

# 1 Nonparametric Bootstrap

Let  $Y_1, \dots, Y_n \sim F$  with pdf  $f(y)$ . Recall, the empirical cdf is defined as

$$F_n(y) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}(Y_i \leq y) \quad y \in \mathbb{R}^d$$

↑

MLE of  $F$  and as  $n \rightarrow \infty$ ,  $F_n \rightarrow F$ .

Theoretical: Sample  $Y \sim F$ , use  $Y_1, \dots, Y_n$  to compute  $F_n$

Bootstrap: Sample  $Y^* \sim F_n$ , use  $Y_1^*, \dots, Y_n^*$  to compute  $F_n^*$

- The idea behind the nonparametric bootstrap is to sample many data sets from  $F_n(y)$ , which can be achieved by resampling from the data **with replacement**. (iid case). ↗ of size  $n$ .

How many possible Bootstrap samples?  $n^n$

Are  $Y_1^*, \dots, Y_n^*$  independent?

$$P(Y_1^* = a, Y_2^* = b) = \frac{\sum_{i=1}^n \mathbb{I}(Y_i^* = a)}{n} + \frac{\sum_{i=1}^n \mathbb{I}(Y_i^* = b)}{n} = P(Y_1^* = a) P(Y_2^* = b) \quad \underline{\underline{\text{Yes}}}$$

Do we always want this?  
(more later...)

```

# observed data
x <- c(2, 2, 1, 1, 5, 4, 4, 3, 1, 2)

# create 10 bootstrap samples
x_star <- matrix(NA, nrow = length(x), ncol = 10)
for(i in 1:10) {
  x_star[, i] <- sample(x, length(x), replace = TRUE)
}
x_star

```

```

##            $\bar{x}^{*(1)}$       [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]  $\bar{x}^{*(10)}$  [,10]
## [1,]      1      2      4      1      2      1      2      3      3      4
## [2,]      4      4      1      1      1      2      2      1      2      1
## [3,]      2      2      2      4      5      4      4      5      1      4
## [4,]      4      4      2      5      2      4      5      5      1      3
## [5,]      2      1      5      1      3      2      4      2      4      4
## [6,]      4      4      2      1      4      4      4      3      1      2
## [7,]      1      1      2      1      2      1      2      2      3      1
## [8,]      4      4      1      3      3      3      5      1      2      4
## [9,]      4      1      2      3      2      1      2      1      4      2
## [10,]     3      4      5      1      5      4      5      2      4      1

```

```

# compare mean of the sample to the means of the bootstrap samples
mean(x)

```

```
## [1] 2.5
```

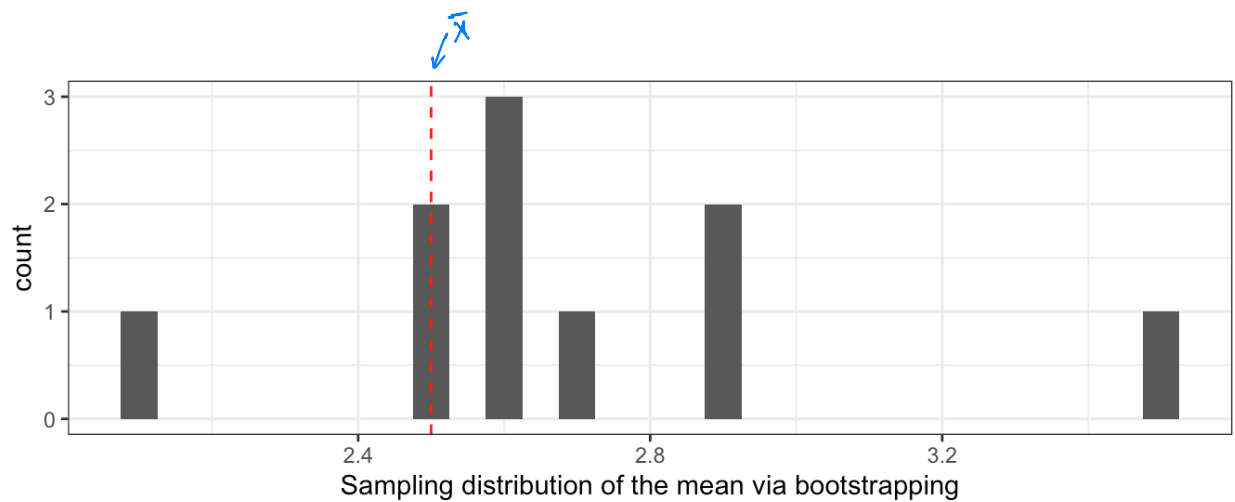
```
colMeans(x_star)
```

```
## [1] 2.9 2.7 2.6 2.1 2.9 2.6 3.5 2.5 2.5 2.6
```

```

ggplot() +
  geom_histogram(aes(colMeans(x_star)), binwidth = .05) +
  geom_vline(aes(xintercept = mean(x)), lty = 2, colour = "red") +
  xlab("Sampling distribution of the mean via bootstrapping")

```



## 1.1 Algorithm

**Goal:** estimate the sampling distribution of a statistic based on observed data  $y_1, \dots, y_n$ .

Let  $\theta$  be the parameter of interest and  $\hat{\theta}$  be an estimator of  $\theta$ . Then,

For  $b=1, \dots, B$

① Sample  $y^{*(b)} = (y_1^{*(b)}, \dots, y_n^{*(b)})$  by sampling w/ replacement from  $(y_1, \dots, y_n)$   
(i.e. Sample from  $F_n$ )

②  $\hat{\theta}^{(b)} = \hat{\theta}(y^{*(b)})$   
 $\uparrow$   
 estimate of  $\theta$  based on  $b^{\text{th}}$  bootstrap sample.

