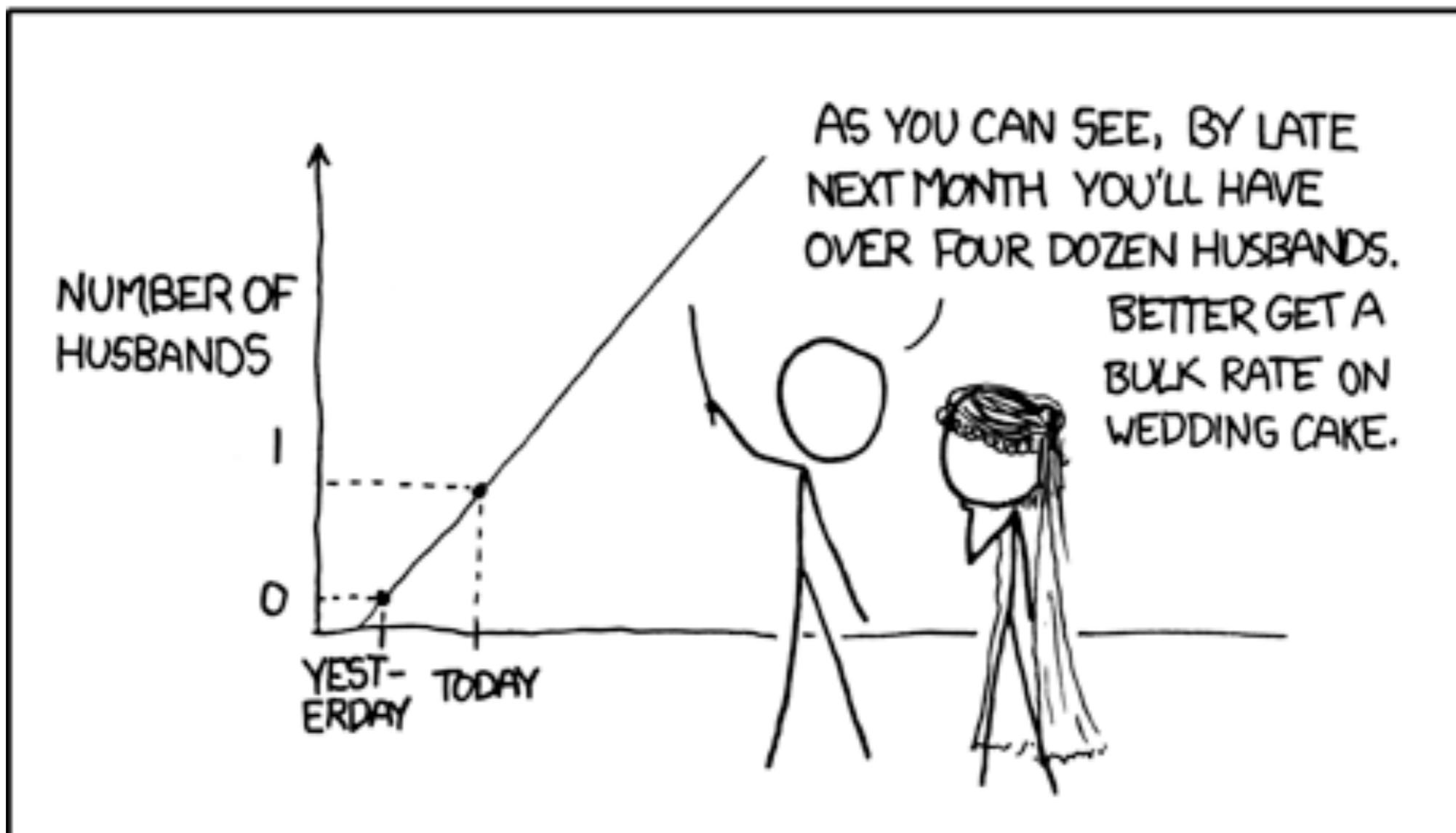


# Linear model 1

MY HOBBY: EXTRAPOLATING

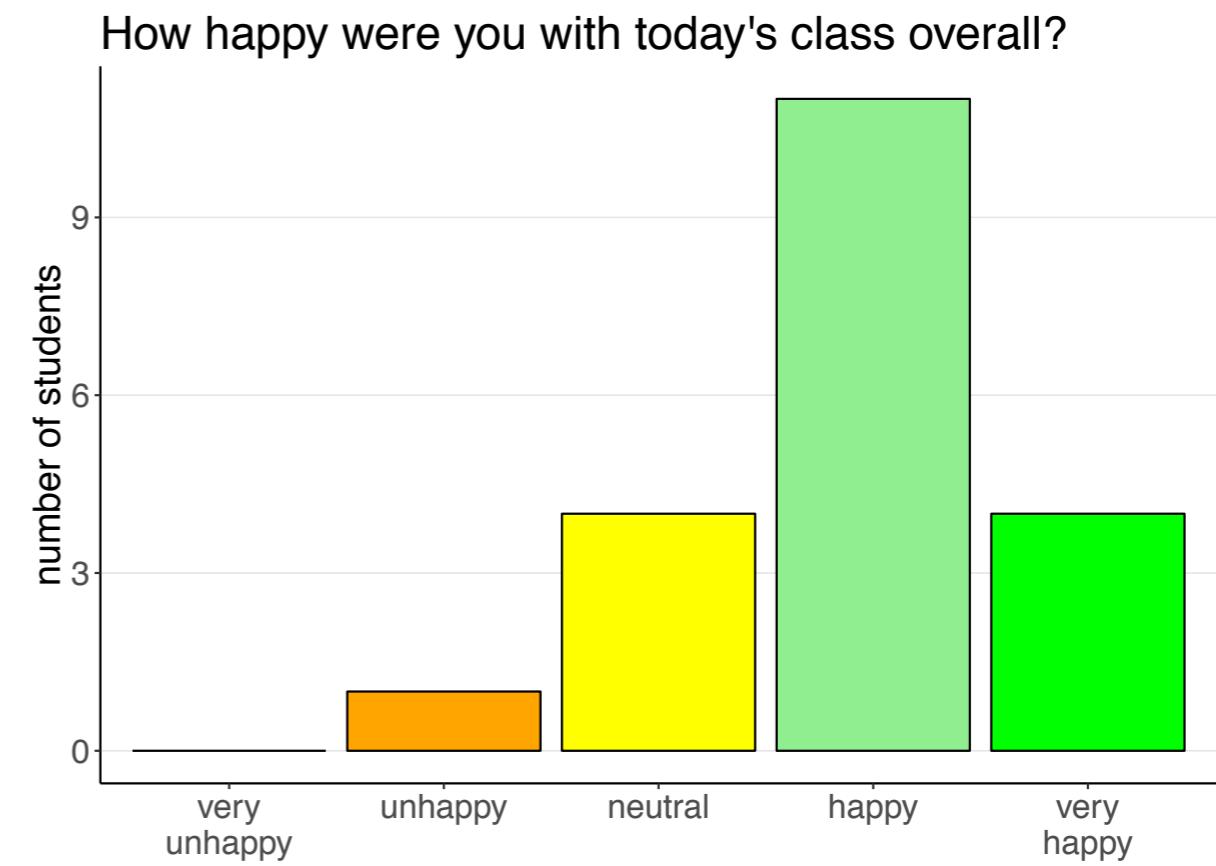
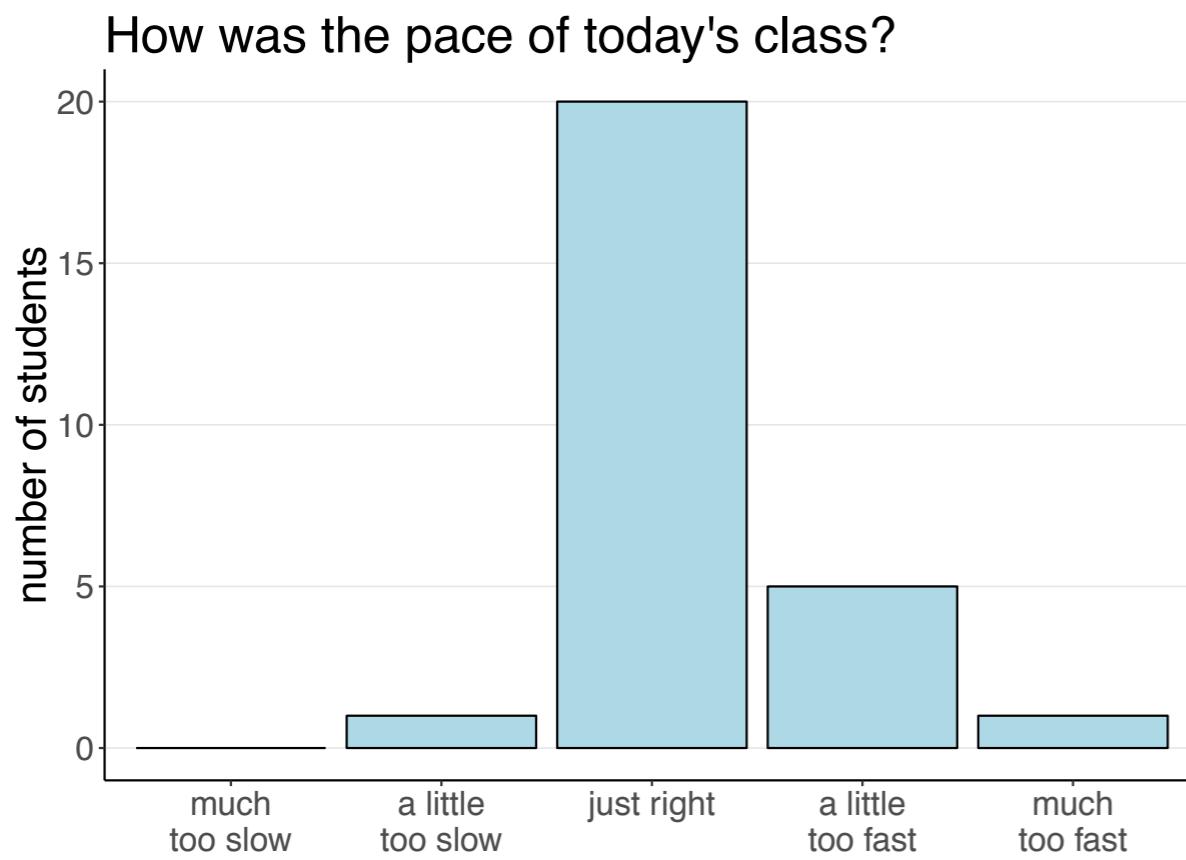


01/30/2019

# **Logistics**

# Your feedback

# Your feedback



# Your feedback

On Homework: I really liked the datacamp class. I found that a different explanation for key concepts complemented your class well. It was very helpful.

**Do you have recommendations for other datacamp lessons we should take weekly?**

I do!

# Your feedback

<https://psych252.github.io/>

## Schedule

Week	Day	Date	Topic	Content	Reading	Resources	Datacamp
1	M	7-Jan	Introduction	<ul style="list-style-type: none"><li>• Course introduction</li></ul>		<ul style="list-style-type: none"><li>• Cheatsheet R Studio</li><li>• Cheatsheet R Markdown 1</li><li>• Cheatsheet R Markdown 2</li><li>• R Markdown for class reports</li></ul>	<ul style="list-style-type: none"><li>• Introduction to R</li><li>• RStudio IDE 1</li><li>• RStudio IDE 2</li><li>• RMarkdown</li></ul>
	W	9-Jan	Visualization I	<ul style="list-style-type: none"><li>• Best practices</li><li>• Introduction to RStudio</li><li>• Introduction to <code>library(ggplot2)</code></li><li>• Reporting results using Rmarkdown</li></ul>	<ul style="list-style-type: none"><li>• Data visualization (#1)</li><li>• Data visualization (#3)</li></ul>	<ul style="list-style-type: none"><li>• Cheatsheet ggplot2</li></ul>	<ul style="list-style-type: none"><li>• ggplot part 1</li><li>• ggplot part 2</li><li>• Reporting</li></ul>
	F	11-Jan	Visualization II	<ul style="list-style-type: none"><li>• Making nice plots</li></ul>	<ul style="list-style-type: none"><li>• Data visualization (#4)</li><li>• Data visualization (#8)</li><li>• R for Data Science (#27)</li></ul>	<ul style="list-style-type: none"><li>• Cheatsheet shiny</li></ul>	<ul style="list-style-type: none"><li>• ggplot part 3</li><li>• Shiny 1</li><li>• Shiny 2</li></ul>

# Your feedback

I was a little confused at the end about how we got the distribution on possible values for the F statistic? I understand that we can generate the F statistic from our data, but then I don't understand how we knew that was an unlikely one.

**I will give it another shot  
(and you'll have to do  
something similar for  
your homework)**

# Your feedback

Some times the material can be hard to fully understand, but I thought you helped walk us through it. **The gifs are great**

**playing around with code  
is key to understanding**

# **Homework 4**

# Homework 4

## Part 2

## Part 1

INTERACTIVE COURSE

### Correlation and Regression

[Continue Course](#)

4 hours | 18 Videos | 58 Exercises | 32,959 Participants | 4,200 XP



**Course Description**

Ultimately, data analysis is about understanding relationships among variables. Exploring data with multiple variables requires new, more complex tools, but enables a richer set of comparisons. In this course, you will learn how to describe relationships between two numerical quantities. You will characterize these relationships graphically, in the form of summary statistics, and through simple linear regression models.

**1 Visualizing two variables** FREE

In this chapter, you will learn techniques for exploring bivariate relationships.

[VIEW CHAPTER DETAILS](#) [Continue Chapter](#)

This course is part of these tracks:

- Data Analyst with R
- Data Scientist with R
- Intro to Statistics with R

  
Ben Baumer  
Assistant Professor at Smith College

My name goes here

The names of the people I have worked with go here

2019-01-28 23:42:22

#### Instructions

This homework is due by **Tuesday, February 5th, 8:00pm**.

This homework has **two parts**:

1. Complete the data camp assignment for this week: [Correlation and Regression](#)
2. Complete this problem sheet. Upload **only** the pdf file containing the code and plots to Canvas and name the file: `yourlastname_modeling_homework.pdf`

Make sure to explain what your doing (or trying to do)! Either as comments within the code, or as text underneath the code blocks.

#### Working with probabilities

To get familiar with the different functions for computing probabilities, we'll use the exponential distribution below.

Assuming an exponential distribution with the rate parameter  $\lambda = 1$ :

**What's the probability of observing a value larger than 4?**

# write code here

**What's the probability of observing a value between 2 and 3?**

# write code here

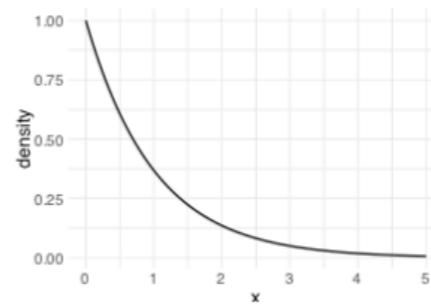


Figure 1: The exponential distribution with  $\lambda = 1$ .

# Plan for today

- Quick review of statistical inference in frequentist statistics
- Correlation
  - Pearson's moment correlation
  - Spearman's rank correlation
- Regression
  - The conceptual tour
  - The R route

# **Quick review of statistical inference in frequentist statistics**

# 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)

# 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

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

Model A has  
15% less error  
than Model C.

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

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

- we have to construct a sampling distribution of PRE assuming that  $H_0$  is true
- and then compare the observed value of PRE to that distribution

## Population distribution

$$Y_i = 75 + \epsilon_i \quad \epsilon_i \sim \mathcal{N}(\mu = 0, \sigma = 5)$$

### Model C

$$Y_i = 75 + e_i$$

0 parameters

### Model A

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

1 parameter

# Sampling distribution of PRE

```
1 # simulation parameters
2 n_samples = 1000
3 sample_size = 50
4 mu = 75 # true mean of the distribution
5 sigma = 5 # true standard deviation of the errors
6
7 # function to draw samples from the population distribution
8 fun.draw_sample = function(sample_size, sigma) {
9   sample = mu + rnorm(sample_size, mean = 0, sd = sigma)
10 }
11
12 # draw samples
13 samples = n_samples %>%
14   replicate(fun.draw_sample(sample_size, sigma)) %>%
15   t() # transpose the resulting matrix (i.e. flip rows and columns)
```

# Sampling distribution of PRE

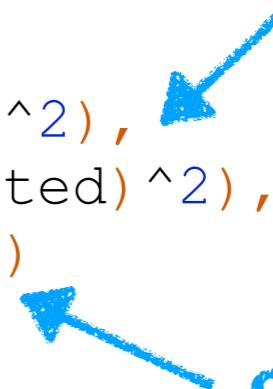
```
17 # put samples in data frame and compute PRE
18 df.samples = samples %>%
19   as_tibble(.name_repair = "unique") %>%
20   mutate(sample = 1:n()) %>%
21   gather("index", "value", -sample) %>%
22   mutate(compact = mu) %>%
23   group_by(sample) %>%
24   mutate(augmented = mean(value))
```

sample	index	value	compact	augmented
1	1	73.43	75	74.75
	2	76.38	75	74.75
	3	79.92	75	74.75
	4	72.33	75	74.75
	5	77.75	75	74.75
2	1	79.84	75	73.92
	2	78.44	75	73.92
	3	79.49	75	73.92
	4	71.81	75	73.92
	5	79.57	75	73.92
3	1	78.99	75	74.93
	2	67.28	75	74.93
	3	77.74	75	74.93
	4	73.73	75	74.93
	5	73.49	75	74.93

# Sampling distribution of PRE

```
17 # put samples in data frame and compute PRE
18 df.samples = samples %>%
19   as_tibble(.name_repair = "unique") %>%
20   mutate(sample = 1:n()) %>%
21   gather("index", "value", -sample) %>%
22   mutate(compact = mu) %>%
23   group_by(sample) %>%
24   mutate(augmented = mean(value)) %>%
25   summarize(sse_compact = sum((value - compact)^2),
26             sse_augmented = sum((value - augmented)^2),
27             pre = 1 - sse_augmented/sse_compact)
```

calculate SSE  
for each model



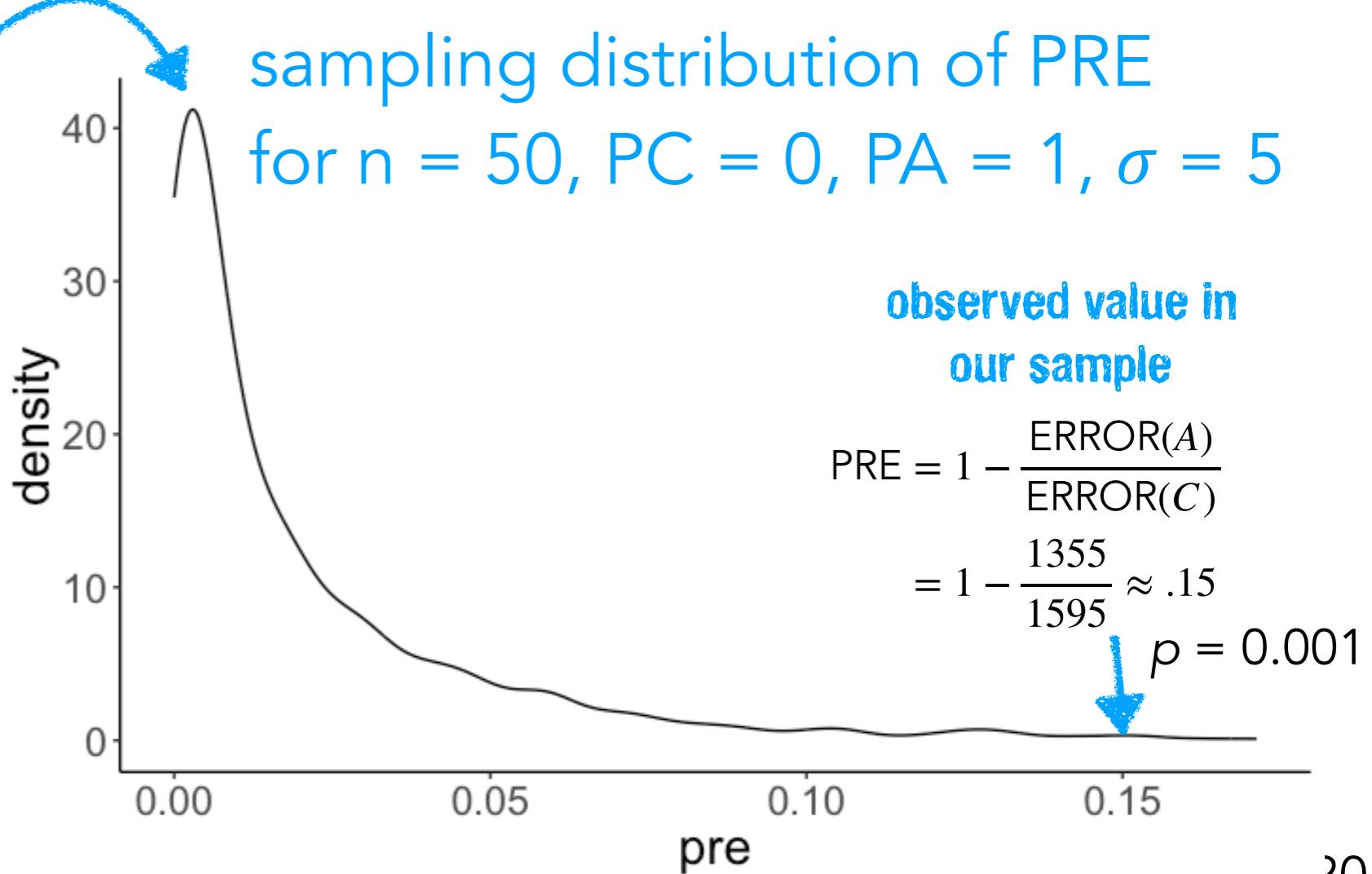
calculate PRE

sample	sse_compact	sse_augmented	pre
1	1244.66	1241.52	0.00
2	953.25	894.95	0.06
3	1083.60	1083.32	0.00
4	888.06	798.19	0.10
5	1246.57	1233.25	0.01

# Sampling distribution of PRE

```
29 # sampling distribution for PRE  
30 ggplot(data = df.samples,  
31         mapping = aes(x = pre)) +  
32         stat_density(geom = "line")  
33  
34 # p-value for our sample  
35 df.samples %>%  
36 summarize(p_value = sum(pre >= df.summary$pre) / n())
```

