

# Studying Social Inequality with Data Science

INFO 3370 / 5371  
Spring 2024

Sampling: Stratified, Clustered, and the Future

# Learning goals for today

By the end of class, you will be able to

- ▶ sample from a population in R
- ▶ write an estimator function
- ▶ apply the function to your sample
- ▶ connect sampling to the replication crisis
- ▶ discuss the future of sampling

# Baseball salaries

BASEBALL

The New York Times

THE NEW YORK TIMES

Sections

Los Angeles Times

SUBSCRIBE

LOG IN



Dodgers news Teoscar Hernández California dreaming Dodgers pitchers rising \$1 billion boon?

## Channeling the Old Steinbrenner Ways, Yankees Stepped Up for Judge

Aaron Judge, who hit 62 home runs in 2022, agreed to a nine-year, \$360 million contract with the Yankees after meeting with at least two other teams.

Share via email Share on Facebook Share on Twitter 128



Aaron Judge set career highs in batting average (.311), home runs (62) and R.B.I. (137) in 2022. Chris DiDonato for The New York Times

DODGERS

## Complete coverage: Shohei Ohtani signs record deal with Dodgers



Shohei Ohtani speaks during his introductory Dodgers news conference at Dodger Stadium on Thursday. (Wafiq Skali / Los Angeles Times)

BY LOS ANGELES TIMES STAFF

PUBLISHED DEC. 9, 2023 | UPDATED DEC. 22, 2023 8:54 AM PT

# Baseball salaries

BASEBALL

The New York Times

Sections

Channeling the Old Steinbrenner Ways, Yankees Stepped Up for Judge

Aaron Judge, who hit 62 home runs in 2022, agreed to a nine-year, \$360 million contract with the Yankees after meeting with at least two other teams.

Share via email

128



Aaron Judge set career highs in batting average (.311), home runs (62) and R.B.I. (137) in 2022. Chris DiDonato for The New York Times

Los Angeles Times

Sections

Subscribe

Log In

Dodgers news Teoscar Hernández California dreaming Dodgers pitchers rising \$1 billion boon?

## Complete coverage: Shohei Ohtani signs record deal with Dodgers



Shohei Ohtani speaks during his introductory Dodgers news conference at Dodger Stadium on Thursday. (Wafiq Skaliq / Los Angeles Times)

BY LOS ANGELES TIMES STAFF

PUBLISHED DEC. 9, 2023 | UPDATED DEC. 22, 2023 8:54 AM PT

## Major League Baseball Minimum: \$720,000

# Baseball salaries

## Major League Baseball Salaries 2023

Major League Baseball salaries based on players on opening day rosters and injured list and restricted list. Figures, compiled by USA TODAY, are based on documents obtained from Major League Baseball, the MLB Players Association, clubs officials and agents, filed with MLB's central office. Deferred payments and incentive clauses are not included. See [more salaries for 2022](#).

