

ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ

**Σχολή Χρηματοοικονομικής και
Στατιστικής**



**Τμήμα Στατιστικής και Ασφαλιστικής Επιστήμης
ΜΕΤΑΠΤΥΧΙΑΚΟ ΠΡΟΓΡΑΜΜΑ ΣΠΟΥΔΩΝ
ΣΤΗΝ ΕΦΑΡΜΟΣΜΕΝΗ ΣΤΑΤΙΣΤΙΚΗ**

**ΕΦΑΡΜΟΓΕΣ ΤΟΥ BOOTSTRAP ΣΤΗΝ
ΑΝΑΛΥΣΗ ΧΡΟΝΟΣΕΙΡΩΝ**

Κωνσταντίνος Καβρός

Διπλωματική Εργασία

που υποβλήθηκε στο Τμήμα Στατιστικής και Ασφαλιστικής
Επιστήμης του Πανεπιστημίου Πειραιώς ως μέρος των
απαιτήσεων για την απόκτηση του Μεταπτυχιακού Διπλώματος
Ειδίκευσης στην *Εφαρμοσμένη Στατιστική*

Πειραιάς
Νοέμβριος 2023

ΠΑΝΕΠΙΣΤΗΜΙΟ ΠΕΙΡΑΙΩΣ
Σχολή Χρηματοοικονομικής και Στατιστικής



Τμήμα Στατιστικής και Ασφαλιστικής Επιστήμης

ΜΕΤΑΠΤΥΧΙΑΚΟ ΠΡΟΓΡΑΜΜΑ ΣΠΟΥΔΩΝ
ΣΤΗΝ ΕΦΑΡΜΟΣΜΕΝΗ ΣΤΑΤΙΣΤΙΚΗ

ΕΦΑΡΜΟΓΕΣ ΤΟΥ BOOTSTRAP ΣΤΗΝ
ΑΝΑΛΥΣΗ ΧΡΟΝΟΣΕΙΡΩΝ

Κωνσταντίνος Καβρός

Διπλωματική Εργασία

που υποβλήθηκε στο Τμήμα Στατιστικής και Ασφαλιστικής
Επιστήμης του Πανεπιστημίου Πειραιώς ως μέρος των
απαιτήσεων για την απόκτηση του Μεταπτυχιακού Διπλώματος
Ειδίκευσης στην *Εφαρμοσμένη Στατιστική*

Πειραιάς
Νοέμβριος 2023

Η παρούσα Διπλωματική Εργασία εγκρίθηκε ομόφωνα από την Τριμελή Εξεταστική Επιτροπή που ορίσθηκε από τη ΓΣΕΣ του Τμήματος Στατιστικής και Ασφαλιστικής Επιστήμης του Πανεπιστημίου Πειραιώς στην υπ' αριθμ. συνεδρίασή του σύμφωνα με τον Εσωτερικό Κανονισμό Λειτουργίας του Προγράμματος Μεταπτυχιακών Σπουδών στην Εφαρμοσμένη Στατιστική.

Τα μέλη της Επιτροπής ήταν:

- Καθηγητής Γ. Ηλιόπουλος (Επιβλέπων)
- Αναπληρωτής Καθηγητής Μ. Μπούτσικας
- Αναπληρωτής Καθηγητής Κ. Πολίτης

Η έγκριση της Διπλωματικής Εργασίας από το Τμήμα Στατιστικής και Ασφαλιστικής Επιστήμης του Πανεπιστημίου Πειραιώς δεν υποδηλώνει αποδοχή των γνώμων του συγγραφέα.

UNIVERSITY OF PIRAEUS
School of Finance and Statistics



Department of Statistics and Insurance Science

POSTGRADUATE PROGRAM IN
APPLIED STATISTICS

BOOTSTRAP APPLICATIONS TO TIME
SERIES ANALYSIS

By
Konstantinos Kavros

MSc Dissertation

submitted to the Department of Statistics and Insurance Science
of the University of Piraeus in partial fulfilment of the
requirements for the degree of Master of Science in Applied
Statistics

Piraeus, Greece
November 2023

Περίληψη

Το bootstrap Efron (1979) είναι μια υπολογιστική στατιστική τεχνική η οποία έχει αποδειχθεί ένα ισχυρό εργαλείο για την εκτίμηση της διακύμανσης αλλά και της δειγματικής κατανομής μιας στατιστικής συνάρτησης είτε παραμετρικά είτε μη παραμετρικά. Αν και το bootstrap αρχικά αναπτύχθηκε για ανεξάρτητα δεδομένα, στη συνέχεια επεκτάθηκε και για πιο σύνθετα προβλήματα όπου τα δεδομένα είναι εξαρτημένα, όπως συμβαίνει στην περίπτωση των χρονοσειρών. Στην παρούσα εργασία θα παρουσιαστούν οι κυριότερες τεχνικές bootstrap για χρονοσειρές και θα εφαρμοστούν σε προσομοιωμένα αλλά και πραγματικά δεδομένα. Για δεδομένα χρονοσειρών υπάρχουν δύο μέθοδοι εφαρμογής bootstrap. Η πρώτη είναι η παραμετρική προσέγγιση όπου υποτίθεται ότι τα δεδομένα προέρχονται από κάποιο παραμετρικό μοντέλο και το bootstrap γίνεται στα κατάλοιπα που προκύπτουν μετά από την εκτίμησή του. Η δεύτερη προσέγγιση είναι η μη παραμετρική, δηλαδή ελεύθερη μοντέλου, όπου η δειγματοληψία με επανάθεση πραγματοποιείται σε blocks παρατηρήσεων της αρχικής χρονοσειράς.

Η παρούσα διπλωματική εργασία αποτελείται από πέντε κεφάλαια. Στα πρώτα τρία κεφάλαια παρουσιάζεται το θεωρητικό υπόβαθρο του αντικειμένου της διπλωματικής εργασίας ενώ στα υπόλοιπα πραγματοποιείται εφαρμογή των μεθόδων. Πιο συγκεκριμένα, στο πρώτο κεφάλαιο παρουσιάζεται η κλασική μέθοδος bootstrap και κάποιες παραλλαγές της που αναφέρονται στο κύριο μέρος της εργασίας. Στο δεύτερο κεφάλαιο γίνεται μια εισαγωγή στις χρονοσειρές και σε μοντέλα που θα χρησιμοποιηθούν στο bootstrap χρονοσειρών. Το τρίτο κεφάλαιο αποτελεί το κύριο μέρος της εργασίας στο οποίο παρουσιάζονται θεωρητικά οι μέθοδοι bootstrap για χρονοσειρές. Στην συνέχεια στο τέταρτο κεφάλαιο γίνεται οπτικοποίηση των αποτελεσμάτων της εφαρμογής των μεθόδων και μια μελέτη προσομοίωσης ώστε να φανεί η συμπεριφορά των μεθόδων bootstrap σε μια αυτοπαλίνδρομη διαδικασία πρώτης τάξεως. Τέλος, στο πέμπτο κεφάλαιο πραγματοποιείται εφαρμογή των μεθόδων σε ένα πραγματικό σύνολο δεδομένων. Για τη μελέτη προσομοίωσης και την ανάλυση των δεδομένων έχει χρησιμοποιηθεί η γλώσσα προγραμματισμού R.

Abstract

The bootstrap, introduced by Efron (1979), is a statistical computational technique that has been proven to be a powerful tool for estimating the variance and the sampling distribution of a statistical function, either parametrically or nonparametrically. Although initially developed for independent data, it has since been extended to more complex problems where the data are dependent, as is the case with time series. In this thesis, the main bootstrap techniques for time series will be presented and applied to both simulated and real data. In time series data, there are two methods of applying the bootstrap. The first is the parametric approach, where it is assumed that the data come from some parametric model, and the bootstrap is performed on the residuals obtained after its estimation. The second approach is nonparametric, i.e. model-free, where resampling is done on blocks of observations from the original time series.

This thesis consists of five chapters. The first three chapters present the theoretical background of the thesis topic, while the remaining chapters focus on the application of the methods. Specifically, the first chapter introduces the classical bootstrap method and some variations relevant to the main part of the work. The second chapter provides an introduction to time series and models that will be used in the time series bootstrap. The third chapter is the main part of the work, presenting the theoretical methods of the bootstrap for time series. Subsequently, in the fourth chapter, the results of applying the methods are visualized and a simulation study is conducted to demonstrate the behavior of bootstrap methods in a first order autoregressive process. Finally, in the fifth chapter, the methods are applied to real-world data. The R programming language is used for the simulation study and data analysis.