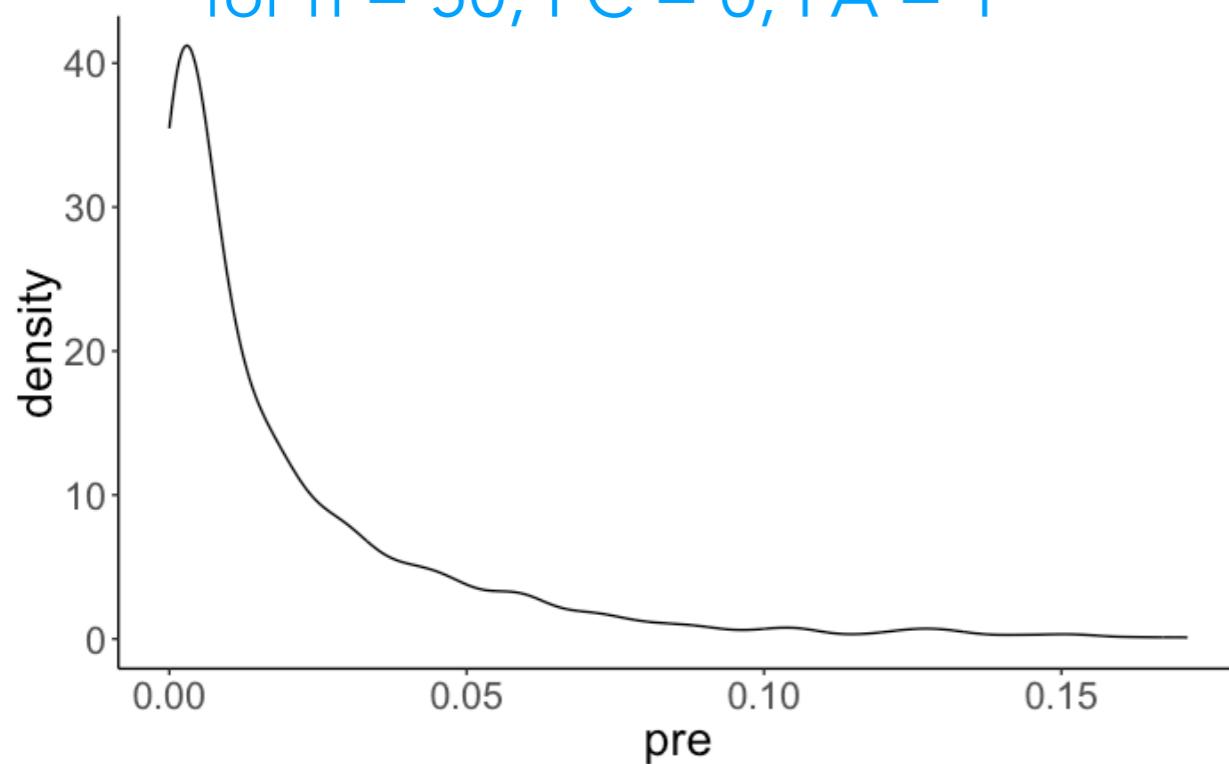
sample	sse_compact	sse_augmented	pre
1	1244.66	1241.52	0.00
2	953.25	894.95	0.06
3	1083.60	1083.32	0.00
4	888.06	798.19	0.10
5	1246.57	1233.25	0.01



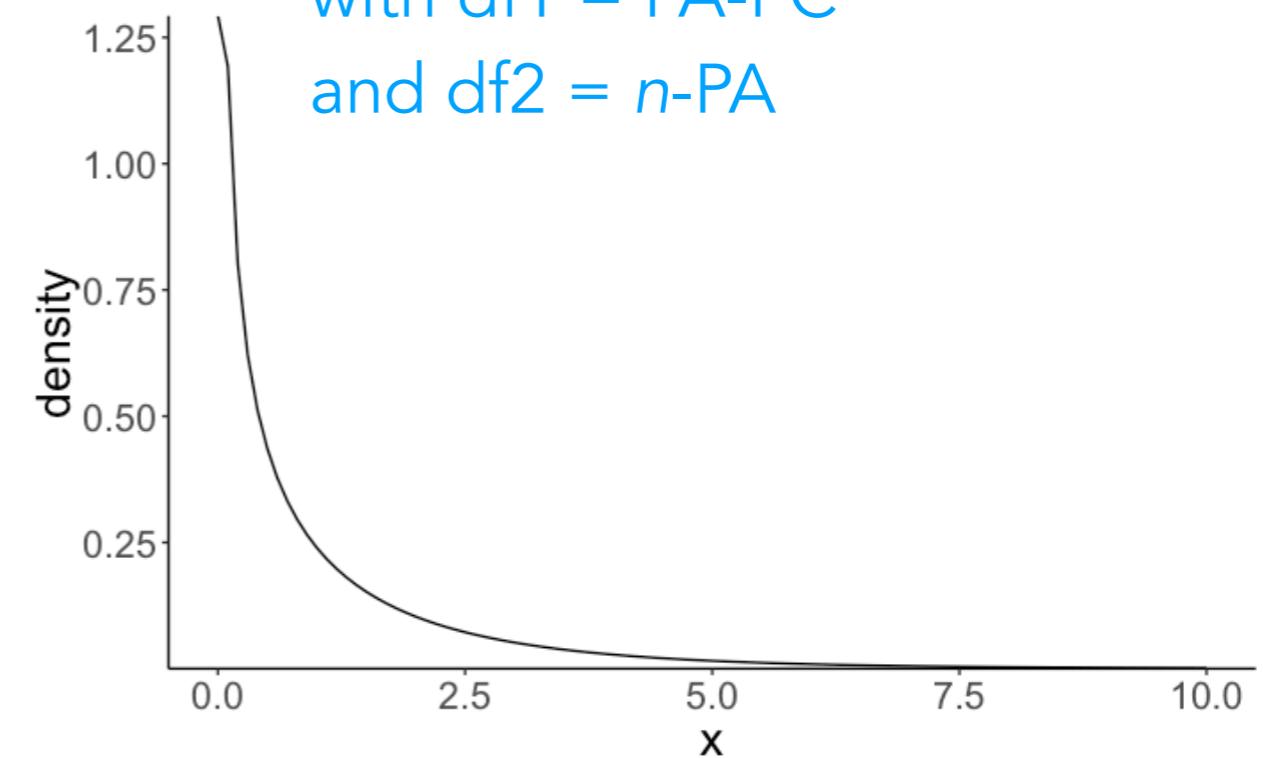
# Sampling distribution of PRE

deterministic mapping

sampling distribution of PRE  
for  $n = 50$ ,  $PC = 0$ ,  $PA = 1$



$F(df1, df2)$  distribution  
with  $df1 = PA - PC$   
and  $df2 = n - PA$



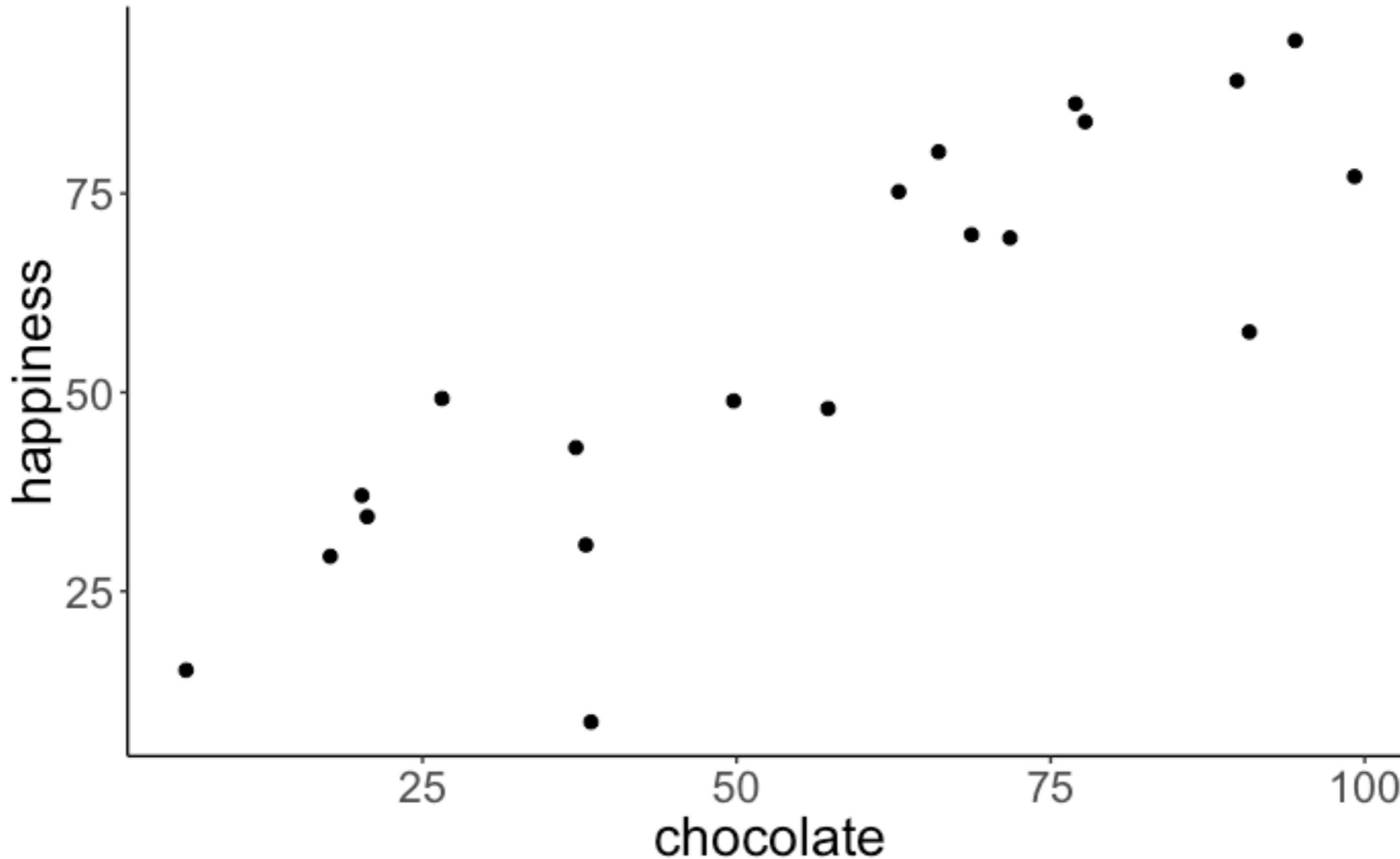
we use the F-distribution since it comes with R (and is the standard statistic to report)

# 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)

# **Correlation**

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

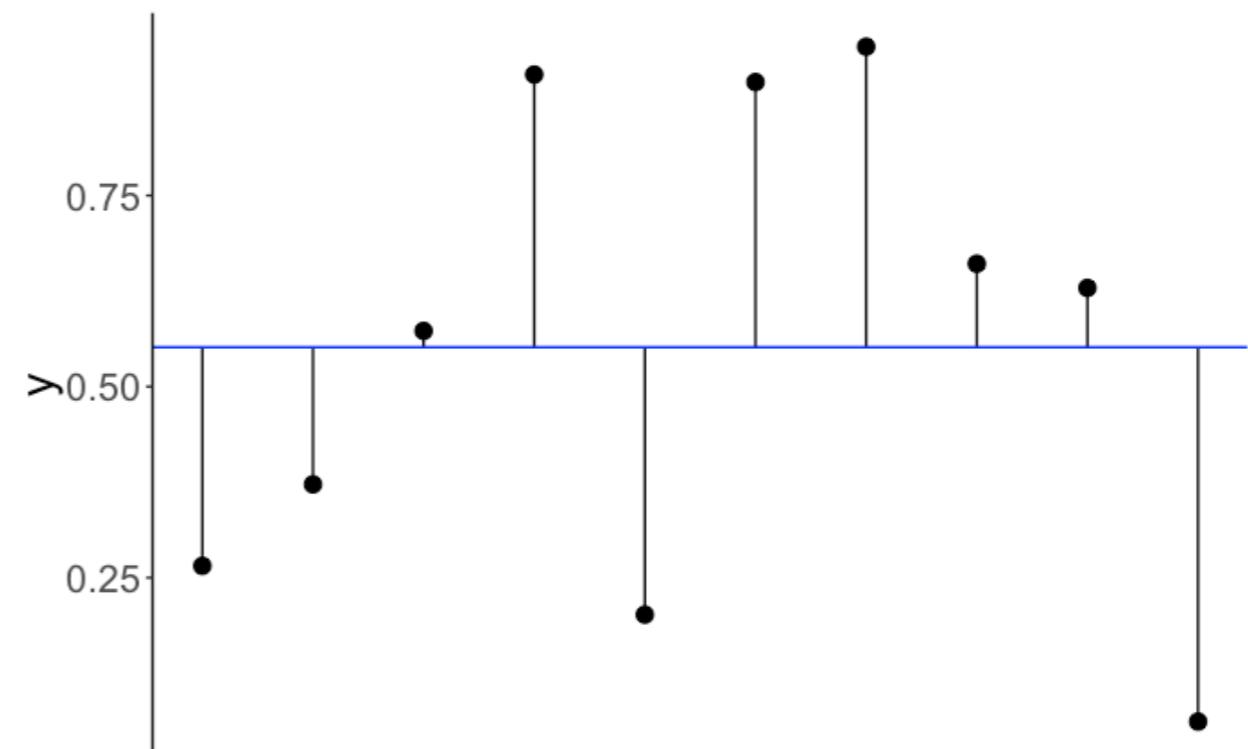
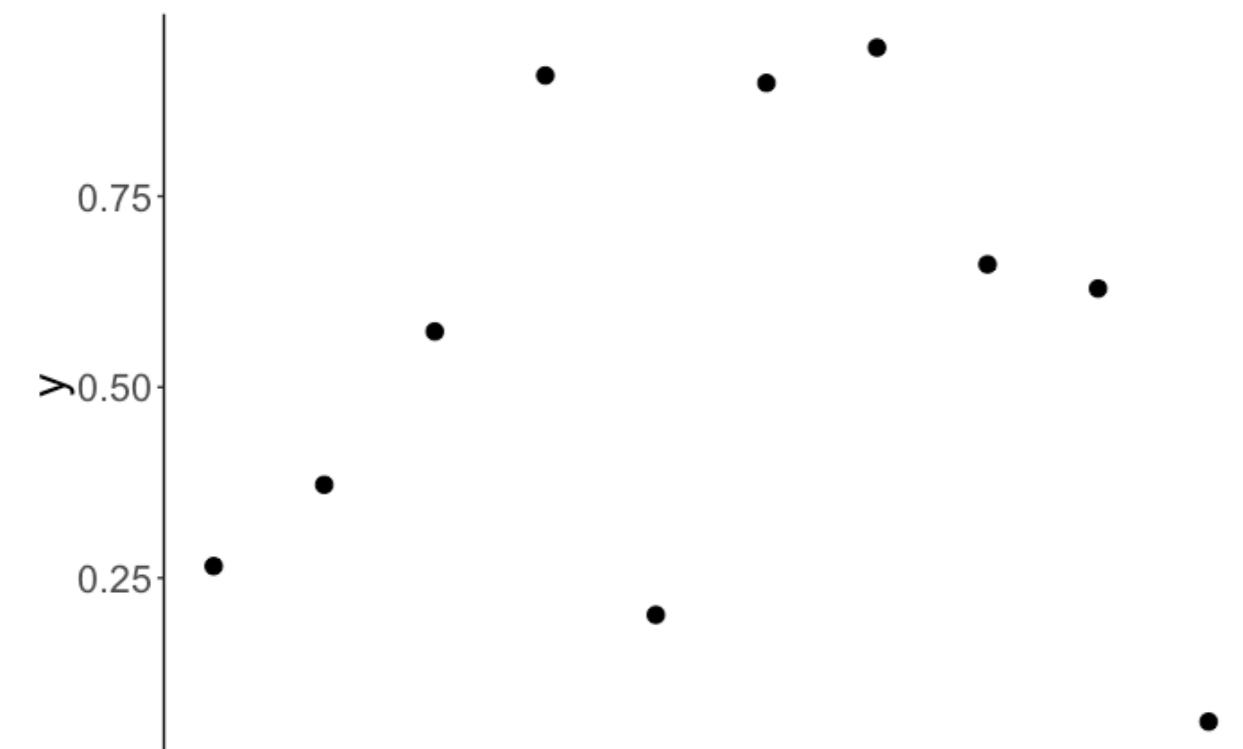


correlation = a measure of the relationship  
between two variables

## variance

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

sum of squared errors

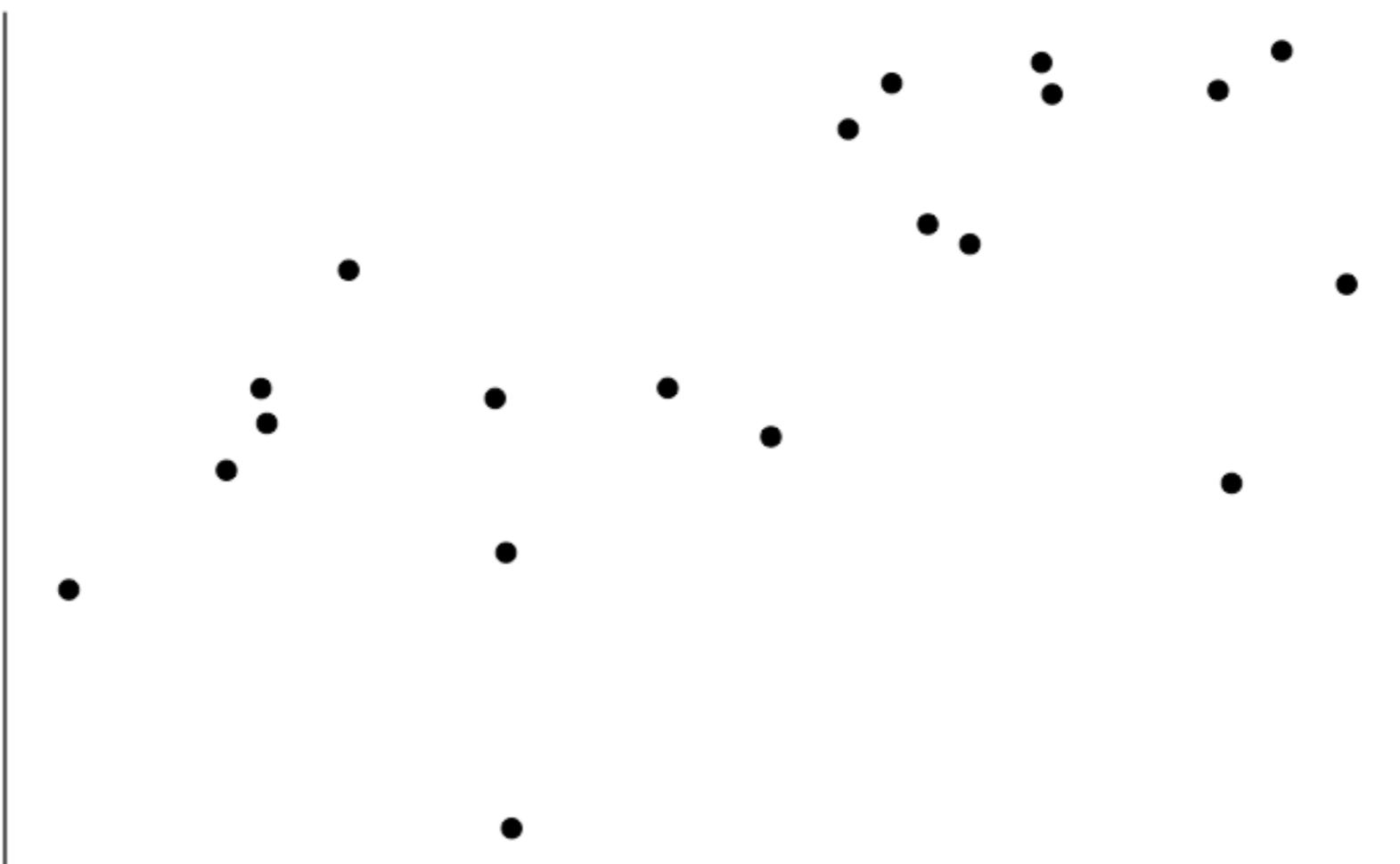


(I was too lazy to draw rectangles ...)

How well does the mean capture the data?

