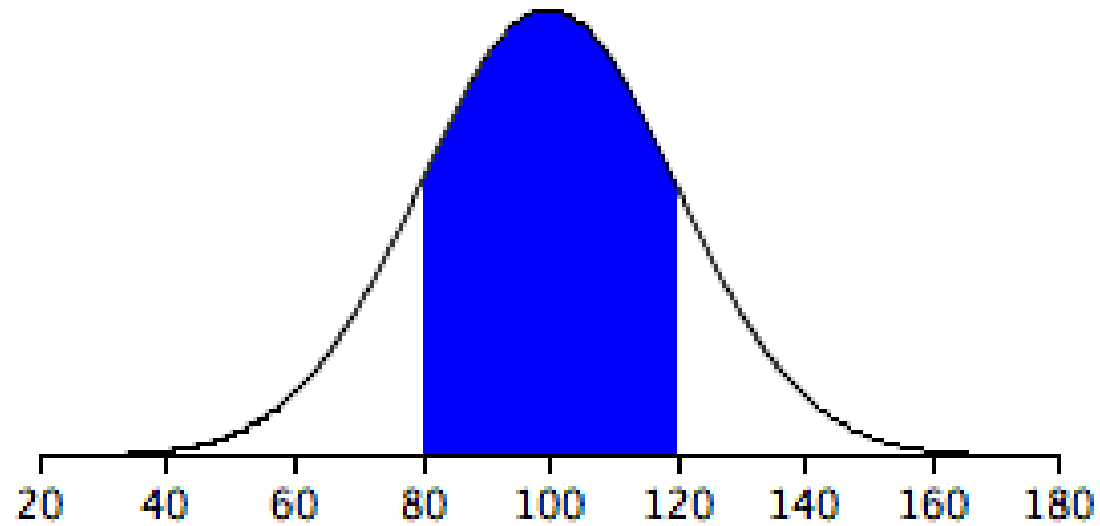


Sampling distributions



Overview

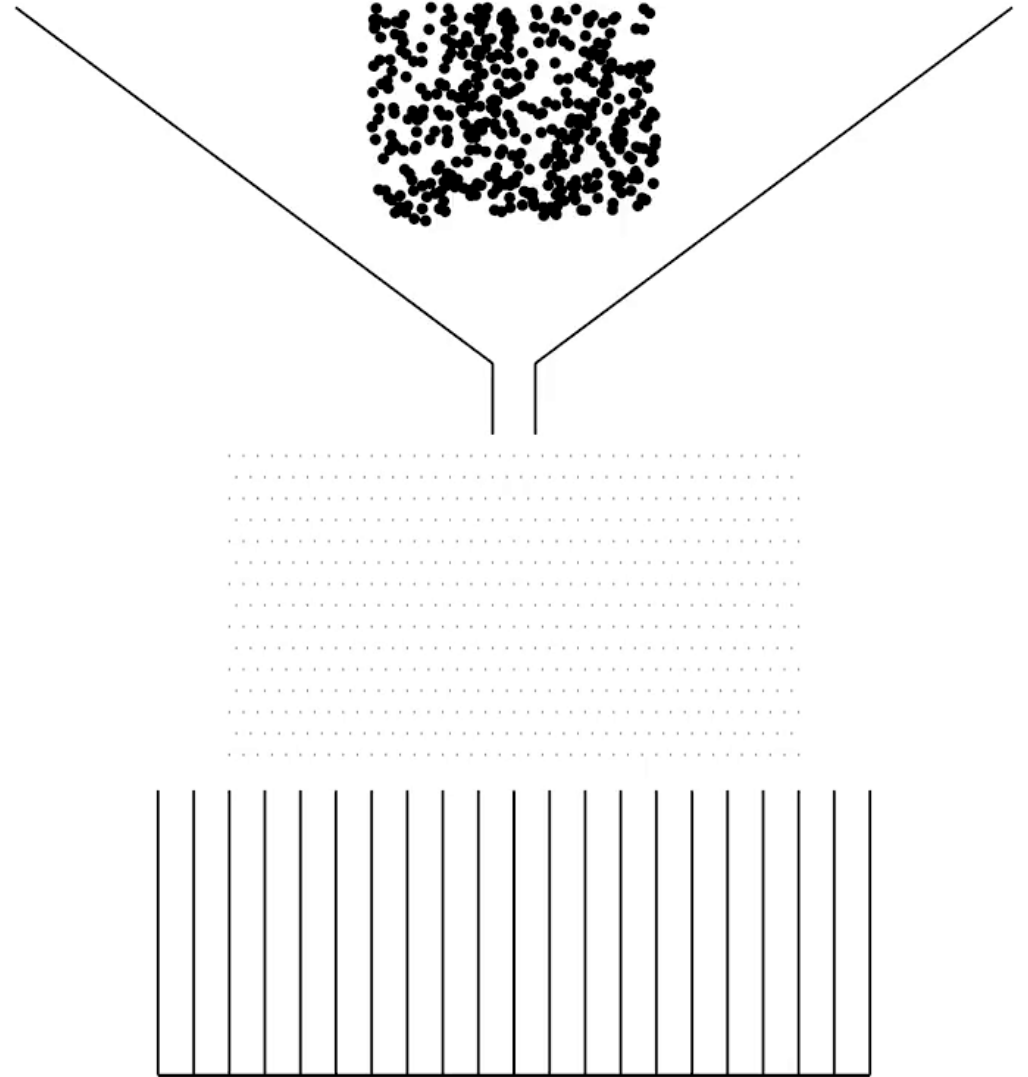
Very quick review

For loops

Generating random numbers and
selecting random samples

Sampling distributions

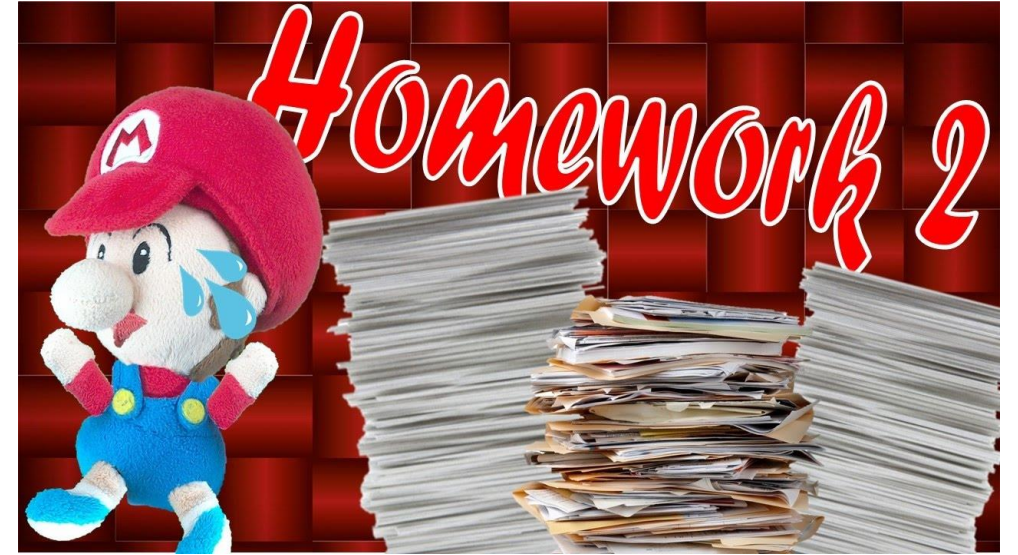
If there is time: confidence intervals



Announcements

Homework 2 has been posted

- Due Sunday (9/17) at 11pm
- Start early on it!
 - You can do problems 1 and 2 after today's class
- How was homework 1?



Where we are in the plan for the semester

			<u>Analysis</u>	<u>R</u>
1	Sep 2	Course overview, introduction to R, descriptive statistics		
2	Sep 7-9	Review of central statistical concepts and exploratory analysis using R		
