

## Ch. 10: Two Quantitative Variables

1 / 74

### 10.1: Two Quantitative Variables

#### Scatterplots and Correlation

3 / 74

## Navigation

### By Section

- 10.1: start - end
- 10.2: start - end
- 10.3: start - end
- 10.4: start - end
- 10.5: start - end

2 / 74

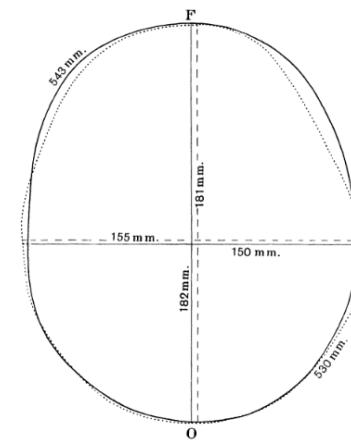
## Graphical Summaries of Quantitative Variables

R.J. Gladstone (1905). "A Study of the Relations of the Brain to the Size of the Head", Biometrika, Vol. 4, p 105-123.

Data collected during 237 autopsies at Middlesex Hospital in London, excluding cases "in which the brain showed a distinctly pathological condition which would have obviously affected its weight"

Variables:

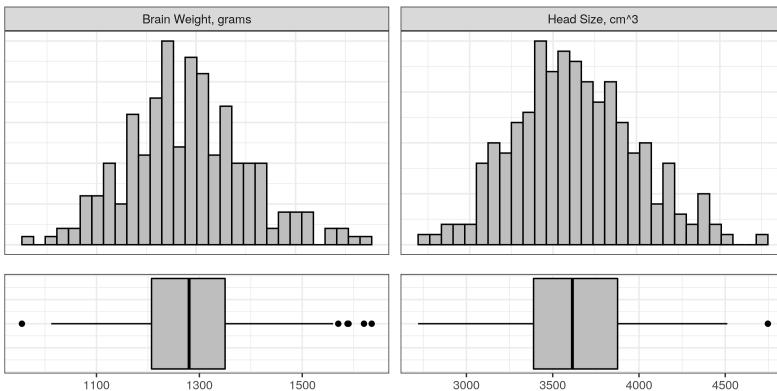
- Gender
- Age (20 - 45 or 46+)
- Brain Weight (g)
- Head Size (cubic cm) the smallest rectangular block which could contain the head



4 / 74

## Graphical Summaries of Quantitative Variables

A single quantitative variable can be summarized visually using a histogram or a bar chart:

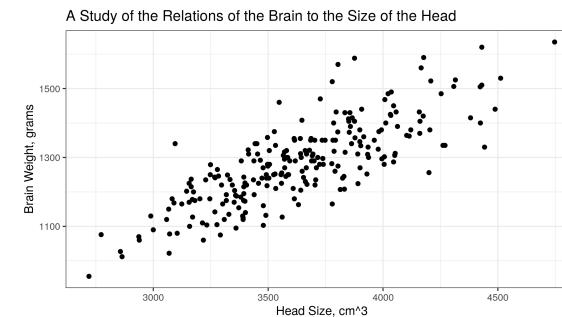


5 / 74

## Graphical Summaries of Quantitative Variables

But, summarizing each variable separately doesn't tell us how the two variables might be related.

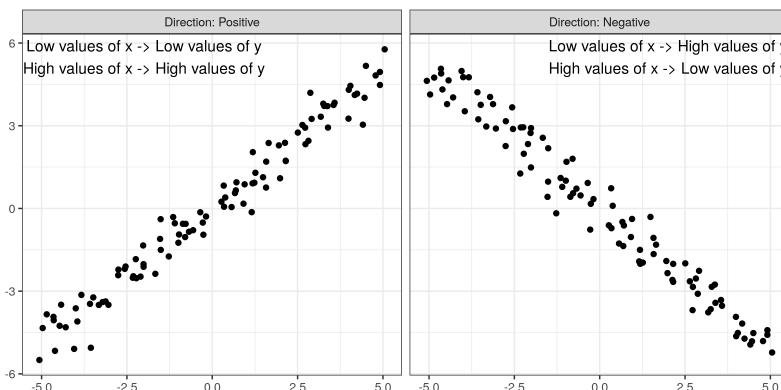
Is there a relationship between brain weight and head size? How do you know?



A **scatterplot** is a plot with the explanatory variable on the x-axis, and the response variable on the y-axis. Observations are shown as points corresponding to a set of quantitative measurements.

6 / 74

## Describing Variable Relationships: Direction

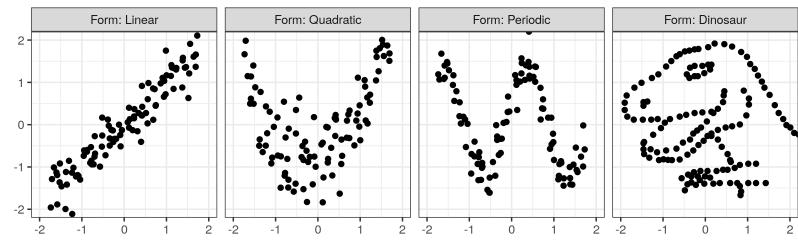


Positive slope: as  $x$  increases,  $y$  increases too.

