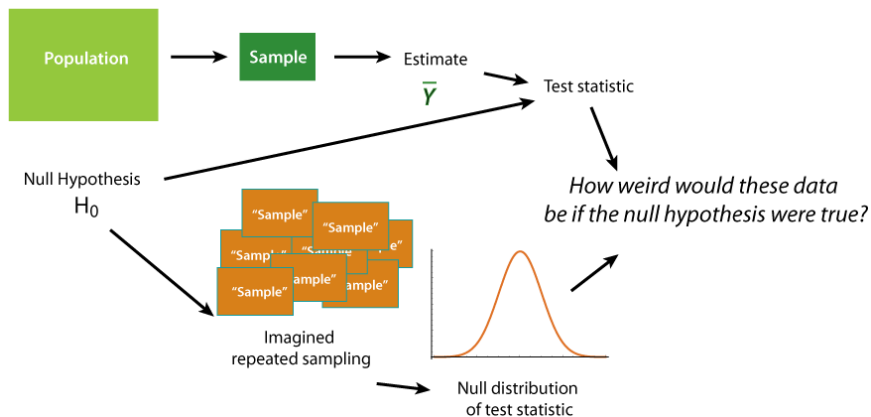# Hypothesis testing

Chapter 6

Hypothesis testing asks how unusual it is to get data that differ from the null hypothesis.

If the data would be quite unlikely under $H_0$, we reject $H_0$.



Population → Sample → Estimate $\bar{Y}$ → Test statistic

Null Hypothesis $H_0$

"Sample" "Sample" "Sample" "Sample" "Sample" "Sample"

Imagined repeated sampling

*How weird would these data be if the null hypothesis were true?*

Null distribution of test statistic

## Hypotheses are about populations, but are tested with data from samples

Hypothesis testing usually assumes that sampling is random.

Null hypothesis: a specific statement about a population parameter made for the purposes of argument.

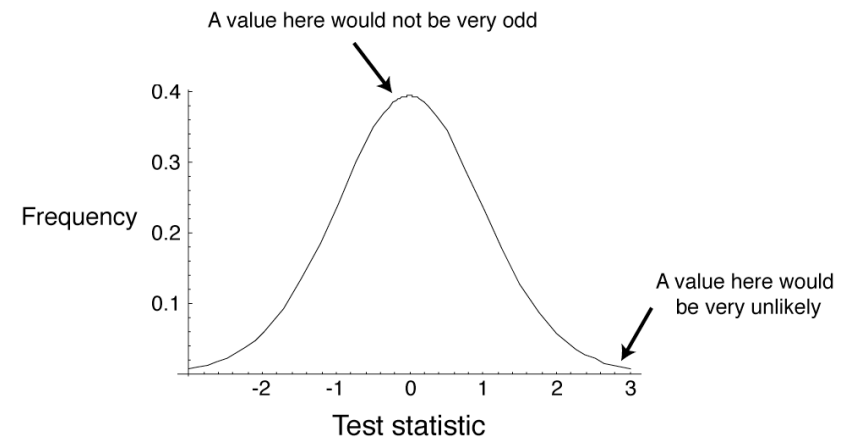Alternate hypothesis: represents all other possible parameter values except that stated in the null hypothesis.

The *null hypothesis* is usually the simplest statement, whereas the *alternative hypothesis* is usually the statement of greatest interest.

A good null hypothesis would be interesting if proven wrong.

A null hypothesis is specific;
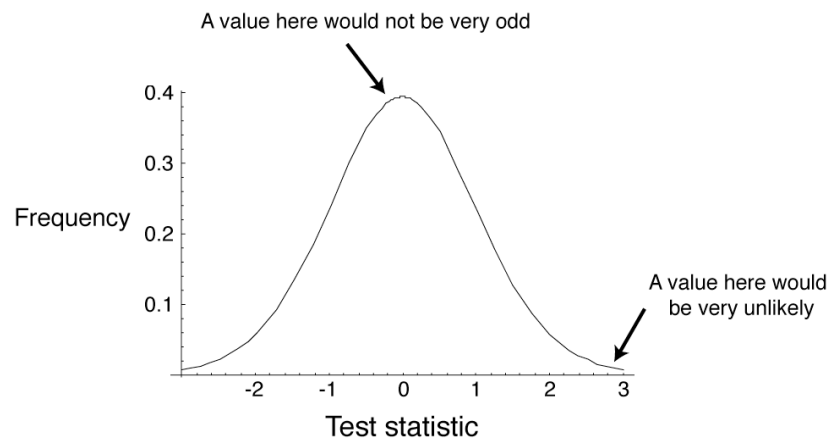an alternate hypothesis is not.

