



**DATA 1202**  
**Spring 2024**

# Lecture 20

---

Interpreting Confidence

# Announcements

---

- **Homework 10** due Wednesday
- **Lab 11** due Friday at 5pm

# Estimation

# Inference: Estimation

---

- **Parameter:** Fixed quantity in the population
- How can we figure out the value of an unknown parameter?
- **If you don't have a census:**
  - Take a random sample from the population
  - Use a statistic as an **estimate** of the parameter
- **Problem:** One sample → One estimate
  - But the random sample could have come out differently
  - And so the estimate could have been different

**We need to know the variability of our estimate**

---

# Where to Get Another Sample?

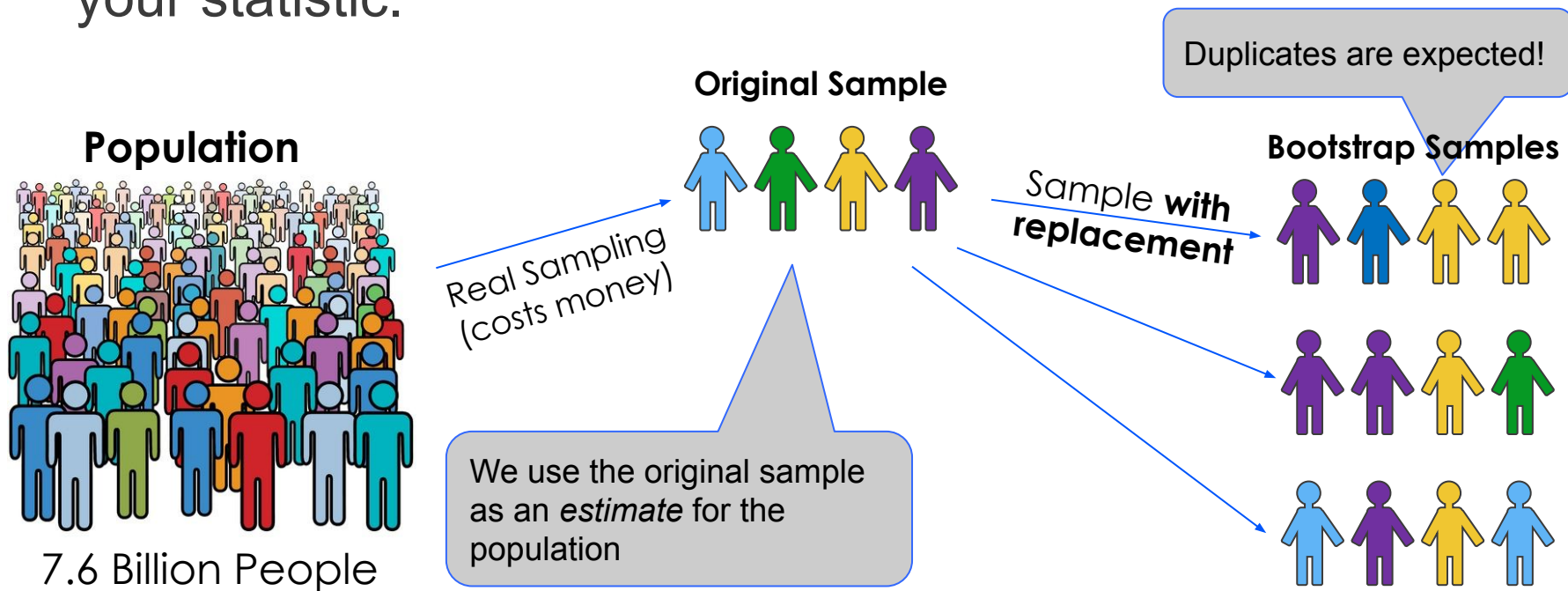
---

- **We want to understand variability of our estimate**
  - We only have the **sample**
  - To get many values of the estimate, we need many random samples
  - We can't go back and sample again from the population
-

