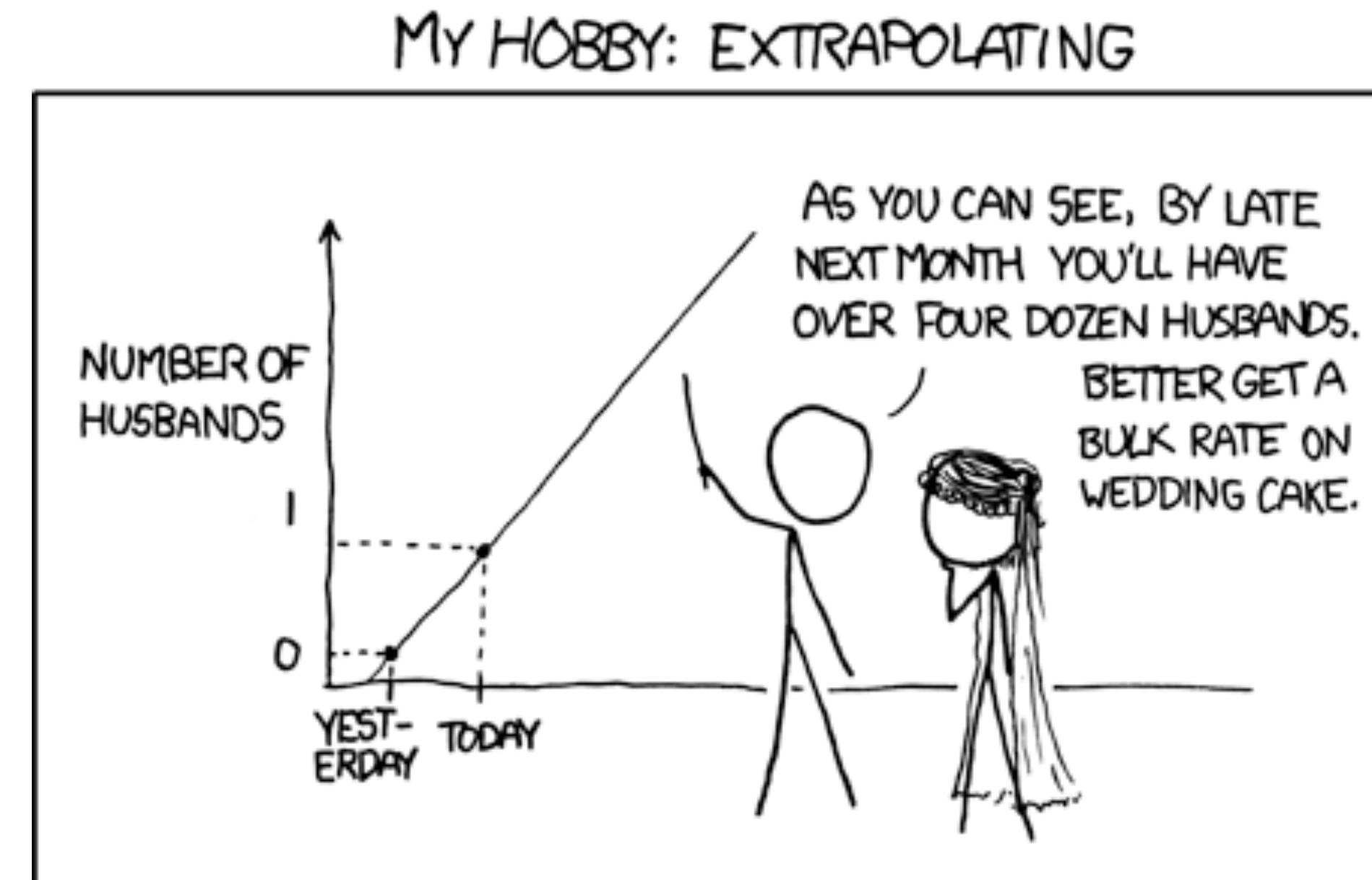


# Linear model 1



# **Things that came up**

# Ed discussion

## Recording for Jan 26 #11



Anonymous  
21 hours ago in Lectures



27

PIN

STAR

WATCH

VIEWS



Hi professors, can you share the recording link for Jan 26th's? It is not yet updated on Canvas > Pages.  
Thank you!

Comment Edit Delete Endorse ...

## 1 Answer



Tobi Gerstenberg STAFF

3 hours ago



Unfortunately, I forgot to record this lecture. I'm sorry about that! We've now set it up on zoom that the recording starts automatically so that this doesn't happen again in the future.



Comment Edit Delete Endorse ...

Add comment

# Ed discussion

## interpreting regression tables #12



Maya Provencal

1 hour ago in R/RStudio



PIN



STAR



WATCH

2

VIEWS



a friend just shared this with me and i thought it was a nice little guide for how to interpret your regression output in r :)

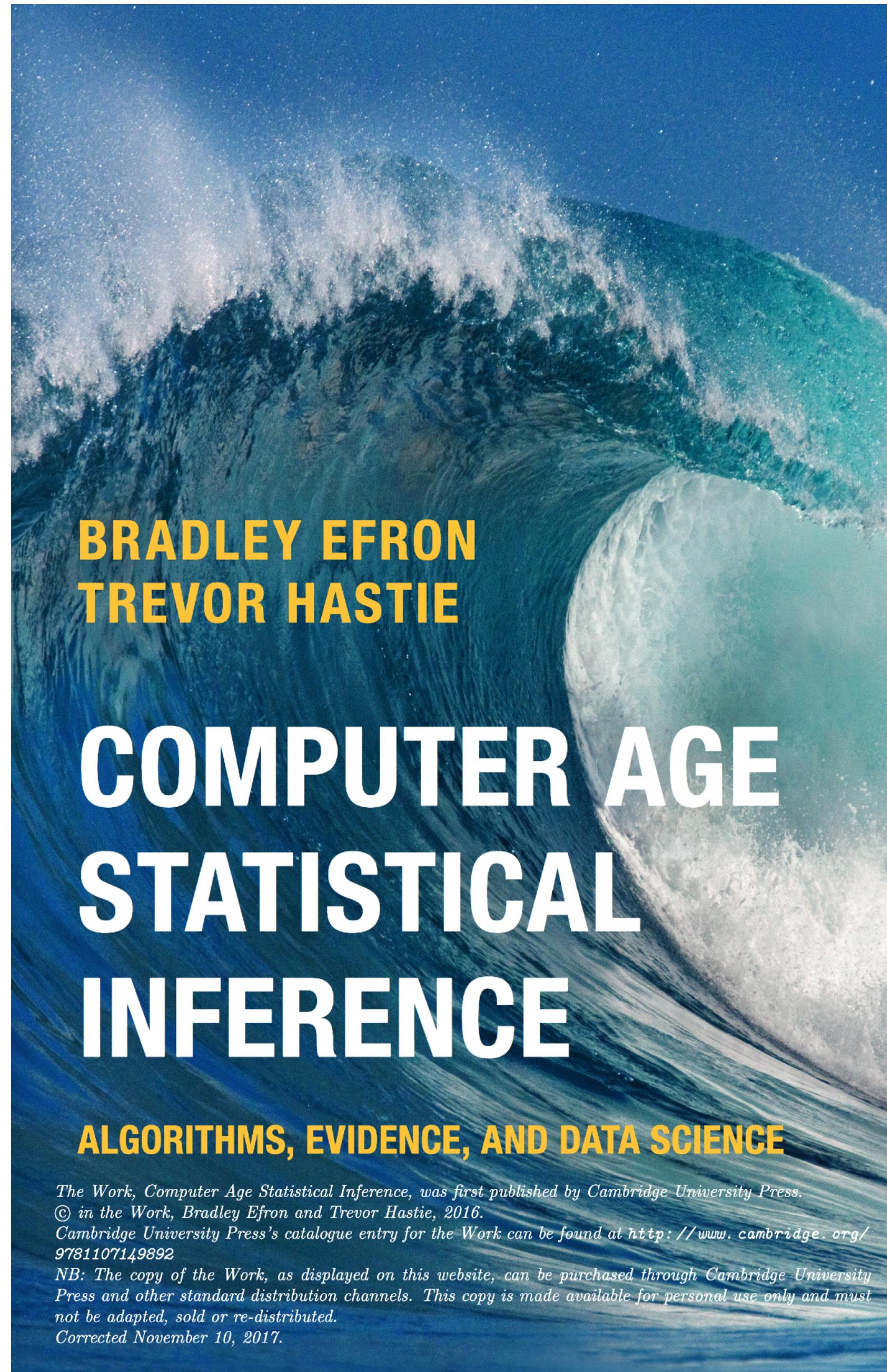
link: [https://jonathanold.github.io/pdf/teaching/Interpreting\\_regression\\_tables.pdf](https://jonathanold.github.io/pdf/teaching/Interpreting_regression_tables.pdf)

Comment Edit Delete Endorse ...

Add comment

[https://jonathanold.github.io/pdf/teaching/Interpreting\\_regression\\_tables.pdf](https://jonathanold.github.io/pdf/teaching/Interpreting_regression_tables.pdf)

# More on bootstrapping



## Tobi Gerstenberg

**Pronouns:** he/him

**Department:** Psychology

**Interests:**

- running the [Causality in Cognition Lab](#)
- computational modeling

**Fun facts:**

- I like surfing and R



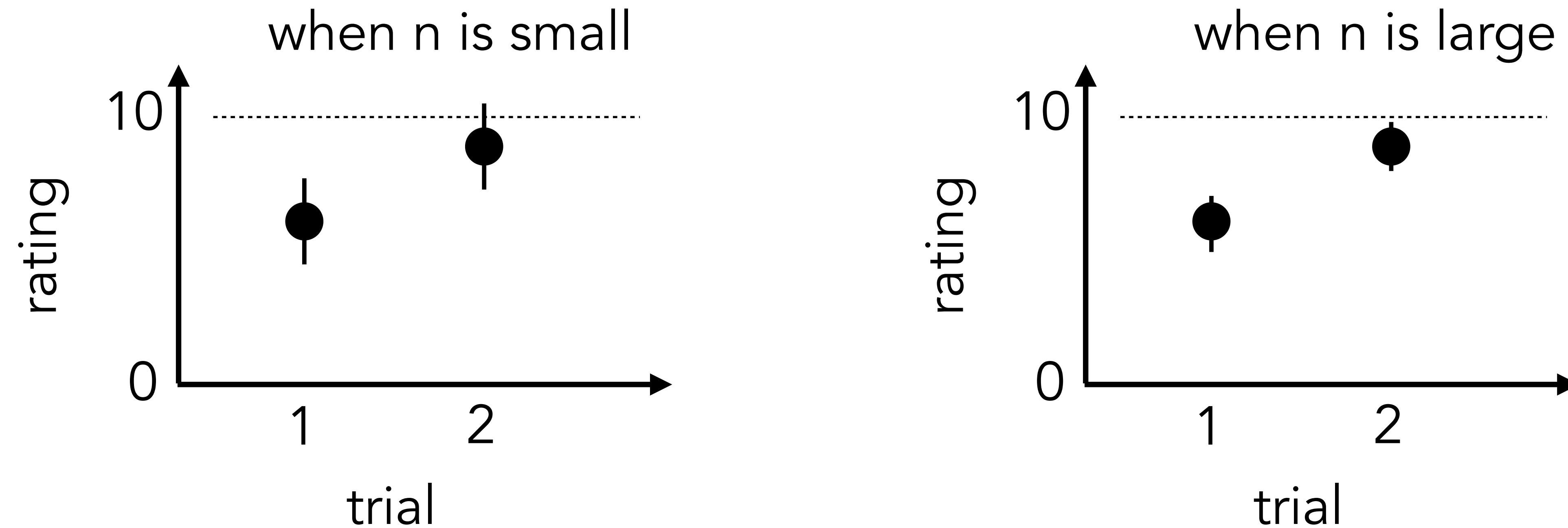
# Confidence intervals: Bootstrap vs. analytic

| analytic   | bootstrap  |
|--|--|
| <b>computationally efficient</b>   | <b>no distributional assumption</b><br>(You don't need normality. For skewed data, heavy tails, or <b>bounded variables</b> , bootstrap intervals adapt to the actual shape of your data. The normal-based CI assumes symmetric errors, which can be wrong.) |
| <b>exact coverage (no approximation)</b>   | <b>works for almost any statistic</b><br>(median, ratio, correlation, regression coefficient, ... ; bootstrap handles them all with the same procedure)  |
| <b>better small sample performance for means</b> (t-distribution for $n < 15$ ); empirical distribution is a crude approximation to the population)  | <b>captures the right asymmetry</b><br>(if the sampling distribution of your statistic is skewed, the bootstrap interval will be asymmetric; normal-based intervals are always symmetric.)   |
| <b>analytical tractability</b><br>(formulas let you derive properties mathematically—how the interval width scales with $n$ , how it responds to variance changes, etc. this aids study design and power analysis) |  |

# Central limit theorem and bootstrapping

If you take a sample of  $n$  independent observations from **\*any\*** distribution with finite mean  $\mu$  and finite variance  $\sigma^2$ , then **as  $n$  grows large**, the **distribution of the sample mean** approaches a normal distribution.

The original distribution can be **anything**—skewed, bimodal, discrete, bounded, whatever. Yet the distribution of the mean of many such observations tends toward the same bell-shaped normal.



The CLT is the bridge between "I have data from some unknown process" and "I can use normal-distribution-based tools to make inferences about it."

# Plan for today

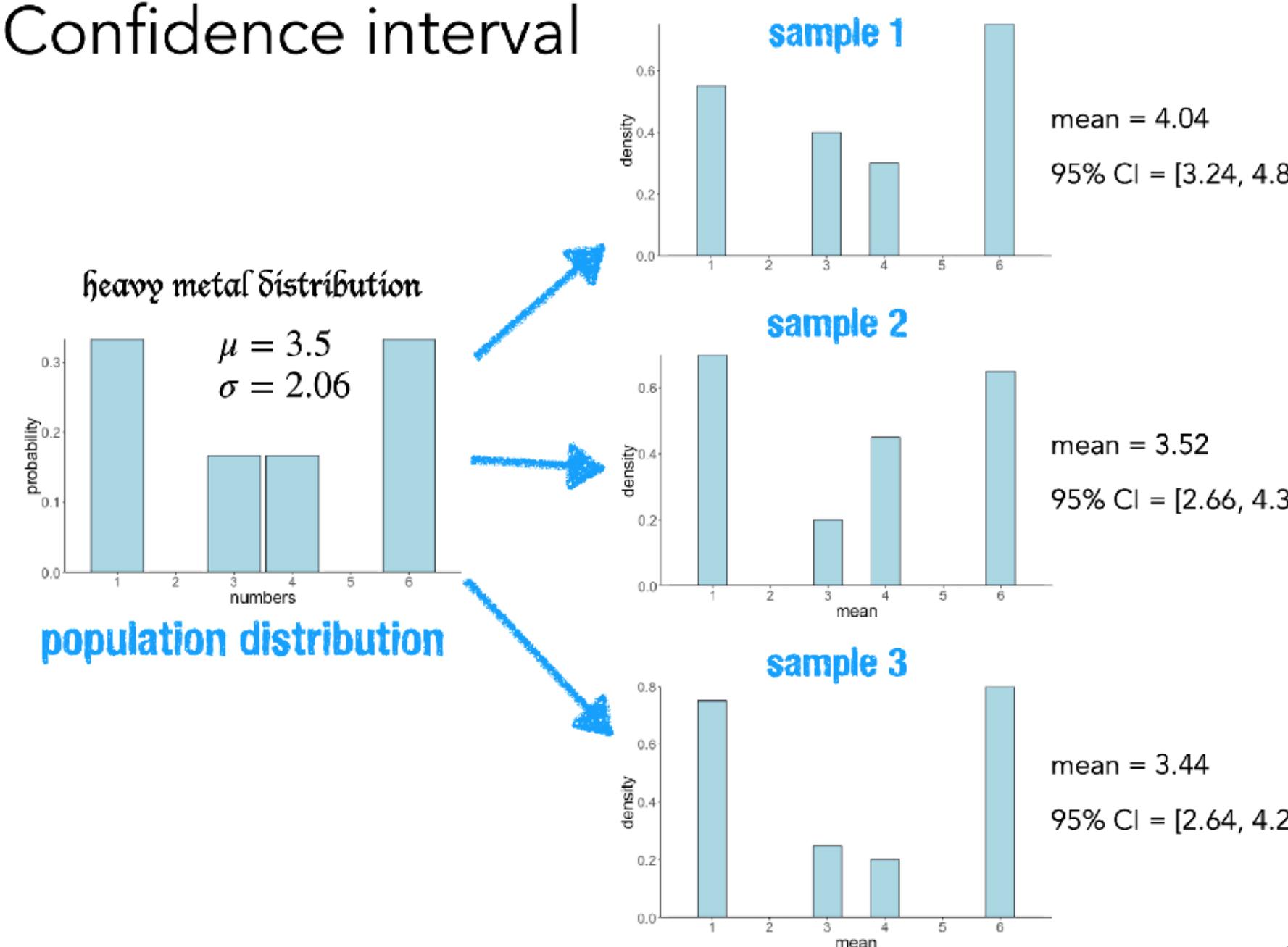
- Quick recap
- Modeling data
- Hypothesis testing as model comparison
- Correlation
  - Pearson's moment correlation
  - Spearman's rank correlation
- Regression

# Quick recap

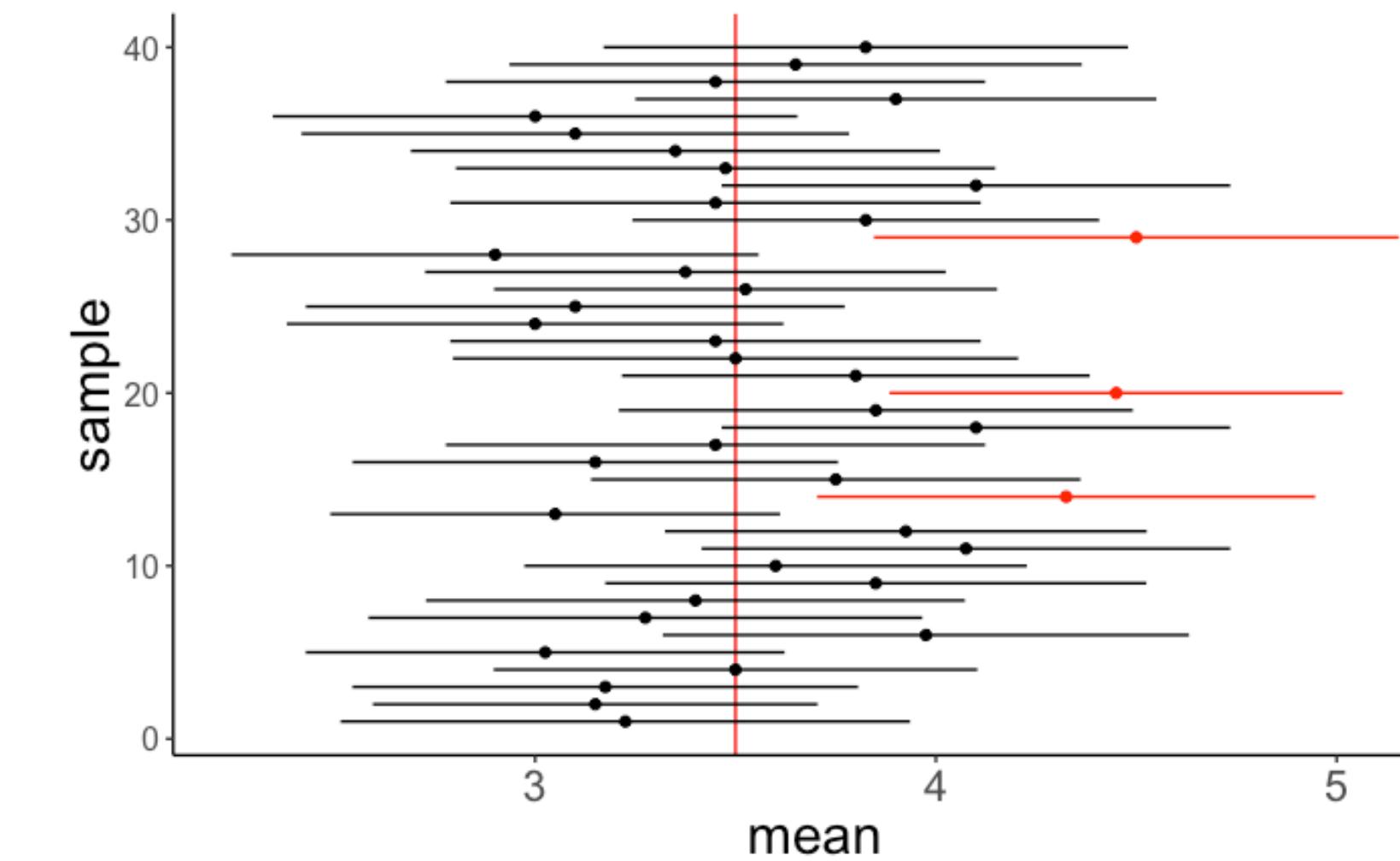
# Quick recap: Confidence intervals

"If we were to repeat the experiment over and over, then 95% of the time the confidence intervals contain the estimate of interest."

Confidence interval

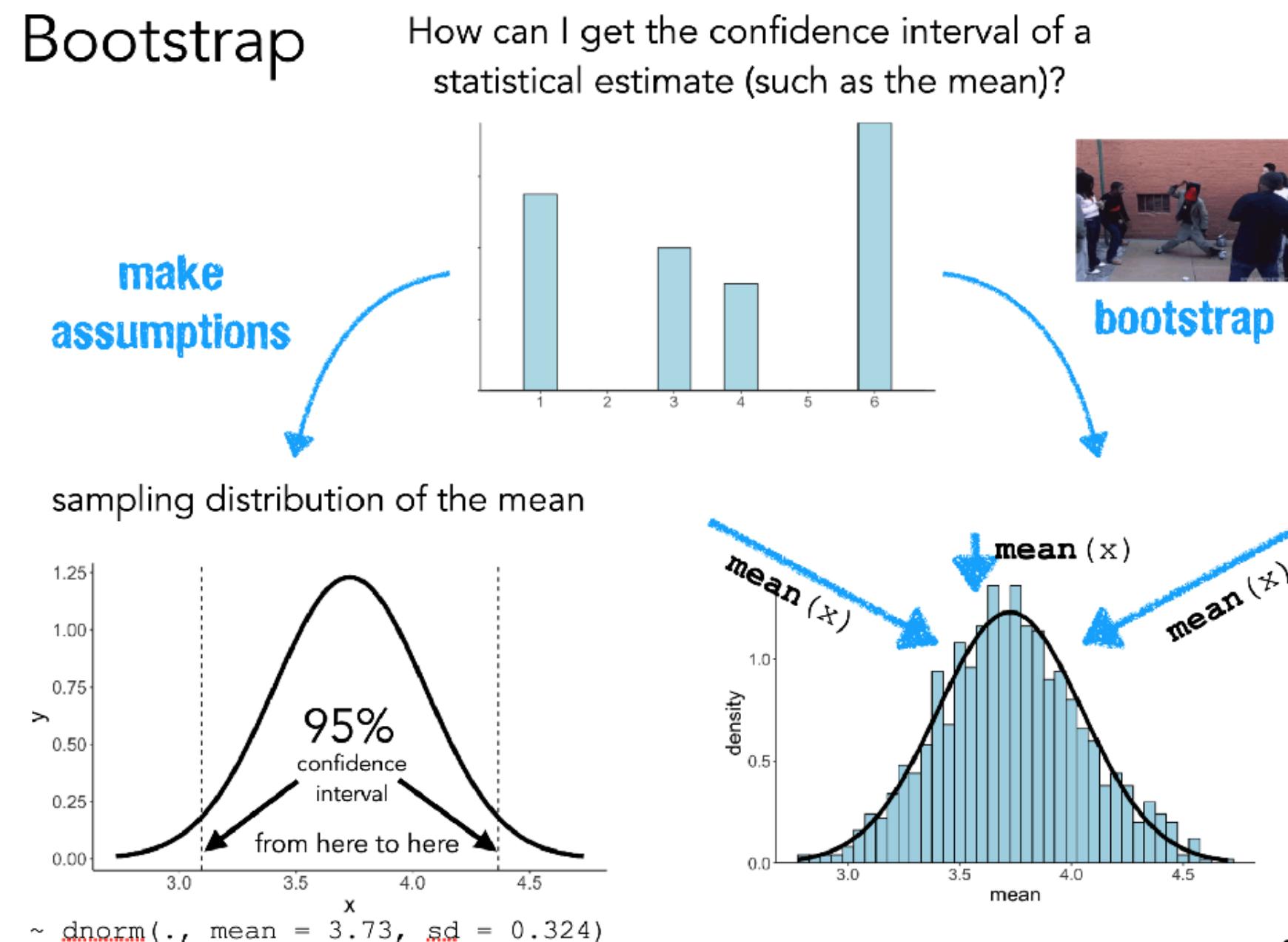
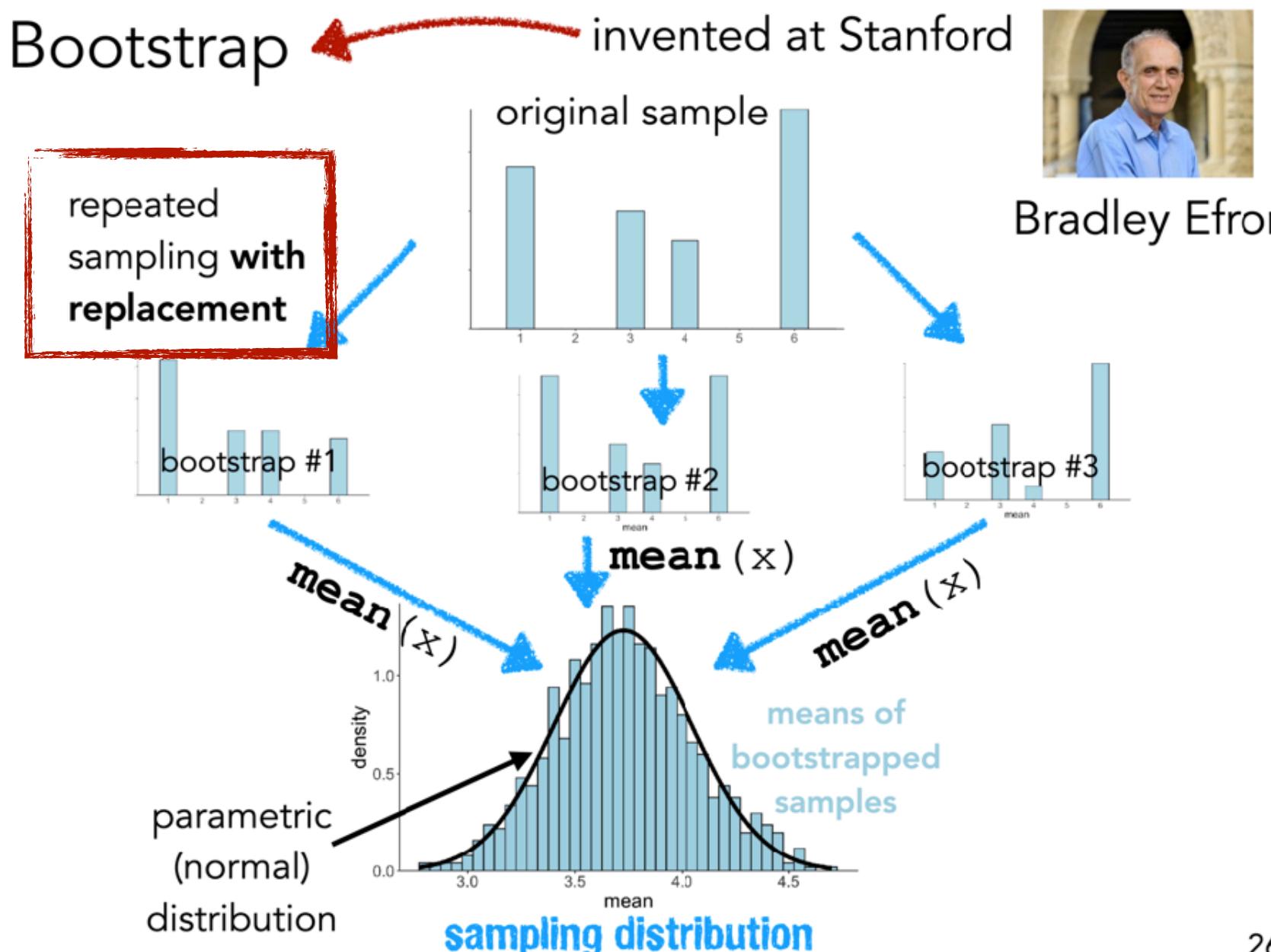


