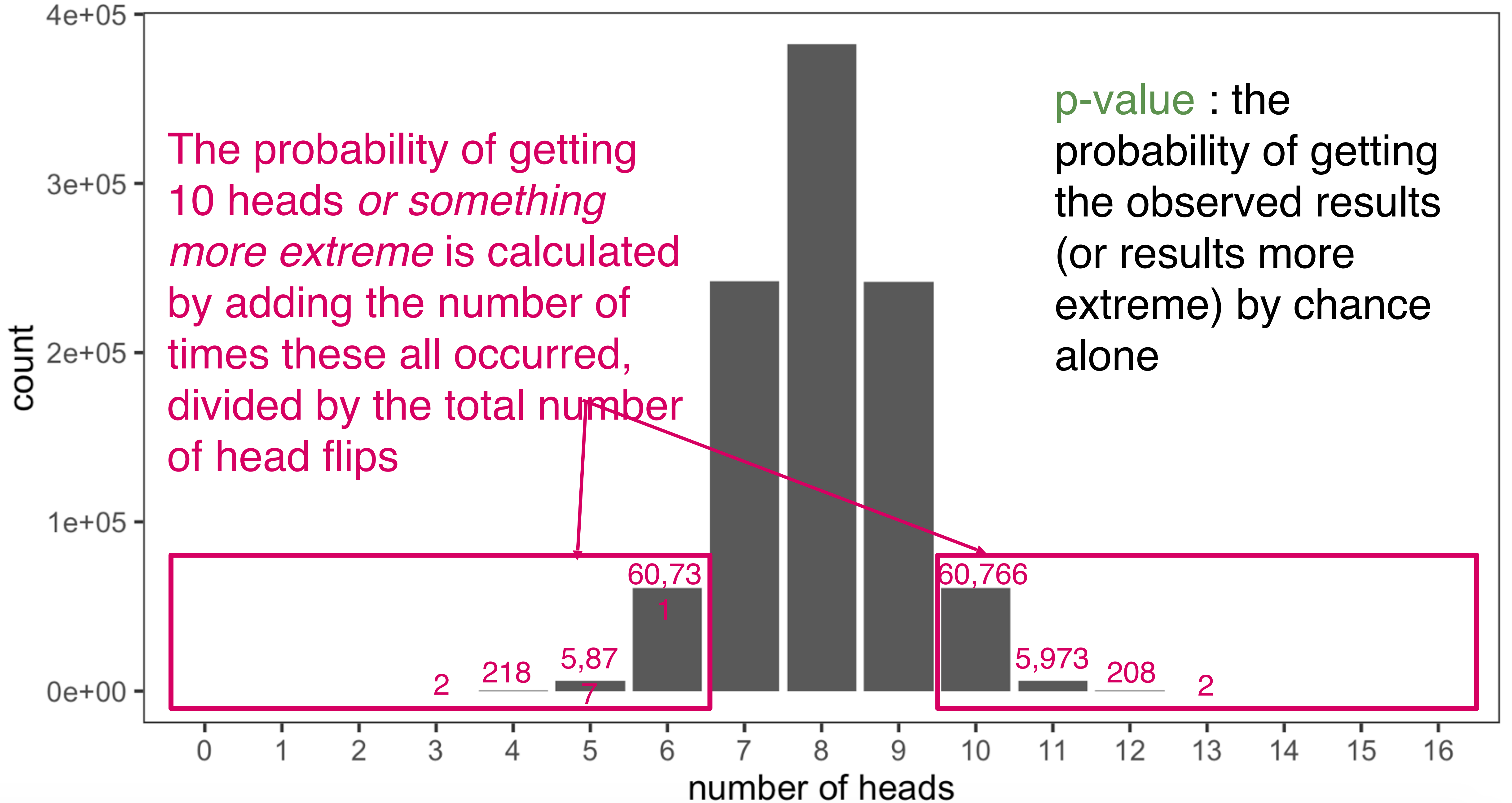


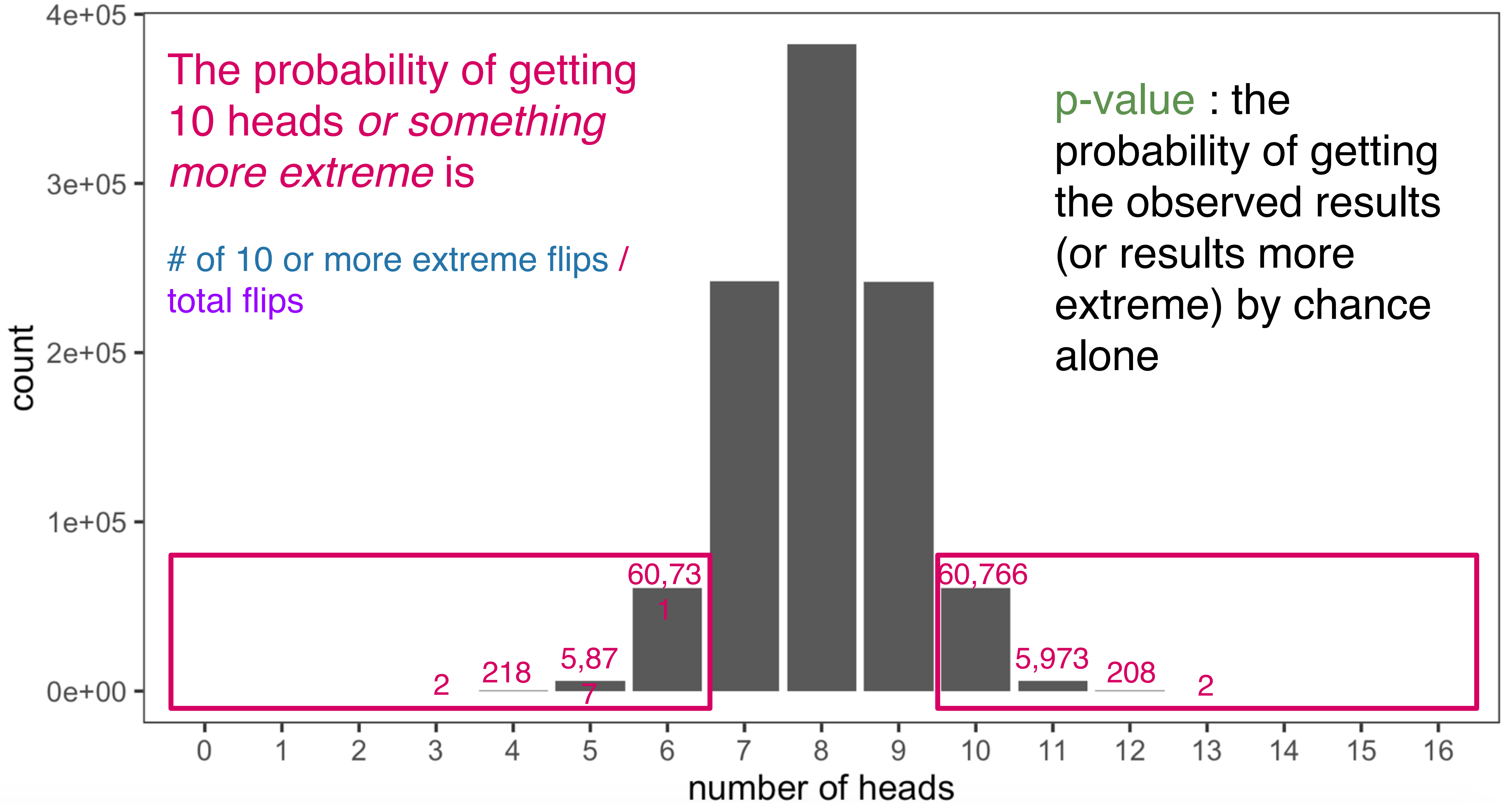




