

Lecture 12: Sept 28, 2018

Bootstrap

- *Non-parametric Bootstrap*
- *Parametric Bootstrap*

James Balamuta
STAT 385 @ UIUC

Announcements

- **hw04** is due **Friday, Sep 28th, 2018** at **6:00 PM**
- **hw05** will be released on Saturday evening
 - Due **Friday, Oct 5th, 2018** at **6:00 PM**
- **Quiz 06** covers Week 5 contents @ [CBTF](#).
 - Window: Oct 2nd - 4th
 - Sign up: <https://cbtf.engr.illinois.edu/sched>
- Want to review your homework or quiz grades?
Schedule an appointment.
- Got caught using GitHub's web interface in hw01 or hw02? Let's chat.

Recap

- **Iteration**

- Forms of repeating the same instruction
- Common structures: *for*, *while*, and *repeat*
- Special controls for inside an iteration structure
 - *break*: exit out of loop
 - *next*: go to the next value in loop.

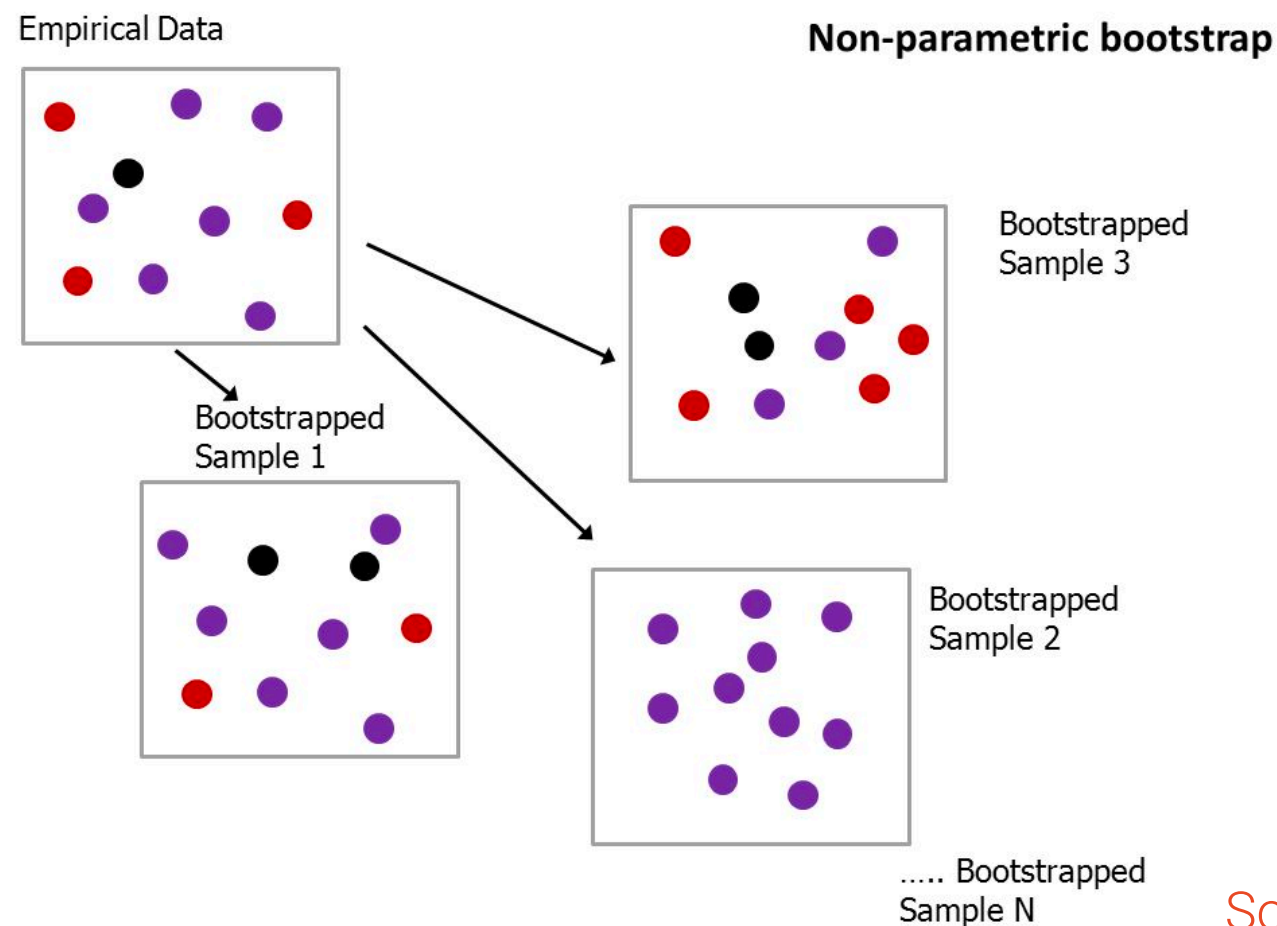
Lecture Objectives

- **Understanding scenarios** to use bootstrap
- **Importance of resampling** in the bootstrap procedure.
- **Differences** between **parametric** and **non-parametric** bootstrap.
- **Creating** and **interpreting quantile** confidence intervals.