20



<https://www.the100.ci/2024/12/05/why-you-are-not-allowed-to-say-that-your-95-confidence-interval-contains-the-true-parameter-with-a-probability-of-95/>

# Quick recap: Bootstrapping



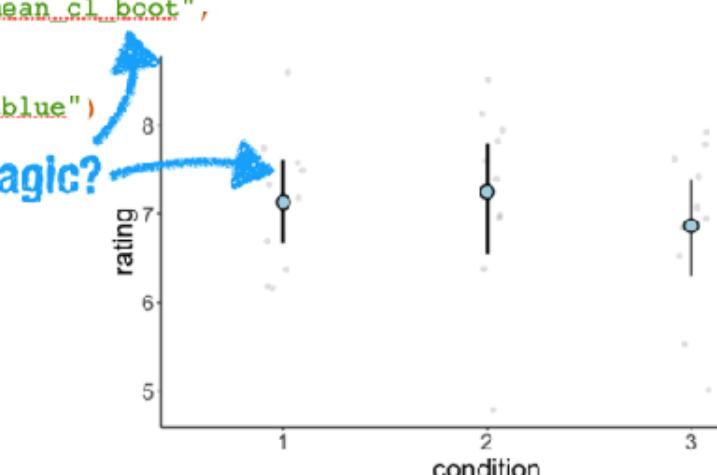
make sure to set the seed!

mean\_cl\_boot() explained

```
1 set.seed(1)
2
3 n = 10 # sample size per group
4 k = 3 # number of groups
5
6 df.data = tibble(participant = 1:(n*k),
7   condition = as.factor(rep(1:k, each = n)),
8   rating = rnorm(n*k, mean = 7, sd = 1))
9
10 ggplot(data = df.data,
11   mapping = aes(x = condition,
12     y = rating)) +
13   geom_point(alpha = 0.1,
14     position = position_jitter(width = 0.1, height = 0)) +
15   stat_summary(fun.data = "mean_cl_boot",
16     shape = 21,
17     size = 1,
18     fill = "lightblue")
```

| participant | condition | rating |
|-------------|-----------|--------|
| 1           | 1         | 6.37   |
| 2           | 1         | 7.18   |
| 3           | 1         | 6.16   |
| 4           | 1         | 6.60   |
| 5           | 1         | 7.33   |

what is this magic?



# Quick recap: Modeling data

$$\text{Data} = \text{Model} + \text{Error}$$

↑  
what makes for  
a good model?

- we build models with parameters, and fit those parameters to **minimize error**
- adding additional parameters to the model will always improve the model fit and reduce error
- fundamental trade-off between **simplicity** and **accuracy**

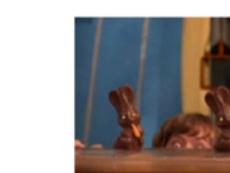
Assumption of normal distribution

$$\text{Error} = \text{Data} - \text{Model}$$

↑  
assumed to be  
normally  
distributed

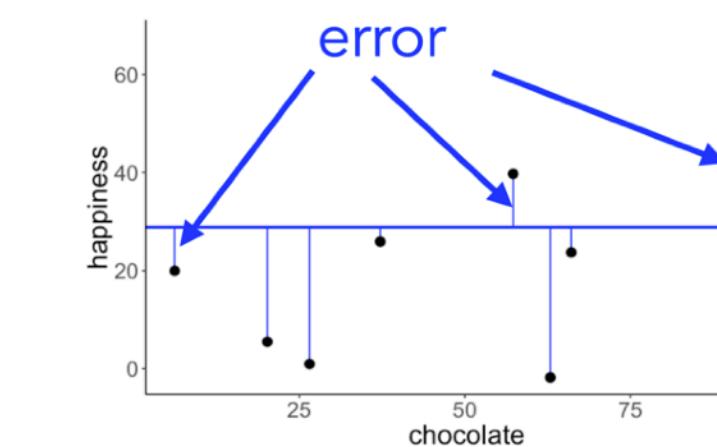
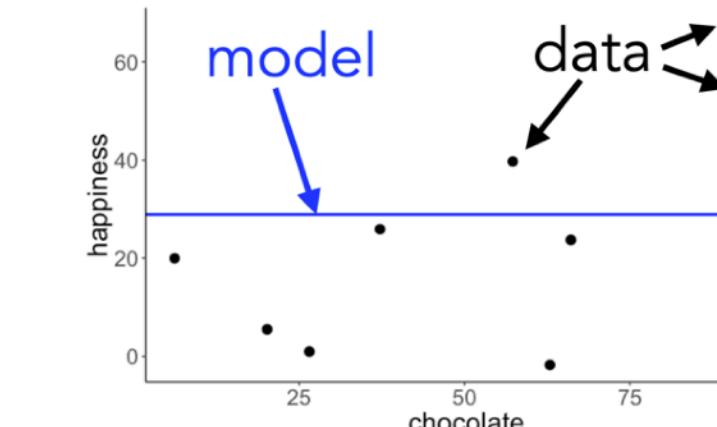
↑  
don't need to  
be normally  
distributed!!

very common misconception!!!



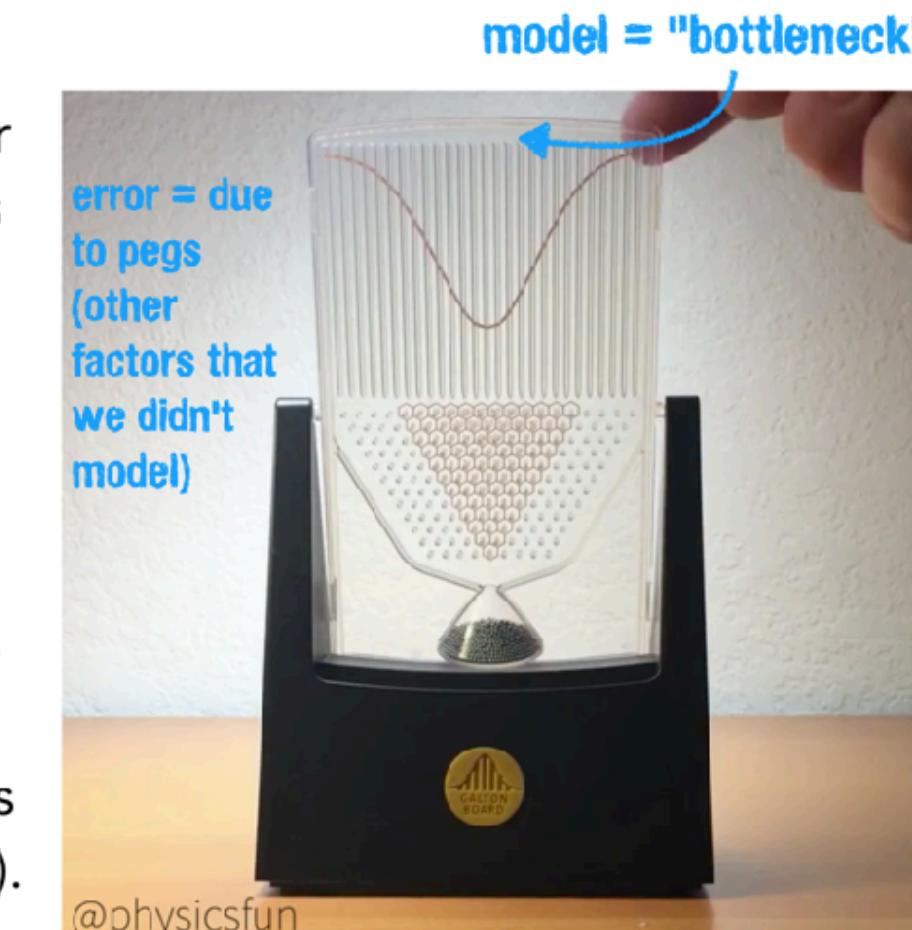
$$\text{Data} = \text{Model} + \text{Error}$$

$H_0$ : Chocolate consumption and happiness are unrelated.



## ERROR

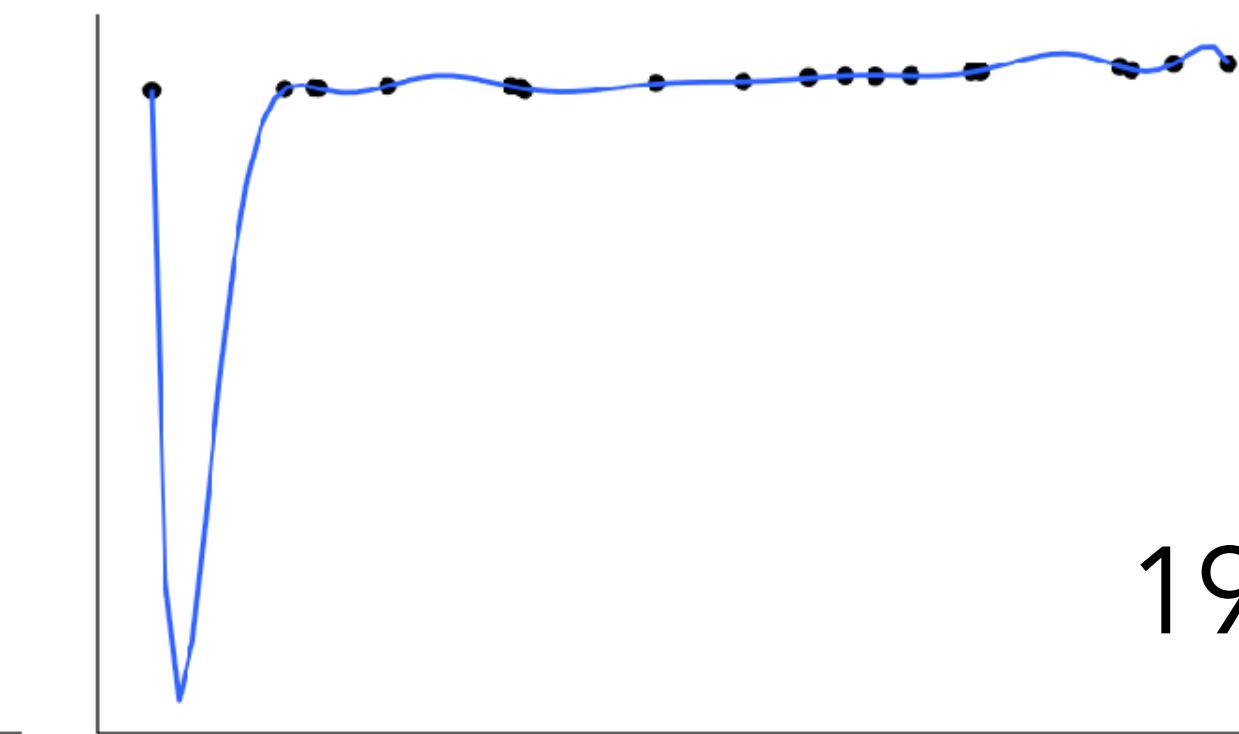
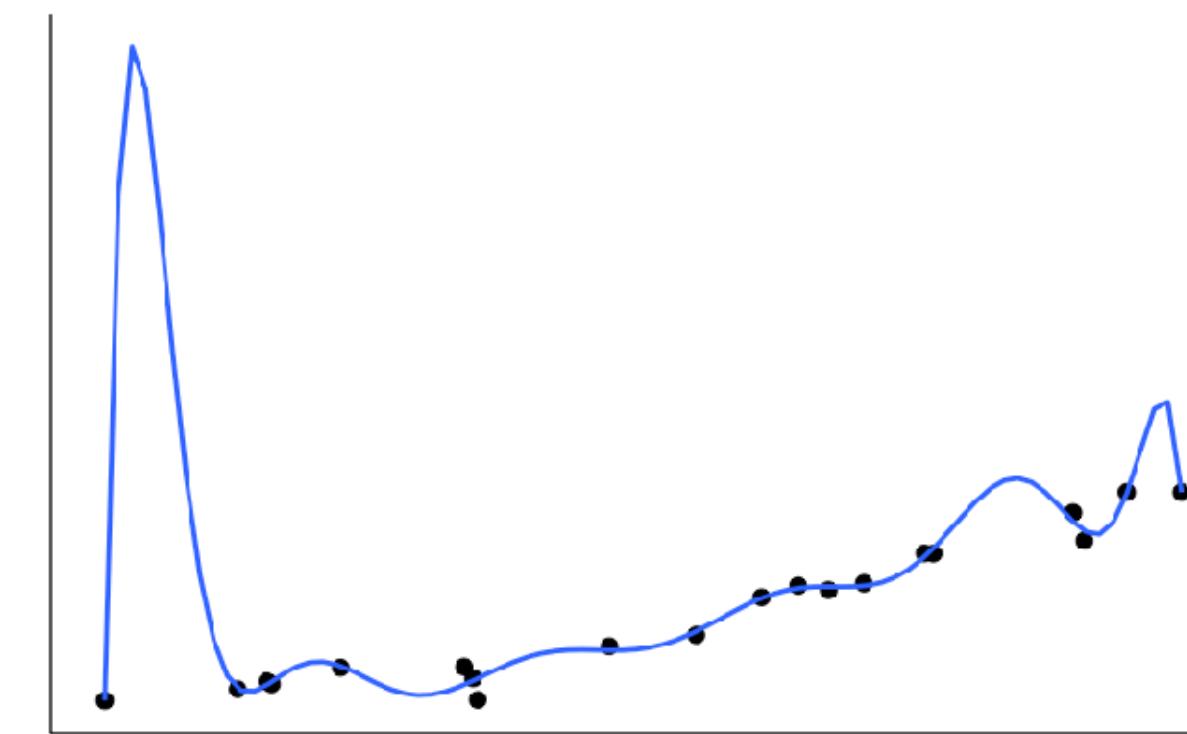
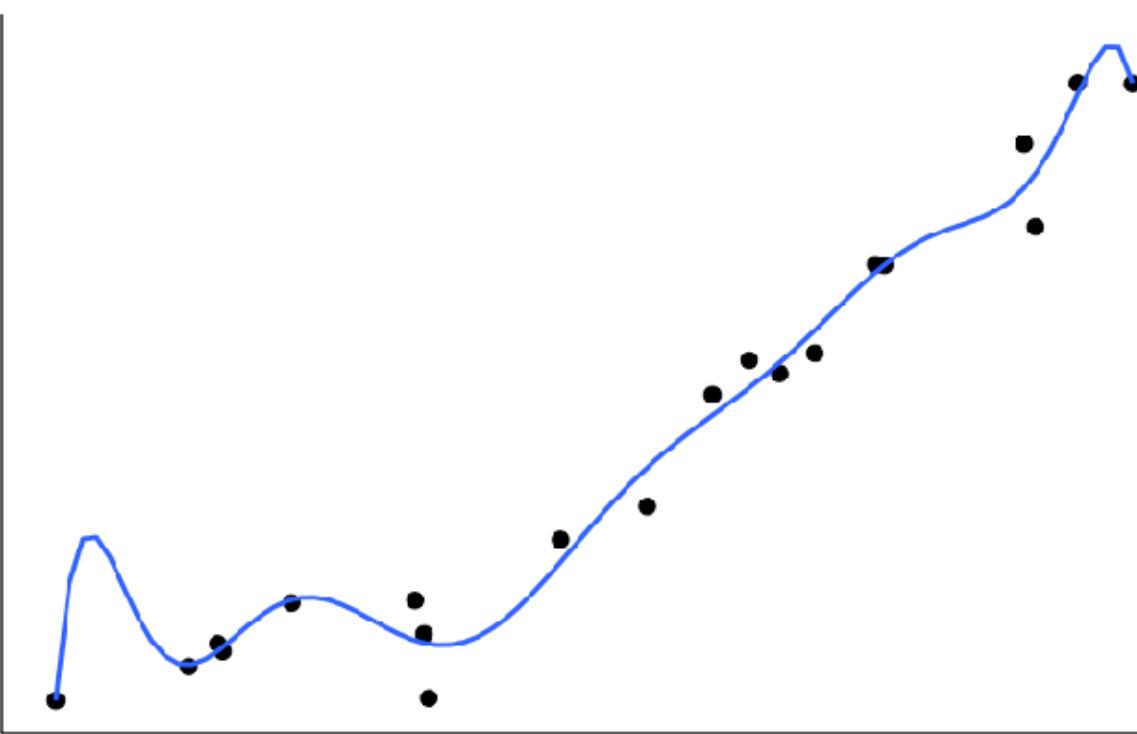
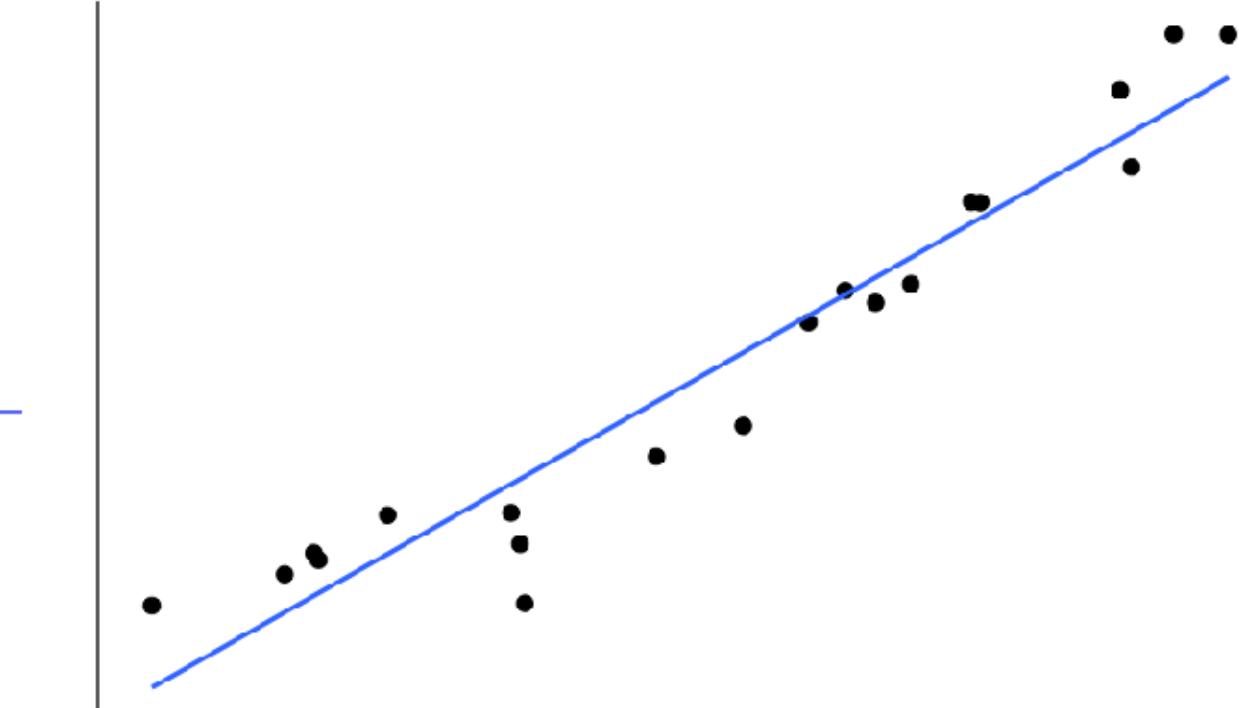
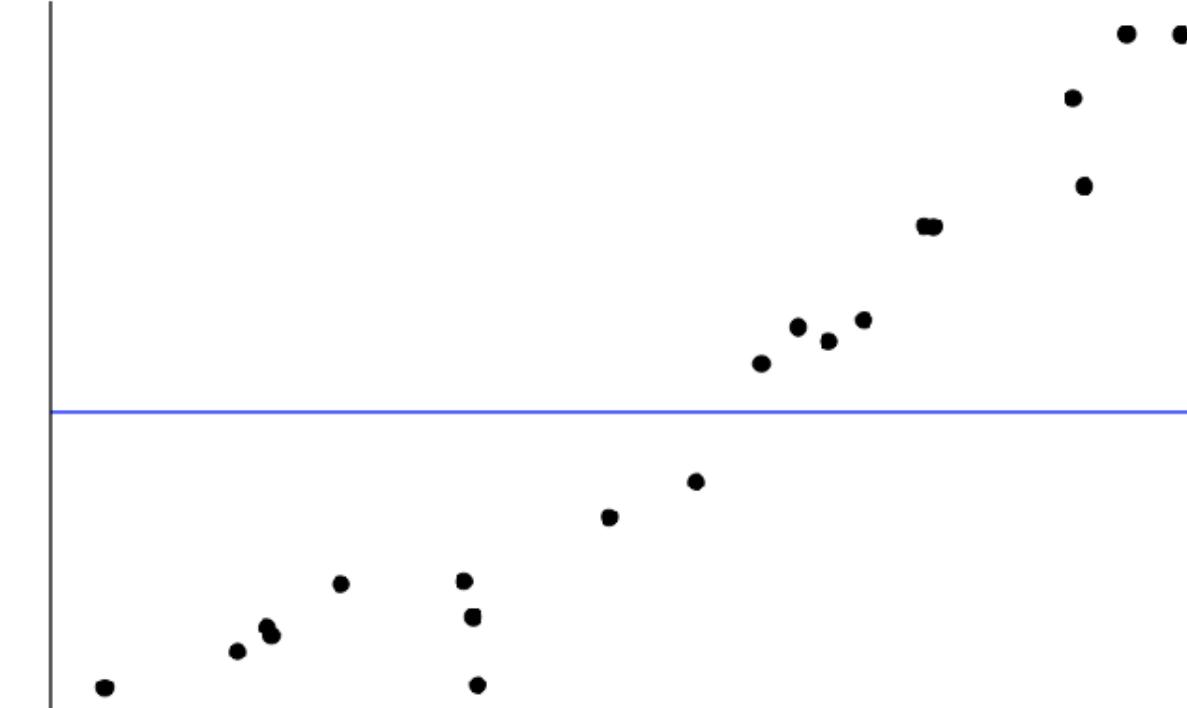
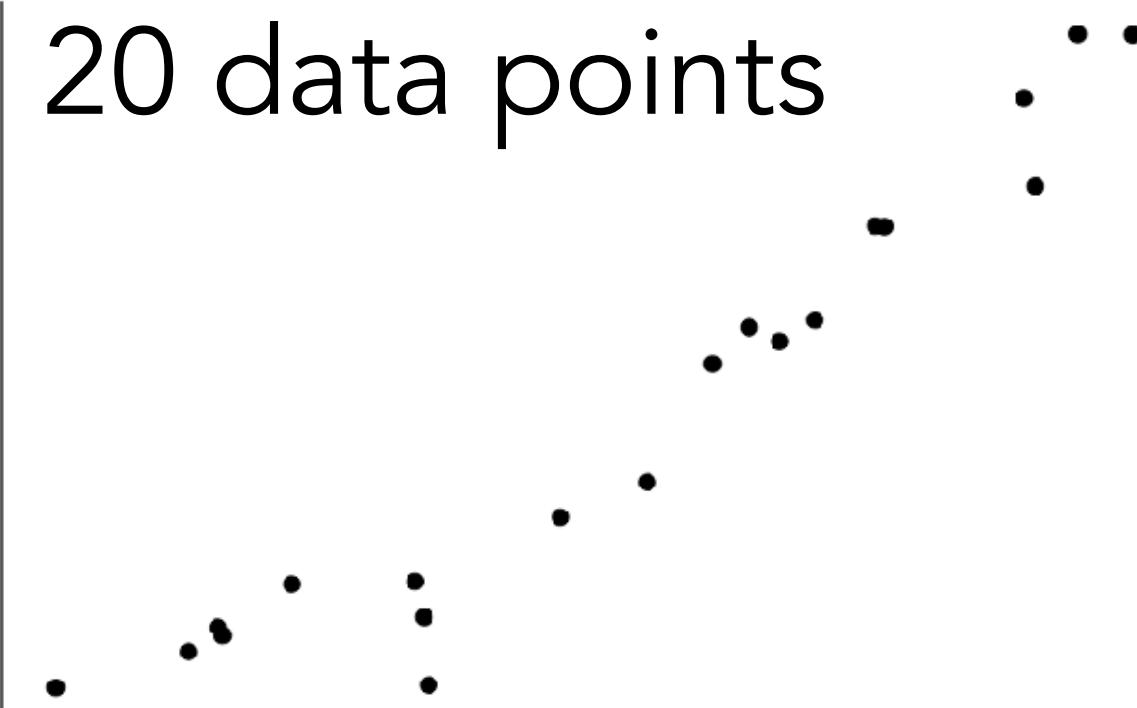
1. We assume that the error between model and data is due to (a potentially large number of) factors that we didn't take into account.
2. We assume that each of these factors influences the data in **an additive way** (some pulling in one, others pulling in another direction).