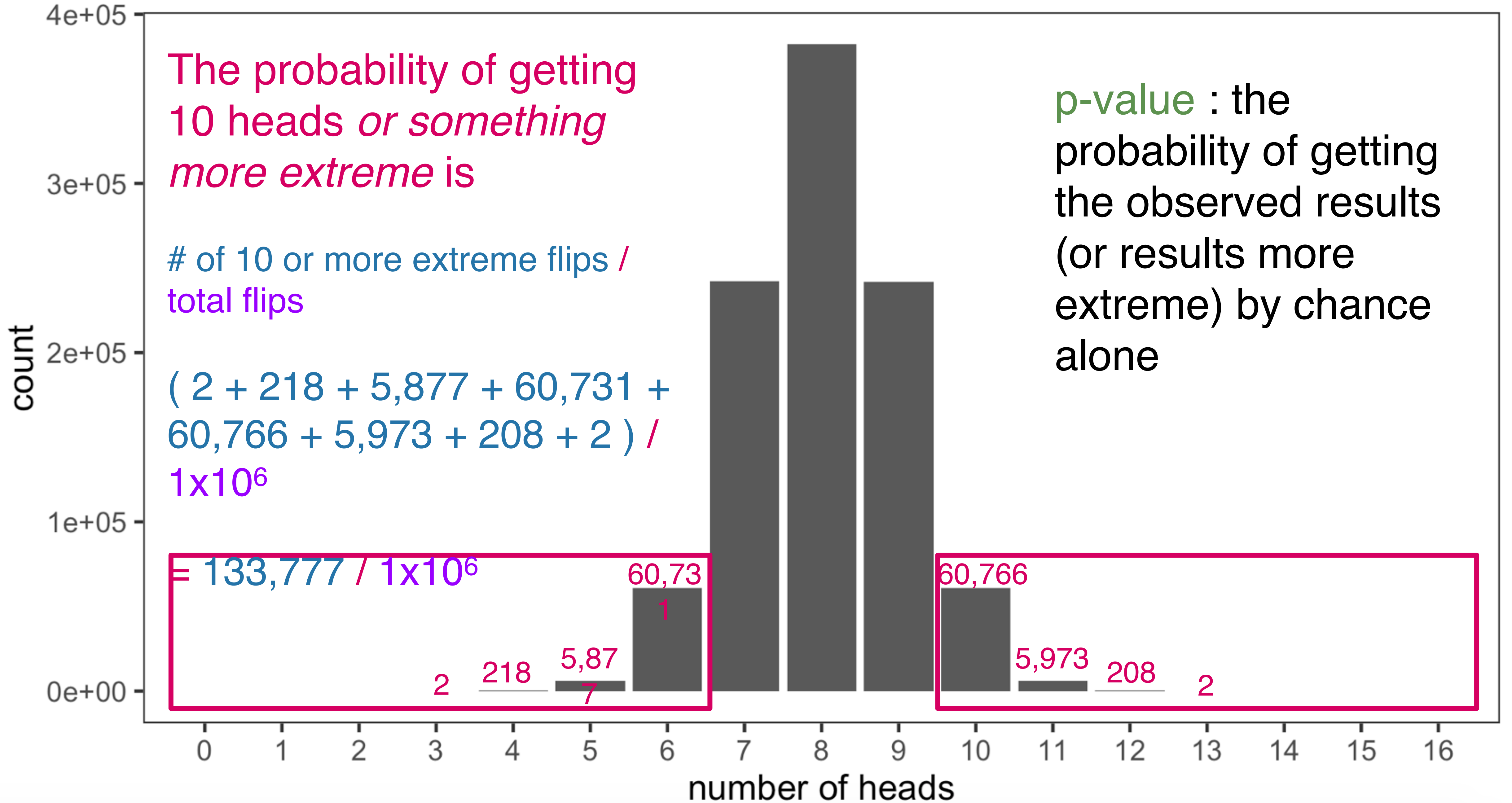