# The Bootstrap

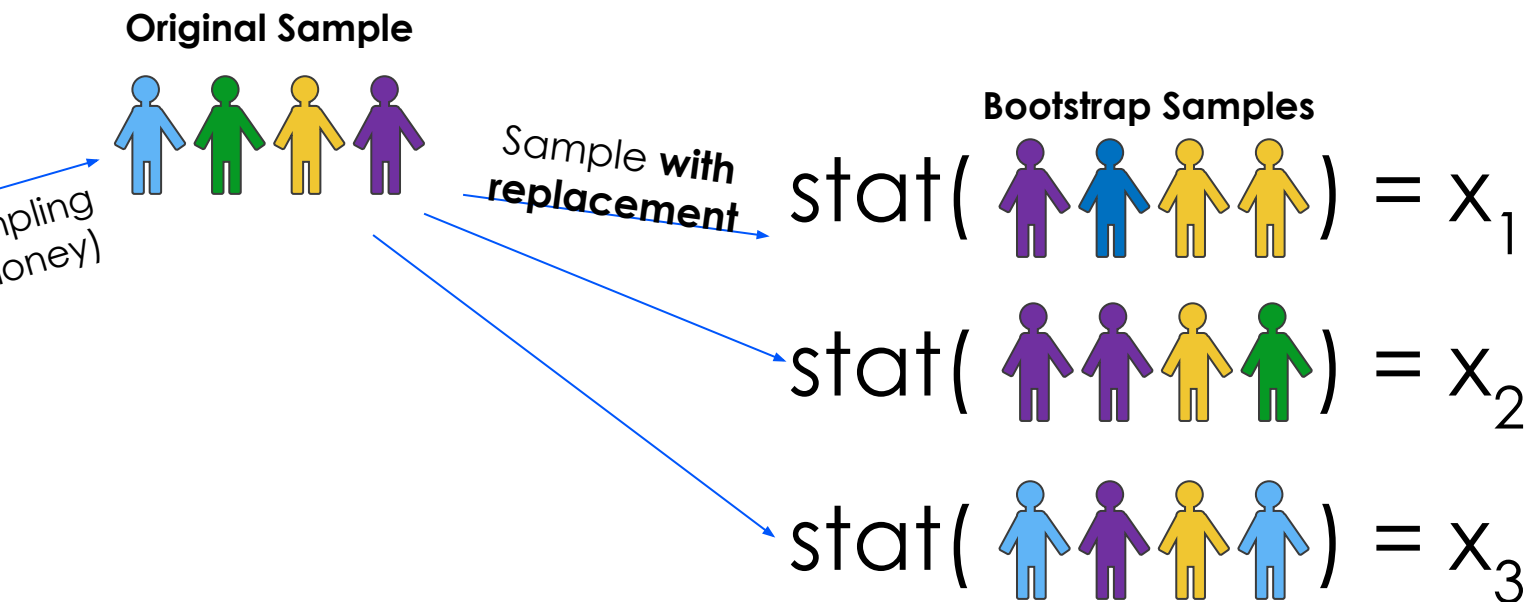
# Bootstrap the Distribution of a Statistic

Simulation method to estimate the sample distribution of your statistic.



# Bootstrap the Distribution of a Statistic

Simulation method to estimate the sample distribution of your statistic.





# Bootstrap the Distribution of a Statistic

Simulation method to estimate the sample distribution of your statistic.