# Test Statistic

A number calculated to represent the match between a set of data and the null hypothesis

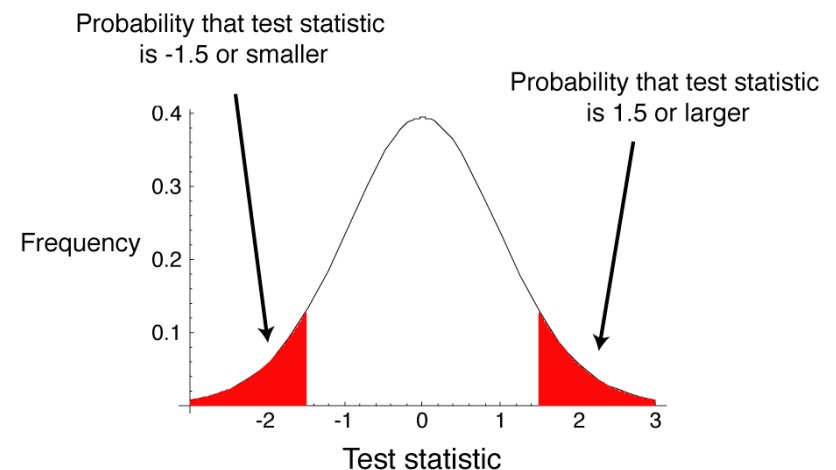Can be compared to a general distribution to infer probability

A value here would not be very odd

A value here would be very unlikely

Test statistic

Possible outcomes from samples under null hypothesis

A value here would not be very odd

A value here would be very unlikely

Test statistic

A test statistic summarizes the match between the data and the null hypothesis

# *P*-value

Probability that test statistic is -1.5 or smaller

Probability that test statistic is 1.5 or larger
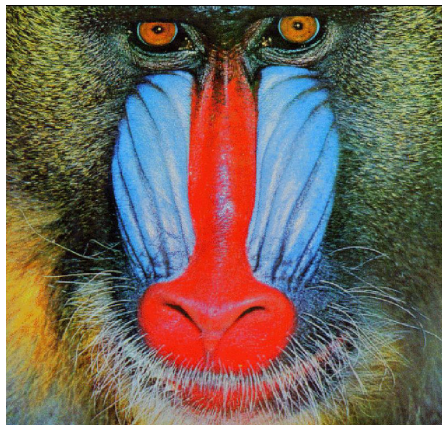
Test statistic

## How to find *P*-values

Simulation

Parametric tests

Permutation

A *P*-value is the probability of getting the data, or something as or more unusual, if the null hypothesis were true.

## Hypothesis testing: an example

**Does a red shirt help win wrestling?**

## The experiment and the results

Animals use red as a sign of aggression

Does red influence the outcome of wrestling, taekwondo, and boxing?

- 16 of 20 rounds had more red-shirted than blue-shirted winners in these sports in the 2004 Olympics

- Shirt color was randomly assigned

Hill, RA, and RA Burton 2005. Red enhances human performance in contests Nature 435:293.

## Stating the hypotheses

$H_0$: Red- and blue-shirted athletes are <u>equally likely</u> to win (*proportion* = 0.5).

$H_A$: Red- and blue-shirted athletes are <u>not equally likely</u> to win (*proportion* ≠ 0.5).

Is this discrepancy by chance alone?:
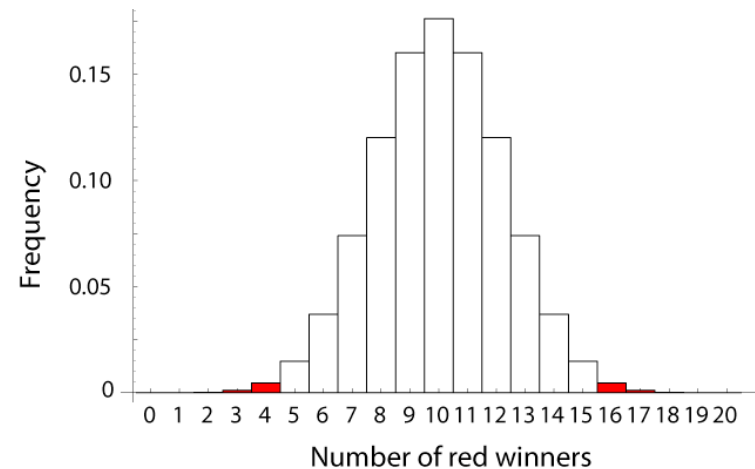
Estimating the probability of such an extreme result

