

Block Bootstrap Methods

Agnieszka Tracz

October 2023

1 Moving Block Bootstrap

1.1 Introduction

The Moving Block Bootstrap (MBB) is suitable for dependent data without making parametric model assumptions. MBB resamples consecutive blocks of observations not only one observation at a time as before and in this way it preserves the original observations dependence structure within each block.

Moreover, as the sample size increases, the block size also increases. Consequently, when dealing with data generated by a weakly dependent process, MBB asymptotically replicates the underlying dependence structure of the process.

1.2 Description of the MBB

We denote:

X_1, X_2, \dots - sequence of stationary random variables

$X_n = \{X_1, \dots, X_n\}$ - observations

$\hat{\theta}_n = T(F_n)$ - MBB version of estimators

F_n - the empirical distribution function of X_1, \dots, X_n

$T(\cdot)$ - (real-valued) functional of F_n .

$l \equiv l_n \in [1, n]$ - an integer

For dependent data, we typically require that:

$l \rightarrow \infty$ and $n^{-1}l \rightarrow 0$ as $n \rightarrow \infty$

However, a description of the MBB can be given without this restriction.

The process of obtaining MBB samples:

1. Take the observations $\{X_1, \dots, X_n\}$ and choose the integer l - length of blocks.
2. Place the observations in blocks $B_i = (X_i, \dots, X_{i+l-1})$, $1 \leq i \leq N$ where $N = n - l + 1$ (each block has length l)

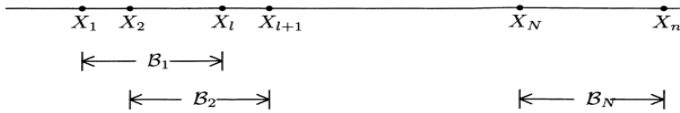


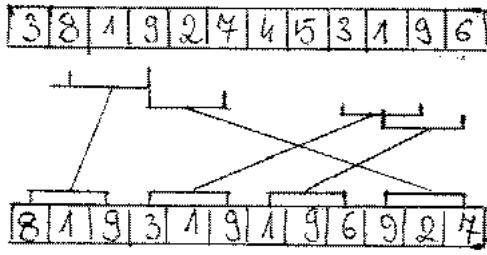
FIGURE 2.1. The collection $\{B_1, \dots, B_N\}$ of overlapping blocks under the MBB.

3. Fix number k and take B_1^*, \dots, B_k^* - simple random sample drawn with replacement from $\{B_1, \dots, B_N\}$
4. We get $(X_{(i-1)l+1}^*, \dots, X_{il}^*)$ - elements in $B_i^*, i = 1, \dots, k$
5. Then X_1^*, \dots, X_m^* constitute the MBB sample of size $m = kl$

The MBB version $\theta_{m,n}^*$ of $\hat{\theta}_n$ is defined as $\theta_{m,n}^* = T(F_{m,n}^*)$, where $F_{m,n}^*$ denotes the empirical distribution of X_1^*, \dots, X_m^*

MBB sample size is typically chosen to be of the same order as the original sample size. If b_1 denotes the smallest integer such that $b_1 l \geq n$, then one may select $k = b_1$ blocks to generate the MBB samples,

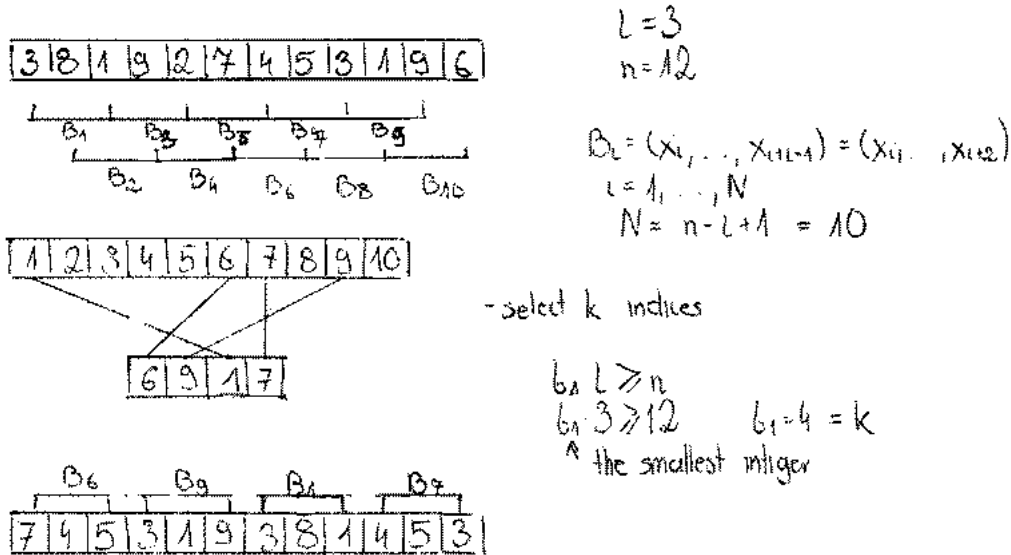
and use only the first n values to define the bootstrap version of T_n



1.3 An alternative formulation of the MBB

Selecting the blocks B_i^* 's randomly from $\{B_1, \dots, B_N\}$ is equivalent to selecting k indices at random from the set $\{1, \dots, N\}$.