Non-parametric Bootstrap

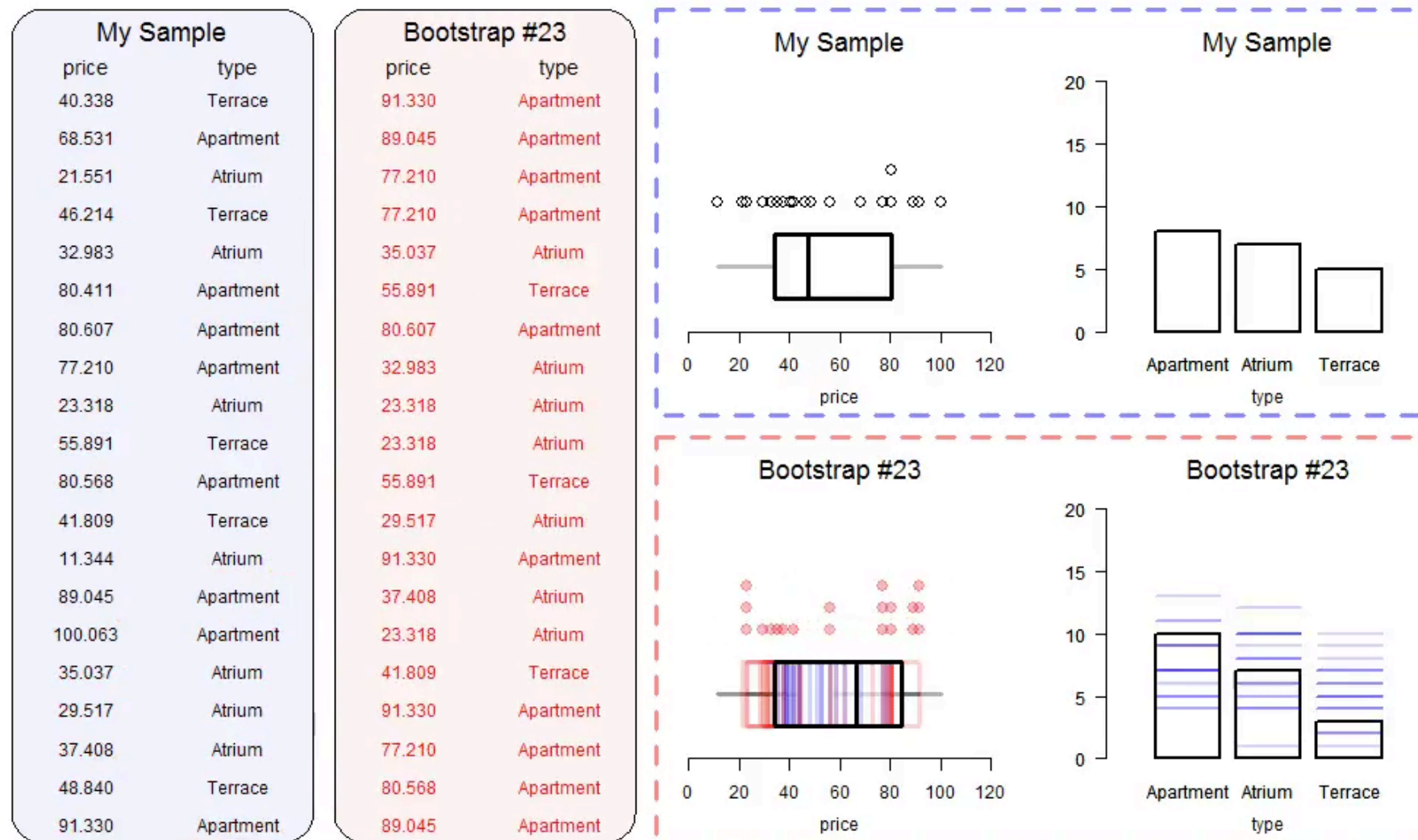
Definition:

Non-parametric Bootstrap seeks to estimate an underlying probability distribution by resampling observed values.