The *null distribution* for a test statistic is the probability distribution of alternative outcomes when a random sample is taken from a population corresponding to the null expectation.

## Estimating the value

16 of 20 is a proportion of *proportion* = 0.8

This is a discrepancy of 0.3 from the proportion proposed by the null hypothesis, *proportion* = 0.5

## The null distribution of the *sample* number of red wins

## Calculating the *P*-value from the null distribution

The *P*-value is calculated as

$P = 2 \times [\text{Pr}(16) + \text{Pr}(17) + \text{Pr}(18) + \text{Pr}(19) + \text{Pr}(20)] = 0.012$.

$\alpha$ is often 0.05

## Statistical significance

The significance level, $\alpha$, is a probability used as a criterion for rejecting the null hypothesis.

If the *P*-value for a test is less than or equal to $\alpha$, then the null hypothesis is rejected.

## Significance for the red shirt example

$P = 0.012$

$P < \alpha$, so we can reject the null hypothesis

Athletes in red shirts were more likely to win.

# Larger samples give more information

A larger sample will tend to give and estimate with a smaller confidence interval

A larger sample will give more power to reject a false null hypothesis

# Sample R code for doing this simulation (Note: This is not the most efficient code for this!)

```
binarySample = function(n, prob){
  results = rep(NA,n)
  for(i in 1:n){
    if(runif(1) < prob) results[i] = "red"
  else
    results[i] = "blue"
  }
  length(which(results=="red"))
}

numreps=10000
resultsDF = data.frame(numberRedWins =
    replicate(numreps, binarySample(20,.5)))
```

# Hypothesis testing: another example

Do dogs resemble their owners?



# Common wisdom holds that dogs resemble their owners. Is this true?

41 dog owners approached in parks; photos taken of dog and owner separately

Photo of owner and dog, along with another photo of dog, shown to students to match

Roy, M.M., & Christenfeld, N.J.S. (2004). Do dogs resemble their owners? *Psychological Science*, **15**, 361–363

## Hypotheses

$H_0$: The proportion of correct matches is *proportion* $= 0.5$.

$H_A$: The proportion of correct matches is different from *proportion* $= 0.5$.
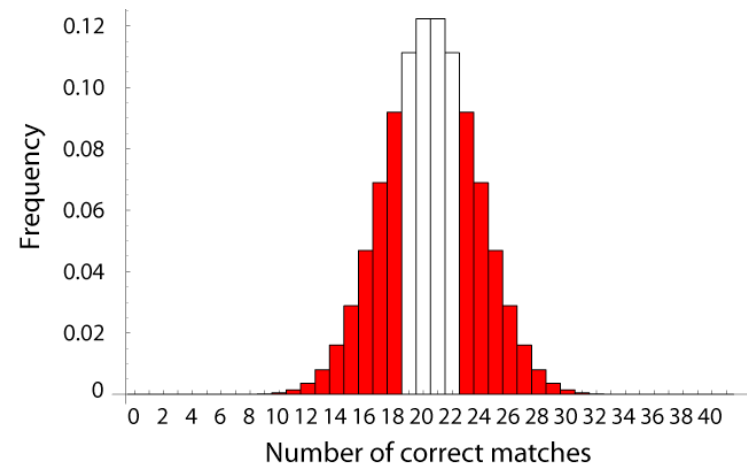
## Data

Of 41 matches, 23 were correct and 18 were incorrect.

## Estimating the proportion

$$sample\ proportion = \frac{23}{41} = 0.56$$

## Null distribution for dog/owner resemblance

## The *P*-value:

$P = 0.53$

We do not reject the null hypothesis that dogs do not resemble their owners.

## Jargon

## Significance level

The acceptable probability of rejecting a true null hypothesis

Called $\alpha$

For many purposes, $\alpha = 0.05$ is acceptable. $\alpha$ is somewhat arbitrarily chosen by researchers.