3	Sep 14-16	Confidence Intervals and the bootstrap		
4	Sep 21-23	Review of hypothesis tests and permutation tests in R	resampling methods	
5	Sep 28-30	Parametric, non-parametric and theories of hypothesis testing		base R

We will be using some simulations to justify and validate methods we use throughout the semester



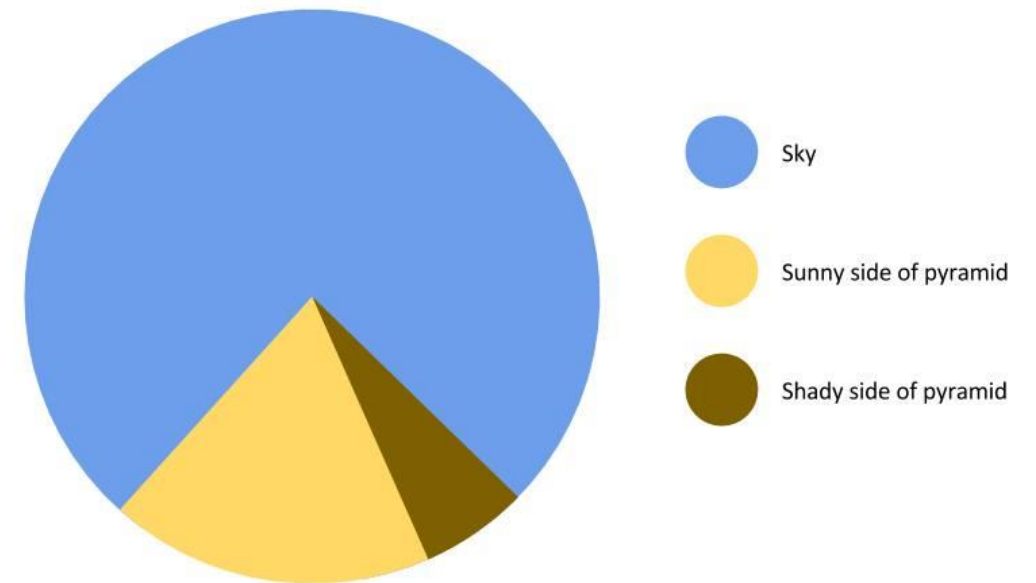
Quick review

Basics of R

```
> my_vec <- c(5, 28, 19)  
> my_vec[3]  
> my_vec[3] <- 7
```

How to plot categorical data

```
> drinks_table <- table(profiles$drinks)  
> barplot(drinks_table)  
> pie(drinks_table)
```

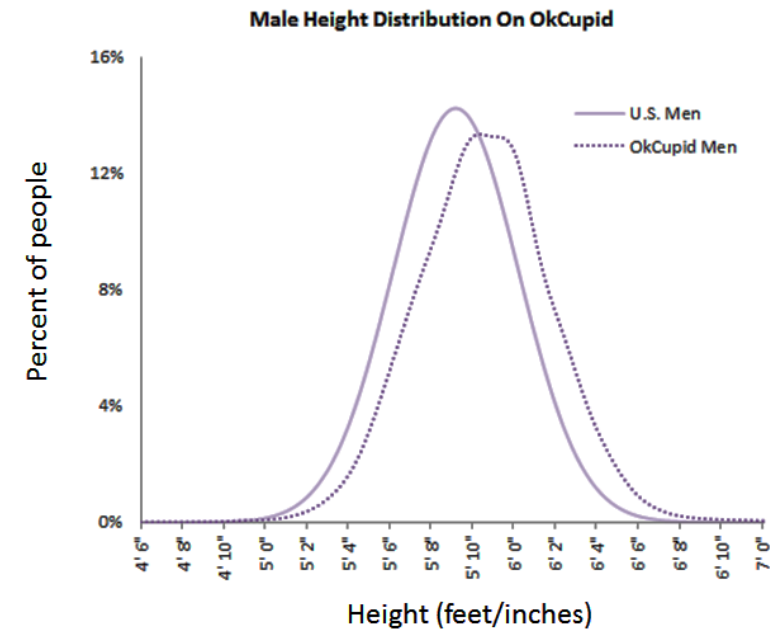
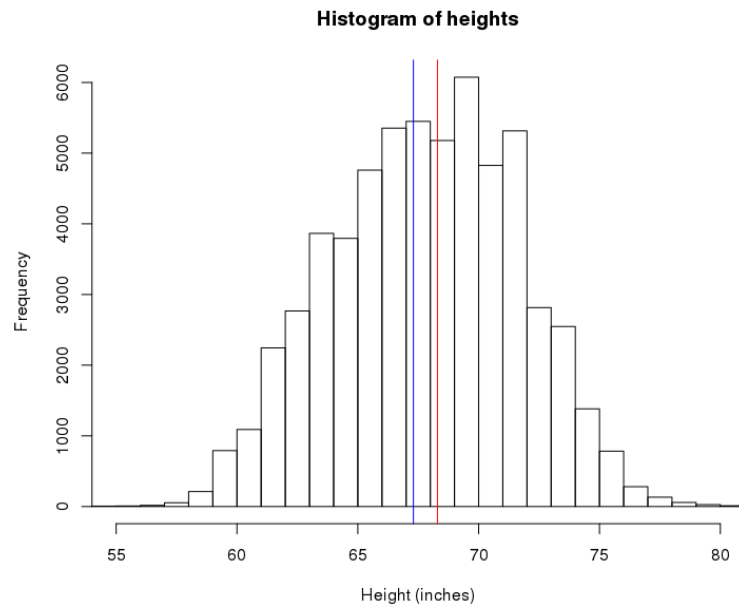