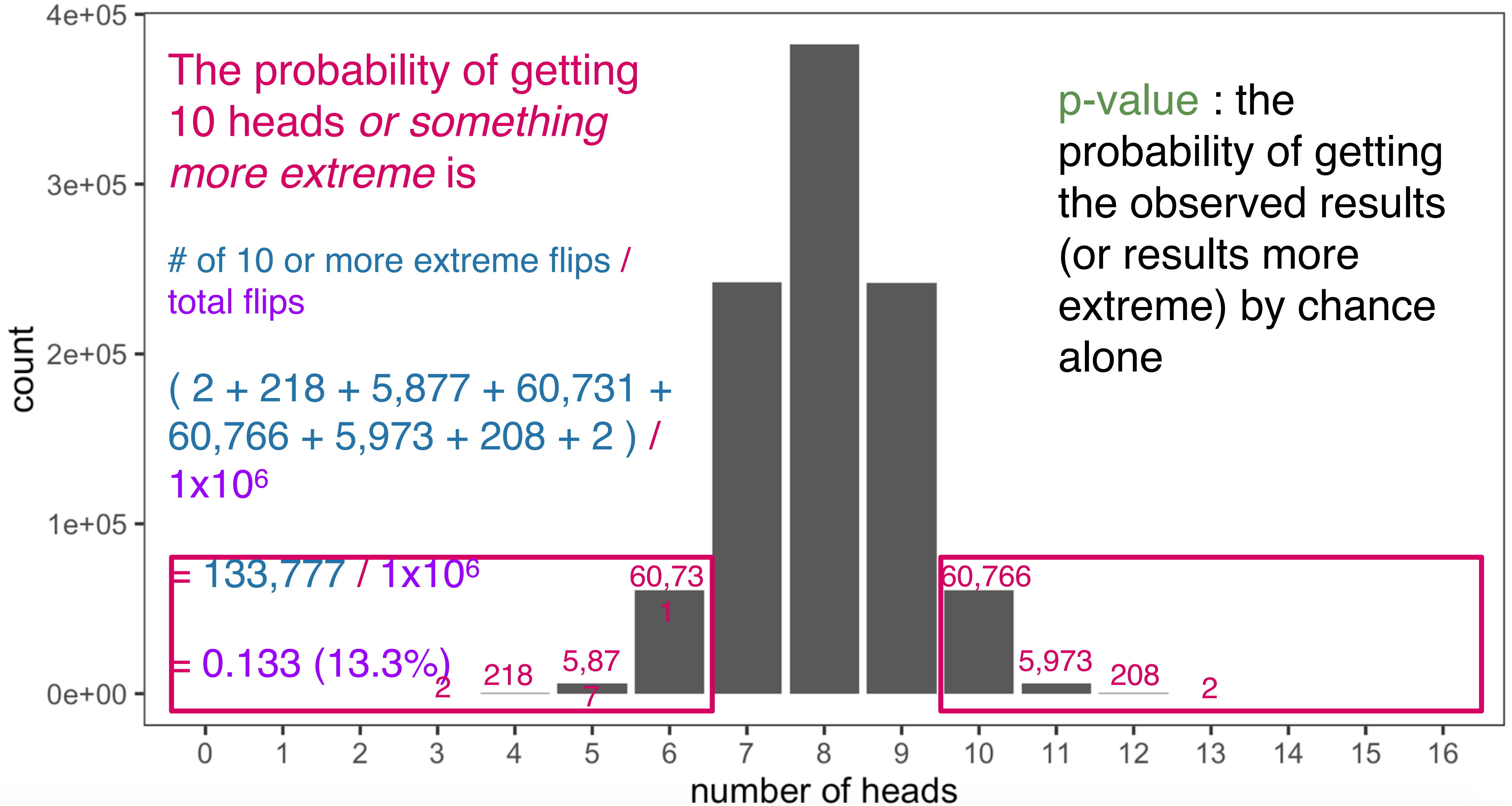