Let I_1, \dots, I_k be iid random variables with the *discrete uniform distribution* on $\{1, \dots, N\}$. If we set $B_i^* = B_{I_i}$ for $i = 1, \dots, k$, then B_1^*, \dots, B_k^* represent a simple random sample drawn with replacement from $\{B_1, \dots, B_N\}$. The bootstrap sample X_1^*, \dots, X_m^* can be defined using the resampled blocks B_1^*, \dots, B_k^* as before.



Note: Conditional on the data X_n , the resampled blocks of observations $(X_1^*, \dots, X_l^*)', (X_{l+1}^*, \dots, X_{2l}^*)', \dots, (X_{(k-1)l+1}^*, \dots, X_{kl}^*)'$ are iid l -dimensional random vectors with:

$$P_*((X_1^*, \dots, X_l^*)' = (X_j, \dots, X_{j+l-1})') = P_*(I_1 = j) = N^{-1}, \text{ for } 1 \leq j \leq N$$

P_* denotes the conditional probability given X_n .

1.4 General version of the MBB

1.4.1 Intorduction

Estimators of the form $\hat{\theta}_n = T(F_n)$ considered above include many commonly used estimators, e.g., the sample mean, but they are not sufficiently rich for applications in the time series context.

($\hat{\theta}_n$ depends only on the one-dimensional marginal empirical distribution F_n , and hence does not cover standard statistics like the sample lag correlations)

Given the observations X_n , let $F_{p,n}$ denote the p -dimensional empirical measure:

$$F_{p,n} = (n-p+1)^{-1} \sum_{j=1}^{n-p+1} \gamma_{Y_j},$$

where $Y_j = (X_j, \dots, X_{j+p-1})$ and where for any $y \in \mathbf{R}^p$, γ_y denotes the probability measure on \mathbf{R}^p putting unit mass on y .

The general version of the MBB concerns estimators of the form :

$$\hat{\theta}_n = T(F_{p,n}),$$

where $T(\cdot)$ is now a functional defined on a (rich) subset of the set of all probability measures on \mathbf{R}^p . Here, $p \geq 1$ may be a fixed integer, or it may tend to infinity with n suitably.

1.4.2 "ordinary" approach of the MBB

To define the "ordinary" MBB version we fix a block size l , $1 < l < n - p + 1$, and define the blocks in terms of Y_i 's as

$$\tilde{B}_j = (Y_j, \dots, Y_{j+l-1}), 1 \leq j \leq n - p - l + 2.$$

For $k \geq 1$, select k blocks randomly from the collection: $\tilde{B}_i : 1 \leq i \leq n - p - l + 2$, to generate the MBB observations: $Y_1^*, \dots, Y_l^*, Y_{l+1}^*, \dots, Y_m^*$, where $m = kl$.

The MBB version of $\hat{\theta}_n = T(F_{p,n})$ is now defined as

$$\theta_{m,n}^* = T(\tilde{F}_{m,n}^*),$$

where $\tilde{F}_{m,n}^* \equiv m^{-1} \sum_{j=1}^m \gamma_{Y_j^*}$ denotes the empirical distribution of Y_1^*, \dots, Y_m^* . Thus, for estimators of the form $\hat{\theta}_n = T(F_{p,n})$, the MBB version is defined by resampling from blocks of Y -values instead of blocks of X -values themselves.

3	8	1	9	2	4	4	5	3	1	9	6
x_1	x_2	x_3	x_4								x_{12}

$$Y_j = (X_{j_1}, \dots, X_{j_{p-1}})$$

$$p=3, l=3$$

Y_1	Y_2	Y_3	Y_4	Y_5	Y_6	Y_7	Y_8	Y_9	Y_{10}
3 x_1	8 x_2	1	9	2	7	4	5	3	1
8 x_2	1 x_3	9	2	7	4	5	3	1	9
1 x_3	9 x_4	2	7	4	5	3	1	9	6

$$\tilde{B}_j = (Y_j, \dots, Y_{j+l-1})$$

$$1 \leq j \leq n - p - l + 2 = 11 - 3 - 3 = 5$$

$$\begin{aligned} \tilde{B}_1 &= (Y_1, Y_2, Y_3) \\ \tilde{B}_2 &= (Y_2, Y_3, Y_4) \\ \tilde{B}_3 &= (Y_3, Y_4, Y_5) \\ \tilde{B}_4 &= (Y_4, Y_5, Y_6) \\ \tilde{B}_5 &= (Y_5, Y_6, Y_7) \\ \tilde{B}_6 &= (Y_6, Y_7, Y_8) \\ \tilde{B}_7 &= (Y_7, Y_8, Y_9) \\ \tilde{B}_8 &= (Y_8, Y_9, Y_{10}) \end{aligned}$$