Source: USA TODAY Sports

### Quick Search

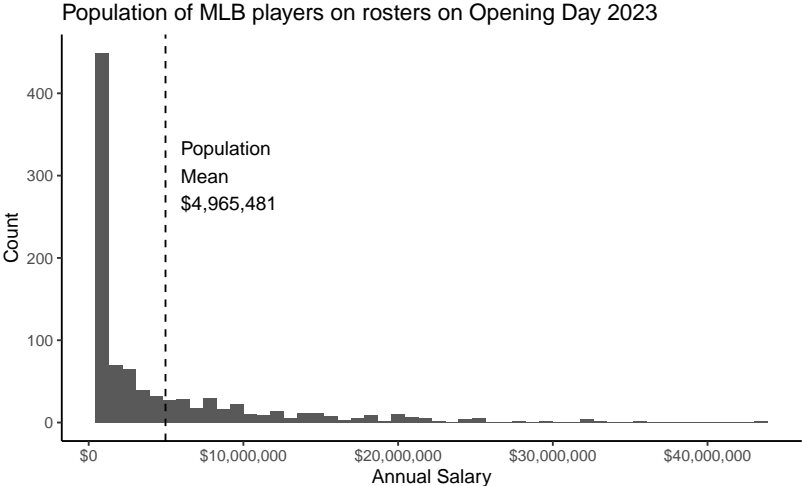
Player  Team  Position

Show/Hide Columns

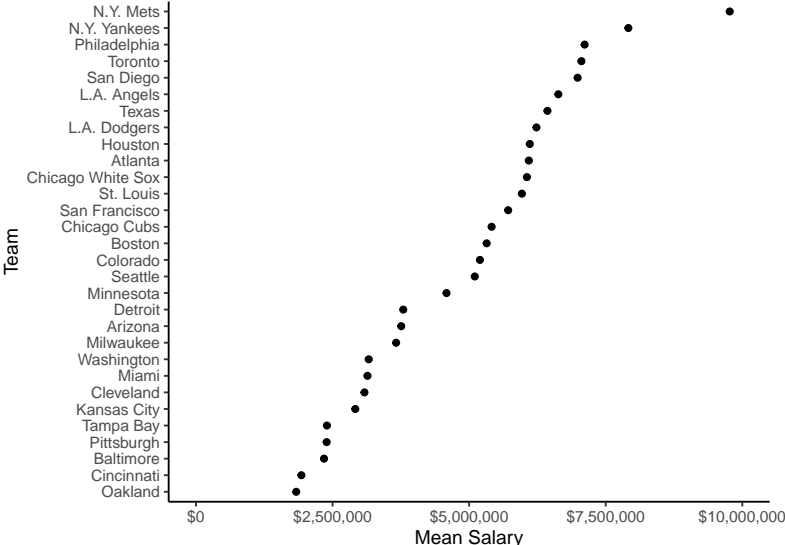
Player	Team	Position	Salary	Years	Total Value
<a href="#">Scherzer, Max</a>	N.Y. Mets	RHP	\$43,333,333	3	\$130,000,000
<a href="#">Verlander, Justin</a>	N.Y. Mets	RHP	\$43,333,333	2	\$86,666,666
<a href="#">Judge, Aaron</a>	N.Y. Yankees	OF	\$40,000,000	9	\$360,000,000
<a href="#">Rendon, Anthony</a>	L.A. Angels	3	\$38,571,429	7	\$245,000,000
<a href="#">Trout, Mike</a>	L.A. Angels	OF	\$37,116,667	12	\$426,500,000

[databases.usatoday.com/major-league-baseball-salaries-2023/](https://databases.usatoday.com/major-league-baseball-salaries-2023/)

# Baseball salaries



# Baseball salaries



# Draw a Sample to Estimate the Mean Salary

```
baseball <- read_csv("https://info3370.github.io/data/baseball.csv")
```

How would you design:

- ▶ Simple random sample of 60 players
- ▶ Random sample stratified by team
- ▶ Random sample clustered by team

and why would you do it each way?

Stuck? See last week's [reading](#)



# Draw a Sample to Estimate the Mean Salary

simple random sampling	60 players chosen at random
stratified sampling	2 players on each of the 30 teams
clustered sampling	20 players on 3 sampled teams

