

# Advanced Statistical Methods Hw8

Do Hyup Shin

2021-11-22

## Problem 10.4

Verify formula (10.38) for the number of distinct bootstrap samples.

### Solution

We'll show that the number of distinct bootstrap samples  $= \binom{2n-1}{n}$ . This problem is a duplicate combination problem. Let  $(x_1, x_2, \dots, x_n)$  be the sample and the size of sample is  $n$ . Let the number of times each observation is chosen is  $a_i \forall i = 1, 2, \dots, n$ . Then,  $\sum_{i=1}^n a_i = n$  with  $\forall 0 \leq a_i \leq n$  and  $\forall a_i$  are nonnegative integer. We should find the number of combination  $a_i$  satisfying above condition. This problem is the same as following problem. Suppose that there exist  $n-1$  bars(= |) and  $n$  dots(= ·). Let's arrange the two types of symbols in a row. Then, we can express the arranged line in this way  $\_\_\_ | \_\_\_ | \_\_\_ \dots \_\_\_ | \_\_\_ | \_\_\_$  and  $\_\_\_$  means where  $\cdot$  can enter. There exists  $n$  separation which is  $\_\_\_$ .

Thus, we can correspond  $\forall a_i$  to the number of  $\cdot$  in  $i$ th  $\_\_\_$ . We know that the number of permutations  $n-1$  bars(= |) and  $n$  dots(= ·) is  $\frac{(2n-1)!}{(n-1)!n!} = \binom{2n-1}{n}$ .

Therefore, the number of distinct bootstrap samples is  $\binom{2n-1}{n}$ .

## Problem 10.5

A normal theory least squares model (7.28)-(7.30) yields  $\hat{\beta}$  (7.32). Describe the parametric bootstrap estimates for the standard errors of the components of  $\hat{\beta}$ .

### Solution

The distribution of  $\hat{\beta}$  is  $\hat{\beta} \sim N(\beta, (X^T X)^{-1} \sigma^2)$ . If we know the  $\sigma^2$ , the standard errors of components of  $\hat{\beta}$  are  $se(\hat{\beta}_i) = \sigma(e_i^T (X^T X)^{-1} e_i)^{1/2} \quad \forall i = 1, 2, \dots, p$  where  $e_i$  is the standard basis vector with  $i$ th element zero.

But, if we don't know the  $\sigma^2$ , then we replace  $s^2 = MSE = \frac{1}{n-p-1} y^T (I - H) y = \frac{1}{n-p-1} \sum_{i=1}^n (y_i - \hat{y}_i)^2$  instead of  $\sigma^2$ . Let the design matrix  $X$  be fixed.

## Problem 10.7

Verify formula (10.70).

### Solution

We'll show that the variance of sample mean of bootstrap sample  $X^* = (x_1^*, x_2^*, \dots, x_n^*)$  is  $\sum_{i=1}^n (x_i - \bar{x})^2 / n^2$ . Let  $X = (x_1, x_2, \dots, x_n)$  be random sample from population  $F$  and define  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ . Let  $\hat{F}$  be the empirical probability distribution that puts probability  $1/n$  on each point  $x_i$ . The bootstrap sample with replace from  $\{x_1, x_2, \dots, x_n\}$  is  $X^* = (x_1^*, x_2^*, \dots, x_n^*)$ , i.e  $x_i^* \stackrel{iid}{\sim} \hat{F}$ . Then,  $P(x_i^* = x_j) = \frac{1}{n} \quad \forall 1 \leq i, j \leq n$ . So the expectation of  $x_i^*$  is  $E_{\hat{F}}(x_i^*) = \sum_{j=1}^n x_j P(x_i^* = x_j) = \sum_{j=1}^n x_j \frac{1}{n} = \frac{1}{n} \sum_{j=1}^n x_j = \bar{x}$ .

Define  $\bar{x}^* = \frac{1}{n} \sum_{j=1}^n x_j^*$  which is sample mean of bootstrap sample. Then, the variance of  $\bar{x}^*$  is

$$\begin{aligned}
\text{var}_{\hat{F}}(\bar{x}^*) &= E_{\hat{F}}((\bar{x}^* - E(\bar{x}^*))^2) = E_{\hat{F}}((\bar{x}^* - \bar{x})^2) \\
&= E_{\hat{F}}\left(\sum_{j=1}^n \frac{1}{n} (x_j^* - \bar{x})\right)^2 = \frac{1}{n^2} E_{\hat{F}}\left(\sum_{j=1}^n (x_j^* - \bar{x})\right)^2 \\
&= \frac{1}{n^2} E_{\hat{F}}\left(\sum_{j=1}^n (x_j^* - \bar{x})^2 + \sum_{i \neq j} (x_i^* - \bar{x})(x_j^* - \bar{x})\right) \\
&= \frac{1}{n^2} E_{\hat{F}}\left(\sum_{j=1}^n (x_j^* - \bar{x})^2\right) \quad (\because E_{\hat{F}}(x_i^* - \bar{x})(x_j^* - \bar{x}) = 0 \quad \forall i \neq j) \\
&= \frac{1}{n^2} n E_{\hat{F}}(x_1^* - \bar{x})^2 \quad (\because \forall (x_j^* - \bar{x}) \text{ are following independently identical distribution}) \\
&= \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2 P(x_1^* = x_j) = \frac{1}{n} \sum_{j=1}^n \frac{1}{n} (x_j - \bar{x})^2 \\
&= \frac{1}{n^2} \sum_{j=1}^n (x_j - \bar{x})^2
\end{aligned}$$