Sample



Sample *with replacement*

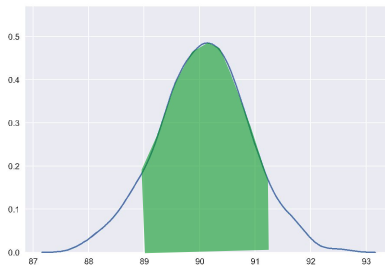
Bootstrap Samples

$$\text{stat}(\text{purple, blue, yellow, yellow}) = x_1$$

$$\text{stat}(\text{purple, purple, yellow, green}) = x_2$$

$$\text{stat}(\text{light blue, purple, yellow, light blue}) = x_3$$

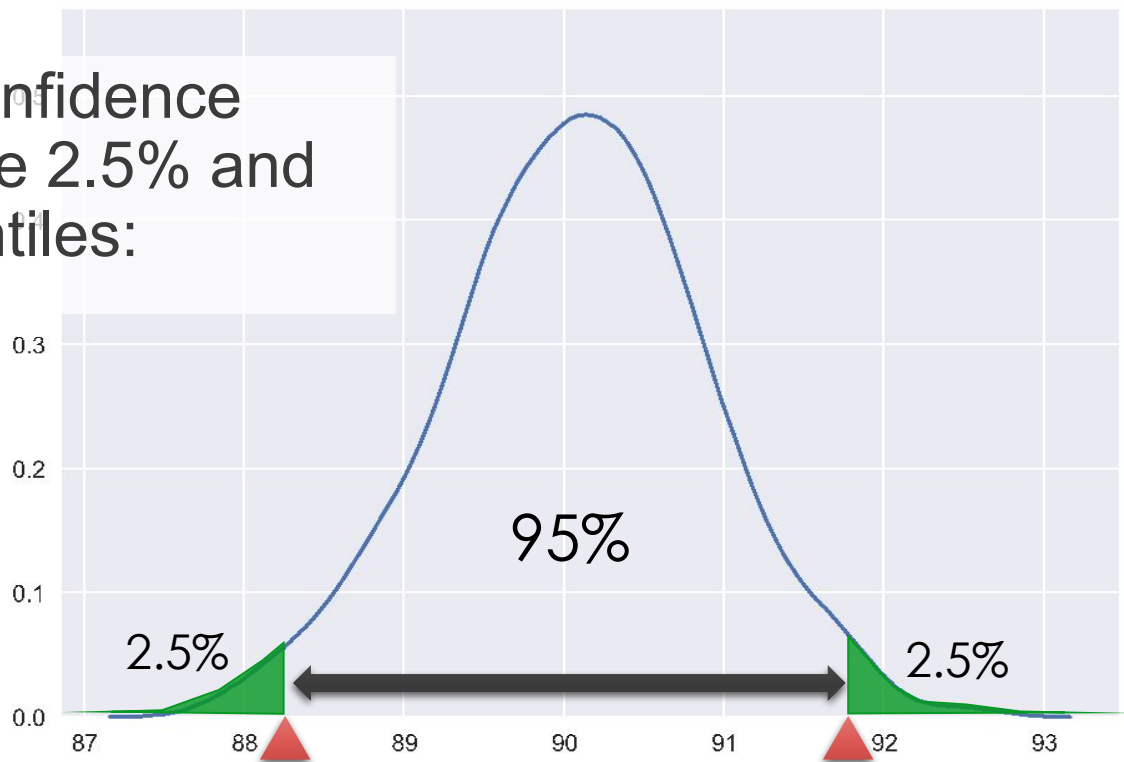
Empirical Distribution of the **Statistic**



Confidence Interval

# Bootstrap Confidence Interval

Construct a 95% confidence interval by taking the 2.5% and (100 - 2.5)% percentiles:



# The Bootstrap in words

---

- From the original sample,
  - draw at random
  - **with replacement**
    - Otherwise you would always get the same sample
  - **Use the same sample size** as the original sample
    - The size of the new sample has to be the same as the original one, so that estimates are comparable
- For each sample, **compute the statistic**
- Compute **empirical distribution of the statistics**