Using  $\hat{\theta}^{(1)}, \dots, \hat{\theta}^{(B)}$  we can

- estimate the sampling dsn of  $\hat{\theta}_n$  (via histogram, density estimator)
- estimate SE of  $\hat{\theta}$
- estimate bias of  $\hat{\theta}$
- estimate a CI (many ways).

etc.

## 1.2 Justification for iid data

Suppose  $Y_1, \dots, Y_n$  are iid with  $EY_1 = \mu \in \mathbb{R}$ ,  $\text{Var}(Y_1) = \sigma^2 \in (0, \infty)$ . Let's approximate the distribution of  $T_n = \sqrt{n}(\bar{Y}_n - \mu)$  via the bootstrap.

**Theorem:** If  $Y_1, Y_2, \dots$  are iid with  $\text{Var}(Y_1) = \sigma^2 \in (0, \infty)$ , then  $\sup_{y \in \mathbb{R}} |P(T_n \leq y) - P_*(T_n^* \leq y)| \equiv \Delta_n \rightarrow 0$  as  $n \rightarrow \infty$  almost surely (a.s.).

Given  $\underline{y} = \{y_1, \dots, y_n\}$  draw  $Y_1^*, \dots, Y_n^*$  bootstrap sample. Then,

bootstrap probability  $\rightarrow P_*(Y_i^* = y_i) = P(Y_i = y_i | \underline{y}) = \frac{1}{n} \quad 1 \leq i \leq n.$

The bootstrap version of our statistic  $T_n$  is  $T_n^* = \sqrt{n}(\bar{Y}_n^* - E_* Y_i^*) = \sqrt{n}(\bar{Y}_n^* - \bar{y}_n)$

bootstrap expected value.  $\rightarrow$  where  $E_*(Y_i^*) = E[Y_i^* | \underline{y}] = \sum_{i=1}^n \frac{1}{n} y_i = \bar{y}_n$  also  $E_*(\bar{Y}_n^*) = E_*(\frac{1}{n} \sum_{i=1}^n Y_i^*) = \frac{1}{n} \sum_{i=1}^n E_* Y_i^* \stackrel{iid}{=} \bar{y}_n$

Also  $P_*(T_n^* \leq y) = P(T_n^* \leq y | \underline{y})$  approximates  $P(T_n \leq y) \quad y \in \mathbb{R}$  (Theorem).  $\leftarrow$  exists because dsn of  $Y_1^*, \dots, Y_n^*$  exists, but hard to compute directly b/c  $n^n$  possible bootstrap samples.  $\Rightarrow$  use simulation.

The proof of this theorem requires two facts:

- i. (Berry-Esseen Lemma) Let  $Y_1, \dots, Y_n$  be independent with  $EY_i = 0$  and  $E|Y_i|^3 < \infty$  for  $i = 1, \dots, n$ . Let  $\sigma_n^2 = n\text{Var}(\bar{Y}_n) = n^{-1} \sum_{i=1}^n EY_i^2 > 0$ . Then,

$$\sup_{y \in \mathbb{R}} \left| P\left(\frac{\sqrt{n}\bar{Y}_n}{\sigma_n} \leq y\right) - \Phi(y) \right| = \sup_{x \in \mathbb{R}} \left| P(\sqrt{n}\bar{Y}_n \leq x) - \Phi\left(\frac{x}{\sigma_n}\right) \right| \leq \frac{2.75}{n^{3/2}\sigma_n^3} \sum_{i=1}^n E|Y_i|^3.$$

$\uparrow$  standard normal cdf.

- ii. (Marcinkiewicz-Zygmund SLLN) Let  $X_i$  be a sequence of iid random variables with  $E|X_i|^p < \infty$  for  $p \in (0, 2)$ . Then, for  $S_n = \sum_{i=1}^n X_i$ ,

$$\frac{1}{n^{1/p}}(S_n - nc) \rightarrow 0 \text{ as } n \rightarrow \infty \text{ almost surely } (*)$$

for any  $c \in \mathbb{R}$  if  $p \in (0, 1)$  and for  $c = EX_1$  if  $p \in [1, 2)$ . If  $(*)$  holds for some  $c \in \mathbb{R}$ , then  $E|X_1|^p < \infty$ .