Therefore,  $\text{var}_{\hat{F}}(\bar{x}^*) = \frac{1}{n^2} \sum_{j=1}^n (x_j - \bar{x})^2$ . Suppose that there exist B bootstrap samples. Then, We can calculate  $\bar{x}^{*j}$  for each jth bootstrap sample. So the estimate of  $\text{var}_{\hat{F}}(\bar{x}^*)$  is  $\hat{\text{var}}_{boot}(\bar{x}^*) = \frac{1}{B-1} \sum_{j=1}^B (\bar{x}^{*j} - \bar{x}_{(\cdot)})^2$  where  $\bar{x}_{(\cdot)} = \frac{1}{B} \sum_{j=1}^B \bar{x}^{*j}$ .

In conclusion,  $\hat{\text{var}}_{boot}(\bar{x}^*) = \frac{1}{B-1} \sum_{j=1}^B (\bar{x}^{*j} - \bar{x}_{(\cdot)})^2 \rightarrow \text{var}_{\hat{F}}(\bar{x}^*) = \frac{1}{n^2} \sum_{j=1}^n (x_j - \bar{x})^2$  as  $B \rightarrow \infty$ .

### Problem 10.9

A survey in a small town showed incomes  $x_1, x_2, \dots, x_m$  for men and  $y_1, y_2, \dots, y_n$  for women. As an estimate of the differences,

$$\hat{\theta} = \text{median}\{x_1, x_2, \dots, x_m\} - \text{median}\{y_1, y_2, \dots, y_n\}$$

was computed.

- How would you use nonparametric bootstrapping to assess the accuracy of  $\hat{\theta}$  ?
- Do you think your method makes full use of the bootstrap replications?

### Solution

(a)

Let  $X = (x_1, x_2, \dots, x_m)$  and  $Y = (y_1, y_2, \dots, y_n)$  be the samples of men and women, respectively. Some large number B of bootstrap samples are independently drawn. Let  $X^{*j} = (x_1^{*j}, x_2^{*j}, \dots, x_m^{*j})$  and  $Y^{*j} = (y_1^{*j}, y_2^{*j}, \dots, y_n^{*j})$  be the Bth bootstrap sample of X and Y, respectively. The corresponding bootstrap replications are calculated, say  $\hat{\theta}^{*j} = \text{median}\{x_1^{*j}, x_2^{*j}, \dots, x_m^{*j}\} - \text{median}\{y_1^{*j}, y_2^{*j}, \dots, y_n^{*j}\}$ . Then, the bootstrap estimate of standard error for  $\hat{\theta}$  is  $\hat{\text{se}}_{boot}(\hat{\theta}) = \left(\frac{1}{B-1} \sum_{i=1}^B (\hat{\theta}^{*i} - \hat{\theta}_{(\cdot)})^2\right)^{1/2}$ . So, we can assess the accuracy of  $\hat{\theta}$  by above process.

(b)

### Problem 11.1

We observe  $y \sim \lambda G_{10}$  to be  $y = 20$ . Here  $\lambda$  is an unknown parameter while  $G_{10}$  represents a gamma random variable with 10 degrees of freedom ( $y \sim G(10, \lambda)$  in the notation of Table 5.1). Apply the Neyman constructions as in Figure 11.1 to find the confidence limit endpoints  $\hat{\lambda}(0.025)$  and  $\hat{\lambda}(0.975)$ .

**Solution**

**Problem 11.3**

Suppose  $\hat{G}$  in (11.33) was perfectly normal, say  $\hat{G} \sim N(\hat{\mu}, \hat{\sigma}^2)$ . What does  $\hat{\theta}_{BC}(\alpha)$  reduce to in this case, and why does this make intuitive sense?

**Solution**

**Problem 11.5**

Suppose  $\hat{\theta} \sim \text{Poisson}(\theta)$  is observed to equal 16. Without employing simulation, compute the 95% central BCa interval for  $\theta$ . (You can use the good approximation  $z_0 = a = 1/(6\hat{\theta}^{1/2})$ .)

**Solution**

**Problem 11.6**

Use the R program `bcajack` (available with its help file from [efron.web.stanford.edu](http://efron.web.stanford.edu) under “Talks”) to find BCa confidence limits for the student score eigenratio statistic as in Figure 10.2.

**Solution**