Περιεχόμενα

Κεφάλαιο 1	1
1.1 Εισαγωγή στη μέθοδο Bootstrap	1
1.2 Η μέθοδος Monte Carlo.....	2
1.3 Γενική παρουσίαση Bootstrap για εκτίμηση τυπικού σφάλματος.....	5
1.4 Εμπειρική συνάρτηση κατανομής	8
1.5 Εκτιμητής αντικατάστασης	9
1.6 Η μέθοδος Jackknife.....	12
1.6.1 Η μέθοδος delete-d Jackknife.....	14
1.6.2 Σχέση μεταξύ μεθόδου Jackknife και Bootstrap	15
1.6.3 Η μέθοδος Cross-Validation	15
1.7 Βήματα αλγορίθμου Bootstrap για εκτίμηση τυπικού σφάλματος και μεροληψίας.....	16
1.7.1 Jackknife after Bootstrap.....	18
1.8 Παραμετρικό Bootstrap.....	19
1.9 Αποτυχία μεθόδου Bootstrap.....	20
1.10 M-out-of-n Bootstrap	22
1.11 Διάστημα εμπιστοσύνης Percentile Bootstrap.....	22
Κεφάλαιο 2.....	25
2.1 Δεδομένα χρονοσειρών	25
2.2 Στασιμότητα	25
2.3 Λευκός Θόρυβος.....	27
2.4 Τυχαίος περίπατος.....	28
2.5 Γραμμικές στοχαστικές διαδικασίες.....	29
2.6 Αυτοπαλίνδρομη χρονοσειρά (AR).....	30
2.6.1 Αυτοπαλίνδρομη χρονοσειρά τάξης 1	31
2.7 Χρονοσειρά κινητού μέσου (MA).....	32
2.7.1 Χρονοσειρά κινητού μέσου τάξης 1.....	34
2.8 Αυτοπαλίνδρομα μοντέλα κινητού μέσου ARMA.....	34
Κεφάλαιο 3.....	38
3.1 Βάσει μοντέλου Bootstrap στα κατάλοιπα.....	38
3.1.2 Βάσει μοντέλου Bootstrap στα κατάλοιπα ενός AR(p).....	42
3.1.3 Ασταθής αυτοπαλίνδρομη διαδικασία.....	44
3.2 Bootstrap σε εξαρτημένα δεδομένα.....	45
3.3 Block Bootstrap.....	46

3.4 Moving Block Bootstrap (MBB).....	49
3.5 Circular Block Bootstrap (CBB)	52
3.6 Stationary Bootstrap (SB)	53
3.7 Post-blackening Bootstrap.....	54
3.8 Matched Block Bootstrap (MaBB) και Tapered Block Bootstrap (TBB)	55
3.9 Επιλογή μεγέθους block	57
3.10 AR-Sieve Bootstrap.....	58
Κεφάλαιο 4.....	62
4.1 Περιγραφή μελέτης προσομοίωσης.....	62
4.2 Οπτικοποίηση εφαρμογής μεθόδων Bootstrap.....	62
4.3 Προσομοίωση για $n = 100$	69
4.4 Προσομοίωση για $n = 200$	76
Κεφάλαιο 5.....	80
5.1 Εφαρμογή μεθόδων Bootstrap σε πραγματικά δεδομένα.....	80
Παραρτήματα	91
Βιβλιογραφία.....	113

Κεφάλαιο 1

1.1 Εισαγωγή στη μέθοδο Bootstrap

Η βασική ιδέα της στατιστικής ανάλυσης είναι η εξαγωγή όλης της πληροφορίας από τα δεδομένα (Rao, 1989). Αυτό που μας ενδιαφέρει είναι η κατανομή του πληθυσμού F και συγκεκριμένα οι παράμετροι της κατανομής. Οι παράμετροι είναι συγκεκριμένες ποσότητες που προσδιορίζουν τα χαρακτηριστικά ενός πληθυσμού και μπορεί να είναι είτε γνωστές είτε άγνωστες. Τα αντίστοιχα χαρακτηριστικά του πληθυσμού αν υπολογιστούν από το τυχαίο δείγμα, δηλαδή δειγματικά, που λαμβάνουμε από τον πληθυσμό καλούνται στατιστικές συναρτήσεις ή στατιστικά. Υπολογίζοντας την τιμή μιας στατιστικής συνάρτησης στο δείγμα μπορούμε να κάνουμε γενίκευση στον πληθυσμό. Η στατιστική ανάλυση βασίζεται κυρίως σε στατιστικές συναρτήσεις οι οποίες είναι συναρτήσεις του τυχαίου δείγματός μας. Αξίζει να σημειωθεί, ότι σε ένα τυχαίο δείγμα (iid παρατηρήσεις), οι τυχαίες μεταβλητές είναι ανεξάρτητες και ισόνομες. Πριν τη συλλογή δεδομένων μια στατιστική συνάρτηση είναι μια τυχαία μεταβλητή η οποία έχει κατανομή πιθανότητας που ονομάζεται δειγματική κατανομή του στατιστικού. Ας υποθέσουμε ότι έχουμε δείγμα $\mathbf{X} = (X_1, X_2, \dots, X_n)$ ανεξάρτητων και ισόνομων, δηλαδή iid, παρατηρήσεων από πληθυσμό με κατανομή F . Έχοντας συλλέξει τις παρατηρήσεις του δείγματος για κάθε τυχαία μεταβλητή $X_i, i = 1, 2, \dots, n$ έχουμε τιμές x_1, x_2, \dots, x_n από τις οποίες μπορεί να εκτιμηθεί κάποια παράμετρος $\vartheta(F)$ της πληθυσμιακής κατανομής F της οποίας η μορφή είναι άγνωστη. Είναι φανερό λοιπόν ότι η παράμετρος ϑ είναι μια συνάρτηση της F . Χρησιμοποιώντας μια στατιστική συνάρτηση $g_n(\mathbf{X})$, δηλαδή συνάρτηση των τιμών του δείγματος ή αλλιώς εκτιμητήρια συνάρτηση, μπορούμε να εκτιμήσουμε την παράμετρο $\vartheta(F)$. Η τιμή αυτής της στατιστικής συνάρτησης για τις παρατηρήσεις του δείγματος $g_n(x_1, \dots, x_n)$ αποτελεί την εκτίμηση $\hat{\vartheta}_n$ της παραμέτρου ϑ , δηλαδή είναι μια τιμή που υπολογίζεται με βάση τα δεδομένα του δείγματος μας και προσεγγίζει την πραγματική τιμή της αντίστοιχης παραμέτρου του πληθυσμού.

Ας υποθέσουμε ότι $\hat{\vartheta}_n$ είναι ένα γενικό στατιστικό που υπολογίζεται από το δείγμα μας (δείγμα μεγέθους n).

$$\hat{\vartheta}_n = g_n(x_1, \dots, x_n)$$

για κάποια συνάρτηση $g_n(\cdot)$.

Σε ένα πρόβλημα εκτίμησης είναι απαραίτητο να υπάρχει μια ένδειξη της ακρίβειας του εκτιμητή μέσω μιας εκτίμησης του τυπικού σφάλματος του εκτιμητή ή ενός διαστήματος εμπιστοσύνης. Για να το κάνουμε αυτό χρειαζόμαστε να εκτιμήσουμε τη δειγματική κατανομή του $\hat{\theta}_n$. Αξίζει να σημειωθεί, ότι για διαφορετικά δείγματα, δηλαδή διαφορετικές τιμές $\{x_1, \dots, x_n\}$, η εκτιμήτρια συνάρτηση της παραμέτρου $\hat{\theta}_n$ παίρνει διαφορετικές τιμές. Δηλαδή, το $\hat{\theta}_n$ παρουσιάζει μεταβλητότητα μεταξύ διαφορετικών δειγμάτων. Αν είχαμε πολλαπλά δείγματα θα μπορούσαμε να πάρουμε πολλές τιμές του εκτιμητή και συνεπώς να μετρήσουμε την μεταβλητότά του εμπειρικά. Σύμφωνα με τους Politis and McElroy (2020), αυτή είναι και η προσέγγιση της μεθόδου bootstrap που ανήκει στην ευρύτερη κατηγορία μεθόδων επαναδειγματοληψίας (resampling methods) και η οποία δημιουργεί πολλαπλά δείγματα με σκοπό να μετρηθεί η μεταβλητότητα του εν λόγω στατιστικού εμπειρικά. Η μέθοδος bootstrap είναι μέθοδος που προσεγγίζει την δειγματική κατανομή ενός στατιστικού και εκτιμάει τα χαρακτηριστικά της (Shao & Tu, 1995). Η ιδέα του bootstrap είναι να χρησιμοποιήσουμε μόνο την πληροφορία που διαθέτουμε από τα δεδομένα που έχουμε στην διάθεσή μας χωρίς να κάνουμε κάποια υπόθεση για την κατανομή του πληθυσμού.