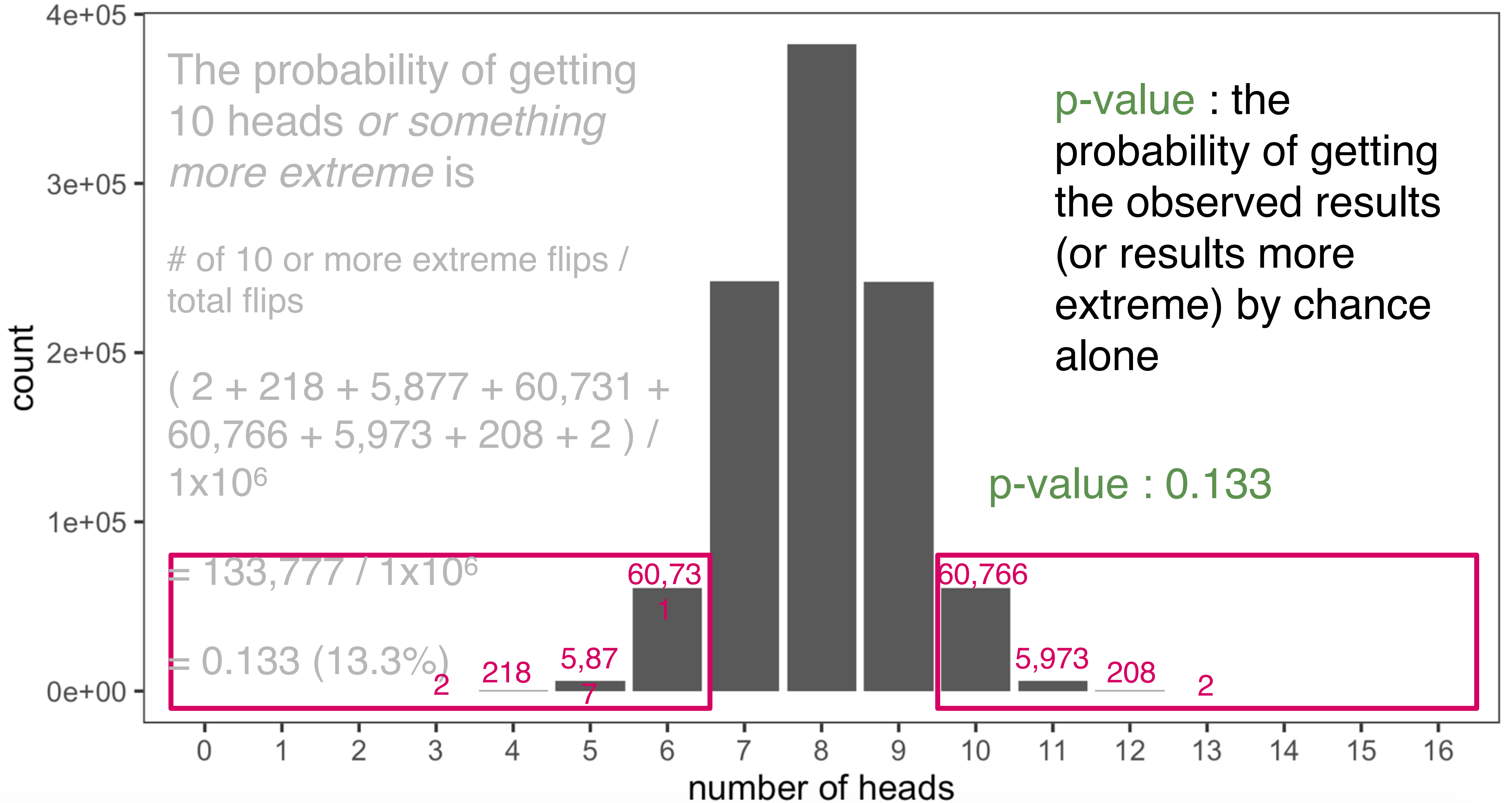


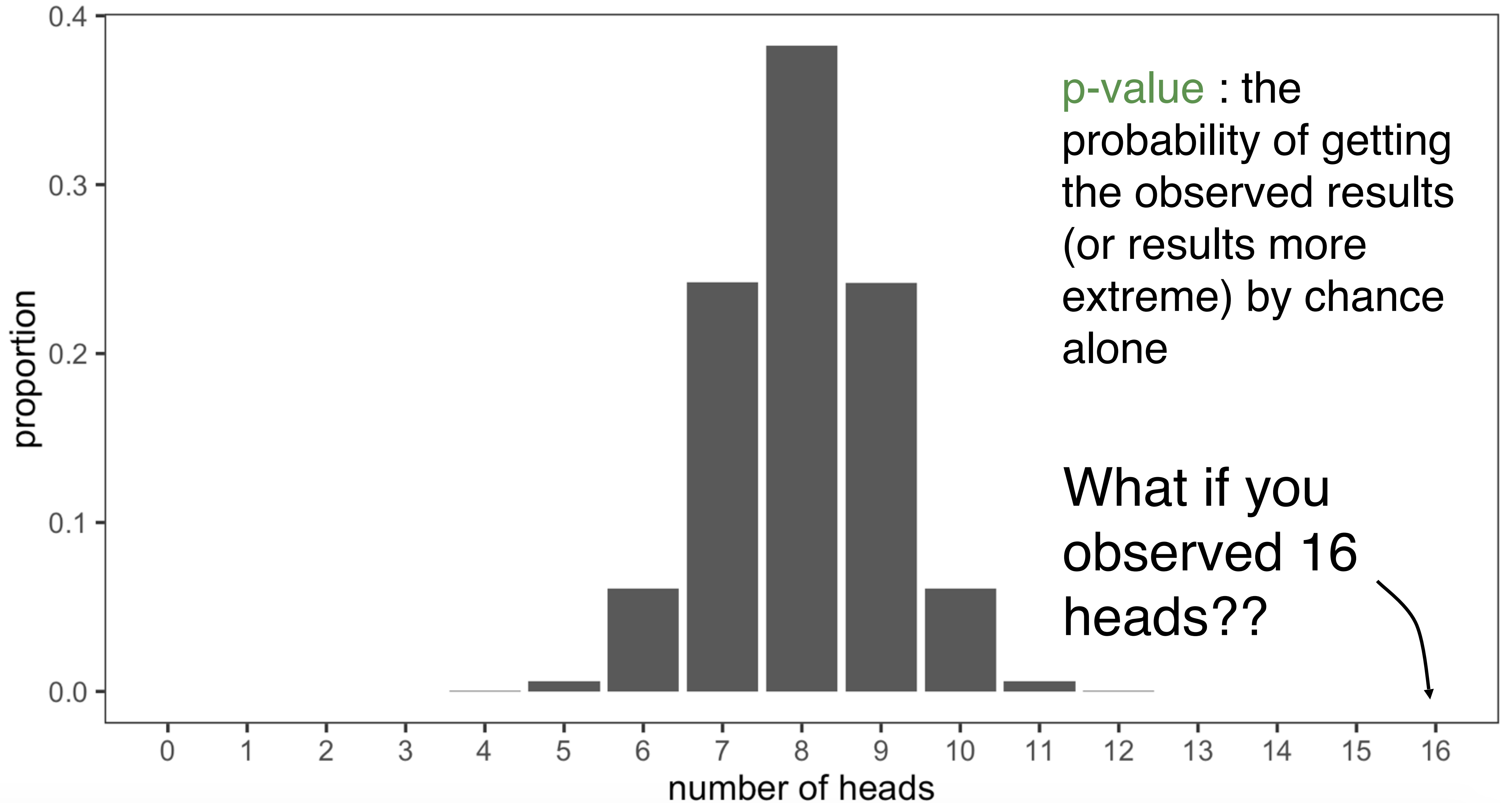


