

3 Bootstrapping Dependent Data

Suppose we have dependent data $\mathbf{y} = (y_1, \dots, y_n)$ generated from some unknown distribution $F = F_{\mathbf{Y}} = F_{(Y_1, \dots, Y_n)}$.

y_1, \dots, y_n no longer independent, could be time series for example. (or spatial)

Goal:

To approximate the dsn of a statistic
 $\theta = T(\mathbf{y})$.

Challenge:

Since y_i 's dependent it is inappropriate to use the bootstrap for iid data! Bootstrapped samples would no longer reproduce the data generating process.

We will consider 2 approaches

- ① model-based
- ② block bootstrap (2 types)

3.1 Model-based approach

Example 3.1 Suppose we observe a time series $\mathbf{Y} = (Y_1, \dots, Y_n)$ which we assume is generated by an AR(1) process, i.e., "Auto-regressive"

$$Y_t = \alpha Y_{t-1} + \varepsilon_t \quad t=1, \dots, n$$

iid

$|\alpha| < 1$ and $\varepsilon_1, \dots, \varepsilon_n \sim (0, \sigma^2)$

↑ mean 0 ↑ variance = σ^2

turn our problem into an iid bootstrap!

If we assume an AR(1) model for the data, we can consider a method similar to bootstrapping residuals for linear regression.

- ① Estimate $\hat{\alpha}$ from data (fit the model)
 - ② Define estimated Innovations $\hat{\varepsilon}_t = Y_t - \hat{\alpha} Y_{t-1}, t=2, \dots, n$
 - and $\bar{\hat{\varepsilon}} = \frac{1}{n-1} \sum_{t=2}^n \hat{\varepsilon}_t$ their mean.
 - ③ Define the residuals of the model as the centered innovations
- $$\hat{\varepsilon}_t = \hat{\varepsilon}_t - \bar{\hat{\varepsilon}} \quad (E\varepsilon_i = 0)$$
- ④ For $r=1, \dots, R$,
 - a) create the bootstrap sample $\hat{\varepsilon}_0^*, \dots, \hat{\varepsilon}_n^*$ by independently sampling $n-1$ values from the $n-1$ values $\hat{\varepsilon}_t, t=2, \dots, n$.
 - b) construct pseudo-data $\mathbf{Y}^* = (Y_1^*, \dots, Y_n^*)$ from $Y_0^* = \hat{\varepsilon}_0^*$,
 $Y_t^* = \hat{\alpha} Y_{t-1}^* + \hat{\varepsilon}_t^*, t=1, \dots, n$.
 - c) define $\hat{\alpha}_r^*$ as the estimate of α from Y_1^*, \dots, Y_n^* . - ⑤ The dist of $\hat{\alpha}_1^*, \dots, \hat{\alpha}_R^*$ is used to estimate the sampling dist of $\hat{\alpha}$.

Model-based – the performance of this approach depends on the model being appropriate for the data.

This may not always be a good assumption.

3.2 Nonparametric approach

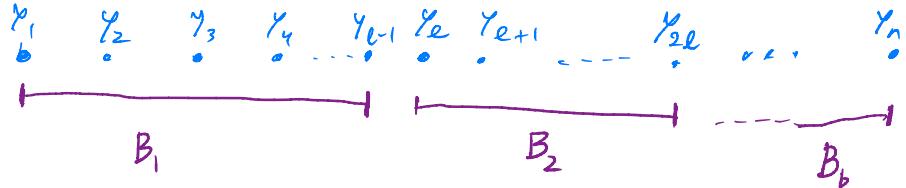
To deal with dependence in the data, we will employ a nonparametric *block* bootstrap.

Idea:

resample data in blocks to preserve the dependence structure
within the blocks.

3.2.1 Nonoverlapping Blocks (NBB)

Consider splitting $\mathbf{Y} = (Y_1, \dots, Y_n)$ in b consecutive blocks of length ℓ .