- select k blocks

$$k = b_1$$

b_1 - the smallest integer :

$$b_1 \cdot l \geq n$$

$$b_1 \cdot 3 \geq 12$$

$$b_1 = 4 = k$$

\tilde{B}_6	\tilde{B}_2	\tilde{B}_{10}	\tilde{B}_7
$Y_6 Y_7 Y_8$	$Y_6 Y_7 Y_8$	$Y_1 Y_2 Y_3$	$Y_1 Y_8 Y_9$
7 4 5	2 7 4	3 8 1	4 5 3
4 5 3	7 4 5	8 1 9	5 3 1
5 3 1	4 5 3	1 9 2	3 1 9

- MBB sample

1.4.3 "naive" approach of the MBB

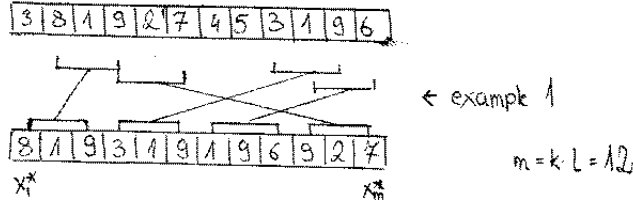
For the more general class of statistics $\hat{\theta}_n$ given by $\hat{\theta}_n = T(F_{p,n})$ for some $p \geq 2$, there is an alternative way of defining the bootstrap version of $\hat{\theta}_n$.

Since the estimator $\hat{\theta}_n$ can always be expressed as a function of the given observations X_1, \dots, X_n , we define the bootstrap version of $\hat{\theta}_n$ by resampling from X_1, \dots, X_n directly.

Suppose that the block bootstrap observations X_1^*, \dots, X_m^* are generated by resampling from the blocks $B_i = (X_i, \dots, X_{i+l-1})$, $i = 1, \dots, N$ of X -values. Define bootstrap "analogs" of the p -dimensional variable $Y_i \equiv (X_i, \dots, X_{i+p-1})'$ in terms of X_1^*, \dots, X_m^* as $Y_i^{**} \equiv (X_i^*, \dots, X_{i+p-1}^*)'$, $i = 1, \dots, m - p + 1$. Then, the bootstrap version of $\hat{\theta}_n$ under this alternative approach is defined as:

$$\theta_{m,n}^{**} = T(\tilde{F}_{m,n}^{**}),$$

where $\tilde{F}_{m,n}^{**} = \sum_{i=1}^{m-p+1} \gamma_{Y_i^{**}}$.



$$\begin{aligned}
 Y_1^{**} &= (8, 1, 9) = (X_1^*, X_2^*, X_3^*) \\
 Y_2^{**} &= (1, 9, 3) = (X_2^*, X_3^*, X_4^*) \\
 Y_3^{**} &= (9, 3, 1) \\
 Y_4^{**} &= (3, 1, 9) \quad \dots \\
 Y_5^{**} &= (1, 9, 1) \\
 Y_6^{**} &= (9, 1, 9) \\
 Y_7^{**} &= (1, 9, 6) \\
 Y_8^{**} &= (9, 6, 9) \\
 Y_9^{**} &= (6, 9, 2) \\
 Y_{10}^{**} &= (9, 2, 7)
 \end{aligned}$$

Y_1^{**}	Y_2^{**}	Y_3^{**}	Y_4^{**}	Y_5^{**}	Y_6^{**}	Y_7^{**}	Y_8^{**}	Y_9^{**}	Y_{10}^{**}
8	1	9	3	1	9	1	9	6	9
1	9	3	1	9	1	9	6	9	2
9	3	1	9	1	9	6	9	2	7

2 Nonoverlapping Block Bootstrap

In this section, we consider the blocking rule due to Carlstein. For simplicity, here we consider estimators given by $\theta_n = T(F_{p,n})$, with $p = 1$ only.

The key feature of Carlstein's blocking rule is to **use nonoverlapping segments of the data** to define the blocks. The corresponding block bootstrap method will be called the **nonoverlapping block bootstrap (NBB)**

Notation:

$l \equiv l_n \in [1, n]$ - an integer (length of blocks)

$b \geq 1$ - the largest integer satisfying: $lb \leq n$

The process of obtaining NBB samples:

1. Take the observations $\{X_1, \dots, X_n\}$ and choose the integer l - length of blocks.
2. Place the observations in blocks $B_i^{(2)} = (X_{(i-1)l+1}, \dots, X_{il})'$, $1 \leq i \leq b$ (each block has length l)

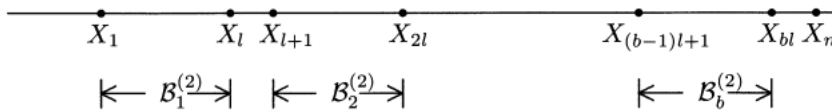


FIGURE 2.2. The collection $\{B_1^{(2)}, \dots, B_b^{(2)}\}$ of nonoverlapping blocks under Carlstein's (1986) rule.