1.2 Η μέθοδος Monte Carlo

Μια ακολουθία τυχαίων μεταβλητών X_1, X_2, \dots , που είναι ορισμένες στον ίδιο δειγματικό χώρο Ω , λέμε ότι συγκλίνει στην τυχαία μεταβλητή X , η οποία είναι επίσης ορισμένη στον Ω , σχεδόν βεβαίως ή με πιθανότητα 1, αν ισχύει

$$P\left(\lim_{n \rightarrow \infty} X_n = X\right) = 1$$

όπου με $P\left(\lim_{n \rightarrow \infty} X_n = X\right)$ εννοούμε την πιθανότητα του ενδεχομένου

$$A = \left\{ \omega \in \Omega: \lim_{n \rightarrow \infty} X_n(\omega) = X(\omega) \right\}.$$

Η σχεδόν βέβαια σύγκλιση συνεπάγεται ότι, σε ένα σύνολο του οποίου η πιθανότητα εμφάνισης είναι ίση με 1, για κάθε δειγματικό σημείο του ω , η διαφορά $X_n(\omega) - X(\omega)$, όταν το n μεγαλώνει απεριόριστα, γίνεται αυθαίρετα μικρή κατ' απόλυτη τιμή.

Σύμφωνα με τον ισχυρό νόμο των μεγάλων αριθμών αν

$$X_1, X_2, \dots \sim iid F$$

τότε

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{\sigma. \beta} E_F(X) \text{ καθώς } n \uparrow \infty$$

με την προϋπόθεση ότι $E_F(X)$ υπάρχει και είναι πεπερασμένη. Αυτό πρακτικά σημαίνει ότι για αρκετά μεγάλο n ισχύει $\bar{X} \approx E_F(X)$. Αξίζει να σημειωθεί, ότι για οποιαδήποτε συνάρτηση g που είναι συνεχής στο σημείο $E_F(X)$ ισχύει

$$g(\bar{X}) \xrightarrow{\sigma. \beta} g(E_F(X)) \text{ καθώς } n \uparrow \infty$$

Γενικότερα, για κάθε συνάρτηση h για την οποία υπάρχει και είναι πεπερασμένη η $E_F\{h(X)\}$,

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h(X_i) \xrightarrow{\sigma. \beta} E_F\{h(X)\} \text{ καθώς } n \uparrow \infty$$

κάτι το οποίο πρακτικά σημαίνει ότι για αρκετά μεγάλο n έχουμε $\bar{h} \approx E_F\{h(X)\}$

Στις μεθόδους προσομοίωσης για τη μελέτη στοχαστικών φαινομένων, το στοχαστικό φαινόμενο αναπαρίσταται εικονικά με τη βοήθεια ενός ηλεκτρονικού υπολογιστή στο οποίο παρακολουθείται η εξέλιξή του και καταγράφονται εκείνα τα χαρακτηριστικά που μας ενδιαφέρουν. Τα παραπάνω βήματα επαναλαμβάνονται με σκοπό να προκύψουν στατιστικά συμπεράσματα. Ας υποθέσουμε τώρα ότι θέλουμε να βρούμε την μέση τιμή

$$E_F\{h(X)\} = \int h(x)dF(x) = \begin{cases} \sum h(x)f(x), & F \text{ διακριτή με } \sigma. \mu. \pi f \\ \int h(x)f(x)dx, & F \text{ απολύτως συνεχής με } \sigma. \pi. \pi f \end{cases}$$

Προσομοιώνοντας ένα πολύ μεγάλο τυχαίο δείγμα από την F τότε, σύμφωνα με τα προηγούμενα, θα μπορούσαμε να εκμεταλλευτούμε τον ισχυρό νόμο των μεγάλων αριθμών και να την προσεγγίσουμε με τον μέσο όρο

$$\bar{h} = \frac{1}{n} \sum_{i=1}^n h(X_i)$$

Συνεπώς, έχοντας δεδομένη την κατανομή F , αυτό μπορούμε να το κάνουμε μέσω ενός ηλεκτρονικού υπολογιστή. Η παραπάνω διαδικασία που περιγράψαμε είναι γνωστή ως ολοκλήρωση Monte Carlo.

Γενικότερα, οι μέθοδοι προσομοίωσης είναι χρήσιμες σε περίπλοκα προβλήματα όταν οι αναλυτικές μέθοδοι αποτυγχάνουν. Στις αναλυτικές μεθόδους πραγματοποιείται η κατάλληλη μαθηματική μοντελοποίηση του στοχαστικού φαινομένου και αυτό μελετάται αναλυτικά. Όταν δεν είμαστε σε θέση να βρούμε την κατανομή μιας πολύπλοκης στατιστικής συνάρτησης όταν αυτή προέρχεται από πληθυσμό με γνωστή κατανομή τότε μπορούμε να καταφύγουμε στη λύση

της προσομοίωσης Monte Carlo. Ας υποθέσουμε τώρα ότι γνωρίζουμε την κατανομή F του πληθυσμού, μια πολύ δύσκολη υπόθεση σε πραγματικά δεδομένα. Σε αυτήν την περίπτωση η μεροληψία και η διακύμανση ενός εκτιμητή θα μπορούσε να υπολογιστεί ακριβώς από αναλυτικές μεθόδους ή κατά προσέγγιση από προσομοίωση Monte Carlo σε περίπτωση που οι αναλυτικοί υπολογισμοί είναι δύσκολοι. Αξίζει να σημειωθεί πως αν καταφέρουμε να προσομοιώσουμε δείγμα από την κατανομή του πληθυσμού τότε αυτόματα μπορούμε να έχουμε μια πολύ καλή εικόνα για την κατανομή της συνάρτησης που μας ενδιαφέρει αφού χρησιμοποιήσουμε τα χαρακτηριστικά αυτού του δείγματος. Η ιδέα πίσω από την προσομοίωση Monte Carlo είναι η εξής: Δεδομένου ότι ο πληθυσμός είναι γνωστός μπορούμε να προσομοιώσουμε οποιοδήποτε πλήθος iid παρατηρήσεων από αυτόν. Συνεπώς, δεδομένου ότι η κατανομή F είναι γνωστή και θέλουμε να εκτιμήσουμε την διακύμανση $\text{var}(\hat{\theta})$ του εκτιμητή $\hat{\theta}$ μιας παραμέτρου του πληθυσμού θ για κάποιο μέγεθος δείγματος n , τότε προσομοιώνουμε πολλαπλά ανεξάρτητα αντίγραφα της τυχαίας μεταβλητής $\hat{\theta}$ τα οποία υπολογίζονται από τα προσομοιωμένα n -διάστατα τυχαία δείγματα, δηλαδή από τα δείγματα n παρατηρήσεων, από τον γνωστό πληθυσμό F .

Προσομοίωση:

$$X_1^{(1)}, X_2^{(1)}, \dots, X_n^{(1)} \sim iid F$$

$$X_1^{(2)}, X_2^{(2)}, \dots, X_n^{(2)} \sim iid F$$

...

$$X_1^{(M)}, X_2^{(M)}, \dots, X_n^{(M)} \sim iid F$$

Συνεπώς, προσομοιώνουμε M δείγματα $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(M)}$ όπου κάθε δείγμα περιέχει n σε αριθμό iid παρατηρήσεις από τον γνωστό πληθυσμό F .

Για $j = 1, 2, \dots, M$ υπολογίζουμε τη τιμή $\hat{\theta}_n^{(j)} = T_n(x_1^{(j)}, \dots, x_n^{(j)})$. Σημειώνεται ότι ένα ιστόγραμμα των τιμών μπορεί να μας υποδείξει τη μορφή της κατανομής του $\hat{\theta}_n$.

Αν το M είναι αρκετά μεγάλο τότε σύμφωνα με τον ισχυρό νόμο των μεγάλων αριθμών:

$$E_F g(T(\mathbf{X})) \approx \frac{1}{M} \sum_{i=1}^M g(T(\mathbf{X}^{(i)}))$$

όπου $g(\cdot)$ είναι κάποια συνάρτηση. Επίσης, θα έχουμε

$$\text{Bias}_F(T) \approx \frac{1}{M} \sum_{i=1}^M T(\mathbf{X}^{(i)}) - \theta(F)$$

$$\text{Var}_F(T) \approx \frac{1}{M} \sum_{i=1}^M T^2(\mathbf{X}^{(i)}) - \left[\frac{1}{M} \sum_{i=1}^M T(\mathbf{X}^{(i)}) \right]^2$$

Η μεροληψία και η διακύμανση του εκτιμητή μας δίνουν μια εικόνα για τη μεταβλητότητα των τιμών που θα έχει ο εκτιμητής της παραμέτρου στα διάφορα δείγματα. Πιο συγκεκριμένα, η μεροληψία ενός εκτιμητή μας δείχνει πόσο διαφέρει «κατά μέσο όρο» από την πραγματική τιμή της παραμέτρου: αν όχι τότε είναι αμερόληπτος, διαφορετικά είναι ένας μεροληπτικός εκτιμητής. Η διακύμανση ενός εκτιμητή μας δείχνει την μεταβλητότητα που έχει ο εκτιμητής. Ένας εκτιμητής με μεγάλη διακύμανση και συνεπώς μεγάλο τυπικό σφάλμα μας υποδεικνύει ότι για διαφορετικά δείγματα ίδιου μεγέθους από τον ίδιο πληθυσμό θα παίρνουμε έντονα διαφορετικές τιμές του εκτιμητή. Με τη μέθοδο Monte Carlo μπορούμε να εκτιμήσουμε την μεροληψία και την διακύμανση του εκτιμητή εμπειρικά από την μεροληψία και την διακύμανση των προσομοιωμένων δειγμάτων. Δηλαδή, εκτιμούμε τη μεταβλητότητα που έχει ο εκτιμητής στα διαφορετικά δείγματα από την μεταβλητότητα που έχει ο εκτιμητής στα προσομοιωμένα δείγματα. Ωστόσο, σε πρακτικά προβλήματα η πληθυσμιακή κατανομή F δεν είναι γνωστή συνεπώς εδώ έρχεται να μας βοηθήσει η μέθοδος bootstrap. Η ιδέα είναι απλή: αφού δεν έχουμε ολόκληρο τον πληθυσμό κάνουμε ό,τι μπορούμε με αυτά που έχουμε δηλαδή με το παρατηρούμενο δείγμα (Politis, 1998).