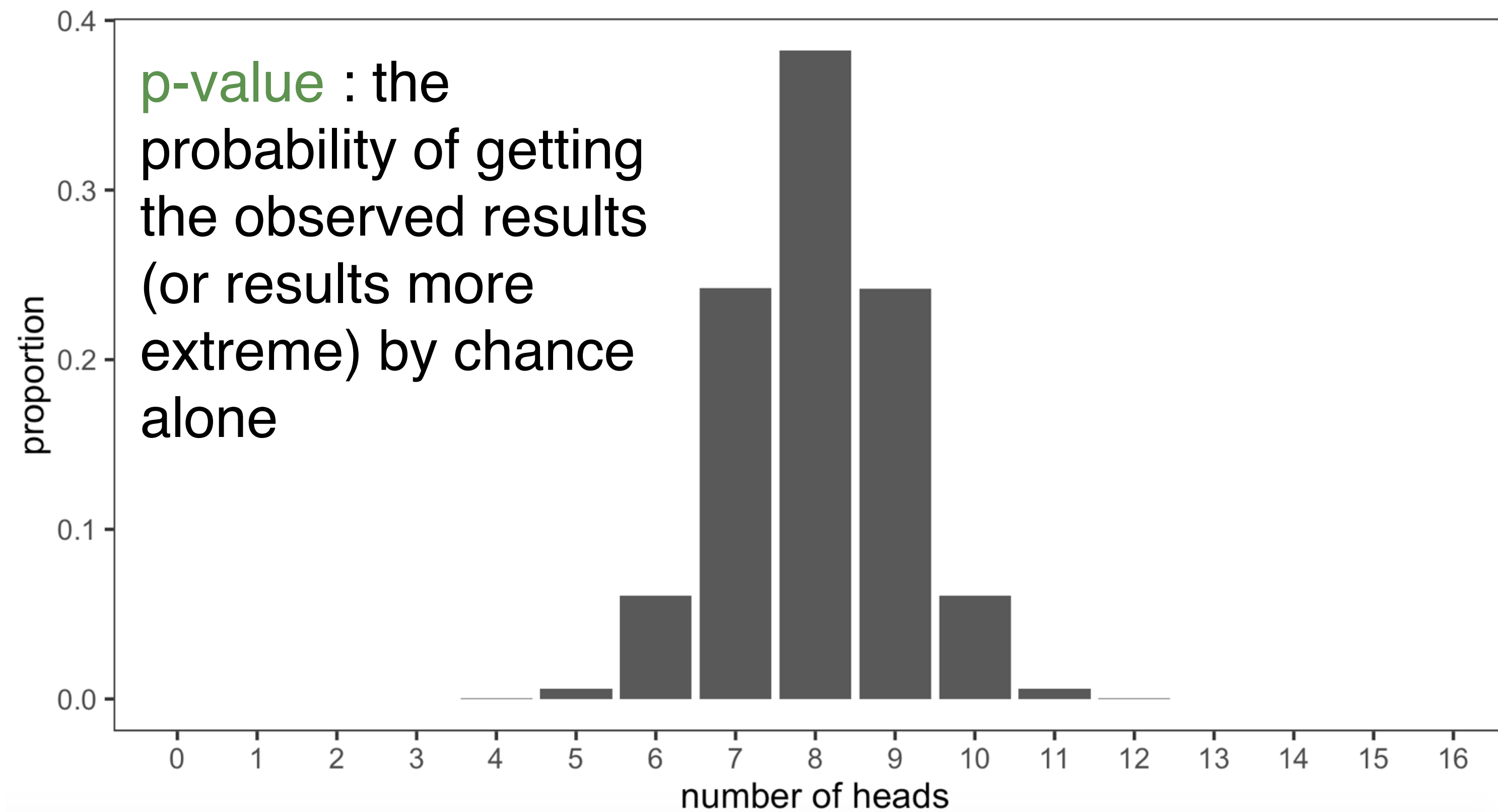