Quick review

How to plot quantitative data:

> `hist(profiles$height)`

> `abline(v = 67)`



Staying organized

It is useful to create separate folders for different homework and even for the different pieces of class code.

Be sure to set your working directory properly so that R can find the relevant files.



A little more R...

Things that
begin with

Rr



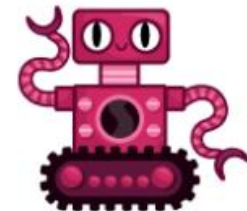
rabbit



rocket



rain



robot



ribbon



rat

For loops

For loops are useful when you want to repeat a piece of code many times under similar conditions

The syntax for a for loop is:

```
for (i in 1:100) {
```

```
    # do something
```

```
}
```



This is repeated 100 times
i is incremented by 1 each time

For loops

For loops are particularly useful in conjunction with vectors...

```
my_results <- NULL    # create an empty vector to store the results
for (i in 1:100) {
  my_results[i] <- i^2
}
```

Try this at home!: Use a for loop to create a vector that holds the values at multiples of 3 from 3 to 300

- i.e., 3, 6, 9, ..., 300

Let's try it in R!

Generating random data

R has built in functions to generate data from different distributions

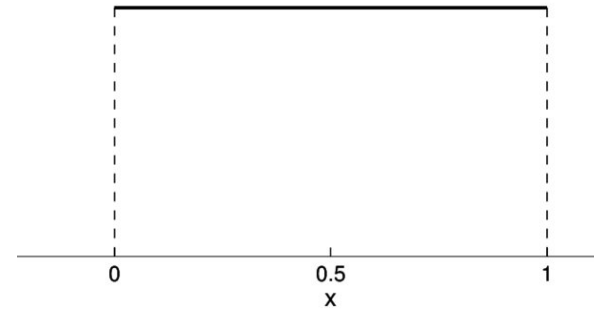
- All these functions start with the letter *r*

The uniform distribution

generate $n = 100$ points from $U(0, 1)$

```
> rand_data <- runif(100)
```

```
> hist(rand_data)
```

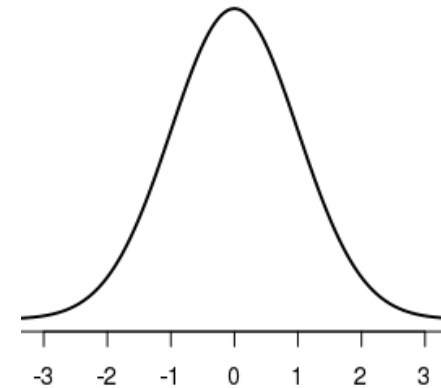


The normal distribution

generate $n = 100$ points from $N(0, 1)$

```
> rand_data <- rnorm(100)
```

```
> hist(rand_data)
```



Generating random data

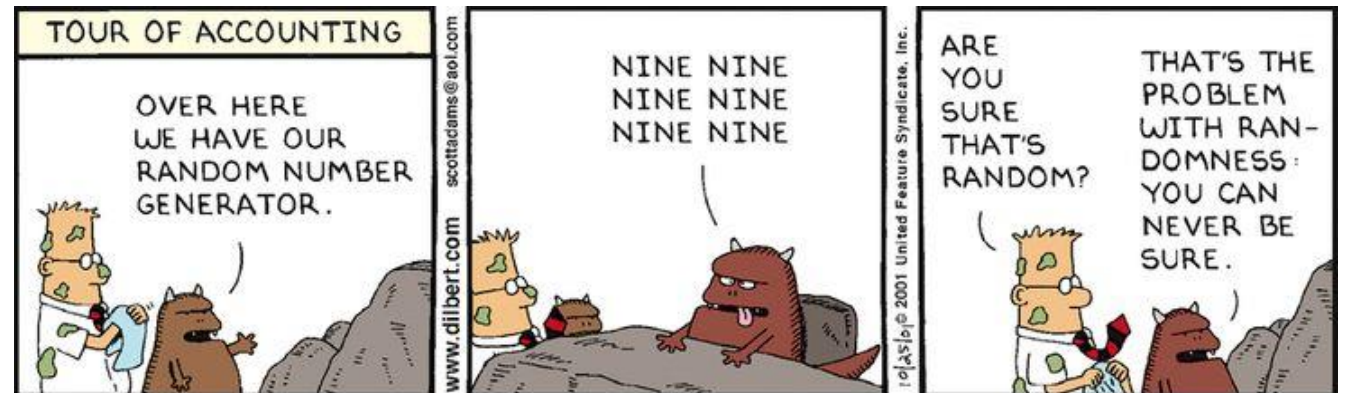
If we want the same sequence of random numbers we can set the random number generating seed

```
> set.seed(123)
```

```
> runif(100)
```

Q: Why would we want the same sequence of random number?

A: Reproducibility!



Sampling data

The `sample(v, n)` function samples `n` random points from a vector `v`

For example, suppose we had a vector with the ages of all US citizens in a vector called `pop_ages`

We could sample the ages of 100 random people using:

- `rand_sample <- sample(pop_ages, 100)`

We can sample with replacement using the `replace = TRUE` argument:

- `rand_sample_replace <- sample(pop_ages, 100, replace = TRUE)`

Let's try it in R!

Questions?