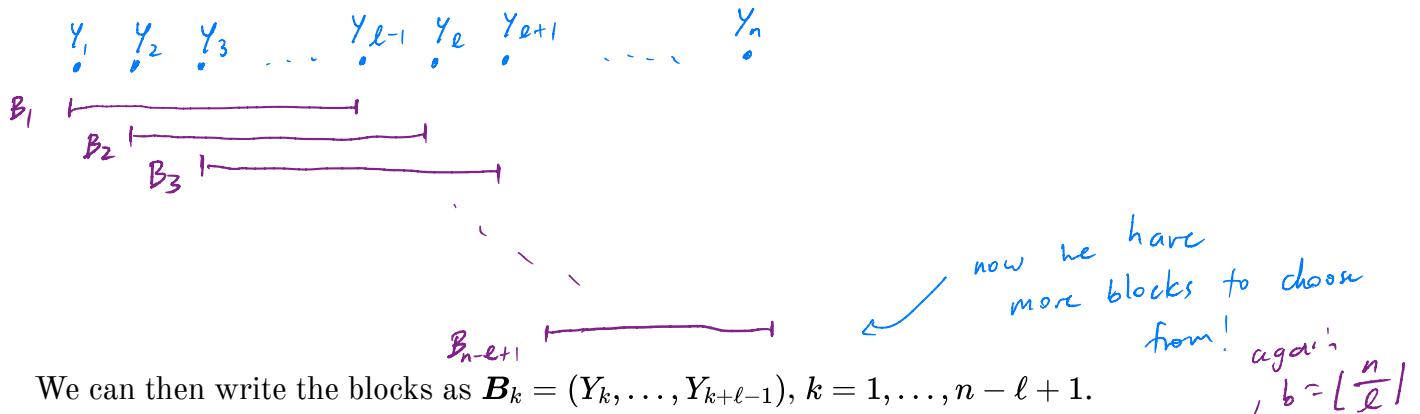
We can then rewrite the data as $\mathbf{Y} = (\mathbf{B}_1, \dots, \mathbf{B}_b)$ with $\mathbf{B}_k = (Y_{(k-1)\ell+1}, \dots, Y_{k\ell})$,
 $k = 1, \dots, b$. $b = \left\lfloor \frac{n}{\ell} \right\rfloor$ "floor function" = round down.

- ① Sample nonoverlapping blocks B_1^*, \dots, B_b^* independently from B_1, \dots, B_b (with replacement) to form pseudo dataset $\mathbf{Y}^* = (B_1^*, \dots, B_b^*)$.
- ② Statistic of interest θ is estimated from \mathbf{Y}^* to create $\hat{\theta}^*$
- ③ Repeat R times to obtain $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$ with which to estimate the dsn of $\hat{\theta}$.

Note, the order of data within the blocks must be maintained, but the order of the blocks that are resampled does not matter.

3.2.2 Moving Blocks (MBB)

Now consider splitting $\mathbf{Y} = (Y_1, \dots, Y_n)$ into overlapping blocks of adjacent data points of length ℓ .



We can then write the blocks as $\mathbf{B}_k = (Y_k, \dots, Y_{k+\ell-1})$, $k = 1, \dots, n - \ell + 1$.

- ① Create a pseudo dataset by sampling $\mathbf{y}^* = (B_1^*, \dots, B_b^*)$ from the blocks independently w/ replacement from $B_1, \dots, B_{n-\ell+1}$
- ② Calculate $\hat{\theta}^*$ from \mathbf{y}^*
- ③ Repeat 1-2 R times to obtain $\hat{\theta}_1^*, \dots, \hat{\theta}_R^*$.

3.2.3 Choosing Block Size

If the block length is too short,

the resampling cannot capture the dependence ($\ell=1$ is the iid bootstrap!)

If the block length is too long,

not many blocks to sample. (does not resample data generation).