## Type I error

Rejecting a true null hypothesis

Probability of Type I error is $\alpha$ (the significance level)

## Type II error

Not rejecting a false null hypothesis

The probability of a Type II error is $\beta$.

The smaller $\beta$, the more *power* a test has.

## Power

The ability of a test to reject a false null hypothesis

Power $= 1 - \beta$

## Power increases with more information (i.e. with larger sample size)

## One- and two-tailed tests

Most tests are *two-tailed tests.*

This means that a deviation in either direction would reject the null hypothesis.

Normally $\alpha$ is divided into $\alpha/2$ on one side and $\alpha/2$ on the other.

Test statistic

## One-tailed tests

Only used when the other tail is nonsensical

For example, comparing grades on a multiple choice test to that expected by random guessing

## Critical value

The value of a test statistic beyond which the null hypothesis can be rejected

We never "accept the null hypothesis"

"Statistically significant"
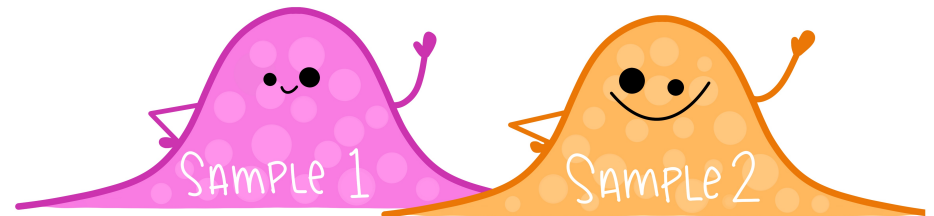
$P < \alpha$

We can "reject the null hypothesis"
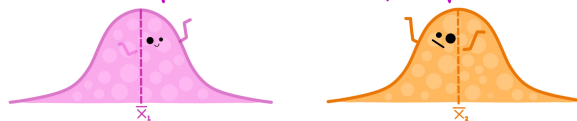


2-SAMPLE T-TESTS

@allison_horst

teaching assistants:

Sample 1

Sample 2

LET'S START HERE: if random samples are drawn from populations w/ the same mean...

Then it is more likely that the 2 sample means will be close together... (i.e. the same population)
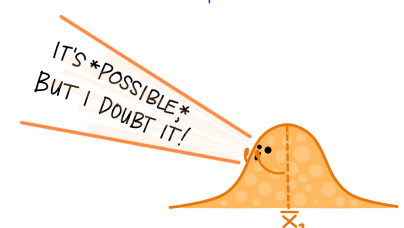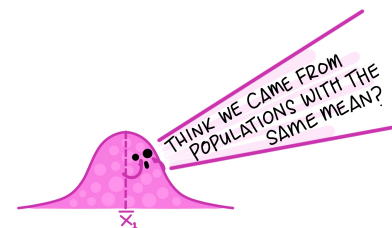
...and it is less likely (but always possible!) that the sample means will be far apart.

$\overline{x}_1$     $\overline{x}_2$

@allison_horst

in OTHER WORDS...: The more different the sample means are,* the less likely it is they were drawn from populations w/ the same mean.

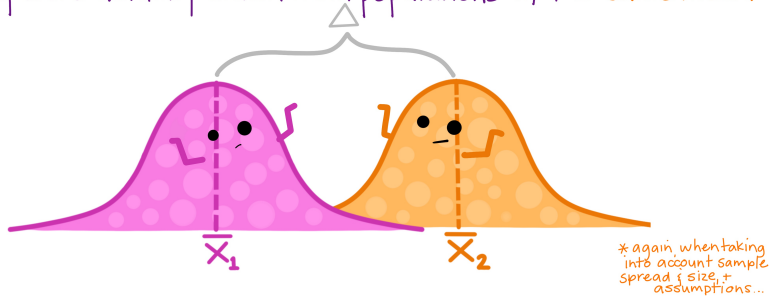*(when taking into account sample spread + size) & assuming we've randomly sampled

THINK WE CAME FROM POPULATIONS WITH THE SAME MEAN?

IT'S *POSSIBLE* BUT I DOUBT IT!