data = where the balls land

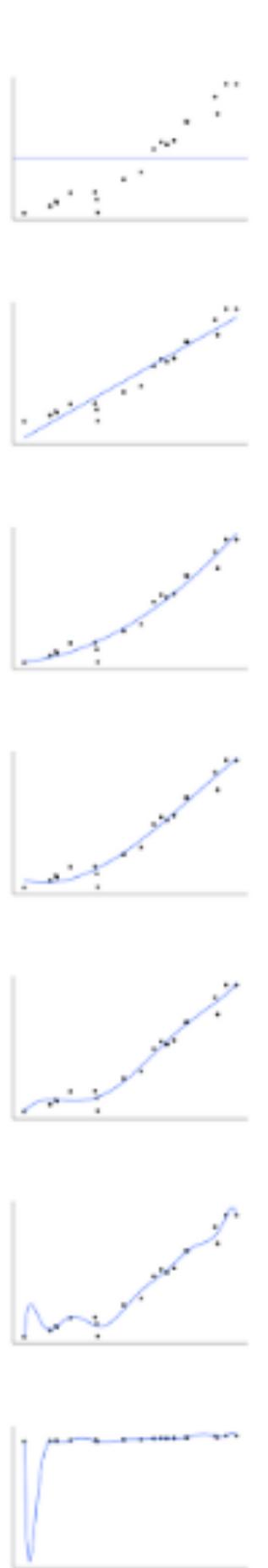
Result: normal distribution

# Modeling data



19 parameters

# Which model describes the data best





# Example

Percentage of households that had internet access in the year 2013 by US state

| i  | internet | state | college | auto | density |
|----|----------|-------|---------|------|---------|
| 1  | 79.0     | AK    | 28.0    | 1.2  | 1.2     |
| 2  | 63.5     | AL    | 23.5    | 1.3  | 94.4    |
| 3  | 60.9     | AR    | 20.6    | 1.7  | 56.0    |
| 4  | 73.9     | AZ    | 27.4    | 1.3  | 56.3    |
| 5  | 77.9     | CA    | 31.0    | 0.8  | 239.1   |
| 6  | 79.4     | CO    | 37.8    | 1.0  | 48.5    |
| 7  | 77.5     | CT    | 37.2    | 1.0  | 738.1   |
| 8  | 74.5     | DE    | 29.8    | 1.1  | 460.8   |
| 9  | 74.3     | FL    | 27.2    | 1.2  | 350.6   |
| 10 | 72.2     | GA    | 28.3    | 1.1  | 168.4   |

set before looking at the data

$$\text{model}_1: Y_i = 75 + \text{ERROR}$$

worth it?

fit to the data

$$\text{model}_2: Y_i = \beta_0 + \text{ERROR}$$

worth it?

additional predictor

$$\text{model}_3: Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$$

fit to the data

## Compact model

$\text{model}_C: Y_i = \beta_0 + \text{ERROR}$

## Augmented model

$\text{model}_A: Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$

$$\text{ERROR}(C) \geq \text{ERROR}(A)$$

## Proportional reduction in error (PRE)

$$\text{PRE} = \frac{\text{ERROR}(C) - \text{ERROR}(A)}{\text{ERROR}(C)}$$

## Compact model

$$\text{model}_C: Y_i = \beta_0 + \text{ERROR}$$

$$\text{ERROR}(C) = 50$$

## Augmented model

$$\text{model}_A: Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$$

$$\text{ERROR}(A) = 30$$

## Proportional reduction in error (PRE)

$$\text{PRE} = 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)}$$

$$= 1 - \frac{30}{50} = .40$$

Increasing the complexity of the model reduced the error by 40%. **worth it?**

# worth it?

## Compact model

model<sub>C</sub>:  $Y_i = \beta_0 + \text{ERROR}$

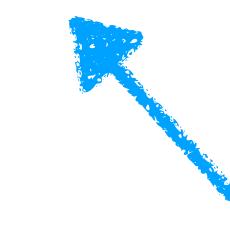
## Augmented model

model<sub>A</sub>:  $Y_i = \beta_0 + \beta_1 X_i + \text{ERROR}$

## Proportional reduction in error (PRE)

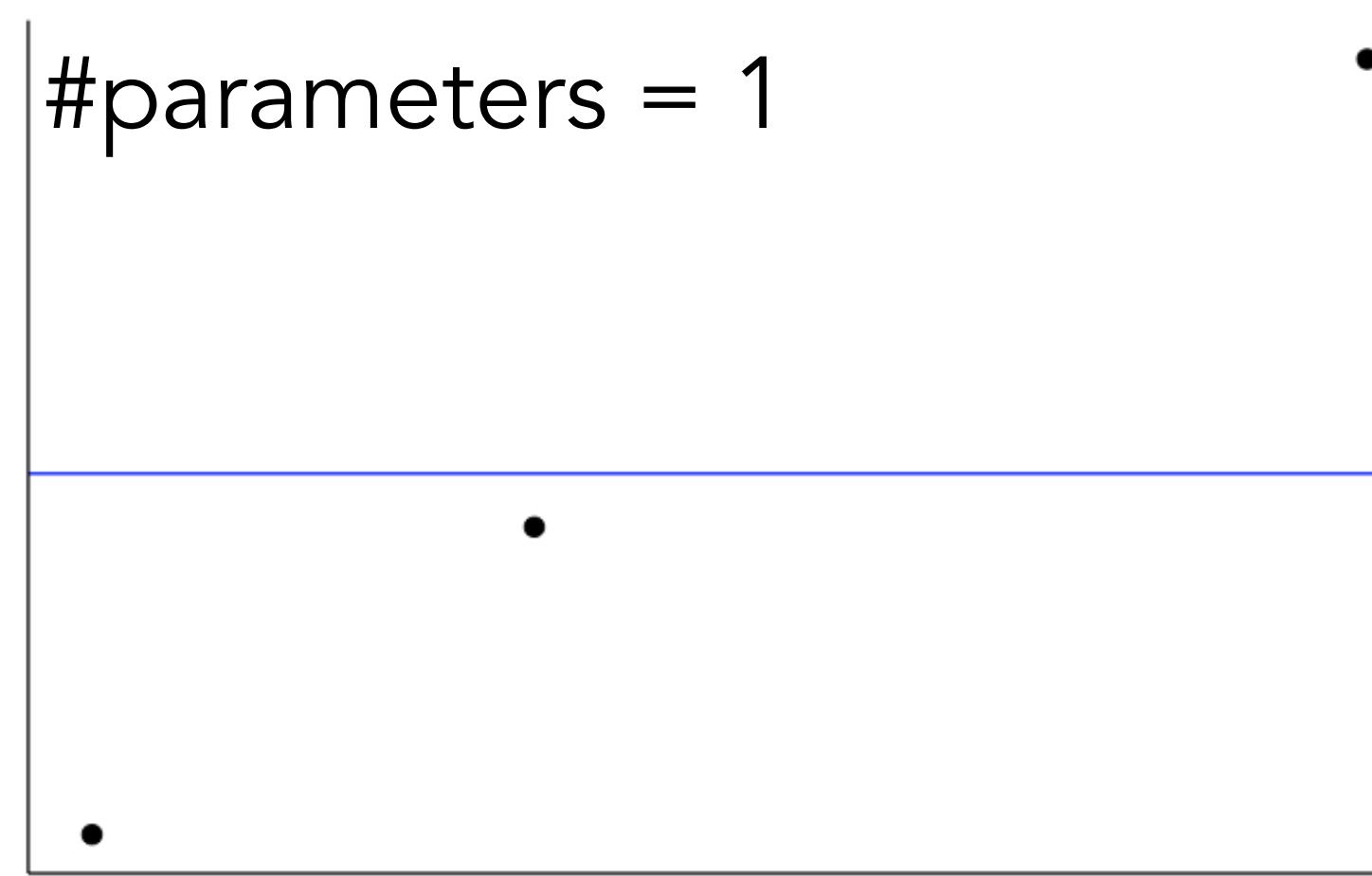
$$\begin{aligned}\text{PRE} &= 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)} \\ &= 1 - \frac{30}{50} = .40\end{aligned}$$

- to answer the **worth it?** question we need inferential statistics (define ERROR, sampling distributions, etc.)
- more likely to be **worth it** if:
  1. **PRE** is high
  2. the number of additional parameters in A compared to C is low
  3. the number of parameters that could have been added to model<sub>C</sub> to create model<sub>A</sub> but were not is high

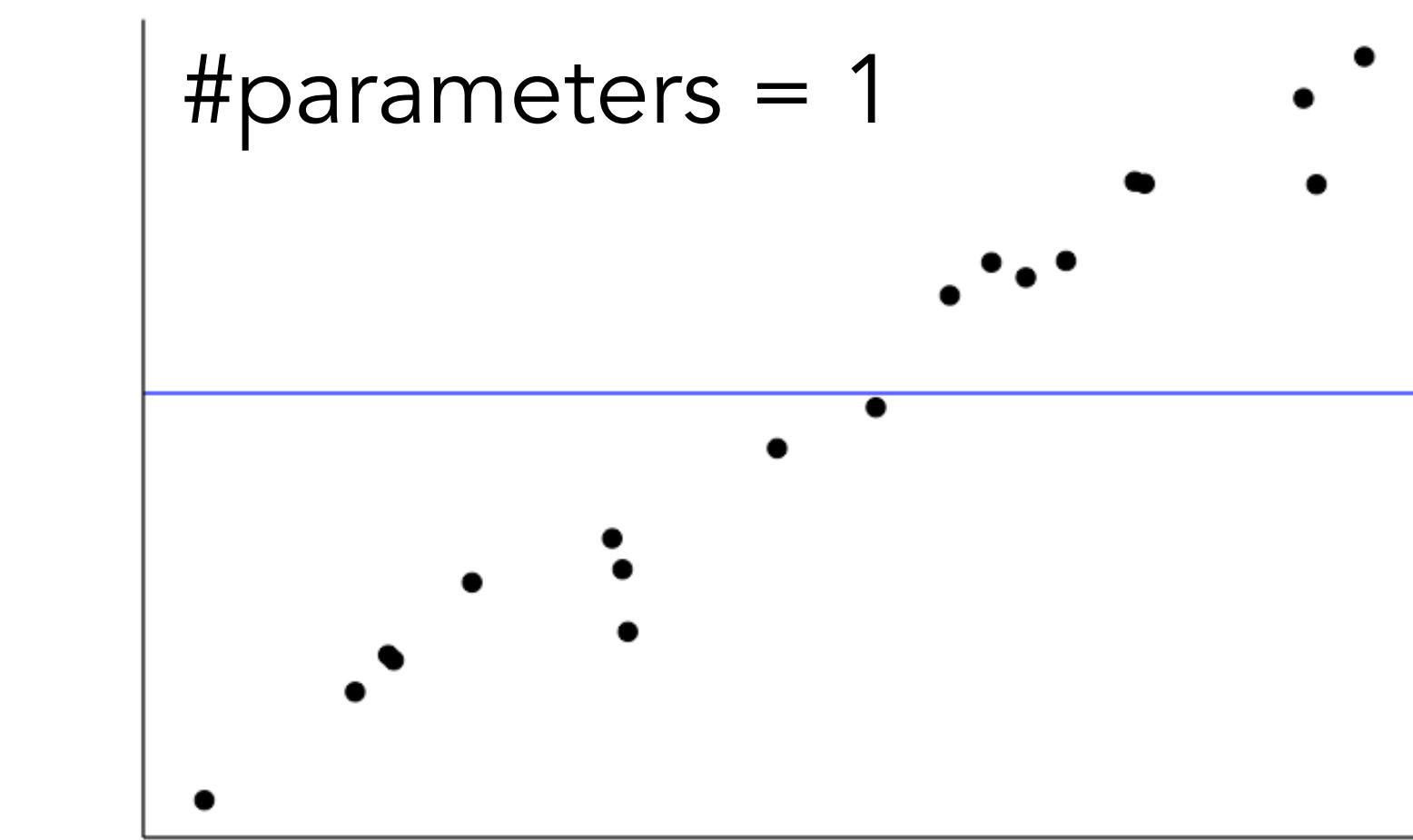
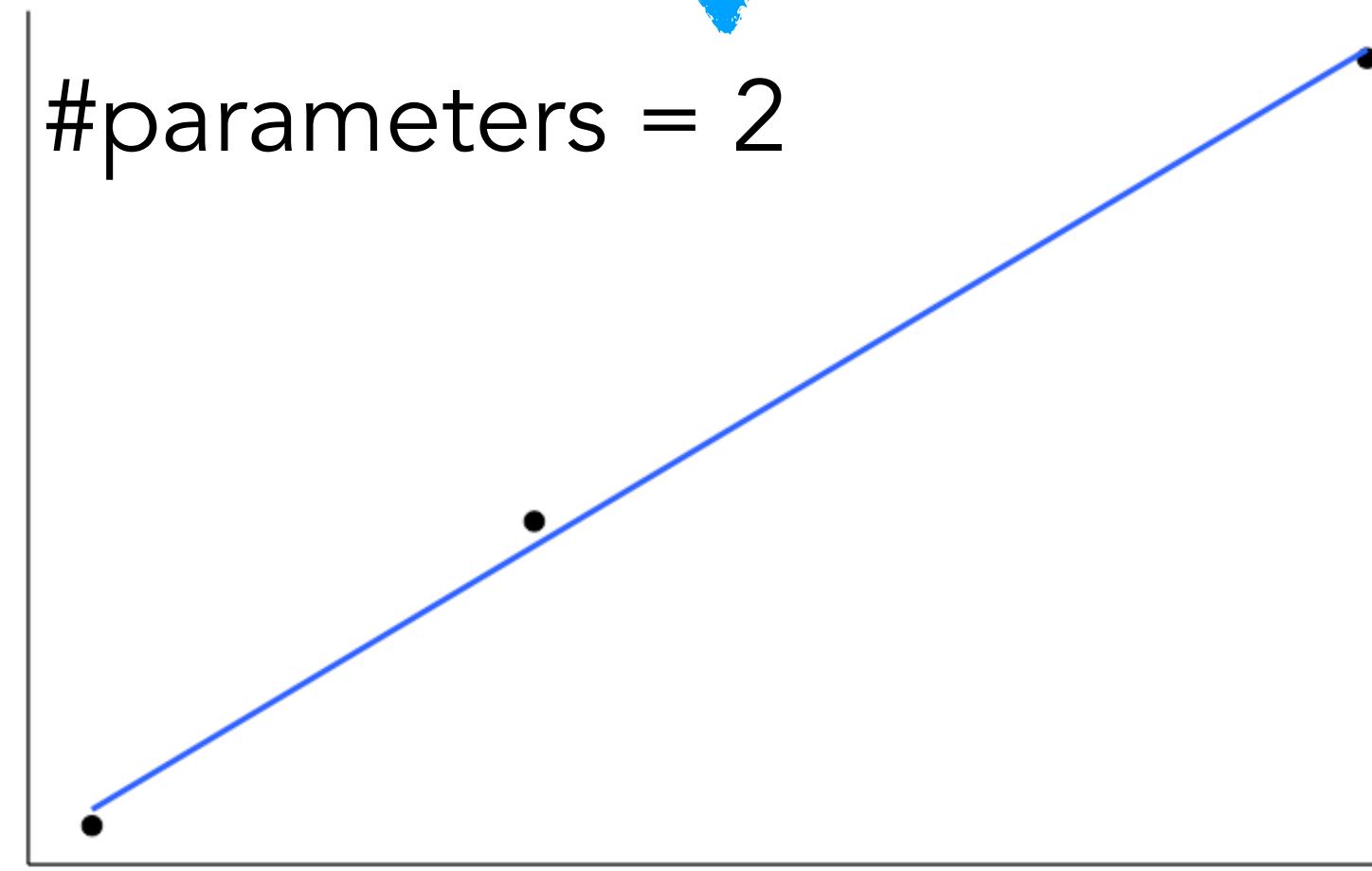


more impressed if the number of observations n is much greater than the number of parameters

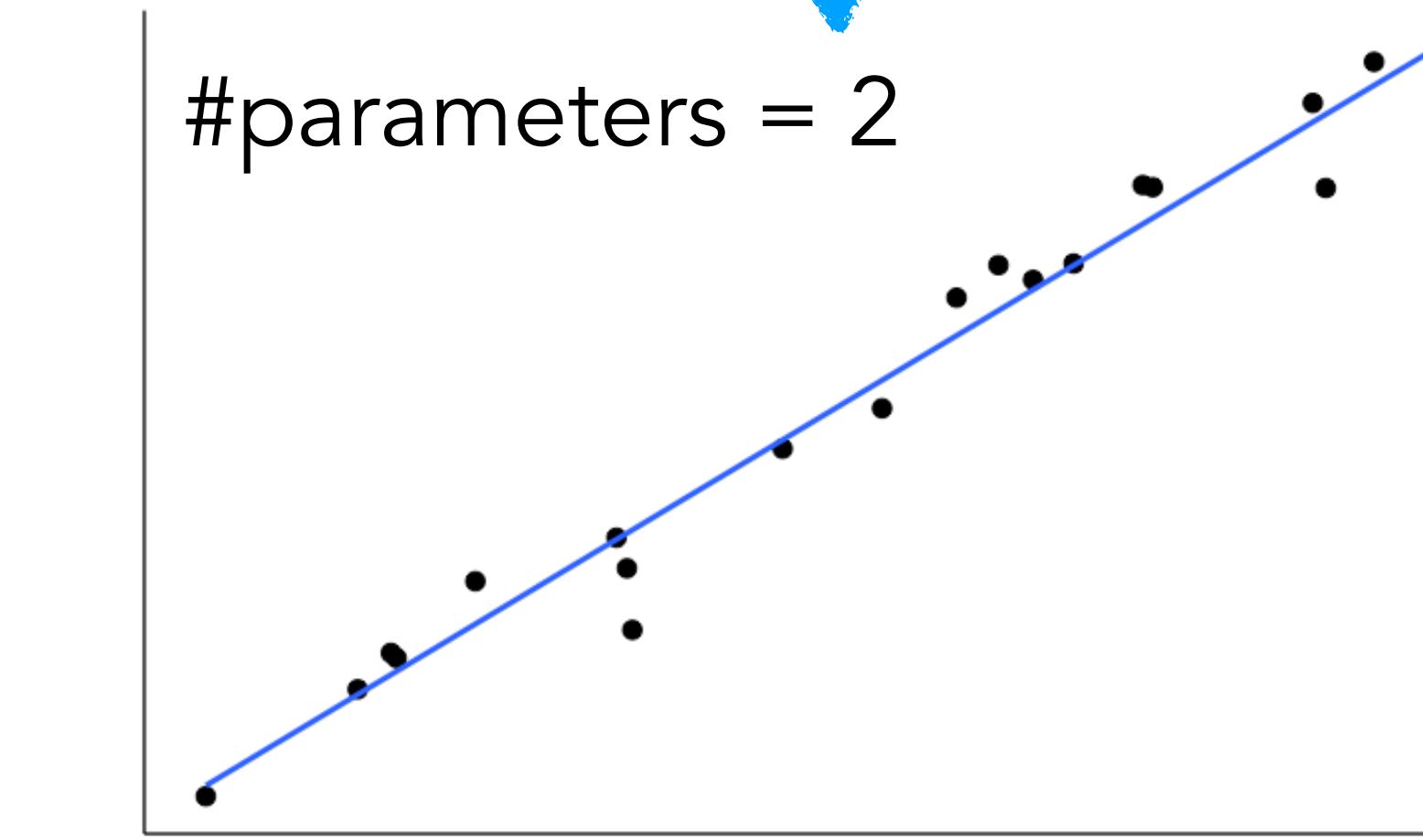
# PRE per parameter for different $n$



↓ neato!



↓ impressive!



# General procedure

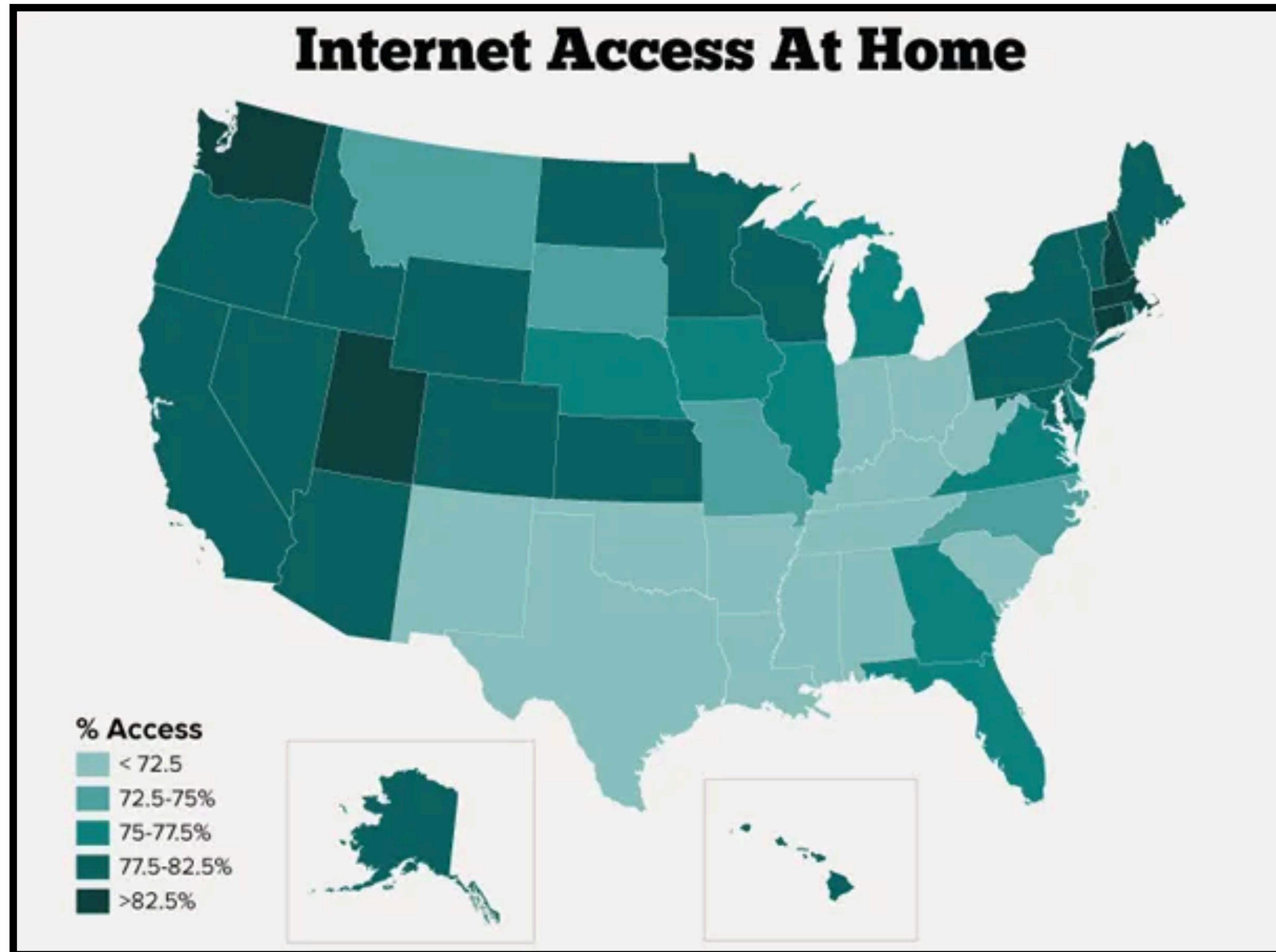
- for any question we want to ask about our DATA
    - we define model<sub>C</sub> and model<sub>A</sub>
    - compare the models using PRE
    - determine whether PRE is **worth it**
  - in standard frequentist lingo:
    - model<sub>C</sub> =  $H_0$  (null hypothesis)
    - model<sub>A</sub> =  $H_1$  (alternative hypothesis)
  - hypothesis test:
    - $H_0$ : **all** the parameters that are included in model<sub>A</sub> but not in model<sub>C</sub> are 0
    - $H_1$ : **not all** the parameters that are included in model<sub>A</sub> but not in model<sub>C</sub> are 0
- 
- model comparison**

# Hypothesis testing as model comparison

# Procedure

1. Start with research question
2. Formulate hypothesis as a comparison between compact and augmented model
3. Fit parameters in each model
4. Calculate the proportional reduction of error (PRE)
5. Decide whether PRE is **worth it**

# Internet Access At Home



<http://thedataviz.com/2012/07/19/mapping-internet-access-in-u-s-homes/>

# Notation

DATA = MODEL + ERROR

$$Y_i = \beta_0 + \epsilon_i \quad \text{simple model (true parameters)}$$

$$Y_i = b_0 + e_i \quad \text{simple model (estimated parameters)}$$

$$Y_i = b_0 + b_1 X_{i1} + e_i \quad \text{more complex model}$$

college



Greek letters  $\beta$  or  $\epsilon$  represent the true but unknowable parameters in the population.

Roman letters  $b$  or  $e$  represent estimates of these parameters using our DATA.

Percentage of households that had internet access in the year 2013 by US state

| i  | internet | state | college | auto | density |
|----|----------|-------|---------|------|---------|
| 1  | 79.0     | AK    | 28.0    | 1.2  | 1.2     |
| 2  | 63.5     | AL    | 23.5    | 1.3  | 94.4    |
| 3  | 60.9     | AR    | 20.6    | 1.7  | 56.0    |
| 4  | 73.9     | AZ    | 27.4    | 1.3  | 56.3    |
| 5  | 77.9     | CA    | 31.0    | 0.8  | 239.1   |
| 6  | 79.4     | CO    | 37.8    | 1.0  | 48.5    |
| 7  | 77.5     | CT    | 37.2    | 1.0  | 738.1   |
| 8  | 74.5     | DE    | 29.8    | 1.1  | 460.8   |
| 9  | 74.3     | FL    | 27.2    | 1.2  | 350.6   |
| 10 | 72.2     | GA    | 28.3    | 1.1  | 168.4   |

# Research question and hypotheses

Is the average percentage of internet users per state significantly different from 75%?