1.3 Γενική παρουσίαση Bootstrap για εκτίμηση τυπικού σφάλματος

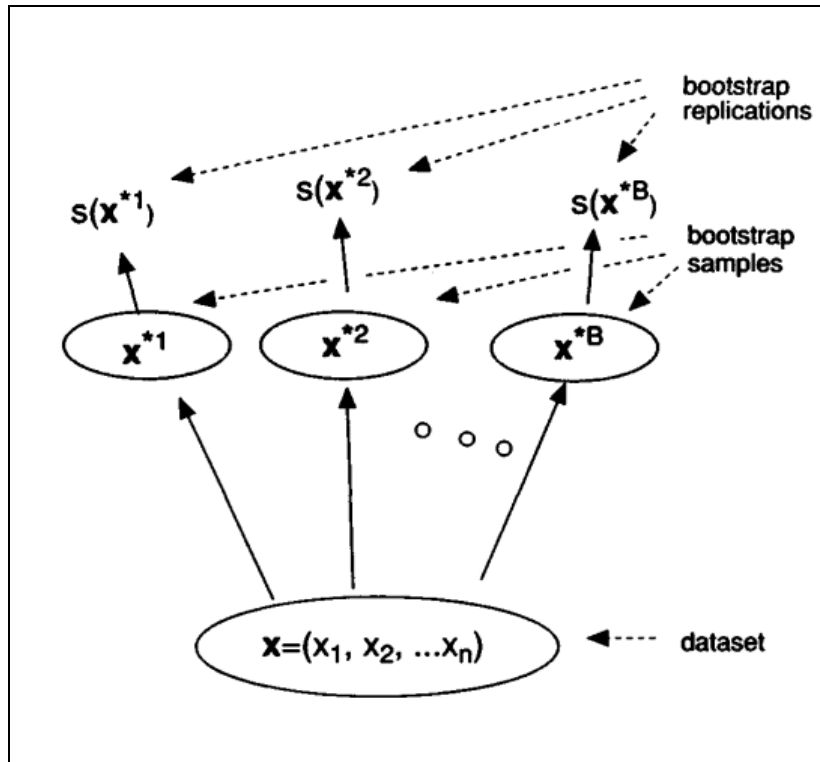
Αν γνωρίζουμε ότι οι τυχαίες μεταβλητές X_i ακολουθούν κανονική κατανομή, δηλαδή προέρχονται από κανονικό πληθυσμό, με μέση τιμή μ και διακύμανση σ^2 και η παράμετρος που μας ενδιαφέρει θ είναι η μέση τιμή του πληθυσμού, δηλαδή $\theta = \mu$, τότε ο εκτιμητής $\hat{\theta}$ θα είναι η δειγματική μέση τιμή $\bar{x} = \sum_{i=1}^n x_i/n$. Η εκτίμηση του τυπικού σφάλματος του εκτιμητή \bar{x} γνωρίζουμε ότι δίνεται από το $\sqrt{s^2/n}$, όπου $s^2 = \sum_{i=1}^n (x_i - \bar{x})^2/(n - 1)$ είναι η δειγματική διασπορά. Επίσης, σε αυτήν την περίπτωση είναι εύκολο να κατασκευάσουμε διαστήματα εμπιστοσύνης για το θ είτε όταν το σ είναι γνωστό είτε όταν είναι άγνωστο. Σύμφωνα με το κεντρικό οριακό θεώρημα ακόμα και αν δεν γνωρίζουμε την πληθυσμιακή κατανομή από την οποία προέρχονται τα δεδομένα, για αρκετά μεγάλο μέγεθος δείγματος, καθώς $n \rightarrow \infty$, $\bar{x} \approx N(\mu, \sigma^2/n)$. Πιο συγκεκριμένα, το κεντρικό οριακό θεώρημα αναφέρει ότι, το άθροισμα και επομένως η μέση τιμή, μεγάλου αριθμού ανεξάρτητων παρατηρήσεων, ακολουθεί κατά προσέγγιση κανονική κατανομή ανεξαρτήτως από ποια κατανομή προέρχονται οι παρατηρήσεις (υπό την προϋπόθεση ότι η διασπορά της είναι πεπερασμένη). Δηλαδή με λίγα λόγια υποδεικνύει ότι οι δειγματικοί μέσοι συμπεριφέρονται για μεγάλο n περίπου όπως όταν

έχουμε κανονική κατανομή. Το τυπικό σφάλμα ενός εκτιμητή δείχνει την ακρίβεια του εκτιμητή. Ωστόσο, το τυπικό σφάλμα είναι δύσκολο να εκτιμηθεί σε περίπτωση άλλων εκτιμητών πέραν του μέσου και ειδικότερα σε πολύπλοκες στατιστικές συναρτήσεις που η εύρεση της συνάρτησης κατανομής τους είναι σχεδόν αδύνατη. Αξίζει να σημειωθεί ότι σε αρκετά παραμετρικά αλλά ειδικότερα σε μη παραμετρικά μοντέλα στα οποία η κατανομή F είναι άγνωστη δεν μπορούμε εύκολα να προσδιορίσουμε μια κατάλληλη στατιστική συνάρτηση που στην συνέχεια θα εκτιμήσουμε τα χαρακτηριστικά της. Συνεπώς, εδώ έρχεται και μας βοηθάει η μέθοδος bootstrap. Η μεθοδολογία Monte Carlo που είδαμε παραπάνω βασίστηκε στην πολύ σημαντική υπόθεση ότι η κατανομή του πληθυσμού είναι γνωστή. Η μέθοδος bootstrap ξεπερνάει τα προβλήματα της κλασικής μεθόδου Monte Carlo καθώς δεν χρειάζεται να γνωρίζουμε την κατανομή από την οποία προέρχονται τα δεδομένα μας. Δηλαδή, είναι μια μη παραμετρική μέθοδος. Σύμφωνα με τους Efron and Hastie (2016) ο 21^{ος} αιώνας χαρακτηρίζεται από ταυτόχρονη διαθεσιμότητα μεγάλου όγκου δεδομένων (Big Data) και ισχυρής υπολογιστικής δύναμης, επομένως οι επαγγελματίες του χώρου έχουν αρχίσει σταδιακά να αμφισβητούν την λογική της έναρξης με ένα περιοριστικό και συχνά αβάσιμο σύνολο παραμετρικών υποθέσεων, όπως για παράδειγμα αυτό της κανονικότητας. Ως αποτέλεσμα η οπτική της μη παραμετρικής στατιστικής γίνεται συνεχώς και πιο σημαντική, και υπολογιστικές μέθοδοι όπως η μέθοδος bootstrap γίνονται ολοένα και πιο σχετικές.

Η μέθοδος bootstrap χρησιμοποιεί την εμπειρική κατανομή \hat{F} , δηλαδή την κατανομή που δίνει πιθανότητα $1/n$ σε κάθε μια από τις n παρατηρήσεις του δείγματος μας και 0 σε οποιαδήποτε άλλη τιμή. Δηλαδή, η απλή ιδέα είναι να αντικαταστήσει την άγνωστη πληθυσμιακή κατανομή με την γνωστή εμπειρική κατανομή, συνεπώς θα χρησιμοποιήσουμε μια ποσότητα του δείγματος για να εκτιμήσουμε την αντίστοιχη ποσότητα του πληθυσμού. Οι ιδιότητες τώρα ενός εκτιμητή, όπως το τυπικό σφάλμα του, προσδιορίζονται από την εμπειρική. Μερικές φορές αυτές οι ιδιότητες μπορούν να προσδιοριστούν αναλυτικά, αλλά πιο συχνά προσεγγίζονται χρησιμοποιώντας την μέθοδο Monte Carlo, δηλαδή κάνοντας δειγματοληψία με επανάθεση από την εμπειρική κατανομή (Chernick, 2007).