In Action

... what's happening ...



Source

Why Bootstrap?

... where's the infinite data ???

1

- Sampling data takes both **time** and **money**

2

- No ability to make an assumption about the sample's underlying distribution...

3

- Inference done with asymptotic theory on samples may not make sense **if values have a restriction.**
(e.g. height cannot be negative or zero)

Bootstrap Terms

... describing non-parametric bootstrap statistically ...

Real World

Unknown
Distribution

Observed
Values

$$F \rightarrow X = (X_1, \dots, X_{\boxed{n}})$$

$$\hat{\theta} = s(X)$$

Test Statistic on
Observed Values

Bootstrap World

Empirical
Distribution

Bootstrap
Sample

$$\hat{F}_n \rightarrow X^{*,i} = (X_1^{*,i}, \dots, X_{\boxed{n}}^{*,i})$$

$$\hat{\theta}^{*,i} = s(X^{*,i})$$

Test Statistic on
Bootstrapped Values

$\hat{\theta}^*$

to

$\hat{\theta}$

is like

$\hat{\theta}$

to

θ

Test Statistic on
Observed



Test Statistic on
Bootstrapped Values

Population Statistic

Resampling

generating fictional data from real data ...

Original Data

$$F \rightarrow X = (X_1, \dots, X_n)$$

	id	sex	height
X_1	1	M	6.1
X_2	2	F	5.5
X_3	3	F	5.2
X_4	4	M	5.6
X_5	5	M	5.9

Resampled Data

$$\hat{F}_n \rightarrow X^* = (X_1^*, \dots, X_n^*)$$

id	sex	height
2	F	5.5
3	F	5.2
3	F	5.2
1	M	6.1
4	M	5.6

$X_1^{*,1}$
 $X_2^{*,1}$
 $X_3^{*,1}$
 $X_4^{*,1}$
 $X_5^{*,1}$

...
...

New
Samples

id	sex	height
4	M	5.6
4	M	5.6
2	F	5.5
1	M	6.1
2	F	5.5

$X_1^{*,n}$
 $X_2^{*,n}$
 $X_3^{*,n}$
 $X_4^{*,n}$
 $X_5^{*,n}$

Sample **with** Replacement

Strategy

... how to roll a **non-parametric** bootstrap ...

Step 1: Obtain a sample of the population

Step 2: Using resampling technique, construct

$$\hat{F}_n^{iid} \sim X^{*,i} = \left(X_1^{*,i}, X_2^{*,i}, \dots, X_n^{*,i} \right)$$

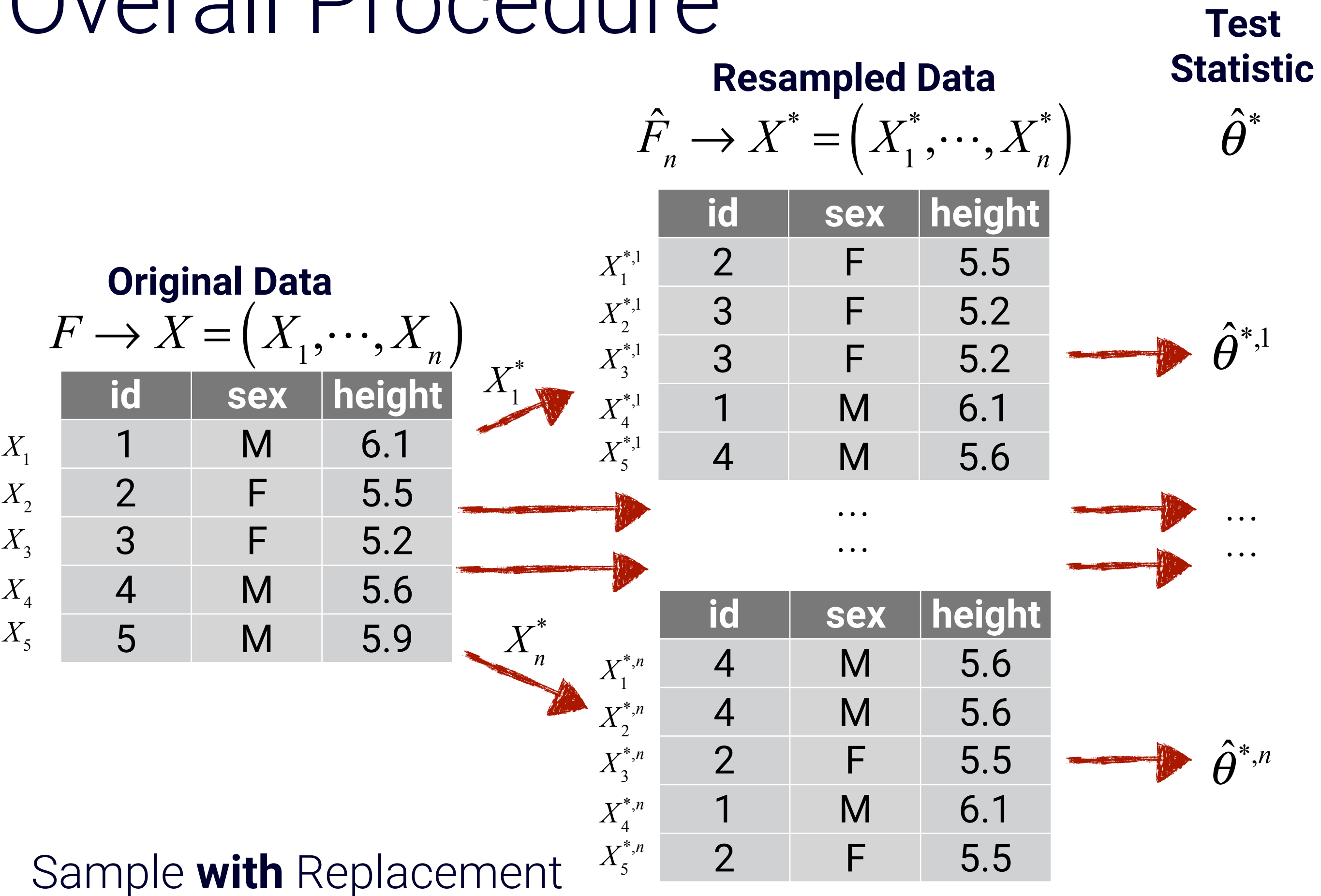
that contains the same number of observations as

