# Quick overview of basic concepts in Statistics

(and a few bad jokes)

# Motivation

Statistical concepts can be subtle and
[difficult to explain](#)

I thought it would be useful to review some central ideas in Statistics

# Overview

Descriptive statistics

- Distributions, sampling
- Descriptive statistics and plots
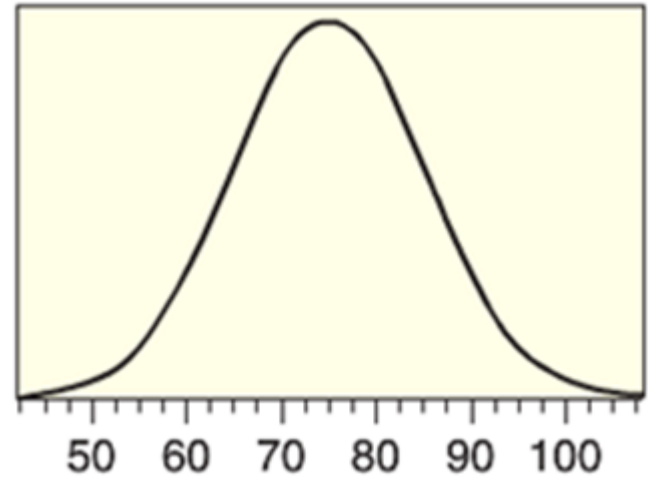- Sampling distributions

Inferential statistics

- Point and interval estimates
- Confidence intervals and the bootstrap
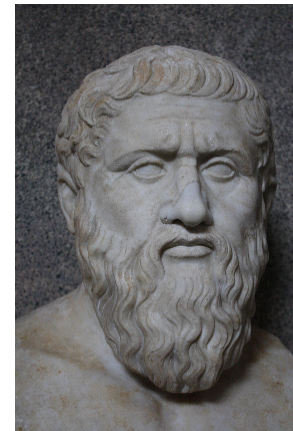- Hypothesis tests

If there is interest: neural data analysis

# Where do data come from?



**Distributions!**



Data ☹

Truth!

# How do we get data?



**Simple random sample**: each member in the population is equally likely to be in the sample

- Random selection

*Soup analogy!*



**Q:** Why is this good?

# An Example Dataset (flight delays)

Variables

Cases

| | flight | date | carrier | origin | dest | air_time | arr_delay |
|---|---|---|---|---|---|---|---|
| 1 | 1545 | 1-1-2013 | UA | EWR | IAH | 227 | 11 |
| 2 | 1714 | 1-1-2013 | UA | LGA | IAH | 227 | 20 |
| 3 | 1141 | 1-1-2013 | AA | JFK | MIA | 160 | 33 |
| 4 | 725 | 1-1-2013 | B6 | JFK | BQN | 183 | -18 |
| 5 | 461 | 1-1-2013 | DL | LGA | ATL | 116 | -25 |
| 6 | 1696 | 1-1-2013 | UA | EWR | ORD | 150 | 12 |
| 7 | 507 | 1-1-2013 | B6 | EWR | FLL | 158 | 19 |

# An Example Dataset (flight delays)

Categorical Variable

Quantitative Variable

Cases

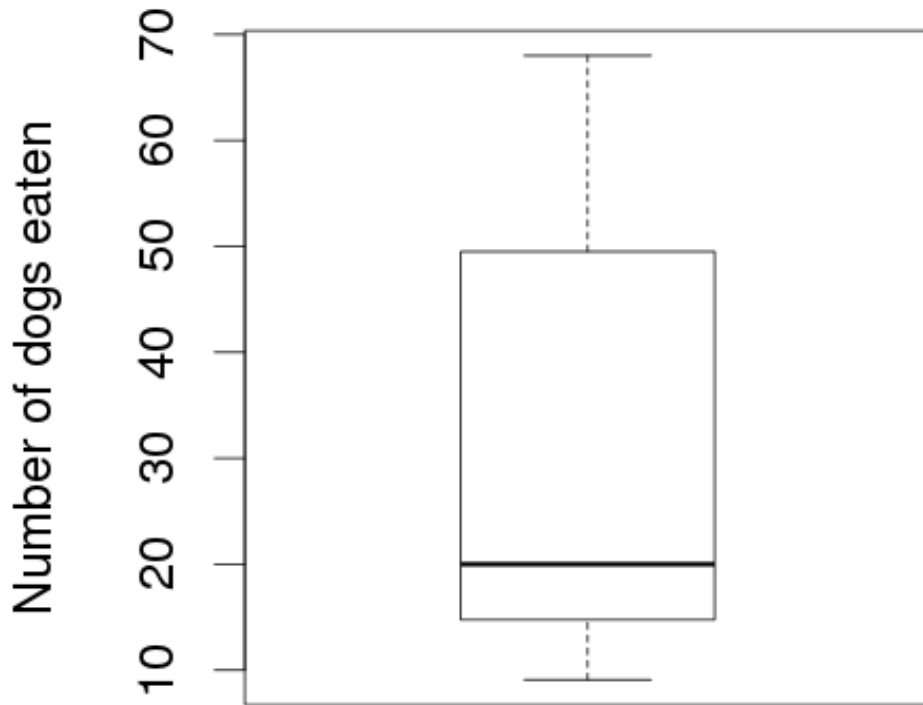| | flight | date | carrier | origin | dest | air_time | arr_delay |
|---|---|---|---|---|---|---|---|
| 1 | 1545 | 1-1-2013 | UA | EWR | IAH | 227 | 11 |
| 2 | 1714 | 1-1-2013 | UA | LGA | IAH | 227 | 20 |
| 3 | 1141 | 1-1-2013 | AA | JFK | MIA | 160 | 33 |
| 4 | 725 | 1-1-2013 | B6 | JFK | BQN | 183 | -18 |
| 5 | 461 | 1-1-2013 | DL | LGA | ATL | 116 | -25 |
| 6 | 1696 | 1-1-2013 | UA | EWR | ORD | 150 | 12 |
| 7 | 507 | 1-1-2013 | B6 | EWR | FLL | 158 | 19 |

# What is a good first step when analyzing data?

What is a useful plot
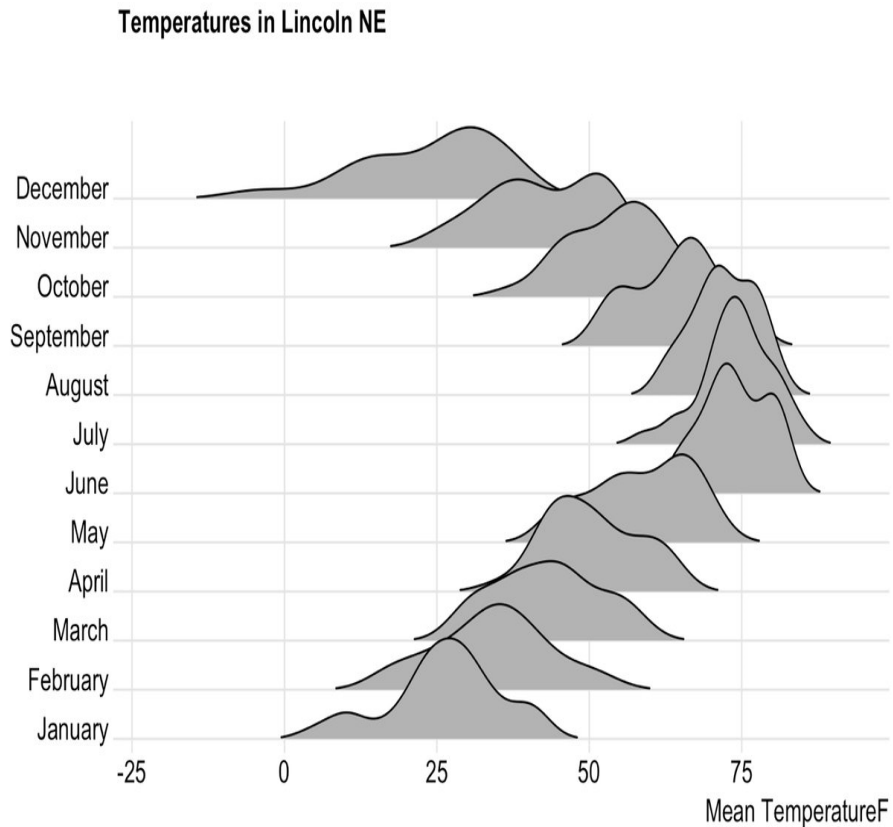for categorical data?