Note: Blocks in the MBB overlap, but the blocks $B_i^{(2)}$'s under the NBB do not. The collection of blocks from which the bootstrap blocks are selected is smaller than the collection for the MBB.

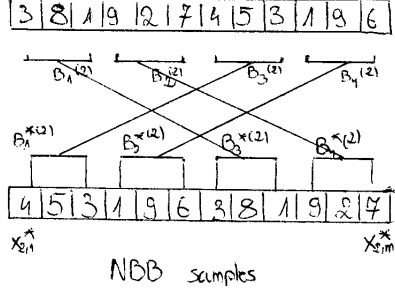
3. Calculate the number k and take $B_1^{*(2)}, \dots, B_k^{*(2)}$ - simple random sample drawn with replacement from $\{B_1^{(2)}, \dots, B_b^{(2)}\}$

4. We get $(X_{(i-1)l+1}^*, \dots, X_{il}^*)$ - elements in $B_i^*, i = 1, \dots, k$

5. Then $X_{2,1}^*, \dots, X_{2,l}^*, \dots, X_{2,(b-1)l+1}^*, \dots, X_{2,m}^*$ constitute the NBB sample of size $m = kl$

$F_{m,n}^{*(2)}$ denote the empirical distribution of the bootstrap sample $X_{2,1}^*, \dots, X_{2,m}^*$, obtained by writing the elements of $B_1^{*(2)}, \dots, B_k^{*(2)}$ in a sequence.

Then, the bootstrap version of estimator $\hat{\theta}_n$ is given by: $\theta_{m,n}^{*(2)} = T(F_{m,n}^{*(2)})$



Even though the definition of the bootstrapped estimators are very similar for the MBB and for the NBB, the resulting bootstrap versions $\theta_{m,n}^*$ and $\theta_{m,n}^{*(2)}$ have different distributional properties.

We illustrate the point with the simplest case, where $\hat{\theta}_n = n^{-1} \sum_{j=1}^n X_j$ is the sample mean. The bootstrap version of $\hat{\theta}_n$ under the two methods are respectively given by:

$$\theta_{m,n}^* = m^{-1} \sum_{j=1}^m X_j^*, \text{ and } \theta_{m,n}^{*(2)} = m^{-1} \sum_{j=1}^m X_{2,j}^*$$

$$\begin{aligned} E_*(\theta_{m,n}^*) &= E_*(\ell^{-1} \sum_{i=1}^{\ell} X_i^*) \\ &= N^{-1} \sum_{j=1}^N \left(\ell^{-1} \sum_{i=1}^{\ell} X_{j+i-1} \right) \\ &= N^{-1} \left\{ n\bar{X}_n - \ell^{-1} \sum_{j=1}^{\ell-1} (\ell-j)(X_j + X_{n-j+1}) \right\}. \end{aligned}$$

To obtain a similar expression for $E_*(\theta_{m,n}^{*(2)})$, note that under the NBB, the bootstrap variables: $(X_{2,1}^*, \dots, X_{2,l}^*), \dots, (X_{2,(b-1)l+1}^*, \dots, X_{2,m}^*)$ are iid, with common distribution:

$$P_*((X_{2,1}^*, \dots, X_{2,l}^*) = (X_{(j-1)l+1}, \dots, X_{jl})) = b^{-1},$$

for $j = 1, \dots, b$

$$\begin{aligned} E_*(\theta_{m,n}^{*(2)}) &= E_*(\ell^{-1} \sum_{i=1}^{\ell} X_{2,i}^*) \\ &= b^{-1} \sum_{j=1}^b \left(\ell^{-1} \sum_{i=1}^{\ell} X_{(j-1)\ell+i} \right) \\ &= (b\ell)^{-1} \left\{ n\bar{X}_n - \sum_{i=b\ell+1}^n X_i \right\}, \end{aligned}$$

which equals \bar{X}_n if n is a multiple of l .

The bootstrapped estimators have different (conditional) means under the two methods.

If the process $\{X_n\}_{n \geq 1}$ satisfies some standard moment and mixing conditions, then $E\{E_*(\theta_{m,n}^*) - E_*(\theta_{m,n}^{*(2)})\}^2 = O(\frac{1}{n^2})$.

