

Nonparametric Bootstrap: rationale

Nonparametric Bootstrap : A statistical method that involves repeatedly resampling the observed data with replacement to generate a large number of bootstrap samples. These samples are used to estimate the sampling distribution of a statistic or to make inferences about population parameters without assuming a specific parametric distribution for the data.

1. Nonparametric bootstrap is advantageous for small sample sizes, where traditional parametric methods can be unreliable.
2. It provides a robust approach when data doesn't conform to specific distributional assumptions.
3. It is useful for estimating parameters, constructing confidence intervals, and conducting hypothesis tests without relying on distribution assumptions.
4. It is useful in scenarios involving complex data or when modeling assumptions are uncertain.

Nonparametric Bootstrap: introduction (1/2)

So suppose that we observe a sample x_1, \dots, x_n which are realizations of i.i.d. r.v. X_1, \dots, X_n , with distribution L_x and with CDF F_x unknown.. The empirical Cumulative Distribution Function (eCDF) is defined as

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{x_i \leq x}$$

and is an estimator of the true CDF F_x . Moreover, we have that

$$E[F_n(x)] = F_x$$

and

$$\text{var}(F_n(x)) = \frac{1}{n} \left((F_x)(1 - F_x) \right)$$

Nonparametric Bootstrap: introduction (2/2)

Besides, from the Law of Large Numbers (LLN), we have that

$$F_n(x) \xrightarrow{a.s.} F_x \quad \text{as } n \text{ goes to } \infty.$$

And from the Central Limit Theorem (CLT), we have that

$$F_n(x) \xrightarrow{L} N\left(F_x, \frac{1}{n}F_x(1 - F_x)\right) \quad \text{as } n \text{ goes to } \infty.$$

The underlying idea of nonparametric bootstrap is that bootstrap samples can be generated from the $F_n(x)$ when F_x is unknown since $F_n(x)$ is a consistent estimator for F_x .

General algorithm

1. B independent samples $x_{b1}^*, \dots, x_{bn}^*$, $b = 1, \dots, B$ are drawn with replacement from the data x_1, \dots, x_n . This is important, we resample with replacement from the data to obtain bootstrap samples.
2. each of the B independent samples, an estimate $\hat{\theta}_b^*$ is computed (the bootstrap replicates of $\hat{\theta}$).
3. The expectation of $\hat{\theta}$ (given $F_n(x)$) is estimated as follows:

$$E_{boot}[\hat{\theta}] = \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^*$$

4. The variance of $\hat{\theta}$ (given $F_n(x)$) is estimated as follows:

$$var_{boot}(\hat{\theta}) = \frac{1}{B-1} \sum_{b=1}^B \left(\hat{\theta}_b^* - \frac{1}{B} \sum_{b=1}^B \hat{\theta}_b^* \right)^2$$

Example

We observe a sample $x_1, \dots, x_n = 8.26, 6.33, 10.4, 5.27, 5.35, 5.61, 6.12, 6.19, 5.2, 7.01, 8.74, 7.78, 7.02, 6, 6.5, 5.8, 5.12, 7.41, 6.52, 6.21, 12.28, 5.6, 5.38, 6.6, 8.74$.

Estimate the bias and the standard error of $\hat{\theta}$, the estimate of the coefficient of variation. The coefficient of variation cv is equal to

$$cv = sd(x)/mean(x)$$

Let's perform the analysis in R and Python.