$$F^{iid} \sim X = \left(X_1, X_2, \dots, X_n \right)$$

Step 3: Compute a statistic on the resampled data as $\hat{\theta}^* = s(X^{*,i})$

Step 4: Repeat **Steps 2 - 3** until i matches required number of iterations.

Overall Procedure



Implementation

... of **non**-parametric bootstrap ...

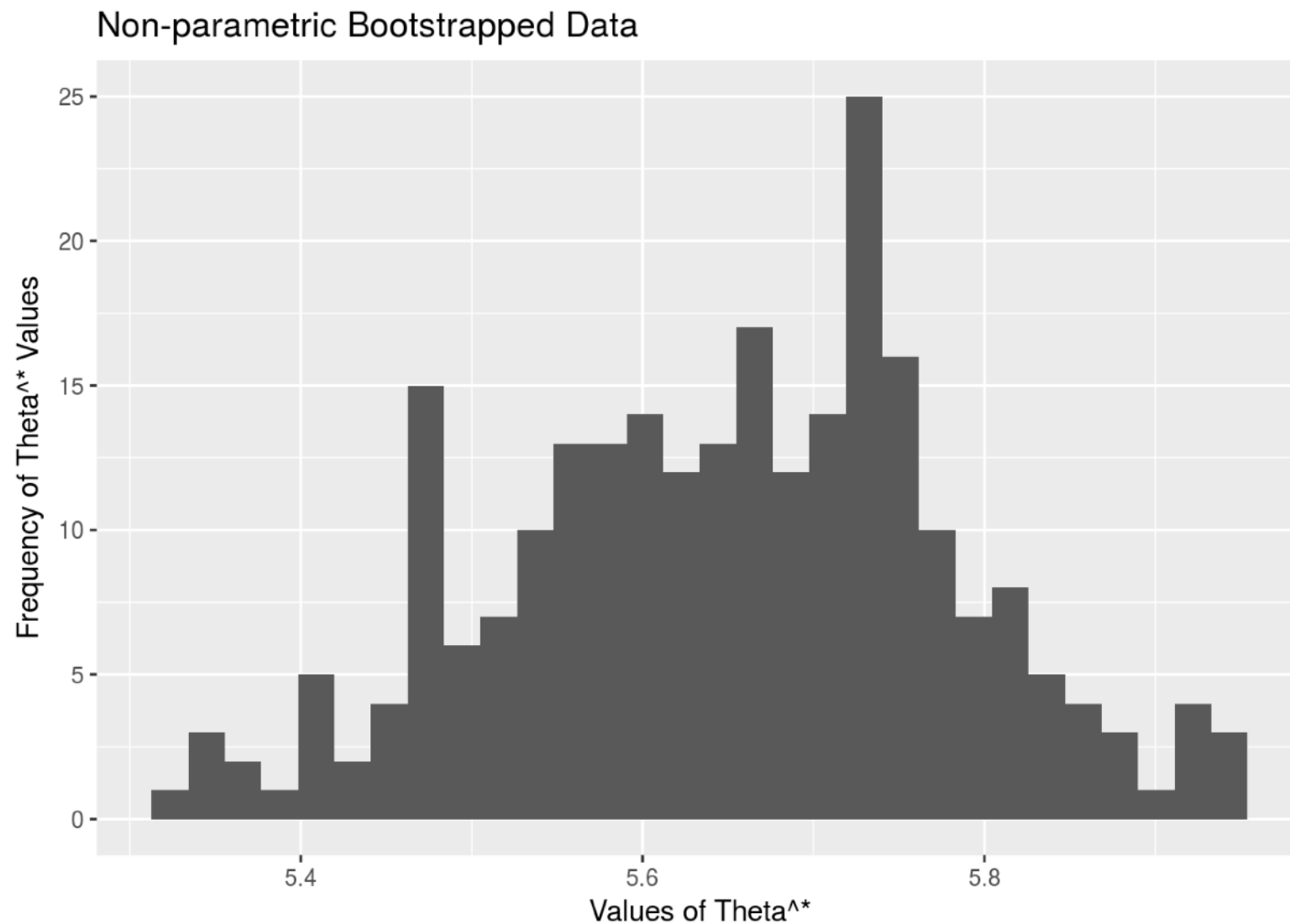
```
sample_data = ???                                # Step 1: Obtain samples from population
theta_hat = mean(sample_data$var)                # Compute the mean for the data
n_obs = nrow(sample_data)                        # Length of data
boot_iter = 250L                                # Number of bootstrap iterations
theta_star = rep(NA, boot_iter)                  # Bootstrapped estimate of theta

for (i in seq_len(boot_iter)) {
  set.seed(11882 + i)                            # Set seed for reproducibility
  # Step 2: Randomly sample observations positions from 1 to n_obs
  indexes = sample(n_obs, n_obs, replace = TRUE)
  # Extract out the observation positions
  sample_data_star = sample_data[indexes,, drop = FALSE]

  # Step 3: Compute the desired statistic on the bootstrapped values
  theta_star[i] = mean(sample_data_star$var)
} # Step 4: Repeat until i matches boot_iter
```