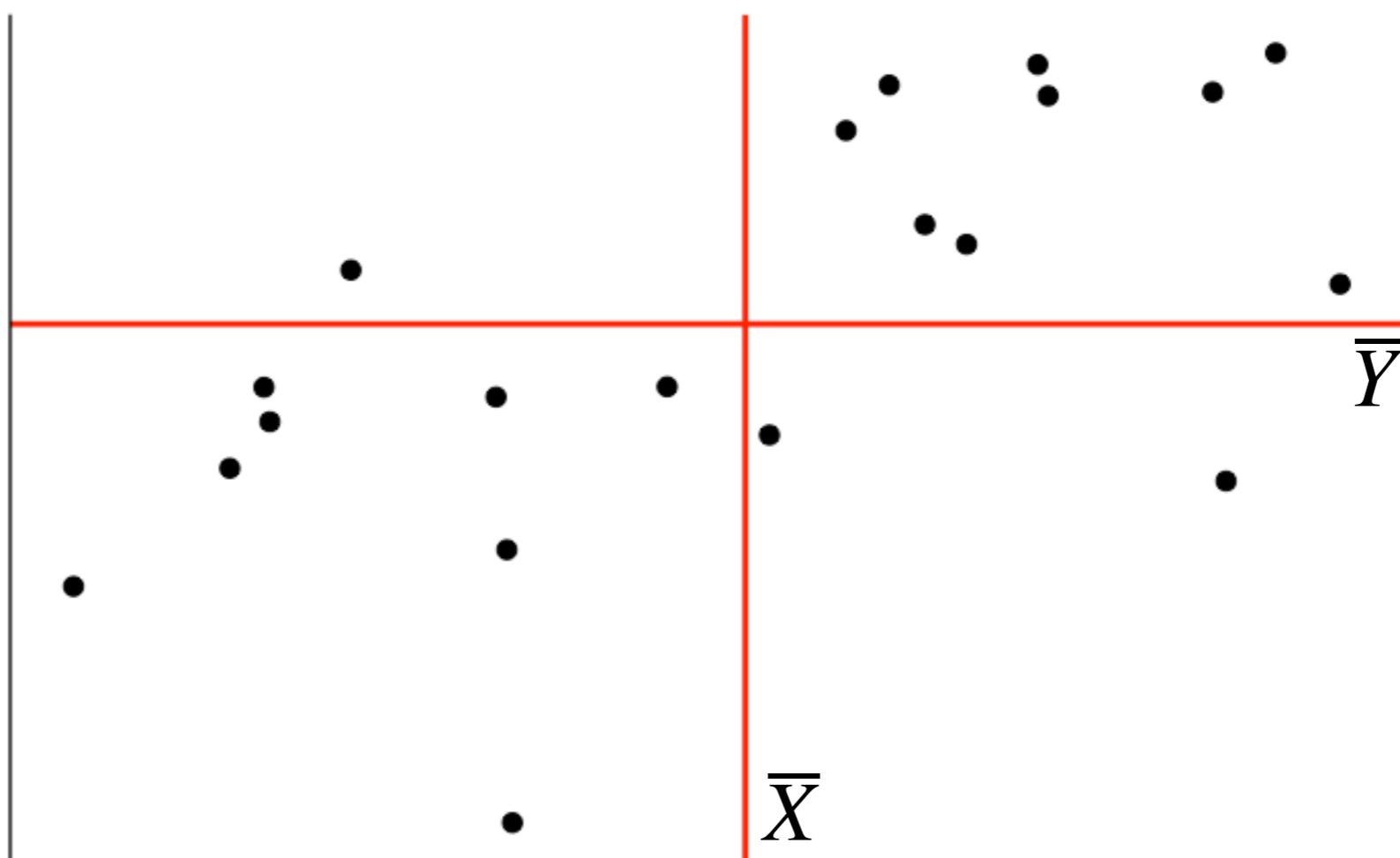
## covariance

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



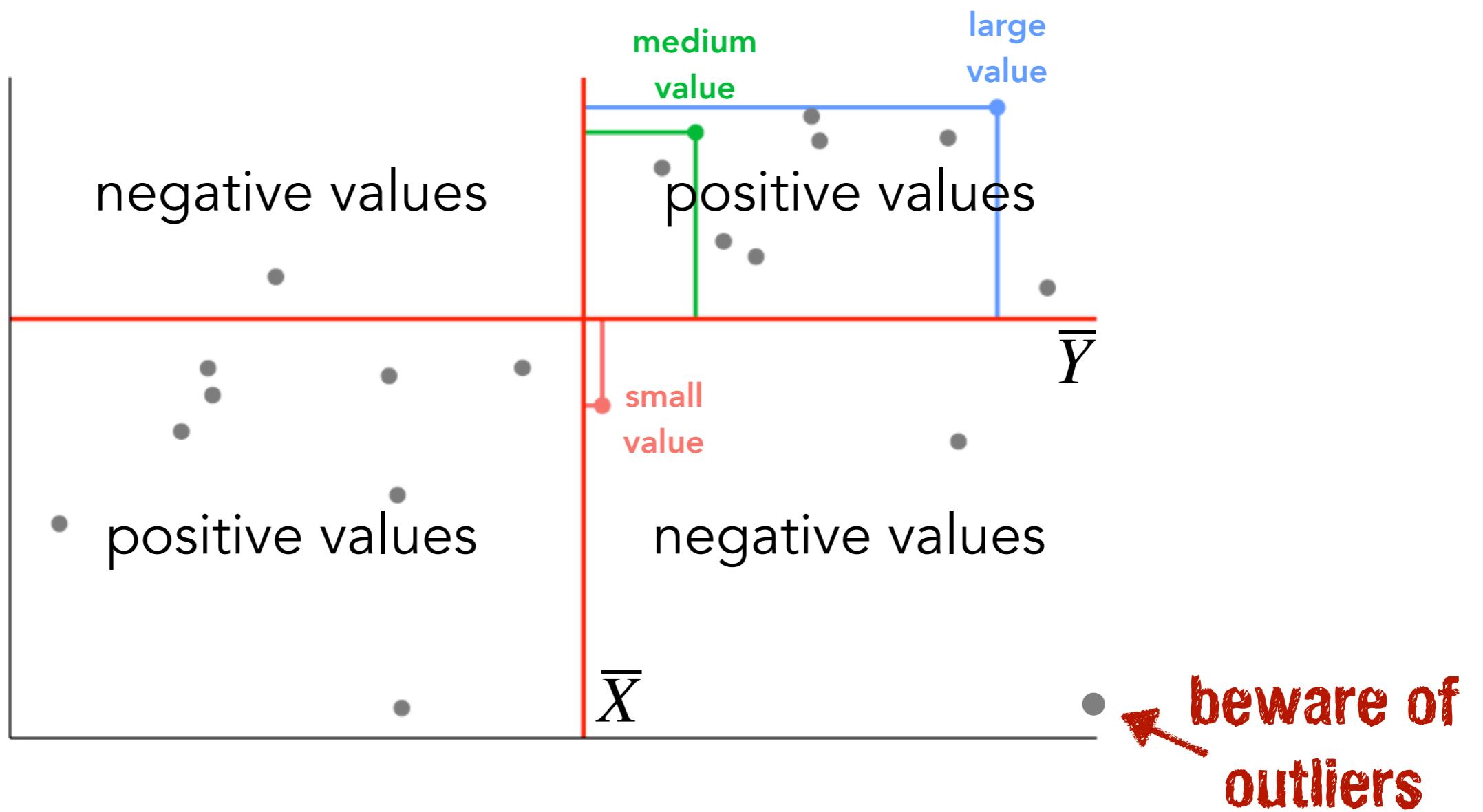
## covariance

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



## covariance

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



## 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

## Pearson correlation

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

standardized covariation  
(dividing by the standard deviations)

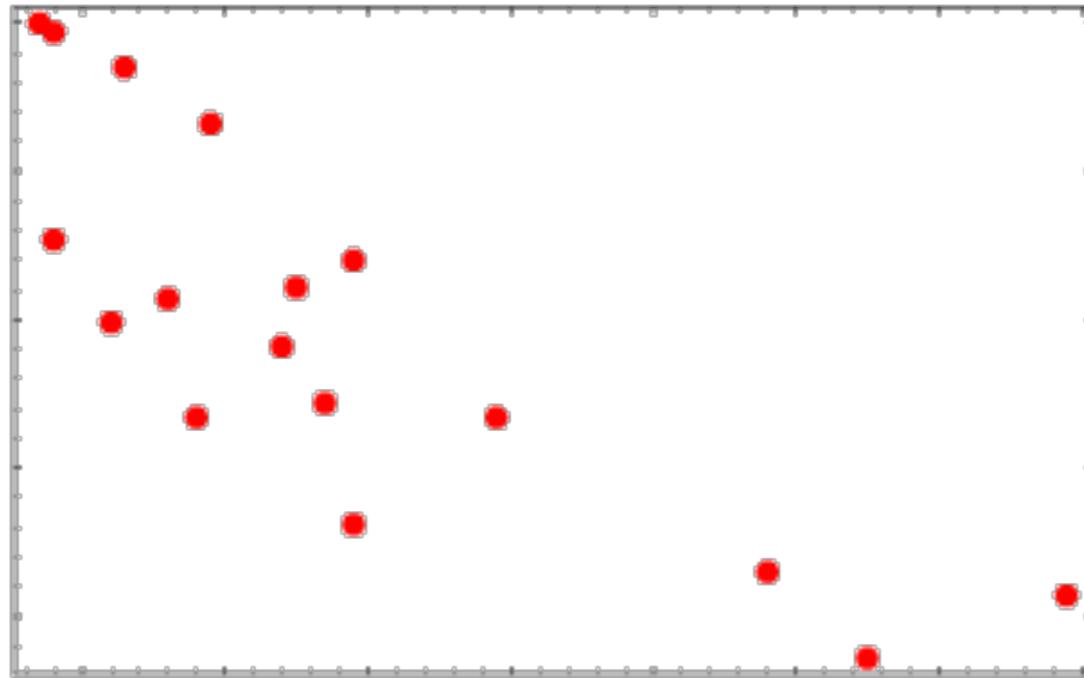
# Properties of the Pearson correlation

- standardized:  $-1 \leq r \leq 1$
- scale independent (for both X and Y)
- commutativity:  $r(X, Y) = r(Y, X)$
- sign determines the direction of dependence
- captures **linear dependence** only

association not  
causation



# Who is the correlation champion?



The faster you  
respond the more  
points you get!

-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

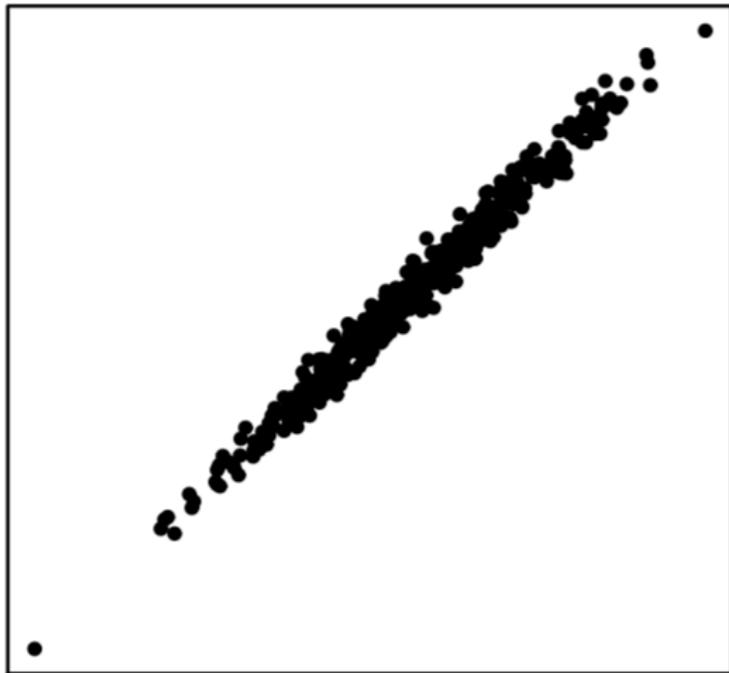
0.75 : 1

# **Who is the correlation champion?**

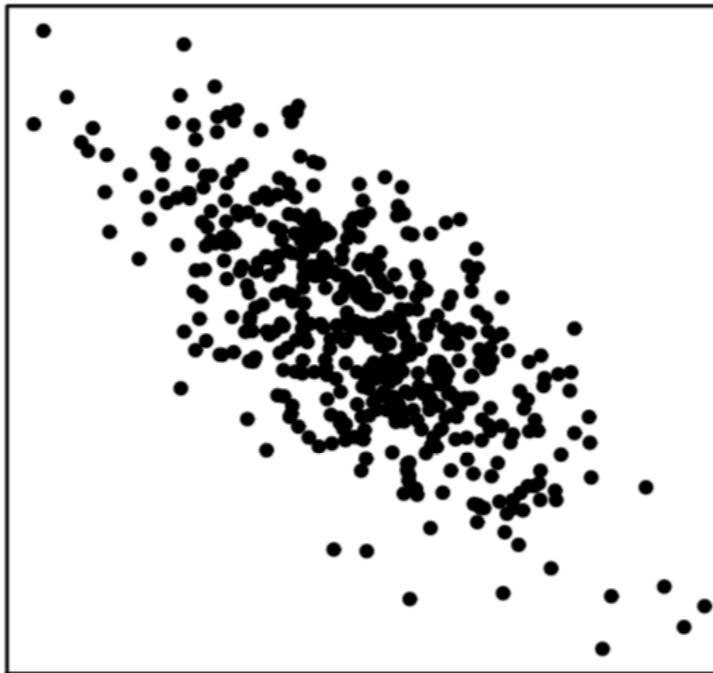
Get ready to compete!