# Apply an Estimator

Write a function that I like to call `estimator()`

- ▶ input is a sample
- ▶ output is an estimate

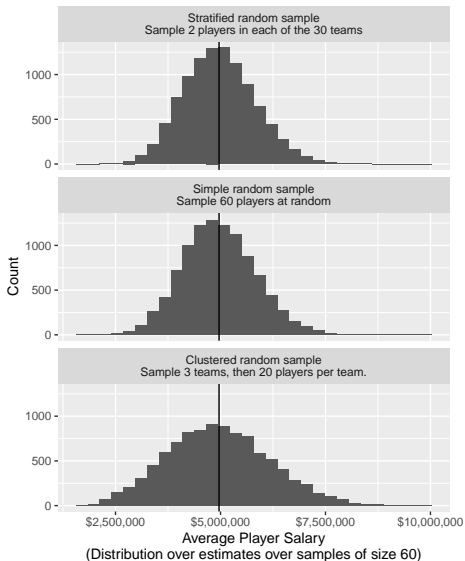
# Evaluate performance

We will first calculate the population mean

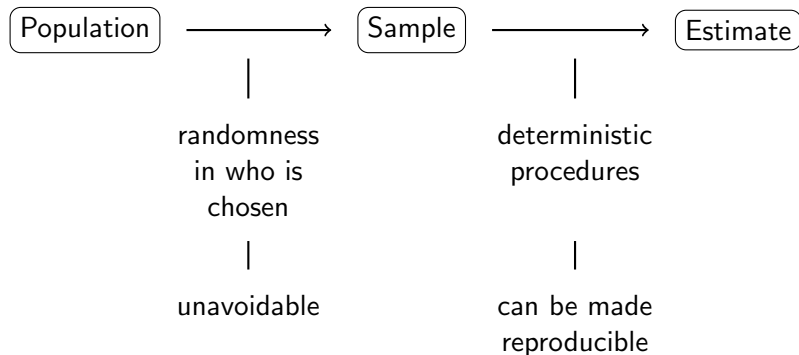
Then we will repeatedly

- ▶ draw a sample
- ▶ apply the estimator
- ▶ store the result

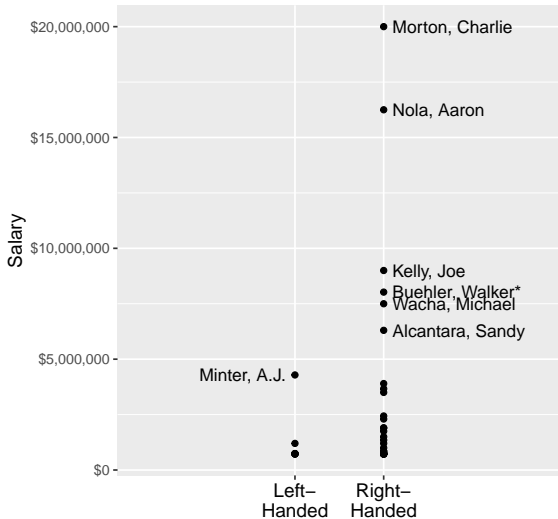
# Three sampling strategies

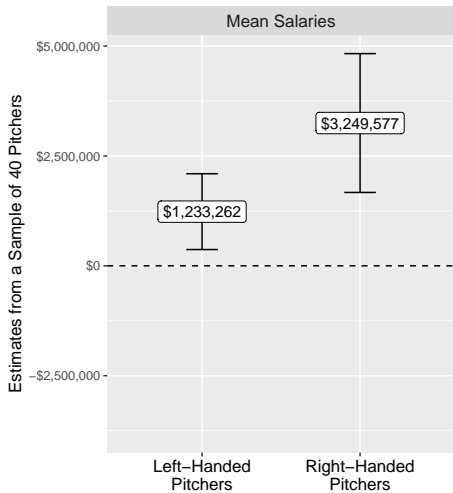