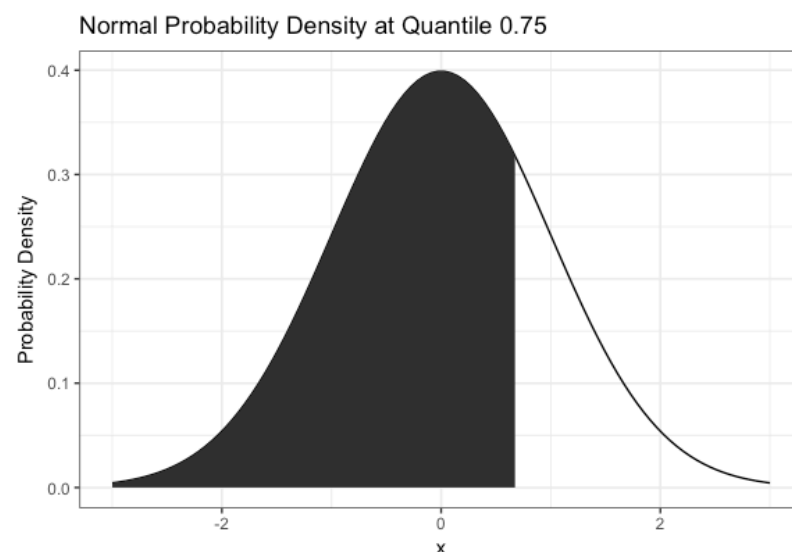
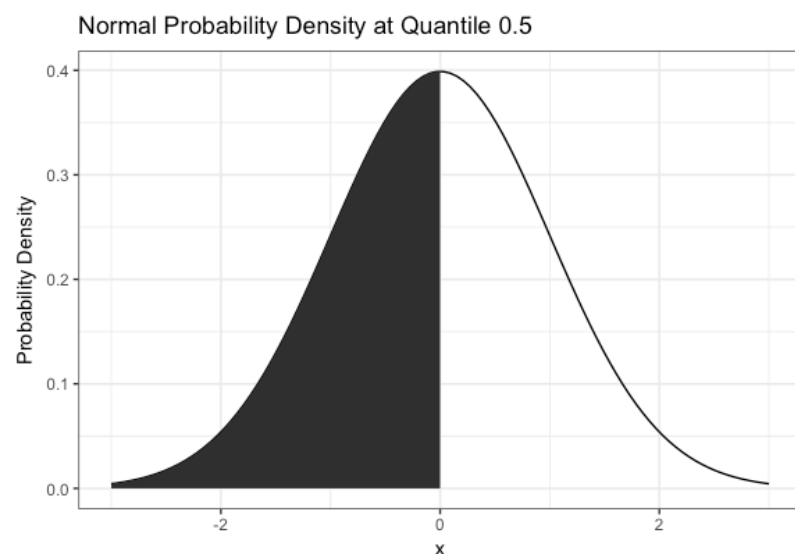
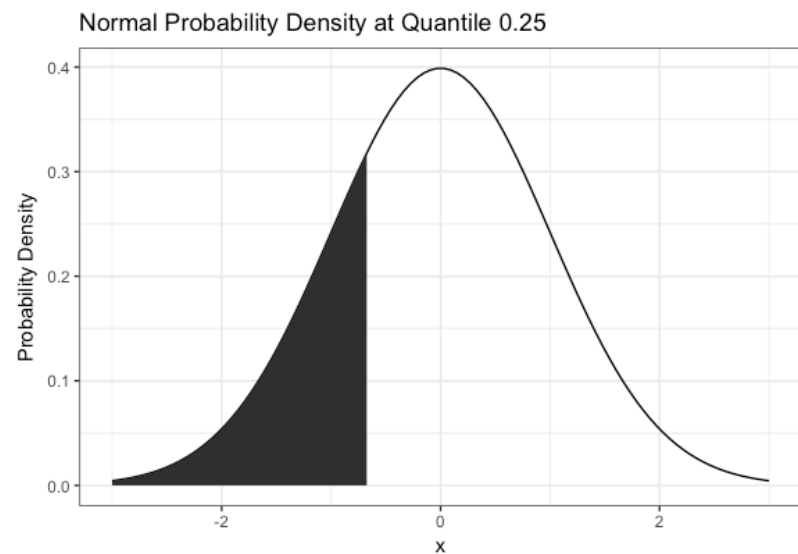
Bootstrapped Distribution

... values of bootstrapped statistic ...



Quantiles / Percentiles

... value of the distribution at p -th location ...



```
# Sample Data
```

```
x = c(1, 2, 3, 4, 5, 6)
```

```
# Value at ordered point in  
# distribution
```

```
quantile(x,  
  probs = c(0.25, 0.5, 0.75, 1)  
)
```

```
# 25% 50% 75% 100%
```

```
# 2.25 3.50 4.75 6.00
```

```
# Median is the 50% quantile
```

```
median(x)
```

```
# [1] 3.5
```


Percentile CIs

Computes a custom confidence interval with data...