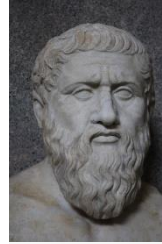


Review and extension of statistical concepts

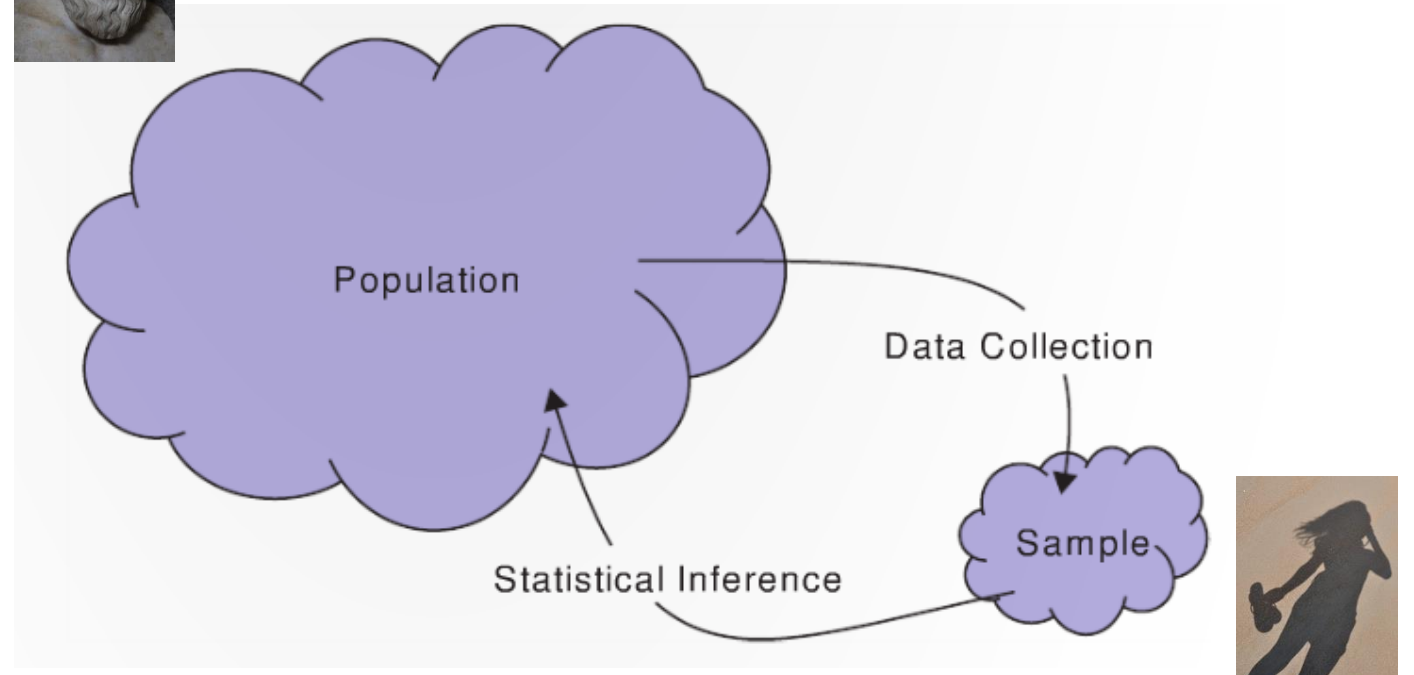
Where does data come from?



DATA SCIENCE!!!



Population: all individuals/objects of interest



Sample: A subset of the population

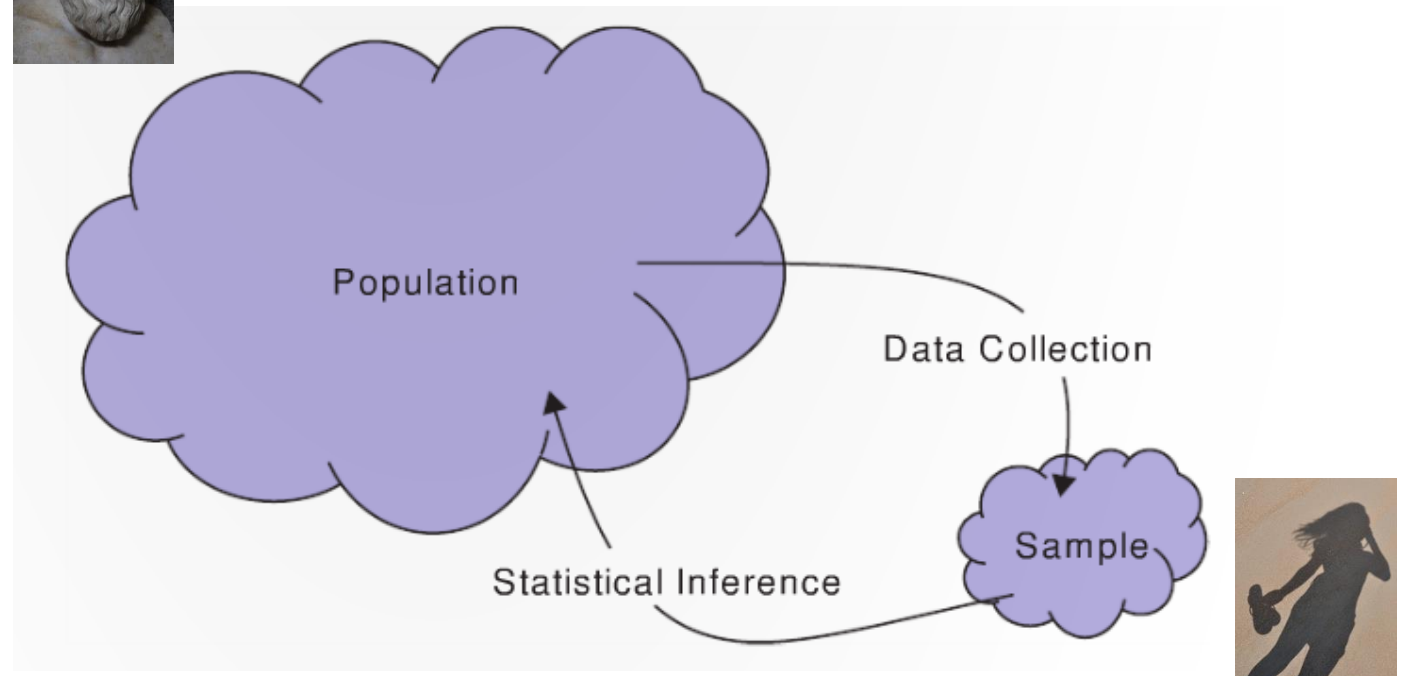
Where does data come from?