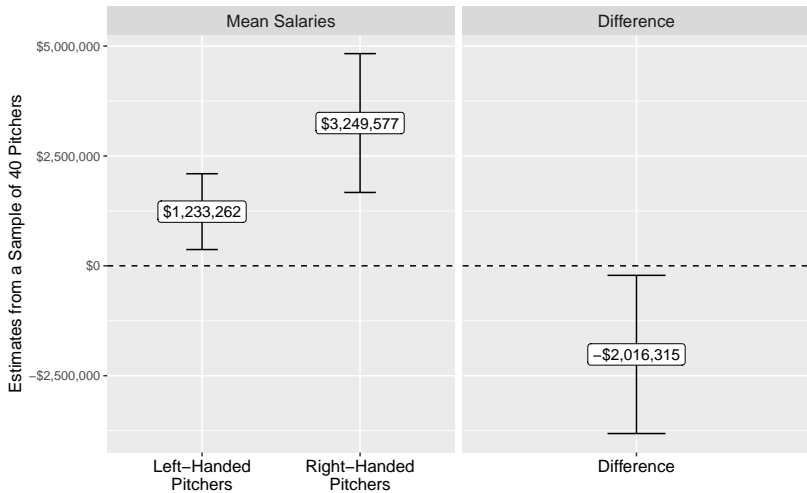
# Danger of One Sample



## Sample of 40 Pitchers from Opening Day 2023

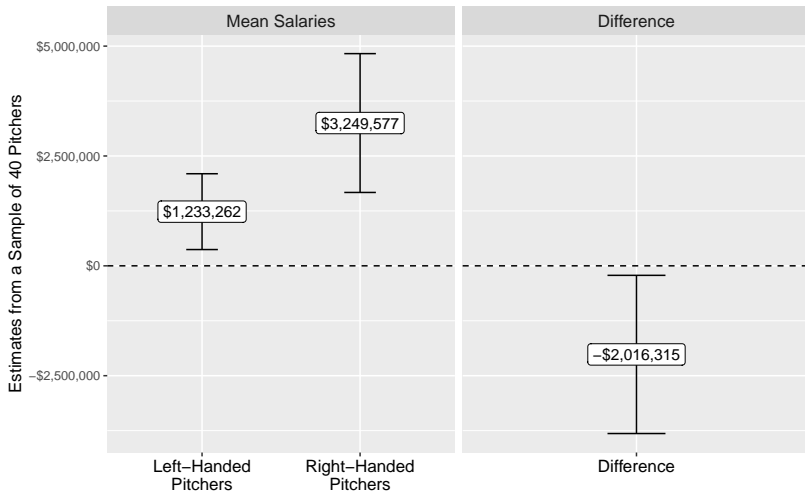








# Why might right-handed pitchers earn more?



# Your turn

- ▶ load the data
- ▶ take a sample of size 40
- ▶ group by position
- ▶ summarize the mean salary

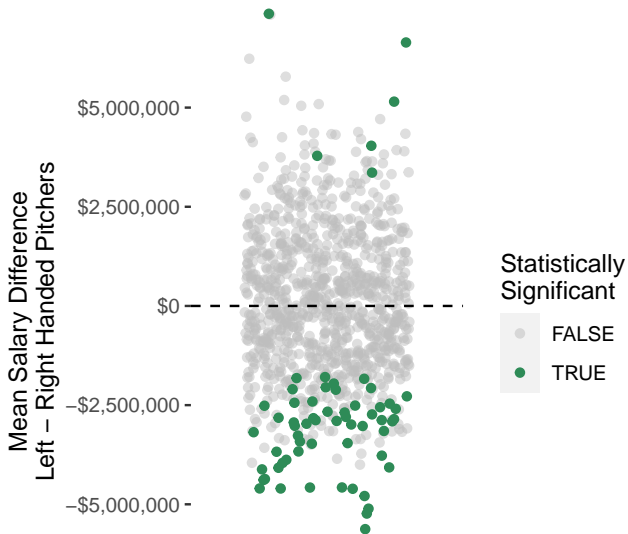
Who has higher average salary in your sample?

