

# Week 10: Bootstrap

## MATH-517 Statistical Computation and Visualization

Tomas Masak

November 25th 2022

- population  $F$
- random sample  $\mathcal{X} = \{X_1, \dots, X_N\}$  from  $F$

**Goal of Statistics:** Extract information about  $F$  using  $\mathcal{X}$ .

- characteristic of interest  $\theta = \theta(F)$

**Running Ex.:** The mean  $\theta = \mathbb{E}X_1 = \int x dF(x)$ .

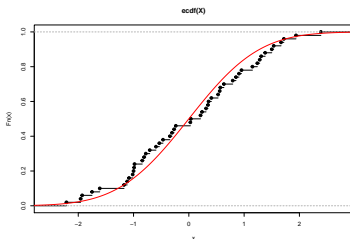
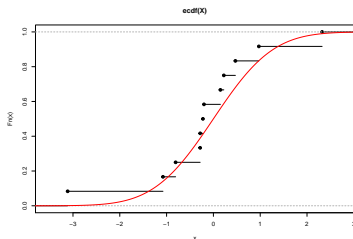
$\Delta$

$F$  can be estimated:

- parametrically
  - assuming  $F \in \{F_\lambda \mid \lambda \in \Lambda \subset \mathbb{R}^p\}$  for some integer  $p$ , take  $\hat{F} = F_{\hat{\lambda}}$  for an  $\hat{\lambda}$  estimator of the parameter vector  $\lambda$
- non-parametrically
  - by the ECDF, i.e.  $\hat{F} = \hat{F}_N$  where  $\hat{F}_N(x) = \frac{1}{N} \sum_{n=1}^N \mathbb{I}_{[X_n \leq x]}$

# ECDF

```
edf_plot <- function(N){  
  X <- rnorm(N)  
  EDF <- ecdf(X)  
  plot(EDF)  
  x <- seq(-4,4,by=0.01)  
  points(x,pnorm(x),type="l",col="red")  
}  
set.seed(517)  
edf_plot(12)  
edf_plot(50)
```



- population  $F$
- random sample  $\mathcal{X} = \{X_1, \dots, X_N\}$  from  $F$
- characteristic of interest  $\theta = \theta(F)$

**Running Ex.:** The mean  $\theta = \mathbb{E}X_1 = \int x dF(x)$ .

- parametrically: MLE
- non-parametrically:  $\hat{\theta} := \int x d\hat{F}_N(x) = \frac{1}{N} \sum_{n=1}^N X_n$

$\Delta$

- population  $F$
- random sample  $\mathcal{X} = \{X_1, \dots, X_N\}$  from  $F$
- characteristic of interest  $\theta = \theta(F)$
- sample characteristic  $\hat{\theta} = \theta(\hat{F})$
- **sampling distribution** of  $\hat{\theta}$ 
  - quantiles of sampling distribution needed for CIs or testing
  - bias or MSE needed to rate the estimator - all characteristics of sampling distr.

**Running Ex.:** The mean  $\theta = \mathbb{E}X_1 = \int x dF(x)$ .

- non-parametrically:  $\hat{\theta} := \int x d\hat{F}_N(x) = \frac{1}{N} \sum_{n=1}^N X_n$
- if  $F$  is Gaussian,  $\hat{\theta} \sim \mathcal{N}(\theta, \frac{\sigma^2}{N})$  is the sampling distribution
  - without Gaussianity, there is still a sampling distribution, we just don't know what it is  $\Delta$

# Intro

Statistics is about the **sampling distribution**, which is given by the sampling process (part of which is  $F$  itself)

- if we controlled the sampling process, we would approximate the sampling distribution by Monte Carlo

**The Bootstrap Idea:** Resampling process from  $\hat{F}$  can mimic the sampling process from  $F$  itself.

- since  $\hat{F}$  is known, the resampling distribution can be studied
  - or approximated by Monte Carlo

Sampling (real world):  $F \implies X_1, \dots, X_N \implies \hat{\theta} = \theta(\hat{F})$

Resampling (bootstrap world):  $\hat{F} \implies X_1^*, \dots, X_N^* \implies \hat{\theta}^* = \theta(\hat{F}^*)$

## Running Ex.

- $X_1, \dots, X_N \stackrel{\text{i.i.d.}}{\sim} F$  and  $\theta = \theta(F) = \int x dF$
- we want  $\hat{\theta}(\alpha)$  such that  $P(\theta \leq \hat{\theta}(\alpha)) = \alpha$ .