Ένα δείγμα bootstrap $\mathbf{X}^* = (x_1^*, \dots, x_n^*)$ (το * χρησιμοποιείται για να μην υπάρχει σύγχυση με το αρχικό δείγμα) είναι ένα τυχαία επιλεγμένο με επανάθεση δείγμα από τα αρχικά μας δεδομένα x_1, \dots, x_n . Αν $n = 7$ τότε ένα δείγμα bootstrap ίσως είναι το $\mathbf{X}^* = (x_5, x_2, x_3, x_3, x_4, x_7, x_1)$.

Σχήμα 1.1



Πηγή: Efron & Tibshirani (1993)

Από το Σχήμα 1.1 μπορούμε να διακρίνουμε την διαδικασία της μεθόδου bootstrap για την εκτίμηση του τυπικού σφάλματος του στατιστικού $s(\mathbf{X})$.

Συνεπώς, στον αλγόριθμο bootstrap από το αρχικό δείγμα δημιουργούνται B τυχαία δείγματα, με επανάθεση, μεγέθους n όπως το αρχικό. Συνήθως, η τιμή του B για την εκτίμηση ενός τυπικού σφάλματος κυμαίνεται ανάμεσα σε 50 και 200 (Efron & Tibshirani, 1993).

Προσομοίωση B bootstrap δειγμάτων από την \hat{F} :

$$X_1^{*(1)}, X_2^{*(1)}, \dots, X_n^{*(1)} \sim iid \hat{F}$$

$$X_1^{*(2)}, X_2^{*(2)}, \dots, X_n^{*(2)} \sim iid \hat{F}$$

...

$$X_1^{*(B)}, X_2^{*(B)}, \dots, X_n^{*(B)} \sim iid \hat{F}$$

Σε κάθε δείγμα bootstrap υπολογίζουμε τη τιμή της στατιστικής συνάρτησης $s(\mathbf{X})$, δηλαδή την $s(\mathbf{X}^{*b})$, και ως αποτέλεσμα από τα B δείγματα bootstrap παίρνουμε στη διάθεση μας τις τιμές $s(\mathbf{X}^{*1}), s(\mathbf{X}^{*2}), \dots, s(\mathbf{X}^{*B})$. Στο τέλος της διαδικασίας η δειγματική τυπική απόκλιση των τιμών $s(\mathbf{X}^{*1}), s(\mathbf{X}^{*2}), \dots, s(\mathbf{X}^{*B})$ είναι η εκτίμηση του τυπικού σφάλματος του $s(\mathbf{X})$.

$$\widehat{se}_{boot} = \{\sum_{b=1}^B [s(\mathbf{X}^{*b}) - s(\cdot)]^2 / (B - 1)\}^{1/2}$$

όπου $s(\cdot) = \sum_{b=1}^B s(\mathbf{X}^{*b}) / B$

Αρχικά από δείγμα \mathbf{X} iid από την F μπορούμε να εκτιμήσουμε την τιμή κάποιας πληθυσμιακής παραμέτρου ϑ .

$$\begin{array}{c} iid \quad s \\ F \rightarrow \mathbf{X} \rightarrow \hat{\vartheta} \end{array}$$

Δεν γνωρίζουμε την F , αλλά μπορούμε να την εκτιμήσουμε από την \hat{F} . Ένα δείγμα bootstrap \mathbf{X}^* είναι ένα iid δείγμα από την \hat{F}

$$\begin{array}{c} iid \quad s \\ \hat{F} \rightarrow \mathbf{X}^* \rightarrow \hat{\vartheta}^* \end{array}$$

Συνεπώς, από κάθε δείγμα bootstrap υπολογίζουμε τον εκτιμητή $\hat{\vartheta}$, δηλαδή:

$$\hat{\vartheta}^*(b) = s(\mathbf{X}^{*b}) \quad b = 1, 2, \dots, B$$

Στον πραγματικό κόσμο θα πάρουμε μόνο μια τιμή για το $\hat{\vartheta}$ από το δείγμα μας, αλλά ο κόσμος του bootstrap είναι πιο γενναιόδωρος και μπορούμε να δημιουργήσουμε πολλές τιμές του εκτιμητή $\hat{\vartheta}^*$ ο οποίος θα υπολογίζεται από κάθε δείγμα bootstrap και επομένως θα μπορούμε να εκτιμήσουμε την μεταβλητότητα αυτών των τιμών απευθείας μέσω του \widehat{se}_{boot} . Επομένως, η βασική ιδέα πίσω από τη μέθοδο bootstrap είναι ότι η μεταβλητότητα των τιμών $\hat{\vartheta}^*$ γύρω από το $\hat{\vartheta}$ θα είναι παρόμοια με την μεταβλητότητα των $\hat{\vartheta}$ γύρω από την πραγματική τιμή της παραμέτρου ϑ . Αξίζει να σημειωθεί ότι σε πολύ λίγες περιπτώσεις μπορούμε να υπολογίσουμε τον εκτιμητή bootstrap χωρίς να εφαρμόσουμε την προσέγγιση της μεθόδου Monte Carlo. Σύμφωνα με τον Efron (1982) στην περίπτωση που ο εκτιμητής είναι ο μέσος της κατανομής μιας τυχαίας μεταβλητής, τότε η εκτίμηση bootstrap του τυπικού σφάλματος του δειγματικού μέσου δίνεται από τον τύπο

$$\widehat{se}_{boot} = [(n - 1)/n]^{1/2} \hat{\sigma}$$

όπου $\hat{\sigma} = \left[\frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2}$.

1.4 Εμπειρική συνάρτηση κατανομής

Έστω $\mathbf{X} = (X_1, \dots, X_n)$ τυχαίο δείγμα, μεγέθους n από κάποια κατανομή F . Σημειώνεται ότι χρησιμοποιούμε το σύμβολο F που συνήθως συμβολίζει αθροιστική συνάρτηση κατανομής

αλλά μπορούμε να συμβολίσουμε έτσι, γενικά, και την κατανομή. Έχοντας παρατηρήσει τα δεδομένα δηλαδή $X_1 = x_1, \dots, X_n = x_n$, η εμπειρική κατανομή είναι η διακριτή κατανομή που δίνει πιθανότητα $1/n$ σε κάθε $x_i, i = 1, 2, \dots, n$ και 0 σε οποιαδήποτε άλλη τιμή. Η εμπειρική συνάρτηση κατανομής \hat{F} , η οποία είναι διακριτή ανεξάρτητα από το τι είναι η κατανομή F , είναι η συνάρτηση

$$\hat{F}(x) = \frac{1}{n} \sum_{i=1}^n I(X_i \leq x) = \#\{X_i \leq x\}/n, \quad x \in \mathbb{R},$$

όπου $I(X_i \leq x)$ είναι η δείκτρια συνάρτηση που λαμβάνει τιμή 1 εάν $X_i \leq x$ και τιμή μηδέν αν $X_i > x$.

Η εμπειρική κατανομή συγκλίνει ισχυρά και ομοιόμορφα στην πραγματική κατανομή F . Το θεώρημα Glivenko-Cantelli μας λέει ότι η απόσταση μεταξύ της εμπειρικής συνάρτησης κατανομής και της συνάρτησης κατανομής F συγκλίνει σχεδόν βεβαίως στο μηδέν:

$$\sup_{x \in \mathbb{R}} |\hat{F}(x) - F(x)| \rightarrow 0 \quad \sigma. \beta, \text{ καθώς } n \uparrow \infty$$

Συνεπώς καθώς μεγαλώνει το μέγεθος του δείγματος η εμπειρική συνάρτηση κατανομής προσεγγίζει την πραγματική συνάρτηση κατανομής. Δηλαδή $\hat{F} \rightarrow F$ όταν $n \rightarrow \infty$. Η μέθοδος bootstrap χρησιμοποιεί την εμπειρική κατανομή αντί για κάποια γνωστή κατανομή όπως κάνει η κλασική μέθοδος Monte Carlo. Για να παράγουμε ένα δείγμα από την εμπειρική κατανομή πραγματοποιούμε δειγματοληψία με επανάθεση στις παρατηρηθείσες τιμές. Σημειώνεται ότι αν το αρχικό δείγμα περιέχει πολλαπλά αντίγραφα κάποιας συγκεκριμένης τιμής τότε κάθε αντίγραφο έχει πιθανότητα $1/n$, δηλαδή σε κάθε διακριτό X_i η εμπειρική συνάρτηση κατανομής δίνει πιθανότητα ανάλογη με τη συχνότητα που εμφανίζεται αυτή η παρατήρηση στο δείγμα (Trosset, 2009). Αξίζει να σημειωθεί, ότι από την σύγκλιση της εμπειρικής κατανομής \hat{F} στην πραγματική κατανομή F προκύπτει και η σύγκλιση των πιο ενδιαφέροντων συναρτησιακών της εμπειρικής κατανομής στα αντίστοιχα της πραγματικής κατανομής:

$$T(\hat{F}) \rightarrow T(F) \text{ συνήθως } \sigma. \beta$$

1.5 Εκτιμητής αντικατάστασης

Η αρχή της αντικατάστασης (plug-in principle) είναι μια απλή μέθοδος εκτίμησης παραμέτρων από το δείγμα. Ο εκτιμητής αντικατάστασης μιας παραμέτρου $\vartheta = T(F)$ είναι ο $\hat{\vartheta} = T(\hat{F})$. Δηλαδή εκτιμάμε την συνάρτηση $\vartheta = T(F)$, χρησιμοποιώντας την ίδια συνάρτηση στην εμπειρική συνάρτηση κατανομής \hat{F} , $\hat{\vartheta} = T(\hat{F})$ (Efron & Tibshirani, 1993).