What is a useful plot
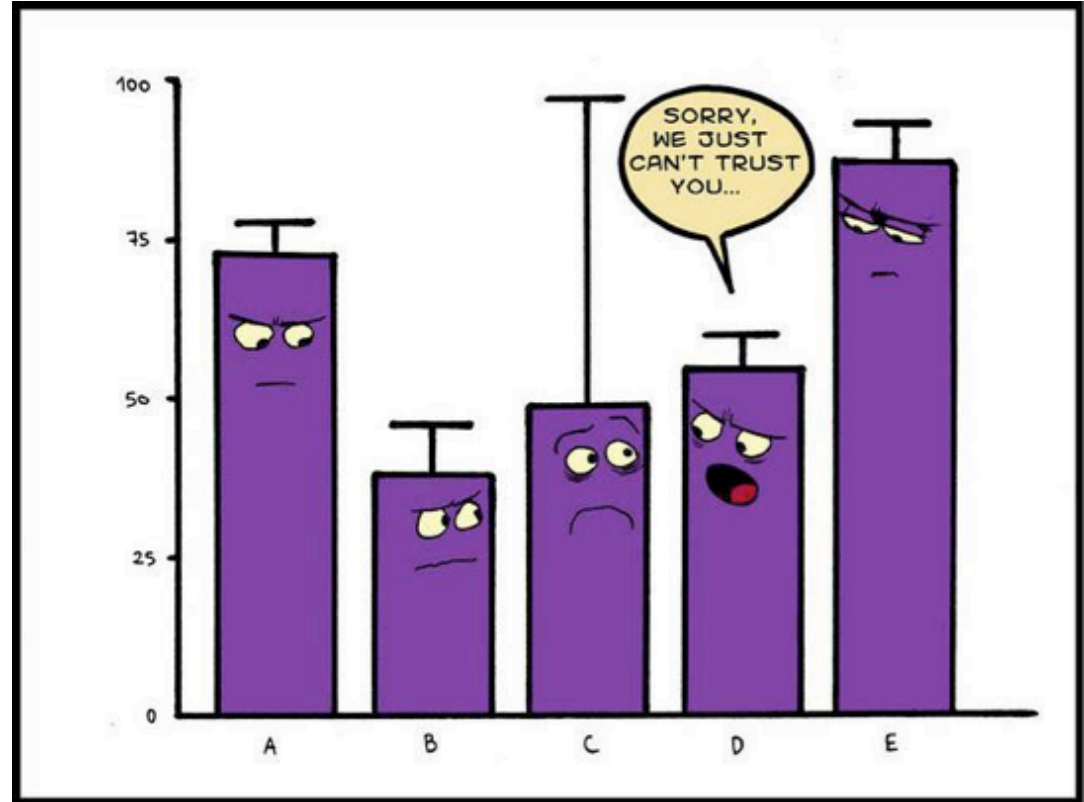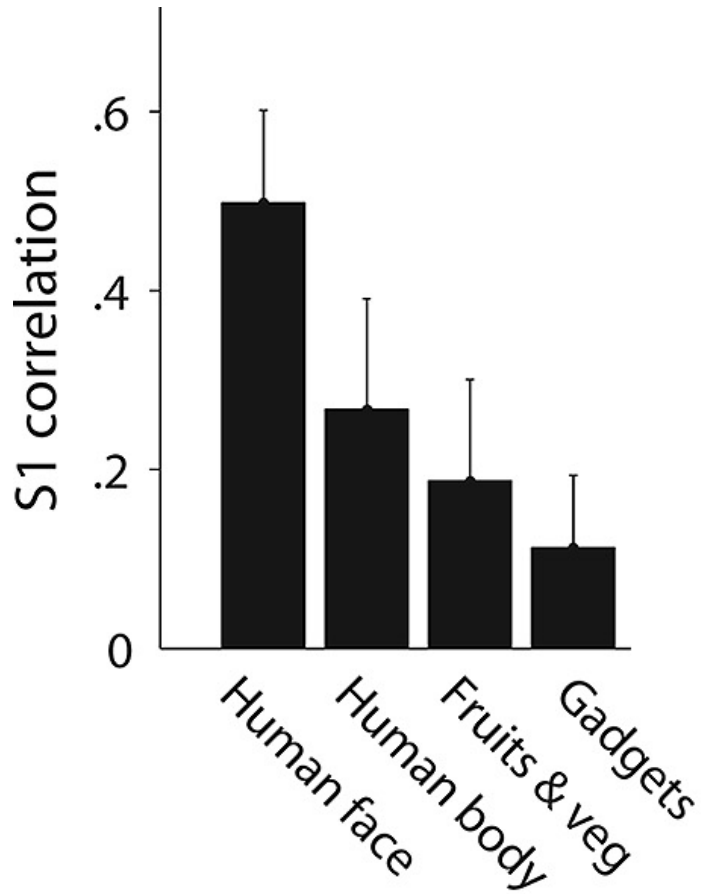for quantitative data?

# Box and violin plots

# Joy plots

Neat fact: Some statisticians find violin plots ugly
- Why do they find them ugly?

**Temperatures in Lincoln NE**

# Dynamite plots



Neat fact: Many statisticians hate dynamite plots
- Why do they hate them?

# What is a statistic?

**A**: A single death is a tragedy; a million deaths is a statistic.
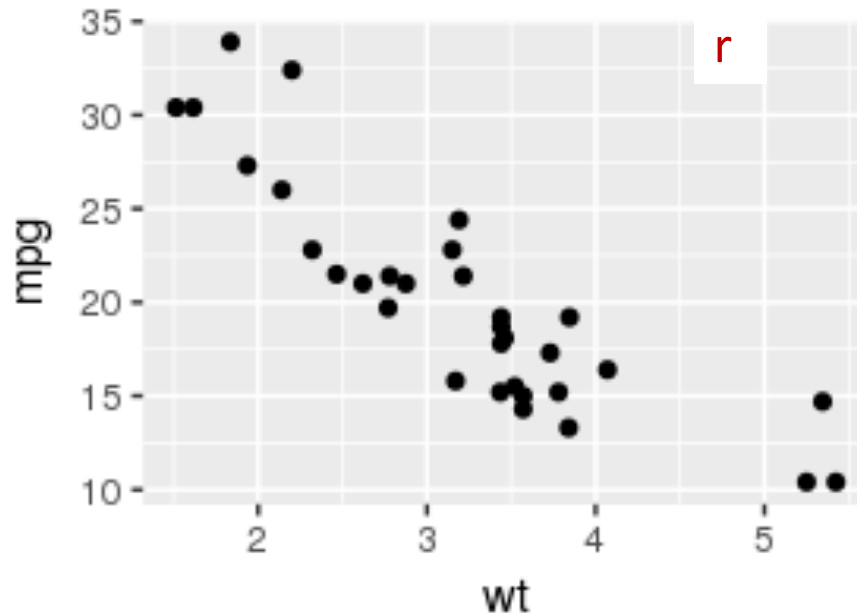
- Joseph Stallin

Descriptive statistics describe the data you have collected

- Usually denoted with roman characters
  - A single categorical variable: proportion
  - A single quantitative variable: the mean
  - A pair of quantitative variables: the correlation

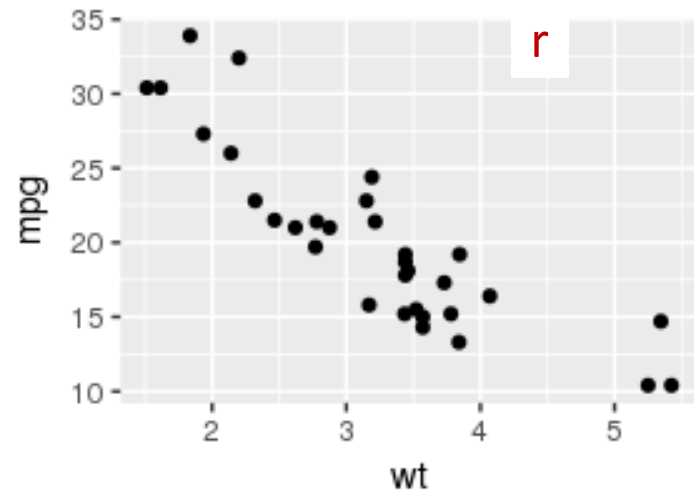# Example statistic: correlation coefficient

The correlation coefficient *r* is a statistic that measures the linear association between two variables X, and Y

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

r is an estimate of ρ
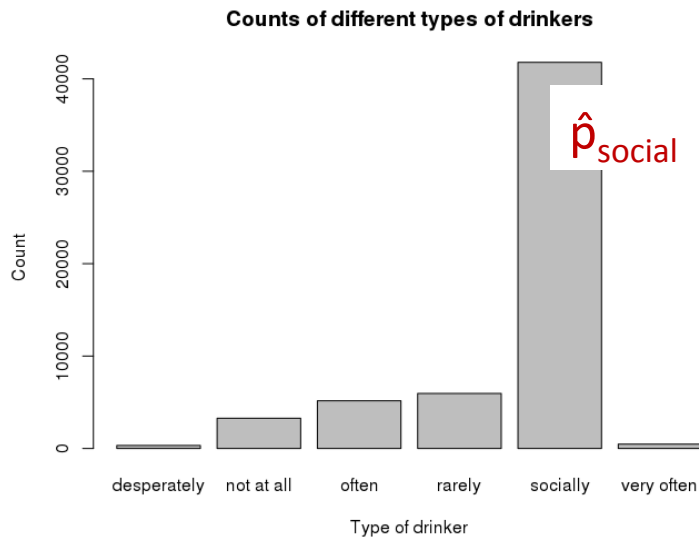
# Side note: correlation ≠ causation

The correlation coefficient $r$ is a statistic that measures the linear association between two variables X, and Y