# Solution

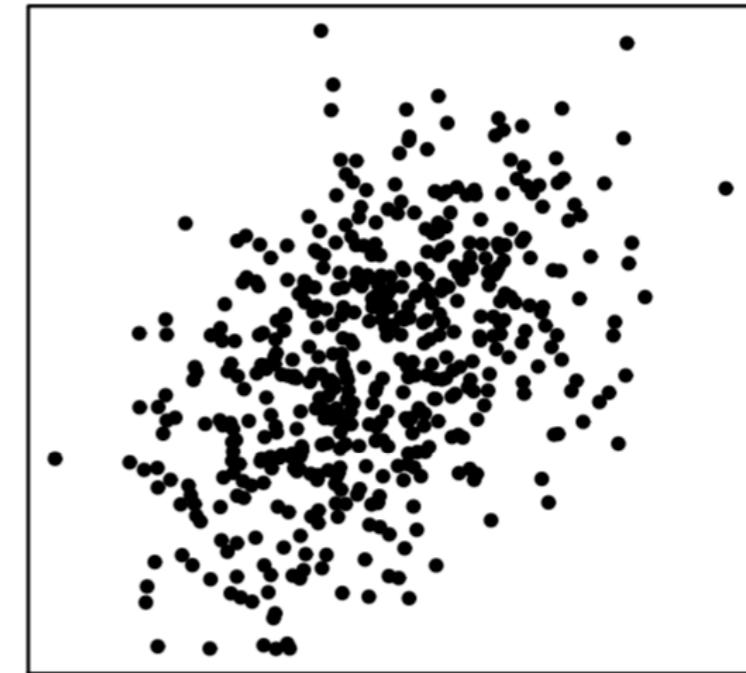
$r = 0.99$



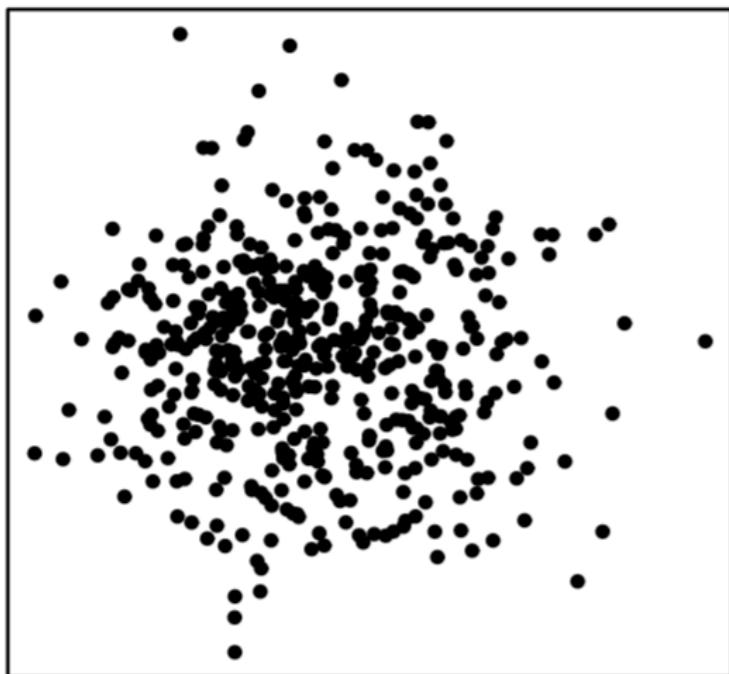
$r = -0.7$



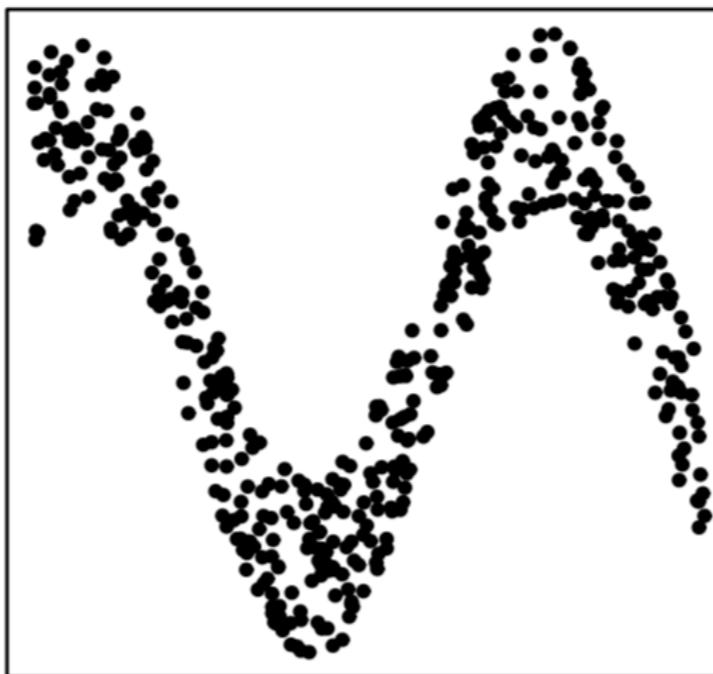
$r = 0.5$



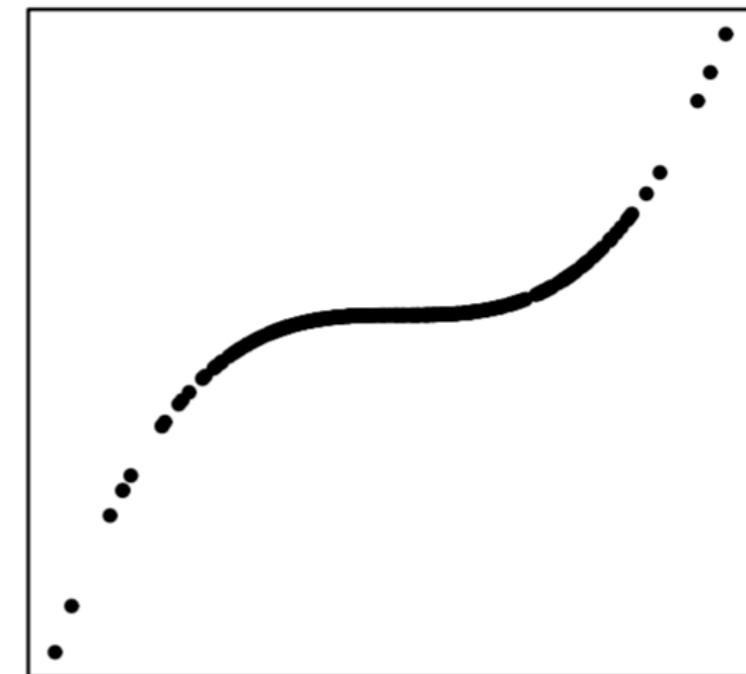
$r = -0.05$



$r = 0.03$



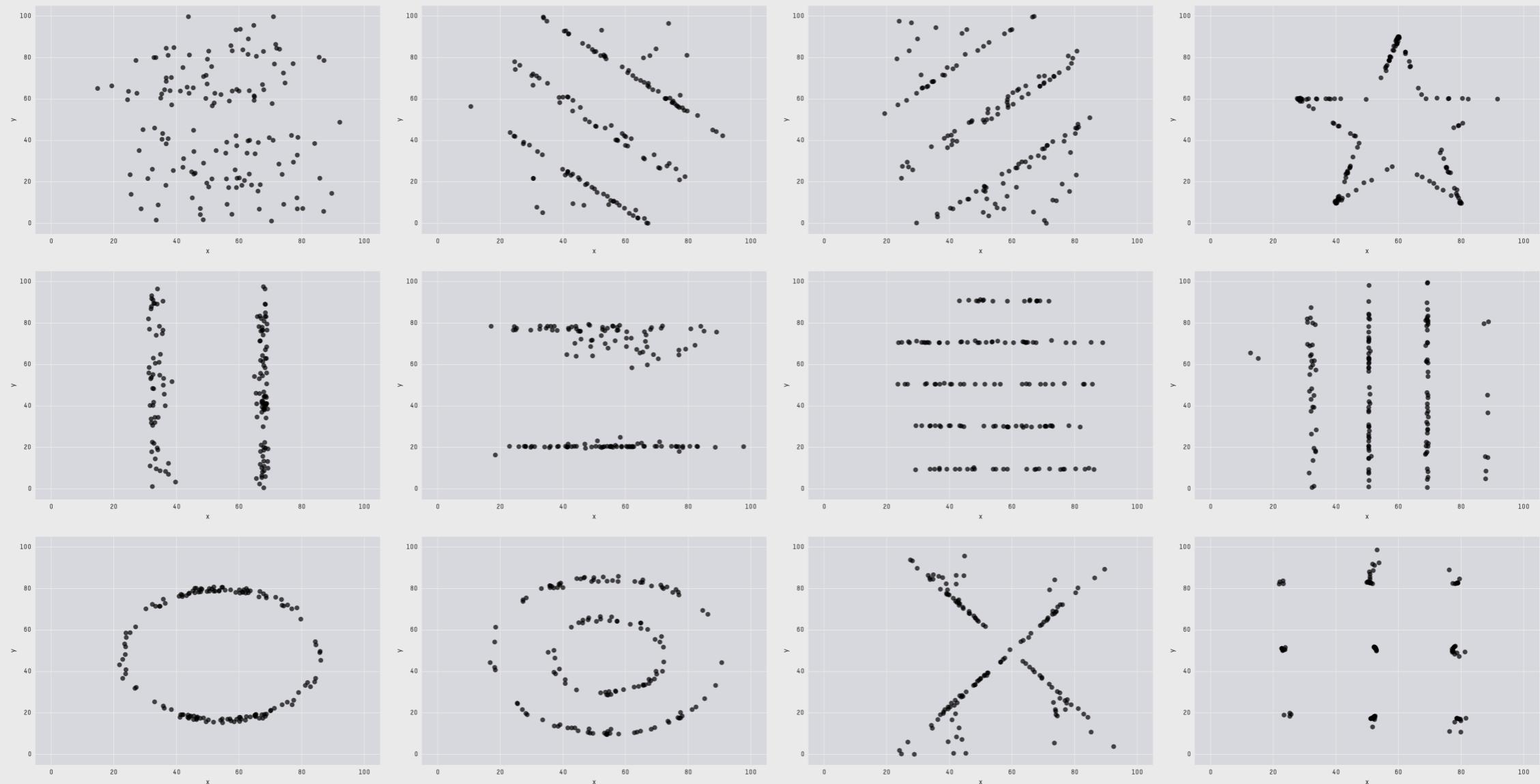
$r = 0.77$



# Be careful about interpreting correlations!



X Mean: 54.26  
Y Mean: 47.83  
X SD : 16.76  
Y SD : 26.93  
Corr. : -0.06

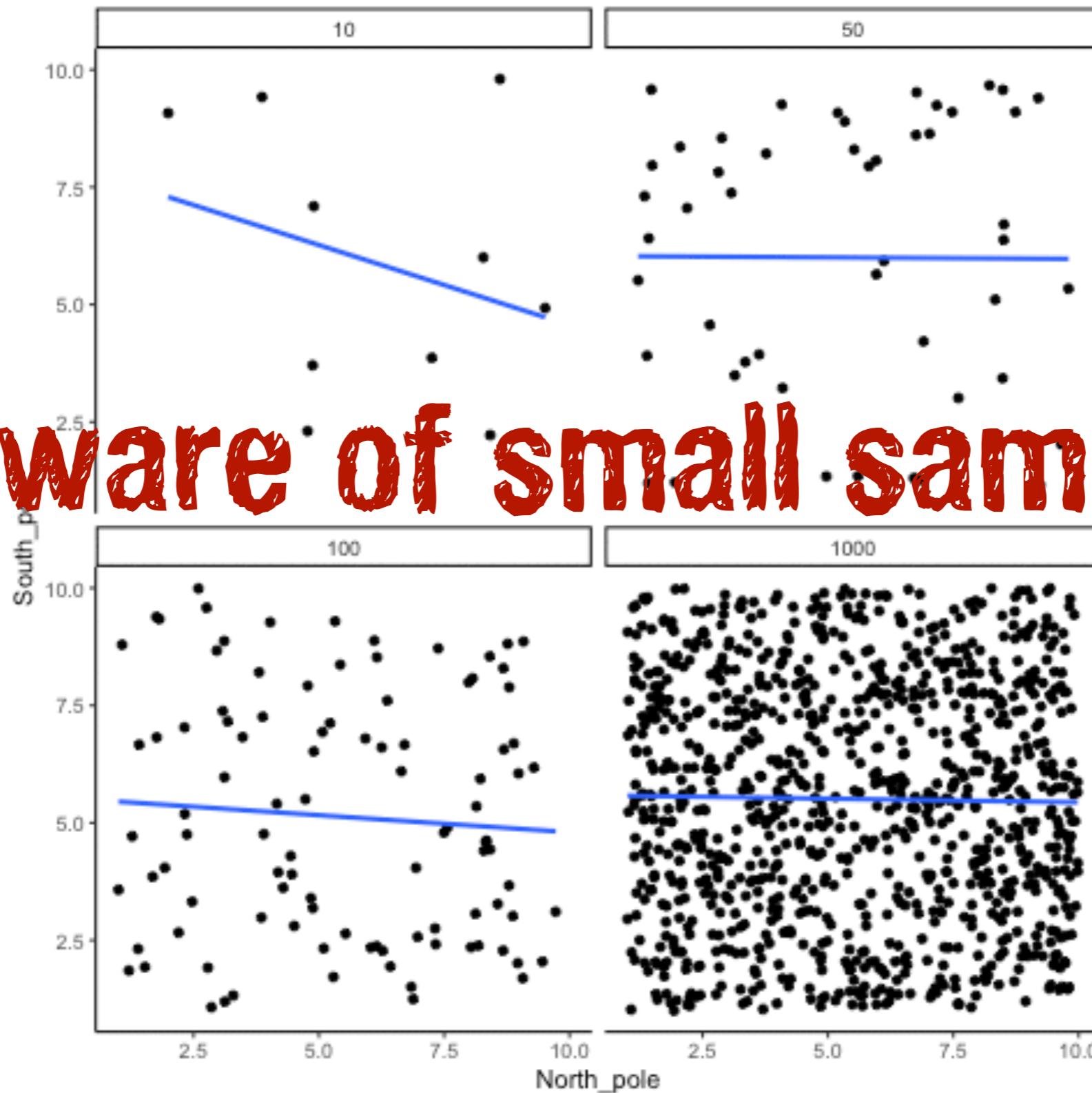


always visualize the data ...

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

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

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

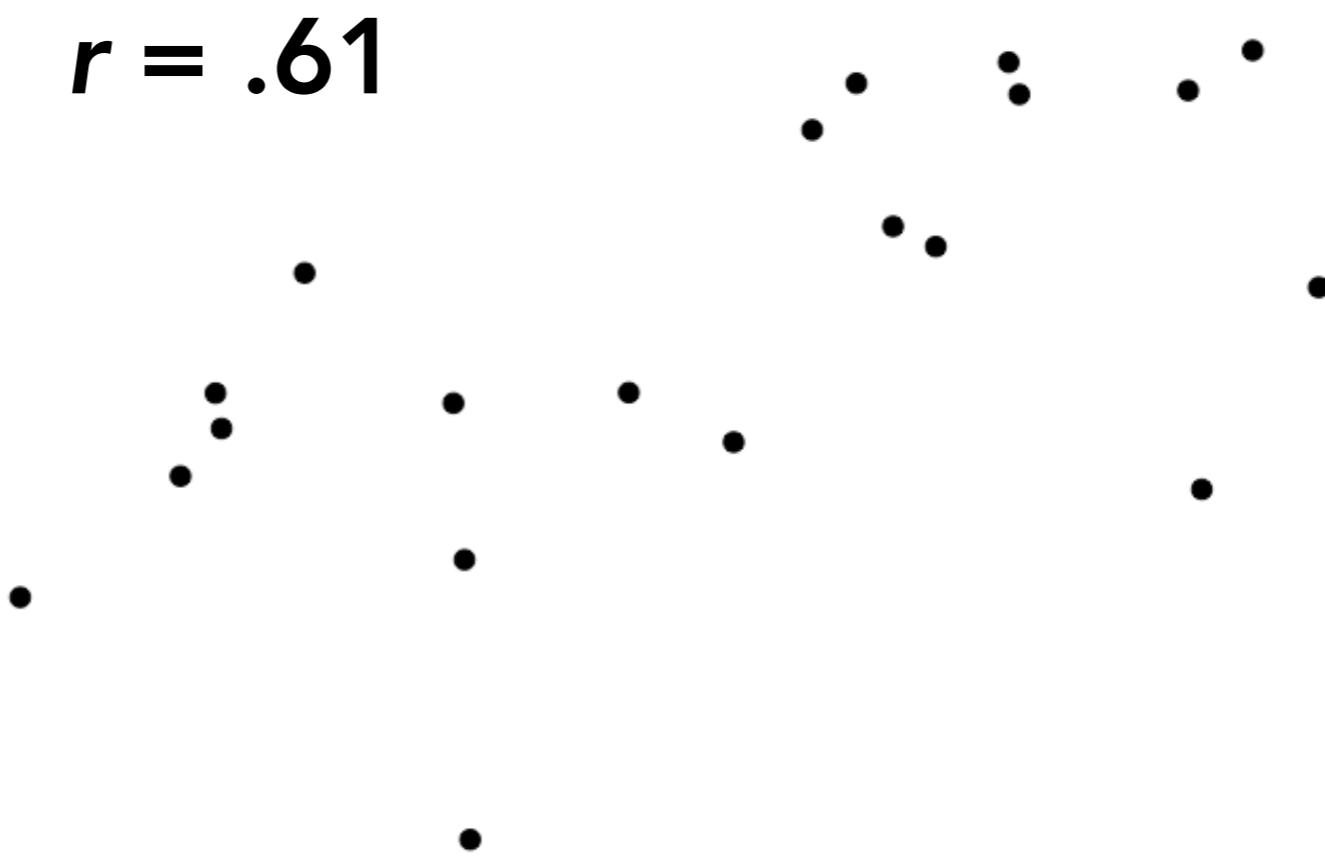


# in R

```
1 # data set
2 df.correlation = tibble(
3   x = runif(20, min = 0, max = 1),
4   y = x + rnorm(x, mean = 0.5, sd = 0.25)
5 )
6
7 # correlation
8 df.correlation %>%
9   summarize(r = cor(x, y))
```

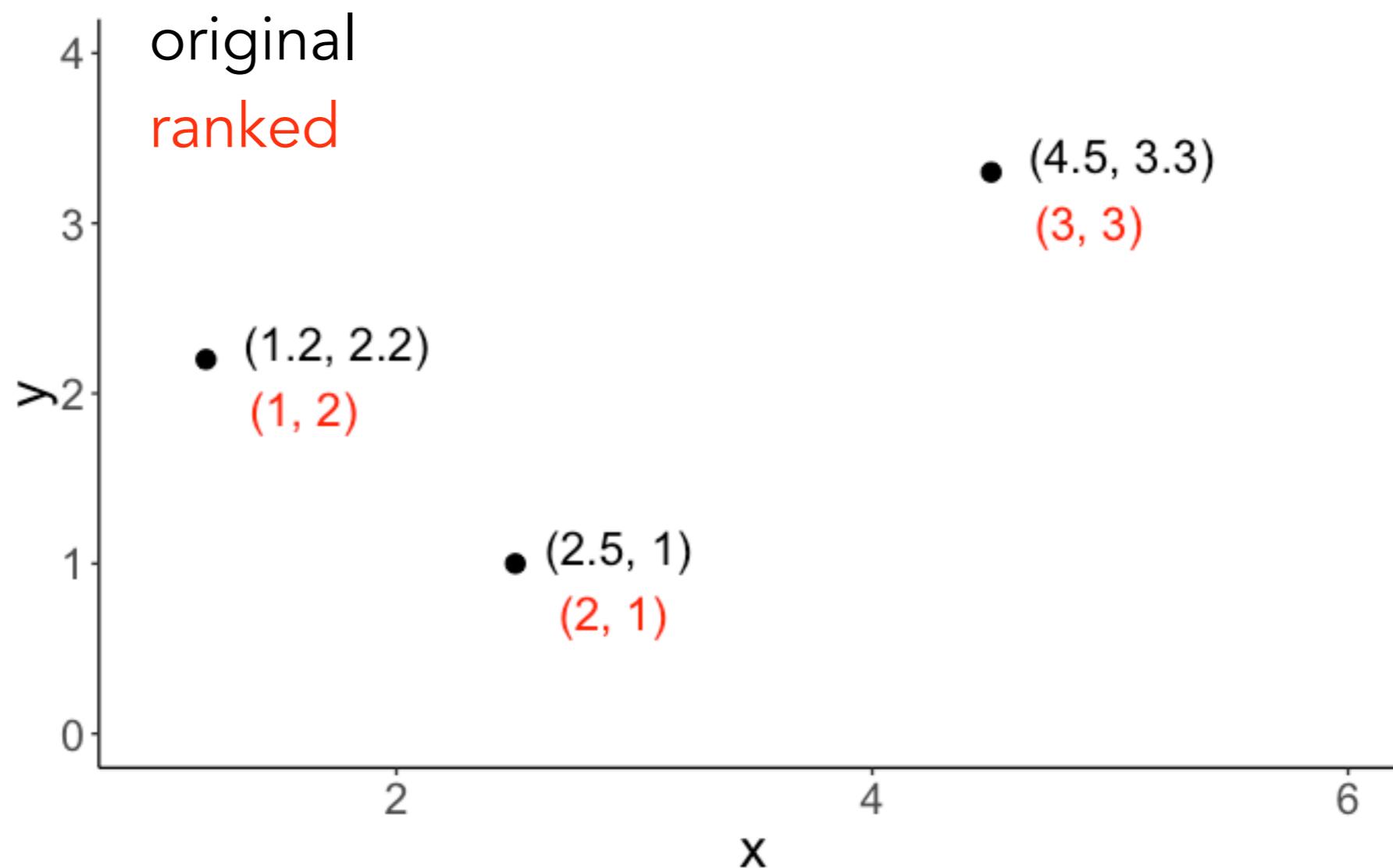
function to  
calculate  
correlation

$$r = .61$$



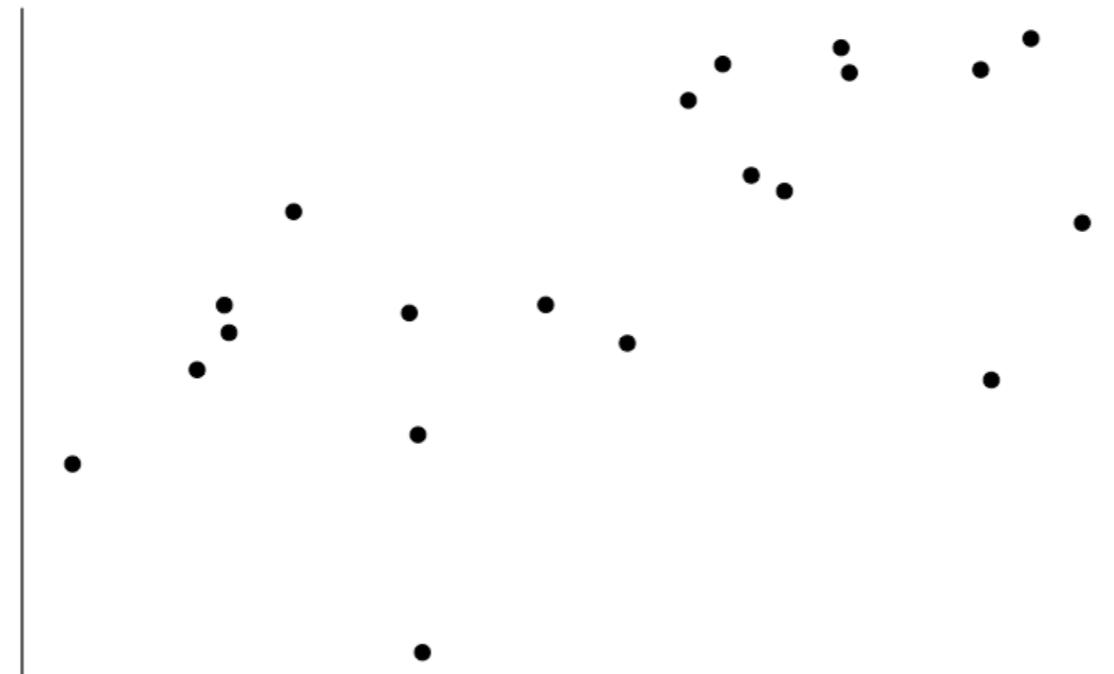
# Spearman rank order correlation

- 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



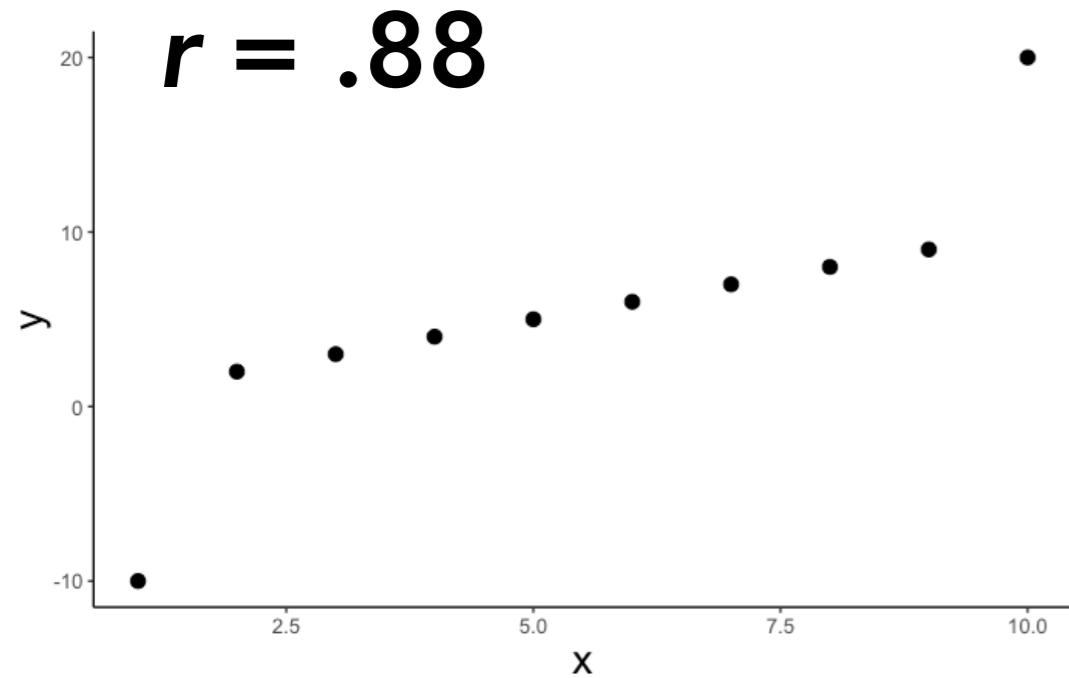
x	y	x_rank	y_rank
0.27	1.14	5	12
0.37	0.97	6	8
0.57	0.92	10	6
0.91	0.85	18	4
0.20	0.98	3	9
0.90	1.39	17	17
0.94	1.44	19	20
0.66	1.40	12	18
0.63	1.33	11	15
0.06	0.71	1	2

r	spearman	r_ranks
0.609	0.595	0.595

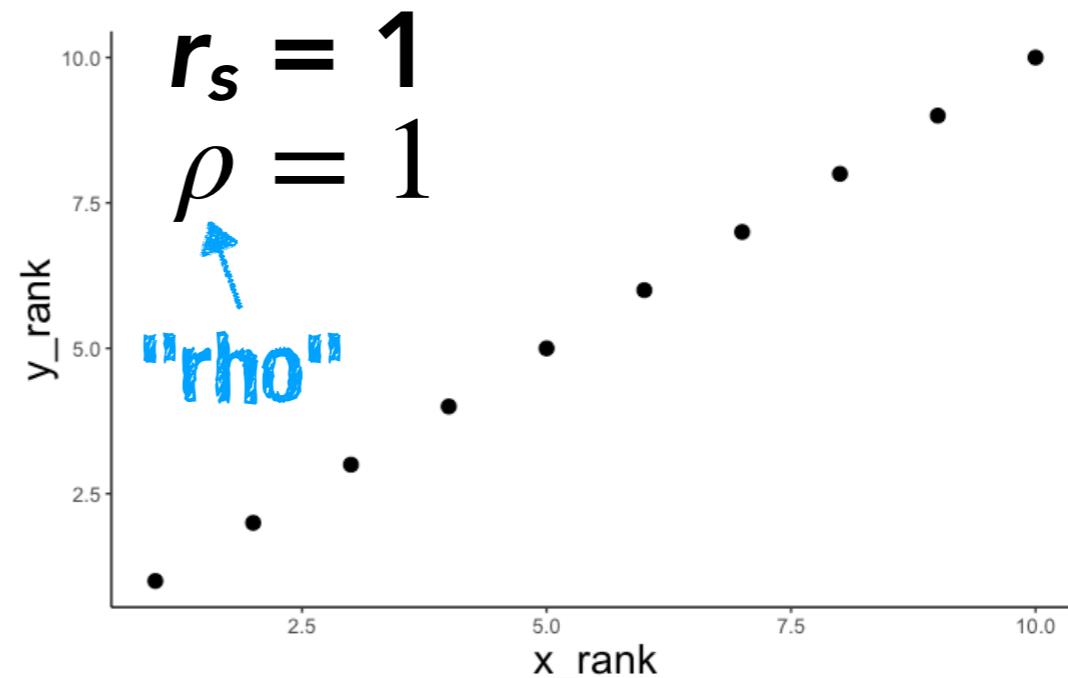
```
1 # correlation
2 df.spearman %>%
3   summarize(r = cor(x, y, method = "pearson"),
4             spearman = cor(x, y, method = "spearman"),
5             r_ranks = cor(x_rank, y_rank))
```