Σε πολλές περιπτώσεις παίρνουμε ως εκτιμητή, το δειγματικό (εμπειρικό) ανάλογο της ποσότητας που μας ενδιαφέρει. Για παράδειγμα, ο εκτιμητής αντικατάστασης (plug-in estimator) της αναμενόμενης μέσης τιμής $\vartheta = E_F(X) = \mu(F)$ είναι ο δειγματικός μέσος $\hat{\vartheta} = E_{\hat{F}}(X) = \mu(\hat{F}) = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$.

Ο εκτιμητής αντικατάστασης της διακύμανσης $\vartheta = Var_F(X) = \sigma^2(F)$ είναι περίπου η δειγματική διασπορά $\hat{\vartheta} = Var_{\hat{F}}(X) = \sigma^2(\hat{F}) = \frac{1}{n} \sum_{i=1}^n x_i^2 - (\bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$, συνεπώς ο εκτιμητής αντικατάστασης της διακύμανσης έχει τη μορφή του εκτιμητή μέγιστης πιθανοφάνειας στην κανονική κατανομή. Είναι περίπου η δειγματική διασπορά με τη διαφορά ότι στον παρονομαστή αντί για $n - 1$ έχει n και επομένως δεν είναι ακριβώς αμερόληπτος εκτιμητής της διασποράς του πληθυσμού. $E \left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right] = \frac{n-1}{n} \sigma^2$

Πιο συγκεκριμένα, αν ϑ είναι η μέση τιμή $\vartheta \equiv E_F(X) = \mu(F)$ ή η διασπορά $\vartheta \equiv Var_F(X) = \sigma^2(F)$ τότε

$$\vartheta = \int_{-\infty}^{\infty} x dF(x) \quad \text{ή} \quad \vartheta = \int_{-\infty}^{\infty} x^2 dF(x) - \left(\int_{-\infty}^{\infty} x dF(x) \right)^2$$

Είναι γνωστό ότι η συνάρτηση κατανομής F μπορεί να εκτιμηθεί χωρίς να κάνουμε κάποια υπόθεση για την μορφή της από την εμπειρική συνάρτηση κατανομής που προκύπτει από το δείγμα. Συνεπώς, η ιδέα εδώ είναι να χρησιμοποιήσουμε την εμπειρική συνάρτηση κατανομής \hat{F} για να εκτιμήσουμε και το ϑ . Αν ϑ είναι η μέση τιμή της κατανομής τότε μπορούμε να θεωρήσουμε ως εκτιμητή τον

$$\hat{\vartheta} = E_{\hat{F}}(X) = \mu(\hat{F}) = \int_{-\infty}^{\infty} x d\hat{F}(x) = E(X^*)$$

όπου X^* είναι μια τυχαία μεταβλητή με κατανομή την εμπειρική \hat{F} . Αν οι τιμές του δείγματος X_1, \dots, X_n είναι x_1, \dots, x_n τότε

$$P(X^* = x_i) = \hat{F}(x_i) - \hat{F}(x_{i-1}) = \frac{1}{n}, i = 1, \dots, n$$

Επομένως, αν ϑ είναι η μέση της κατανομής F τότε

$$\hat{\vartheta} = E_{\hat{F}}(X) = \mu(\hat{F}) = E(X^*) = \sum_{i=1}^n x_i P(X^* = x_i) = \frac{1}{n} \sum_{i=1}^n x_i = \bar{x}$$

Αν ϑ είναι η διασπορά της κατανομής F τότε

$$\begin{aligned}
\hat{\theta} &= \text{Var}_{\hat{F}}(X) = \sigma^2(\hat{F}) = \text{Var}(X^*) \\
&= \int_{-\infty}^{\infty} x^2 d\hat{F}(x) - \left(\int_{-\infty}^{\infty} x d\hat{F}(x) \right)^2 \\
&= \sum_{i=1}^n x_i^2 P(X^* = x_i) - \left(\sum_{i=1}^n x_i P(X^* = x_i) \right)^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \\
&= \frac{1}{n} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{n-1}{n} s^2
\end{aligned}$$

Σε περίπτωση που η F είναι διδιάστατη κατανομή, ο συντελεστής συσχέτισης είναι η ποσότητα

$$\vartheta = \rho = \rho(F) = \frac{\text{Cov}_F(X, Y)}{\sqrt{\text{Var}_F(X)\text{Var}_F(Y)}}$$

σε αυτήν την περίπτωση ο εκτιμητής αντικατάστασης είναι

$$\hat{\theta} = \rho(\hat{F}) = \frac{\text{Cov}_{\hat{F}}(X, Y)}{\sqrt{\text{Var}_{\hat{F}}(X)\text{Var}_{\hat{F}}(Y)}} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}} \equiv r$$

δηλαδή, ο δειγματικός συντελεστής συσχέτισης. Συνεπώς, καταλαβαίνουμε ότι μπορούμε να προτείνουμε έναν εκτιμητή για οποιαδήποτε παράμετρο ϑ μιας άγνωστης κατανομής F χωρίς να κάνουμε καμία υπόθεση για την μορφή της F . Η συνάρτηση κατανομής της στατιστικής συνάρτησης T , εξαρτάται από την F που μπορεί να εκτιμηθεί από την εμπειρική συνάρτηση κατανομής \hat{F} , συνεπώς προκύπτει η βασική ιδέα της μεθόδου bootstrap η οποία είναι να εκτιμήσουμε την κατανομή F_T της στατιστικής συνάρτησης T χρησιμοποιώντας αντί της F την εμπειρική συνάρτηση κατανομής \hat{F} . Ειδικότερα εκτιμάμε την κατανομή της

$$T = T(X_1, \dots, X_n) \text{ όπου } X_i \sim F$$

από την κατανομή της τυχαίας μεταβλητής

$$T^* = T(X_1^*, \dots, X_n^*) \text{ όπου } X_i^* \sim \hat{F}$$

Σαν αποτέλεσμα, όλα τα ζητούμενα χαρακτηριστικά της T μπορούν να εκτιμηθούν από τα αντίστοιχα χαρακτηριστικά της T^* . Για παράδειγμα η μέση τιμή της T εκτιμάται από την

$$\begin{aligned}
E(T^*) &= E(T(X_1^*, \dots, X_n^*)) = \sum_{i_1=1}^n \dots \sum_{i_n=1}^n T(x_{i_1}, \dots, x_{i_n}) P(X_1^* = x_{i_1}) \dots P(X_n^* = x_{i_n}) \\
&= \frac{1}{n^n} \sum_{i_1=1}^n \dots \sum_{i_n=1}^n T(x_{i_1}, \dots, x_{i_n})
\end{aligned}$$

ενώ γενικότερα η μέση τιμή μιας συνάρτησης της T , $E(g(T^*))$ εκτιμάται από

$$E(g(T^*)) = E\left(g(T(X_1^*, \dots, X_n^*))\right) = \frac{1}{n^n} \sum_{i_1=1}^n \dots \sum_{i_n=1}^n g(T(x_{i_1}, \dots, x_{i_n}))$$

Οι παραπάνω εκτιμήσεις καλούνται εκτιμήσεις bootstrap των χαρακτηριστικών της T . Σε ορισμένες ειδικές περιπτώσεις όπως όταν $T = \bar{X}$ η παραπάνω τιμή μπορεί να υπολογιστεί. Η περίπτωση του \bar{X} , και γενικότερα στατιστικών συναρτήσεων που εκφράζονται ως δειγματικοί μέσοι είναι πολύ ειδική, δεδομένου ότι υπάρχει κλειστός τύπος για τη διασπορά τους. Σε αυτήν την περίπτωση ο εκτιμητής bootstrap της διασποράς και του τυπικού σφάλματος του εκτιμητή μπορεί να υπολογιστεί άμεσα και ονομάζεται ιδανικός εκτιμητής bootstrap της διασποράς και ιδανικός εκτιμητής bootstrap του τυπικού σφάλματος αντίστοιχα. Ας υποθέσουμε ότι $X_1, \dots, X_n \sim iid F$ και $\vartheta \equiv \mu = E_F(X) = \mu(F)$. Όπως είδαμε νωρίτερα, ο εκτιμητής αντικατάστασης είναι ο $\hat{\vartheta} = \hat{\mu} = \mu(\hat{F}) = \bar{X}$. Γνωρίζουμε ότι, ισχύει $Var_F(\bar{X}) = Var_F(X)/n$, οπότε από την αρχή της αντικατάστασης, ο ιδανικός εκτιμητής bootstrap της διασποράς του \bar{X} είναι

$$Var_{\hat{F}}(\hat{\vartheta}^*) = Var_{\hat{F}}(\bar{X}^*) = \frac{Var_{\hat{F}}(X)}{n} = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n^2} = \frac{1}{n} \frac{n-1}{n} s^2$$