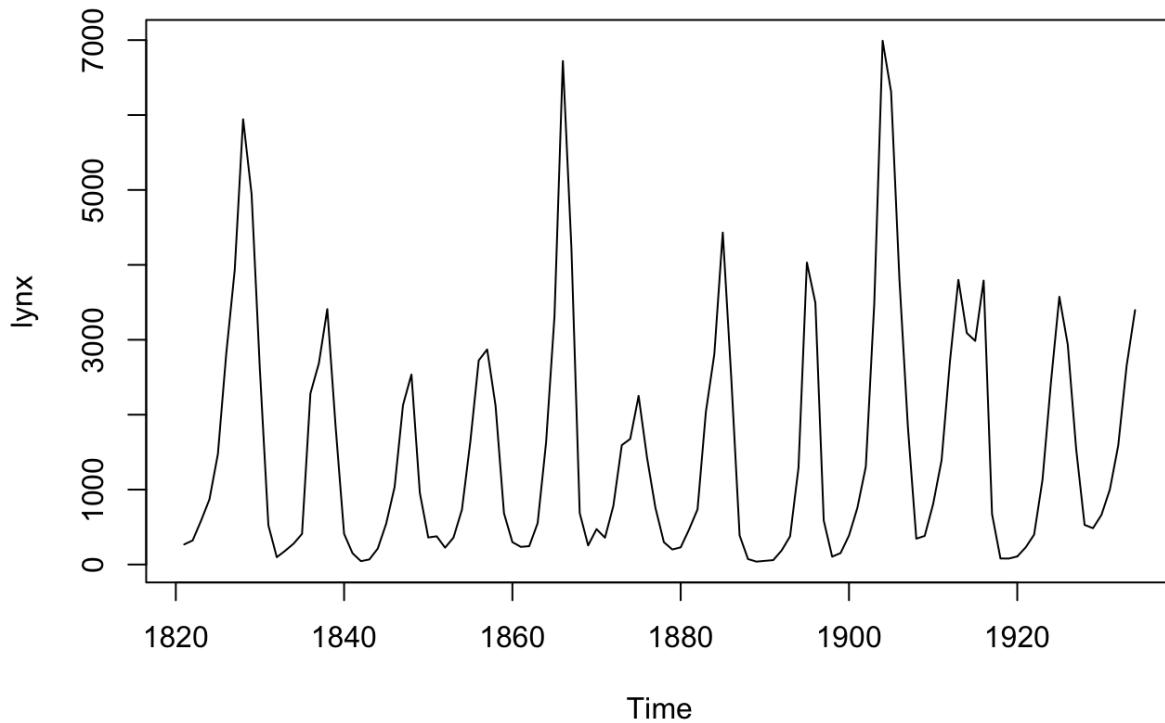
A symptotic result - block length should increase w/ length of the time series. If so, MBB & NBB produce consistent estimators of moments, correct coverage probabilities for CIs and correct error rates for tests.

There are practical methods for choosing ℓ
(Lahiri 2003)

Your Turn

We will look at the annual numbers of lynx trappings for 1821–1934 in Canada. Taken from Brockwell & Davis (1991).

```
data(lynx)
plot(lynx)
```



Goal: Estimate the sample distribution of the mean

$$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n y_i$$

```
theta_hat <- mean(lynx)
theta_hat
```

```
## [1] 1538.018
```