---

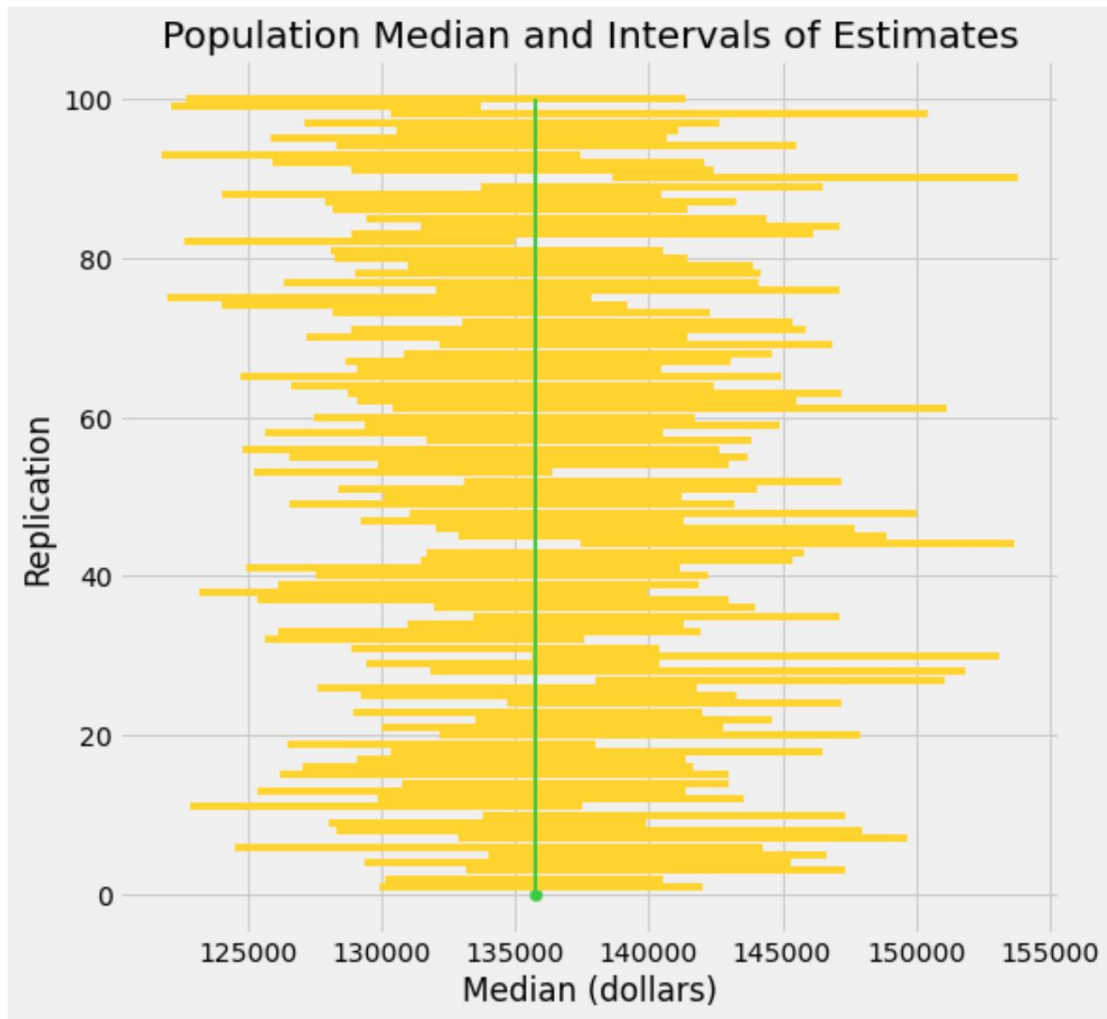
(Demo)

# Confidence Intervals

# 95% Confidence Interval

---

- Interval of **estimates of a parameter**
  - Based on random sampling
  - 95% is called the confidence level
    - Could be any percent between 0 and 100
    - Higher level means wider intervals
  - A “**good**” interval is one that contains the parameter
  - The **confidence is in the process** that creates the interval:
    - It generates a “good” interval about 95% of the time.
-



## The Meaning of 95% confidence

The **green line** is the parameter value.

**It is fixed and unknown.**

(For this demo we we had access to the population but you won't in practice.)

Each **yellow line** is a 95% **confidence interval** based on a **fresh sample** from the population

There are **100 intervals**.  
We expect **roughly 95** to contain the parameter.

(Demo)

**Use Methods Appropriately**

# When *Not* to Use Our Bootstrap Method

---

- If you're trying to estimate any parameter that's greatly affected by rare elements of the population
  - Very high or very low percentiles, or min and max
- If the probability distribution of your statistic is not roughly bell shaped (the shape of the empirical distribution will be a clue)
- If the original sample is very small

(Demo)

---



# Can You Use a CI Like This?

---

By our calculation, an approximate 95% confidence interval for the average age of the mothers in the population is (26.9, 27.6) years.

## True or False:

- About 95% of the mothers in the population were between 26.9 years and 27.6 years old.

**Answer: False.** We're estimating that their **average age** is in this interval.

---

# Is This What a CI Means?

---

An approximate 95% confidence interval for the average age of the mothers in the population is (26.9, 27.6) years.

## True or False:

- There is a 0.95 probability that the average age of mothers in the population is in the range 26.9 to 27.6 years.

**Answer: False.** The parameter is fixed, and the interval (26.9, 27.2) is fixed. The parameter is either in that interval, or not. Once you've picked an interval, there's no probability involved.

---

# 95% Confidence

---

- Interval of estimates of a parameter
  - Based on random sampling
  - The process results in a random interval
  - A “good” interval is one that contains the parameter
  - The **confidence is in the process** that creates the interval:
    - It generates a “good” interval with chance 95%
-

# Confidence Intervals For Testing

# Using a CI for Testing

---

- Null hypothesis: **Population average =  $x$**
- Alternative hypothesis: **Population average  $\neq x$**
- Cutoff for p-value:  $p\%$
- Method:
  - Construct a  $(100-p)\%$  confidence interval for the population average
  - If  $x$  is not in the interval, reject the null
  - If  $x$  is in the interval, can't reject the null

(Demo)

---