$\overline{x}_1$     $\overline{x}_2$
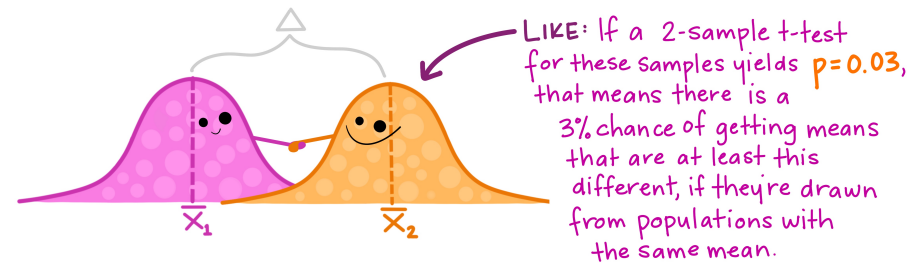
@allison_horst

So for our 2 random samples, we ask:

WHAT IS THE PROBABILITY OF GETTING 2 SAMPLE MEANS THAT ARE AT LEAST THIS DIFFERENT,* if they were actually drawn from populations w/ the same mean?

$\overline{X}_1$   $\overline{X}_2$

*again, when taking into account sample spread & size, + assumptions...
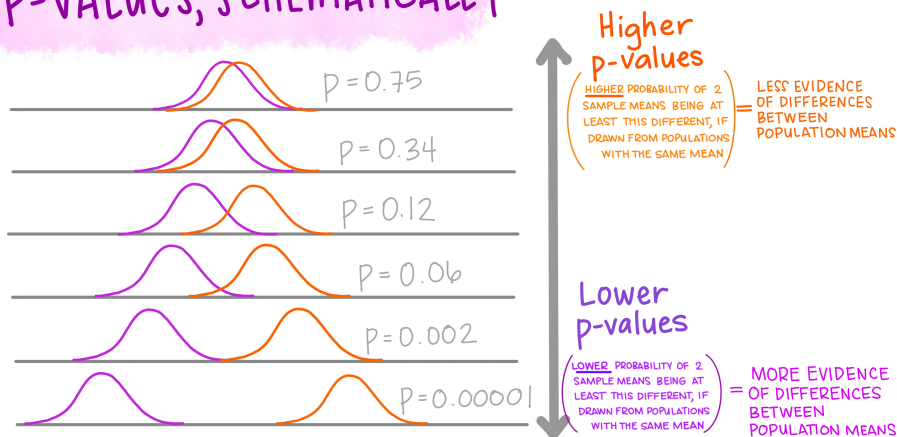
That's our p-value!

WHAT IS THE PROBABILITY OF GETTING 2 SAMPLE MEANS THAT ARE AT LEAST THIS DIFFERENT, if they were actually drawn from populations w/ the same mean?

$\overline{X}_1$   $\overline{X}_2$

LIKE: If a 2-sample t-test for these samples yields p=0.03, that means there is a 3% chance of getting means that are at least this different, if they're drawn from populations with the same mean.

P-VALUES, SCHEMATICALLY:

p = 0.75

p = 0.34

p = 0.12

p = 0.06

p = 0.002

p = 0.00001

Higher p-values

HIGHER PROBABILITY OF 2 SAMPLE MEANS BEING AT LEAST THIS DIFFERENT, IF DRAWN FROM POPULATIONS WITH THE SAME MEAN = LESS EVIDENCE OF DIFFERENCES BETWEEN POPULATION MEANS

Lower p-values

LOWER PROBABILITY OF 2 SAMPLE MEANS BEING AT LEAST THIS DIFFERENT, IF DRAWN FROM POPULATIONS WITH THE SAME MEAN = MORE EVIDENCE OF DIFFERENCES BETWEEN POPULATION MEANS