Specifically, we will use that if  $\{Y_i\}$  are iid w/  $EY_i^2 < \infty$ , then

$$\frac{1}{n^{3/2}} \sum_{i=1}^n |Y_i|^3 \rightarrow 0 \text{ as } n \rightarrow \infty \text{ a.s.}$$

Letting  $X_i = |Y_i|^3$  because  $E|X_i|^p = E|Y_i|^{3p} < \infty$  for  $p = 2/3$ , we may take  $c=0$

Proof:

Write  $\sup_{y \in \mathbb{R}} |P(T_n \leq y) - P_*(T_n^* \leq y)| \leq \underbrace{\sup_{y \in \mathbb{R}} |P(T_n \leq y) - \Phi(y/\sigma)|}_{\tilde{\Delta}_n} + \underbrace{\sup_{y \in \mathbb{R}} |P_*(T_n^* \leq y) - \Phi(y/\sigma)|}_{\Delta_n}$

$\tilde{\Delta}_n \rightarrow 0$  by CLT since  $Y_1, \dots, Y_n$  iid,  $EY_1^2 < \infty$ .

Note that

$$\begin{aligned} \sigma_{n^*}^2 &\equiv n \text{Var}_*(\bar{Y}_n^*) = n \text{Var}_*\left(\frac{1}{n} \sum_{i=1}^n Y_i^*\right) \stackrel{iid}{=} \frac{n}{n^2} \sum_{i=1}^n \text{Var}_*(Y_i^*) = \text{Var}_* Y_i^* \\ &= E_*[(Y_i^*)^2] - [E_* Y_i^*]^2 \quad \text{where } Y_i^* = \begin{cases} Y_1 & \text{w.p. } 1/n \\ \vdots & \vdots \\ Y_n & \text{w.p. } 1/n. \end{cases} \\ &= \frac{1}{n} \sum_{i=1}^n Y_i^2 - (\bar{Y}_n)^2 \end{aligned}$$

So  $\sigma_{n^*}^2 \rightarrow EY_1^2 - (EY_1)^2 = \sigma^2$  as  $n \rightarrow \infty$  w.p. 1 by SLLN since  $EY_1^2 < \infty$ .

By the Berry Esseen Lemma on  $T_n^* = \sqrt{n}(\bar{Y}_n^* - E_* Y_i^*)$  and  $|a-b| \leq 2 \max\{|a|, |b|\} \Rightarrow |a-b|^3 \leq 8 \max\{|a|^3, |b|^3\} \leq 8(|a|^3 + |b|^3)$

$$\begin{aligned} \sup_{y \in \mathbb{R}} |P_*(T_n^* \leq y) - \Phi\left(\frac{y}{\sigma_{n^*}}\right)| &\stackrel{\text{Berry Esseen}}{\leq} \frac{2.75}{n^{3/2} \sigma_{n^*}^3} n E_* |Y_i^* - E_* Y_i^*|^3 \\ &= \frac{2.75}{n^{1/2} \sigma_{n^*}^3} \frac{1}{n} \sum_{i=1}^n |Y_i - \bar{Y}_n|^3 \\ &\leq \frac{2.75}{n^{1/2} \sigma_{n^*}^3} \frac{1}{n} \sum_{i=1}^n 8(|Y_i|^3 + |\bar{Y}_n|^3) \\ &= \frac{8(2.75)}{\sigma_{n^*}^3} \left( \frac{|\bar{Y}_n|^3}{n^{1/2}} + \frac{1}{n^{3/2}} \sum_{i=1}^n |Y_i|^3 \right) \end{aligned}$$

$\rightarrow 0$  as  $n \rightarrow \infty$  w.p. 1 since  $\bar{Y}_n \rightarrow \mu < \infty$  w.p. 1.

$\rightarrow 0$  as  $n \rightarrow \infty$  w.p. 1 by M-Z SLLN (prev. part)

Finally, use  $\sigma_{n^*}^2 \rightarrow \sigma^2$  w.p. 1  $\Rightarrow \Phi\left(\frac{y}{\sigma_{n^*}}\right) \rightarrow \Phi\left(\frac{y}{\sigma}\right)$  as  $n \rightarrow \infty$  for any  $y \in \mathbb{R}$  since  $\Phi(\cdot)$  is continuous.

Thus  $\sup_{y \in \mathbb{R}} |\Phi\left(\frac{y}{\sigma_{n^*}}\right) - \Phi\left(\frac{y}{\sigma}\right)| \rightarrow 0$  as  $n \rightarrow \infty$  w.p. 1 (Polya's theorem).

So,  $\sup_{y \in \mathbb{R}} |P_*(T_n^* \leq y) - \Phi\left(\frac{y}{\sigma}\right)| \leq \sup_{y \in \mathbb{R}} |P_*(T_n^* \leq y) - \Phi\left(\frac{y}{\sigma_{n^*}}\right)| + \sup_{y \in \mathbb{R}} |\Phi\left(\frac{y}{\sigma_{n^*}}\right) - \Phi\left(\frac{y}{\sigma}\right)| \xrightarrow{\text{w.p. 1.}} 0$