Question: Is the okcupid profiles data frame a population or a sample?

Question: If the OkCupid profiles data frame is a sample, what is the population?



Parameters: $\pi, \mu, \sigma, \rho, \beta$



Statistics: $\hat{p}, \bar{x}, s, r, b$

How do we get sample of data?

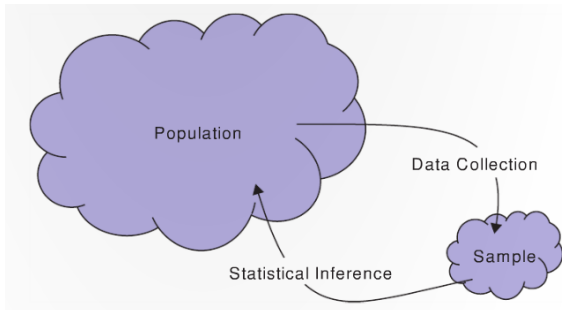
Simple random sample: each member in the population is equally likely to be in the sample

“Random selection”

Q: Why is this good?

A: Allows for generalizations to the population!

- No sampling bias
- Statistic (on average) equal parameter
 - E.g., $E[\bar{x}] = \mu$



Soup analogy!



Questions:

- Is the OkCupid profiles data a simple random sample?
- Would we expect sampling bias from statistics computed from the OkCupid profiles?

Big picture for the week

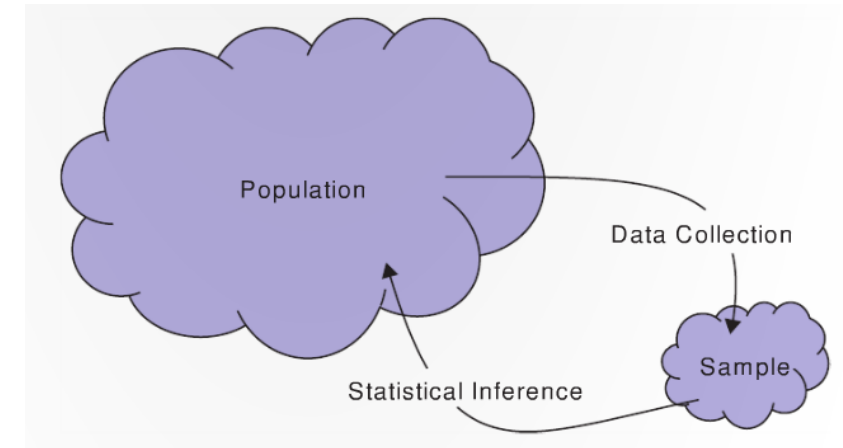
Statistics are point estimates of parameters

We can use sampling distributions (i.e., distributions of statistics) to tell us how much we can trust **any one statistic** to be a good point estimate of a parameter

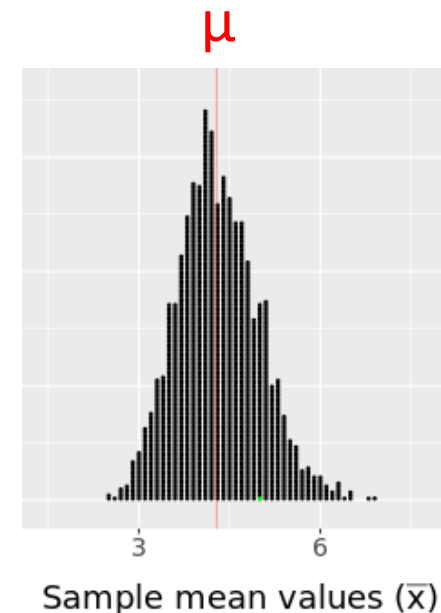
-> confidence interval

Let's start on this now...

parameter: μ



statistic: \bar{x}



**Sampling
distribution of \bar{x}**

Sampling distributions

Sample statistics

Q: What is a statistic?

A: A statistic is number computed from a function on a sample of data

The sample mean \bar{x}

(shadow of the parameter μ)

```
> rand_data <- runif(100)    # generate n = 100 points from U(0, 1)
> mean(rand_data)
```

Q: If we repeat the code above will we get the same statistic?

- A: unlikely

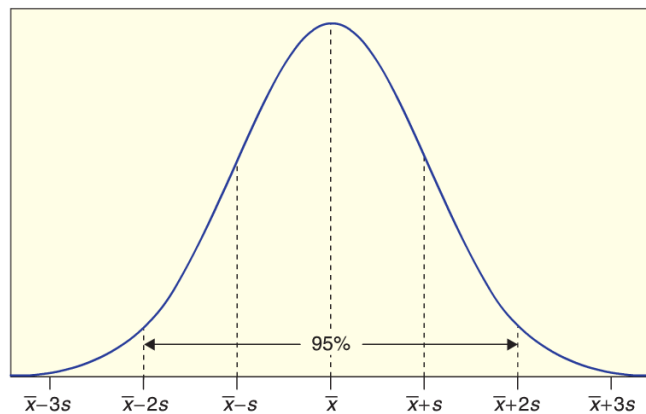
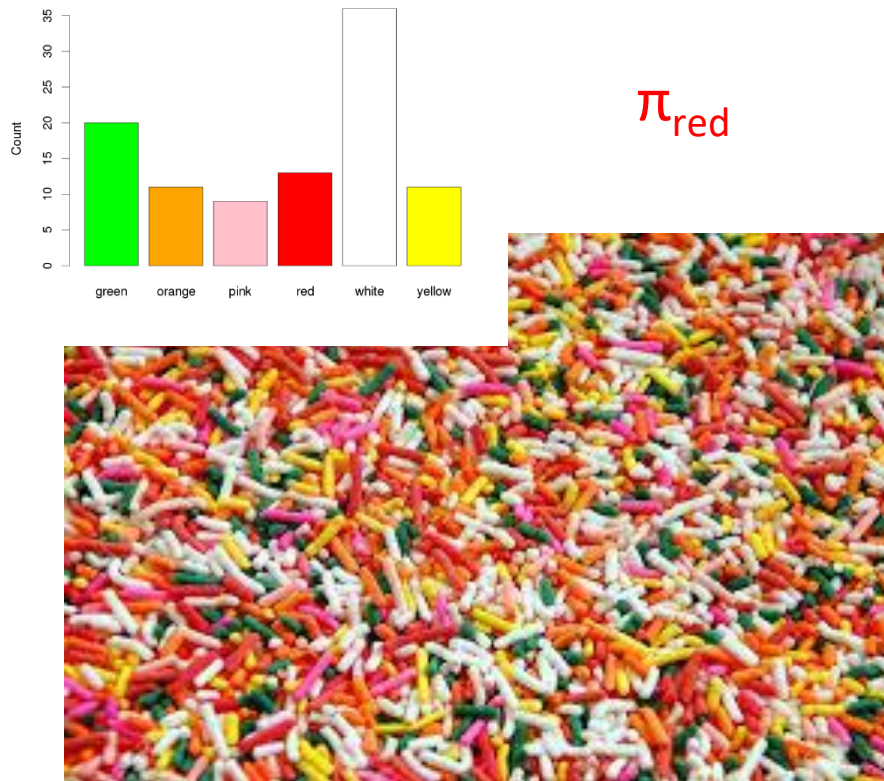
Sampling distributions