### 1 Exact CI. Assuming Gaussianity,

$$T = \sqrt{N} \frac{\bar{X}_N - \theta}{\hat{\sigma}} \sim t_{n-1} \quad \Rightarrow \quad P(-T \leq t_{n-1}(\alpha)) = \alpha$$

and so we get a CI with exact coverage by expressing  $\theta$  from the inequality  $T \leq t_{n-1}(\alpha)$ :

$$\theta \leq \bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}} t_{n-1}(\alpha) =: \hat{\theta}(\alpha).$$

### 2 Asymptotic CI. Assuming only $\mathbb{E}X_1^2 < \infty$ , $T \xrightarrow{d} \mathcal{N}(0, 1)$ and thus

$$P(\theta \leq \hat{\theta}(\alpha)) \approx \alpha \quad \text{for} \quad \hat{\theta}(\alpha) = \bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}} z(\alpha),$$

## Running Ex.

③ **Bootstrap CI.** Let  $\mathbb{E}X_1^2 < \infty$  and  $X_1^*, \dots, X_N^*$  be a resample from the ECDF  $\hat{F}_N$

- set up the bootstrap statistic  $T^* = \sqrt{N} \frac{\bar{X}_N^* - \bar{X}_N}{\hat{\sigma}^*}$
- denote by  $q^*(\alpha)$  the quantile of  $T^*$
- instead of  $\hat{\theta}(\alpha) = \bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}} z(\alpha)$ , consider  $\hat{\theta}^*(\alpha) = \bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}} q^*(\alpha)$

From Edgeworth expansions (complicated!):

$$P_F(T \leq x) = \Phi(x) + \frac{1}{\sqrt{N}} a(x) \phi(x) + \mathcal{O}\left(\frac{1}{N}\right)$$
$$P_{\hat{F}_N}(T^* \leq x) = \Phi(x) + \frac{1}{\sqrt{N}} \hat{a}(x) \phi(x) + \mathcal{O}\left(\frac{1}{N}\right)$$

where  $\hat{a}(x) - a(x) = \mathcal{O}(N^{-1/2})$ .



## Running Ex. - Coverage Comparison

- ② **Asymptotic CI.** By the Berry-Essen theorem

$$P(T \leq x) - P(\mathcal{N}(0, 1) \leq x) = \mathcal{O}\left(\frac{1}{\sqrt{N}}\right) \quad \text{for all } x$$
$$\Rightarrow P\left(\theta \leq \underbrace{\bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}}z(\alpha)}_{=\hat{\theta}(\alpha)}\right) = \alpha + \mathcal{O}\left(\frac{1}{\sqrt{N}}\right).$$

I.e. the coverage of the asymptotic CI is exact up to  $\mathcal{O}(N^{-1/2})$ .

- ③ **Bootstrap CI.** From Edgeworth expansions

$$P_F(T \leq x) - P_{\hat{F}_N}(T^* \leq x) = \mathcal{O}\left(\frac{1}{N}\right)$$
$$\Rightarrow P\left(\theta \leq \underbrace{\bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}}q^*(\alpha)}_{=\hat{\theta}^*(\alpha)}\right) = \alpha + \mathcal{O}\left(\frac{1}{N}\right),$$

I.e. the coverage of the bootstrap CI is exact up to  $\mathcal{O}(N^{-1})$ .

# How is this possible?

- we got a better interval than that from CLT by resampling our data once
  - resampling once  $\equiv$  discarding information

# How is this possible?

- we got a better interval than that from CLT by resampling our data once
  - resampling once  $\equiv$  discarding information
- however, we did “theoretical” resampling
- in practice, we don't know  $q^*(\alpha)$ , we have to approximate it
  - e.g. by Monte Carlo  $\equiv$  resampling many times
  - but still, how can we gain information by resampling?

Baron Munchausen (half-fictional character)

- rode a cannonball
- traveled to the Moon (18th century)
- got out from the bottom of the lake by pulling his *bootstraps*



# Another Example

- $X_1, \dots, X_N$  i.i.d. with  $\mathbb{E}|X_1|^3 < \infty$
- characteristic of interest:  $\theta = \mu^3$ , where  $\mu = \mathbb{E}X_1$
- empirical estimator:  $\hat{\theta} = (\int x d\hat{F}_N)^3 = (\bar{X}_N)^3$  is biased
- bootstrap: estimate the bias  $b := \text{bias}(\hat{\theta}) = \mathbb{E}\hat{\theta} - \theta$  as  $\hat{b}^*$
- bias-corrected estimator:  $\hat{\theta}_b^* = \hat{\theta} - \hat{b}^*$  ... provably smaller bias?