- ▶ RHP: right-handed pitchers
- ▶ LHP: left-handed pitchers

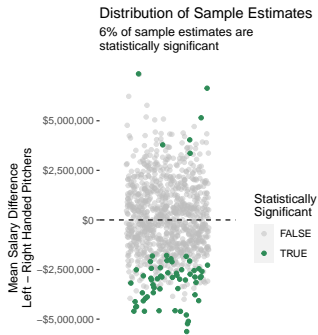
I did this 1,000 times

# Distribution of Sample Estimates

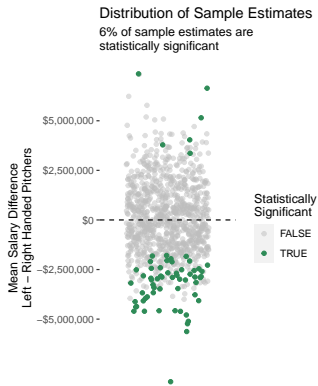
6% of sample estimates are statistically significant



# The replication crisis

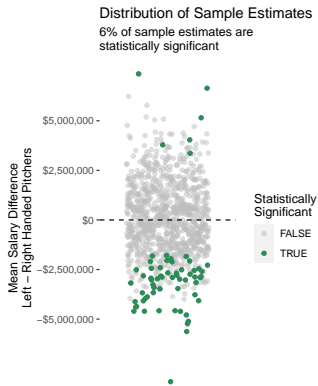


# The replication crisis



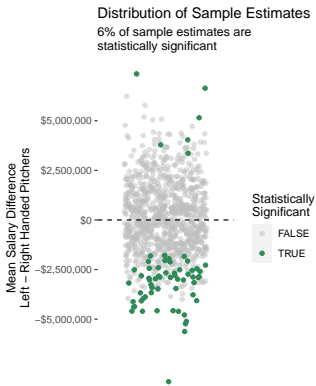
- ▶ unless we see the population, all estimates involve noise

# The replication crisis



- ▶ unless we see the population, all estimates involve noise
- ▶ surprising findings yield big rewards

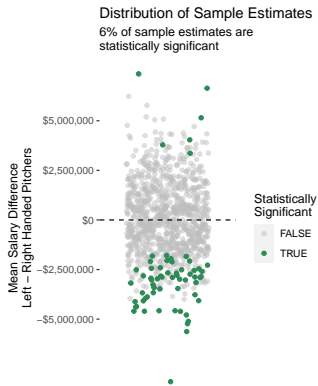
# The replication crisis



- ▶ unless we see the population, all estimates involve noise
- ▶ surprising findings yield big rewards
- ▶ unsurprising findings get ignored



# The replication crisis



- ▶ unless we see the population, all estimates involve noise
- ▶ surprising findings yield big rewards
- ▶ unsurprising findings get ignored
- ▶ science is just discovering noise

## Feeling the Future: Experimental Evidence for Anomalous Retroactive Influences on Cognition and Affect

Daryl J. Bem  
Cornell University

SCIENCE

# Daryl Bem Proved ESP Is Real

Which means science is broken.

BY DANIEL ENGBER

JUNE 07, 2017 • 2:57 PM

[Slate link.](#)

## Evaluating the replicability of social science experiments in *Nature* and *Science* between 2010 and 2015

Colin F. Camerer<sup>1,3\*</sup>, Anna Dreber<sup>2,3\*</sup>, Felix Holzmeister<sup>1,3,4</sup>, Teck-Hua Ho<sup>4,5</sup>, Jürgen Huber<sup>3,6</sup>, Magnus Johannesson<sup>2,3\*</sup>, Michael Kirchner<sup>1,3,4</sup>, Gideon Nave<sup>2,3</sup>, Brian A. Nosek<sup>1,3,4,6\*</sup>, Thomas Pfeiffer<sup>1,3\*</sup>, Adam Altmeld<sup>1</sup>, Nick Buttrick<sup>1,3</sup>, Taizan Chan<sup>7</sup>, Yiling Chen<sup>8</sup>, Eskil Forsell<sup>9</sup>, Anup Gampa<sup>10</sup>, Emma Heikensten<sup>7</sup>, Lily Hummer<sup>8</sup>, Taisuke Imai<sup>11</sup>, Siri Isaksson<sup>7</sup>, Dylan Manfredi<sup>8</sup>, Julia Rose<sup>7</sup>, Eric-Jan Wagenmakers<sup>12</sup> and Hang Wu<sup>13</sup>