A ***sampling distribution*** is a distribution of ***statistics***

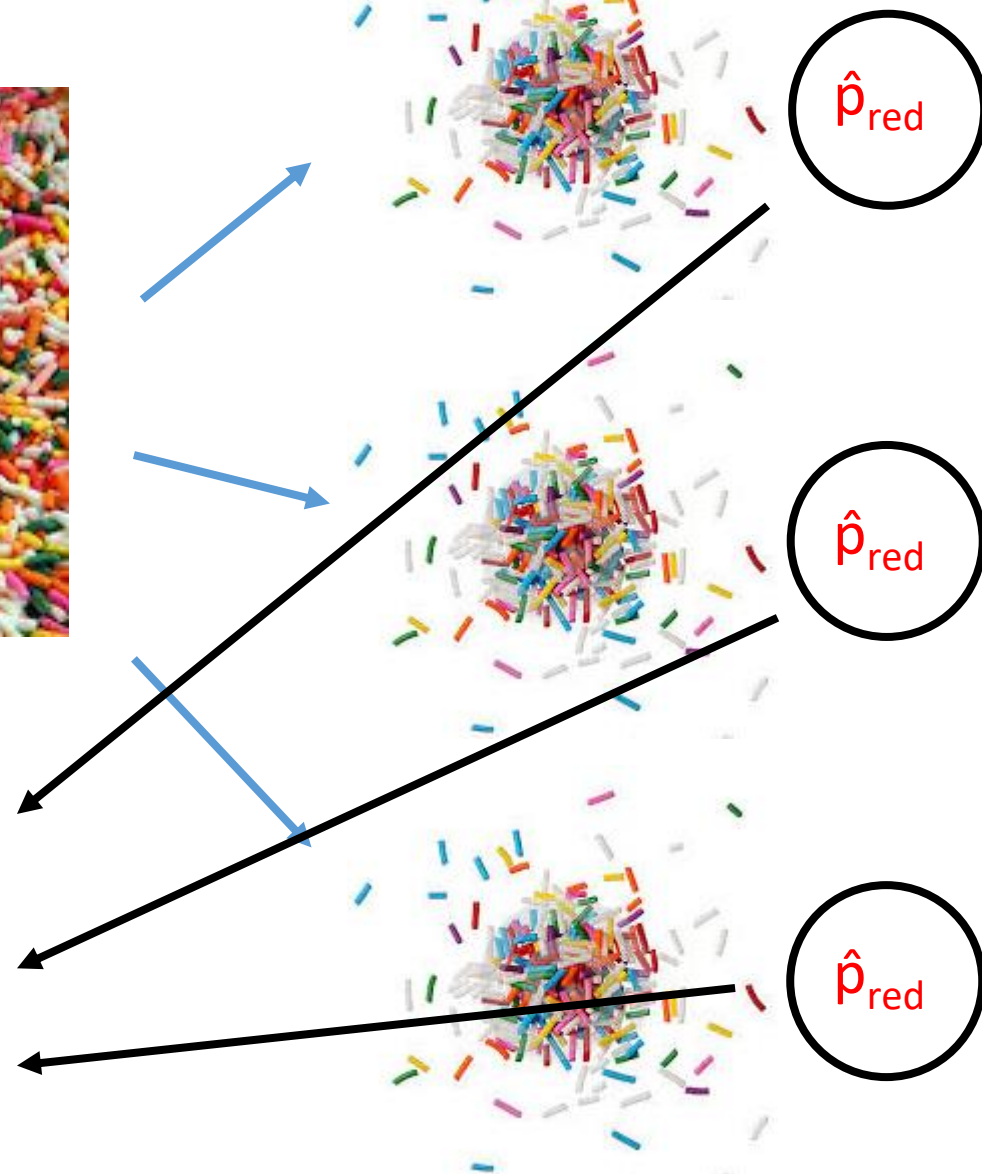
Reminder: For a *single ***categorical variable****, the main statistic of interest is the ***proportion*** (\hat{p}) in each category

- (shadow of the parameter π)

$$\hat{p} = \text{Proportion in a category} = \frac{\text{number in that category}}{\text{total number}}$$



Sampling distribution!



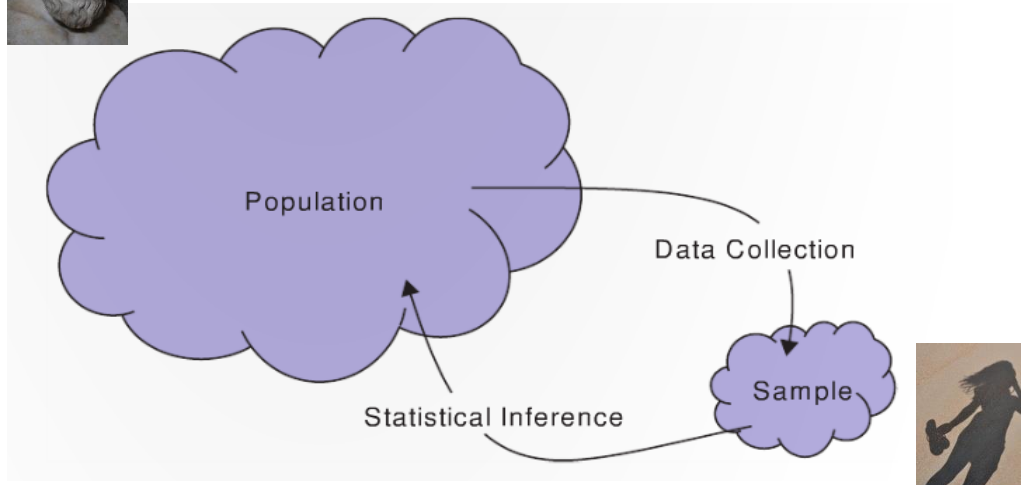
Sampling distribution

Why would we be interested in the sampling distribution?