Όμως, γενικότερα η εκτίμηση bootstrap $E(g(T^*))$ του $E(g(T))$ θα πρέπει να υπολογίζεται από το πολλαπλό άθροισμα που αποτελείται από n^n όρους.

$$E(g(T^*)) = \frac{1}{n^n} \sum_{i_1=1}^n \dots \sum_{i_n=1}^n g(T(x_{i_1}, \dots, x_{i_n}))$$

Συνεπώς, για μικρό σχετικά δείγμα ο ακριβής υπολογισμός του παραπάνω αθροίσματος είναι πρακτικά αδύνατος. Για το λόγο αυτό καταφεύγουμε σε υπολογιστικές μεθόδους και συγκεκριμένα χρησιμοποιούμε προσομοίωση Monte Carlo ώστε να υπολογίσουμε ή τουλάχιστον να προσεγγίσουμε την παραπάνω μέση τιμή. Η προσέγγιση με την μέθοδο Monte Carlo, δειγματοληψία με επανάθεση από την εμπειρική αναφέρεται και ως ομοιόμορφη επαναδειγματοληψία σύμφωνα με τον Hall (1992).

1.6 Η μέθοδος Jackknife

Ο Quenouille (1949) παρουσίασε μια μέθοδο η οποία αργότερα ονομάστηκε μέθοδος jackknife για να εκτιμήσει τη μεροληψία ενός εκτιμητή μη παραμετρικά. Η μέθοδος διέγραφε κάθε φορά μια παρατήρηση από το αρχικό σετ δεδομένων και επαναυπολόγιζε τον εκτιμητή

από τα εναπομείναντα δεδομένα. Ας φανταστούμε δειγματοληψία χωρίς επανάθεση από τα αρχικά μας δεδομένα $\mathbf{X} = (X_1, \dots, X_n)$. Θα πάρουμε ένα δείγμα μεγέθους b όπου $b < N$. Αν $b = N - 1$ ακριβώς τότε είμαστε στην περίπτωση του κλασικού jackknife και μπορούμε να πάρουμε N διαφορετικά μικρότερα δείγματα (subsamples) από το αρχικό δείγμα.

Η μεροληψία ενός εκτιμητή $\hat{\theta}$ μιας ποσότητας θ ορίζεται ως εξής:

$$E(\hat{\theta}) - \theta$$

Ο Tukey (1958) χρησιμοποίησε την ιδέα προκειμένου να εκτιμήσει μη παραμετρικά τη διασπορά και κατ' επέκταση το τυπικό σφάλμα ενός εκτιμητή. Ας υποθέσουμε ότι έχουμε δείγμα \mathbf{X} και έναν εκτιμητή $\hat{\theta} = s(\mathbf{X})$. Θέλουμε να εκτιμήσουμε τη μεροληψία και το τυπικό σφάλμα του εκτιμητή. Η μέθοδος jackknife παίρνει «δείγματα» αφήνοντας σε κάθε δείγμα μια παρατήρηση εκτός κάθε φορά και από τις υπόλοιπες παρατηρήσεις υπολογίζει το $\hat{\theta}$. Τα

$$\mathbf{X}_{(i)} = (x_1, x_2, \dots, x_{i-1}, x_{i+1}, \dots, x_n)$$

για $i = 1, 2, \dots, n$ ονομάζονται δείγματα jackknife. Το i δείγμα jackknife περιέχει τα αρχικά δεδομένα εκτός της i παρατήρησης. Η τιμή του εκτιμητή $\hat{\theta}$ στο i δείγμα jackknife $n - 1$ παρατηρήσεων είναι η παρακάτω.

$$\hat{\theta}_{(i)} = s(\mathbf{X}_{(i)})$$

Συνεπώς στο τέλος της διαδικασίας έχουμε n τιμές $\hat{\theta}_{(i)}$. Ο μέσος όρος αυτών των τιμών των εκτιμητών είναι ο $\hat{\theta}_{(\cdot)}$.

$$\hat{\theta}_{(\cdot)} = \sum_{i=1}^n \hat{\theta}_{(i)} / n$$

Η εκτίμηση jackknife της μεροληψίας ενός εκτιμητή δίνεται από τον παρακάτω τύπο:

$$\widehat{bias}_{jack} = (n - 1)(\hat{\theta}_{(\cdot)} - \hat{\theta})$$

Ο αντίστοιχος διορθωμένος ως προς την μεροληψία εκτιμητής του θ είναι ο παρακάτω:

$$\tilde{\theta} = \hat{\theta} - (n - 1)(\hat{\theta}_{(\cdot)} - \hat{\theta}) = n\hat{\theta} - (n - 1)\hat{\theta}_{(\cdot)}$$

ενώ, η εκτίμηση jackknife του τυπικού σφάλματος δίνεται από τον παρακάτω τύπο:

$$\widehat{se}_{jack} = \left[\frac{n-1}{n} \sum_{i=1}^n (\hat{\theta}_{(i)} - \hat{\theta}_{(\cdot)})^2 \right]^{1/2}$$

1.6.1 Η μέθοδος delete-d Jackknife

Η μέθοδος jackknife δίνει καλά αποτελέσματα για εκτιμήσεις μεροληψίας και τυπικών σφαλμάτων, ωστόσο μπορεί να αποτύχει εάν η στατιστική συνάρτηση δεν είναι λεία (smooth). Η ιδέα των λείων συναρτήσεων μπορεί να διατυπωθεί ότι μικρές αλλαγές στα δεδομένα προκαλούν μικρές αλλαγές στο στατιστικό. Ένα παράδειγμα μη λείου στατιστικού είναι η διάμεσος. Αν έχουμε τις παρακάτω $n = 9$ διατεταγμένες παρατηρήσεις:

10, 27, 31, 40, 46, 50, 52, 104, 146

Παρατηρούμε ότι η διάμεσος είναι η τιμή 46, όμως αν αρχίσουμε και αυξάνουμε την τιμή της τέταρτης μεγαλύτερης παρατήρησης, δηλαδή το 40, η διάμεσος θα αλλάξει αν η τιμή της τέταρτης μεγαλύτερης παρατήρησης ξεπεράσει την τιμή 46. Η τιμή αυτή θα είναι η νέα διάμεσος όσο είναι μεγαλύτερη από 46 και μικρότερη του 50. Αυτό το παράδειγμα μπορεί να υποδείξει ότι η διάμεσος δεν είναι λεία συνάρτηση (Efron & Tibshirani, 1993). Μια μη λεία στατιστική συνάρτηση δημιουργεί στον jackknife εκτιμητή του τυπικού σφάλματος για τον μέσο μια ασυνέπεια. Για τον λόγο αυτό μπορούμε να χρησιμοποιήσουμε την μέθοδο delete-d jackknife που παρουσιάστηκε αρχικά από τον Wu (1986) και αργότερα από τους Shao and Wu (1989) η οποία ξεπερνά το παραπάνω πρόβλημα που παρουσιάζεται σε μη λείες στατιστικές συναρτήσεις. Η μέθοδος αυτή αντί να αφήνει κάθε φορά εκτός μια παρατήρηση από το δείγμα, παραλείπει d παρατηρήσεις, όπου $n = r \cdot d$ με r κάποιον ακέραιο αριθμό. Αν $n^{1/2}/d \rightarrow 0$ και $n - d \rightarrow \infty$, τότε η μέθοδος delete-d jackknife δίνει συνεπείς εκτιμήσεις για την διάμεσο. Γενικότερα, για να επιτευχθεί η απαιτούμενη συνέπεια στην εκτίμηση του τυπικού σφάλματος jackknife πρέπει να αφήσουμε εκτός δείγματος περισσότερες από $d = n^{1/2}$ παρατηρήσεις και λιγότερες από n . Ο τύπος για τον υπολογισμό της delete-d jackknife εκτίμησης του τυπικού σφάλματος είναι

$$\left\{ \frac{r}{\binom{n}{d}} \Sigma (\hat{\vartheta}_{(s)} - \hat{\vartheta}_{(\cdot)})^2 \right\}^{1/2}$$

όπου $\hat{\vartheta}_{(\cdot)} = \Sigma \hat{\vartheta}_{(s)} / \binom{n}{d}$ με $\hat{\vartheta}_{(s)}$ να είναι η εκτίμηση του ϑ στα δεδομένα έχοντας αφαιρεθεί ένα υποσύνολο παρατηρήσεων μεγέθους d , δηλαδή το παραπάνω άθροισμα είναι το άθροισμα είναι όλων των υποδειγμάτων s μεγέθους $n - d$ που έχουν επιλεγεί χωρίς επανάθεση από τα δεδομένα. Αν το n είναι μεγάλο και $n^{1/2} < d < n$ τότε το πλήθος $\binom{n}{d}$ των δειγμάτων jackknife μπορεί να είναι πολύ μεγάλο και η μέθοδος να είναι υπολογιστικά πιο δύσκολη.