Model<sub>C</sub>: 
$$Y_i = B_0 + \epsilon_i$$
  
**0 parameters**

$$Y_i = 75 + e_i$$

Model<sub>A</sub>: 
$$Y_i = \beta_0 + \epsilon_i$$
  
**1 parameter**

$$Y_i = b_0 + e_i$$

$$= \bar{Y} + e_i$$

| i  | internet | state | college | auto | density |
|----|----------|-------|---------|------|---------|
| 1  | 79.0     | AK    | 28.0    | 1.2  | 1.2     |
| 2  | 63.5     | AL    | 23.5    | 1.3  | 94.4    |
| 3  | 60.9     | AR    | 20.6    | 1.7  | 56.0    |
| 4  | 73.9     | AZ    | 27.4    | 1.3  | 56.3    |
| 5  | 77.9     | CA    | 31.0    | 0.8  | 239.1   |
| 6  | 79.4     | CO    | 37.8    | 1.0  | 48.5    |
| 7  | 77.5     | CT    | 37.2    | 1.0  | 738.1   |
| 8  | 74.5     | DE    | 29.8    | 1.1  | 460.8   |
| 9  | 74.3     | FL    | 27.2    | 1.2  | 350.6   |
| 10 | 72.2     | GA    | 28.3    | 1.1  | 168.4   |

# Fit parameters and calculate PRE

$$\mathbf{C: } Y_i = 75 + e_i \quad \mathbf{A: } Y_i = \bar{Y} + e_i$$

| i  | state | internet | compact_b | compact_se | augmented_b | augmented_se |
|----|-------|----------|-----------|------------|-------------|--------------|
| 1  | AK    | 79.0     | 75        | 16.00      | 72.81       | 38.37        |
| 2  | AL    | 63.5     | 75        | 132.25     | 72.81       | 86.60        |
| 3  | AR    | 60.9     | 75        | 198.81     | 72.81       | 141.75       |
| 4  | AZ    | 73.9     | 75        | 1.21       | 72.81       | 1.20         |
| 5  | CA    | 77.9     | 75        | 8.41       | 72.81       | 25.95        |
| 6  | CO    | 79.4     | 75        | 19.36      | 72.81       | 43.48        |
| 7  | CT    | 77.5     | 75        | 6.25       | 72.81       | 22.03        |
| 8  | DE    | 74.5     | 75        | 0.25       | 72.81       | 2.87         |
| 9  | FL    | 74.3     | 75        | 0.49       | 72.81       | 2.23         |
| 10 | GA    | 72.2     | 75        | 7.84       | 72.81       | 0.37         |

$$\text{SSE(C)} = 1595 \quad \text{SSE(A)} = 1355$$

$$\begin{aligned}\text{PRE} &= 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)} \\ &= 1 - \frac{1355}{1595} \approx .15\end{aligned}$$

Model A has  
15% less error  
than Model C.

# Decide whether it's **worth it**

- we need a sampling distribution of PRE
  - a distribution of what PRE would look like if Model C (our  $H_0$ ) were true
- PRE is closely related to the *F* statistic!

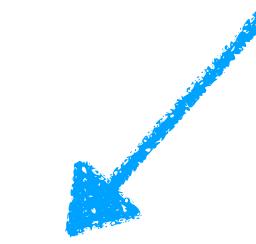
# Decide whether it's **worth it**

- To compute the  $F$  statistic, we need:
  - PRE
  - number of parameters in Model C (PC) and Model A (PA)
  - number of observations  $n$

more likely to be **worth it** if:

1. PRE is high
2. the number of additional parameters in A compared to C is low
3. the number of parameters that could have been added to  $\text{model}_C$  to create  $\text{model}_A$  but were not

**difference in parameters  
between models A and C**



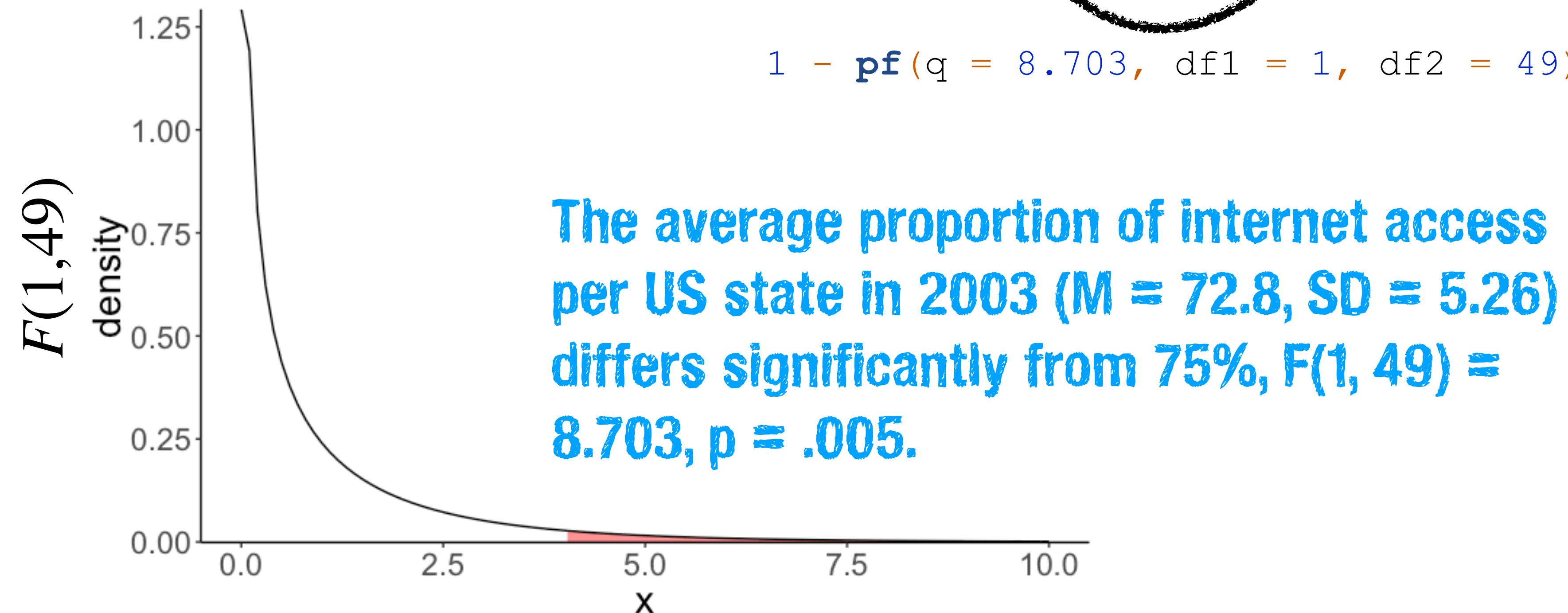
$$F = \frac{\text{PRE}/(\text{PA} - \text{PC})}{(1 - \text{PRE})/(n - \text{PA})}$$



**number of observations vs.  
parameters in Model A**

# Decide whether it's **worth it**

$$F = \frac{\text{PRE}/(\text{PA} - \text{PC})}{(1 - \text{PRE})/(n - \text{PA})}$$
$$= \frac{.15/(1 - 0)}{(1 - .15)/(50 - 1)} = 8.703 \quad p = .00486$$



**Note:** I've rounded PRE here (PRE = 0.15), but I'm reporting the F value based on the exact PRE value.

we just performed a one sample t-test ...

```
t.test(df.internet$internet, mu = 75)
```

## One Sample t-test

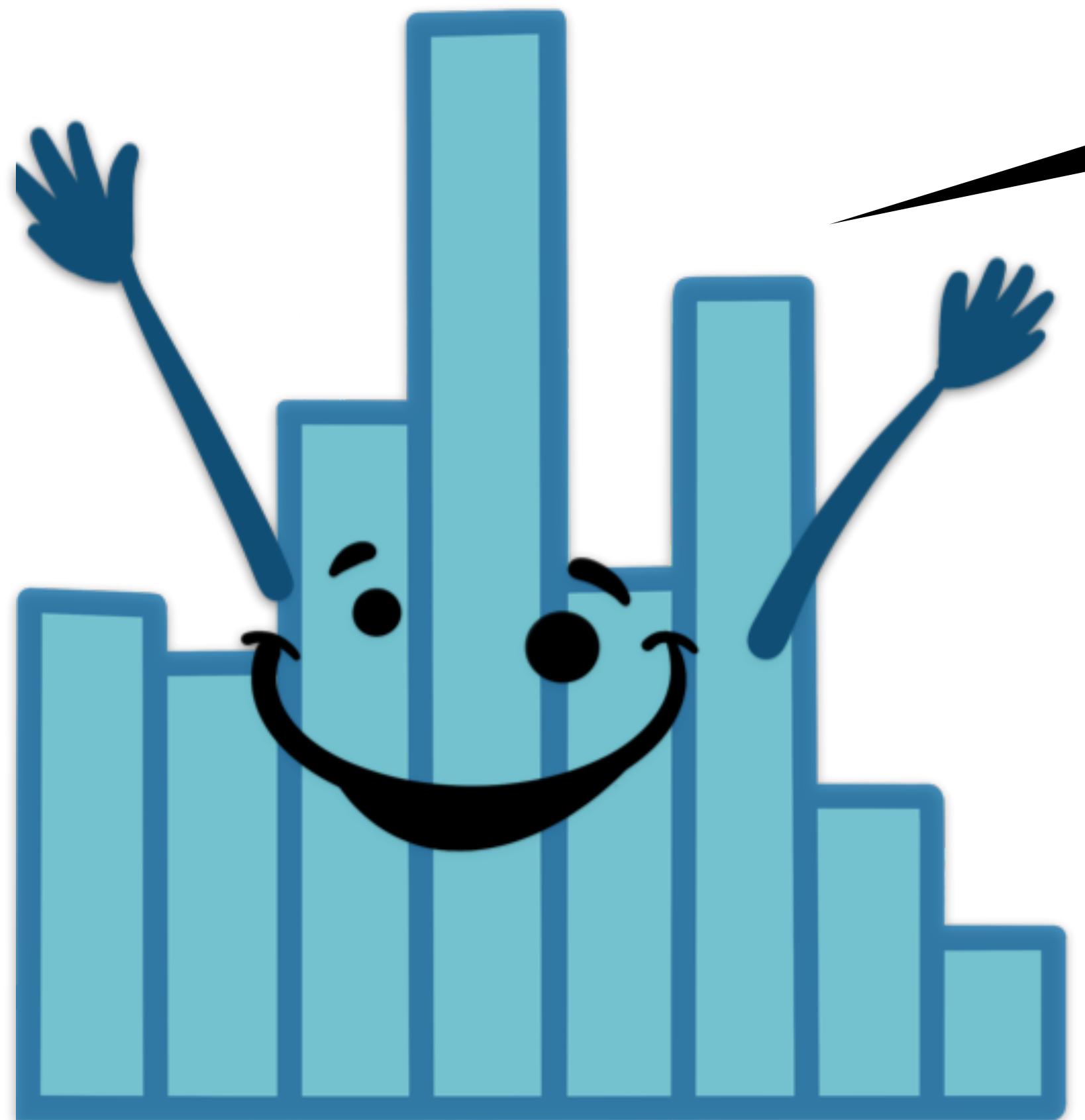
```
data: df.internet$internet  
t = -2.9502, df = 49, p-value = 0.00486  
alternative hypothesis: true mean is not equal to 75
```

but ...

- we will apply the general procedure we went through to day to a large range of different situations
- thinking of hypothesis testing as model comparison is (hopefully more) intuitive
- this approach still works even if we can't find a recipe in our favorite stats cook book

02:00

stretch break!

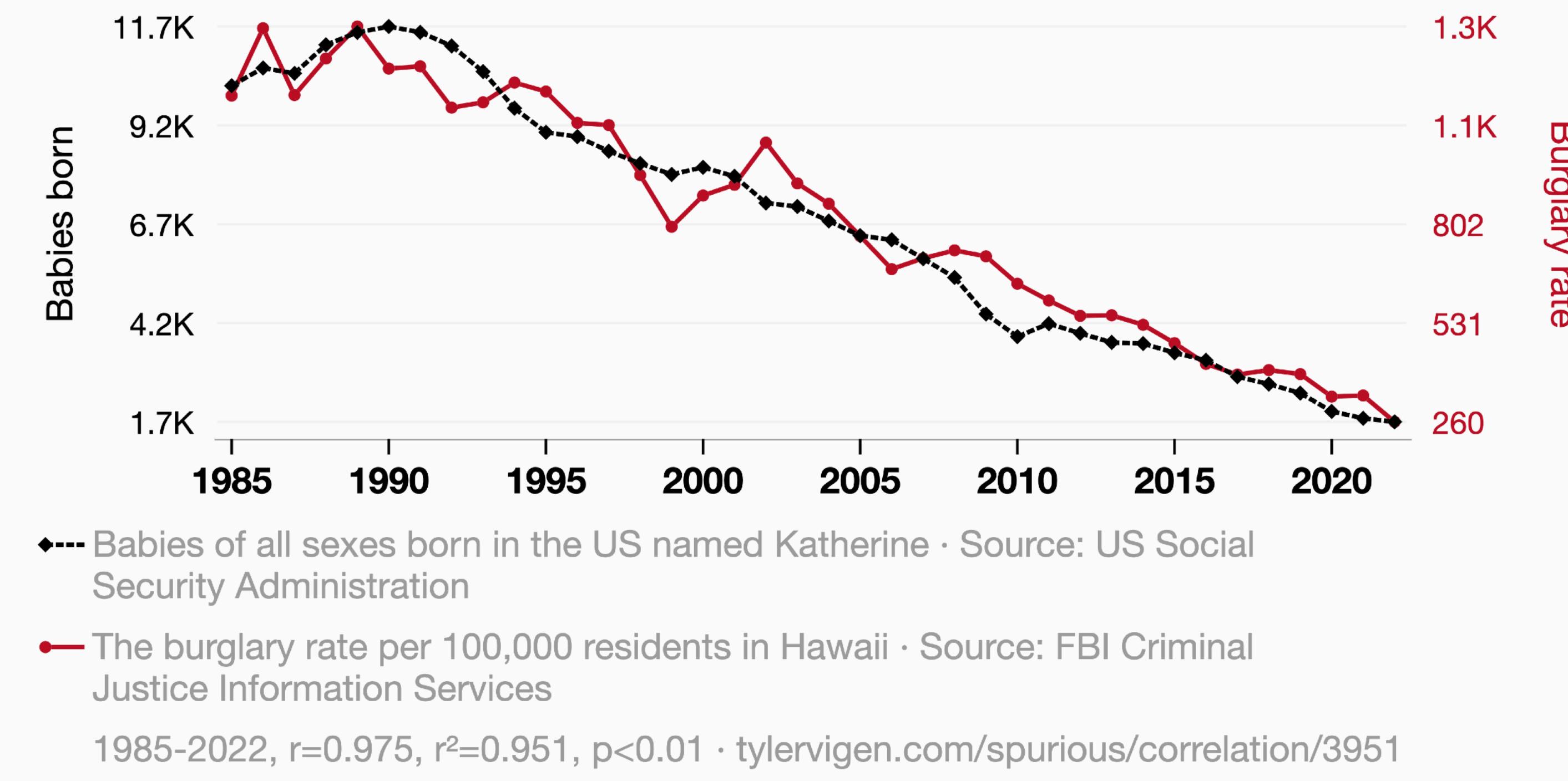


# **Correlation**

# Popularity of the first name Katherine

correlates with

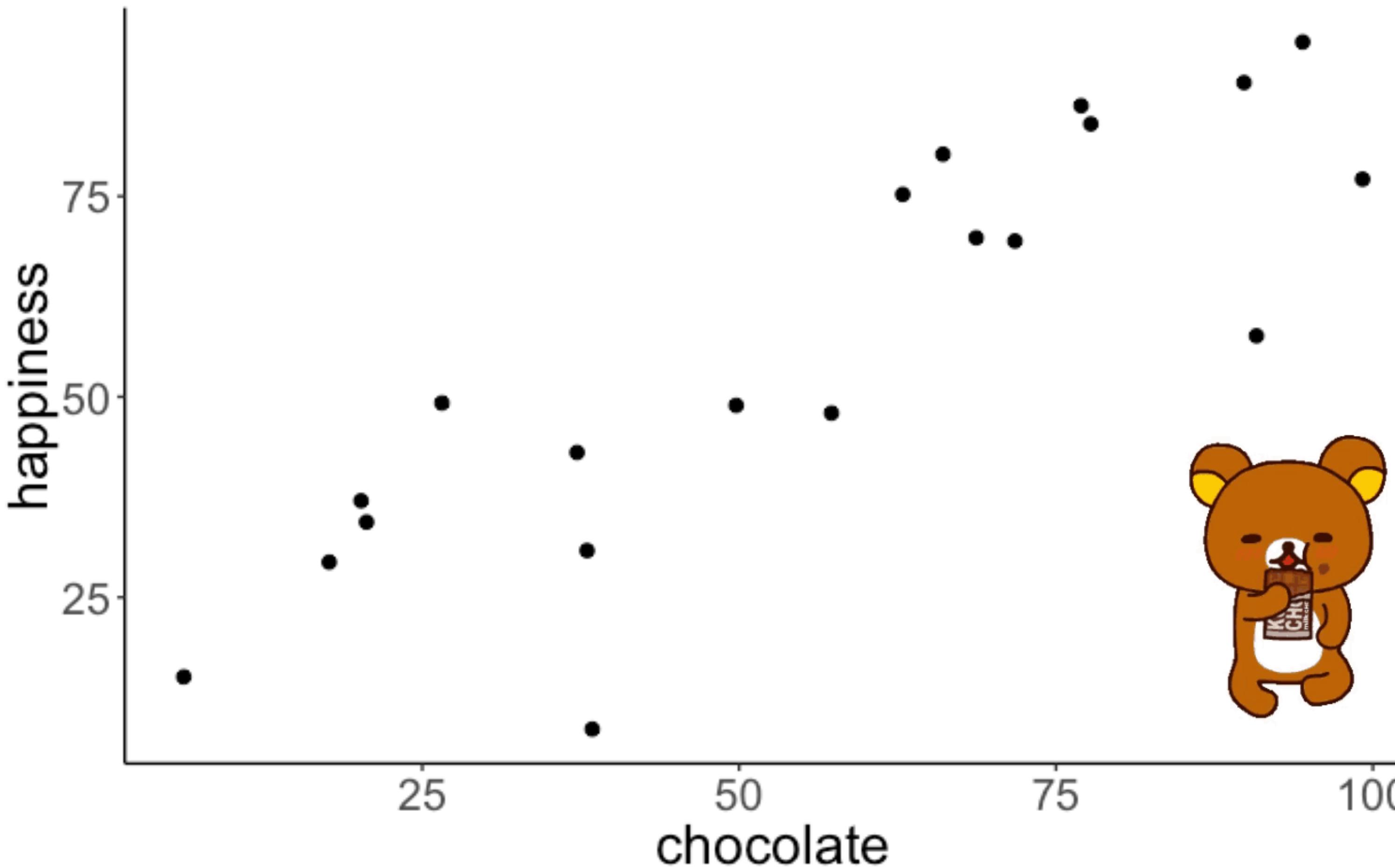
## Burglaries in Hawaii



AI generated

As the saying goes, "Kat's out of the bag," and it seems that also applies to burglars in Hawaii! With fewer Katherines around, there were less Kat burglars trying to pull off heists in the sunny state. It appears that the name Katherine was previously a common alias for cat burglars with a penchant for pilfering pineapples. However, with this name falling out of favor, it seems the purrpetrators have also disappeared, leading to a decrease in burglaries. It's a feline mystery, but it looks like Hawaii can rest easy knowing that the Katherine connection has been pawsitively purvented!

# How to best characterize the relationship between x and y by a single number?

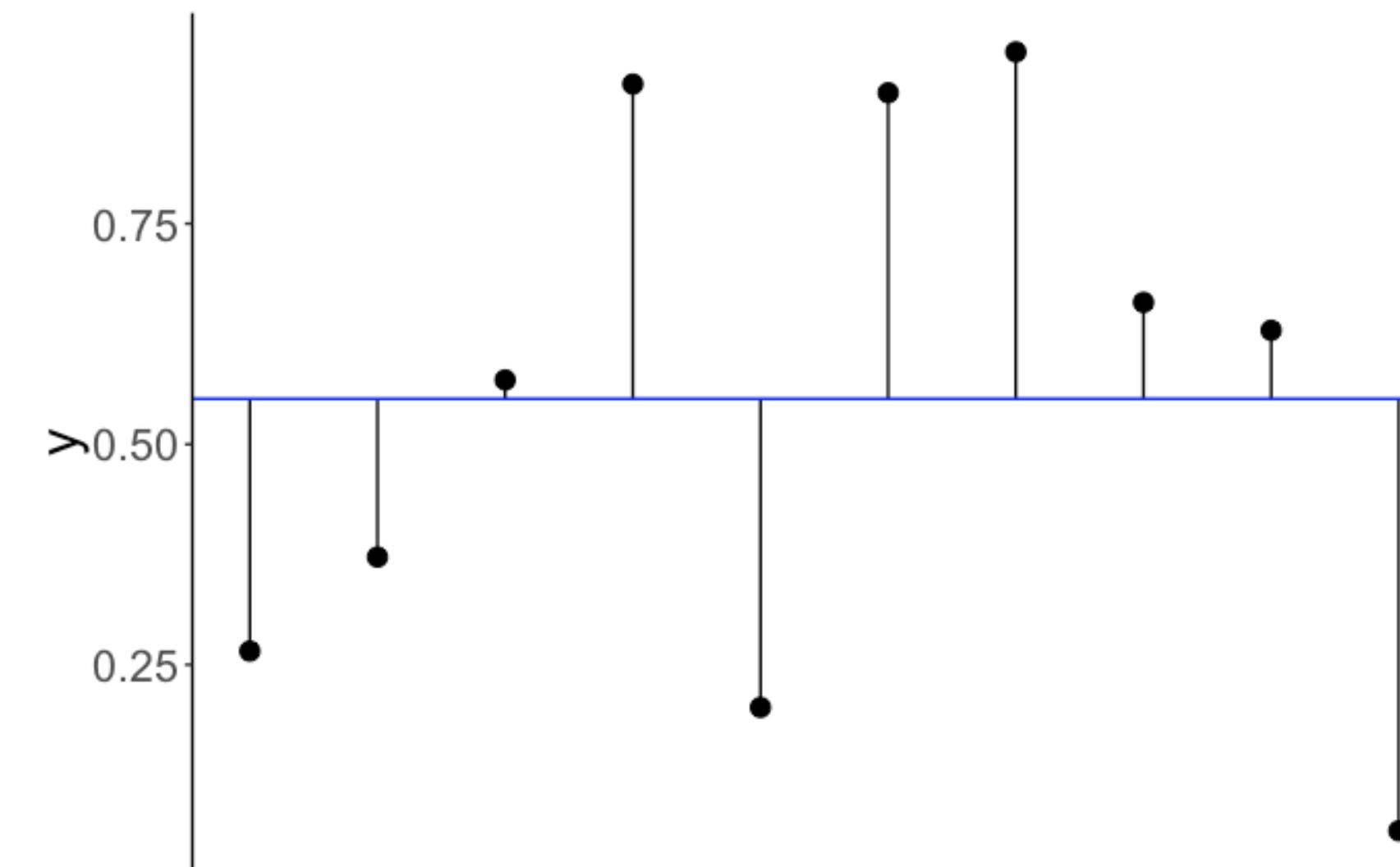
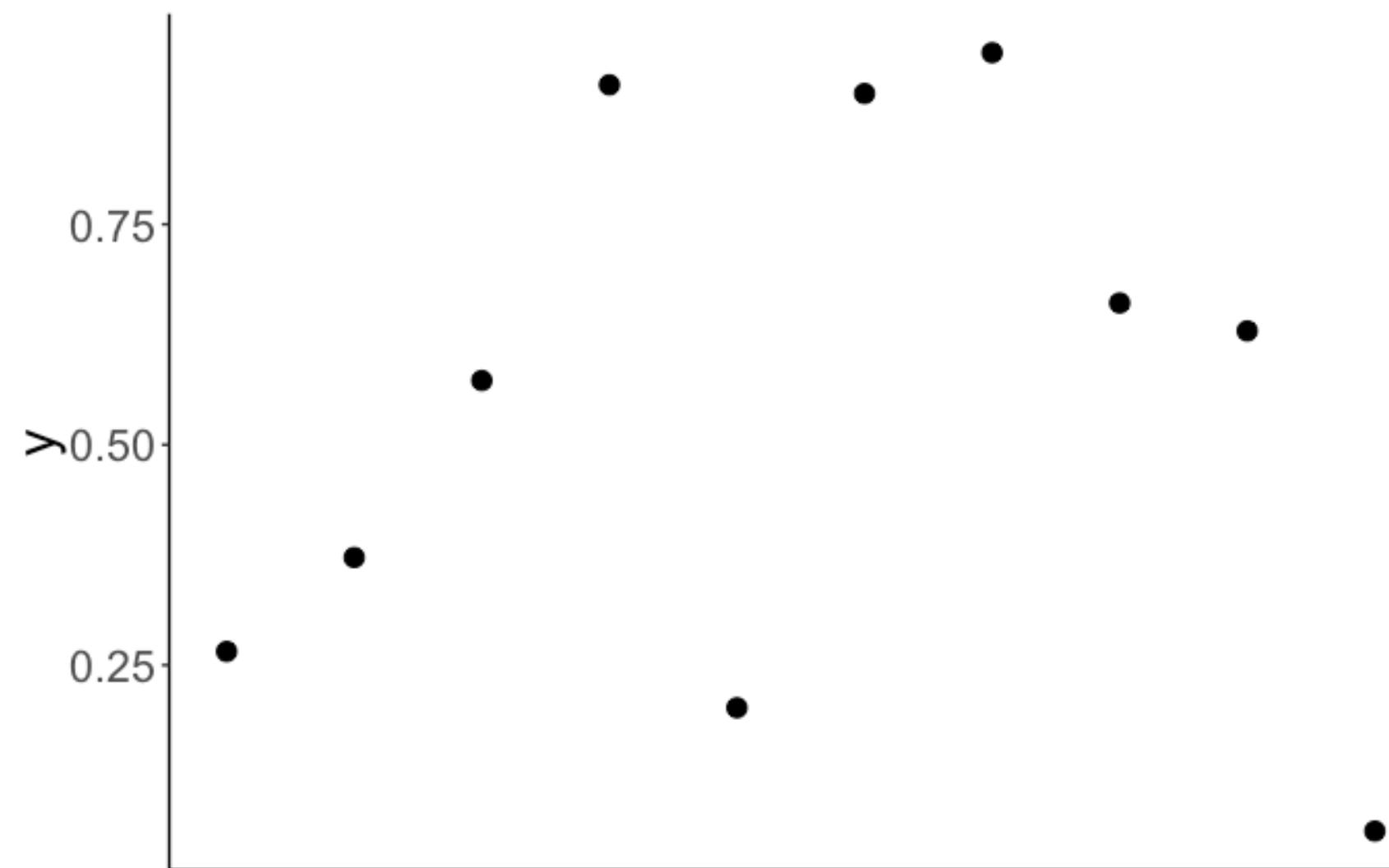