Question:
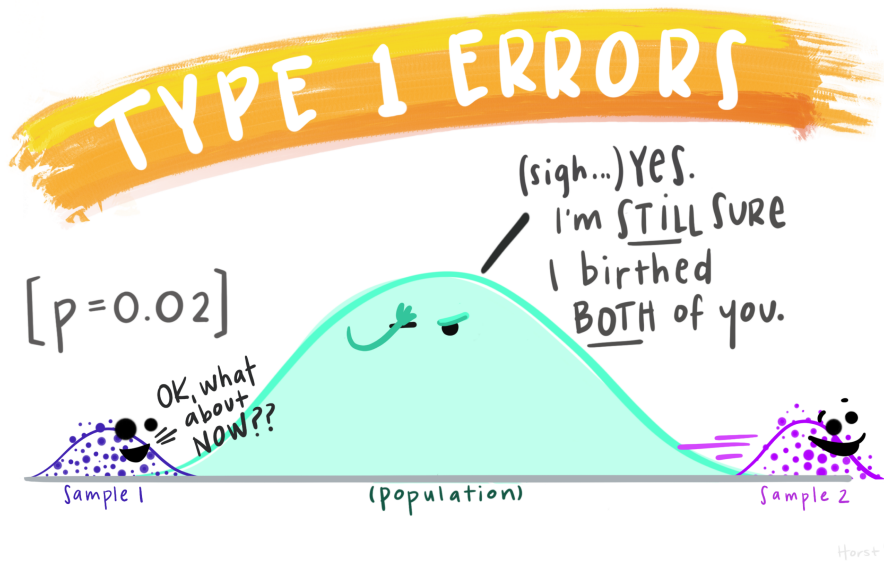WHEN DO WE HAVE ENOUGH EVIDENCE TO SAY THERE IS A SIGNIFICANT DIFFERENCE?

Answer:
WHEN OUR P-VALUE IS BELOW OUR SELECTED SIGNIFICANCE LEVEL ($\alpha$), USUALLY (BUT NOT ALWAYS) = 0.05.
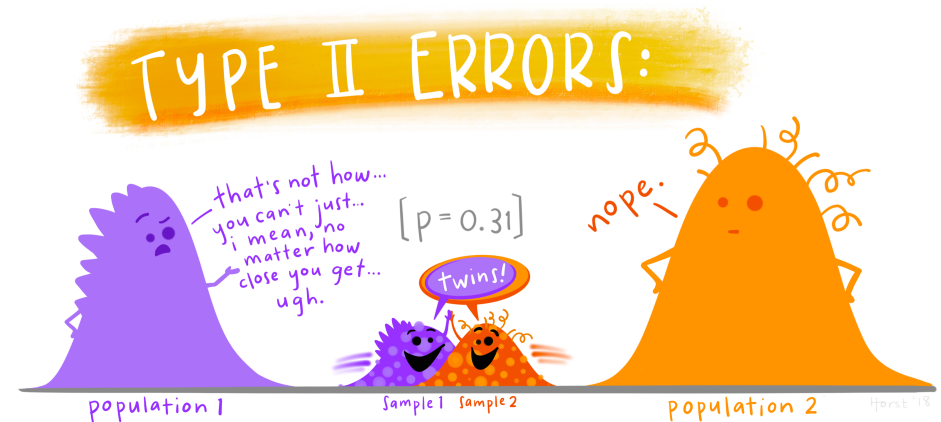
Which means:
IF THE PROBABILITY (p-value) OF FINDING AT LEAST OUR DIFFERENCE IN SAMPLE MEANS (IF THEY WERE DRAWN FROM POPULATIONS WITH THE SAME MEANS) IS LESS THAN 5%, THAT'S ENOUGH EVIDENCE FOR US TO DECIDE THEY ARE LIKELY FROM POPULATIONS WITH UNEQUAL MEANS.

TYPE 1 ERRORS

$[p = 0.02]$

OK, what about NOW??

(sigh...) YES. I'm STILL SURE I birthed BOTH of you.

Sample 1    (population)    Sample 2

TYPE II ERRORS:

that's not how... you can't just... i mean, no matter how close you get... ugh.

$[p = 0.31]$

nope.

twins!

population 1    Sample 1    Sample 2    population 2

Statistical significance ≠ Biological importance

|  | Important | Unimportant |
|---|---|---|
| Significant | Polio vaccine reduces incidence of polio | Things you don't care about, *or* already well known things:<br><br>BRIEFS<br>**Study Shows Frequent Sex Enhances Pregnancy Chances** |
| Insignificant | Small study shows a possible effect, leading to larger study which finds significance.<br>*or*<br>Large study showing no effect of drug that was thought to be beneficial. | Studies with small sample size and high *P*-value<br>*or*<br>Things you don't care about |

# Correlation does not automatically imply causation

# Life expectancy by country:



60

# Confounding variable

An unmeasured variable
that may be the cause of
both *X* and *Y*

Observations vs.
Experiments