The difference between the two is negligible for large sample sizes.

$$\begin{aligned}
\Theta_{min}^* &= m^{-1} \sum_{j=1}^m X_j^* \\
E_*(\Theta_{min}^*) &= E_*(m^{-1} \sum_{j=1}^m X_j^*) = E_*(l^{-1} \sum_{i=1}^l X_i^*) \underset{\uparrow}{=} N^{-1} \sum_{j=1}^N (l^{-1} \sum_{i=1}^l X_{j+i-1}) \\
&= \frac{1}{Nl} \sum_{j=1}^N \sum_{i=1}^l X_{j+i-1} = \frac{1}{Nl} \left(\sum_{i=1}^l X_i + \sum_{i=2}^l X_{i+1} + \dots + \sum_{i=1}^l X_{i+N-1} \right) = \dots \quad N = n-l+1 \\
&= \frac{1}{Nl} \left(\sum_{i=1}^l X_i + \sum_{i=1}^l X_{i+1} + \dots + \sum_{i=1}^l X_{i+N-1} \right) = \dots \\
&= \frac{1}{Nl} \begin{bmatrix} X_1 + \dots + X_l + X_{l+1} + \dots + X_{n-l+1} + \dots + X_n \\ X_l + X_{l+1} + \dots + X_{2l} + \dots + X_{n-l+2} + \dots + X_n \\ \vdots \\ X_{n-l+1} + X_{n-l+2} + \dots + X_n \end{bmatrix} \\
&= \frac{1}{Nl} \left(\sum_{i=1}^l X_i - X_{n-l+2} - \dots - X_n + \sum_{i=l}^n X_i - X_1 - X_{n-l+3} - \dots - X_n \right. \\
&\quad \left. + \sum_{i=1}^n X_i - X_1 - \dots - X_{n-1} \right) \\
&= \frac{1}{Nl} \left(l \sum_{i=1}^l X_i - (l-1)(X_1 + X_n) - (l-2)(X_2 + X_{n-1}) - \dots - 1(X_{l-1} + X_{n-l+2}) \right) \\
&= \frac{1}{Nl} \left(l \cdot n \cdot \frac{1}{n} \sum_{i=1}^n X_i - \sum_{j=1}^{l-1} (l-j)(X_j + X_{n-j+1}) \right) \\
&= \frac{1}{N} \left(n \cdot \bar{X}_n - \frac{1}{l} \sum_{j=1}^{l-1} (l-j)(X_j + X_{n-j+1}) \right)
\end{aligned}$$

$$\Theta_{min}^{*(2)} = m^{-1} \sum_{j=1}^m X_{2j}^*$$

Under the NBB, the bootstrap variables $(X_{2,1}^*, \dots, X_{2,m}^*), \dots, (X_{2,(n-l+1)}^*, \dots, X_{2,m}^*)$ are iid with common distribution:

$$P_*(X_{2,j}^*, \dots, X_{2,l}^*) = (X_{j-nl+1}, \dots, X_{jl}) = 1/b, \quad j = 1, \dots, b$$

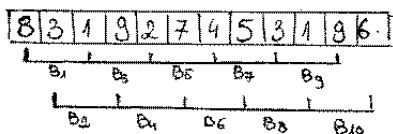
$$\begin{aligned}
E_*(\Theta_{min}^{*(2)}) &= E_*(m^{-1} \sum_{j=1}^m X_{2j}^*) = E_*(l^{-1} \sum_{i=1}^l X_{2i}^*) \underset{\downarrow}{=} b^{-1} \sum_{j=1}^b (l^{-1} \sum_{i=1}^l X_{j-1}l+i) = \\
&= \frac{1}{b \cdot l} \sum_{j=1}^b \sum_{i=1}^l X_{(j-1)l+i} = \\
&= \frac{1}{b \cdot l} \left(\sum_{i=1}^l X_i + \sum_{i=1}^l X_{i+l} + \dots + \sum_{i=1}^l X_{i+(b-1)l} \right) \\
&= \frac{1}{b \cdot l} (X_1 + X_{l+1} + \dots + X_{2l} + \dots + X_{l+(b-1)l} + \dots + X_{bl}) \\
&= \frac{1}{b \cdot l} \left(\sum_{i=1}^{bl} X_i \right) = \frac{1}{b \cdot l} \left(\sum_{i=1}^n X_i - \sum_{i=bl+1}^n X_i \right) = \\
&= \frac{1}{b \cdot l} (n \bar{X}_n - \sum_{i=bl+1}^n X_i) \\
&\quad \uparrow \text{ equals } \bar{X}_n \text{ if } n \text{ is a multiple of } l
\end{aligned}$$

3 Problems with MBB and NBB

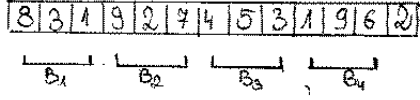
The MBB resampling scheme suffers from an undesirable boundary effect as it assigns lesser weights to the observations toward the beginning and the end of the data set than to the middle part.