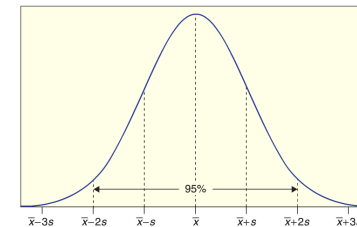
- If we knew what the sampling distribution was, then we could evaluate how much we should trust individual statistics



Parameters: π , μ , σ , ρ , β



Sampling distribution



Statistics: \hat{p} , \bar{x} , s , r , b

Simulating sampling distributions

```
sampling_dist <- NULL
for (i in 1:1000) {
    rand_data <- runif(100)  # generate n = 100 points from U(0, 1)
    sampling_dist[i] <- mean(rand_data)  # save the mean
}

hist(sampling_dist)
```

Simulating sampling distributions

Distribution of OkCupid user's heights $n = 100$

```
heights <- profiles$height
```

```
# get one random sample of heights from 100 people
```

```
height_sample <- sample(heights, 100)
```

```
# get the mean of this sample
```

```
mean(height_sample)
```

Simulating sampling distributions

Distribution of OkCupid user's heights $n = 100$

```
sampling_dist <- NULL
for (i in 1:1000) {
    height_sample <- sample(heights, 100)  # sample 100 random heights
    sampling_dist[i] <- mean(height_sample) # save the mean
}

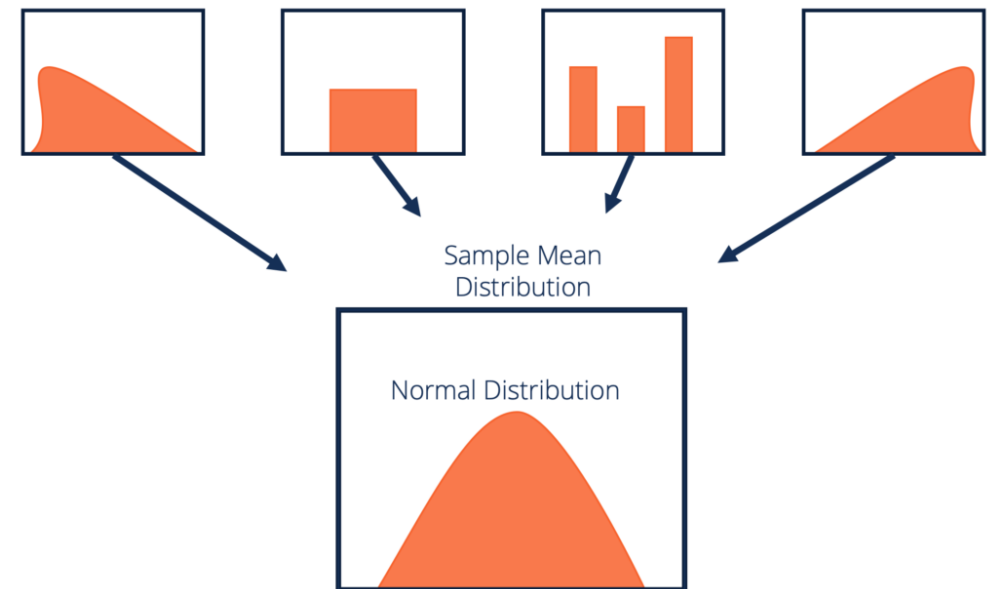
hist(sampling_dist)
```

The central limit theorem

The **central limit theorem** establishes that, in many situations, when independent random variables are summed up, their properly normalized sum tends toward a normal distribution.

Since many statistics we use are the sum of randomly data, many of our sampling distributions will be approximately normal

- You will explore this more on homework 2



Statistics: \hat{p} , \bar{x} , s , r , b

Some would say this sidewalk is broken, but it's actually normal



Confidence intervals

Point Estimate

We use the statistics from a sample as a **point estimate** for a population parameter

- \bar{x} is a point estimate for...? μ

A recent [YouGov poll](#) of 2,335 adults showed Biden's approval rating at 40.2%

Symbols:

π : Biden's approval for all voters

\hat{p} : Biden's approval for those voters in our sample

CBS News Poll – September 5-8, 2023 Adults in the U.S.

YouGov

Sample 2,335 Adults in the U.S.
Margin of Error $\pm 2.7\%$

1. Generally speaking, do you feel things in America today are going...

Very well	5%
Somewhat well	23%
Somewhat badly	34%
Very badly	38%

2. How would you rate the condition of the national economy today?

Very good	8%
Fairly good	21%
Fairly bad	31%
Very bad	35%
Not sure	5%

3. Do you approve or disapprove of the way Joe Biden is handling his job as president?

Approve	40%
Disapprove	60%

Interval estimate based on a margin of error

An **interval estimate** give a range of plausible values for a population parameter

One common form of an interval estimate is:

$$\textit{Point estimate} \pm \textit{margin of error}$$

Where the **margin of error** is a number that reflects the precision of the sample statistic as a point estimate for this parameter

Example: YouGov poll

40.2% of American approve of Biden's job performance, with a margin of error of 2.7%

- i.e., plus or minus 2.7%

How do we interpret this?