# Spearman rank order correlation

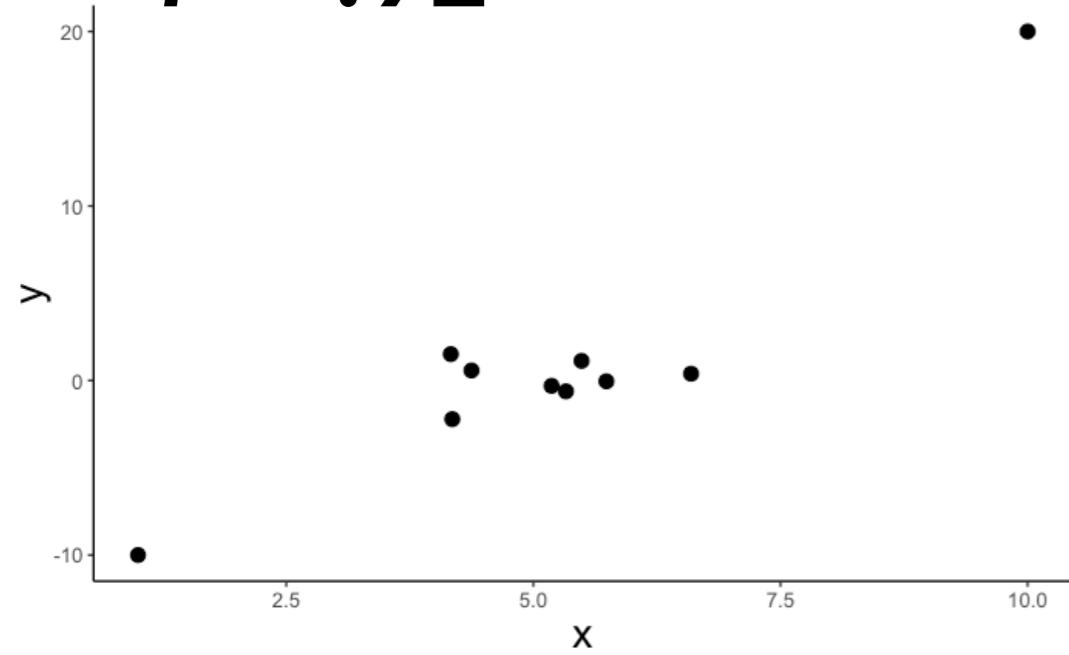
original



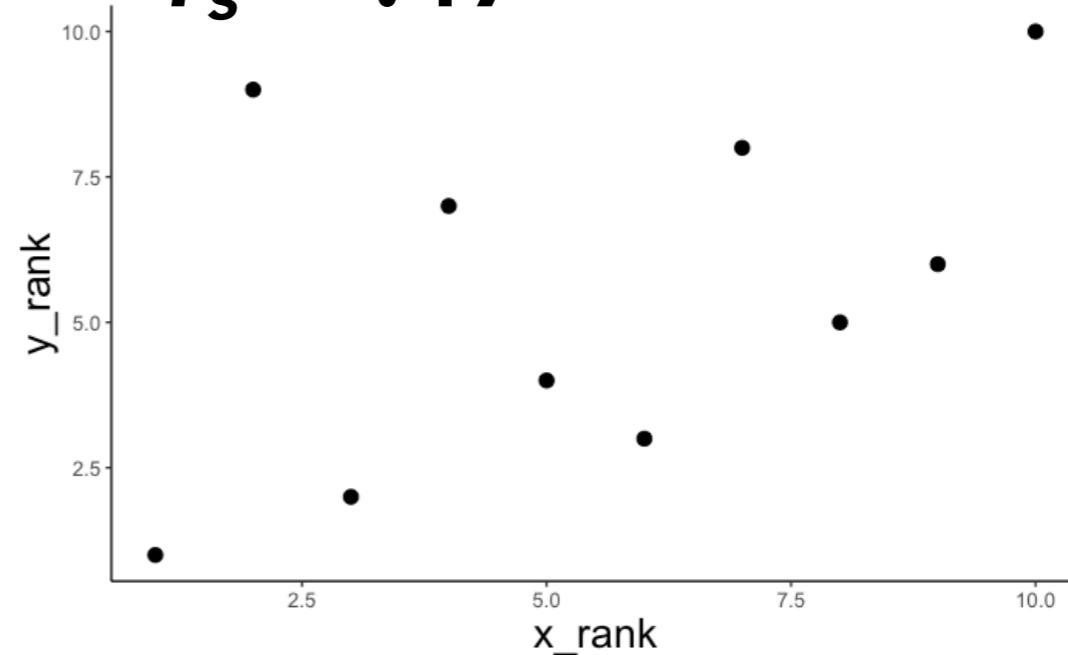
ranked



$r = .92$



$r_s = .47$



# 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

# Regression

# The conceptual tour

# Linear model: Simple regression

Data = Model + Error

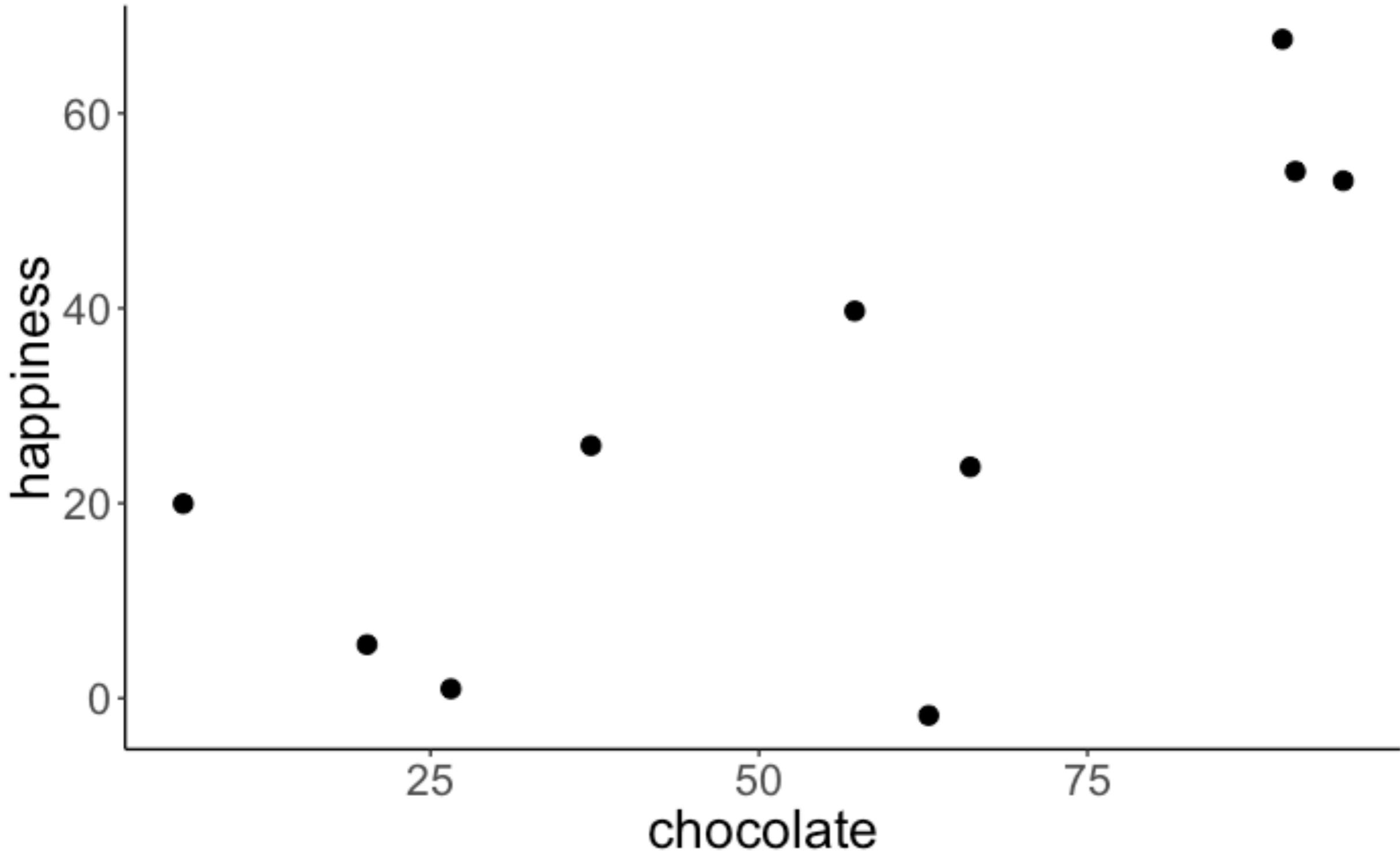
$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$$

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

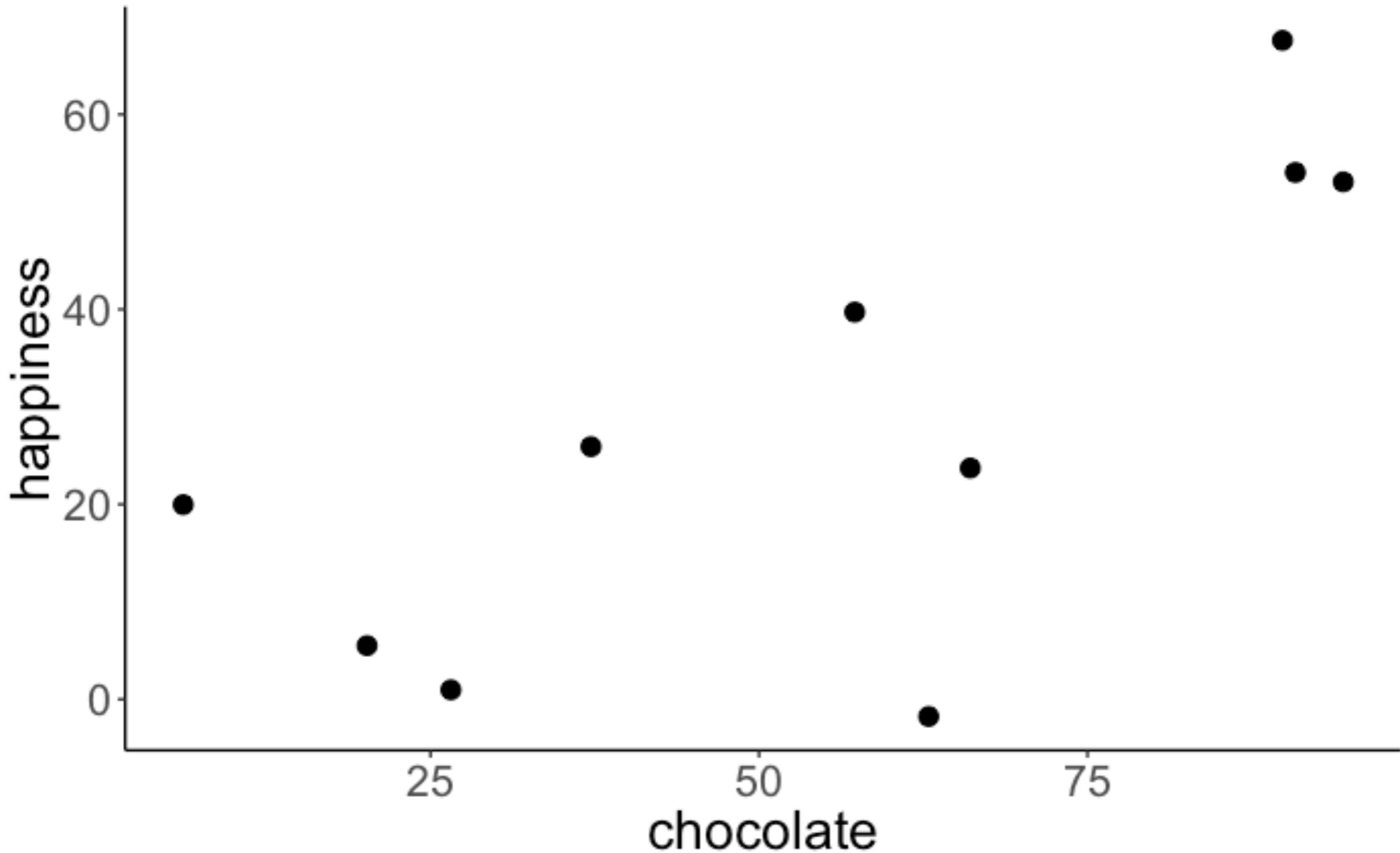


the model is a linear  
combination of predictors

# Does chocolate make us happy?



# 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)

$H_0$ : 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

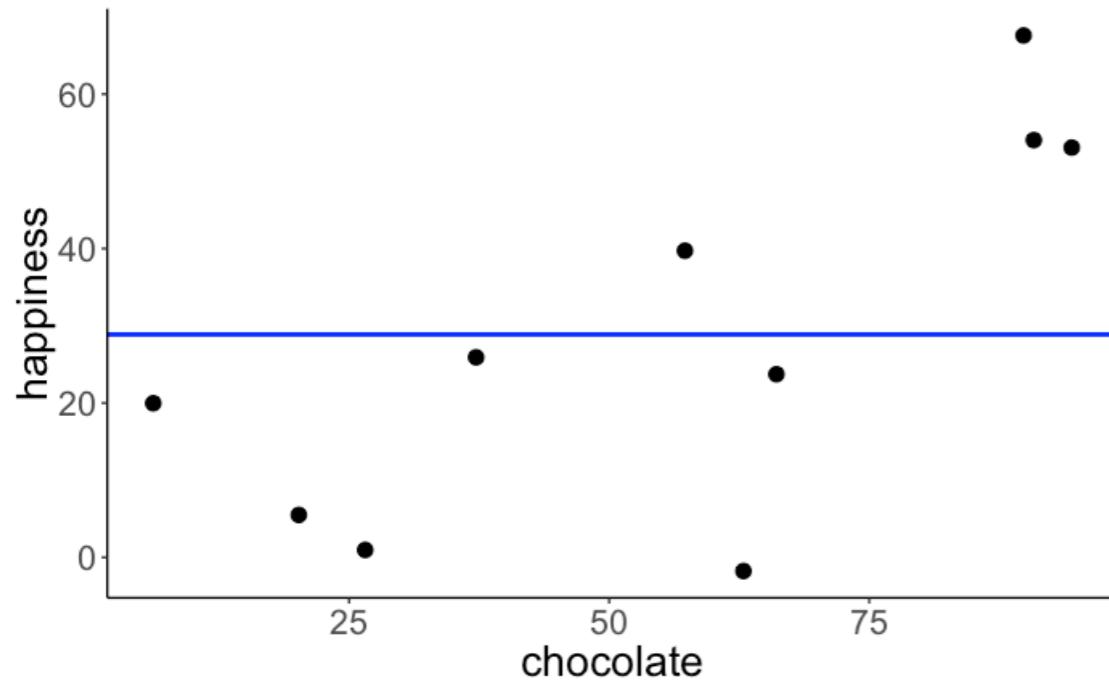
1. Define  $H_0$  as Model C (compact) and  $H_1$  as Model A (augmented)
- 2. Fit model parameters to the data**
3. Calculate the proportional reduction of error (PRE) in our sample
4. Decide whether the augmented model is **worth it** by comparing the observed PRE in our sample to the sampling distribution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

$$Y_i = 28.88 + e_i$$

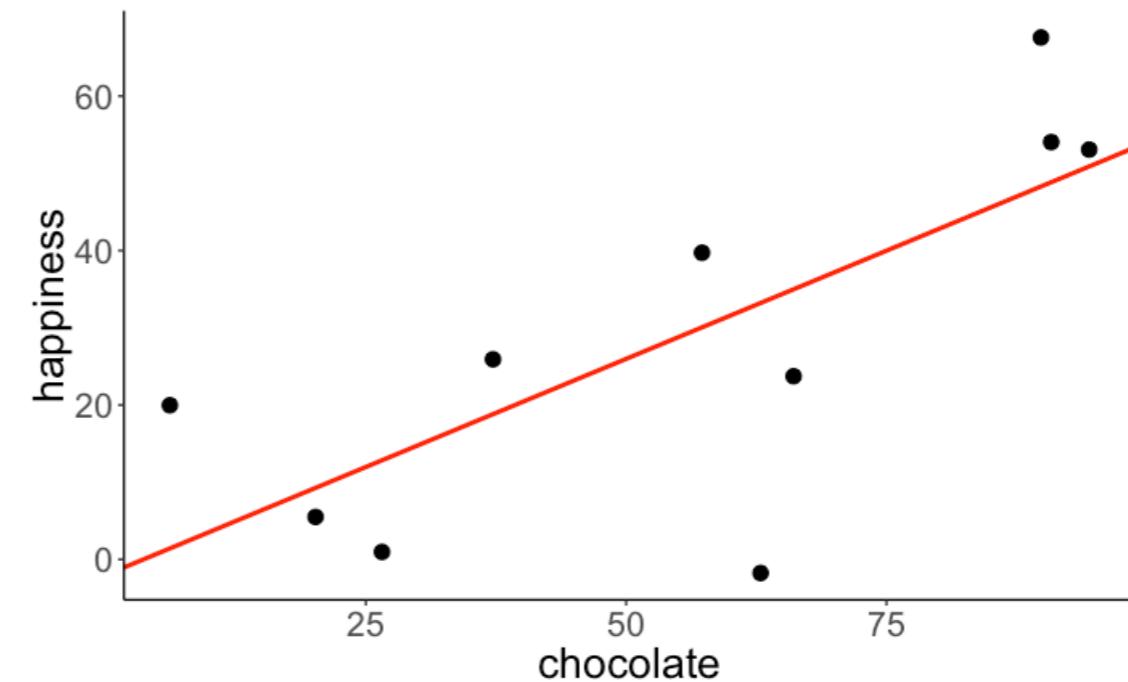
$H_1$ : Chocolate consumption and happiness are related.

### 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)

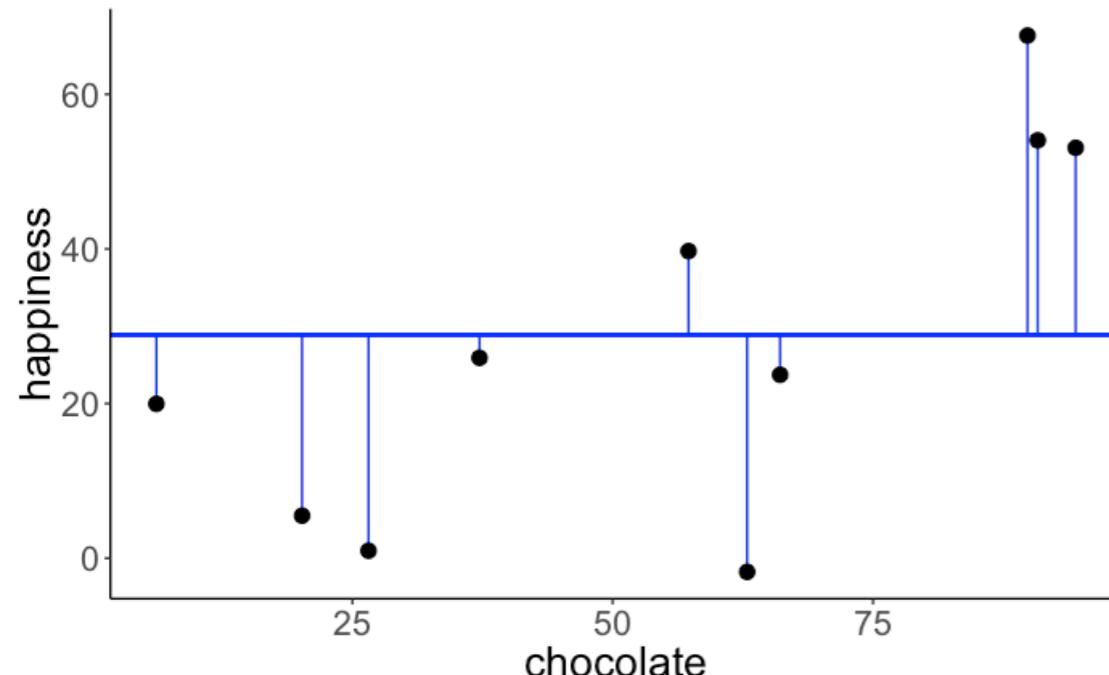
# Calculate PRE

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

Both models were fit to minimize the sum of squared errors

OLS = Ordinary **least squares** regression

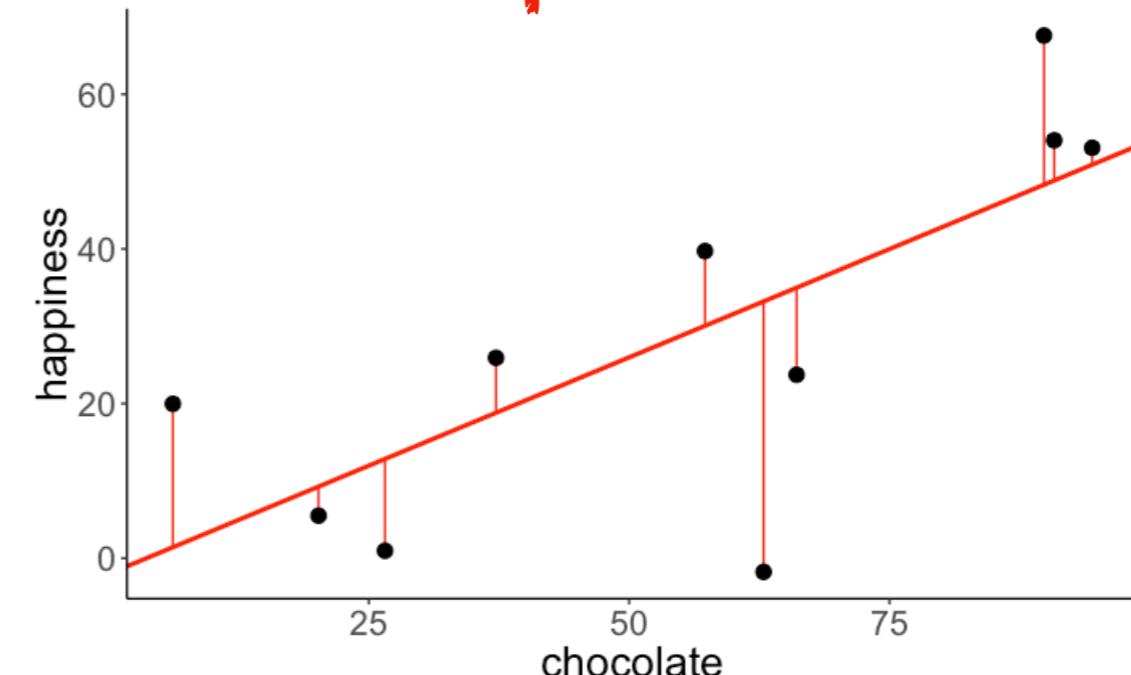
## Sum of squared errors



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

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

## Sum of squared errors



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

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  $\text{model}_C$  to create  $\text{model}_A$  but were not

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

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

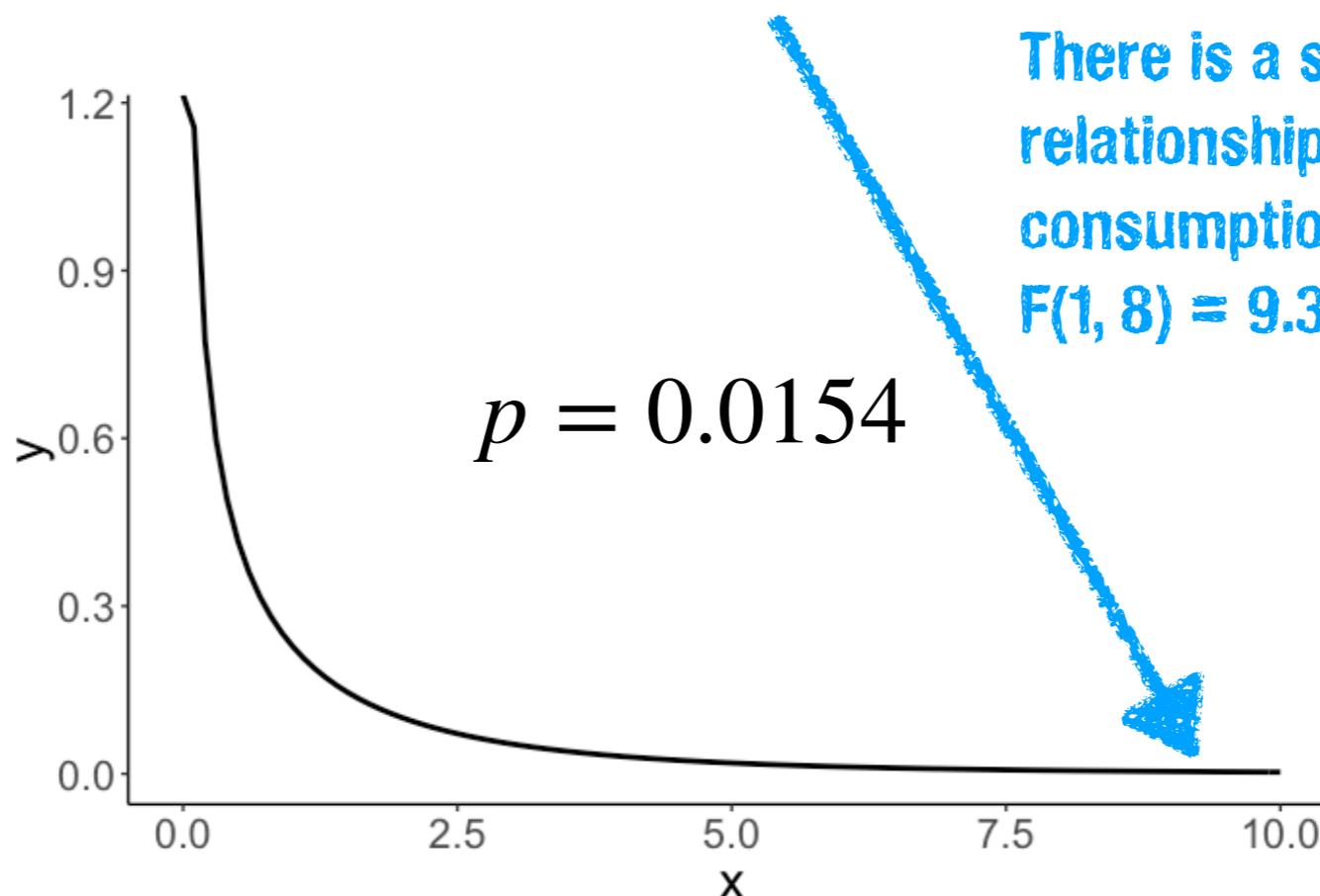
**number of observations  
vs. parameters in Model A**