Indeed, for $l \leq j \leq n-l$, the j th observation X_j appears in exactly l of the blocks $\{B_1, \dots, B_N\}$, whereas for $1 \leq j \leq l-1$, X_j and X_{n-j+1} appear only in j blocks.

Since there is no observation beyond X_n (or prior to X_1), we cannot define new blocks to get rid of this boundary effect.



A similar problem also exists under the NBB with the observations near the end of the data sequence when n is not a multiple of l .



Politis and Romano suggested a simple way out of this boundary problem by **wrapping the data around a circle and forming additional blocks using the 'circularly defined' observations**. They put forward two resampling schemes based on circular blocks, called the 'circular block bootstrap' (CBB) and the 'stationary bootstrap' (SB).

3.1 Circular Block Bootstrap

The Circular Block Bootstrap (CBB) method resamples overlapping and periodically extended blocks of a given length l satisfying $1 \leq l \leq n$.

The process of obtaining CBB samples:

1. Given the variables $X_n = \{X_1, \dots, X_n\}$, define a new time series $Y_{n,i}$, $i \geq 1$ by **periodic extension**. (for any $i \geq 1$, there are integers $k_i \geq 0$ and $j_i \in [1, n]$ such that $i = k_i n + j_i$, then $i = j_i$ (modulo n)).

Define the variables $Y_{n,i}$, $i \geq 1$ by the relation:

$$Y_{n,i} = X_{j_i}$$

(This is equivalent to writing the variables X_1, \dots, X_n repeatedly on a line and labeling them serially as $Y_{n,i}$, $i \geq 1$)

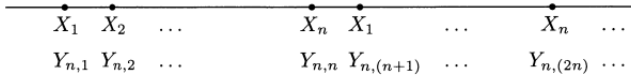


FIGURE 2.3. The periodically extended time series $Y_{n,i}$, $i \geq 1$.

3. Define the blocks:

$$B(i, j) = (Y_{n,i}, \dots, Y_{n,(i+j-1)}) \text{ for } i \geq 1, j \geq 1,$$

In this method we will use the parametr $j : j = l$

4. Then, the CBB resamples blocks from the collection $\{B(i, l) : i \geq 1, \}$

5. Let $(I_1, J_1), (I_2, J_2), \dots$ be a sequence of random vectors with conditional joint distribution. The blocks selected by the CBB are given by $B(I_1, J_1), B(I_2, J_2), \dots$.

6. $X_{C,1}^*, X_{C,2}^*, \dots$ denote the elements of these resampled blocks.

The bootstrap version of an estimator $\hat{\theta}_n = T(F_n)$ under the CBB is defined as $\theta_{m,n}^{*(C)} = T(F_{m,n}^{*(C)})$ for a suitable choice of $m \geq 1$, where $F_{m,n}^{*(C)}$ denotes the empirical distribution of $X_{C,1}^*, \dots, X_{C,m}^*$.

We denote the resampling block indices for the CBB by $I_{3,1}, I_{3,2}, \dots$ (the variables I_i 's in the collection $(I_1, J_1), (I_2, J_2), \dots$

Variables $I_{3,1}, I_{3,2}, \dots$ are **conditionally iid** with:

$$P_*(I_{3,1} = i) = n^{-1} \text{ and } P_*(J_i = l) = 1 \text{ for all } i = 1, \dots, n.$$

Since each X_i appears exactly l times in the collection of blocks $\{B(1, l), \dots, B(n, l)\}$, and since the CBB resamples the blocks from this collection with equal probability, each of the original observations X_1, \dots, X_n receives equal weight under the CBB.

GBB:

8	3	1	9	2	7
---	---	---	---	---	---

 x_1 x_6 $n=6, l=2$ $k=b_1=3$
 $b_1 \cdot l \geq n, b_1 \cdot 2 \geq 6$

$$Y_{n,i} = X_{j_i} \pmod{n}$$

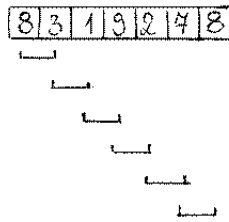
$$\begin{array}{ccccccccc} 8 & 3 & 1 & 9 & 2 & 7 & 8 & 3 & 1 & 9 & 2 & 7 & 8 & 3 & \dots \\ Y_{6,1} & Y_{6,2} & & Y_{6,5} & Y_{6,6} & & Y_{6,12} & Y_{6,13} & & \dots \end{array}$$

$$B(i,j) = (Y_{n,i}, \dots, Y_{n,(i+j-1)}) \quad i \geq 1, j \geq 1$$

CBB We take blocks of given length:

$$B(i,l) = (Y_{n,i}, \dots, Y_{n,(i+l-1)}) = (Y_{n,i}, \dots, Y_{n,i+l-1})$$

$$\begin{aligned} B(1,2) &= (Y_{6,1}, Y_{6,2}) = (8,3) \quad (*) \\ B(2,2) &= (Y_{6,2}, Y_{6,3}) = (3,1) \quad (*) \\ B(3,2) &= (Y_{6,3}, Y_{6,4}) = (1,9) \\ B(4,2) &= (Y_{6,4}, Y_{6,5}) = (9,2) \quad (*) \\ B(5,2) &= (Y_{6,5}, Y_{6,6}) = (2,7) \\ B(6,2) &= (Y_{6,6}, Y_{6,7}) = (7,8) \end{aligned}$$