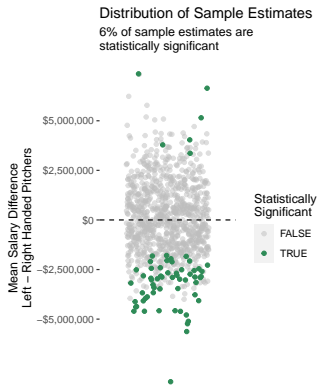
## Essay: The Experiments Are Fascinating. But Nobody Can Repeat Them.

Science is mired in a “replication” crisis. Fixing it will not be easy.

Camerer et al. in *Nature Human Behavior*.

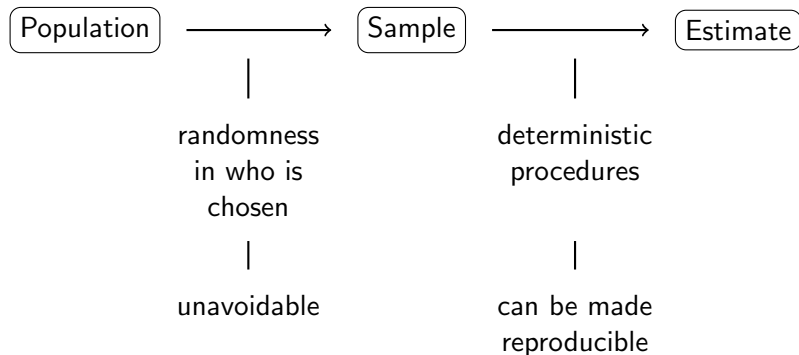
Gelman in *NYTimes*.

# The replication crisis



- ▶ unless we see the population, all estimates involve noise
- ▶ surprising findings yield big rewards
- ▶ unsurprising findings get ignored
- ▶ science is just discovering noise

# Danger of One Sample



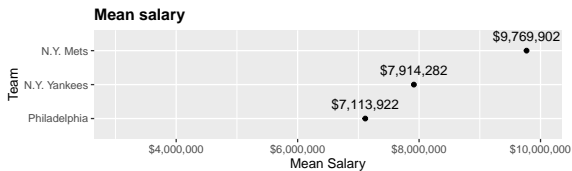
# Reproducibility

What is a typical salary in the three highest-paying teams in American baseball?

- ▶ use the whole population
- ▶ summarize salary grouped by team
- ▶ be ready to tell use your estimates and how you got them

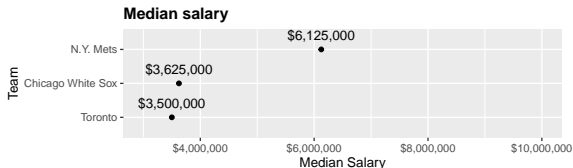
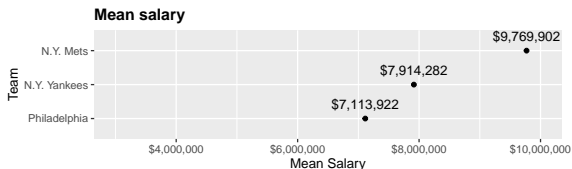
What is a typical salary in the three highest-paying teams in American baseball?

What is a typical salary in the three highest-paying teams in American baseball?