7 / 74

## Describing Variable Relationships: Form

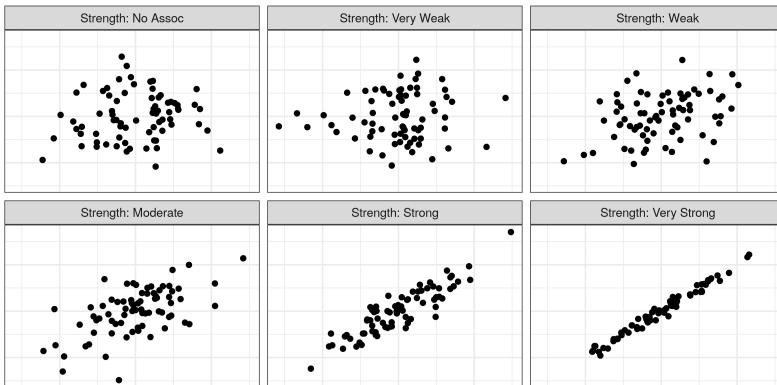
The **form** of an association is whether it follows a linear pattern, or some sort of more complicated pattern - periodic, polynomial (quadratic, cubic, etc.)



8 / 74

## Describing Variable Relationships: Strength

The **strength** of an association indicates how well the value of one variable can be predicted if you know the value of the other variable.

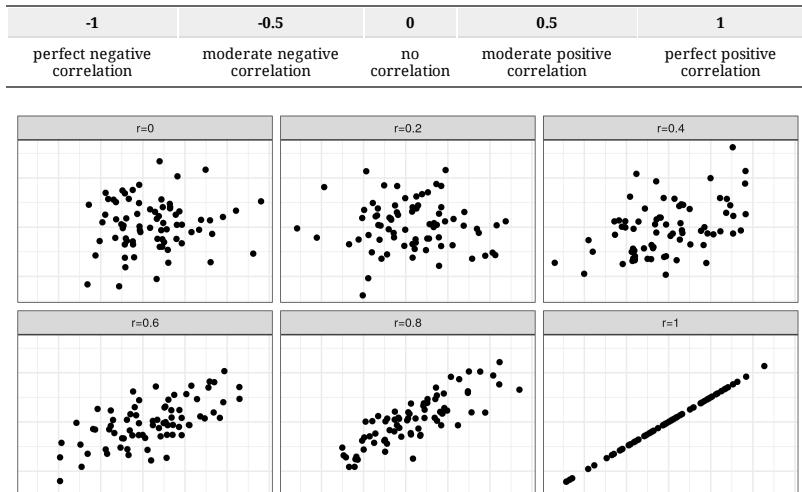


We can describe the strength and direction of a *linear* relationship using the **correlation coefficient**

9 / 74

## Correlation Coefficient

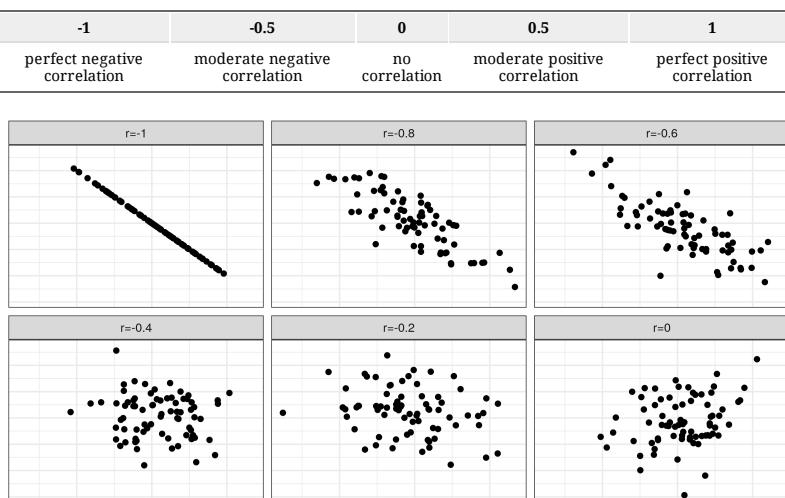
The **correlation coefficient**,  $r$ , is always between -1 and 1.



10 / 74

## Correlation Coefficient

The **correlation coefficient**,  $r$ , is always between -1 and 1.



11 / 74

## Correlation Coefficient

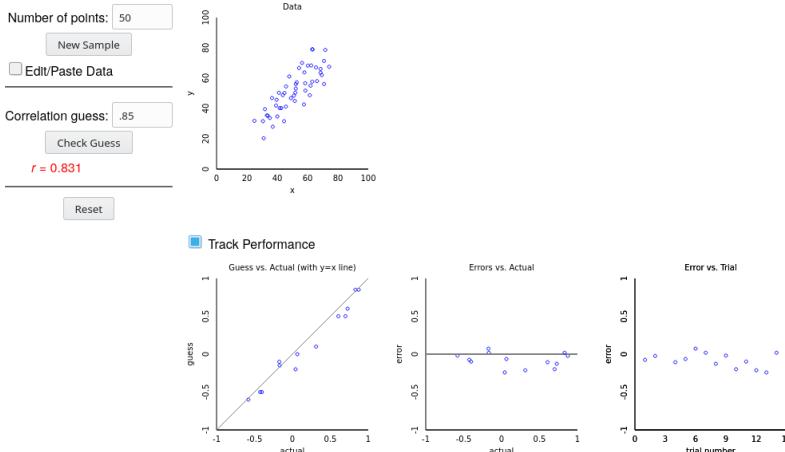
$R^2=0.06$   
REXTHOR, THE DOG-BEARER

I DON'T TRUST LINEAR REGRESSIONS WHEN IT'S HARDER TO GUESS THE DIRECTION OF THE CORRELATION FROM THE SCATTER PLOT THAN TO FIND NEW CONSTELLATIONS ON IT.