correlation = a measure of the relationship  
between two variables

## sample variance

$$Var(Y) = \sum_{i=1}^n \frac{(Y_i - \bar{Y})^2}{n - 1}$$

sum of squared errors

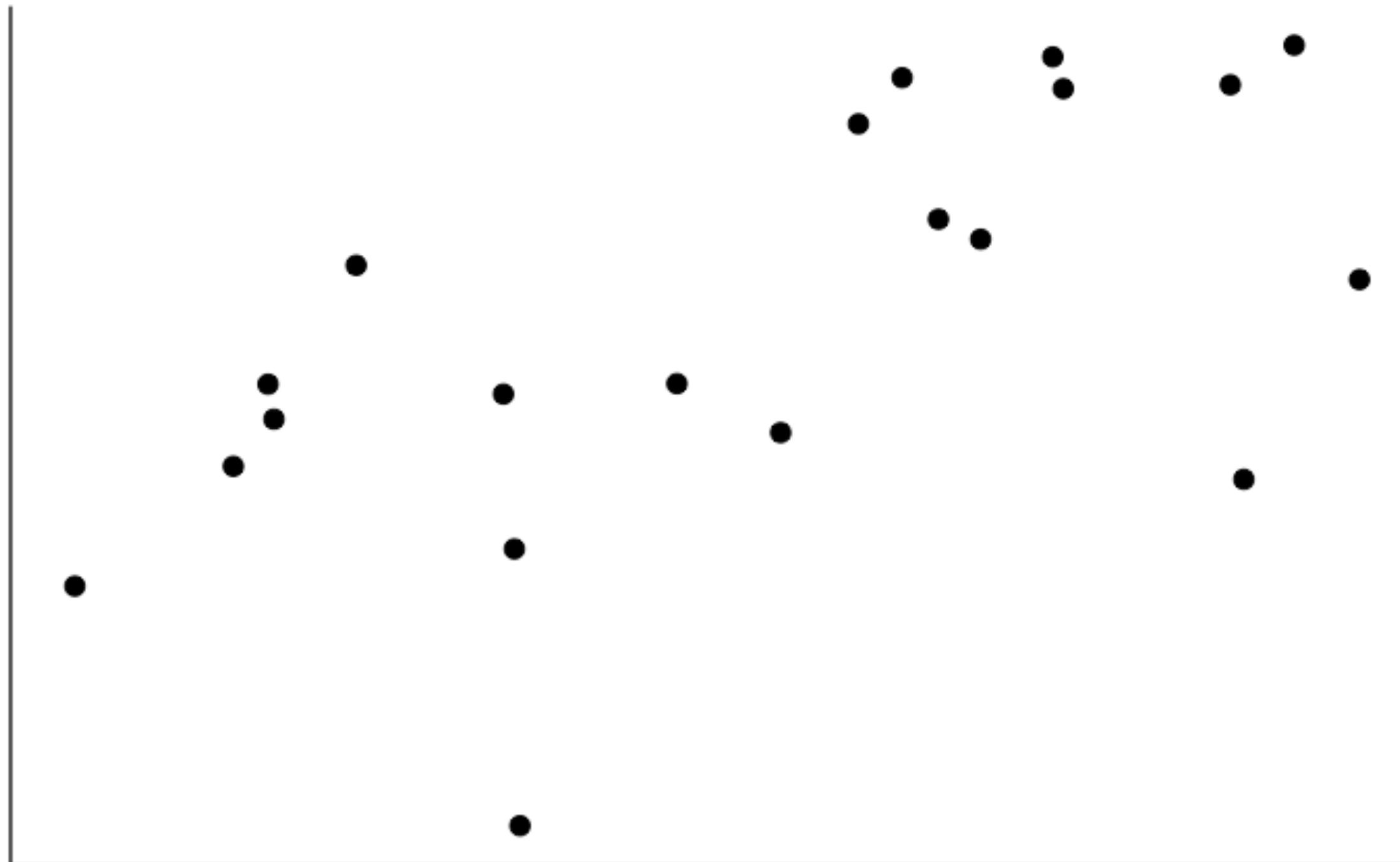


(I was too lazy to draw rectangles ...)

How well does the mean capture the data?

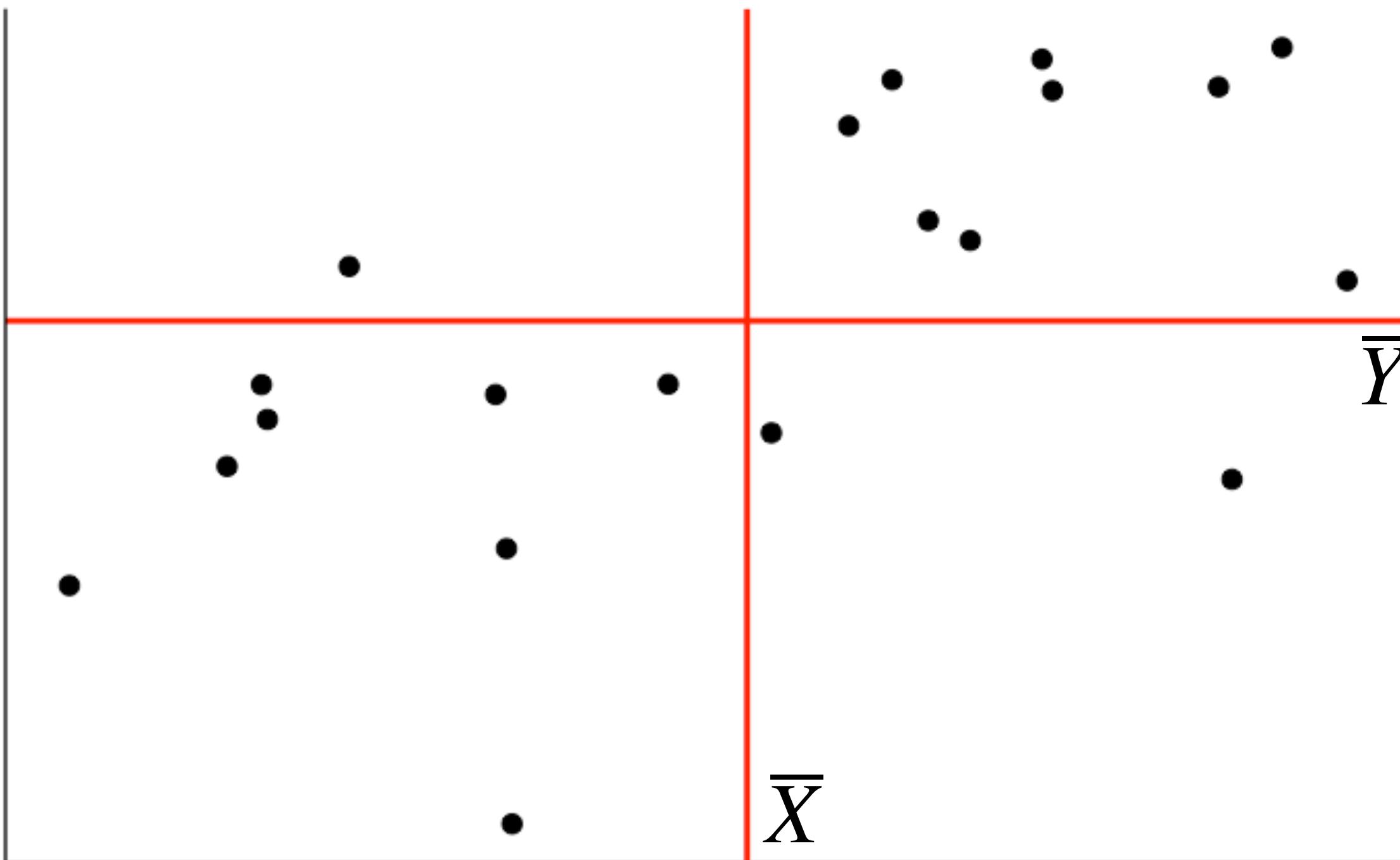
## sample covariance

$$Cov(X, Y) = \sum_{i=1}^n \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$



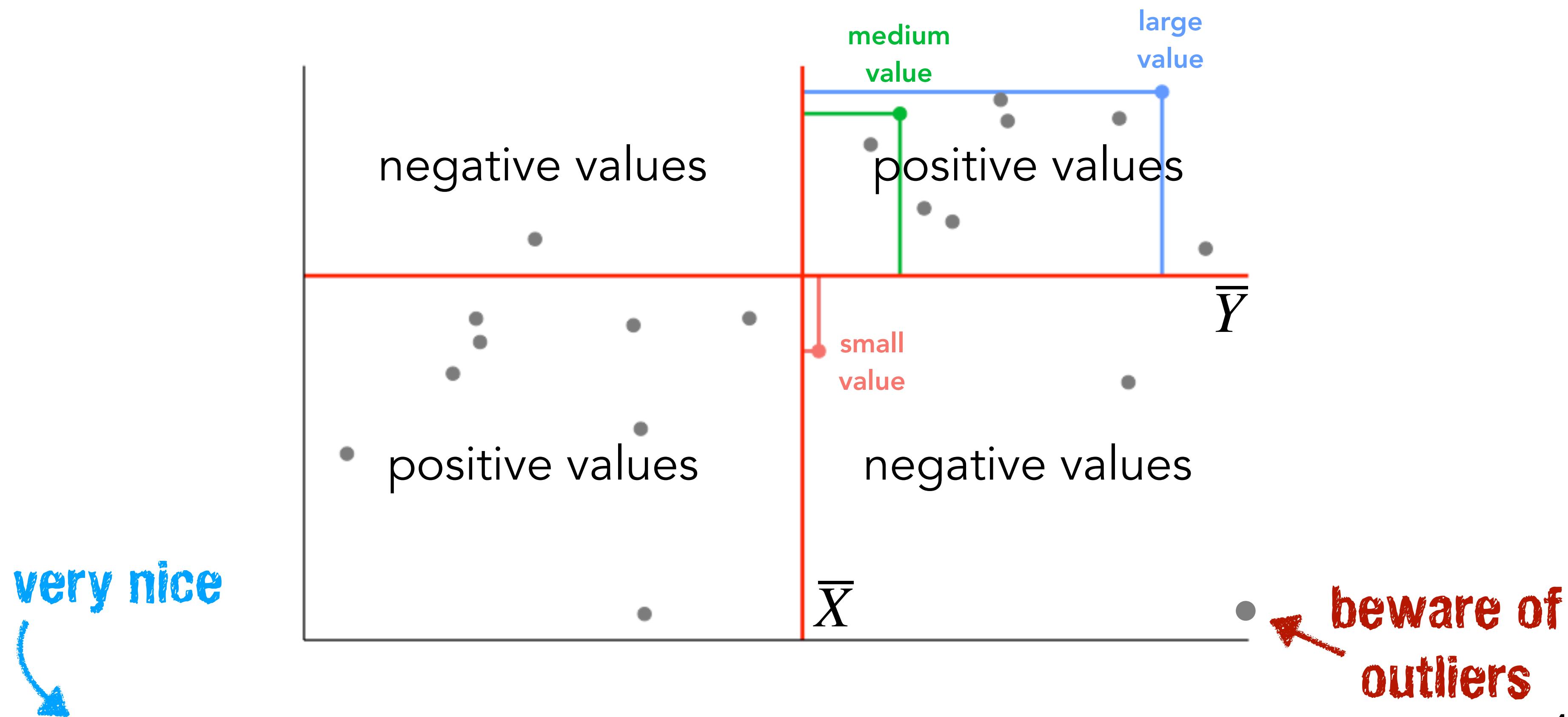
## sample covariance

$$Cov(X, Y) = \sum_{i=1}^n \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$



## sample covariance

$$Cov(X, Y) = \sum_{i=1}^n \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$



## sample covariance

$$Cov(X, Y) = \sum_{i=1}^n \frac{(X_i - \bar{X})(Y_i - \bar{Y})}{n - 1}$$

depends on the scale of the variables

## sample correlation coefficient

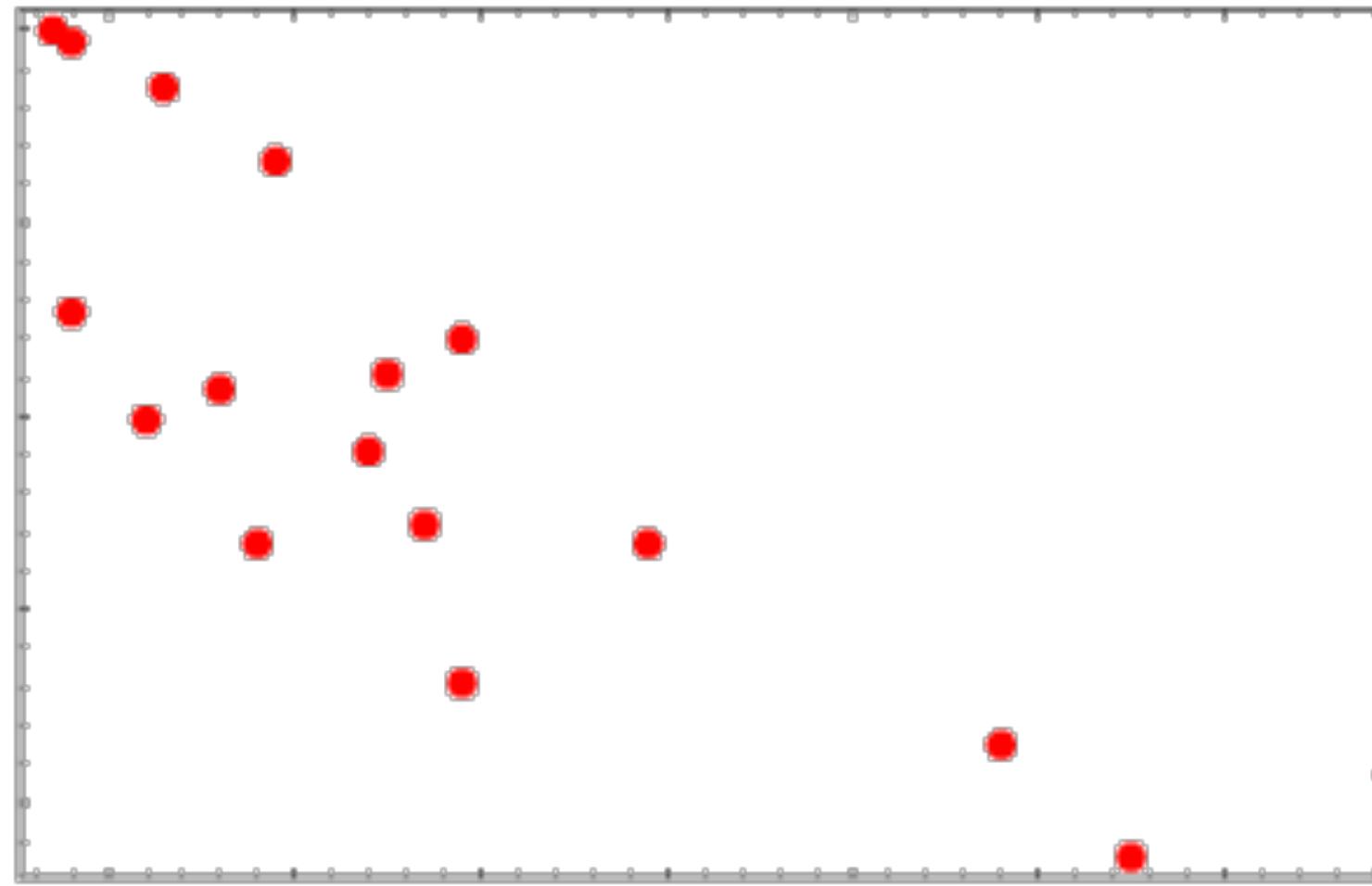
$$r(X, Y) = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2} \cdot \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}$$

standardized covariation  
(dividing by the standard deviations)

# Properties of the Pearson correlation

- standardized:  $-1 \leq r \leq 1$
- scale independent (for both X and Y)
- commutativity:  $r(X, Y) = r(Y, X)$    
**association not causation**
- sign determines the direction of dependence
- captures **linear dependence** only

# Who is the correlation champion?



Winner gets chocolate!

-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

-0.25 : 0

0 : 0.25

0.25 : 0.5

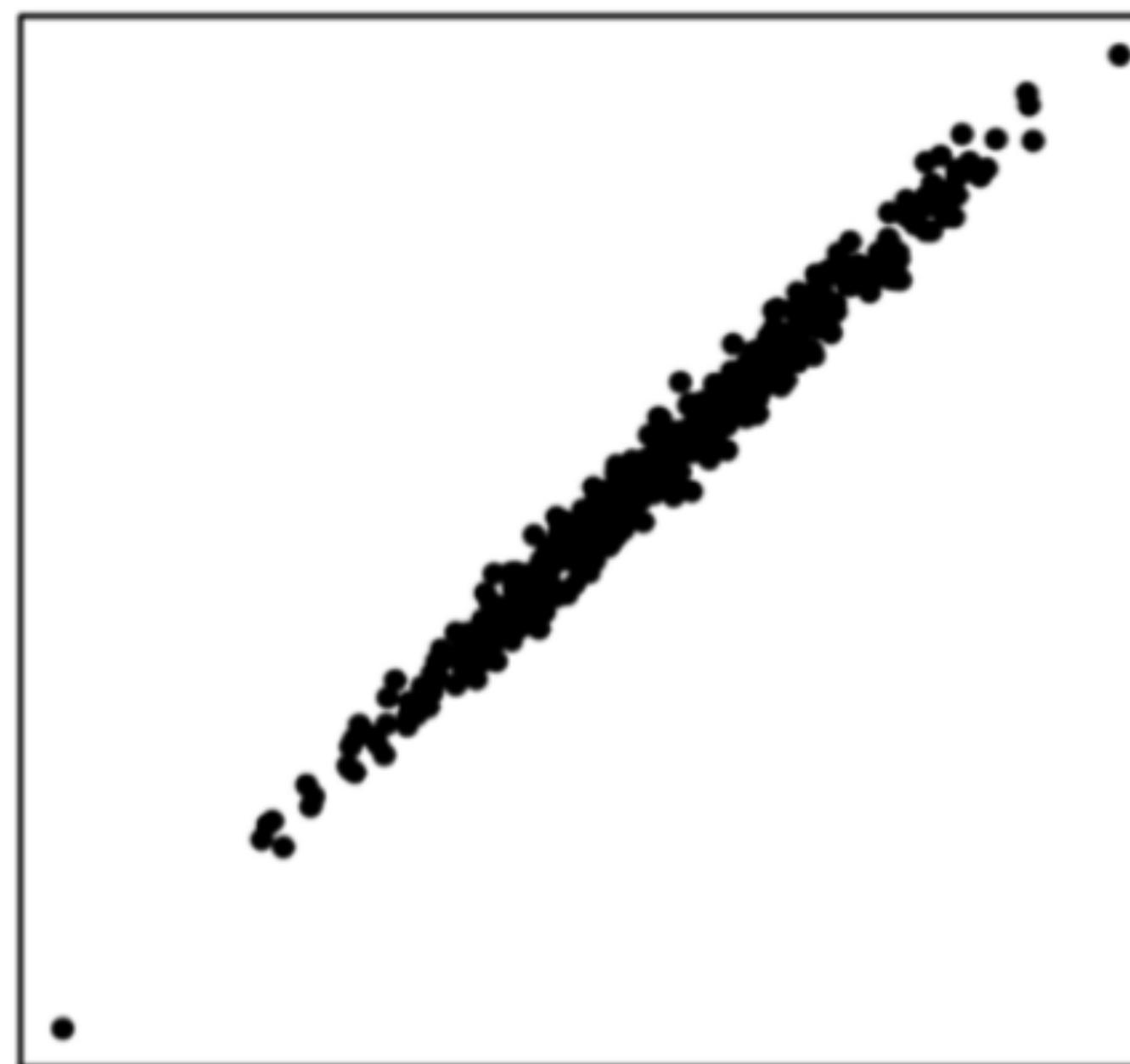
0.5 : 0.75

0.75 : 1

# Who is the correlation champion?

Win up to 1,000 points per answer

# In what range is the correlation coefficient?



-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

-0.25 : 0

0 : 0.25

0.25 : 0.5

0.5 : 0.75

0.75 : 1

-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

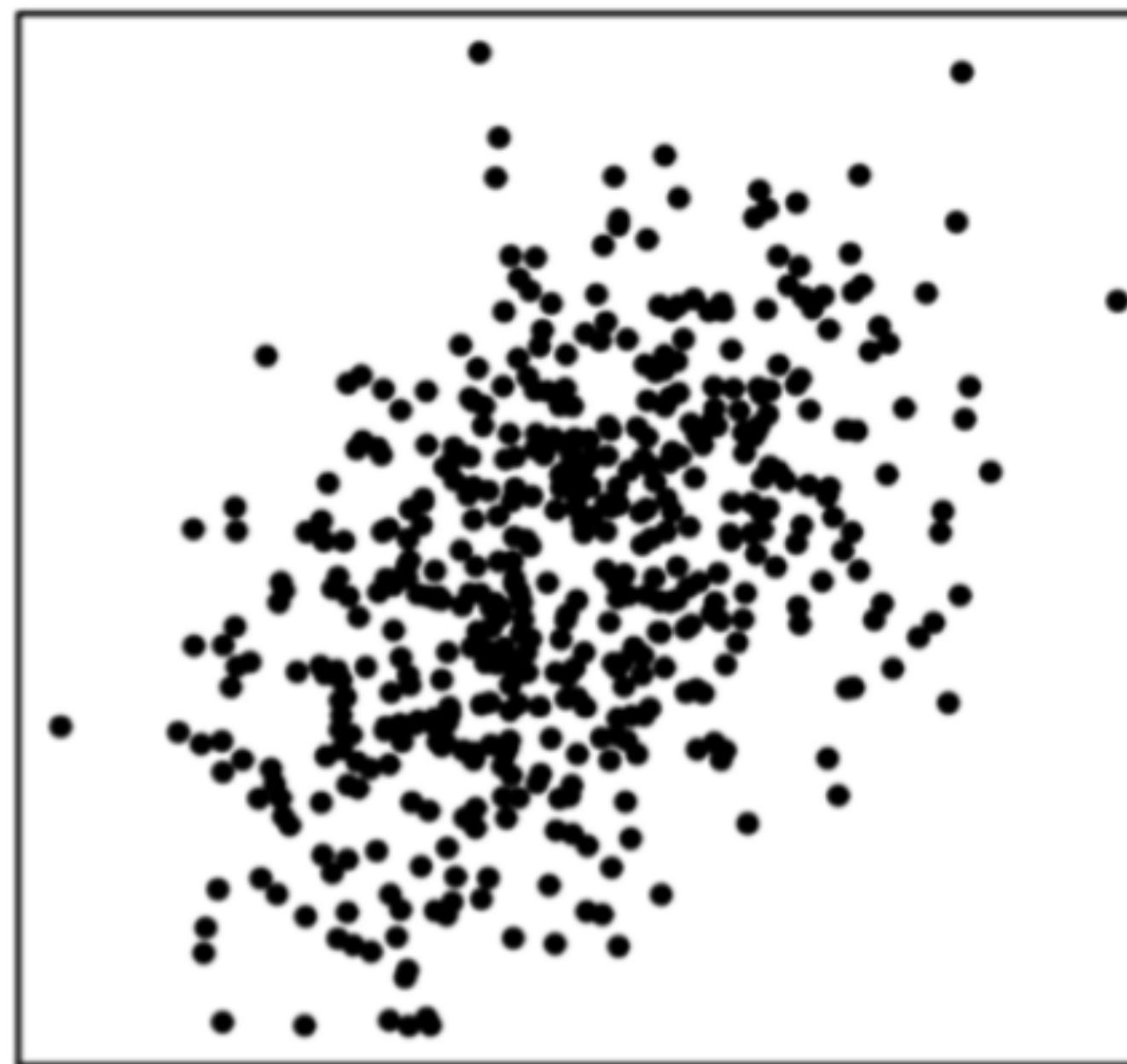
-0.25 : 0

0 : 0.25

0.25 : 0.5

0.5 : 0.75

0.75 : 1



-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

-0.25 : 0

0 : 0.25

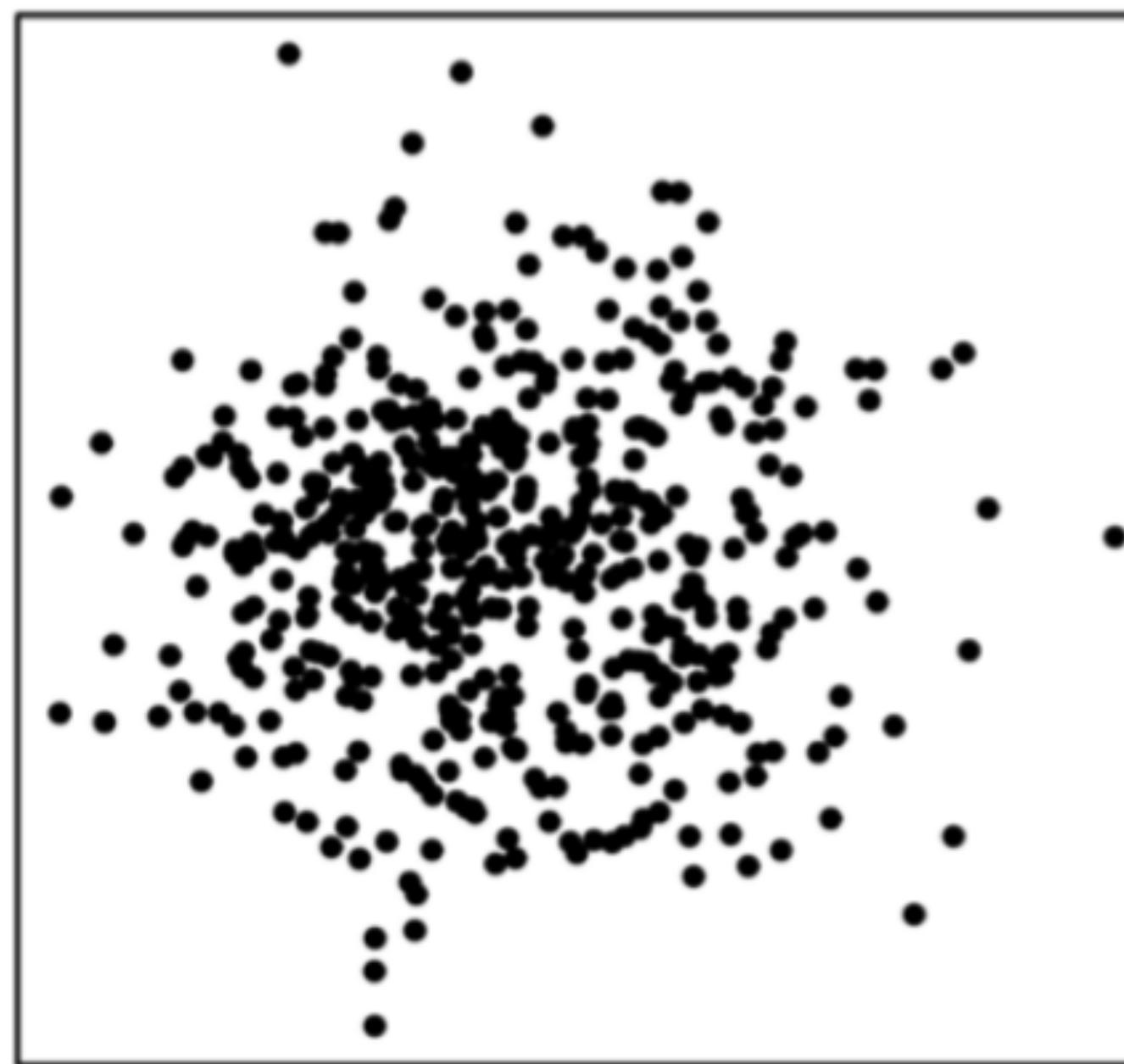
0.25 : 0.5

0.5 : 0.75

0.75 : 1

# Leaderboard

Nobody has responded yet.