# A variety of statistics than can be used to describe data

Descriptive **statistics** are usually denoted with roman characters

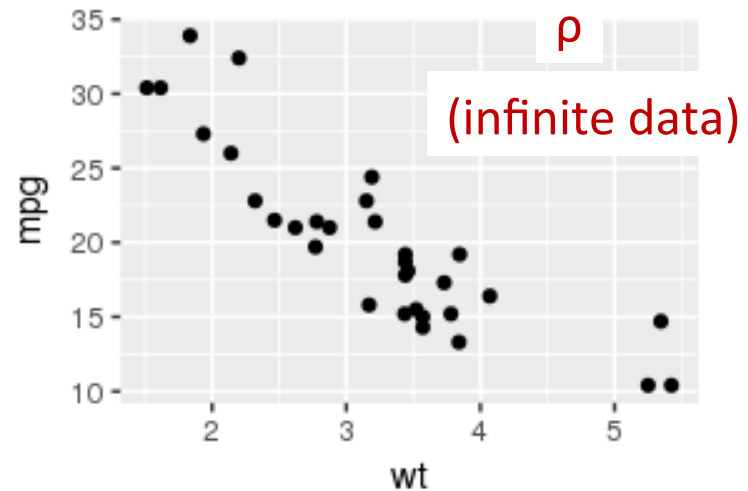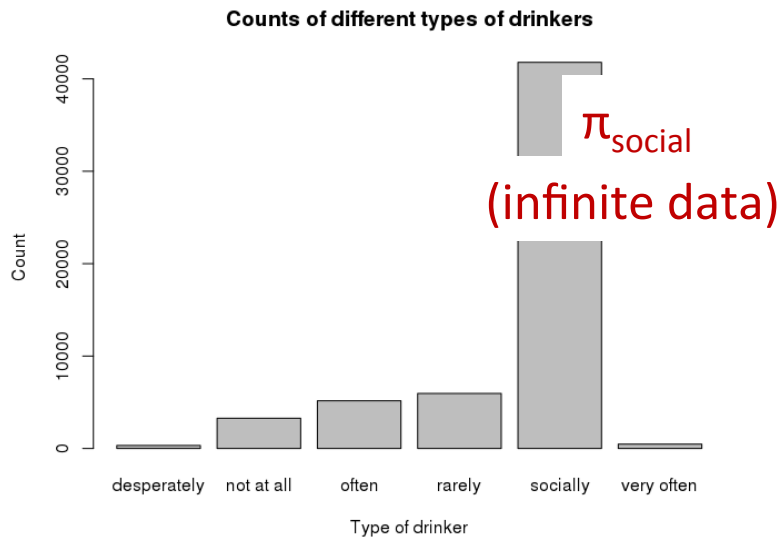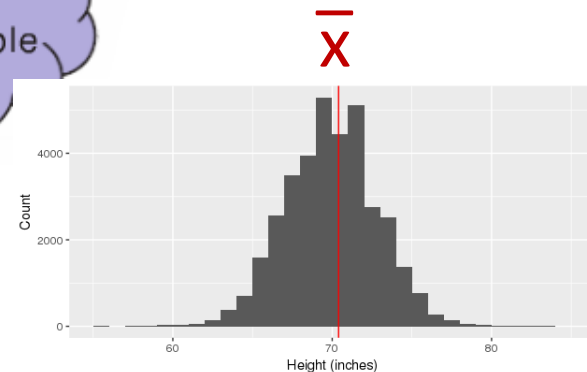- A single quantitative variable: the mean
- A single categorical variable: proportion
- A pair of quantitative variables: the correlation



Counts of different types of drinkers

$\hat{p}_{social}$

$r$

# Statistics have corresponding parameters

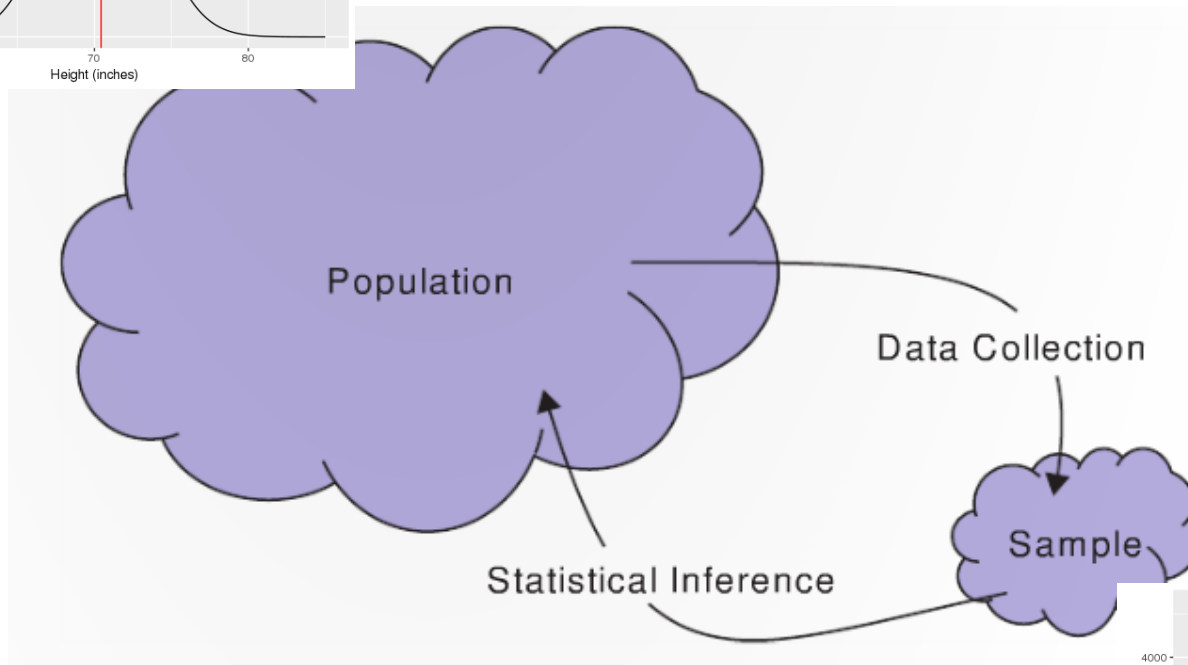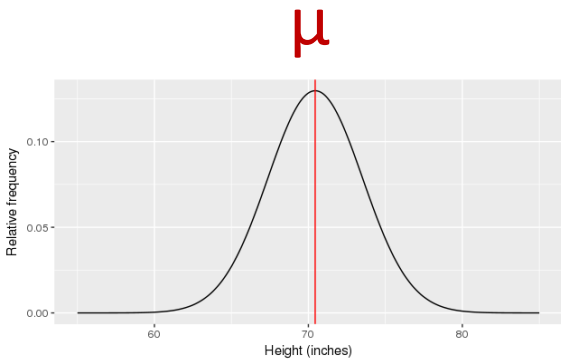Population/process ***parameters*** are usually denoted with greek characters
- A single quantitative variable: the mean
- A single categorical variable: proportion
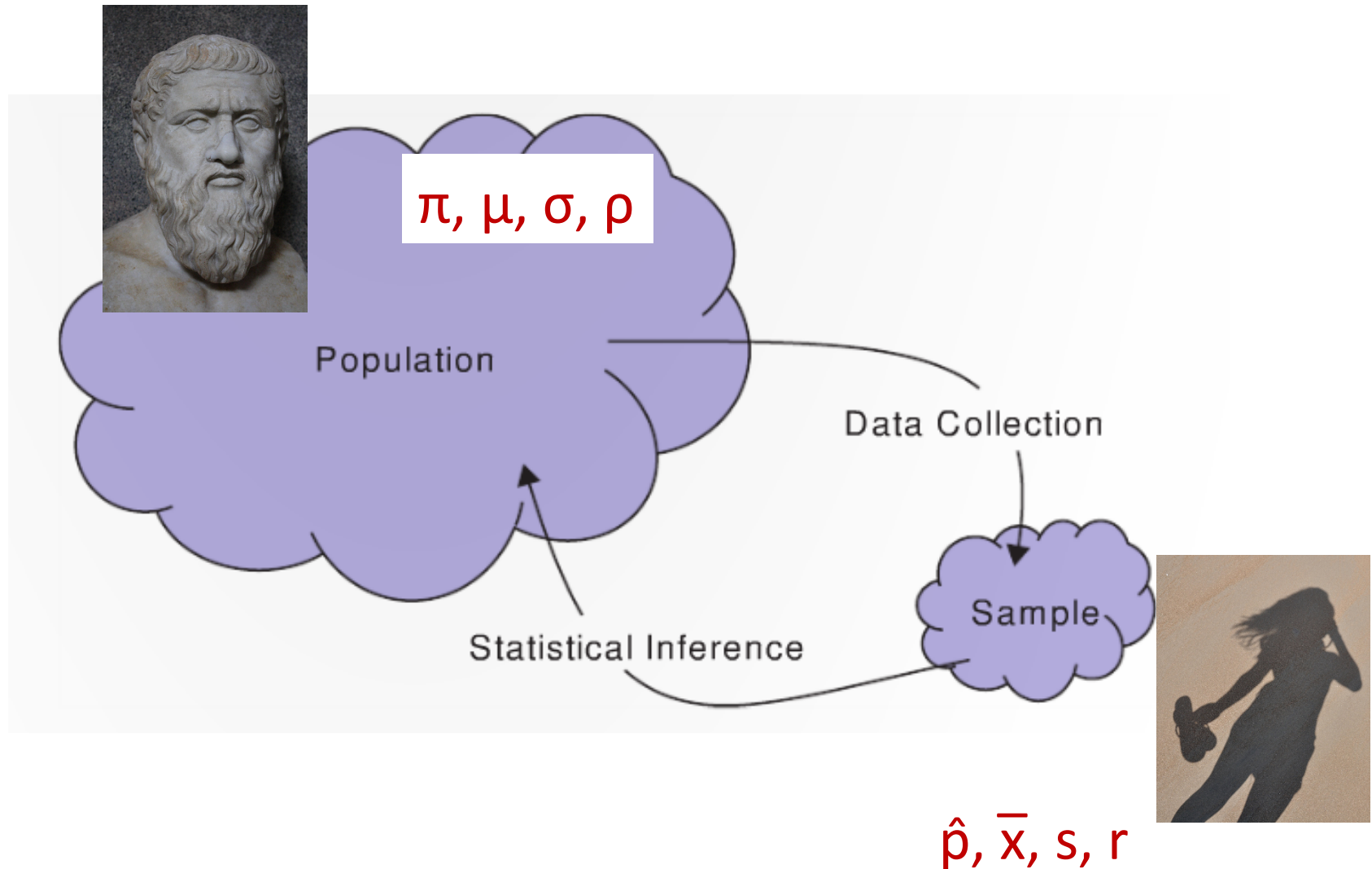- A pair of quantitative variables: the correlation