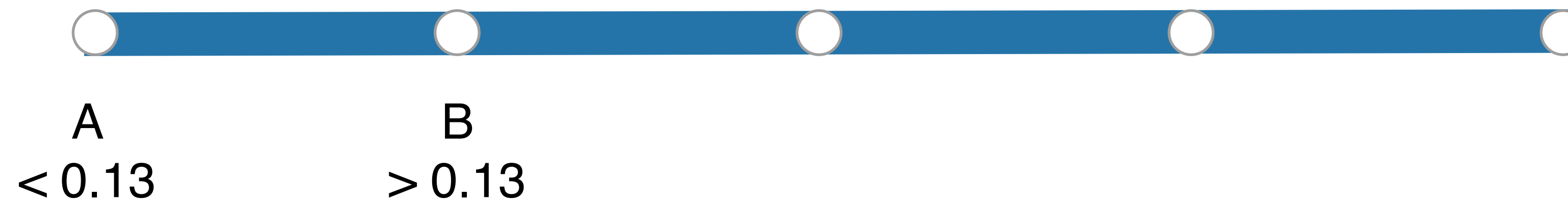


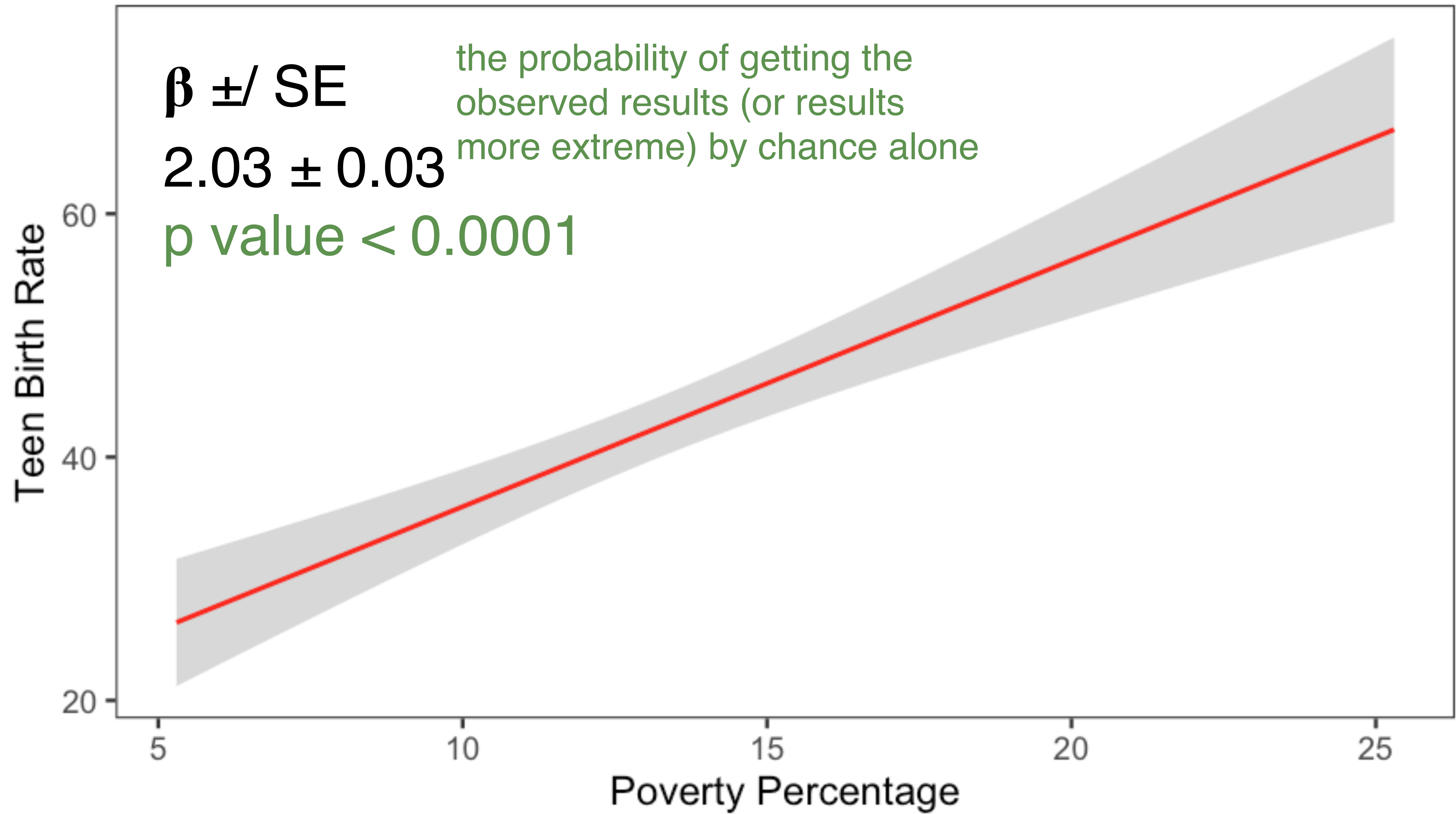






What would be the p-value of you flipping 16 heads?