12 / 74

## Correlation Coefficient

Get a feel for it by [playing the correlation guessing game!](#)

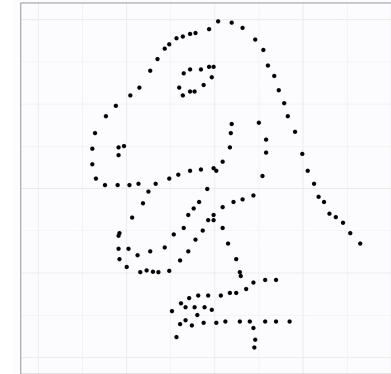


13 / 74

## Correlation Coefficient

The correlation coefficient is only useful for showing the strength of linear relationships.

X: 47.26 (SD = 16.77), Y: 47.83 (SD = 26.94), r = -0.06



All of these plots have essentially the same correlation coefficient, but in some cases there are very clear associations between  $x$  and  $y$

14 / 74

## Correlation Coefficient

(From Calculation Details in the Appendix)

$$r = \frac{1}{n-1} \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s_x} \right) \left( \frac{y_i - \bar{y}}{s_y} \right)$$

15 / 74

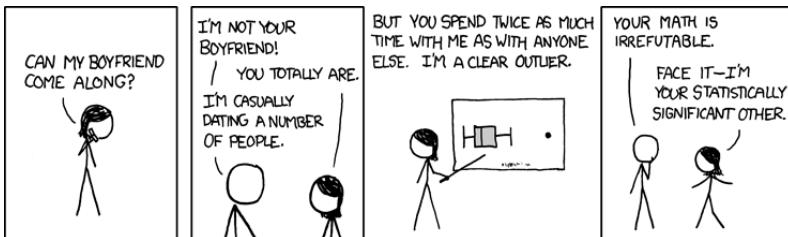
## Correlation is not Causation

Just in case you haven't heard this chant yet, "correlation is not causation". Say it a few times.

It's important to remember that correlation is a measure of association, but that doesn't mean there's any causal factors involved. In some cases, the choice of explanatory and response variables are arbitrary.

16 / 74

## Outliers and Influential Observations



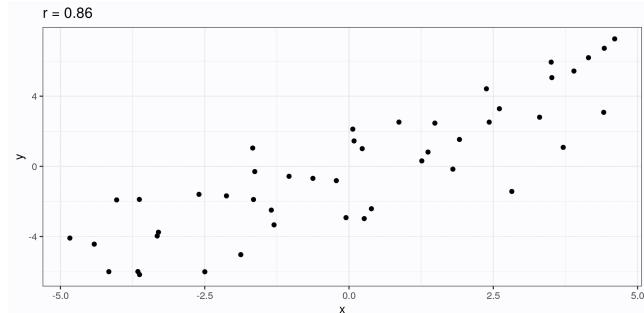
With one variable, outliers are fairly easy to spot

When there are two variables, we don't just have to worry about outliers in one dimension; we also have to worry about **influential observations**

17 / 74

## Outliers and Influential Observations

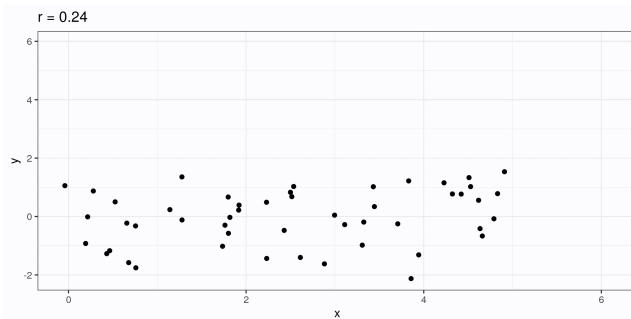
Influential observations are observations which, if included, change our understanding of the relationship between two variables.



18 / 74

## Outliers and Influential Observations

Influential observations are observations which, if included, change our understanding of the relationship between two variables.



19 / 74

## Exploration 10.1

Work through Exploration 10.1 to get a chance to put the material in this section into practice. You can turn it in for 10 points of extra credit in the "Assignment" category.

20 / 74

## 10.2: Inference for the Correlation Coefficient

Simulation Based Approach

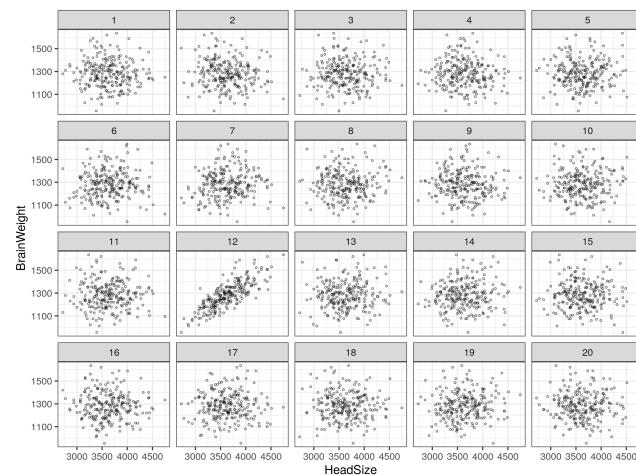
21 / 74

One of these things is not like the others



22 / 74

Which one of these things is not like the others?



23 / 74

Simulation-based Inference for Correlation Coefficient

Our null hypothesis is  $H_0$  : No relationship between  $x$  and  $y$

How can we simulate this?

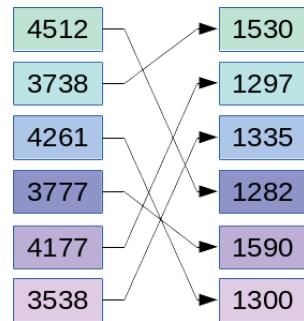
4512	→	1530
3738	→	1297
4261	→	1335
3777	→	1282
4177	→	1590
3538	→	1300

24 / 74

## Simulation-based Inference for Correlation Coefficient

Our null hypothesis is  $H_0$  : No relationship between  $x$  and  $y$

How can we simulate this?