Says that the population parameter (π) lies somewhere between:

$$40.2 - 2.7 \text{ to } 40.2 + 2.7 = 37.5 \text{ to } 42.9$$

i.e., if they sampled all voters the true population proportion (π) would be likely be in this range

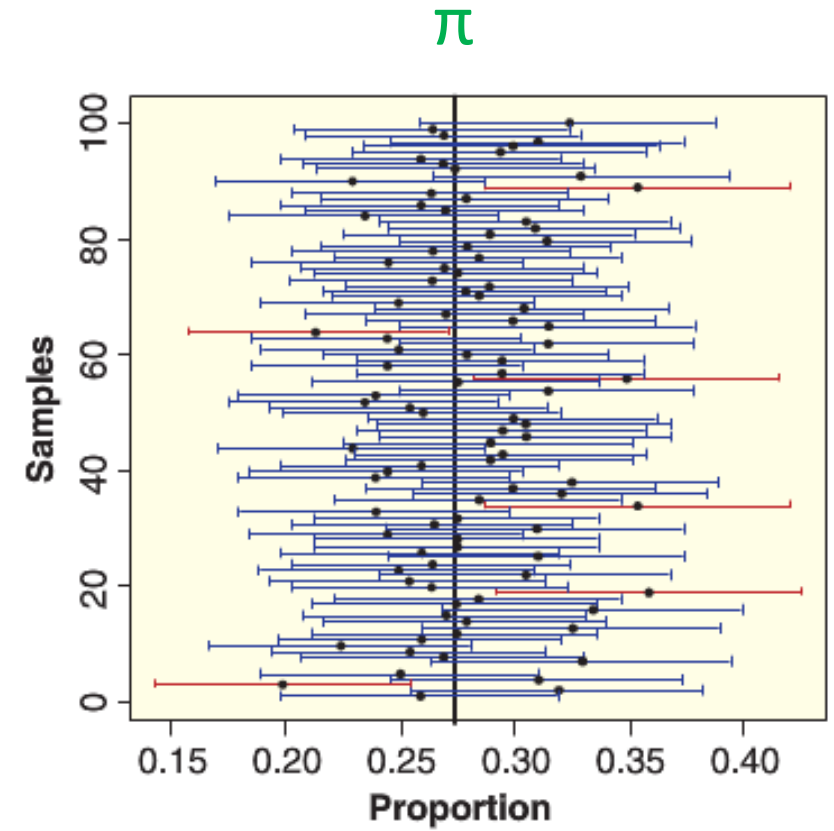


Confidence Intervals

A **confidence interval** is an interval computed by a method that will contain the ***parameter*** a specified percent of times

- i.e., if the interval was calculated repeatedly from many different random samples, the parameter will be in p% of these intervals

The **confidence level** is the percent of all intervals that contain the parameter



Think ring toss...

Parameter exists in the ideal world

We toss intervals at it

95% of those intervals capture the parameter

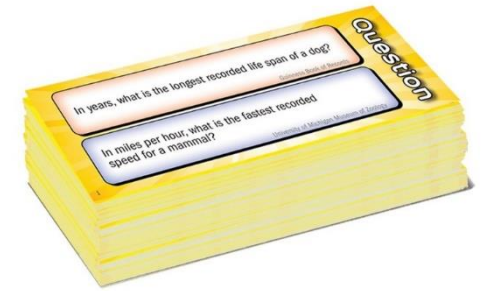


Wits and Wagers: 90% confidence intervals estimators

I am going to ask you 10 questions

You need to produce an **interval range** that contains the true answer for 9 out of the 10 questions I ask

Please write down your answers on a piece of paper



Wits and Wagers...

Question 1: What year was Yale University founded?

Question 2: In what year did Benjamin Franklin prove that lightning was electricity, after flying his kite in a thunderstorm?

Question 3: How many floors does the leaning tower of Pisa have?

Wits and Wagers...

Question 4: In feet, how tall was the tallest giraffe ever recorded?

Question 5: In years, what is the longest recorded life span of a dog?

Question 6: How many pounds does one gallon of whole milk weigh?

Question 7: In pounds, what was the weight of the heaviest domesticated cat ever recorded?

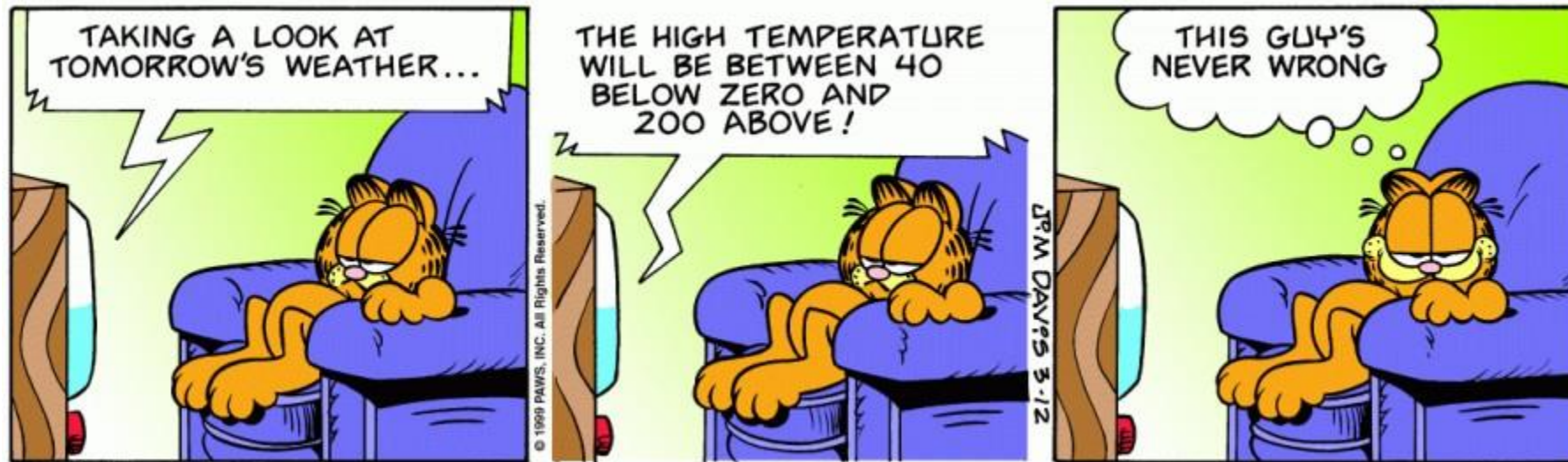
Wits and Wagers...

Question 8: If a person weighs 100 pounds on Earth, how many pounds would they weigh on the surface of the moon?

Question 9: What percentage of American adults say that reading is their favorite leisure-time activity?

Question 10: How many cups of coffee does the average American drink per year?

100% confidence intervals



There is a tradeoff between the **confidence level** (percent of times we capture the parameter) and the **confidence interval size**



Note

For any given confidence interval we compute, we don't know whether it has really captured the parameter

But we do know that if we do this 100 times, 90 of these intervals will have the parameter in it

(for a 90% confidence interval)

Next class

Computing confidence intervals using the bootstrap...