$\pi_{social}$

(infinite data)

$\rho$

(infinite data)

# Statistical inference

$\mu$

We use *sample statistics* to make judgements about *population parameters*

Population

Data Collection

Sample

Statistical Inference

$\bar{x}$

# Statistical inference
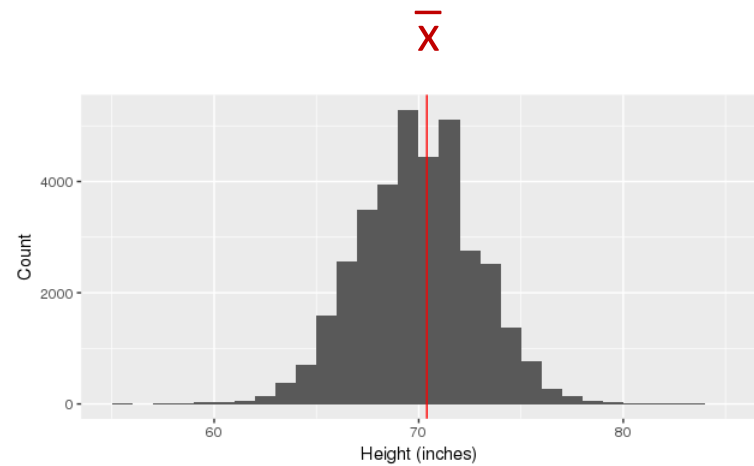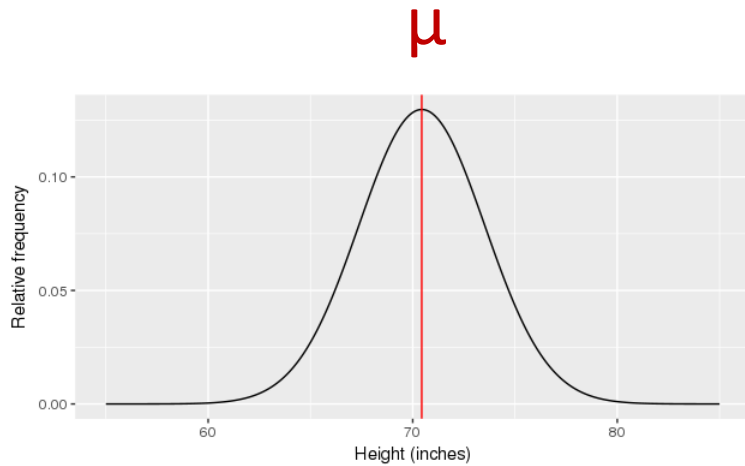


π, μ, σ, ρ

Population
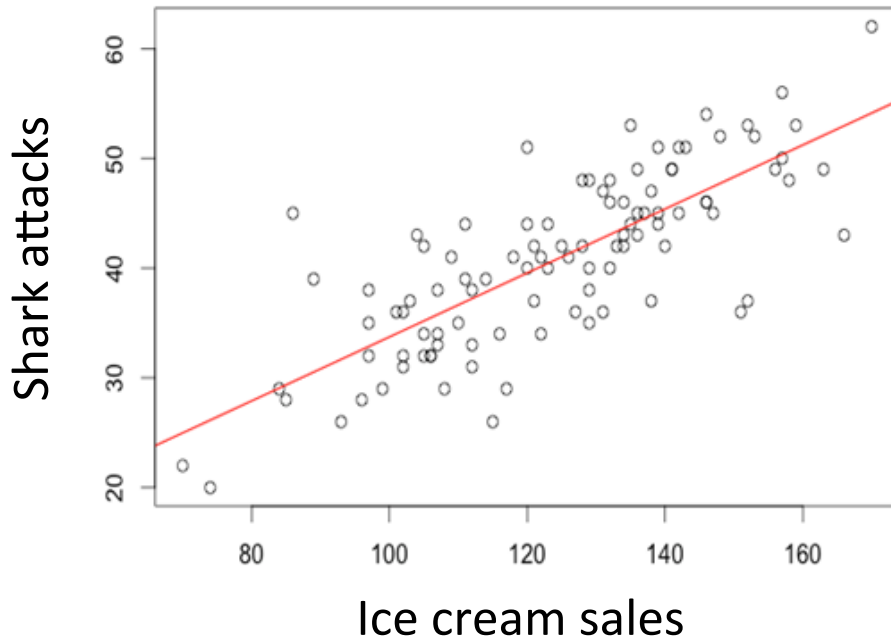
Data Collection

Sample

Statistical Inference

p̂, x̄, s, r

# Estimation

**Point estimation**: $\bar{x}$ is a point estimate for μ

# Regression

In linear regression, we try to fit a line to our data



Truth:  $y = \beta_0 + \beta_1 \cdot x$

Fit:  $\hat{y} = b_0 + b_1 \cdot x$

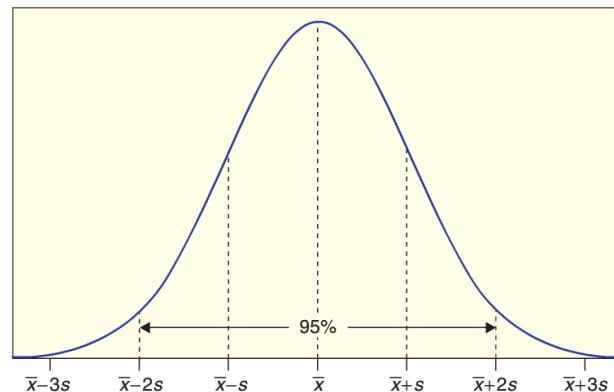The $b_i$'s are usually chosen to minimize the MSE:

$$MSE = \frac{1}{n} \sum_i^n (\hat{y}_i - y_i)^2$$

# Regression caution # 1

Avoid trying to apply the regression line to predict values far from those that were used to create the line

# Statistics have distributions too!

A **sampling distribution** is the distribution of sample statistics computed for different samples of the same size from the same population $\overline{x}_1, \overline{x}_2, \overline{x}_3, ..., \overline{x}_k$

95%

$\overline{x}-3s$  $\overline{x}-2s$  $\overline{x}-s$  $\overline{x}$  $\overline{x}+s$  $\overline{x}+2s$  $\overline{x}+3s$

Sampling distribution are often approximately Normal
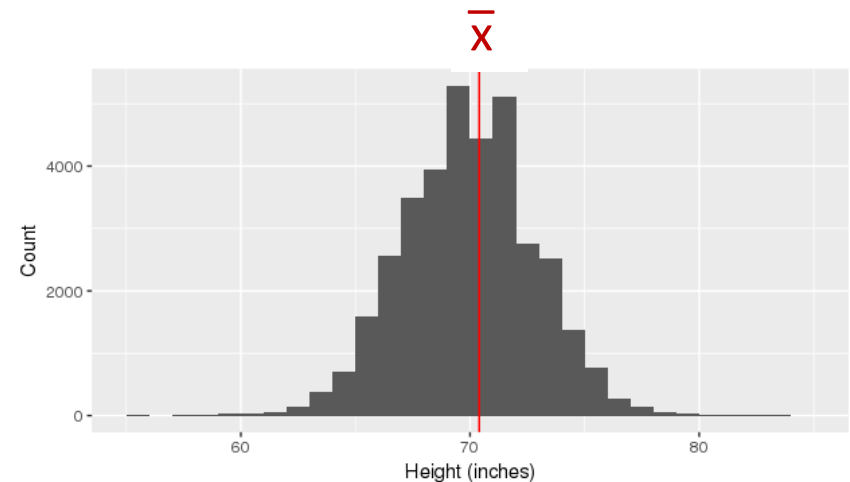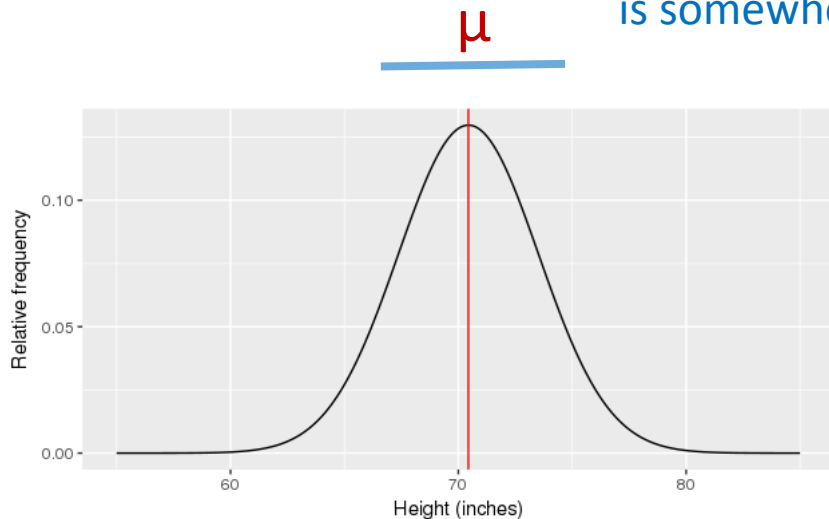- Due to the central limit theorem: sums of many random variates lead to a Normal distribution