-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

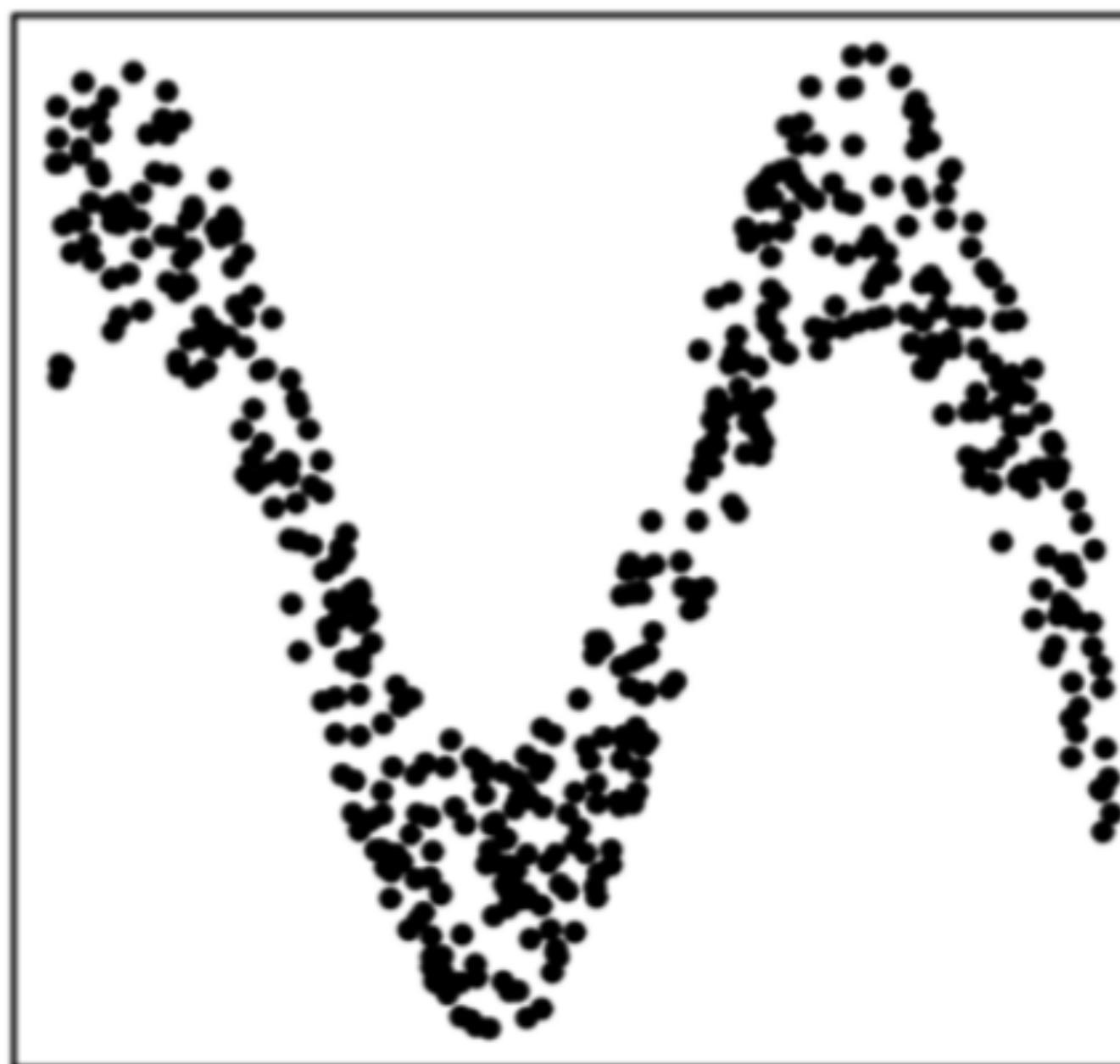
-0.25 : 0

0 : 0.25

0.25 : 0.5

0.5 : 0.75

0.75 : 1



-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

-0.25 : 0

0 : 0.25

0.25 : 0.5

0.5 : 0.75

0.75 : 1



-1 : -0.75

-0.75 : -0.5

-0.5 : -0.25

-0.25 : 0

0 : 0.25

0.25 : 0.5

0.5 : 0.75

0.75 : 1

# Leaderboard

Nobody has responded yet.

# Solution

XX

# Be careful about interpreting correlations!

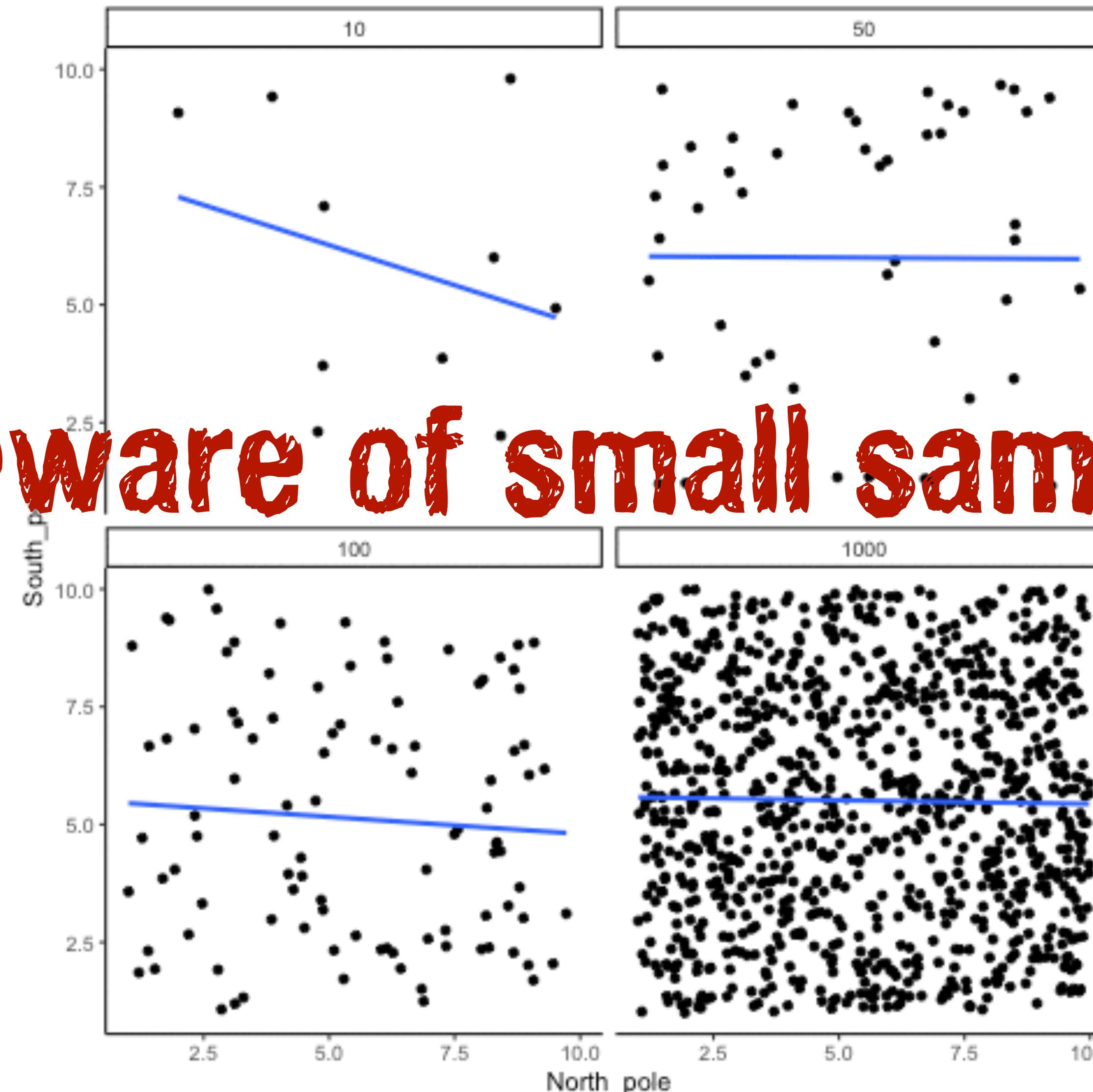


always visualize the data ...

$$n = [10, 50, 100, 1000]$$

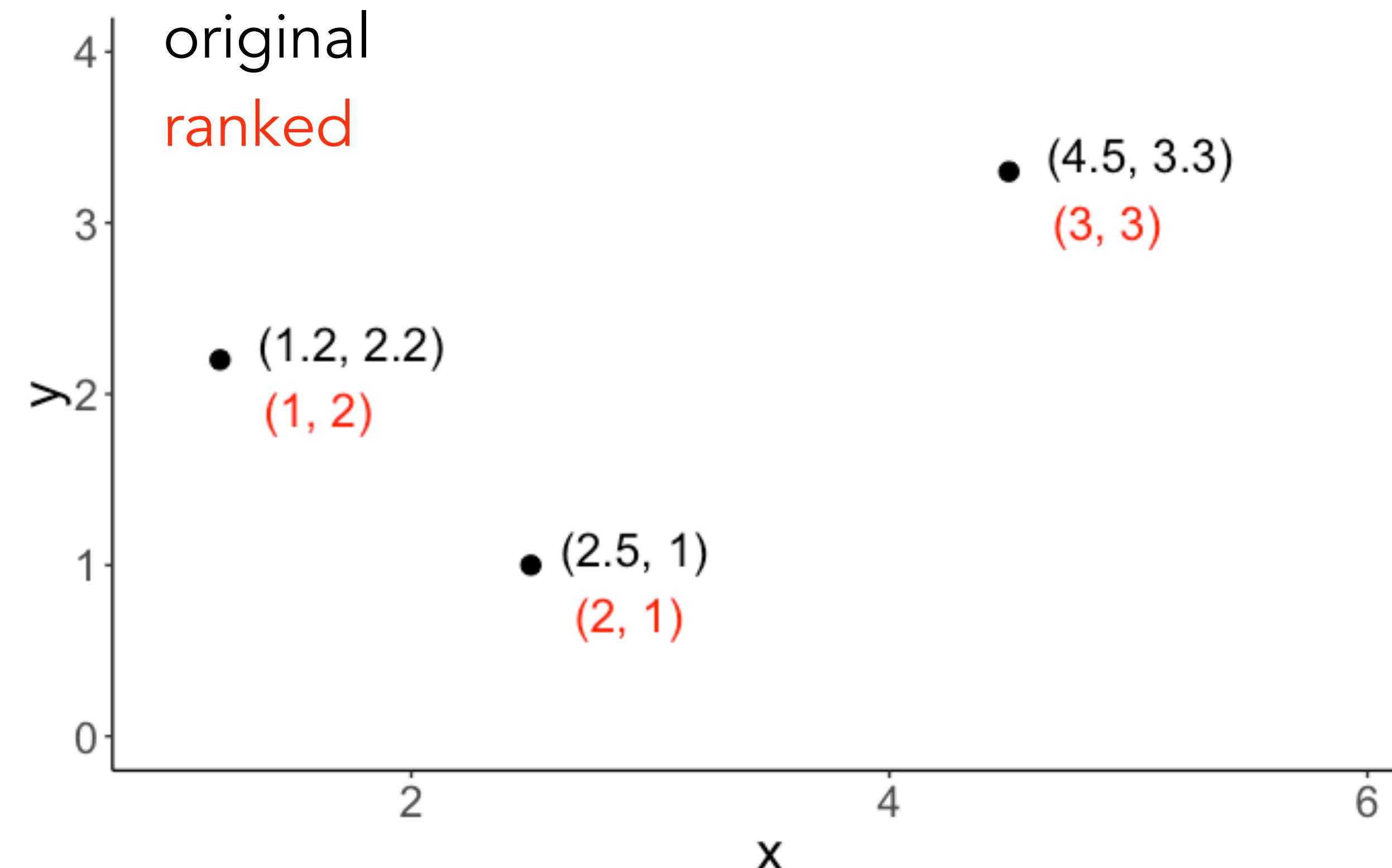
$$X \sim \mathcal{U}(\min = 0, \max = 10)$$

$$Y \sim \mathcal{U}(\min = 0, \max = 10)$$



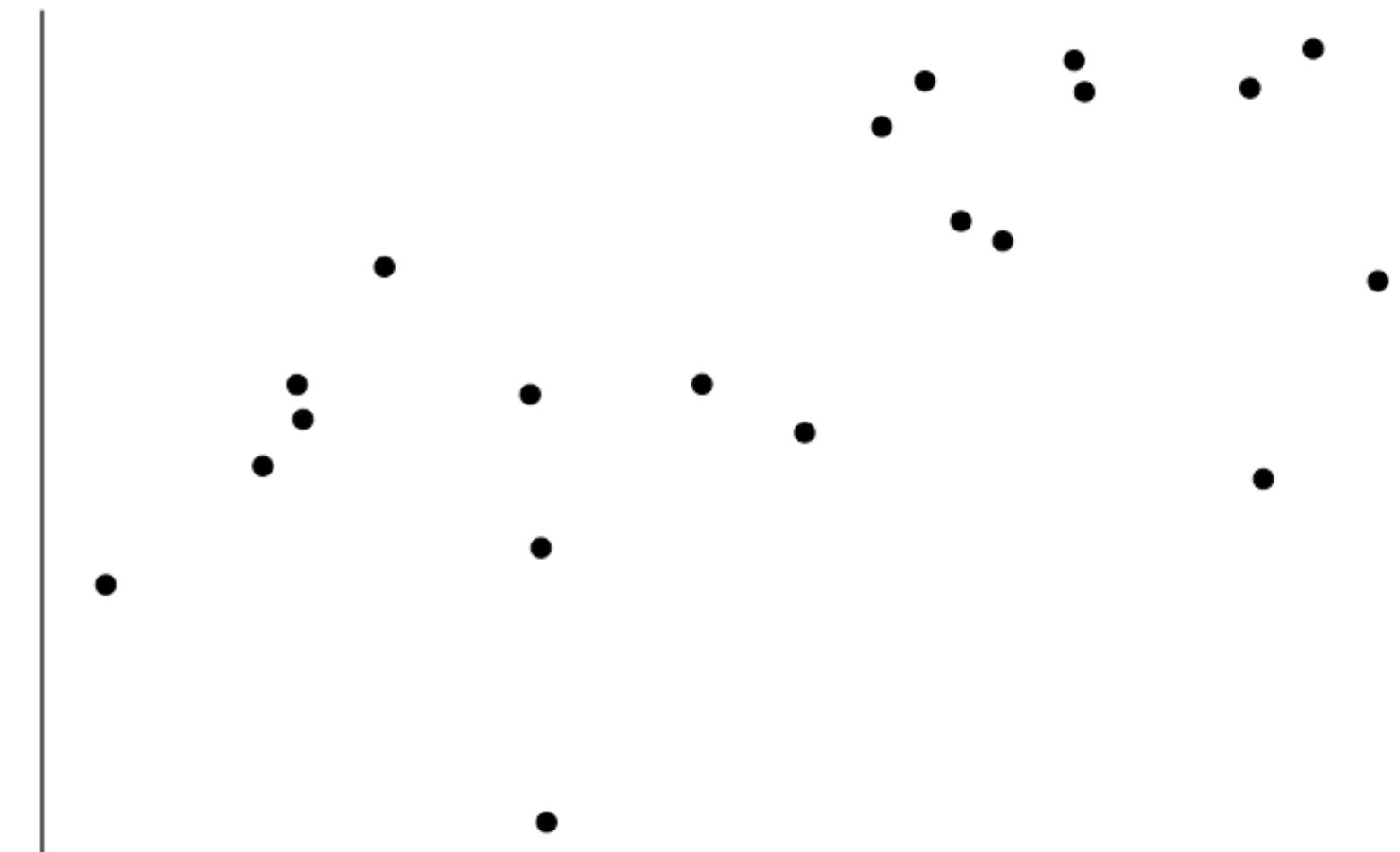
# Spearman rank order correlation

- transform original data into ranks
- calculate correlation on the ranked data



# Spearman rank order correlation

- transform original data into ranks
- calculate correlation on the ranked data



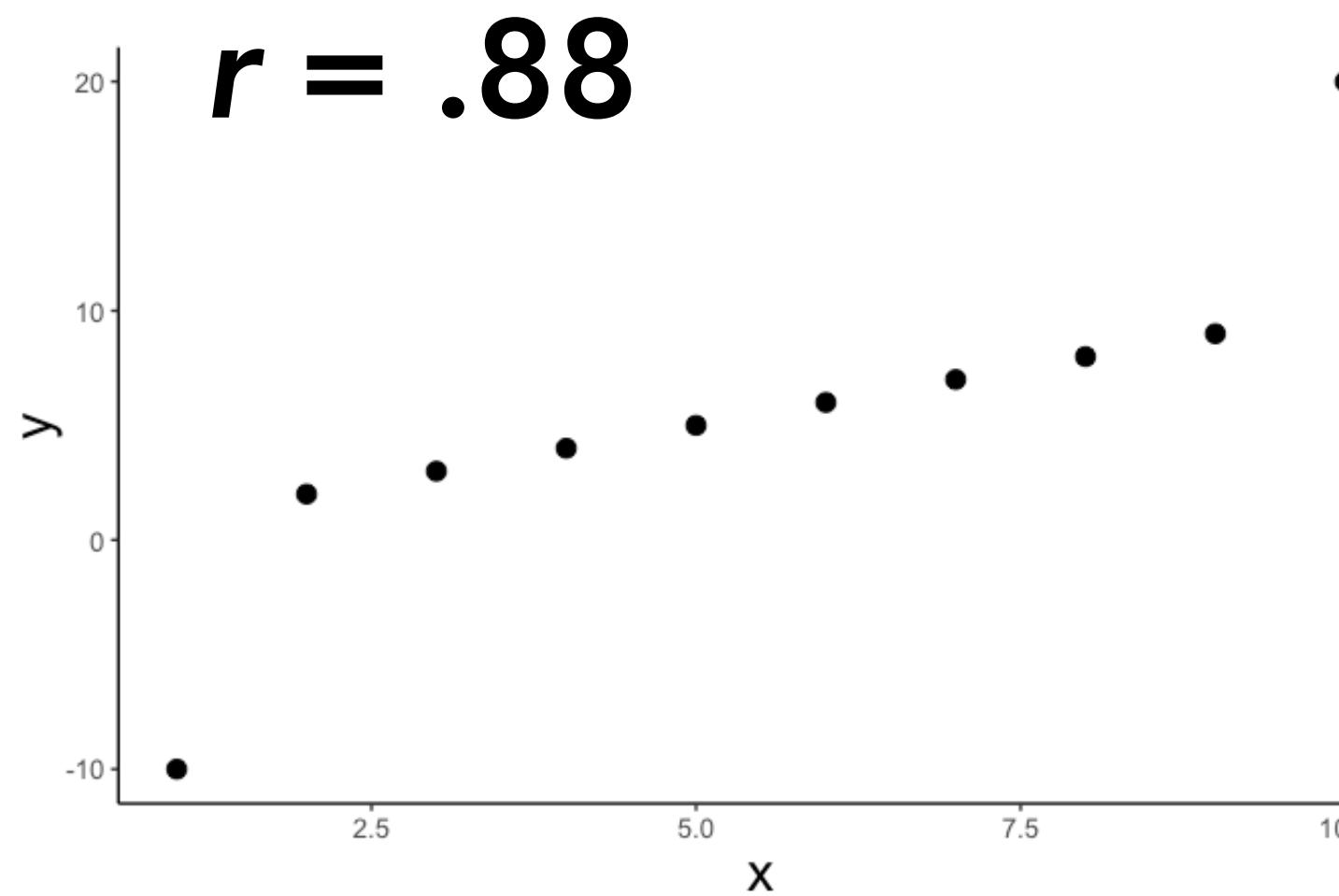
| x    | y    | x_rank | y_rank |
|------|------|--------|--------|
| 0.27 | 1.14 | 5      | 12     |
| 0.37 | 0.97 | 6      | 8      |
| 0.57 | 0.92 | 10     | 6      |
| 0.91 | 0.85 | 18     | 4      |
| 0.20 | 0.98 | 3      | 9      |
| 0.90 | 1.39 | 17     | 17     |
| 0.94 | 1.44 | 19     | 20     |
| 0.66 | 1.40 | 12     | 18     |
| 0.63 | 1.33 | 11     | 15     |
| 0.06 | 0.71 | 1      | 2      |

| r     | spearman | r_ranks |
|-------|----------|---------|
| 0.609 | 0.595    | 0.595   |

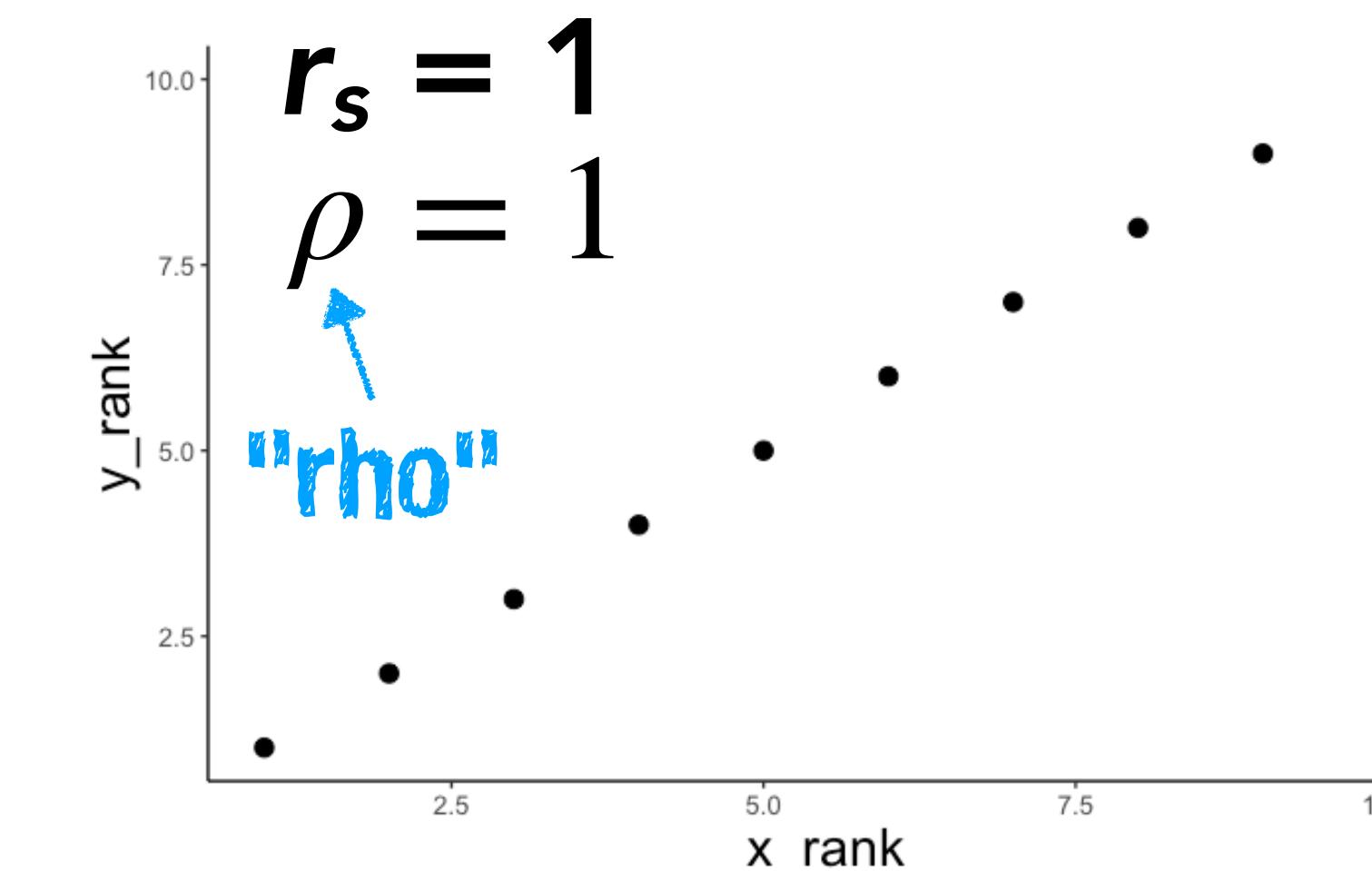
```
1 # correlation
2 df.spearman %>%
3   summarize(r = cor(x, y, method = "pearson"),
4             spearman = cor(x, y, method = "spearman"),
5             r_ranks = cor(x_rank, y_rank, method = "pearson"))
```