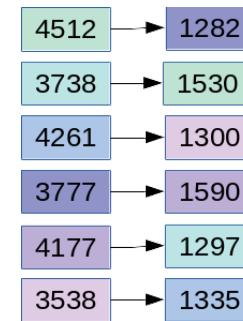
If there's no relationship between  $x$  and  $y$ , then it doesn't really matter what  $x$  value is paired with a given  $y$  value... so we can just change which values are paired together.

25 / 74

## Simulation-based Inference for Correlation Coefficient

Our null hypothesis is  $H_0$  : No relationship between  $x$  and  $y$

How can we simulate this?



This is equivalent to shuffling the order of  $y$  and creating a new regression

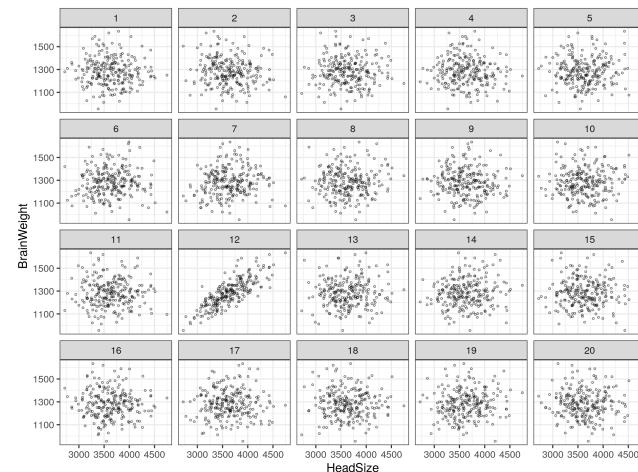
26 / 74

## Summary of Simulation Model

- Null hypothesis: No association between  $x$  and  $y$  variables  
The population symbol for  $r$  is  $\rho$ , so  $H_0 : \rho = 0$
- One repetition: Re randomizing the response outcomes to the explanatory variable values (randomize  $y$  values)
- Statistic: Correlation coefficient,  $r$

27 / 74

## Which one of these things is not like the others?



Here, our  $y$  values have been shuffled for each sub-plot that isn't plot 12.  
Plot 12 contains the original data.

28 / 74

## Hypotheses for Correlation Coefficient

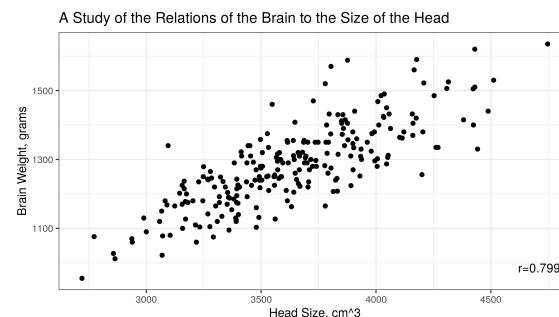
$H_0$  : There is no relationship between  $x$  and  $y$  ( $\rho = 0$ )

$H_A$  : There is a relationship between  $x$  and  $y$  ( $\rho \neq 0$ )

$\rho$  (pronounced "row", spelled "rho") is the population version of  $r$

**Caution**  $\rho$  is not  $p$  and is NOT a p-value.

## Example: Head Size and Brain Weight



$H_0$  : No linear relationship between head size and brain weight ( $\rho = 0$ )

$H_A$  : There is a linear relationship between head size and brain weight ( $\rho \neq 0$ )

29 / 74

30 / 74

## Example: Head Size and Brain Weight

1. **Statistic** - correlation coefficient from the sample

2. **Simulate**

- Assume there is no relationship between head size ( $x$ ) and brain weight ( $y$ )
- Shuffle the values of  $y$
- Calculate the correlation coefficient  $r^*$  from the simulated data

3. **Strength of evidence** - in how many simulated samples did we get a correlation coefficient  $r^*$  with magnitude greater than our sample correlation coefficient  $r$ ?

## Example: Head Size and Brain Weight

Conclusion:

I reject  $H_0$  that there is no linear relationship between head size and brain weight with  $p < 0.001$ . There is very strong evidence that the linear relationship observed between the variables did not occur due to random chance, thus, we must conclude that there is a linear relationship between head size and brain weight, that is, that  $\rho \neq 0$ .

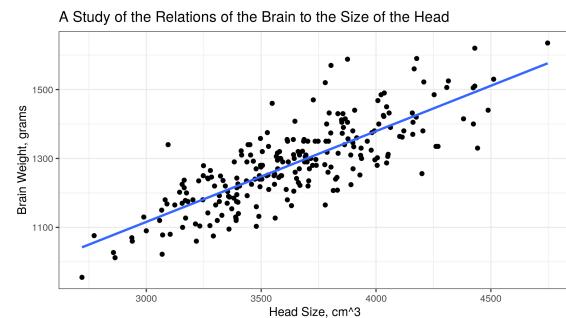
31 / 74

32 / 74

## 10.3: Least Squares Regression

33 / 74

### Regression Line



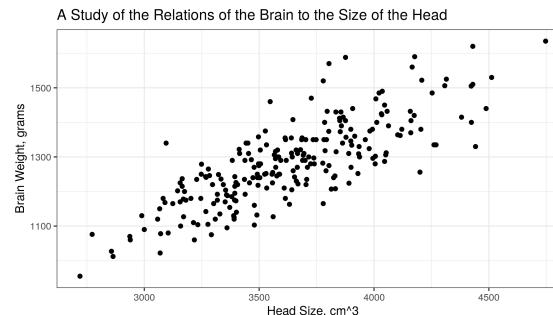
General equation:  $\hat{y} = a + bx$  where