# Estimation

**Point estimation**:  $\bar{x}$  is a point estimate for μ

**Interval estimate**: point estimate + margin of error

The population mean
is somewhere in here

μ

$\bar{x}$

# Confidence Intervals

A **confidence interval** *for a parameter* is an interval computed by ***a method that will capture the parameter a specified proportion of times***

- (if the sampling process were repeated many times)

The success rate (proportion of all samples whose intervals contain the parameter) is known as the **confidence level**
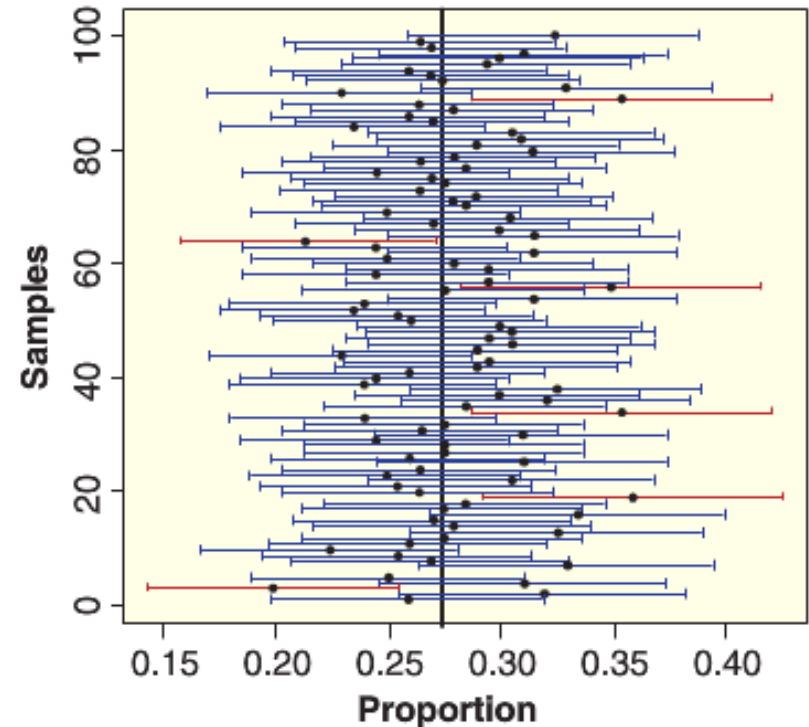
# Think ring toss…



Parameter exists in the ideal world

We toss intervals at it

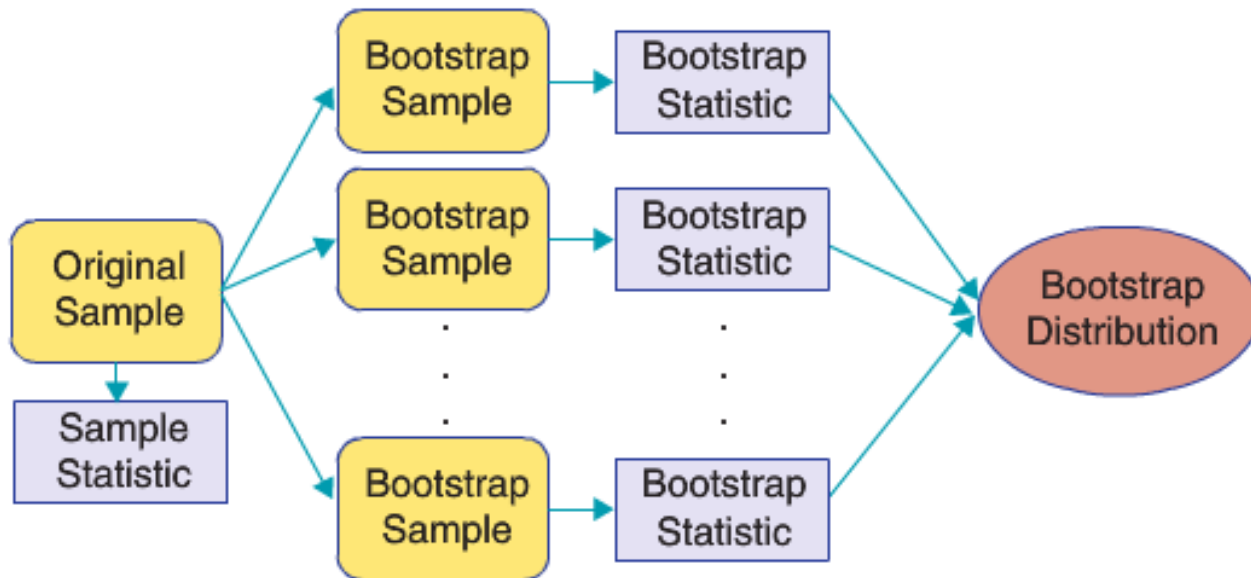95% of those intervals capture the parameter
(for a 95% CI)

Wits and Wagers…

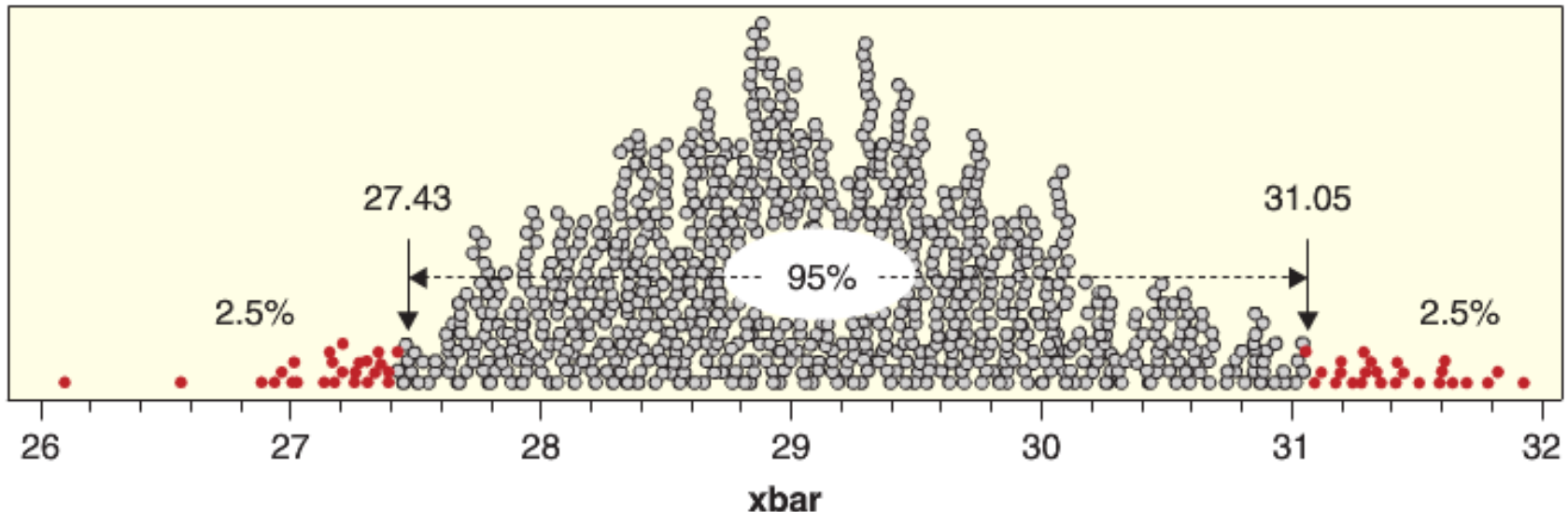# The bootstrap to estimate standard errors

## The plug-in principle

# 95% Confidence Intervals

If the bootstrap distribution is ~symmetric we can use percentiles to calculate a 95% confidence interval
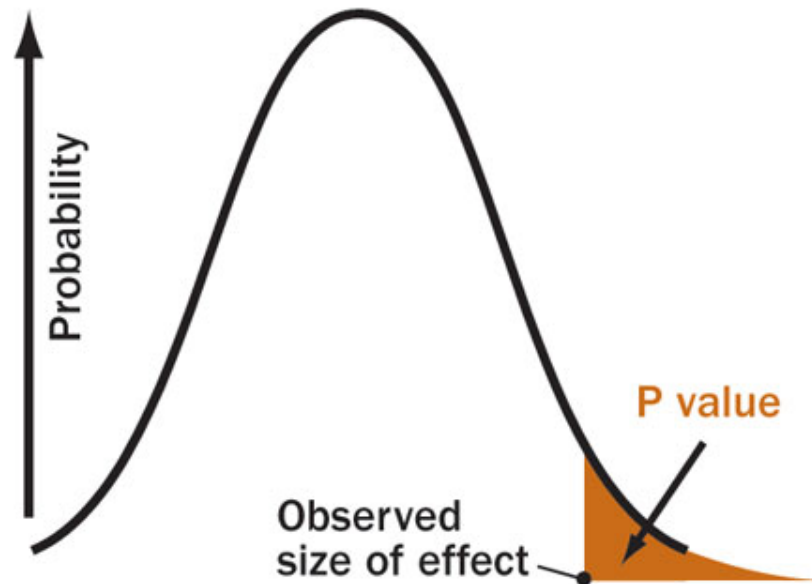