- select k blocks

$$\begin{array}{|c|c|c|c|c|c|} \hline 3 & 1 & 9 & 2 & 8 & 3 \\ \hline \end{array} \quad \begin{array}{c} x_{3,1}^* \\ x_{3,6}^* \end{array} \quad \text{CBB samples}$$

Let $X_{3,1}^*, X_{3,2}^*, \dots$ denote the CBB observations obtained by arranging the elements of the resampled blocks $\{B(I_{3,i}, l) : i \geq 1\}$ and let $\bar{X}_m^{*(3)}$ denote the CBB sample mean based on m bootstrap observations, where $m = kl$ for some integer $k \geq 1$. Then:

$$\begin{aligned} E_* \bar{X}_m^{*(3)} &= E_* \left[m^{-1} \sum_{i=1}^m X_{3,i}^* \right] = E_* \left[l^{-1} \sum_{i=1}^m X_{3,i}^* \right] \\ &= l^{-1} E_* \left[\sum_{i=1}^m X_{3,i}^* \right] = l^{-1} E_* \left[\sum_{j=1}^n \sum_{i=1}^m Y_{n,(j+i-1)} \right] = (*) \\ &\quad \uparrow \begin{array}{l} P(I_{3,1}=i) = \frac{1}{n} \\ \text{from } \{B(I_{3,i}, l) : l \geq 1\} \end{array} \quad \begin{array}{l} P(Y_i=l) = \frac{1}{n} \\ \text{from } \{B(i,j) : j \geq 1\} \end{array} \\ &\quad \quad \quad B(i,j) = (Y_{n,i}, \dots, Y_{n,(i+j-1)}) \\ (*) &= l^{-1} n^{-1} \left[\sum_{i=1}^l Y_{n,i} + \sum_{i=1}^l Y_{n,i+l} + \dots + \sum_{i=1}^l Y_{n,i+(n-l)} \right] = \\ &= \frac{Y_{n,1} + Y_{2n} + \dots + Y_{n,l}}{Y_{2n} + Y_{3n} + \dots + Y_{n,l}} \\ &\quad \quad \quad \begin{array}{cc} Y_{n,n} & Y_{n,(n-l)} \\ \sum_{i=1}^l X_i & \sum_{i=1}^l X_i \end{array} \\ &= l^{-1} n^{-1} \left(l \cdot \sum_{i=1}^n X_i \right) = n^{-1} \sum_{i=1}^n X_i = \bar{X}_n \end{aligned}$$

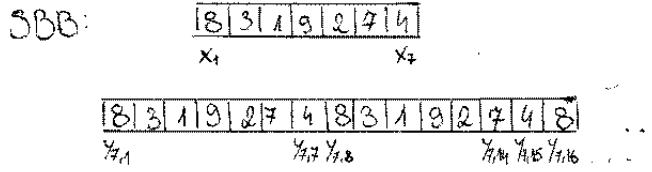
3.2 Stationary Block Bootstrap

The stationary bootstrap (SB) differ from the earlier block bootstrap methods in that it uses blocks of random lengths rather than blocks of a fixed length l .

Let $p \equiv p_n \in (0,1)$ be such that $p \rightarrow 0$ and $np \rightarrow \infty$ as $n \rightarrow \infty$. Then the SB resamples the blocks $B(I_{4,1}, J_{4,1}), B(I_{4,2}, J_{4,2}), \dots$ where the index vectors $(I_{4,1}, J_{4,1}), (I_{4,2}, J_{4,2}), \dots$ are conditionally iid with $I_{4,1}$

having the **discrete uniform distribution** on $\{1, \dots, n\}$, and $J_{4,1}$ having the **geometric distribution** ν_n with parameter p :

$$P_*(J_{4,1} = j) \equiv \nu_n(j) = p(1-p)^{j-1}, \quad j = 1, 2, \dots$$



$$B(l, l) = (y_{n,l}, \dots, y_{n,(l+l-l)})$$

l - not fixed

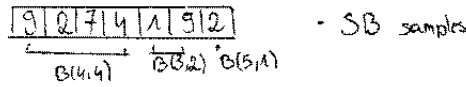
$$B(1, 1) = 8$$

$$B(3, 2) = (1, 9)$$

$$B(6, 3) = (7, 4, 4)$$

etc. ...

$$B(4, 4) = (9, 2, 7, 4)$$



An important property of the SB method is that conditional on \mathbf{X}_n , the bootstrap observations $\{X_{4,i}^*\}_{i \in \mathbf{N}}$ are stationary (which is why it is called the "stationary" bootstrap).