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

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

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

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



# The R route

# Credit data set

## df.credit

index	income	limit	rating	cards	age	education	gender	student	married	ethnicity	balance
1	14.89	3606	283	2	34	11	Male	No	Yes	Caucasian	333
2	106.03	6645	483	3	82	15	Female	Yes	Yes	Asian	903
3	104.59	7075	514	4	71	11	Male	No	No	Asian	580
4	148.92	9504	681	3	36	11	Female	No	No	Asian	964
5	55.88	4897	357	2	68	16	Male	No	Yes	Caucasian	331
6	80.18	8047	569	4	77	10	Male	No	No	Caucasian	1151
7	21.00	3388	259	2	37	12	Female	No	No	African American	203
8	71.41	7114	512	2	87	9	Male	No	No	Asian	872
9	15.12	3300	266	5	66	13	Female	No	No	Caucasian	279
10	71.06	6819	491	3	41	19	Female	Yes	Yes	African American	1350

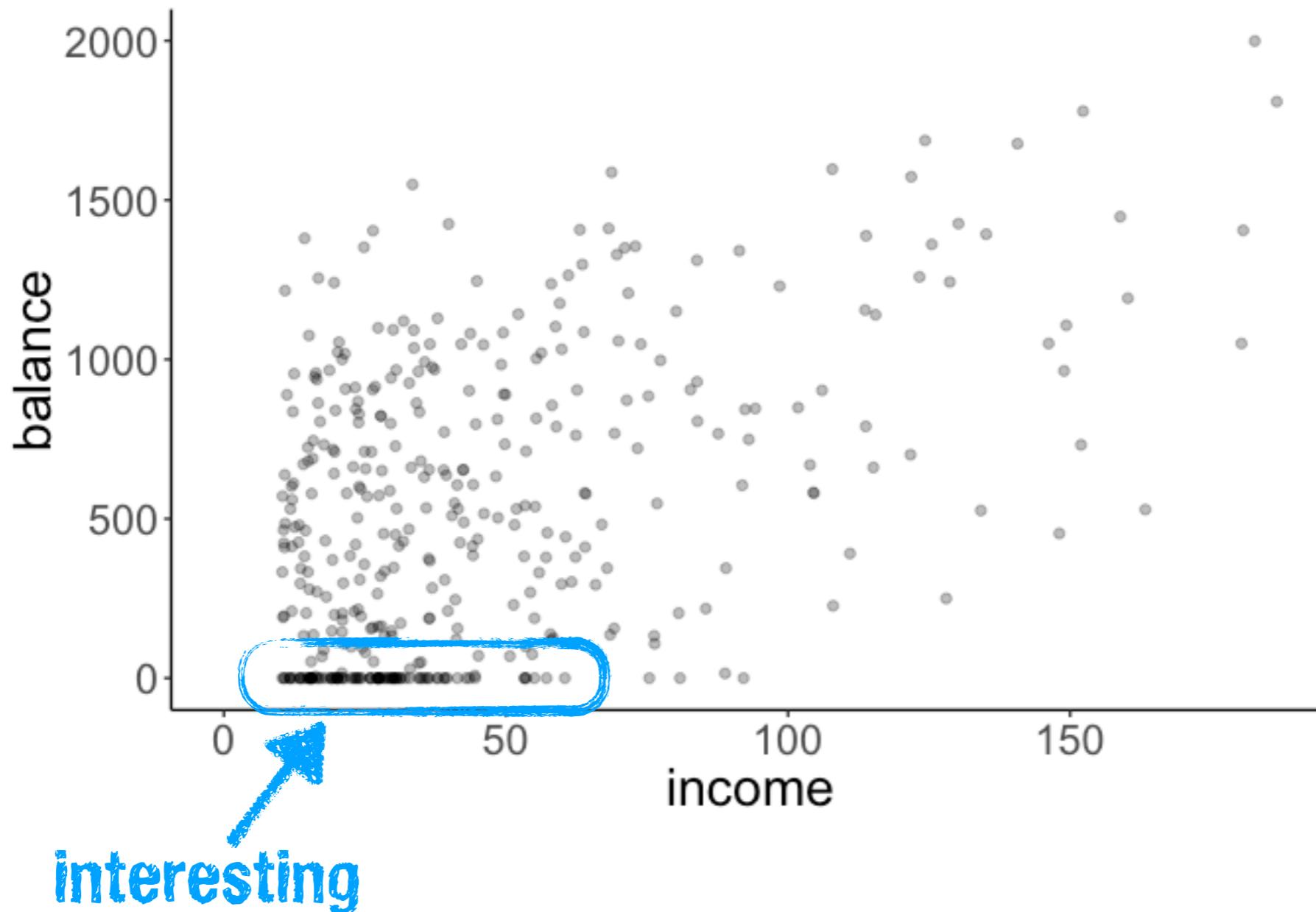
**nrow(df.credit) = 400**

**Is there a relationship between income  
and the average credit card debt?**

variable	description
income	in thousand dollars
limit	credit limit
rating	credit rating
cards	number of credit cards
age	in years
education	years of education
gender	male or female
student	student or not
married	married or not
ethnicity	African American, Asian, Caucasian
balance	average credit card debt in dollars

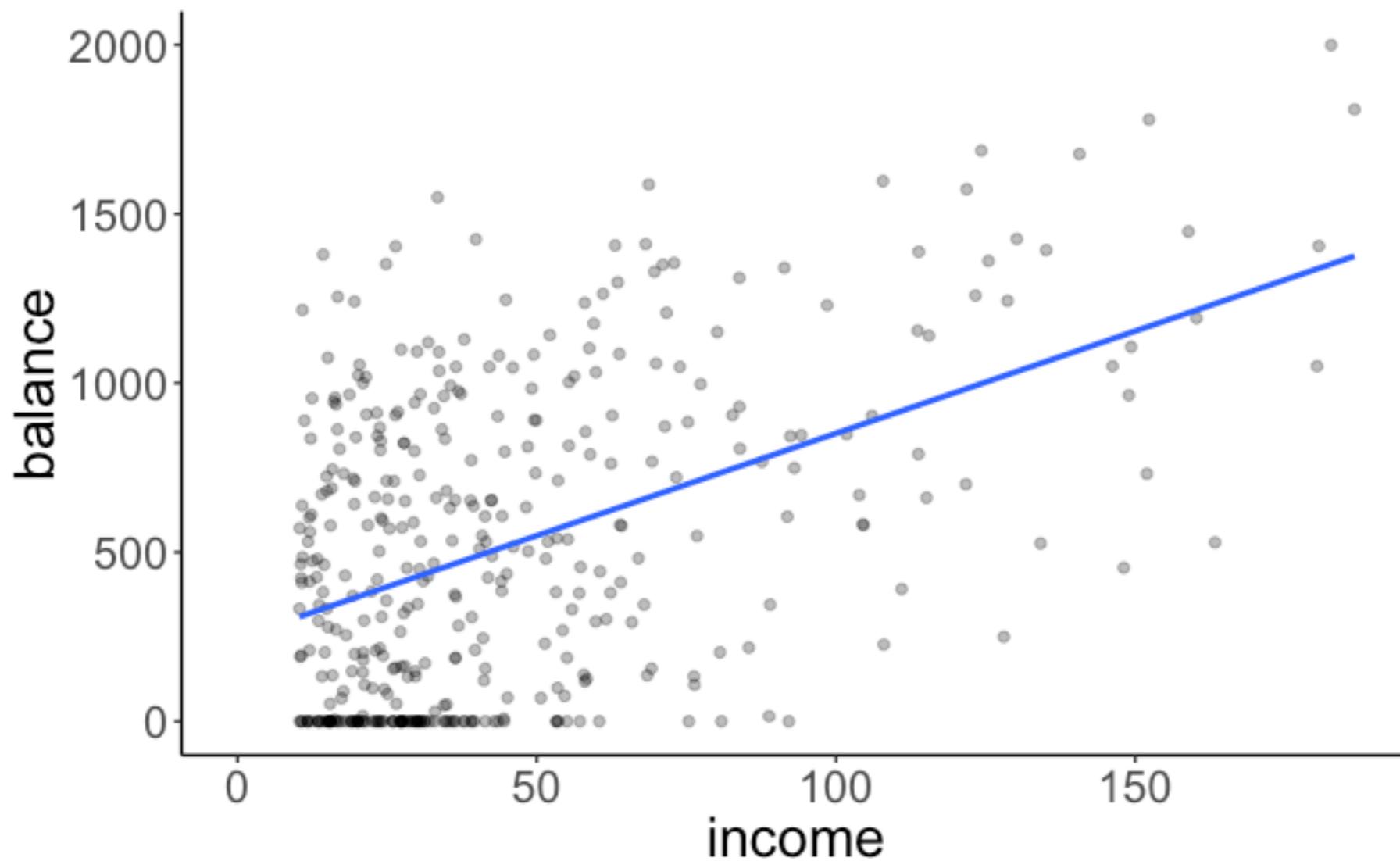
# Always plot the data first ...

```
1 ggplot(data = df.credit,  
2         mapping = aes(x = income,  
3                             y = balance)) +  
4     geom_point(alpha = 0.3)
```



# Always plot the data first ...

```
1 ggplot(data = df.credit,  
2         mapping = aes(x = income,  
3                             y = balance)) +  
4     geom_point(alpha = 0.3) +  
5     geom_smooth(method = "lm", se = F)
```



# Linear model

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

$$\text{balance}_i = b_0 + b_1 \cdot \text{income}_i + e_i$$

```
fit = lm(balance ~ 1 + income, data = df.credit)
```

**outcome**      **intercept**      **predictor**      **data**

(doesn't need to be  
specified explicitly)

# **lm()**

```
fit = lm(balance ~ 1 + income, data = df.credit)
```

```
print(fit)
```

```
Call:  
lm(formula = balance ~ 1 + income, data = df.credit)  
  
Coefficients:  
(Intercept)           income  
      246.515            6.048
```

  
**parameter estimates**

# Interpreting regression parameters

Coefficients:

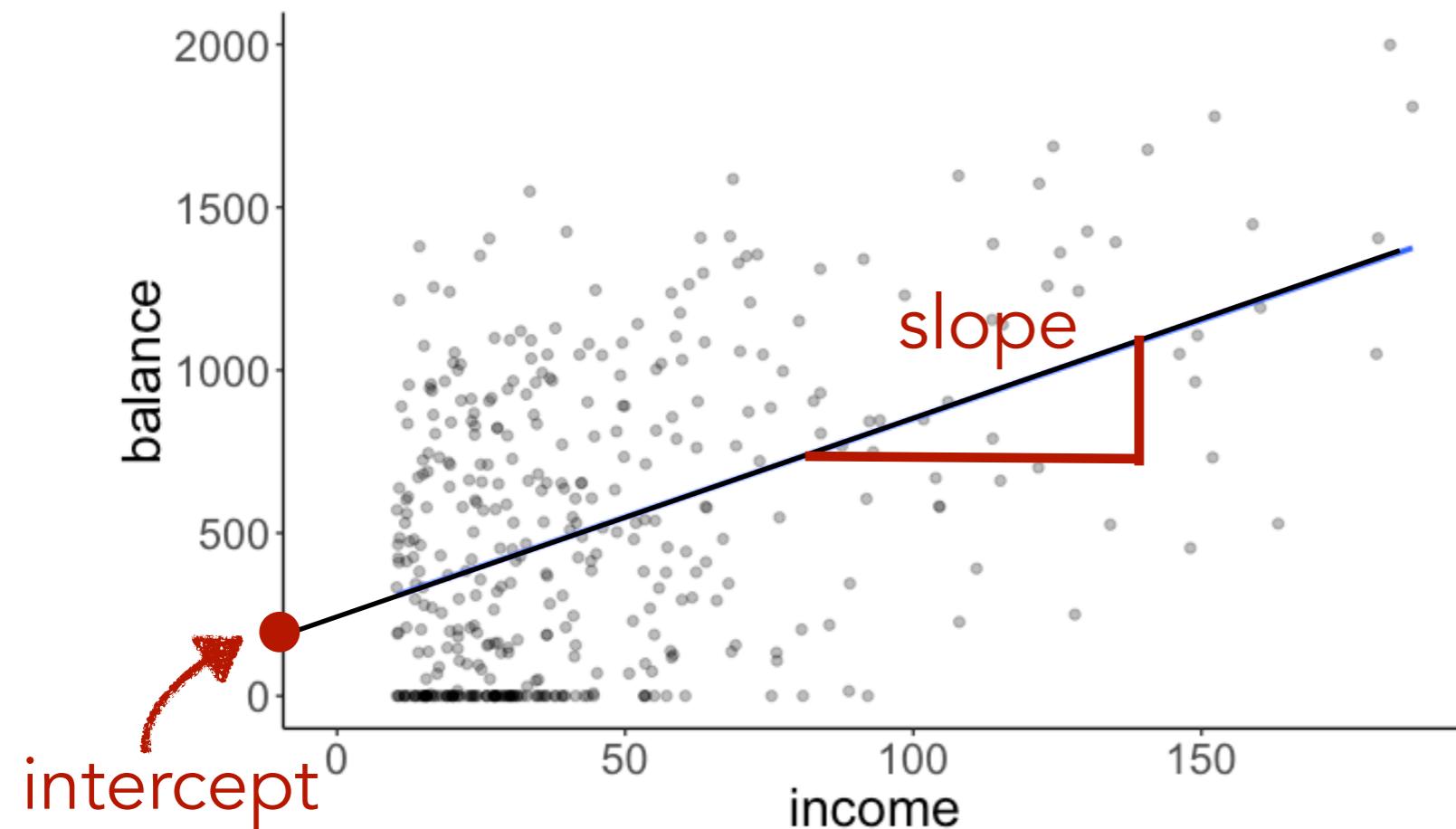
(Intercept) 246.515

income 6.048

variable	description
income	in thousand dollars
balance	average credit card debt in dollars

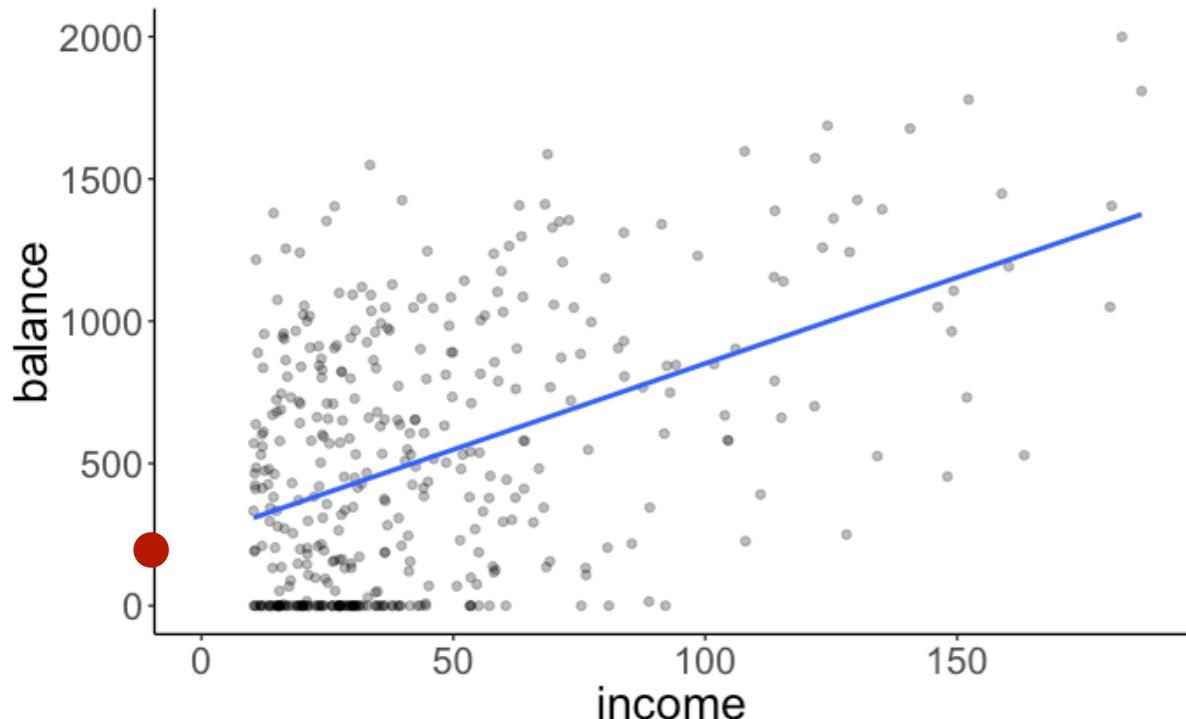
$$\text{balance}_i = b_0 + b_1 \cdot \text{income}_i + e_i$$

$$\text{balance}_i = 246.515 + 6.048 \cdot \text{income}_i + e_i$$



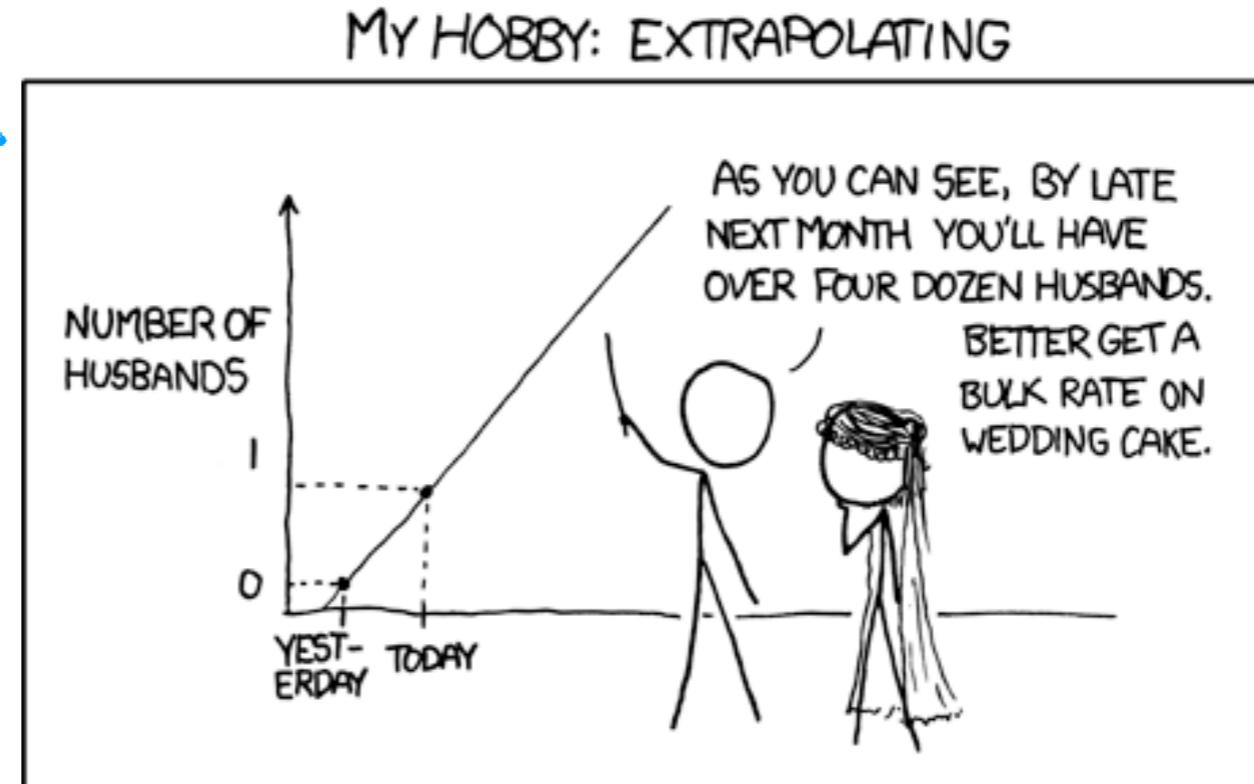
For each additional thousand dollars income, a person's average credit card debt increases by \$6.05.