- $x$  is the explanatory variable
- $\hat{y}$  is the response variable
- $a$  is the **y-intercept** (the predicted value when  $x = 0$ )
- $b$  is the **slope**

35 / 74

## Motivation

In the last two sections, we've talked about the Brain weight vs. Head size data:



34 / 74

### Regression Line

The equation of the regression line in the picture on the last slide is:

$$\hat{y} = 325.573 + 0.263x$$

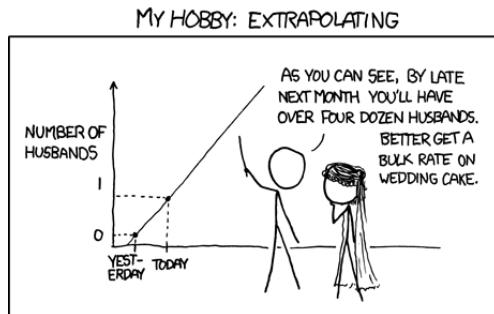
What does that mean?

- With a head size of  $0 \text{ cm}^3$ , we would expect a brain weight of ...  $325.57 \text{ g}$  (on average)  
Sometimes, the interpretation of the intercept doesn't make a lot of sense in context.
- An increase in head size of  $1 \text{ cm}^3$  leads to a predicted change in average brain weight of  $0.2634 \text{ grams}$

Note that the sign of the slope is the same as the sign of the correlation coefficient, because both indicate the direction of the association.

36 / 74

## Extrapolation



Predicting values for the response variable for explanatory variable values outside the range of the original data is known as **extrapolation** and can lead to very misleading predictions.

37 / 74

## Extrapolation

Predicting values for the response variable for explanatory variable values outside the range of the original data is known as **extrapolation** and can lead to very misleading predictions.

For instance, if we try to predict the brain weight of a 2.5 month old child included in the original brain weight publication (but not in our dataset)...

$$x = 1212, \text{ so } \hat{y} = ?$$

$$325.573 + 0.263(1212) = 644.850$$

Our model would predict that the child's brain would weigh about 645 grams, on average. In fact, it weighed 490 grams.

Our prediction had an error of  $490 - 644.850 = -154.850$

This error is called a **residual**

38 / 74

## Extrapolation

Extrapolation is sometimes reasonable, but often is ill-advised.

- I have data on the foot length and height of 58 people, ranging from 4'10 to 6'1. I want to predict the foot length of someone who is
  - 4'2 :
  - 6'3 :
  - 7':
- I have data on the number of ice cream cones sold at a shop and the temperature for January - April. I want to predict how many cones will be sold at a shop per day during
  - July :
  - early May :
  - April 30, in a similarly sized midwestern town :
  - in Saudi Arabia during May :

39 / 74

## Extrapolation

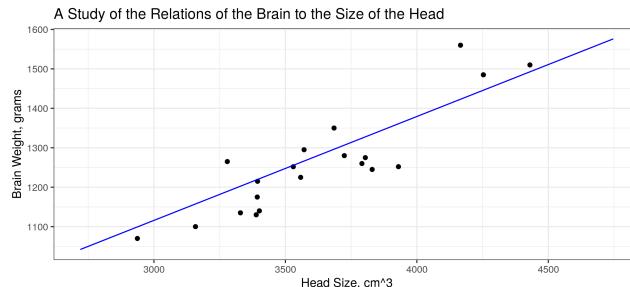
Extrapolation is sometimes reasonable, but often is ill-advised.

- I have data on the foot length and height of 58 people, ranging from 4'10 to 6'1. I want to predict the foot length of someone who is
  - 4'2 : **not a good idea** - most adults are taller, the relationship may not hold for children
  - 6'3 : **probably fine** - we're extrapolating by 2", and someone who is 6'1 is probably not so different from someone who is 6'3
  - 7': **not a good idea**, only a few people have gotten to be that tall, and the general population may not represent them well
- I have data on the number of ice cream cones sold at a shop and the temperature for January - April. I want to predict how many cones will be sold at a shop per day during
  - July : **not a good idea** (much higher avg. temp)
  - early May : **probably fine** (avg temp approximately same as late April)
  - April 30, in a similarly sized midwestern town : **may be reasonable, but use caution**
  - in Saudi Arabia during May : **not a good idea** (much, much higher avg. temp, different location that is not similar)

40 / 74

## Least Squares Regression

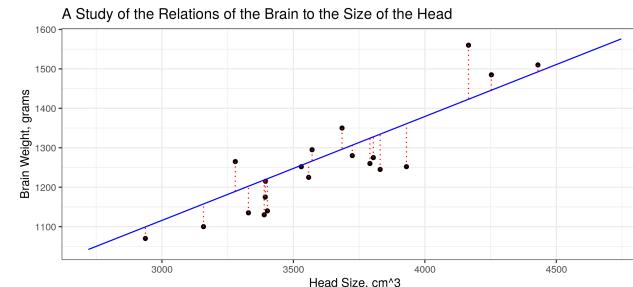
How do we get our regression line?



41 / 74

## Least Squares Regression

How do we get our regression line?



42 / 74

## Least Squares Regression

How do we get our regression line?