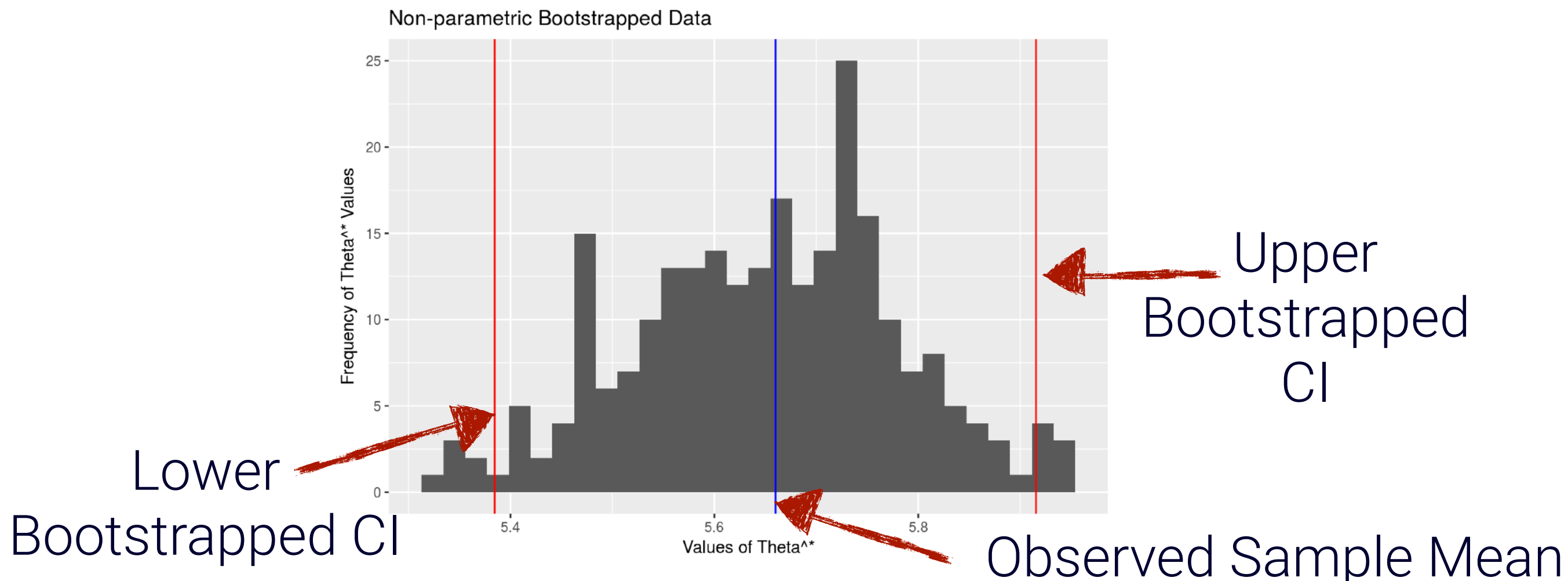
Significance level

alpha = 0.05

Under alpha = 0.05, we are retrieving quantiles for 0.025 and 0.975

```
ci_range = quantile(theta_star, probs = c(alpha / 2, 1 - alpha / 2))
```

[1] low high



Your Turn

Modify the bootstrap so that it computes the **standard deviation** of **Sepal.Width** in the **iris** data set.

Recall: The **sd()** function provides the standard deviation

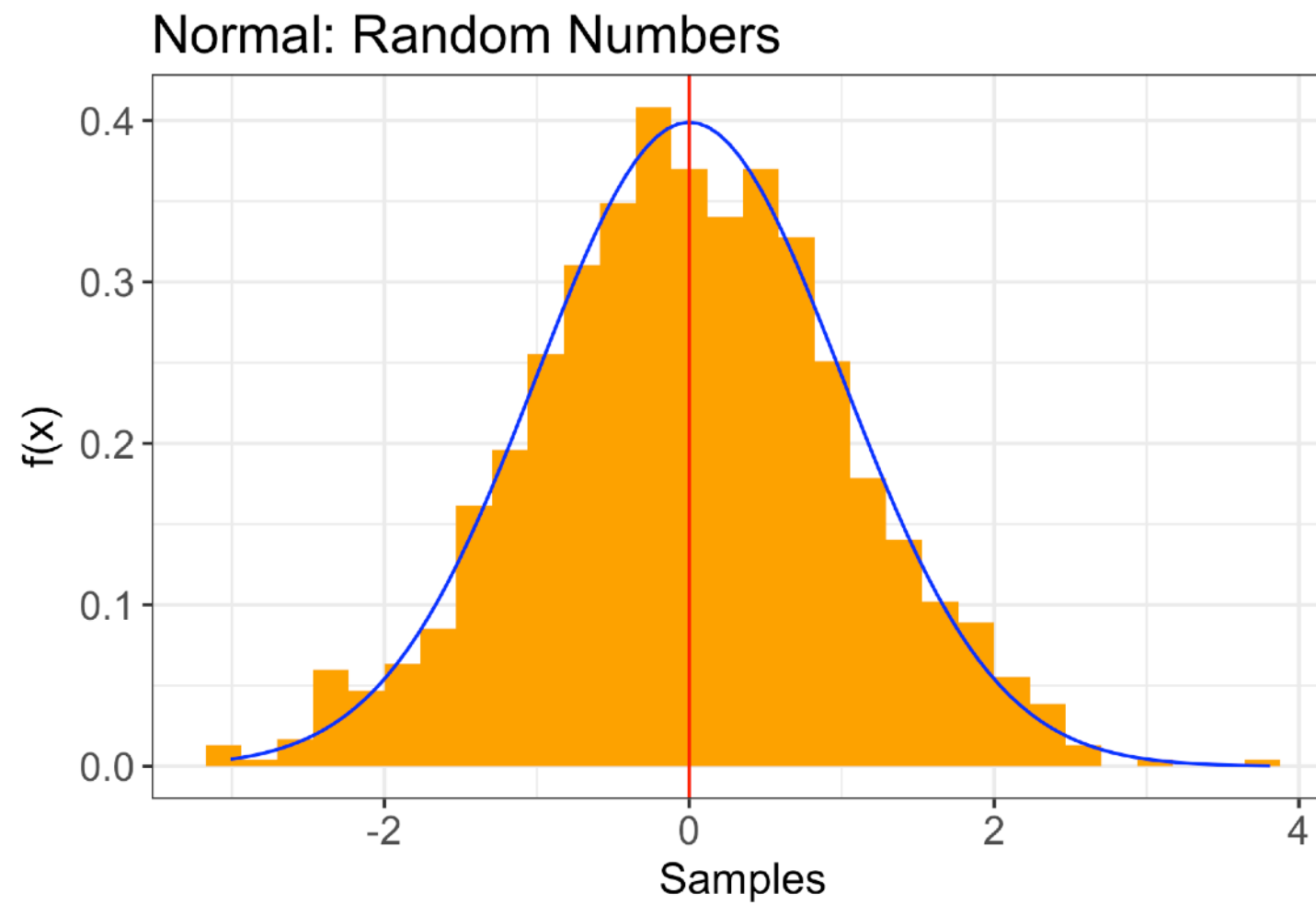
$$\sigma = sd(x) = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\bar{x} = mean(x) = \frac{1}{n} \sum_{i=1}^n x_i$$