# Spearman rank order correlation

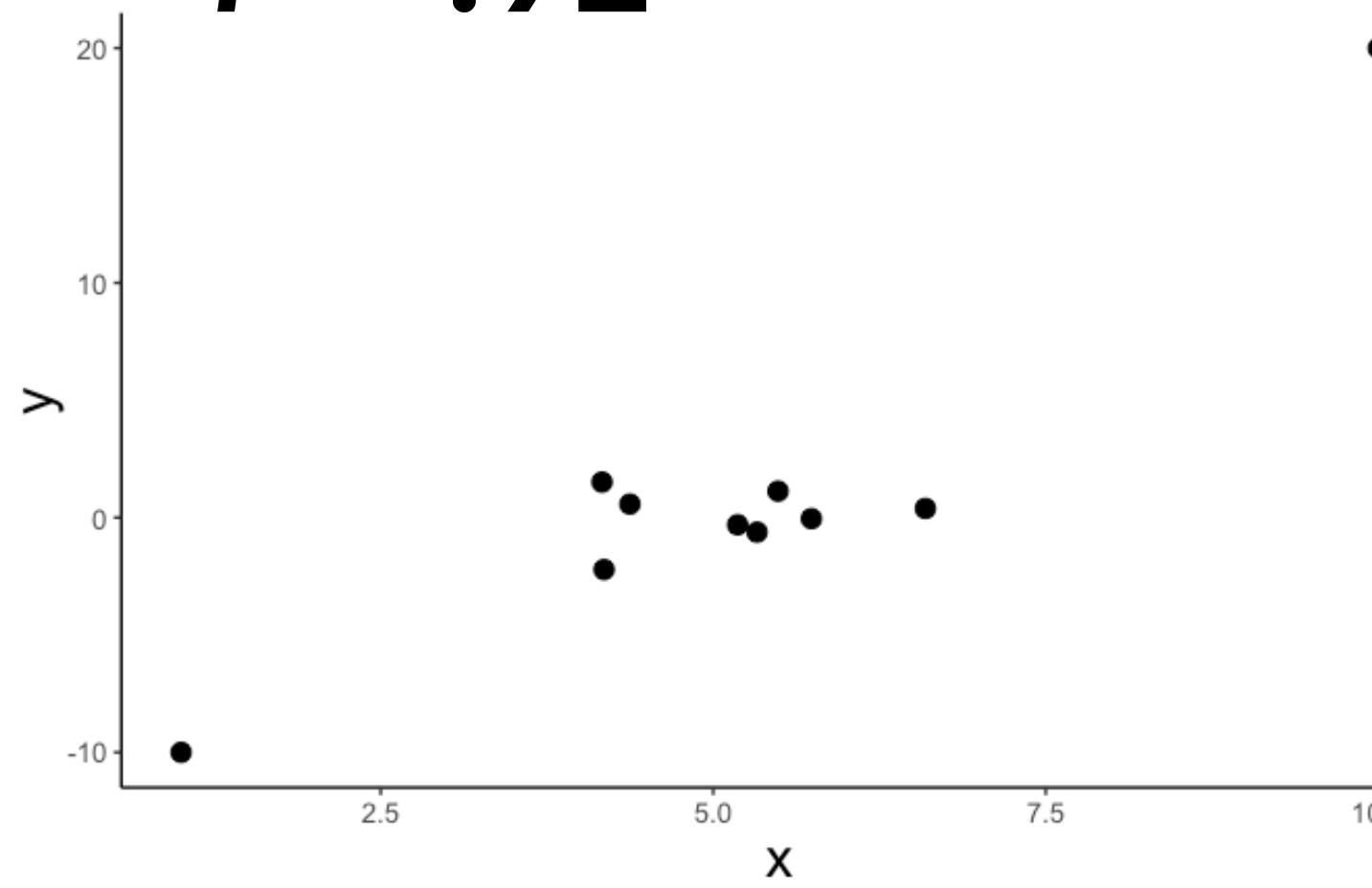
original



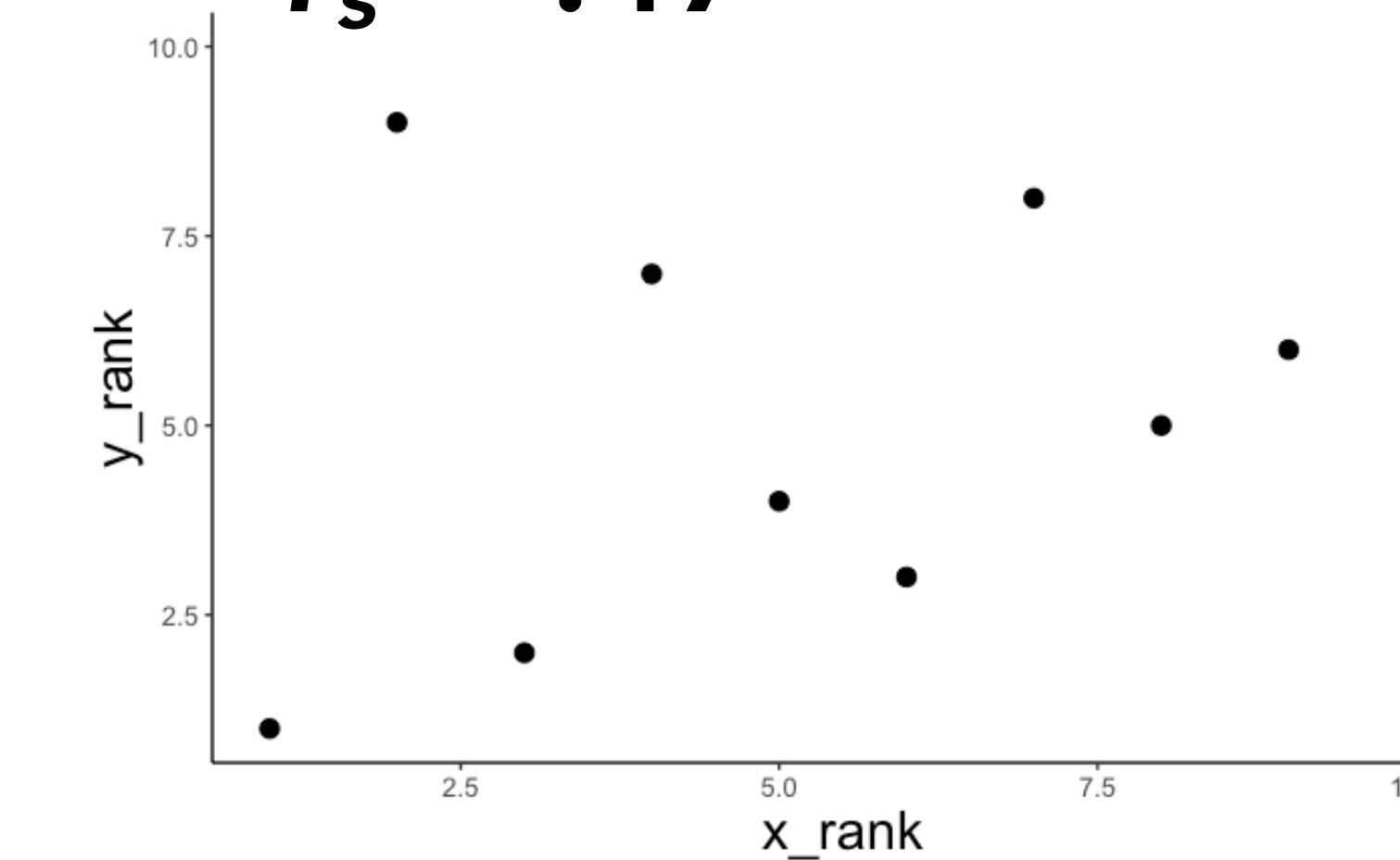
ranked



$r = .92$



$r_s = .47$



# Pearson vs. Spearman

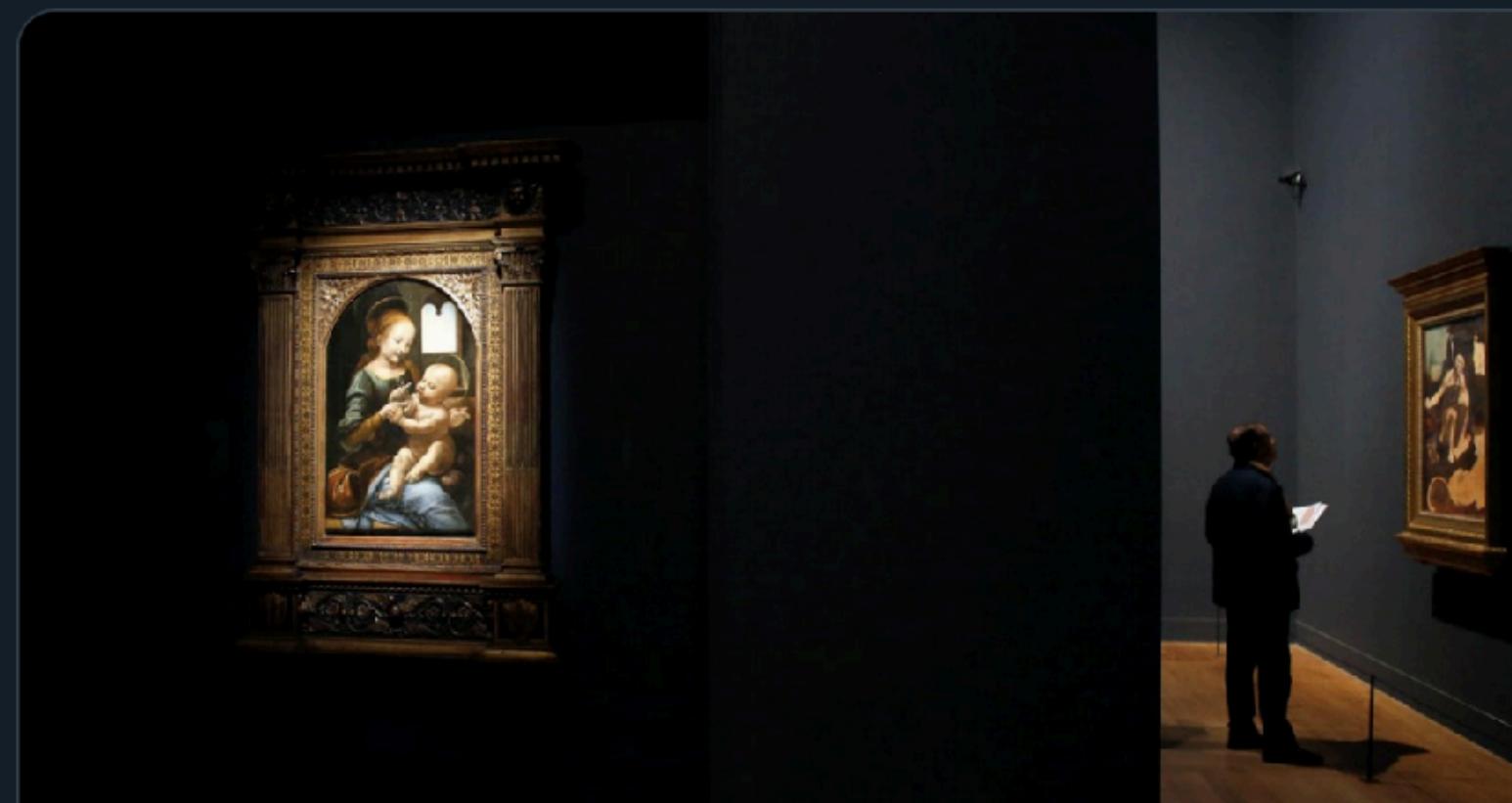
- Pearson's  $r$  captures the extent to which the relationship between two variable is **linear**
- Spearman's  $\rho$  captures the extent to which the relationship between two variables is **monotonic**
- What's better?
  - depends on the context
  - Spearman is robust to outliers, but it throws away (potentially useful) information

# CORRELATION IS NOT CAUSATION



NYT Health  
@NYTHealth

Want to live longer? Try going to the opera. Researchers in Britain have found that people who reported going to a museum or concert even once a year lived longer than those who didn't.



Another Benefit to Going to Museums? You May Live Longer

Researchers in Britain found that people who go to museums, the theater and the opera were less likely to die in the study period than those who didn't.

[nytimes.com](#)

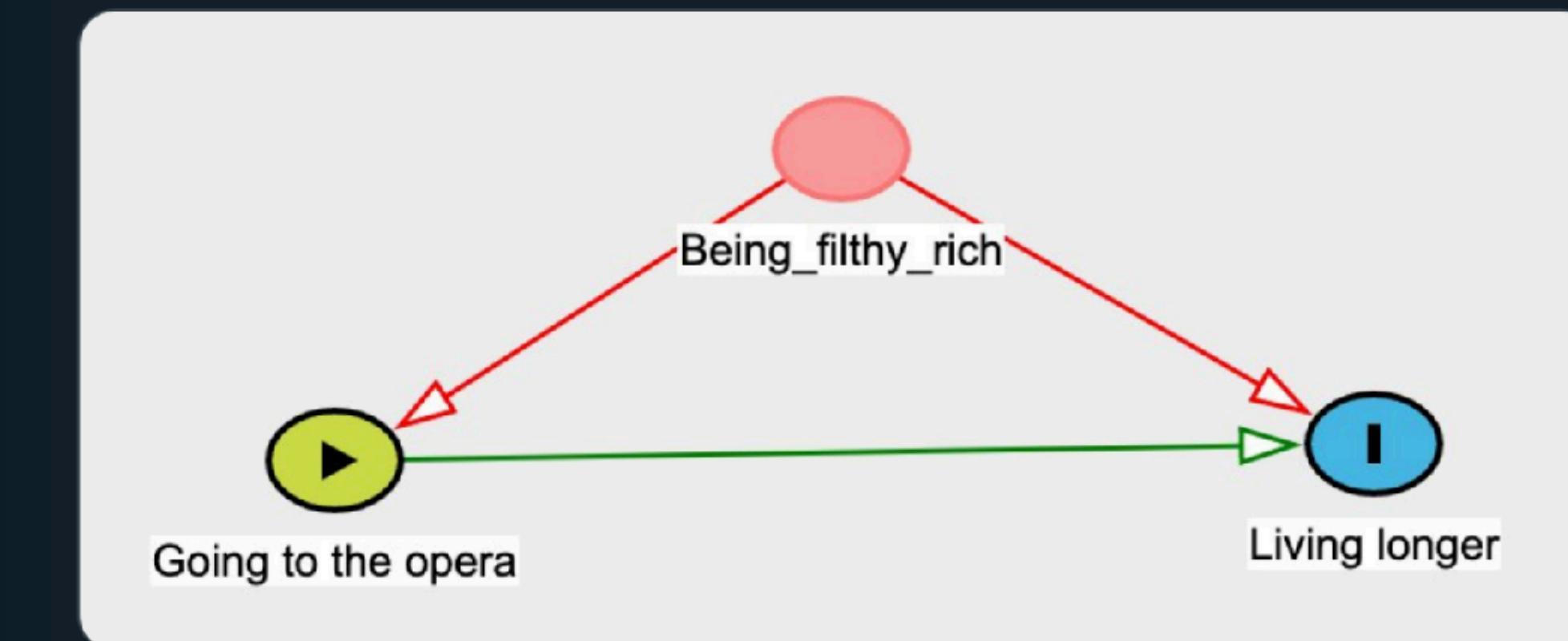
9:19 AM · Dec 22, 2019 · SocialFlow

336 Retweets 1.3K Likes



Andrew Heiss  
@andrewheiss

ooh ooh i can draw the dag for this one!



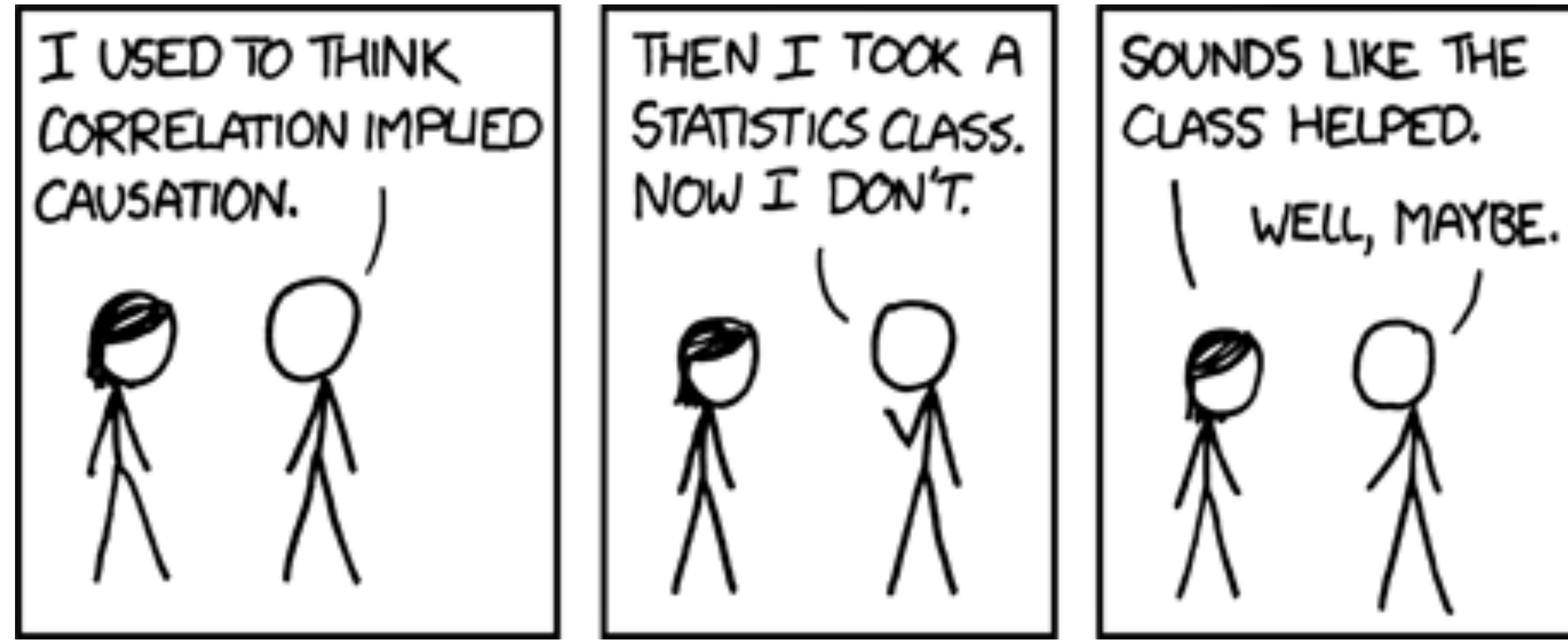
NYT Health @NYTHealth · Dec 22, 2019

Want to live longer? Try going to the opera. Researchers in Britain have found that people who reported going to a museum or concert even once a year lived longer than those who didn't. [nyti.ms/2Q9AmZV](#)

2:47 PM · Dec 22, 2019 · Twitter Web App

[View Tweet activity](#)

837 Retweets 3.9K Likes



- correlations suggest that there is some causal relationship
- but this relationship need not be a direct causal relationship from A to B (or from B to A)

more about causation in a later class

# Regression

# **The conceptual tour**

# Linear model: Simple regression

Data = Model + Error

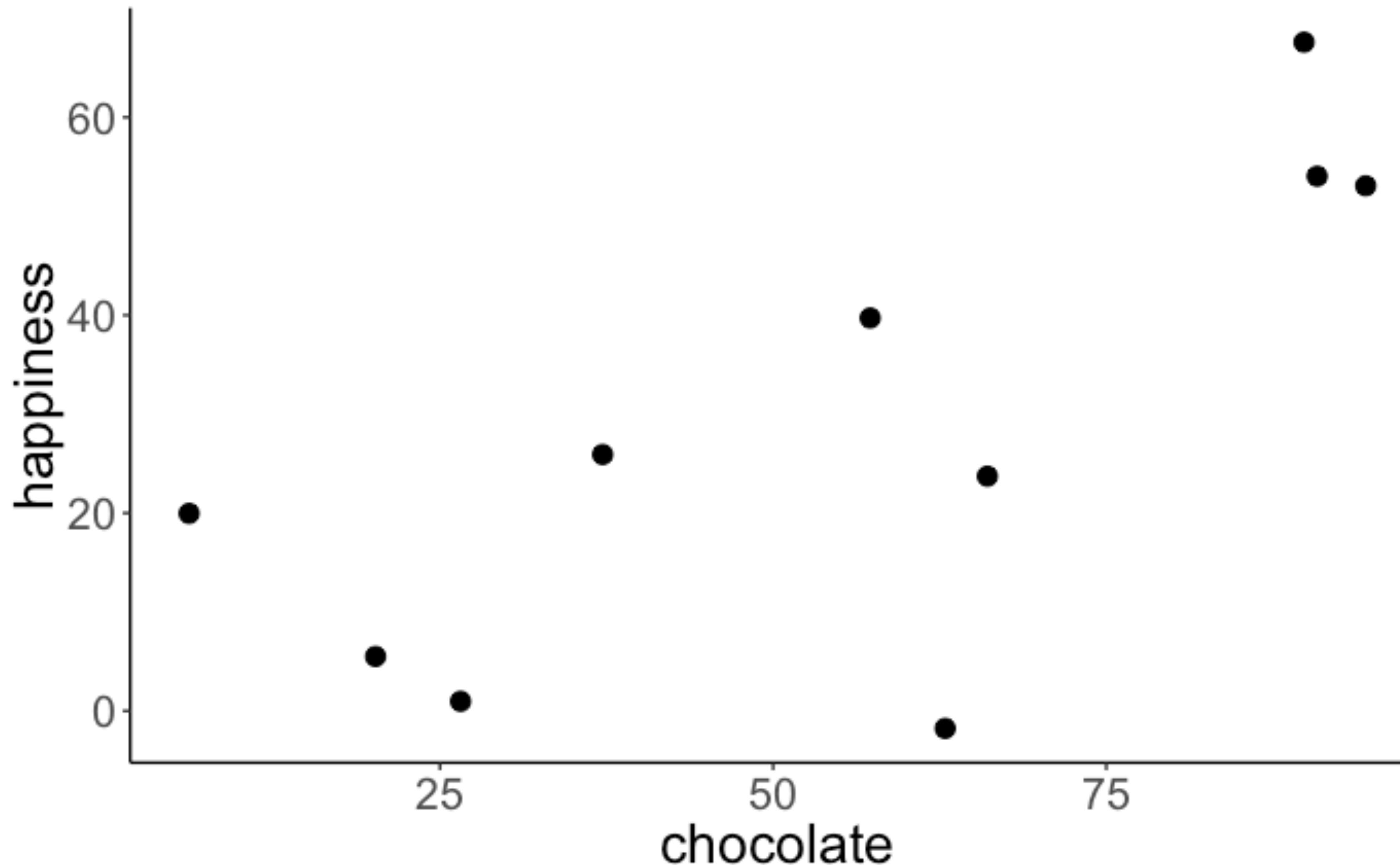
$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

$$\hat{Y}_i = \beta_0 + \beta_1 X_i$$



the model is a linear  
combination of predictors

# Is there a relationship between chocolate consumption and happiness?



# The general procedure

1. Define  $H_0$  as Model C (compact) and  $H_1$  as Model A (augmented)
2. Fit model parameters to the data
3. Calculate the proportional reduction of error (PRE) in our sample
4. Decide whether the augmented model is **worth it** by comparing the observed PRE in our sample to the sampling distribution of PRE (assuming that  $H_0$  is true)

# The general procedure

1. Define  $H_0$  as Model C (compact) and  $H_1$  as Model A (augmented)
2. Fit model parameters to the data
3. Calculate the proportional reduction of error (PRE) in our sample
4. Decide whether the augmented model is **worth it** by comparing the observed PRE in our sample to the sampling distribution of PRE (assuming that  $H_0$  is true)

$H_0$ : Chocolate consumption and happiness are unrelated.

### Model C

$$Y_i = \beta_0 + \epsilon_i$$

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

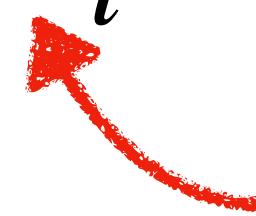
and

$$\beta_1 = 0$$

$H_1$ : Chocolate consumption and happiness are related.

### Model A

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$



chocolate  
consumption

# The general procedure

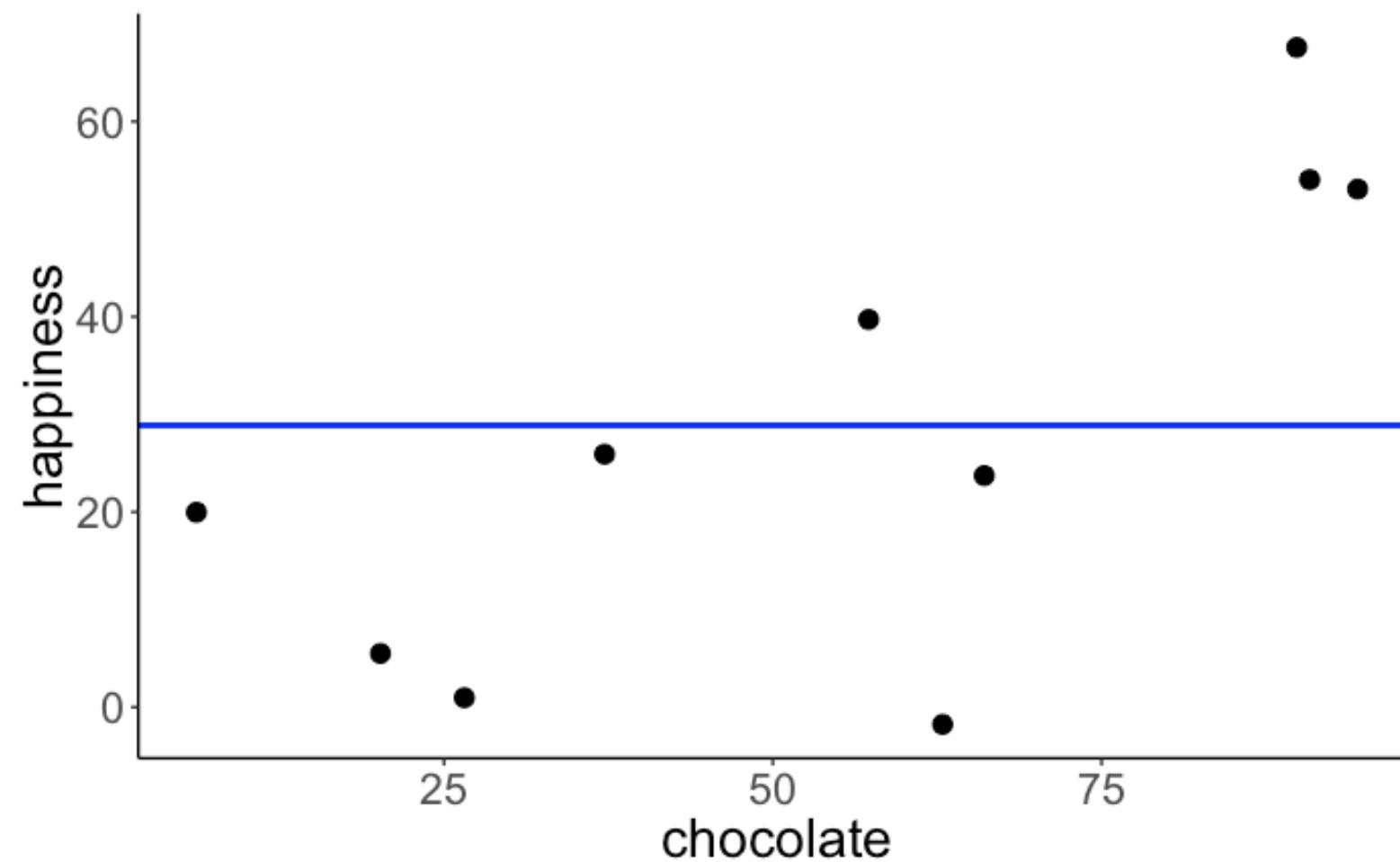
1. Define  $H_0$  as Model C (compact) and  $H_1$  as Model A (augmented)
- 2. Fit model parameters to the data**
3. Calculate the proportional reduction of error (PRE) in our sample
4. Decide whether the augmented model is **worth it** by comparing the observed PRE in our sample to the sampling distribution of PRE (assuming that  $H_0$  is true)

$H_0$ : Chocolate consumption and happiness are unrelated.