## Another Example

- $X_1, \dots, X_N$  i.i.d. with  $\mathbb{E}|X_1|^3 < \infty$
- characteristic of interest:  $\theta = \mu^3$ , where  $\mu = \mathbb{E}X_1$
- estimator:  $\hat{\theta} = (\int x d\hat{F}_N)^3 = (\bar{X}_N)^3$  is biased

$$\mathbb{E}\hat{\theta} = \mathbb{E}\bar{X}_N^3 = \mathbb{E}\left[\mu + N^{-1} \sum_{n=1}^N (X_n - \mu)\right]^3 = \mu^3 + \underbrace{N^{-1}3\mu\sigma^2 + N^{-2}\gamma}_{=b},$$

- bootstrap: estimate the bias  $b := \text{bias}(\hat{\theta}) = \mathbb{E}\hat{\theta} - \theta$  as  $\hat{b}^*$

$$\mathbb{E}_{\hat{F}_N}\hat{\theta}^* = \mathbb{E}_{\hat{F}_N}(\bar{X}_N^*)^3 = \mathbb{E}_{\hat{F}_N}\left[\bar{X}_N + N^{-1} \sum_{n=1}^N (X_n^* - \bar{X}_N)\right]^3 = \bar{X}_N^3 + \underbrace{N^{-1}3\bar{X}_N\hat{\sigma}^2 + N^{-2}\hat{\gamma}}_{=\hat{b}^*},$$

- bias-corrected estimator:  $\hat{\theta}_b^* = \hat{\theta} - \hat{b}^* \dots$  provably smaller bias?

$$\mathbb{E}\hat{\theta}_b^* = \mu^3 + \underbrace{N^{-1}3[\mu\sigma^2 - \mathbb{E}\bar{X}_N\hat{\sigma}^2]}_{\mathcal{O}(N^{-1})} + \underbrace{N^{-2}[\gamma - \mathbb{E}\hat{\gamma}]}_{\mathcal{O}(N^{-1})}.$$

Bootstrap combines

- the plug-in principle, i.e. estimating the unknowns, and
- Monte Carlo principle, i.e. simulation instead of analytic calculations

What are the unknowns?

- parameters  $\Rightarrow$  parametric bootstrap
- the whole  $F$  via ECDF  $\Rightarrow$  the (standard/non-parametric) bootstrap

# The (standard/non-parametric) Bootstrap

- let  $\mathcal{X} = \{X_1, \dots, X_N\}$  be a random sample from  $F$
- characteristic of interest:  $\theta = \theta(F)$
- estimator:  $\hat{\theta} = \theta(\hat{F}_N)$ 
  - write  $\hat{\theta} = \theta[\mathcal{X}]$ , since  $\hat{F}_N$  and thus the estimator depend on the sample
- the distribution  $F_T$  of a scaled estimator  $T = g(\hat{\theta}, \theta) = g(\theta[\mathcal{X}], \theta)$  is of interest
  - e.g.  $T = \sqrt{N}(\hat{\theta} - \theta)$

# The (standard/non-parametric) Bootstrap

- let  $\mathcal{X} = \{X_1, \dots, X_N\}$  be a random sample from  $F$
- characteristic of interest:  $\theta = \theta(F)$
- estimator:  $\hat{\theta} = \theta(\hat{F}_N)$ 
  - write  $\hat{\theta} = \theta[\mathcal{X}]$ , since  $\hat{F}_N$  and thus the estimator depend on the sample
- the distribution  $F_T$  of a scaled estimator  $T = g(\hat{\theta}, \theta) = g(\theta[\mathcal{X}], \theta)$  is of interest
  - e.g.  $T = \sqrt{N}(\hat{\theta} - \theta)$

The workflow of the bootstrap is as follows for some  $B \in \mathbb{N}$  (e.g.  $B = 1000$ ):

Data	Resamples
$\mathcal{X} = \{X_1, \dots, X_N\}$	$\begin{cases} \mathcal{X}_1^* = \{X_{1,1}^*, \dots, X_{1,N}^*\} & \implies & T_1^* = g(\theta[\mathcal{X}_1^*], \theta[\mathcal{X}]) \\ \vdots & & \vdots \\ \mathcal{X}_B^* = \{X_{B,1}^*, \dots, X_{B,N}^*\} & \implies & T_B^* = g(\theta[\mathcal{X}_B^*], \theta[\mathcal{X}]) \end{cases}$

$F_T$  now estimated by  $\hat{F}_{T,B}^*(x) = B^{-1} \sum_{b=1}^B \mathbb{I}_{[T_b^* \leq x]}$