Formulas also exist for calculating CIs
- SE of the mean can be estimated as:
- *CI is:   Statistic ± 2 · SE*

$$SE = \frac{s}{\sqrt{n}}$$

# Question: why is a p-value?

A p-value is the probability of getting a statistic as or more extreme than the observed statistic from a null model (null distribution)
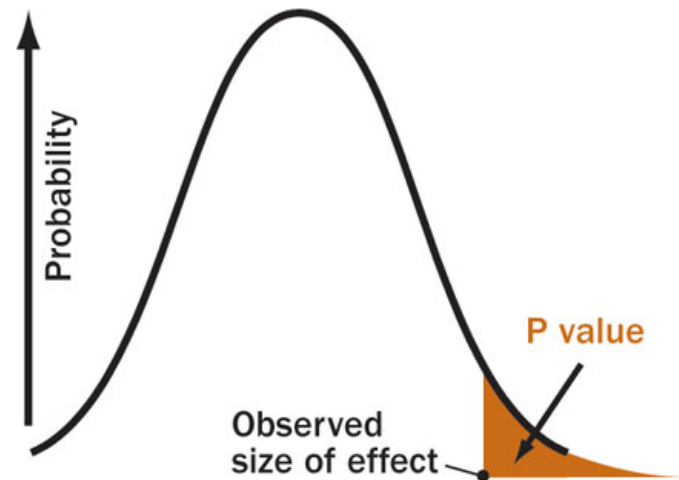
# Hypothesis tests in 2 steps

1. Create a distribution showing what the statistic would like like if there was no effect ($H_0$ = true)
   - Null distribution   (distribution if no effect)


2. Show that the statistic you observed is unlikely to come from this null distribution
   - P-value is:   $Pr(S \geq obs\_stat \mid H_0)$
   - P-value is **<u>not</u>**:   $Pr(H_0 \mid data)$

# Trial metaphor of hypothesis testing

1. Assume innocence: $H_0$ is true
   - State $H_0$ and $H_A$

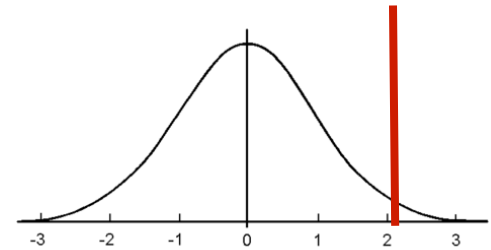2. Gather evidence
   - Calculate the observed statistic

3. Create a distribution of what evidence would look like if $H_0$ is true
   - Null distribution

4. Assess the probability that the observed
   evidence would come from the null distribution
   - p-value

5. Make a judgement
   - Assess whether the results are statistically significant

# Types of hypothesis tests test

**Permutation test**: Create null distribution simulating many random assignments

- 1. Shuffle the condition labels
- 2. Compute statistic on shuffled labels
- 3. Repeat many times to get a null distribution

**Parametric tests**: based on the fact that distribution of statistic is known

- t-tests, ANOVAs, chi-squared tests, etc...
- Less robust (based on Normal approximations)

# Visual hypothesis tests

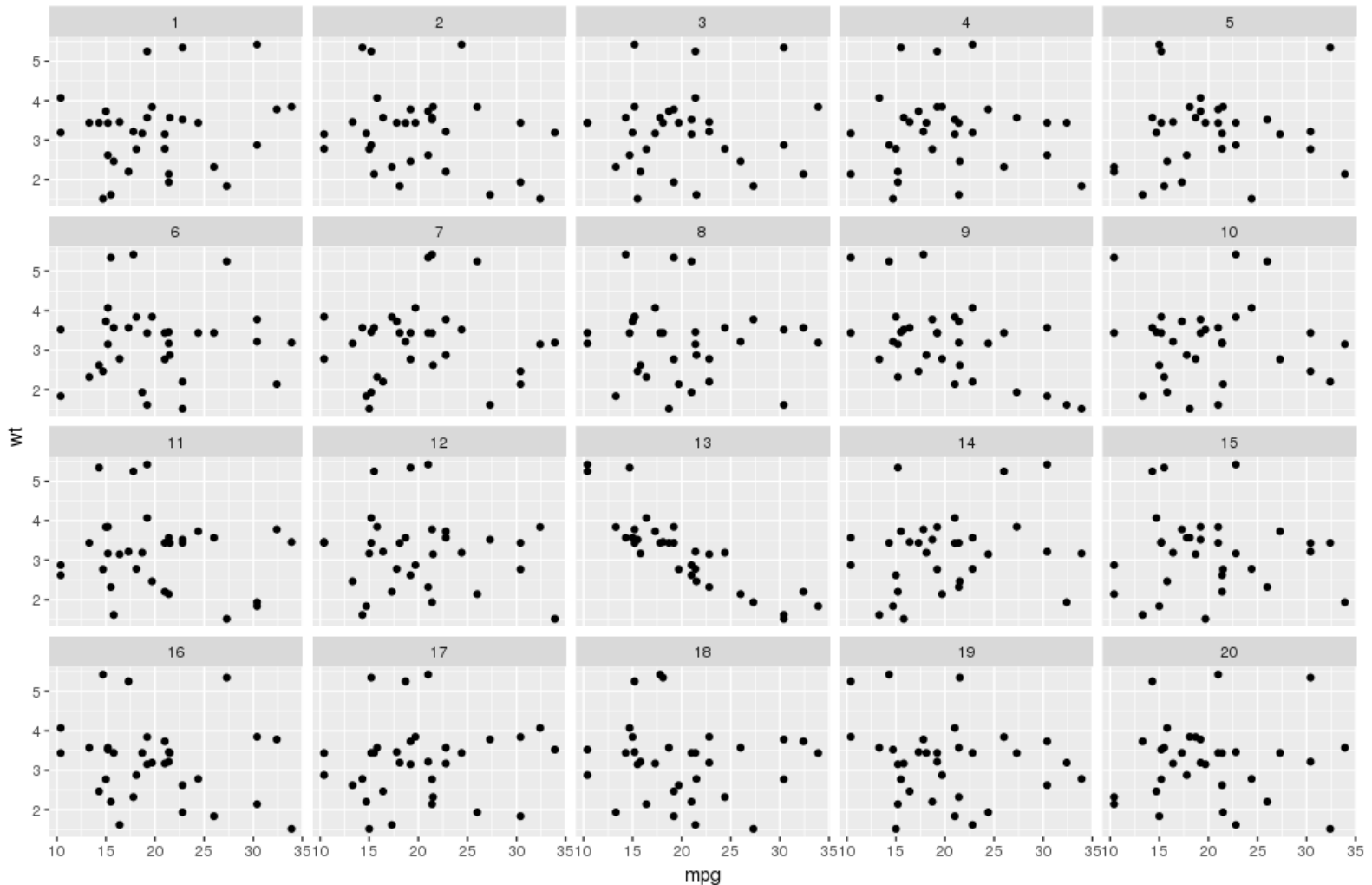"Visual hypothesis tests" create visual plots of randomly rearranged points

- i.e., a set of 'innocent' plots where there is no pattern

If you can tell the real data from the plots of randomly generated data this is equivalent to rejecting the null hypothesis