### Model C

$$Y_i = \beta_0 + \epsilon_i$$

### Model prediction



### Fitted model

$$Y_i = 28.88 + e_i$$

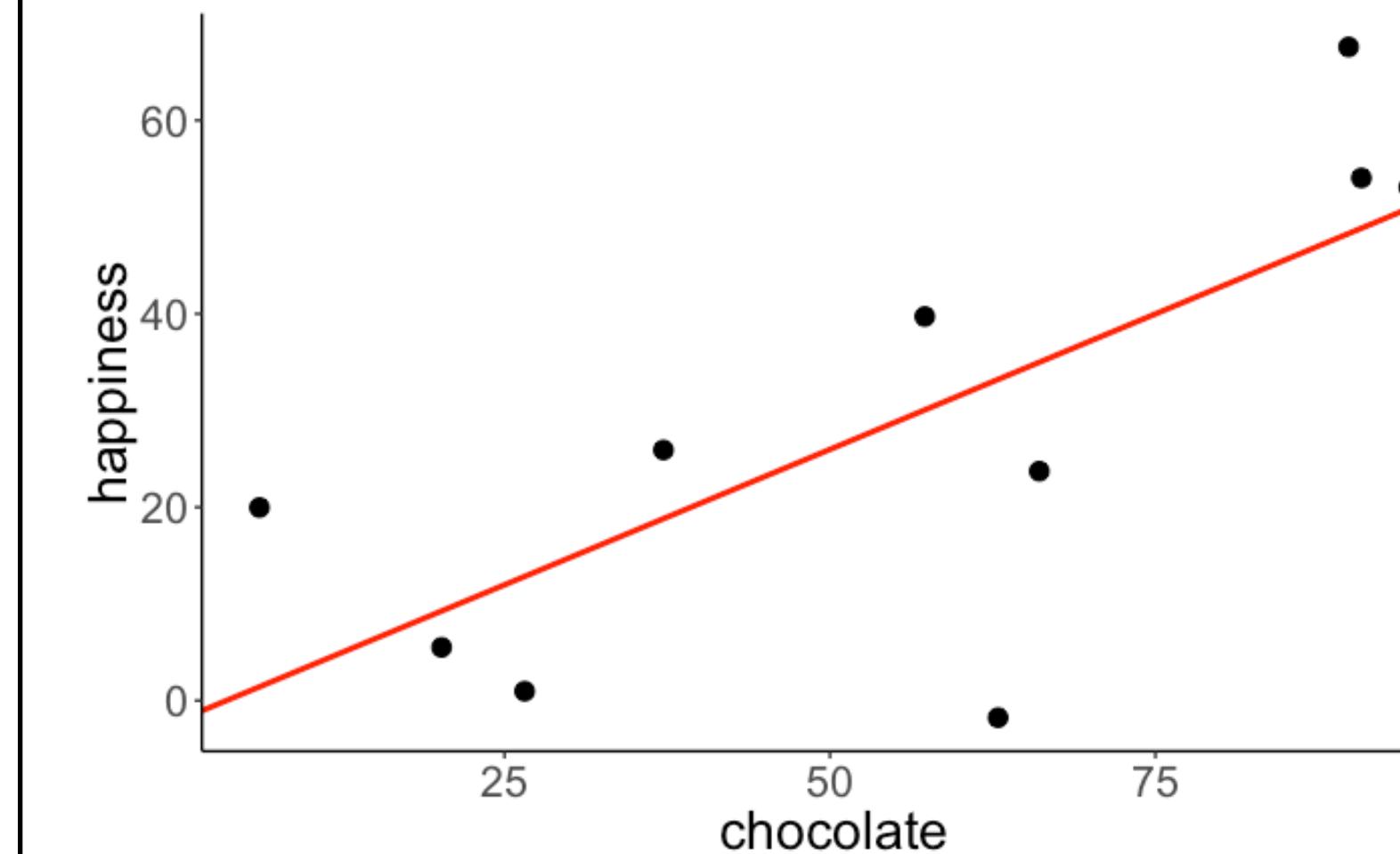
$H_1$ : Chocolate consumption and happiness are related.

### Model A

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

chocolate  
consumption

### Model prediction



### Fitted model

$$Y_i = -2.04 + 0.56X_i + e_i$$

# The general procedure

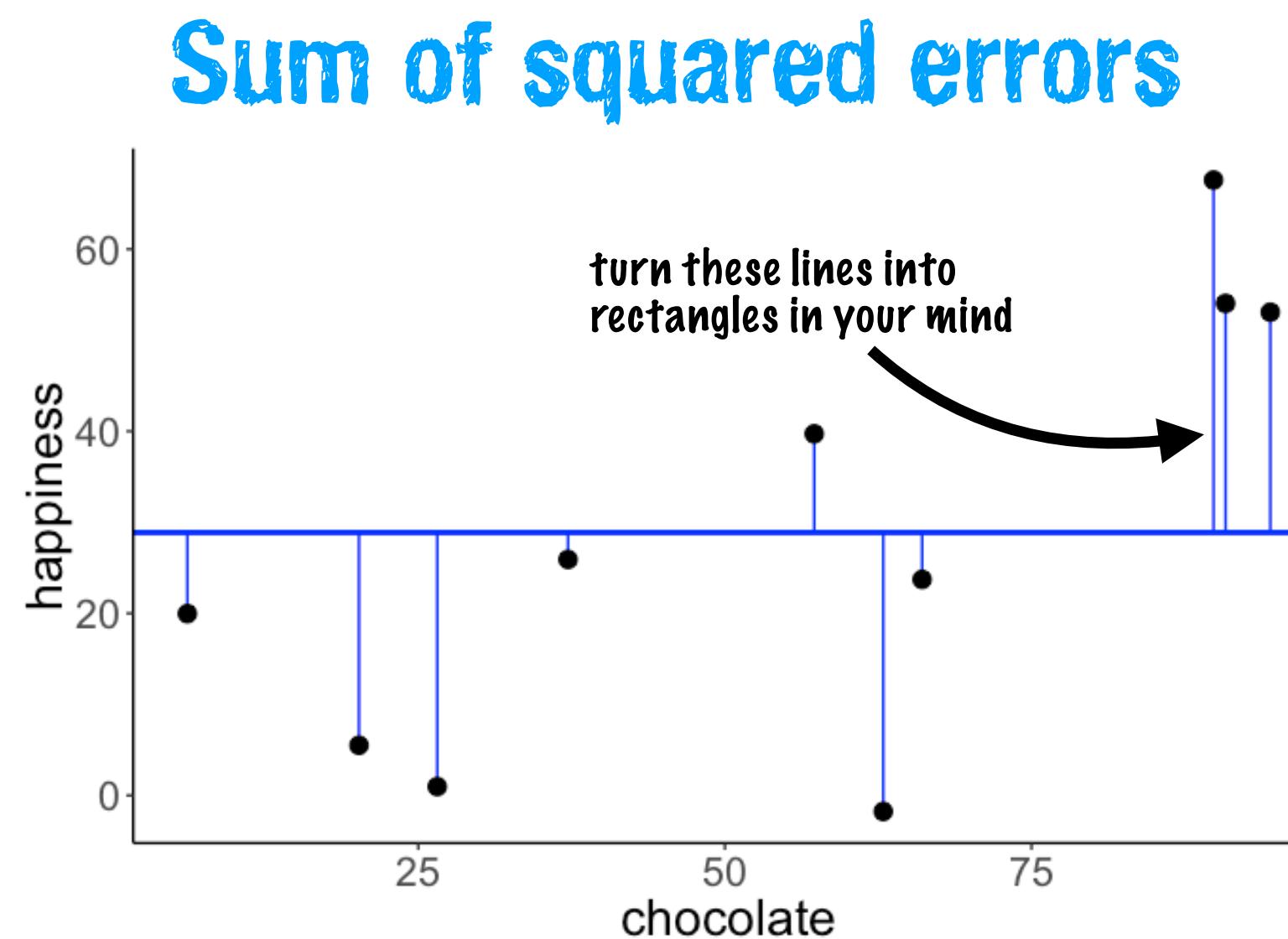
1. Define  $H_0$  as Model C (compact) and  $H_1$  as Model A (augmented)
2. Fit model parameters to the data
3. **Calculate the proportional reduction of error (PRE) in our sample**
4. Decide whether the augmented model is **worth it** by comparing the observed PRE in our sample to the sampling distribution of PRE (assuming that  $H_0$  is true)

# Calculate PRE

$$PRE = 1 - \frac{\text{ERROR}(A)}{\text{ERROR}(C)}$$

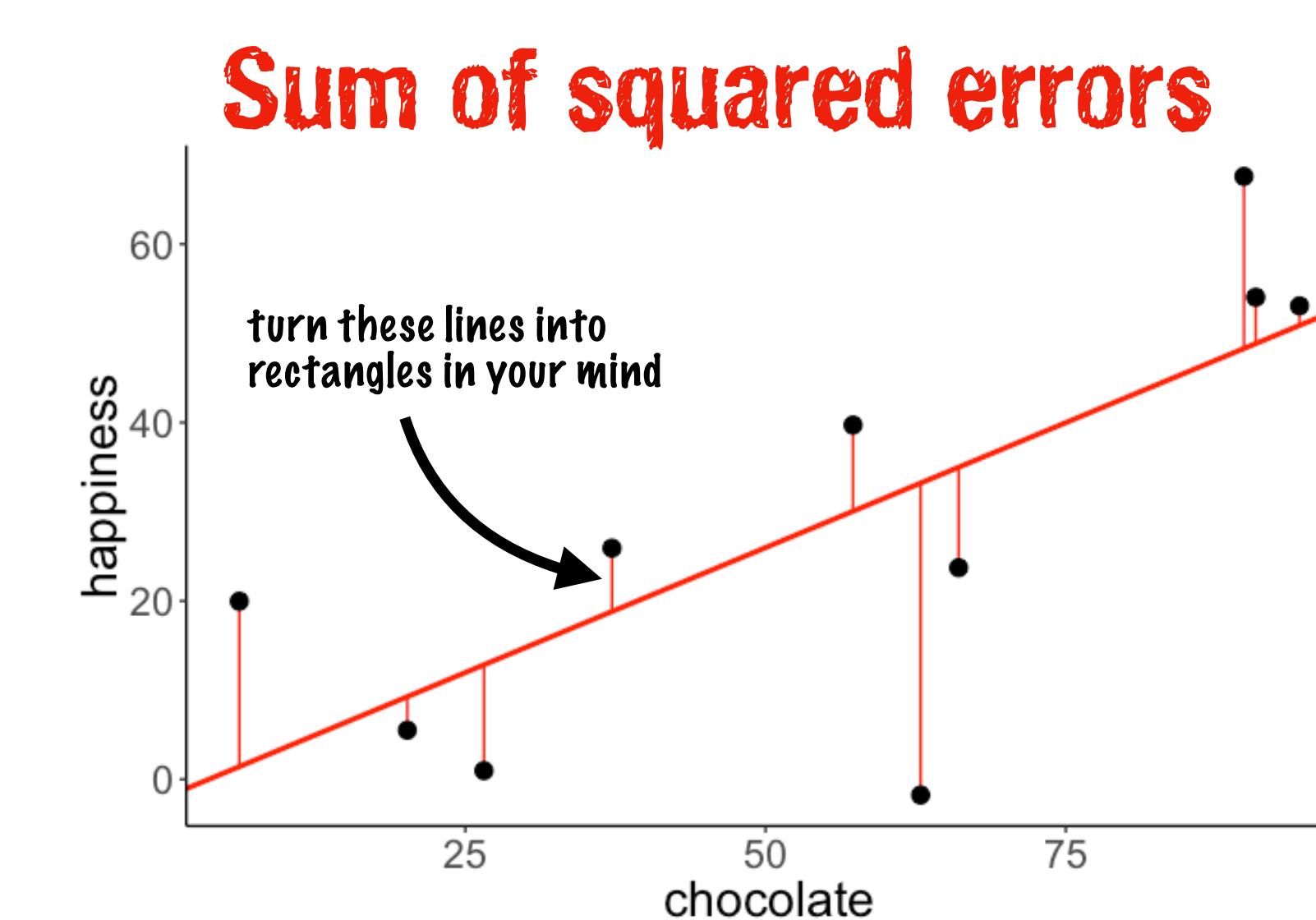
Both models were fit to minimize the sum of squared errors

OLS = Ordinary **least squares** regression



$$\text{SSE}(C) = 5215.016$$

$$\text{PRE} = 1 - \frac{2396.946}{5215.016} \approx 0.54$$



$$\text{SSE}(A) = 2396.946$$

**The augmented model  
reduces the error by 54%.**

# The general procedure

1. Define  $H_0$  as Model C (compact) and  $H_1$  as Model A (augmented)
2. Fit model parameters to the data
3. Calculate the proportional reduction of error (PRE) in our sample
4. Decide whether the augmented model is **worth it** by comparing the observed PRE in our sample to the sampling distribution of PRE (assuming that  $H_0$  is true)

# Decide whether it's **worth it**

- To compute the  $F$  statistic, we need:
  - PRE
  - number of parameters in Model C (PC) and Model A (PA)
  - number of observations  $n$

- more likely to be **worth it** if:
  1. PRE is high
  2. the number of additional parameters in A compared to C is low (**PRE per additional parameter**)
  3. the number of parameters that could have been added to model<sub>C</sub> to create model<sub>A</sub> but were not

**difference in parameters  
between models A and C**

$$F = \frac{\text{PRE}/(\text{PA} - \text{PC})}{(1 - \text{PRE})/(n - \text{PA})}$$

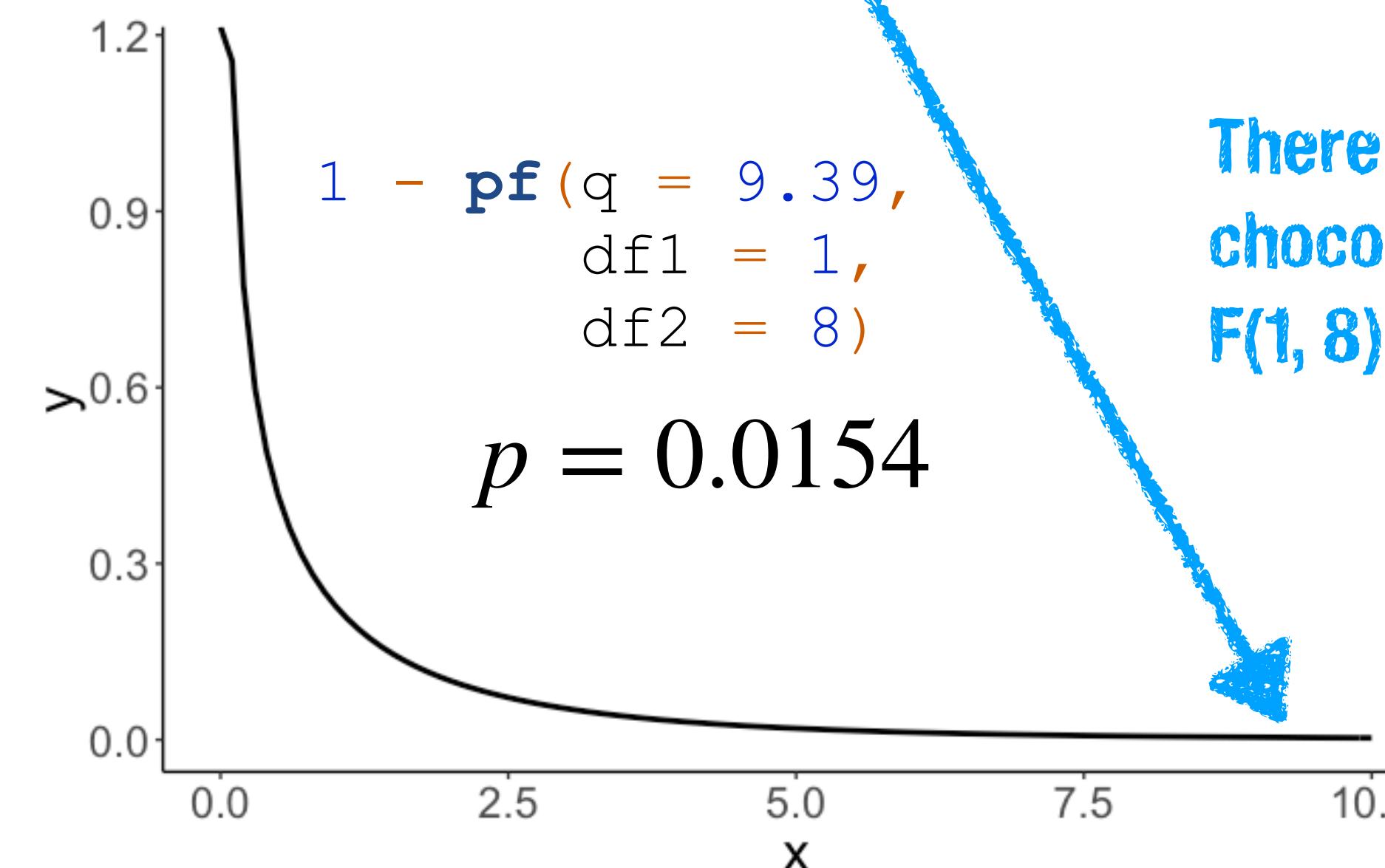
**number of observations  
vs. parameters in Model**

# Decide whether it's **worth it**

- To compute the  $F$  statistic, we need:

- PRE = 0.54
- PC = 1
- PA = 2
- $n$  = 10

$$\begin{aligned} F &= \frac{\text{PRE}/(\text{PA} - \text{PC})}{(1 - \text{PRE})/(n - \text{PA})} \\ &= \frac{0.54/(2 - 1)}{(1 - 0.54)/(10 - 2)} \\ &= 9.39 \end{aligned}$$



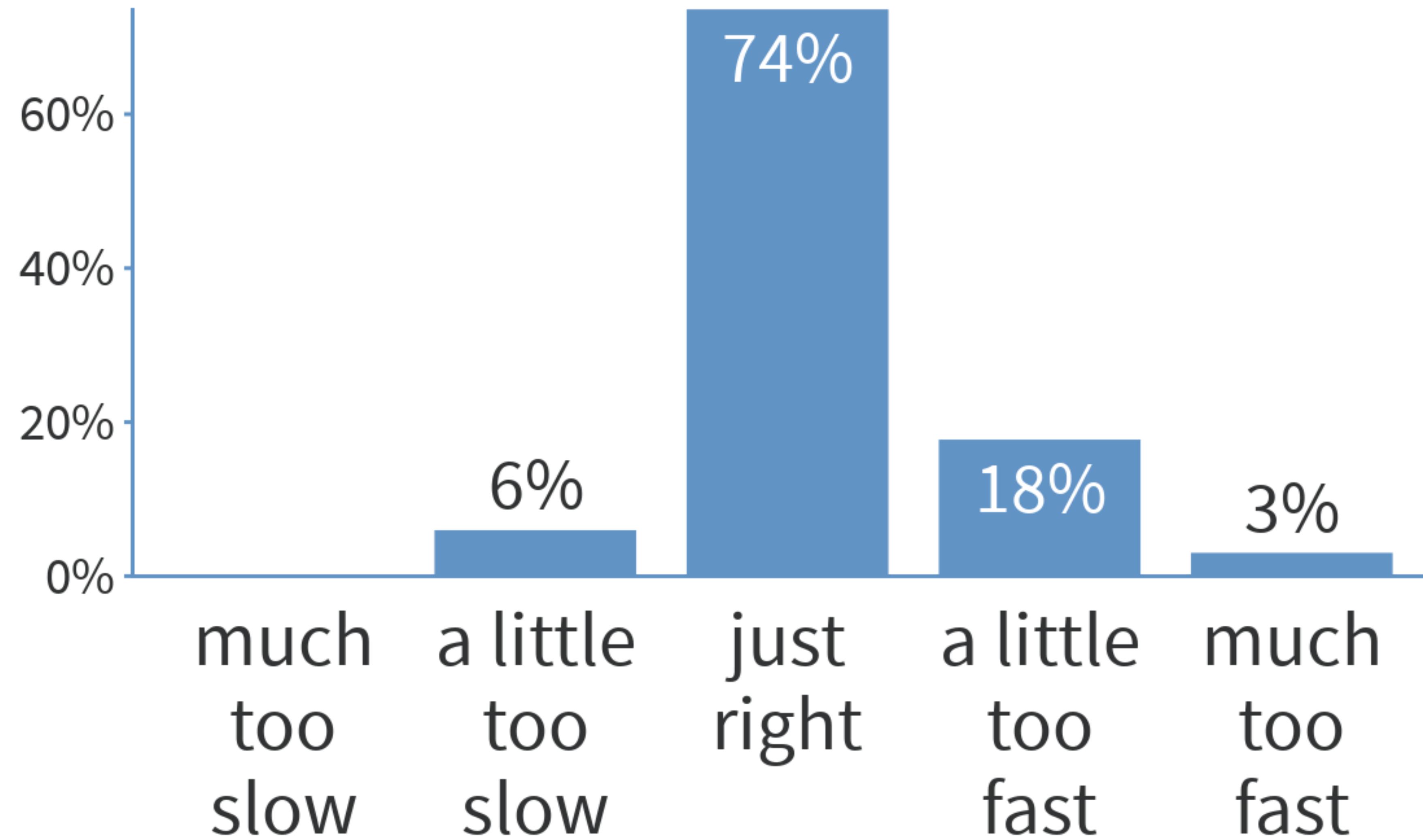
**There is a significant relationship between chocolate consumption and happiness**  
 $F(1, 8) = 9.39, p = .0154.$

# Summary

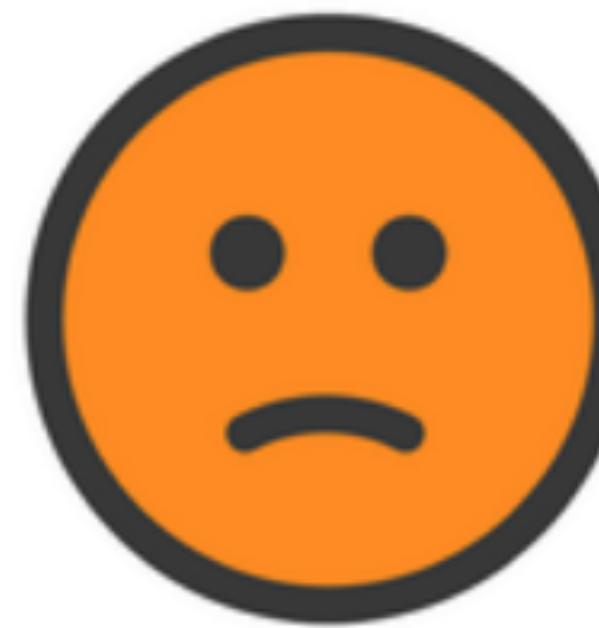
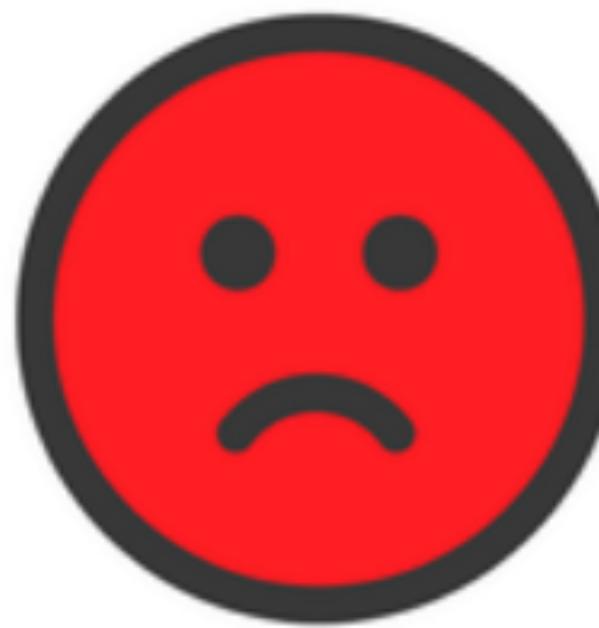
- Quick recap
- Modeling data
- Hypothesis testing as model comparison
- Correlation
  - Pearson's moment correlation
  - Spearman's rank correlation
- Regression

# **Feedback**

# How was the pace of today's class?



# How happy were you with today's class overall?



# What did you like about today's class? What could be improved next time?

A word cloud visualization showing student feedback on today's class. The words are arranged in a large, overlapping cluster. The most prominent words are 'understand', 'really', 'clear', 'well', 'tobi', 'things', 'class', 'concepts', 'background', 'foreground', 'important', 'limit', 'second', 'third', 'class.', 'record', 'explaination', 'lecture', 'interaction', 'time', 'fully', 'fix', 'us', 'assume', 'low', 'may', 'might', 'better', 'didn't', 'let', 'process', 'cover', and 'process'.

**Thank you!**