

## Chapter 5 - Beyond Monte Carlo

Bootstrap Method for Bias and Standard Deviation Estimation.

---

Prof. Alex Alvarez, Ali Raisolsadat

School of Mathematical and Computational Sciences  
University of Prince Edward Island

**Sample Data:**  $X = (X_1, X_2, \dots, X_n)$  independent, identically distributed with  $X_i \sim F$ , where  $F$  is unknown. Usually we observe a realization  $x = (x_1, x_2, \dots, x_n)$  of  $X$ .

We are interested in some statistic  $\theta = \theta(F)$ .

Assume that we have some estimator  $\hat{\theta}_n = \hat{\theta}_n(X)$  of  $\theta$ .

Today we will focus on the study of two problems:

- Estimation of the bias of  $\hat{\theta}_n$
- Estimation of the standard deviation of  $\hat{\theta}_n$

In our previous lecture we defined the empirical cumulative distribution function as:

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n 1_{(-\infty, t]}(x_i)$$

and we saw that the approximation  $F_n \approx F$  may be valid even for not so large values of  $n$ .

We will generate Bootstrap samples  $x^{*(j)} = (x_1^{*(j)}, x_2^{*(j)}, \dots, x_n^{*(j)})$  with  $j = 1, 2, \dots, N$  where the terms  $x_m^{*(j)}$  are i.i.d. random values distributed according to  $F_n$ .

An estimator  $\hat{\theta}_n = \hat{\theta}_n(X)$  of  $\theta$  satisfies the **plug-in principle** if the relation

$$\hat{\theta}_n(x) = \theta(F_n)$$

is valid for all  $x = (x_1, x_2, \dots, x_n)$ .

We know that  $\text{bias}(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta = E_F(\hat{\theta}_n) - \theta(F)$

As we do not know  $F$ , but we do know  $F_n$ , we will approximate  $\text{bias}(\hat{\theta}_n)$  with

$$E_{F_n}(\hat{\theta}_n) - \theta(F_n)$$

If  $\theta_n$  satisfies the plug-in principle, we can replace  $\theta(F_n)$  with  $\hat{\theta}_n(x)$ . For the estimations of  $E_{F_n}(\hat{\theta}_n)$  we will use the sample mean of the Bootstrap sample estimates  $\hat{\theta}_n(x^{*(j)})$  with  $j = 1, 2, \dots, N$ . Then we will have:

$$\widehat{\text{bias}}(\hat{\theta}_n) = \frac{1}{N} \sum_{j=1}^N \hat{\theta}_n(x^{*(j)}) - \hat{\theta}_n(x)$$

The following data are students grades in some course. The grades are distributed according to some unknown distribution. Use the sample median to estimate the median of the unknown distribution and estimate the bias of your estimator.

**Grades:**  $x = (84, 63, 92, 89, 80, 47, 68, 66, 56)$

**Solution:** By ordering the data in ascending order  $(47, 56, 63, 66, 68, 80, 84, 89, 92)$  we can see that the sample median is  $\hat{\theta}_9(x) = 68$ .

For the other term in the bias estimation, we will generate  $N = 1000$  Bootstrap samples as in the following code.

## R Code

```
1 x <- c(84, 63, 92, 89, 80, 47, 68, 66, 56)
2 N <- 1000
3 median_vector <- vector()
4 for(i in c(1:N)) {
5   Y <- sample(x, size=9, replace=TRUE)
6   median_vector[i] = median(Y)
7 }
8
9 average_medians <- mean(median_vector)
```

By running this code once we get

$$\widehat{bias}(\hat{\theta}_n) = \frac{1}{N} \sum_{j=1}^N \hat{\theta}_n(x^{*(j)}) - \hat{\theta}_n(x) = 71.34 - 68 = 3.34$$

The problem of estimating the standard deviation of the estimator  $\hat{\theta}_n$  under the (unknown) distribution  $F$  is simplified by considering the standard deviation under the (known) empirical distribution  $F_n$ :

$$s.d._F(\hat{\theta}_n) \approx s.d._{F_n}(\hat{\theta}_n)$$

Then, we approximate  $s.d._{F_n}(\hat{\theta}_n)$  by the empirical standard deviation of the Bootstrap estimates  $\hat{\theta}_n(x^{*(j)})$  with  $j = 1, 2, \dots, N$

The code corresponding to the estimation of the standard deviation of the median estimator in the previous example is given as follows.

## R Code

```
1 x <- c(84, 63, 92, 89, 80, 47, 68, 66, 56)
2 N <- 1000
3 median_vector <- vector()
4 for (i in c(1:N)) {
5   Y <- sample(x, size=9, replace=TRUE)
6   median_vector[i] = median(Y)
7 }
8 sd_vector <- sd(median_vector)
```

By running this code once we get that  $\widehat{s.d}(\widehat{\theta}_n) = 8.38$ .



## Homework:

With the data from the previous example, find the Bootstrap estimate of the standard deviation of the **sample mean** estimator.