From the computational details section of the appendix:

$$b = r \frac{s_y}{s_x}$$

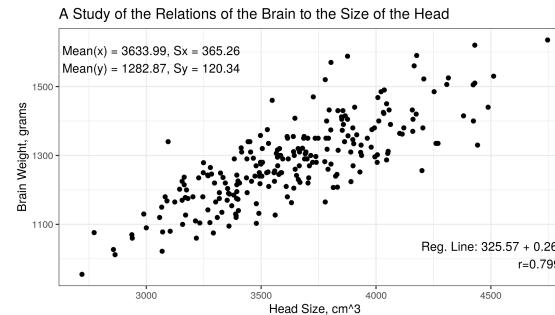
where  $r$  is the correlation coefficient.

Then,

$$a = \bar{y} - b\bar{x}$$

43 / 74

## Example: Head Size and Brain Weight



We can compute  $b$  and  $a$  from the other information given on the graph:

$$b = r \frac{s_x}{s_y} = 0.8 \times \frac{120.34}{365.261} = 0.263$$

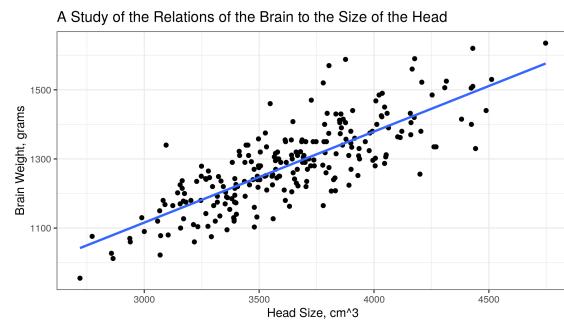
$$a = \bar{y} - b\bar{x} = 1282.873 - 0.263(3633.992) = 325.573$$

44 / 74

## Coefficient of Determination ( $R^2$ )

We can measure how well our data fit our regression line with  $R^2$  (which is literally the correlation value, squared).

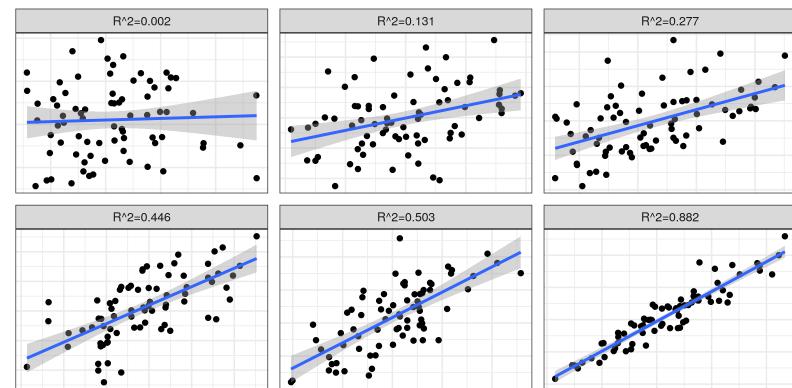
$R^2$  is between 0 and 1, where 1 = perfect predictions and 0 = no relationship at all.



The  $R^2$  for this plot is 0.639

45 / 74

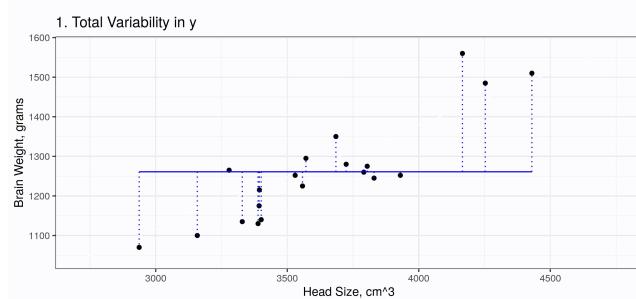
## Coefficient of Determination ( $R^2$ )



46 / 74

## Coefficient of Determination ( $R^2$ )

$R^2$  can be interpreted as the percentage of variability in  $y$  which is explained by the regression line.



47 / 74

## Coefficient of Determination ( $R^2$ )

In math, the last series of graphs can be shown like this:

$$\underbrace{SSE(y - \bar{y})}_{\text{total error}} = \underbrace{SSE(y - \hat{y})}_{\text{residual error}} + \underbrace{SSE(\hat{y} - \bar{y})}_{\text{explained by regression}}$$

This is slightly different than the book's notation, where  $SSE(y - \bar{y}) = SSE(\bar{y})$  and  $SSE(\hat{y} - \bar{y})$  is written as SSE(regression line)

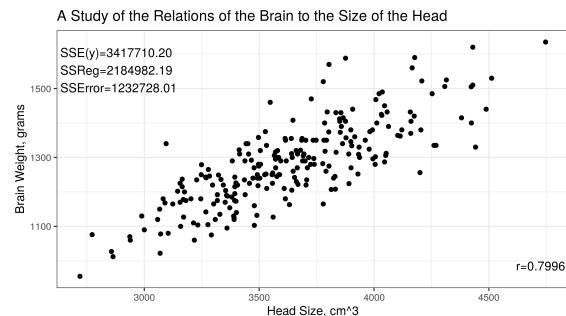
In math, then, we can write  $R^2$  as

$$R^2 = 100 \times \left( 1 - \frac{SSE(y - \hat{y})}{SSE(y - \bar{y})} \right) = 100 \times \left( 1 - \frac{\text{residual std error}}{\text{total std error}} \right)$$

Notice that this equation works out to equal the proportion of the total error in  $y$  explained by the regression.

48 / 74

## Example: Head Size and Brain Weight



Let's compute  $R^2$  using the SSE formula and check that it's the same as the value of  $r$ , squared.

## Example: Head Size and Brain Weight

$$R^2 = \frac{SS(\hat{y} - \bar{y})}{SS(y - \bar{y})} = \frac{SSReg}{SSE(y)} = \frac{2.185 \times 10^6}{3.418 \times 10^6} = 0.639$$

$$R^2 = r^2 = (.7996)^2 = 0.639$$

49 / 74

50 / 74

## Exploration 10.3

Work through Exploration 10.3 (ignoring the parts which say compare with classmates) to review the concepts covered in this lecture.