Parametric Bootstrap

Definition:

Parametric Bootstrap seeks to estimate parameters under the assumption they belong to a specific family of probability distributions.



Strategy

... how to roll a **parametric** bootstrap ...

Step 1: Sample model data under a known distribution with unknown parameters θ .

$$F_{\theta} \overset{iid}{\sim} X = (X_1, X_2, \dots, X_n)$$

Step 2: Compute the statistic under the known distribution $\hat{\theta} = s(X)$

Step 3: Sample model under $F_{\hat{\theta}} \overset{iid}{\sim} X^{*,i} = (X_1^{*,i}, X_2^{*,i}, \dots, X_n^{*,i})$
via $\hat{\theta} = s(X)$

Step 4: Calculate the bootstrapped statistic $\hat{\theta}^{*,i} = s(X^{*,i})$

Step 5: Repeat **Steps 3 - 4** until i matches required number of iterations.

Implementation

... of **parametric** bootstrap ...

```
sample_values = rnorm(100)
```

```
# Step 1: Obtain samples from known  
# population distribution.
```

```
theta_mean_hat = mean(sample_values)  
theta_sd_hat = sd(sample_values)
```

```
# Step 2: Obtain statistics  
# Compute sample mean  
# Compute sample standard deviation
```

```
n_obs = length(sample_values)  
boot_iter = 250L  
theta_mean_star = rep(NA, boot_iter)  
theta_sd_star = rep(NA, boot_iter)
```

```
# Length of data  
# Number of bootstrap iterations  
# Bootstrapped estimate of mean  
# Bootstrapped estimate of standard dev
```

Implementation

... of **parametric** bootstrap ...

See previous slide for setup details...

```
for (i in seq_len(boot_iter)) {  
  set.seed(385 + i)                # Set seed for reproducibility  
  
  # Step 3: Randomly generate observations under distribution  
  sample_values_star = rnorm(n_obs, mean = theta_mean_hat, sd = theta_sd_hat )  
  
  # Step 4: Compute the desired statistic on the bootstrapped values  
  theta_mean_star[i] = mean(sample_values_star )  
  theta_sd_star[i] = sd(sample_values_star )  
  
} # Step 5: Repeat until i matches boot_iter
```

Parametric vs. Nonparametric

... what's the difference ???

Type	Distribution	Parameter
Non-parametric	<i>Unknown</i>	Unknown
Parametric	<i>Known</i>	Unknown

Redux

... highlighting the **difference** ...

Nonparametric

Unknown distribution. Sample values from observed distribution

indexes = **sample**(n_obs, n_obs, replace = TRUE)

Extract out the observation positions

sample_data_star = sample_data[indexes,, drop = FALSE]

Parametric

Known distribution. Sample values underneath estimated parameters.

sample_values_star = **rnorm**(n_obs, mean = theta_mean_hat, sd = theta_sd_hat)

Your Turn

Implement a **parametric** bootstrap that determines the **mean**, **standard deviation**, and **median** of a Poisson distribution with

$$\lambda = 3$$

(lambda)

Recap

- **Non-parametric Bootstrap**

- Resampling from a sample of the population to obtain an empirical distribution.

- **Parametric Bootstrap**

- Sampling under a known probability distribution with estimated values of the initial sample.

This work is licensed under the
Creative Commons
Attribution-NonCommercial-
ShareAlike 4.0 International
License