# Visual hypothesis tests

## Which plot shows the real data?
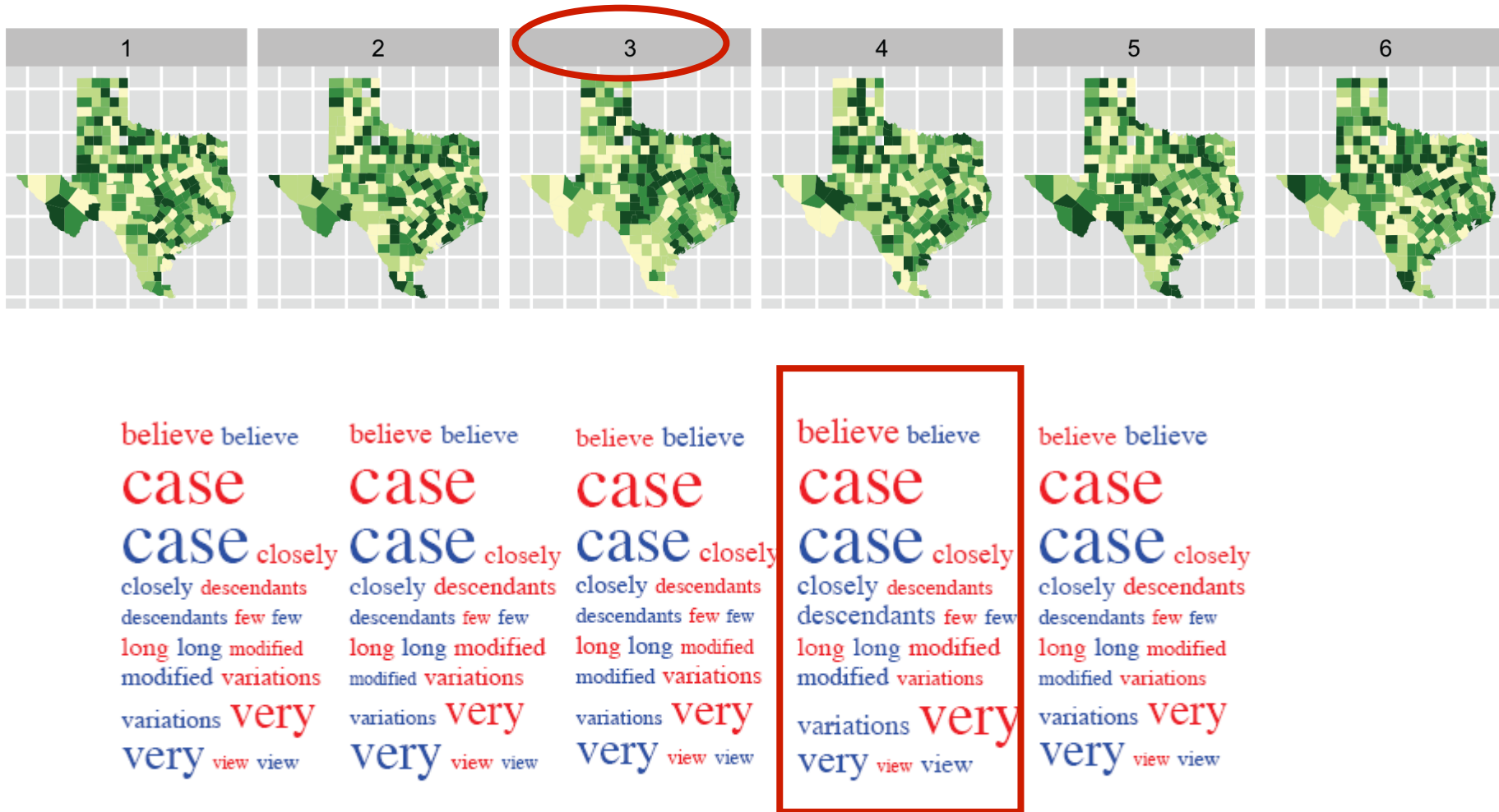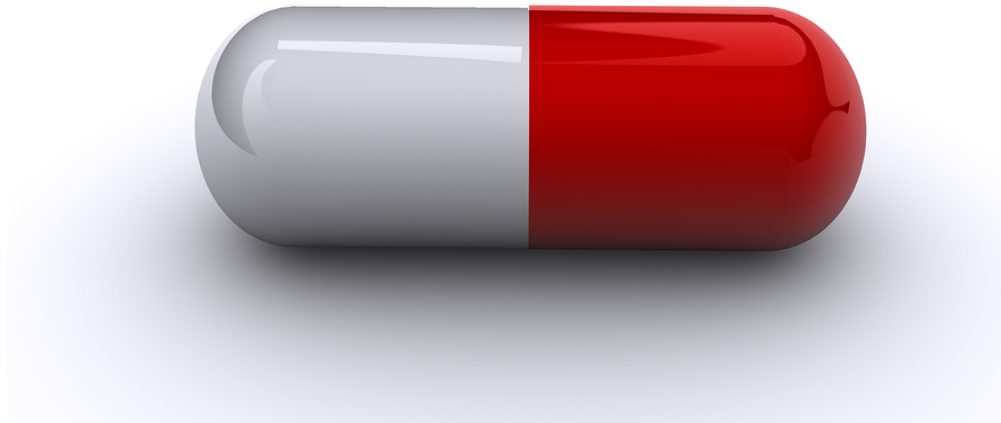
# Visual hypothesis tests



Fig. 5. Five tag clouds of selected words from the 1st (red) and 6th (blue) editions of Darwin's "Origin of Species". Four of the tag clouds were generated under the null hypothesis of no difference between editions, and one is the true data. Can you spot it?

See Wickham, et al, 2010

# Permutation test example: Hypothesis tests for comparing two means
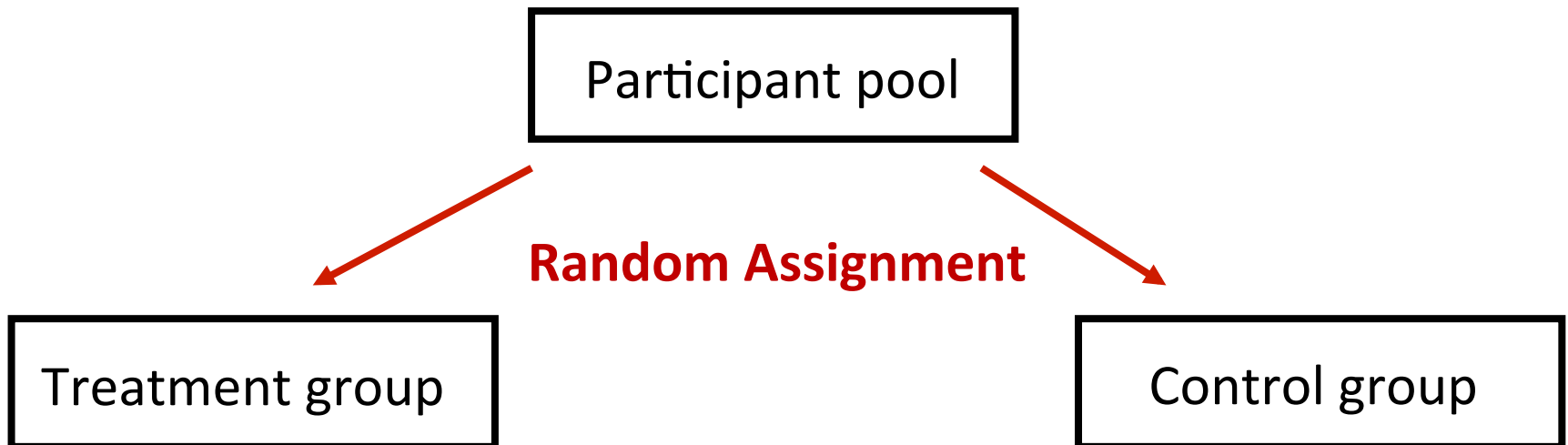


**Question**: Is this pill effective?

# Experimental design

Take a group of participant and **_randomly assign_**:

- Half to a _treatment group_ where they get the pill

- Half in a _control group_ where they get a fake pill (placebo)

- See if there is more improvement in the treatment group compared to the control group

```
                    ┌─────────────────────┐
                    │   Participant pool   │
                    └─────────────────────┘
                      ↙                 ↘
         Random Assignment
┌─────────────────────┐        ┌─────────────────────┐
│  Treatment group    │        │   Control group      │
└─────────────────────┘        └─────────────────────┘
```

# Hypothesis tests for differences in two group means

1. State the null and alternative hypothesis
   - $H_0$: $\mu_{Treatment} = \mu_{Control}$    or    $\mu_{Treatment} - \mu_{Control} = 0$
   - $H_A$: $\mu_{Treatment} > \mu_{Control}$    or    $\mu_{Treatment} - \mu_{Control} > 0$

2. Calculate statistic of interest
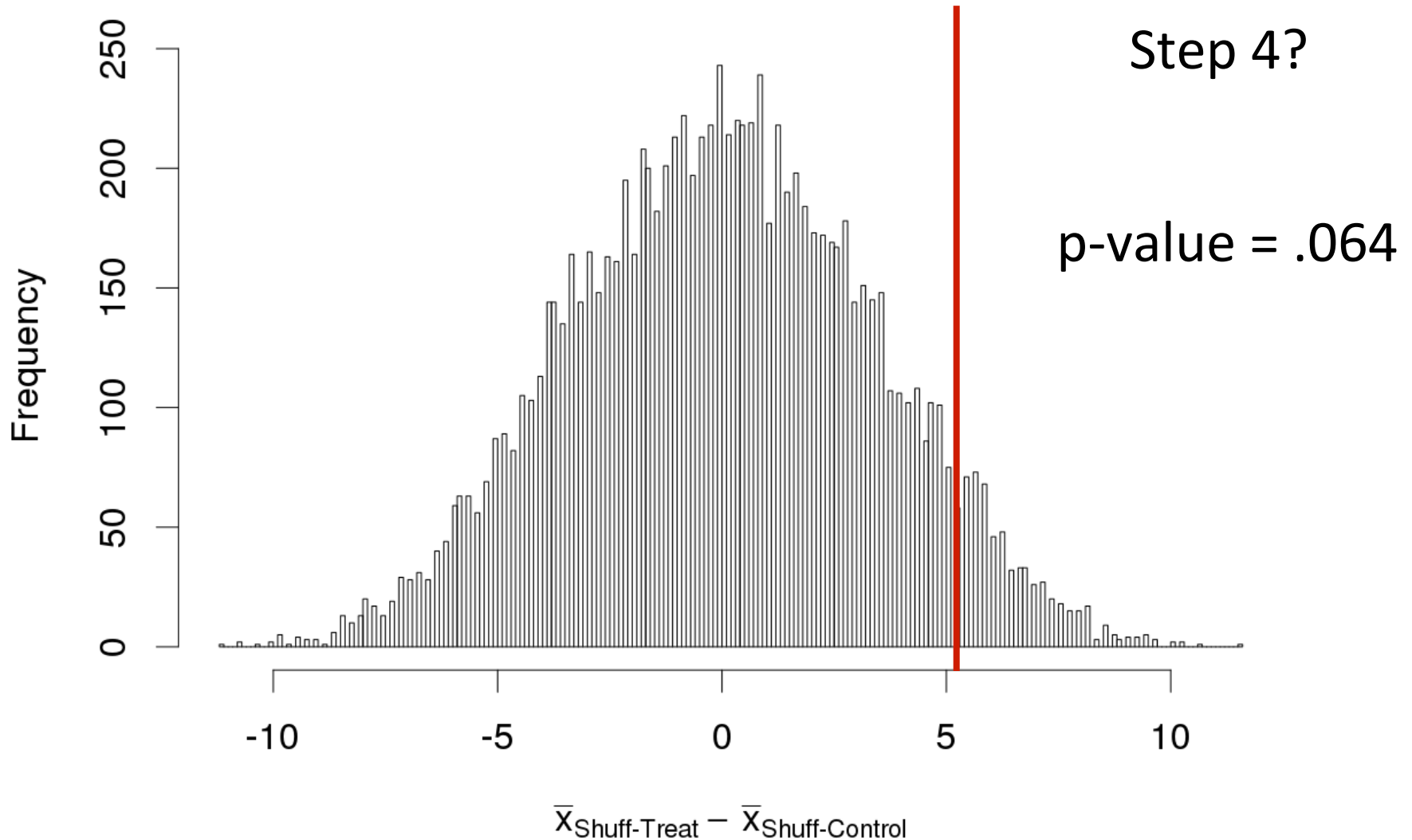   - $\bar{x}_{Treatment} - \bar{x}_{Control}$