Upload your answers to Canvas for extra credit.

## 10.4: Inference for the Regression Slope: Simulation-Based Approach

51 / 74

52 / 74

## Review

Regression equation:  $\hat{y} = a + bx$

- $\hat{y}$  is the
- $a$  is the
- $b$  is the
- $x$  is the

53 / 74

## Review

Regression equation:  $\hat{y} = a + bx$

- $\hat{y}$  is the **response variable**
- $a$  is the **intercept**
- $b$  is the **slope**
- $\hat{x}$  is the **explanatory variable**

54 / 74

## Inference for Regression Slope

To test whether two variables are linearly associated, we usually want to know if  $b = 0$ .

If  $b = 0$  then  $\hat{y} = a + 0x = a$ ... That is,  $x$  doesn't affect  $y$  at all.

What happens if we estimate a  $b \neq 0$  by chance?

55 / 74

## Simulation-based Inference for the Slope

Our null hypothesis is  $H_0$  : No relationship between  $x$  and  $y$ , that is,  $\beta = 0$   
 $\beta$  is the greek letter corresponding to  $b$  - the "true/population slope"

How can we simulate this?

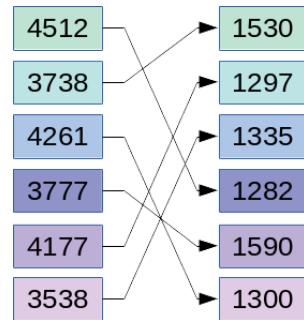
4512	→	1530
3738	→	1297
4261	→	1335
3777	→	1282
4177	→	1590
3538	→	1300

56 / 74

## Simulation-based Inference for the Slope

Our null hypothesis is  $H_0$  : No relationship between  $x$  and  $y$ , that is,  $\beta = 0$   
 $\beta$  is the greek letter corresponding to  $b$  - the "true/population slope"

How can we simulate this?



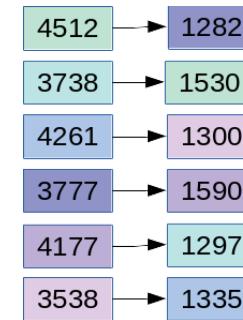
If there's no relationship between  $x$  and  $y$ , then it doesn't really matter what  $x$  value is paired with a given  $y$  value... so we can just change which values are paired together.

57 / 74

## Simulation-based Inference for the Slope

Our null hypothesis is  $H_0$  : No relationship between  $x$  and  $y$ , that is,  $\beta = 0$   
 $\beta$  is the greek letter corresponding to  $b$  - the "true/population slope"

How can we simulate this?



This is equivalent to shuffling the order of  $y$  and creating a new regression

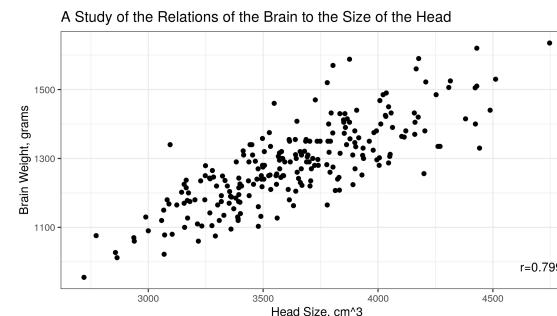
58 / 74

## Summary of Simulation Model

- Null hypothesis: No association between  $x$  and  $y$  variables ( $\beta = 0$ )
- One repetition: Re randomizing the response outcomes to the explanatory variable values (randomize  $y$  values)
- Statistic: Slope,  $b$

59 / 74

## Example: Head Size and Brain Weight



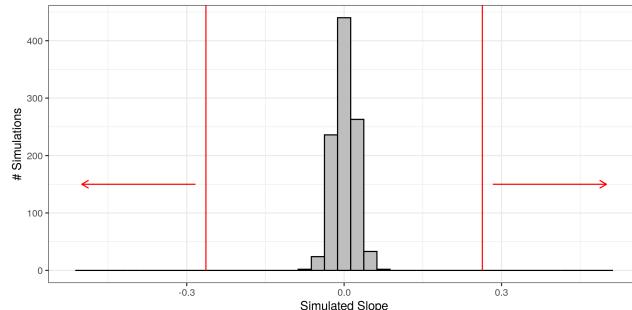
$H_0$  : No linear relationship between head size and brain weight ( $\beta = 0$ )

$H_A$  : There is a linear relationship between head size and brain weight ( $\beta \neq 0$ )

60 / 74

## Example: Head Size and Brain Weight

1. Statistic - slope computed from the sample
2. Simulate
3. Strength of evidence - in how many simulated samples did we get a slope  $b^*$  with magnitude greater than our sample slope  $b$ ?



In 0 samples out of 1000,  $|b^*| > 0.263$ , so  $p < 0.001$

61 / 74

## Example: Head Size and Brain Weight

Conclusion:

I reject  $H_0$  that there is no linear relationship between head size and brain weight with  $p < 0.001$ . There is very strong evidence that the linear relationship observed between the variables did not occur due to random chance, thus, we must conclude that there is a linear relationship between head size and brain weight, that is, that  $\beta \neq 0$ .