# Be careful about extrapolating predictions



- intercept is often outside the range of predictor values
- sometimes doesn't make sense (e.g. age = 0, height = 0, ...)

comic from  
slide 1



# Centering the predictor

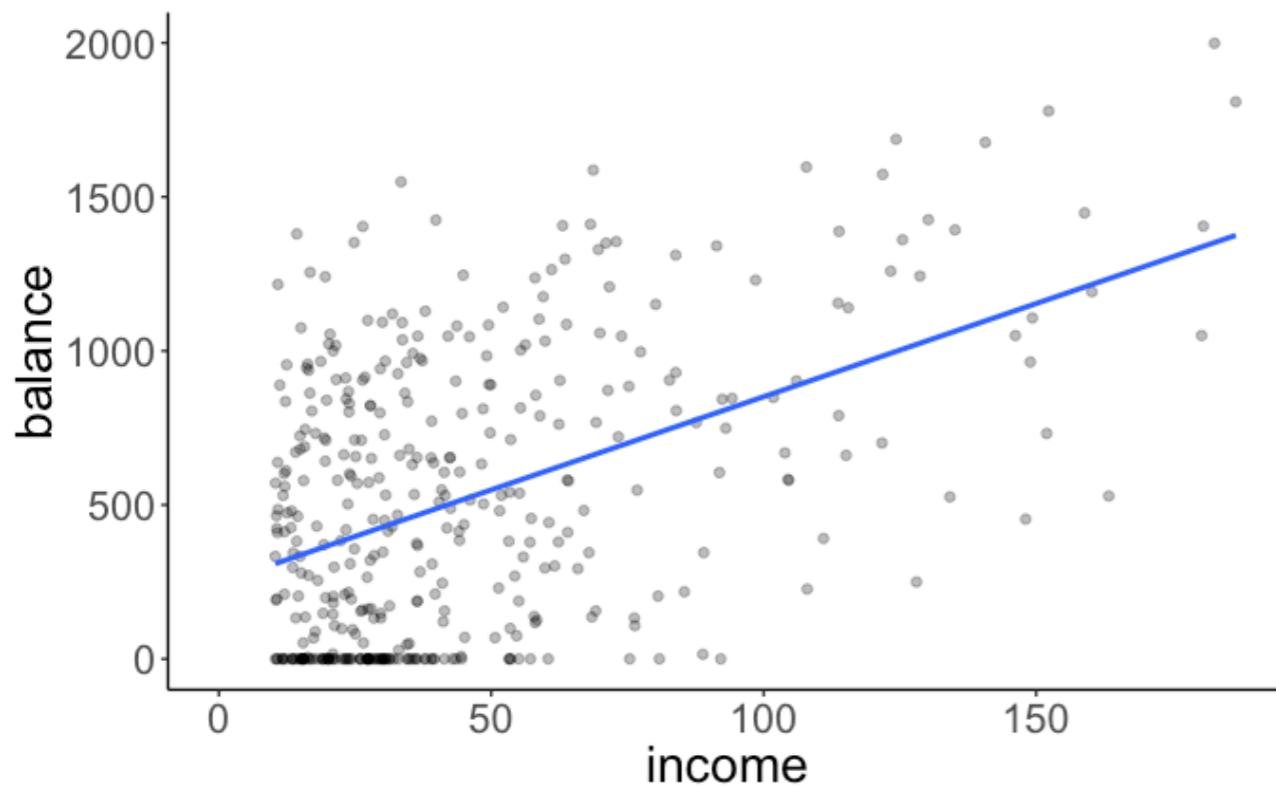
```
1 df.credit %>%
2   mutate(income_centered = income - mean(income)) %>%
3   select(balance, income, income_centered)
```

balance	income	income_centered
333	14.89	-30.33
903	106.03	60.81
580	104.59	59.37
964	148.92	103.71
331	55.88	10.66
1151	80.18	34.96
203	21.00	-24.22
872	71.41	26.19
279	15.12	-30.09
1350	71.06	25.84

# Centering predictors

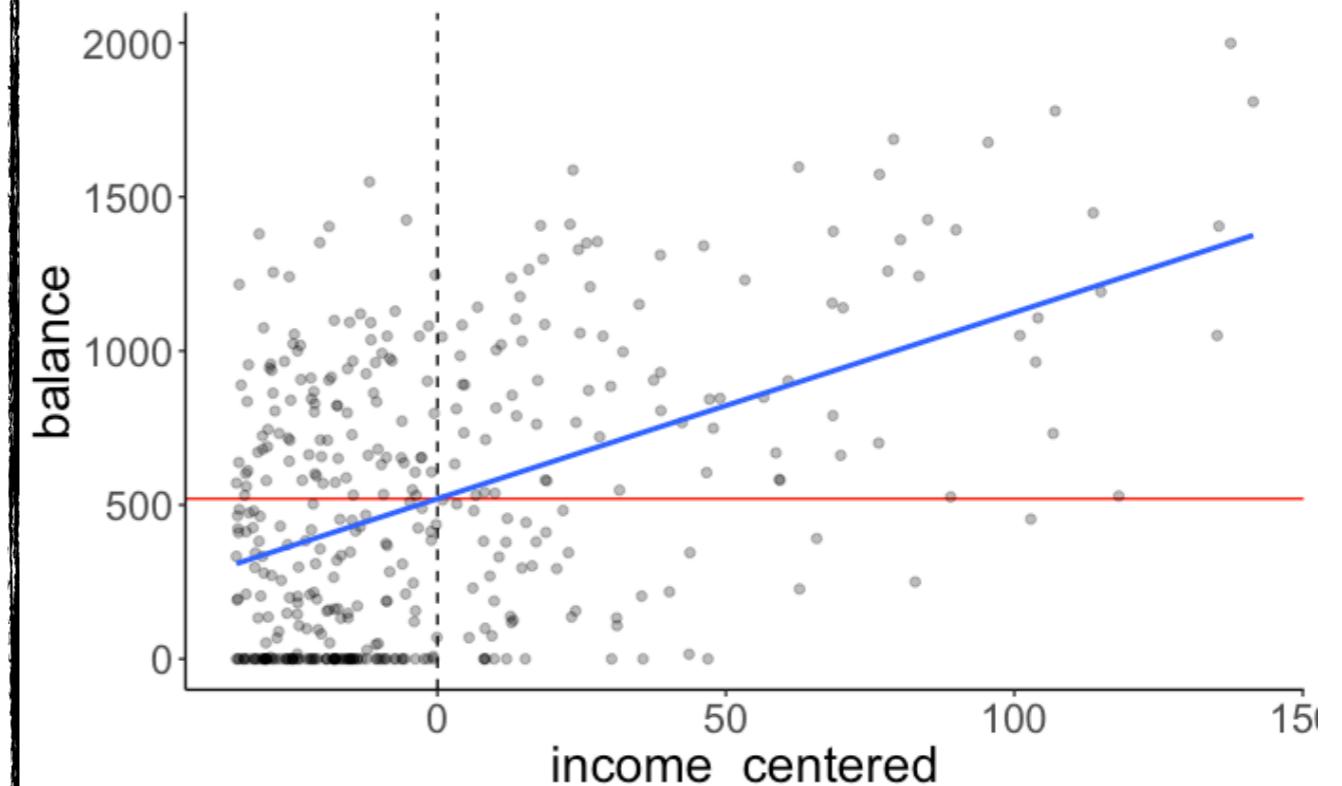
uncentered

$$\text{balance}_i = 246.515 + 6.048 \cdot \text{income}_i + e_i$$



centered

$$\text{balance}_i = 520.015 + 6.048 \cdot \text{income\_centered}_i + e_i$$



intercept = predicted value if  
income is 0

intercept = predicted value if  
income is  
**mean** (income)

# summary( )

```
1 lm(balance ~ 1 + income, data = df.credit) %>%
2   summary()
```

```
Call:
lm(formula = balance ~ 1 + income, data = df.credit)

Residuals:
    Min      1Q  Median      3Q     Max 
-803.64 -348.99 -54.42  331.75 1100.25 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 246.5148   33.1993   7.425  6.9e-13 ***  
income       6.0484    0.5794  10.440 < 2e-16 ***  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 407.9 on 398 degrees of freedom
Multiple R-squared:  0.215, Adjusted R-squared:  0.213 
F-statistic: 109 on 1 and 398 DF,  p-value: < 2.2e-16
```

```
library ("broom")
```



helps with tidying up  
model objects in R

**augment()** adds columns to the original data such as predictions, residuals and cluster assignments

**tidy()** summarizes a model's statistical findings such as coefficients of a regression

**glance()** provides a one-row summary of model-level statistics

# summary()

Residuals:					
Min	1Q	Median	3Q	Max	
-803.64	-348.99	-54.42	331.75	1100.25	

fit %>%

**augment()**

balance	income	.fitted	.se.fit	.resid	.hat	.sigma	.cooksdi	.std.resid
333	14.89	336.58	26.92	-3.58	0.00	408.38	0.00	-0.01
903	106.03	887.79	40.71	15.21	0.01	408.38	0.00	0.04
580	104.59	879.13	39.99	-299.13	0.01	408.10	0.00	-0.74
964	148.92	1147.26	63.45	-183.26	0.02	408.27	0.00	-0.45
331	55.88	584.51	21.31	-253.51	0.00	408.18	0.00	-0.62
1151	80.18	731.47	28.74	419.53	0.00	407.83	0.00	1.03
203	21.00	373.51	24.76	-170.51	0.00	408.29	0.00	-0.42
872	71.41	678.42	25.42	193.58	0.00	408.26	0.00	0.48
279	15.12	338.00	26.83	-59.00	0.00	408.37	0.00	-0.14
1350	71.06	676.32	25.30	673.68	0.00	406.97	0.01	1.65

# summary( )

Residuals:

Min	1Q	Median	3Q	Max
-803.64	-348.99	-54.42	331.75	1100.25

fit %>%

augment( )

balance	income	.fitted	.resid
333	14.89	336.58	-3.58
903	106.03	887.79	15.21
580	104.59	879.13	-299.13
964	148.92	1147.26	-183.26
331	55.88	584.51	-253.51
1151	80.18	731.47	419.53
203	21.00	373.51	-170.51
872	71.41	678.42	193.58
279	15.12	338.00	-59.00
1350	71.06	676.32	673.68

1 fit %>%

2 augment() %>%

3 clean\_names() %>%

4 summarize(

5 min = min(resid),

6 first\_quantile = quantile(resid, 0.25),

7 median = median(resid),

8 third\_quantile = quantile(resid, 0.75),

9 max = max(resid)

10 )

min	first_quantile	median	third_quantile	max
-803.64	-348.99	-54.42	331.75	1100.25

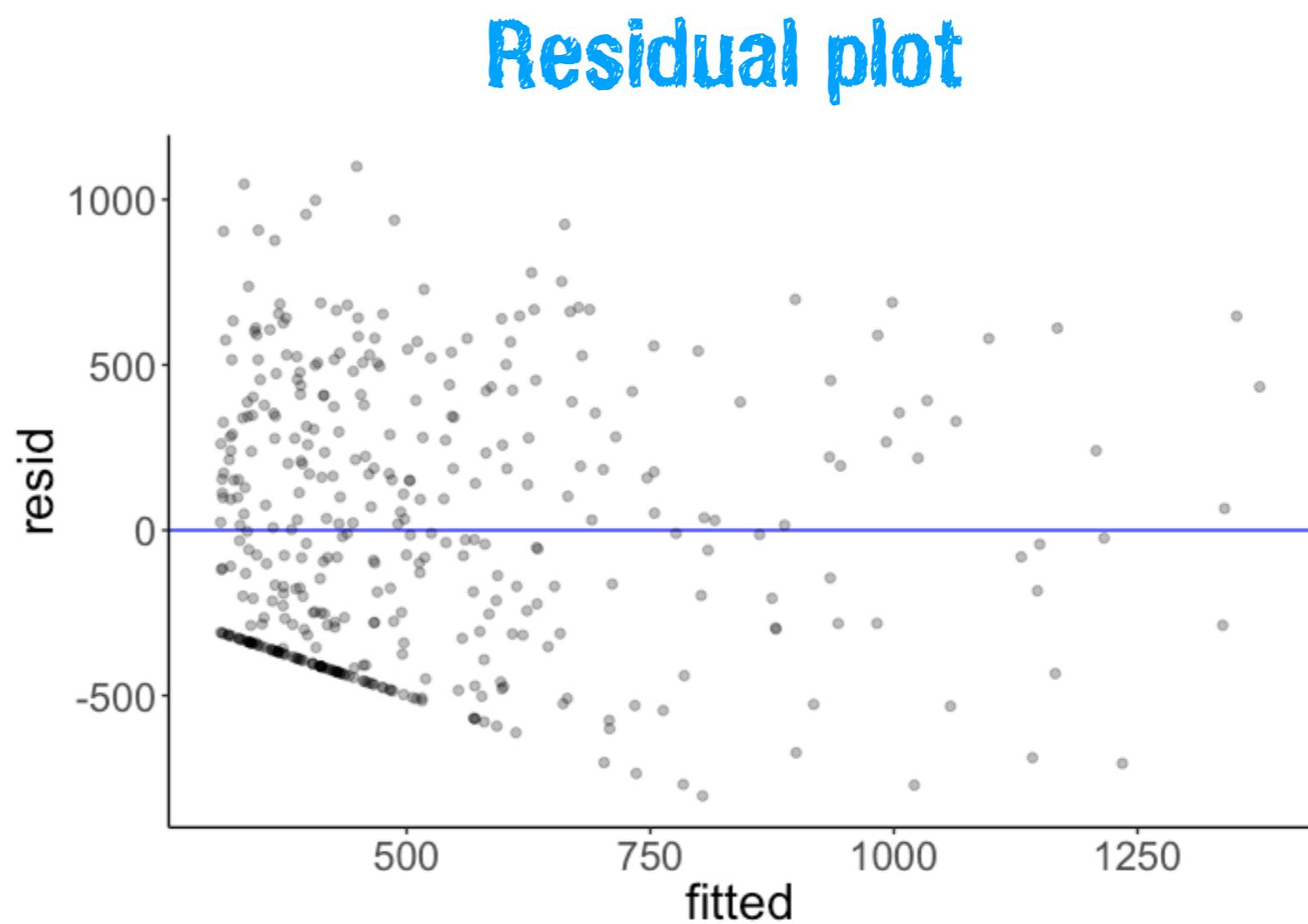
# summary()

Residuals:					
Min	1Q	Median	3Q	Max	
-803.64	-348.99	-54.42	331.75	1100.25	

fit %>%

augment()

balance	income	.fitted	.resid
333	14.89	336.58	-3.58
903	106.03	887.79	15.21
580	104.59	879.13	-299.13
964	148.92	1147.26	-183.26
331	55.88	584.51	-253.51
1151	80.18	731.47	419.53
203	21.00	373.51	-170.51
872	71.41	678.42	193.58
279	15.12	338.00	-59.00
1350	71.06	676.32	673.68



# summary()

```
1 lm(balance ~ 1 + income, data = df.credit) %>%
2   summary()
```

```
Coefficients:
              Estimate Std. Error t value Pr(>|t|)    
(Intercept) 246.5148    33.1993   7.425  6.9e-13 ***
income       6.0484     0.5794  10.440 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
1 fit %>%
2   tidy(conf.int = TRUE)
```

term	estimate	std.error	statistic	p.value	conf.low	conf.high
(Intercept)	246.51	33.20	7.43	0	181.25	311.78
income	6.05	0.58	10.44	0	4.91	7.19

# summary()

```
1 lm(balance ~ 1 + income, data = df.credit) %>%
2   summary()
```

```
Residual standard error: 407.9 on 398 degrees of freedom
Multiple R-squared:  0.215, Adjusted R-squared:  0.213
F-statistic: 109 on 1 and 398 DF,  p-value: < 2.2e-16
```

```
1 fit %>%
2   glance()
```

r.squared	adj.r.squared	sigma	statistic	p.value	df	logLik	AIC	BIC	deviance	df.residual
0.21	0.21	407.86	108.99	0	2	-2970.95	5947.89	5959.87	66208745	398

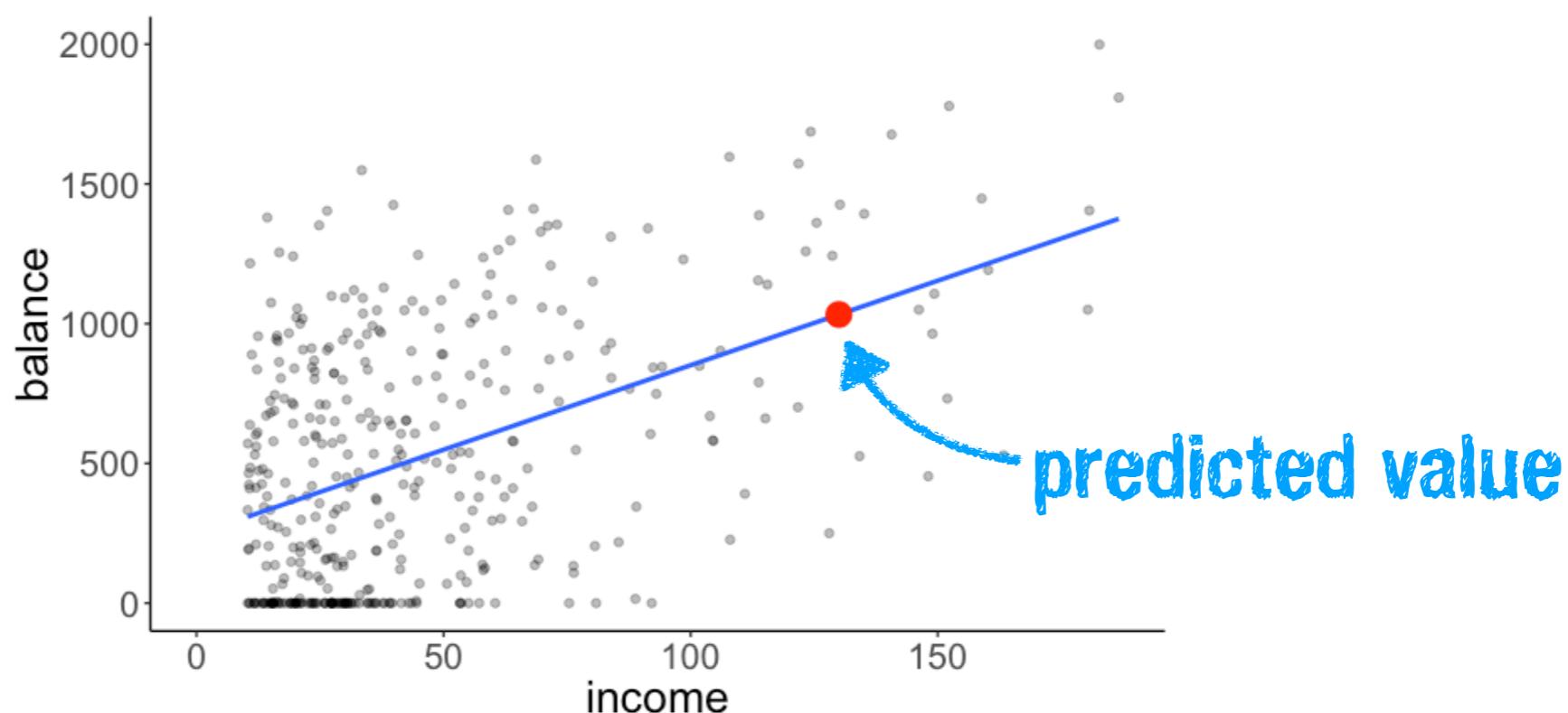
# Making predictions

```
fit = lm(balance ~ 1 + income, data = df.credit)
```

$$\text{balance}_i = 246.515 + 6.048 \cdot \text{income}_i + e_i$$

```
augment(fit, newdata = tibble(income = 130))
```

$$\begin{aligned}\widehat{\text{balance}} &= 246.515 + 6.048 \cdot 130 \\ &= 1032.755\end{aligned}$$



# Hypothesis test

## Compact Model

$$\text{balance}_i = \beta_0 + \epsilon_i$$

```
fit_c = lm(formula = balance ~ 1,  
           data = df.credit)
```

## Augmented Model

$$\text{balance}_i = \beta_0 + \beta_1 \text{outcome}_i + \epsilon_i$$

```
fit_a = lm(formula = balance ~ 1 + income,  
           data = df.credit)
```

**anova**(fit\_c, fit\_a)

Analysis of Variance Table

Model 1: balance ~ 1

Model 2: balance ~ 1 + income

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	399	84339912			
2	398	66208745	1	18131167 108.99 < 2.2e-16 ***	

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

# 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)

```
fit_c = lm(formula = balance ~ 1,  
           data = df.credit)  
  
fit_a = lm(formula = balance ~ 1 + income,  
           data = df.credit)
```

```
anova(fit_c, fit_a)
```

# Hypothesis test

**anova (fit\_c, fit\_a)**

Analysis of Variance Table

Model 1: balance ~ 1

Model 2: balance ~ 1 + income

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	399	84339912				
2	398	66208745	1	18131167	108.99 < 2.2e-16	***
---						

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

$$PRE = 1 - \frac{66208745}{84339912} \approx 0.215$$

The augmented model reduces the error by 21.5%.

```
lm(balance ~ 1 + income, data = df.credit) %>%  
  summary()
```

$R^2$

```
Residual standard error: 407.9 on 398 degrees of freedom  
Multiple R-squared: 0.215, Adjusted R-squared: 0.213  
F-statistic: 109 on 1 and 398 DF, p-value: < 2.2e-16
```

# Hypothesis test

- in the case of a simple regression PRE (proportion of reduced error) is identical to  $R^2$  (variance explained)
- and  $R^2$  is directly related to the correlation coefficient  $r$

```
cor(df.credit$balance,  
df.credit$income)
```

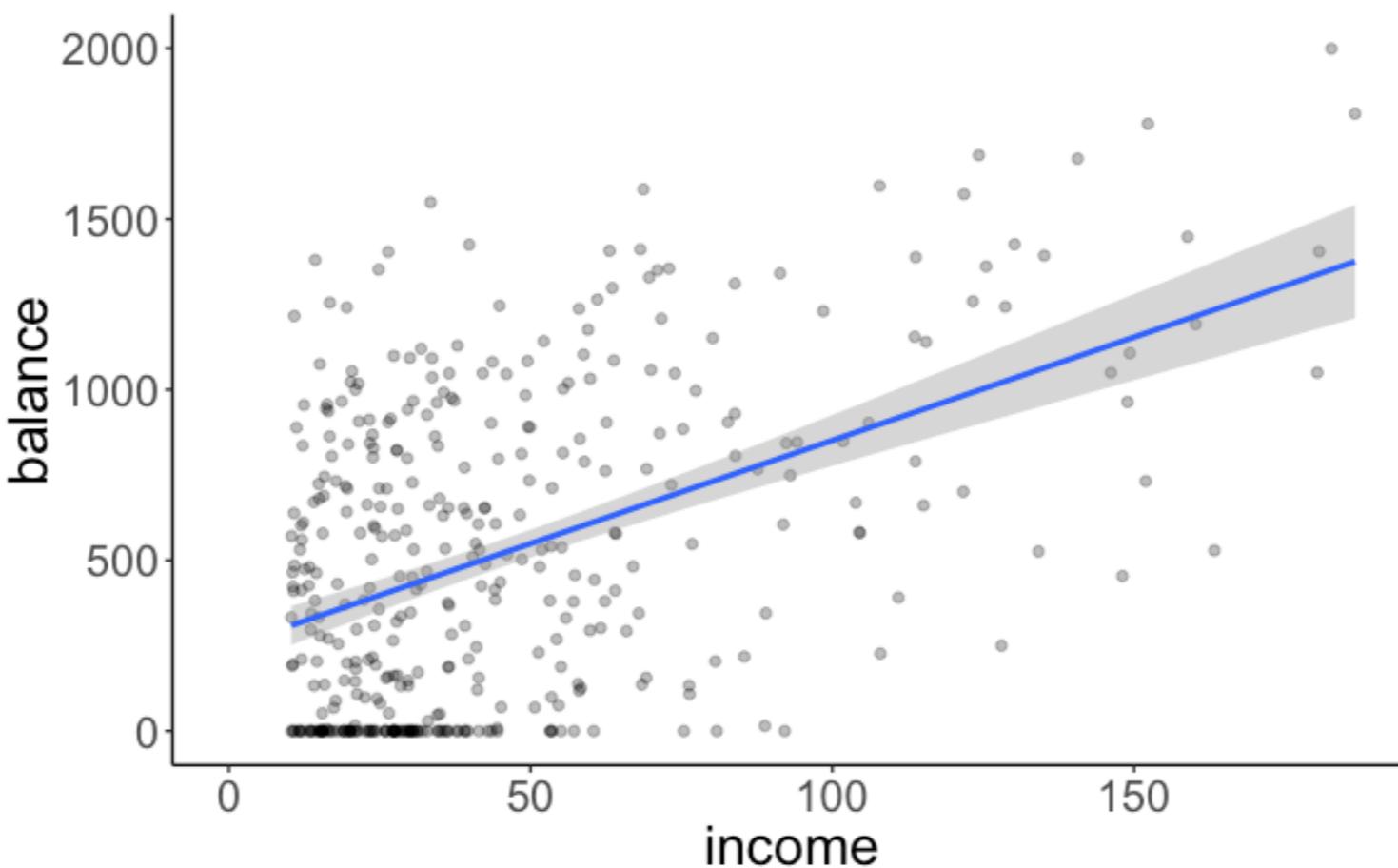
$$R^2 = 0.215$$

$$r = .463$$



effect size measure

# Reporting the results



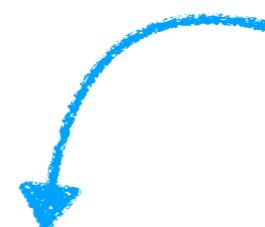
There is a significant relationship between a person's income and the average balance on their credit cards  
 $F(1, 389) = 108.99, p < .001, r = .463$ .