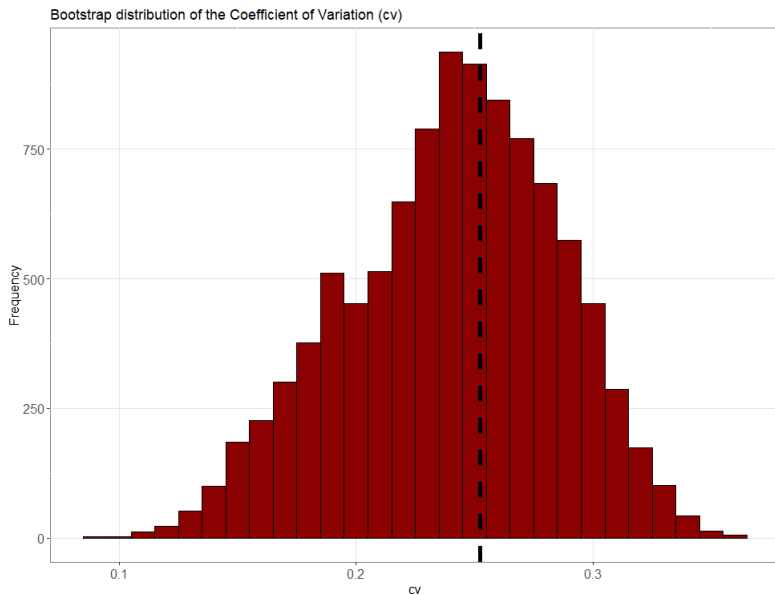
Example in R

```
1 # dataset
2 x <- c(8.26, 6.33, 10.4, 5.27, 5.35, 5.61, 6.12, 6.19, 5.2, 7.01, 8.74, 7.78,
3       7.02, 6, 6.5, 5.8, 5.12, 7.41, 6.52, 6.21, 12.28, 5.6, 5.38, 6.6, 8.74)
4
5 # coefficient of variation (CV)
6 cv <- sd(x) / mean(x)
7 cv # 0.2524712
8
9 # bootstrap
10 num_bootstraps <- 10000
11 bootstrap_cvs <- numeric(num_bootstraps)
12 set.seed(2023)
13 for (i in 1:num_bootstraps) {
14   resample <- sample(x, replace = TRUE)
15   bootstrap_cvs[i] <- sd(resample) / mean(resample)
16 }
17
18 # bias and standard error of the CV estimator
19 bias <- mean(bootstrap_cvs) - cv
20 bias # [1] -0.01216963
21 standard_error <- sd(bootstrap_cvs)
22 standard_error # [1] 0.04484109
```

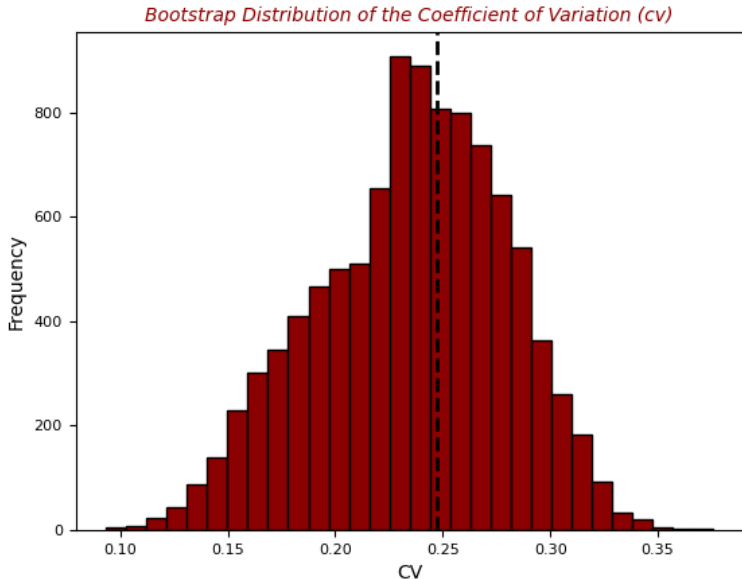
Example in Python

```
1 import numpy as np
2 import pandas as pd
3
4 # dataset
5 x = np.array([8.26, 6.33, 10.4, 5.27, 5.35, 5.61, 6.12, 6.19, 5.2, 7.01, 8.74,
6               7.78,
7               7.02, 6, 6.5, 5.8, 5.12, 7.41, 6.52, 6.21, 12.28, 5.6, 5.38, 6.6,
8               8.74])
9
10 # coefficient of variation (CV)
11 cv = np.std(x) / np.mean(x)
12 cv # 0.24737024423748963
13
14 # bootstrap
15 num_bootstraps = 10000
16 bootstrap_cvs = np.zeros(num_bootstraps)
17 np.random.seed(2023)
18 for i in range(num_bootstraps):
19     resample = np.random.choice(x, size=len(x), replace=True)
20     bootstrap_cvs[i] = np.std(resample) / np.mean(resample)
21
22 # bias and standard error of the CV estimator
23 bias = np.mean(bootstrap_cvs) - cv
24 bias # -0.011917596931186353
25 standard_error = np.std(bootstrap_cvs)
26 standard_error # 0.04395079547052894
```

Visualizing the distribution in R



Visualizing the distribution in Python



References

Rizzo, M. L. (2019), Statistical Computing with R, Second Edition, Chapman Hall/CRC, ISBN 9781466553330

The R Project for Statistical Computing:

<https://www.r-project.org/>

Python:

<https://www.python.org/>

course notes