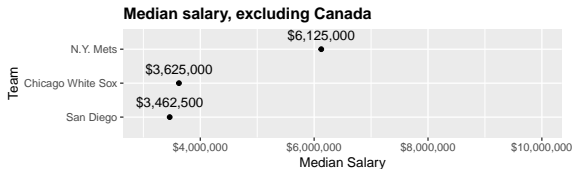
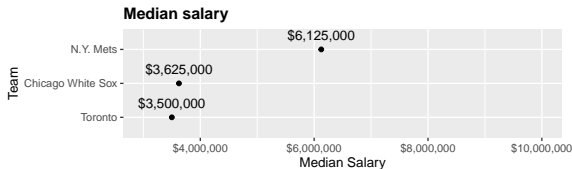
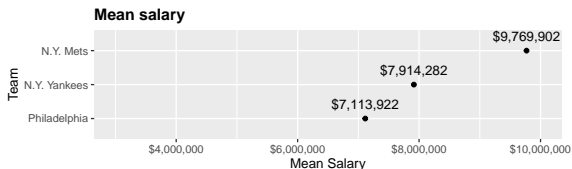




# What is a typical salary in the three highest-paying teams in American baseball?



# What is a typical salary in the three highest-paying teams in American baseball?





```
---  
title: "Problem Set 1: Visualization"  
format: pdf  
---
```

**\*\*Due: 5pm on Wednesday, January 31.\*\***

Student identifier: [type your anonymous identifying number here]

- Use this template to complete the problem set
- In Canvas, you will upload the PDF produced by your .qmd file
- Put your identifier above, not your name! We want anonymous grading to be possible

This problem set involves both data analysis and reading.

### ### Data analysis

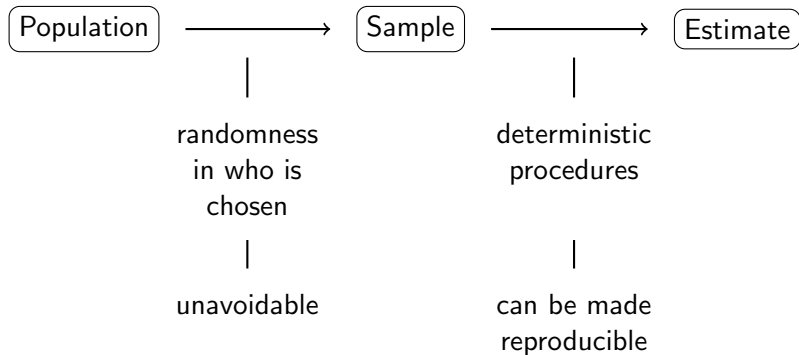
This problem set uses the data

[`lifeCourse.csv`](<https://info3370.github.io/data/lifeCourse.csv>).

```
```{r, comment = F, message = F}  
library(tidyverse)  
library(scales)  
lifeCourse <- read_csv("https://info3370.github.io/data/lifeCourse.csv")  
```
```

The data contain life course earnings profiles for four cohorts of American workers: those born in 1940, 1950, 1960, and 1970. Each row contains a

# Danger of One Sample



# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

# The Future of Sample Surveys

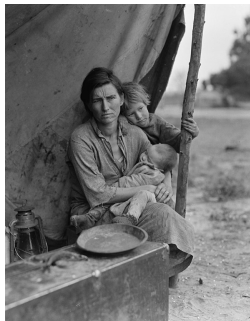
Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1930–1960: Era of Invention

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

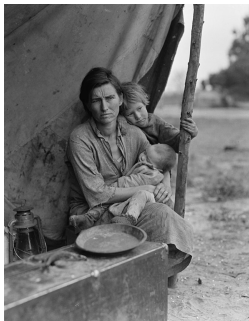
## 1930–1960: Era of Invention



# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1930–1960: Era of Invention



[Skip to Main Content](#)

**USDA** **United States Department of Agriculture**  
National Agricultural Statistics Service

Subscriptions: [Website](#) | [State](#) | [News](#)

Search NASS

Home | [Data & Statistics](#) | [Publications](#) | [Newsroom](#) | [Surveys](#) | [Census](#) | [About NASS](#) | [Contact Us](#) | [Help](#)