For a given resample size m , the conditional expectation of the SB sample mean $\bar{X}_m^{*(4)} \equiv m^{-1} \sum_{i=1}^m X_{4,i}^*$ is given by:

$$E_*(\bar{X}_m^{*(4)}) = E_* X_{4,i}^* = \bar{X}_n$$

4 Subsampling

Subsampling is a widely used technique in statistical analysis, particularly in scenarios with independently and identically distributed (iid) observations. It involves the selection of different subsets of data to approximate the bias and variance of a statistical estimator.

However, subsampling is not limited to iid observations. It can also be applied to subseries of dependent observations to yield valid estimations of not only bias and variance but also the entire sampling distribution of a statistic. Surprisingly, these valid estimations can be achieved under very mild assumptions.

Suppose that $\hat{\theta}_n = t_n(\mathbf{X}_n)$ is an estimator of a parameter θ , such that for some normalizing constant $a_n > 0$, the probability distribution

$$Q_n(x) = P(a_n(\hat{\theta}_n - \theta) \leq x)$$

of the centered and scaled estimator $\hat{\theta}_n$ converges weakly to a limit distribution $Q(x)$:

$$Q_n(x) \rightarrow Q(x)$$

for all continuity points x of Q . Furthermore, assume that $a_n \rightarrow \infty$ as $n \rightarrow \infty$ and that Q is not degenerate at zero ($Q(\{0\}) < 1$).

Let $1 \leq l \leq n$ be a given integer and let:

$$B_i = (X_i, \dots, X_{i+l-1})', \quad 1 \leq i \leq N$$

denote the overlapping blocks of length l where $N = n - l + 1$. (the same as in MBB)

Then, the subsampling estimator of Q_n , based on the overlapping version of the subsampling method, is given by :

$$\hat{Q}_n(x) = N^{-1} \sum_{i=1}^N 1(a_l(\hat{\theta}_{i,l} - \hat{\theta}_n) \leq x), \quad x \in \mathbf{R}$$

where $\hat{\theta}_{i,l}$ is a "copy" of the estimator $\hat{\theta}_n$ on the block B_i , defined by $\hat{\theta}_{i,l} = t_l(B_i)$, $i = 1, \dots, N$.

Note: We used constant a_l instead of a_n and $t_l(\cdot)$ (in place of $t_n(\cdot)$) to define the subsample copy $\hat{\theta}_{i,l}$ as the i th block B_i contains only l observations.

Note: The overlapping version of the subsampling method is a special case of the MBB where a single block is resampled.

The estimator \hat{Q}_n of the distribution function $Q_n(x)$ can be used to obtain subsampling estimators of the bias and the variance of $\hat{\theta}_n$.

The bias of $\hat{\theta}_n$ is given by:

$$\text{Bias}(\hat{\theta}_n) = E\hat{\theta}_n - \theta = a_n^{-1} \int x dQ_n(x)$$

The subsampling estimator of $\text{Bias}(\hat{\theta}_n)$ is then obtained by replacing $Q_n(\cdot)$ by $\hat{Q}_n(\cdot)$

$$\widehat{\text{Bias}}(\hat{\theta}_n) = a_n^{-1} \int x d\hat{Q}_n(x) = a_l a_n^{-1} (N^{-1} \sum_{i=1}^N \hat{\theta}_{i,l} - \hat{\theta}_n)$$

Similarly, the subsampling estimator of the variance of $\hat{\theta}_n$ is given by :

$$\widehat{\text{Var}}(\hat{\theta}_n) = a_l^2 a_n^{-2} [N^{-1} \sum_{i=1}^N \hat{\theta}_{i,l}^2 - (N^{-1} \sum_{i=1}^N \hat{\theta}_{i,l})^2],$$

which is the sample variance of $\hat{\theta}_{i,l}$ multiplied by the scaling factor $a_l^2 a_n^{-2}$. We need to use the correction factors $\frac{a_l}{a_n}$ and $(\frac{a_l}{a_n})^2$ to scale up from the level of $\hat{\theta}_{i,l}$'s, which are defined using l observations, to the level of $\hat{\theta}_n$, which is defined using n -observations.

In applying a bootstrap method, we typically use a resample size that is comparable to the original sample size, and therefore, such explicit corrections of the bootstrap bias and variance estimators are usually unnecessary.

In analogy to the bootstrap methods, one may attempt to apply the subsampling method to a centered variable of the form $T_{1n} \equiv (\hat{\theta}_n - \theta)$. However, this may not be the right thing to do.

(When we don't take the scaling factor we will overestimate the variance)

As noted previously, the subsampling method is a special case of the MBB where the number of resampled blocks is identically equal to 1. Exploiting this fact, we may similarly define other versions of the subsampling method based on nonoverlapping blocks or circular blocks.