3. What should we do next…?

# 3. Create the null distribution!

| Treatment group | Control group |
|---|---|

Reconstructed participant pool data under $H_0$

**Shuffle data for random assignment consistent with $H_0$**

| Shuffled 'treatment group' | Shuffled 'control group' |
|---|---|

One null distribution statistic: $\bar{x}_{Shuff\_Treatment} - \bar{x}_{Shuff\_control}$

# Repeat 10,000 times



**Null distribution**

Step 4?

p-value = .064

$\overline{x}_{Shuff\text{-}Treat} - \overline{x}_{Shuff\text{-}Control}$

# Fisher vs. Neyman

**Question:** should you report exact p-values?
- p < .05
- p = .021

Null-hypothesis significance testing (NHST) is a hybrid of two theories:
- 1. Significance testing of Ronald Fisher
- 2. Hypothesis testing of Jezy Neyman and Egon Pearson
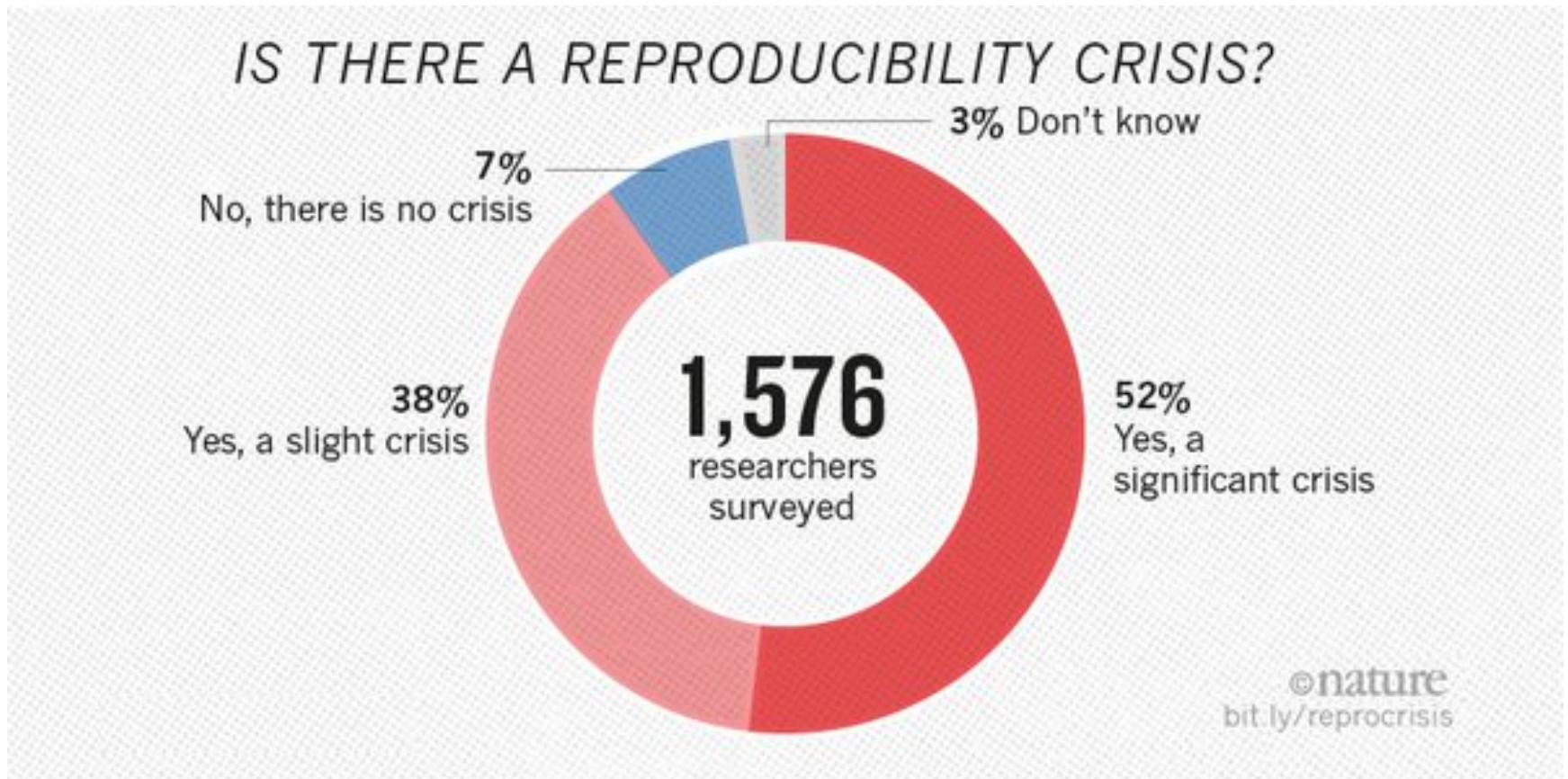
Neat fact: they hated each other

Fisher (1890-1962)        Neyman (1894-1981)

# Type I and Type II Errors



|  | Reject $H_0$ | Do not reject $H_0$ |
|---|---|---|
| $H_0$ is true | Type I error ($\alpha$) (false positive) | No error |
| $H_0$ is false | No error | Type II error ($\beta$) (false negative) |

# Why Most Published Research Findings Are False

John P. A. Ioannidis



IS THERE A REPRODUCIBILITY CRISIS?

3% Don't know

7%
No, there is no crisis

38%
Yes, a slight crisis

1,576
researchers
surveyed

52%
Yes, a
significant crisis

©nature
bit.ly/reprocrisis

# Do reproducible research!

Tools exist that allow you to embed your analyses into your written reports:

- R: R Markdown

- Python/Julia/R: Jupyter

- MATLAB: Live Scripts

# John Tukey

"Far better an approximate answer to the *right* question, which is often vague, than an *exact* answer to the wrong question"

- *The future of data analysis*. Annals of Mathematical Statistics 33 (1), (1962), page 13.
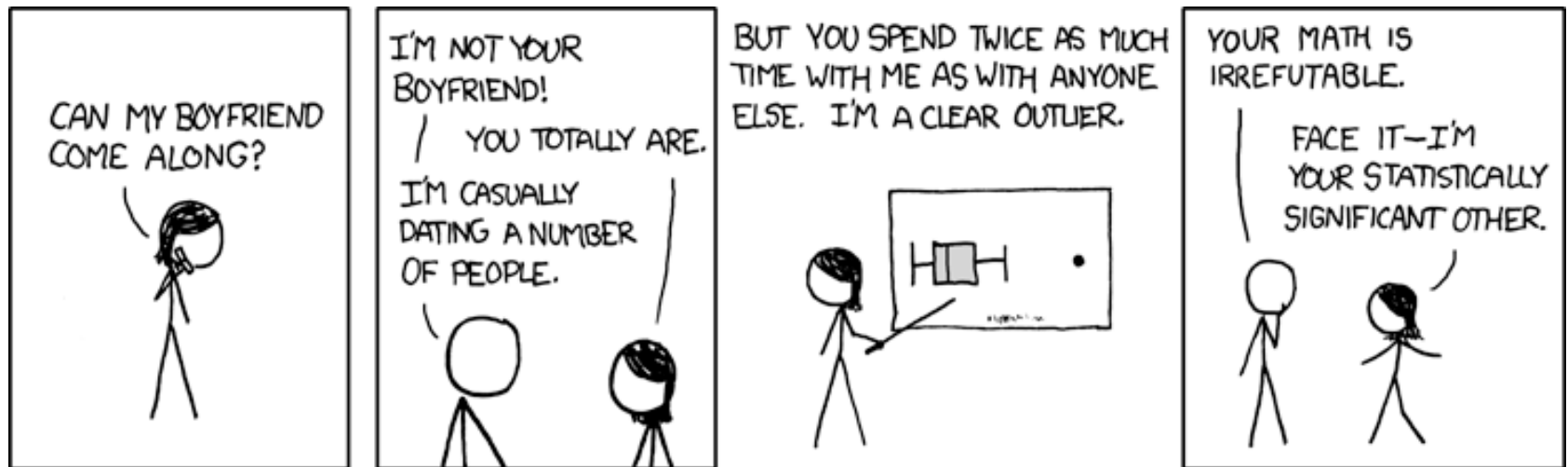
# What is Data Science?



See Donoho, 2017

# Questions

# Interactive single neuron analysis tutorial…

https://neuraldata.net/spike/

Created by Brooke FItzgerald

# Are we feeling ok about hypothesis tests?

# Why Most Published Research Findings Are False

John P. A. Ioannidis

## The Earth Is Round ($p < .05$)

Jacob Cohen

*After 4 decades of severe criticism, the ritual of null hypothesis significance testing—mechanical dichotomous decisions around a sacred .05 criterion—still persists. This article reviews the problems with this practice, including* sure how to test $H_0$, chi-square with Yates's (1951) correction or the Fisher exact test, and wonders whether he has enough power. Would you believe it? And would you believe that if he tried to publish this result without a

American Statistical Association's Statement on p-values