With each additional \$1000 of income, the average balance increases by \$6.05 [4.91, 7.19] (95% CI).

# The general procedure

Test whether the intercept is significantly different from 0

tells R to remove the intercept



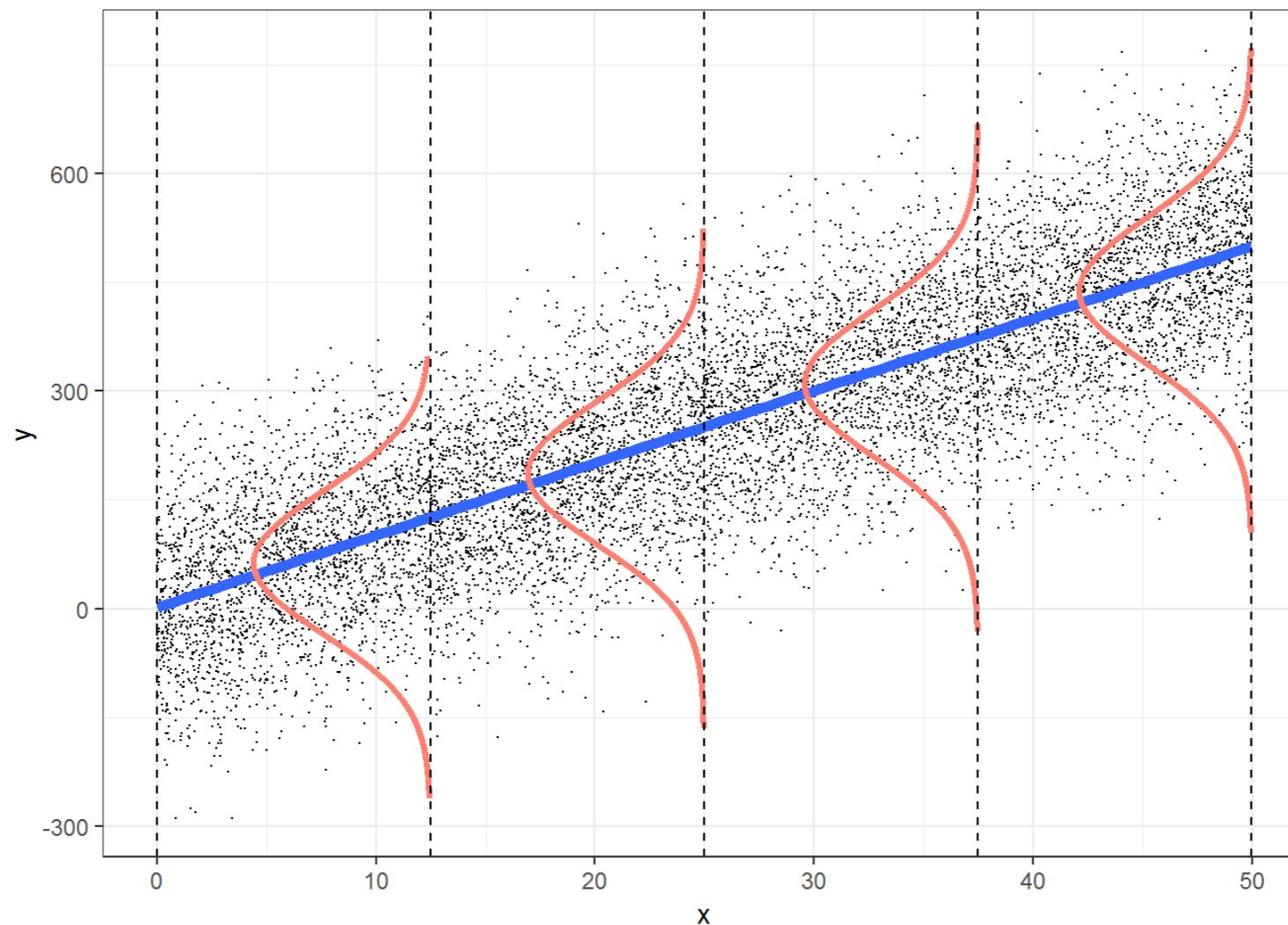
```
fit_c = lm(formula = balance ~ -1 + income,  
           data = df.credit)
```

```
fit_a = lm(formula = balance ~ 1 + income,  
           data = df.credit)
```

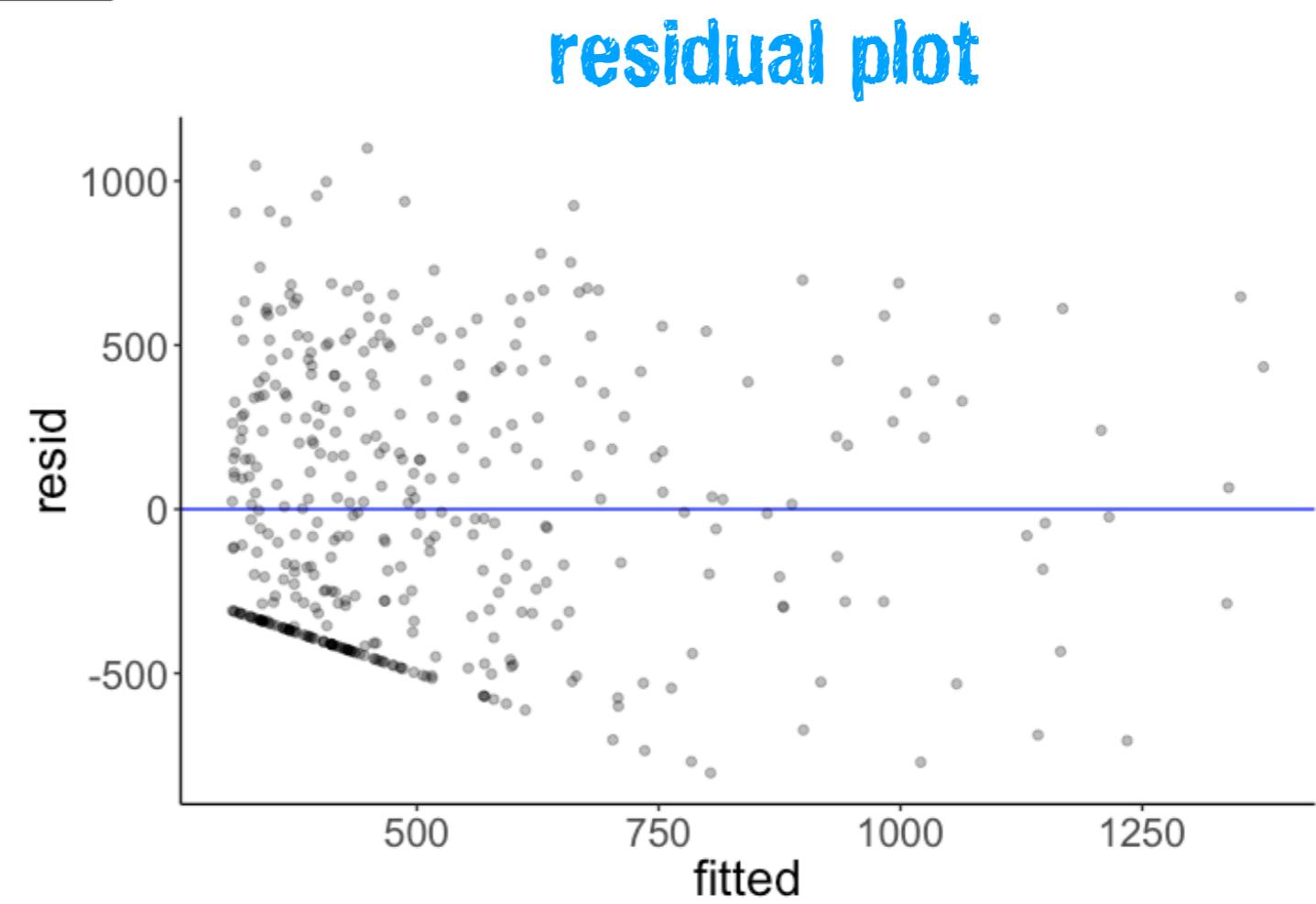
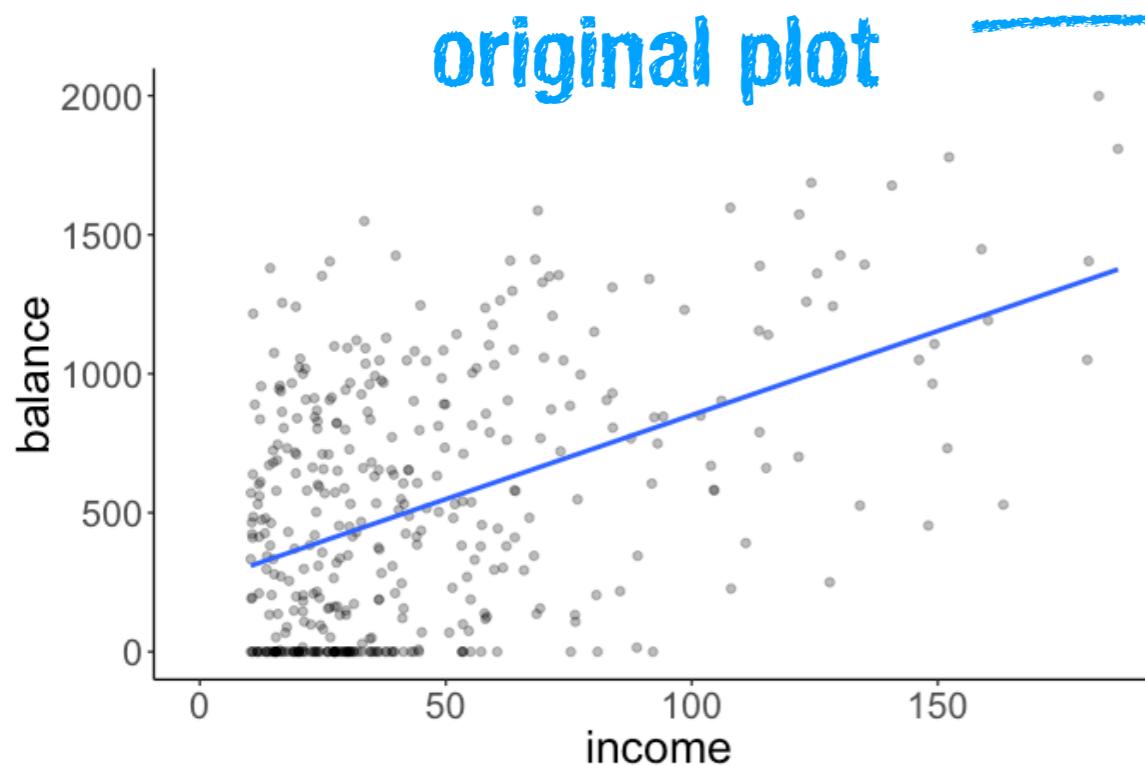
```
anova(fit_c, fit_a)
```

# Model assumptions

- independent observations
- $Y$  is continuous
- errors are normally distributed
- errors have constant variance
- error terms are uncorrelated

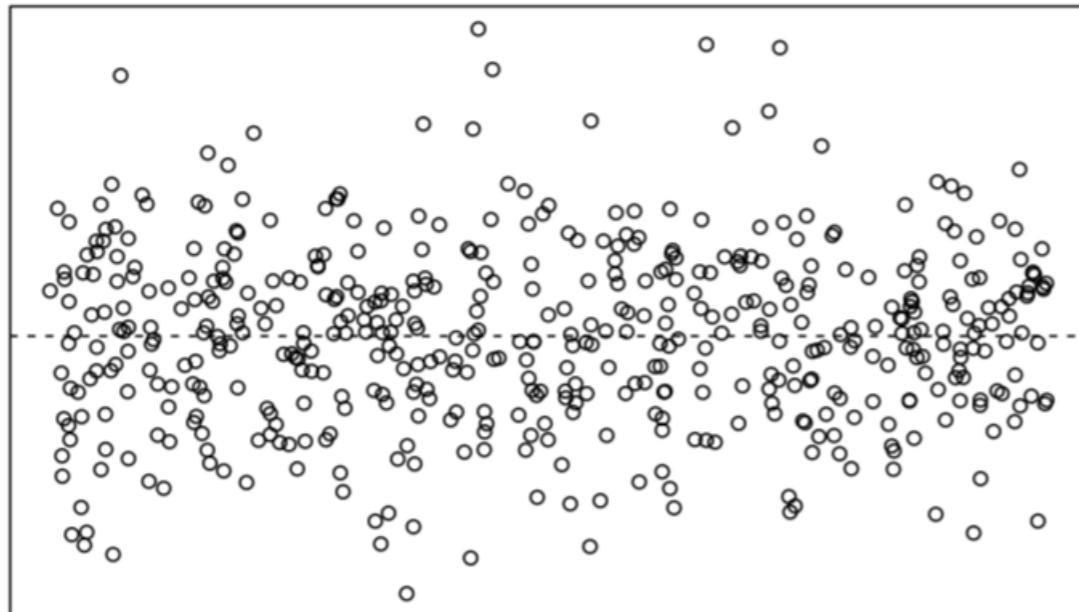


# Model assumptions

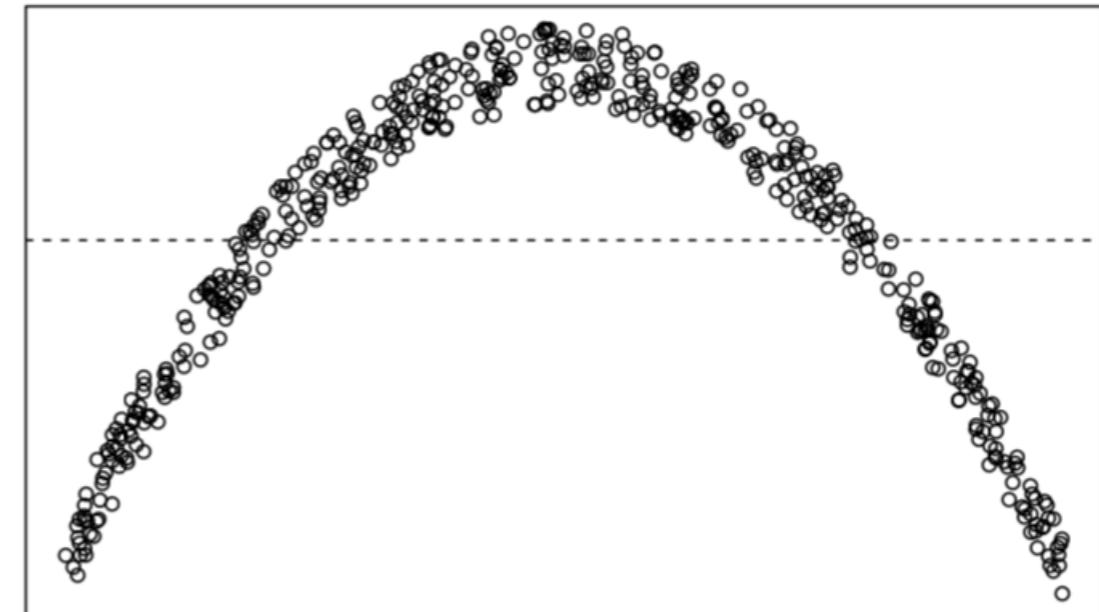


# Model assumptions

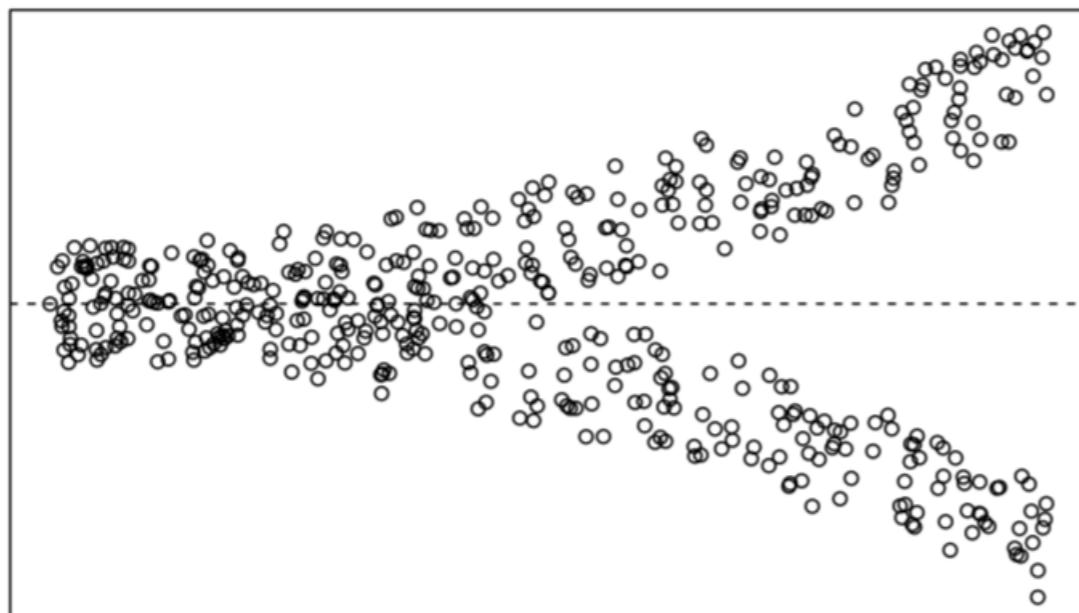
No Violation



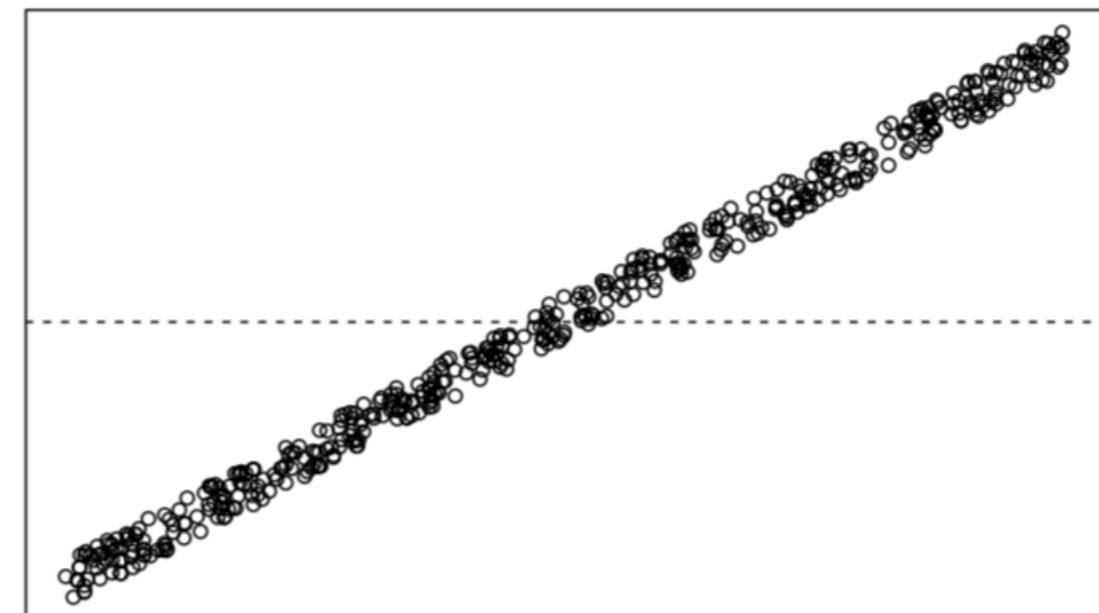
Nonlinear Relationship



Nonconstant Error Variance



Dependent Error Terms



# Summary

- Quick review of statistical inference in frequentist statistics
- Correlation
  - Covariation
  - Pearson's moment correlation
  - Spearman's rank correlation
- Regression
  - The conceptual tour
  - The R route

**Thank you!**