$\beta \pm SE$

$2.03 \pm 0.03$

$p \text{ value} < 0.0001$

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

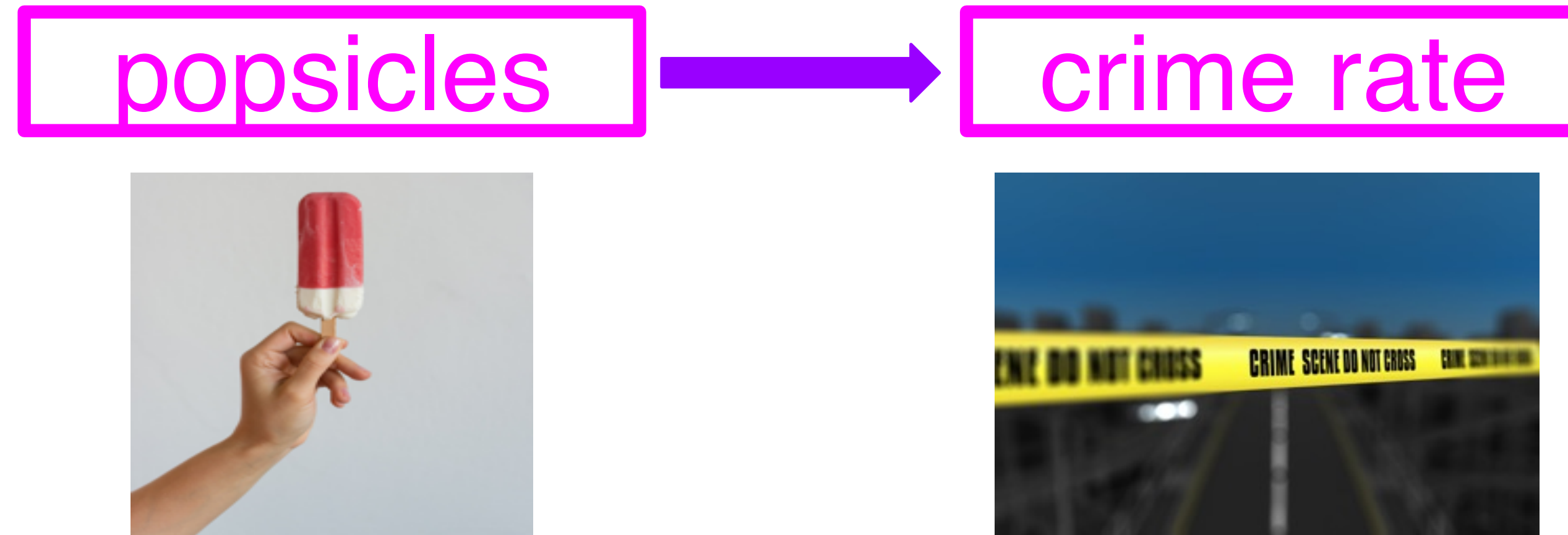


**p-value** : the probability of getting the  
observed results (or results more  
extreme) by chance alone

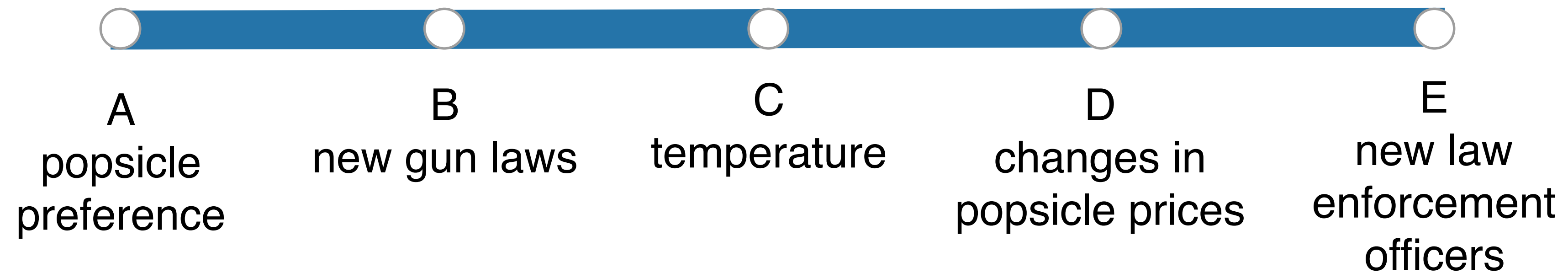
Confounding



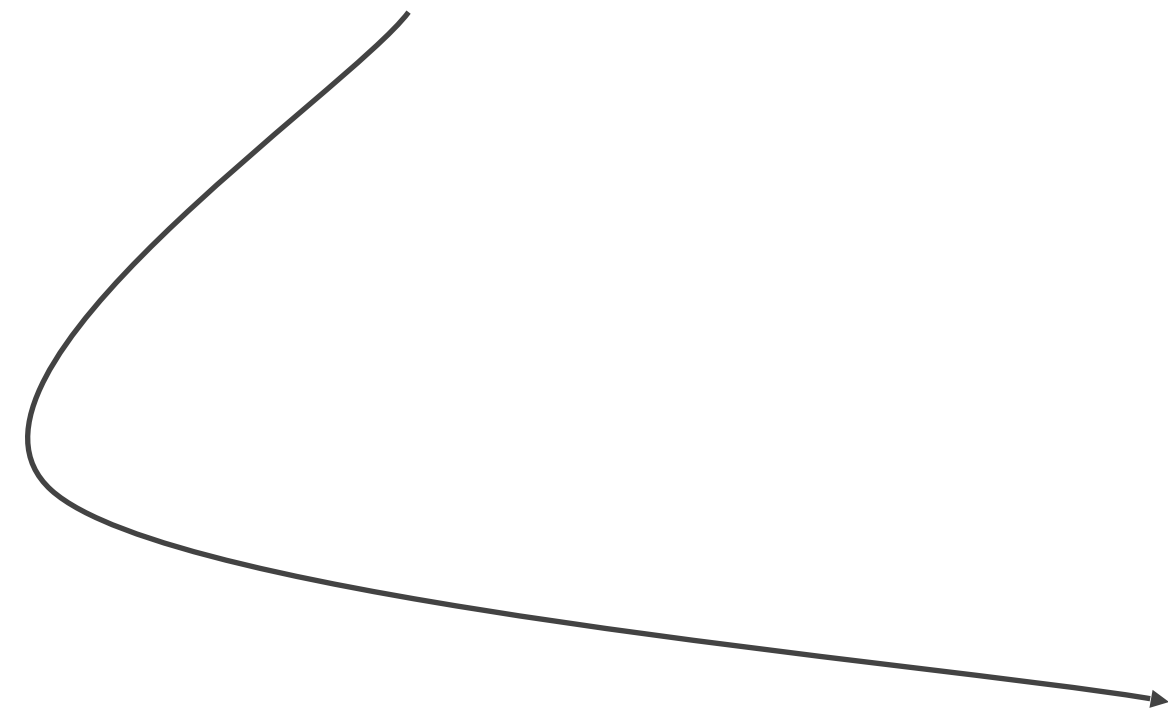
# Confounding



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

<b>Vertebroplasty</b>	<b>Conservative care</b>	<b>Relative risk (95% confidence interval)</b>
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 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