Note that this is the same interpretation (word for word) as we used in Section 10.2 (with the parameter changed)

62 / 74

## 10.5: Inference for the Regression Slope: Theory-Based Approach

### Requirements for Theory-Based Inference

Validity conditions:

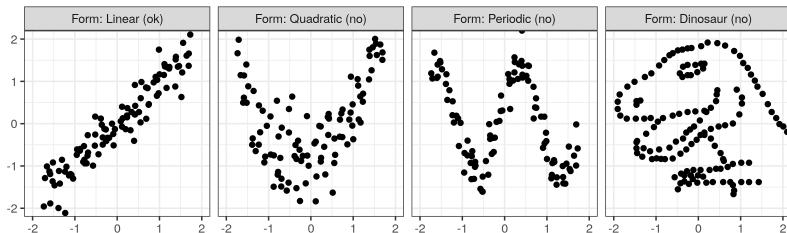
- general pattern of the points should follow a linear trend (no curved or nonlinear patterns)
- about the same number of points above and below the regression line (symmetry)
- variability of points around the regression line should be similar regardless of the value of  $x$  (equal variance)

63 / 74

64 / 74

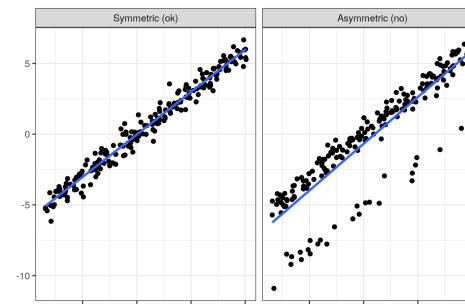
## Requirements for Theory-Based Inference

- general pattern of the points should follow a linear trend (no curved or nonlinear patterns)



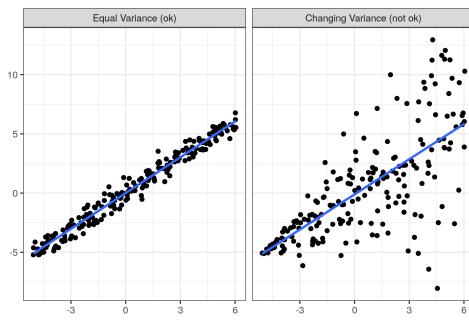
65 / 74

## Requirements for Theory-Based Inference



66 / 74

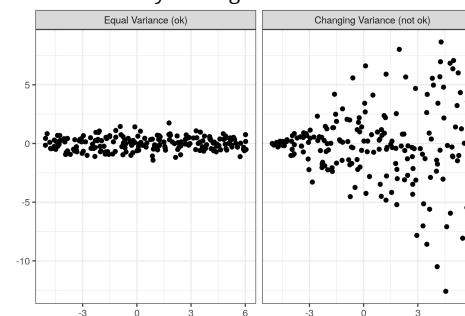
## Requirements for Theory-Based Inference



67 / 74

## Requirements for Theory-Based Inference

Looking at the residuals from your regression can be more useful:



68 / 74

## Theory-Based Inference

We also need a formula for the standard error of our statistic,  $r$

$$SE(r) = \sqrt{\frac{1 - r^2}{n - 2}}$$

Then,

$$t = \frac{r}{SE(r)}$$

Because  $t$  can be calculated from either the slope or the correlation coefficient, the test results will be identical.

We could also make a theory-based confidence interval for  $\rho$  using the standard error.

69 / 74

## Theory-Based Inference

If we work with  $b$ , we will be given  $SE(b)$  by statistical software:

```
##             Estimate Std. Error t value Pr(>|t|)  
## (Intercept) 325.573   47.1409   6.91 4.61e-11  
## HeadSize     0.263     0.0129   20.41 5.96e-54
```

Then,

$$t = \frac{b - 0}{SE(b)}$$

Or, just get  $t$  (and the corresponding p-value) from the table...

Because  $t$  can be calculated from either the slope or the correlation coefficient, the tests are identical.

We could also make a theory-based confidence interval for  $\beta$  using the standard error.

70 / 74

## Interpretations - Slope

Hypothesis test interpretation:

With  $t = ...$ , I have (strength) evidence to suggest that there is a linear relationship between (explanatory variable) and (response variable). I (accept/reject)  $H_0$  and conclude that (conclusion in context of the problem)

Confidence interval interpretation:

I am 95% confident that as explanatory variable increases by 1 unit, the predicted population average value of response variable will change by between (lower bound of b interval) and (upper bound of b interval)

71 / 74

## Example: Head Size and Brain Weight

```
##             Estimate Std. Error t value Pr(>|t|)  
## (Intercept) 325.573   47.1409   6.91 4.61e-11  
## HeadSize     0.263     0.0129   20.41 5.96e-54
```

We can see that the  $t$  value is 20.409, which corresponds to extremely strong evidence that there is a relationship between head size and brain weight.

72 / 74

## Example: Head Size and Brain Weight

```
##             Estimate Std. Error t value Pr(>|t|)  
## (Intercept) 325.573    47.1409   6.91 4.61e-11  
## HeadSize     0.263     0.0129   20.41 5.96e-54
```

If we wanted to make a 95% CI for the slope of the line (e.g. for the increase in brain weight when head size increases by 1 cm<sup>3</sup>), we could use the 2\*SE method:

$$0.263 \pm 2 * 0.0129 = (0.237, 0.289)$$

I am 95% confident that the population average brain weight will increase by between 0.237 and 0.289 grams for each additional cubic centimeter of head size.

73 / 74

## Interpretations - Correlation Coefficient

Hypothesis test interpretation:

With  $t = \dots$ , I have (strength) evidence to suggest that there is a linear relationship between (explanatory variable) and (response variable). I accept/reject  $H_0$  and conclude that conclusion in context of the problem

Confidence interval interpretation:

I am 95% confident that the population correlation between explanatory variable and response variable is between lower bound and upper bound

74 / 74