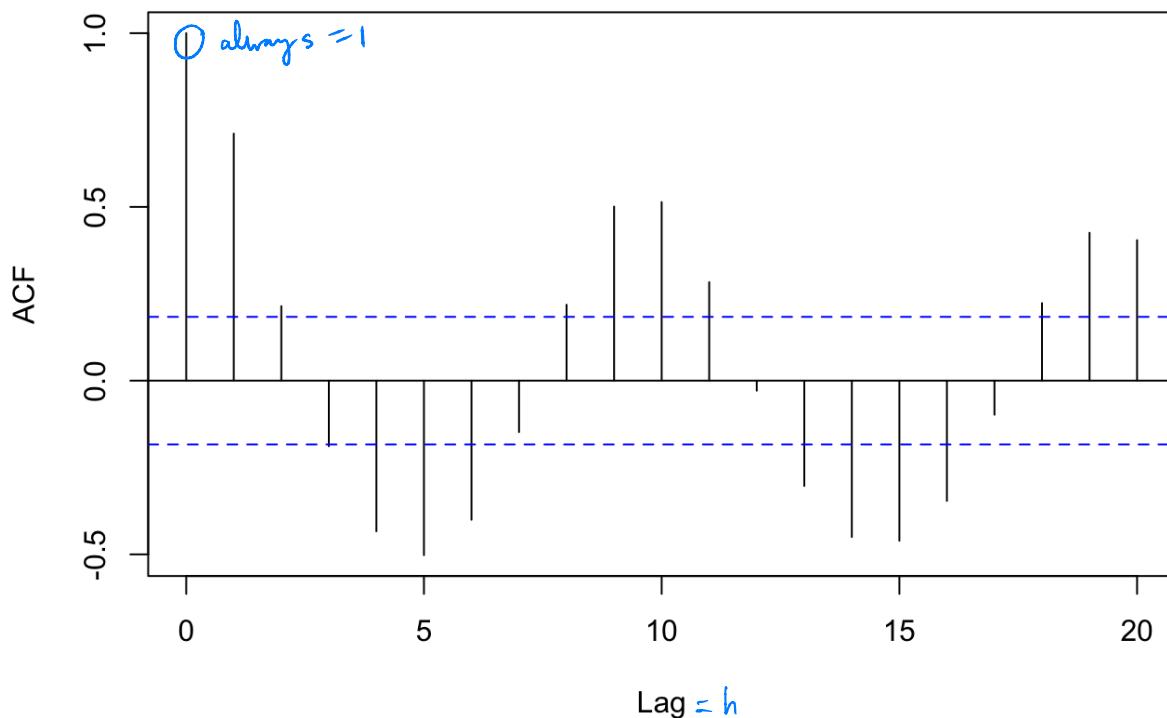
3.2.4 Independent Bootstrap

```
library(simpleboot)
B <- 10000

## Your turn: perform the independent bootstrap
## what is the bootstrap estimate se?
```

We must account for the dependence to obtain a correct estimate of the variance!

acf(lynx)
 Auto correlation function, $\gamma(h) = \text{Cor}(Y_t, Y_{t+h})$
 \uparrow \uparrow
 fraction of lag. Series lynx
 lagged observations



The acf (autocorrelation) in the dominant terms is positive, so we are *underestimating* the standard error.

3.2.5 Non-overlapping Block Bootstrap

```

# function to create non-overlapping blocks
nb <- function(x, b) {
  n <- length(x)
  l <- n %/% b

  blocks <- matrix(NA, nrow = b, ncol = 1)
  for(i in 1:b) {
    blocks[i, ] <- x[((i - 1)*l + 1):(i*l)]
  }
  blocks
}

# Your turn: perform the NBB with b = 10 and l = 11
theta_hat_star_nbb <- rep(NA, B)
nb_blocks <- nb(lynx, 10)
for(i in 1:B) {
  # sample blocks
  # get theta_hat^*
}

# Plot your results to inspect the distribution
# What is the estimated standard error of theta hat? The Bias?

```

3.2.6 Moving Block Bootstrap

```

# function to create overlapping blocks
mb <- function(x, l) {
  n <- length(x)
  blocks <- matrix(NA, nrow = n - l + 1, ncol = 1)
  for(i in 1:(n - l + 1)) {
    blocks[i, ] <- x[i:(i + l - 1)]
  }
  blocks
}

# Your turn: perform the MBB with l = 11
mb_blocks <- mb(lynx, 11)
theta_hat_star_mbb <- rep(NA, B)
for(i in 1:B) {
  # sample blocks
  # get theta_hat^*
}

```

```
}
```

```
# Plot your results to inspect the distribution  
# What is the estimated standard error of theta hat? The Bias?
```

3.2.7 Choosing the Block size

```
# Your turn: Perform the mbb for multiple block sizes l = 1:12  
# Create a plot of the se vs the block size. What do you notice?
```

4 Summary

Bootstrap methods are simulation methods for frequentist inference.

Bootstrap methods are useful for

- many problem types
- esp. when standard assumptions are invalid (like non normal statistics)

Remember Bootstrap principle: the bootstrap dsn should approximate the sampling dsn of the statistic!

Bootstrap methods can fail when

we have extremes or heavy tailed distributions

can be computer intensive (i.e. slow)

need to be careful with dependent data.

as opposed to Bayesian inference, which you should take next semester.