- any characteristic of  $F_T$  can be estimated by the char. of  $\hat{F}_{T,B}^*(x)$



## Running Ex. Again

- $X_1, \dots, X_N \stackrel{\text{d}}{\sim} F$  and  $\theta = \theta(F) = \int x dF$
- we want  $\hat{\theta}(\alpha)$  such that  $P(\theta \leq \hat{\theta}(\alpha)) = \alpha$ .
- ③ **Bootstrap CI.** Let  $\mathbb{E}X_1^2 < \infty$  and  $X_1^*, \dots, X_N^*$  be a resample from the ECDF  $\hat{F}_N$ 
  - set up the bootstrap statistic  $T^* = \sqrt{N} \frac{\bar{X}_N^* - \bar{X}_N}{\hat{\sigma}^*}$
  - denote by  $q^*(\alpha)$  the quantile of  $T^*$
  - take  $\left( -\infty, \bar{X}_N + \frac{\hat{\sigma}}{\sqrt{N}} q^*(\alpha) \right)$

In practice,  $q^*(\alpha)$  approximated by Monte Carlo:

Data

Resamples

$$\mathcal{X} = \{X_1, \dots, X_N\} \implies \begin{cases} \mathcal{X}_1^* = \{X_{1,1}^*, \dots, X_{1,N}^*\} & \implies T_1^* \\ \vdots & \vdots \\ \mathcal{X}_B^* = \{X_{B,1}^*, \dots, X_{B,N}^*\} & \implies T_B^* \end{cases}$$

$\Rightarrow$  take  $q^*(\alpha)$  as the sample quantile of  $T_1^*, \dots, T_B^*$

# Running Ex. Specific

```
lambda <- 2
N <- 100
X <- rexp(N,lambda)

( CI_aspytotic <- mean(X) + qnorm(0.95)*sd(X)/sqrt(N) )

## [1] 0.5645591

Tstar <- function(Xstar,X){
  return( (mean(Xstar)-mean(X))/sd(Xstar)*sqrt(N))
}
B <- 10^3
boot_ind <- sample(1:N, size=N*B, replace=T)
boot_data <- matrix(X[boot_ind],ncol=B)
Tstars <- rep(0,B)
for(b in 1:B){
  Tstars[b] <- Tstar(boot_data[,b],X)
}
( CI_boot <- mean(X) + quantile(Tstars,0.95)*sd(X)/sqrt(N) )

##          95%
## 0.5557836
```

# The Bootstrap

- now we know what the bootstrap is
  - the scheme is very simple, though a bit mysterious, spawning questions:
- when does it work? (“work” = consistency)
- when does it give us something extra? (e.g. faster rates)

# Consistency for Smooth Transformation of the Mean

Bootstrap setup in practice:

- $T$  is the scaled estimator with unknown distribution  $F_T$
- the bootstrap statistic  $T^*$  has  $F_T^*$  also unknown
- the Monte Carlo proxy  $F_{T,B}^*$  is used instead of  $F_T^*$

Glivenko-Cantelli:

$$\sup_x \left| \widehat{F}_{T,B}^*(x) - F_T^*(x) \right| \xrightarrow{a.s.} 0 \quad \text{as } B \rightarrow \infty.$$

**Question:**  $\sup_x \left| F_T^*(x) - F_T(x) \right| \rightarrow 0$  for  $N \rightarrow \infty$ ?

**Theorem:** Let  $\mathbb{E}X_1^2 < \infty$  and  $T = h(\bar{X}_N)$ , where  $h$  is continuously differentiable at  $\mu := \mathbb{E}X_1$  and such that  $h(\mu) \neq 0$ . Then

$$\sup_x \left| F_T^*(x) - F_T(x) \right| \xrightarrow{a.s.} 0 \quad \text{as } N \rightarrow \infty.$$

# Remarks

- bootstrap should not be used blindly
  - verification via theory
  - and/or via simulations
- folk knowledge
  - bootstrap “works” when we have non-degenerate asymptotic normality
  - bootstrap “doesn’t work” when working with order statistics, extremes, non-smooth transformations, non-i.i.d. regimes (e.g. time series), etc.
- bootstrap replaces analytic calculations (in particular the Delta method), but showing that it actually works requires even deeper analytic calculations
- faster rates can be achieved by bootstrap
  - hard to prove, but often happens e.g. when working with a skewed distribution
- how many Monte Carlo draws needed?
  - $B = 10^2$  is enough for variance estimation (next week)
  - $B = 10^3$  is taken most commonly
  - $B = 10^4$  better for small/large quantiles