**Today's Reports** [View previous reports](#)

Feb 01, 2024

**Cotton System**  
Released at 3:00 pm ET [Text](#) | [ESE](#) | [CSV](#)

**Fats & Oils**  
Released at 3:00 pm ET [Text](#) | [ESE](#) | [CSV](#)

**Flour Milling**  
Released at 3:00 pm ET [Text](#) | [ESE](#) | [CSV](#)

**MILK PRODUCTION**  
ENHANCED Visualizations and Interactive Data

DATA ACCESS: We are updating our systems and plan to avoid interruptions. However, NASS data and reports are available in multiple ways in addition to this website - Cornell University Mann Library (a USDA repository) [website](#) and [e-mail report subscription service](#); QuickStats [database](#), [API](#), and downloadable [data files](#); and a [JSON file](#) for principal economic indicator data.



# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1930–1960: Era of Invention

sampling frame

pieces of land

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1930–1960: Era of Invention

sampling frame  
mode

pieces of land  
face-to-face interviews

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1930–1960: Era of Invention

sampling frame

pieces of land

mode

face-to-face interviews

cost

high

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1930–1960: Era of Invention

sampling frame

pieces of land

mode

face-to-face interviews

cost

high

response rate

over 90 percent

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1960–1990: Era of Expansion

Technology helped: Telephones



Source: Wikimedia

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1960–1990: Era of Expansion

Technology helped: Telephones  
— sampling frame



Source: Wikimedia

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1960–1990: Era of Expansion

Technology helped: Telephones

— sampling frame

— mode of data collection



Source: Wikimedia

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1960–1990: Era of Expansion

Technology helped: Telephones

- sampling frame
- mode of data collection
- falling costs



Source: Wikimedia



# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1960–1990: Era of Expansion

Technology helped: Telephones

- sampling frame
- mode of data collection
- falling costs
- falling response rates



Source: Wikimedia

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1990–Present

Technology brought challenges    Technology brought opportunities

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1990–Present

Technology brought challenges    Technology brought opportunities  
— answering machines

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present

Technology brought challenges      Technology brought opportunities

— answering machines

— cell phones

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present

Technology brought challenges      Technology brought opportunities

— answering machines

— cell phones

— caller ID

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present

Technology brought challenges      Technology brought opportunities

- answering machines
- cell phones
- caller ID
- response rates plummeted

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present

Technology brought challenges

- answering machines
- cell phones
- caller ID
- response rates plummeted

Technology brought opportunities

- digital trace data
- internet panels

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1990–Present: Designed and Organic Data



# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

1990–Present: Designed and Organic Data

Designed data

Organic data

**Example**

Census age distribution

**Example**

Web histories

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present: Designed and Organic Data

Designed data

— high cost

Organic data

— almost free

**Example**

Census age distribution

**Example**

Web histories

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present: Designed and Organic Data

Designed data

— high cost

— becoming scarce

Organic data

— almost free

— becoming abundant

**Example**

Census age distribution

**Example**

Web histories

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present: Designed and Organic Data

### Designed data

- high cost
- becoming scarce
- speak to population

### Example

Census age distribution

### Organic data

- almost free
- becoming abundant
- iffy for population

### Example

Web histories

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present: Designed and Organic Data

### Designed data

- high cost
- becoming scarce
- speak to population

### Organic data

- almost free
- becoming abundant
- iffy for population

### Example

Census age distribution

### Example

Web histories

future of **organic data**

future of **designed data**

# The Future of Sample Surveys

Groves, R. M. (2011). [Three eras of survey research](#). Public Opinion Quarterly.

## 1990–Present: Designed and Organic Data

### Designed data

- high cost
- becoming scarce
- speak to population

### Organic data

- almost free
- becoming abundant
- iffy for population

### Example

Census age distribution

### Example

Web histories

the future is **together**

# Learning goals for today

By the end of class, you will be able to

- ▶ sample from a population in R
- ▶ write an estimator function
- ▶ apply the function to your sample
- ▶ connect sampling to the replication crisis
- ▶ discuss the future of sampling