1.6.2 Σχέση μεταξύ μεθόδου Jackknife και Bootstrap

Η εκτίμηση του τυπικού σφάλματος ενός εκτιμητή με τη μέθοδο bootstrap μπορεί να θεωρηθεί γενίκευση της μεθόδου jackknife με την έννοια ότι μετράει τη μεταβλητότητα των τιμών του εκτιμητή ανάμεσα σε διαφορετικά σετ δεδομένων. Η μέθοδος jackknife υπολογίζει το στατιστικό ενδιαφέροντος από n υποσύνολα μεγέθους $n - 1$ του αρχικού δείγματος, ενώ η μέθοδος bootstrap βασίζεται στην παραγωγή πολλών τυχαίων δειγμάτων bootstrap τα οποία προέρχονται από δειγματοληψία με επανάθεση από το αρχικό δείγμα. Οι εκτιμητές των μεθόδων jackknife και bootstrap ονομάζονται εκτιμητές επαναδειγματοληψίας, ενώ οι μέθοδοι αυτοί ανήκουν στην ευρύτερη κατηγορία των μεθόδων επαναδειγματοληψίας (Shao & Tu, 1995). Είναι σημαντικό να αναφερθεί ότι η μέθοδος bootstrap μπορεί να μας δώσει καλά αποτελέσματα εκεί που η κλασική μέθοδος jackknife αποτυγχάνει, δηλαδή σε στατιστικές συναρτήσεις που δεν είναι λείες, όπως στην περίπτωση της διαμέσου και αυτός είναι ένας επί πλέον λόγος που η μέθοδος bootstrap είναι τόσο δημοφιλής. Όμως, αξίζει να σημειωθεί ότι όταν μπορεί να εφαρμοστεί η μέθοδος jackknife για την εκτίμηση ενός τυπικού σφάλματος τότε μπορεί να μην είναι απαραίτητη η χρήση της μεθόδου bootstrap καθώς η jackknife απαιτεί λιγότερη υπολογιστική ισχύ. Επίσης, η μέθοδος bootstrap προσεγγίζει την κατανομή της στατιστικής συνάρτησης που μας ενδιαφέρει και από τις B τιμές $\hat{\theta}^*$ μπορούμε να έχουμε μια εικόνα για τη μορφή της κατανομής, δηλαδή αυτή η μέθοδος μας δίνει περισσότερες δυνατότητες εκτός από την εκτίμηση του τυπικού σφάλματος και της μεροληψίας του εκτιμητή. Η πραγματική υπεροχή όμως του bootstrap είναι στα διαστήματα εμπιστοσύνης. Τα διαστήματα εμπιστοσύνης bootstrap σε πολλά προβλήματα είναι ακριβέστερα από τα κλασικά ασυμπτωτικά διαστήματα, με την έννοια ότι προσεγγίζουν πιο γρήγορα τον ονομαστικό συντελεστή εμπιστοσύνης. Τέλος, αξίζει να σημειωθεί ότι σε κάποιες περιπτώσεις η μέθοδος jackknife μπορεί να χρησιμοποιηθεί για να εκτιμήσει την ακρίβεια των εκτιμήσεων της μεθόδου bootstrap (Efron, 1992). Η μέθοδος αυτή ονομάζεται jackknife after bootstrap και θα σχολιαστεί παρακάτω.

1.6.3 Η μέθοδος Cross-Validation

Οι Allen (1974) και Stone (1974) πρότειναν μια μέθοδο επιλογής μοντέλων η οποία ονομάζεται cross-validation. Η μέθοδος cross-validation μας επιτρέπει να εξετάσουμε την προσαρμογή ενός στατιστικού μοντέλου στα δεδομένα μας. Επίσης, η μέθοδος αυτή μπορεί να χρησιμοποιηθεί για να προσδιοριστεί το βέλτιστο ανάμεσα σε κάποια υποψήφια μοντέλα. Συνεπώς, είναι μέθοδος που χρησιμοποιείται τόσο για τη μελέτη της καλής προσαρμογής ενός

μοντέλου στα δεδομένα όσο και για τη σύγκριση μεταξύ μοντέλων. Αξίζει να αναφερθεί, ότι παρ' όλο που οι μέθοδοι cross-validation και jackknife είναι στενά συνδεδεμένες, η μέθοδος cross-validation δεν είναι μέθοδος που εξετάζουμε την ποιότητα των εκτιμητών όπως η μέθοδος jackknife. Η βασική ιδέα του cross-validation ή αλλιώς διασταυρούμενης επικύρωσης είναι ο τυχαίος διαχωρισμός των δεδομένων μας σε δυο υποσύνολα. Το ένα χρησιμοποιείται για να προσαρμοστεί το μοντέλο και λέγεται δείγμα εκμάθησης, και το άλλο χρησιμοποιείται για να αξιολογηθεί σε "νέα" δεδομένα αφού έχει φυσικά πρώτα προσαρμοστεί και αναφέρεται ως δείγμα επικύρωσης. Στο K-fold cross-validation χωρίζουμε τα αρχικά δεδομένα σε K ίσου μήκους υποσύνολα. Στα $K - 1$ υποσύνολα των δεδομένων τα οποία αποτελούν το δείγμα εκμάθησης προσαρμόζουμε το μοντέλο μας, και στο K υποσύνολο, το οποίο αποτελεί το δείγμα επικύρωσης, το αξιολογούμε χρησιμοποιώντας κάποιο κριτήριο δηλαδή κάποιο μέτρο που μας υποδεικνύει την απόδοση του μοντέλου. Την παραπάνω διαδικασία την πραγματοποιούμε K φορές, δηλαδή για $k = 1, 2, \dots, K$ και στο τέλος συνδυάζουμε τις K αποδόσεις του μοντέλου παίρνοντας τον μέσο όρο τους. Καταλληλότερο θεωρείται το μοντέλο με την βέλτιστη μέση απόδοση. Η ειδική περίπτωση της προσαρμογής του μοντέλου σε όλες τις παρατηρήσεις εκτός μιας και ελέγχου στην μια εναπομένονσα παρατήρηση, μερικές φορές αναφέρεται και ως leave-one-out cross-validation αλλά και ως jackknife. Αυτή η περίπτωση μοιάζει πολύ με την κλασική μέθοδο jackknife, αλλά διαφέρει στο ότι η παρατήρηση που δεν χρησιμοποιήθηκε για την εκτίμηση του μοντέλου, χρησιμοποιείται για τον έλεγχο αυτού ενώ στην μέθοδο jackknife όπως έχουμε δει δεν χρησιμοποιείται.

1.7 Βήματα αλγορίθμου Bootstrap για εκτίμηση τυπικού σφάλματος και μεροληψίας

1. Επιλέγουμε B ανεξάρτητα τυχαία δείγματα bootstrap $\mathbf{X}^{*1}, \mathbf{X}^{*2}, \dots, \mathbf{X}^{*B}$ μεγέθους n με επανάθεση από το αρχικό μας δείγμα.
2. Εκτιμάμε από κάθε δείγμα τον εκτιμητή $\hat{\theta}$, δηλαδή:

$$\hat{\theta}^*(b) = s(\mathbf{X}^{*b}) \quad b = 1, 2, \dots, B$$

3. Υπολογίζουμε το τυπικό σφάλμα από την τυπική απόκλιση των τιμών $\hat{\theta}^*(b)$ των B δειγμάτων bootstrap.

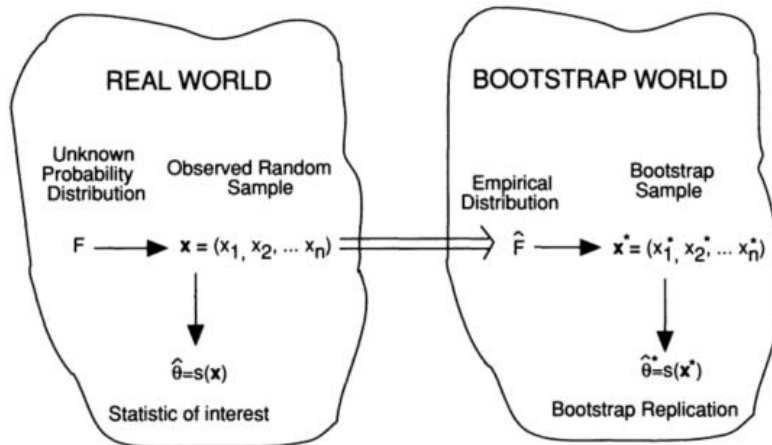
$$\widehat{se}_{boot} = \left\{ \sum_{b=1}^B [\hat{\theta}^*(b) - \hat{\theta}^*(\cdot)]^2 / (B - 1) \right\}^{1/2}$$

$$\text{όπου } \hat{\theta}^*(\cdot) = \sum_{b=1}^B \hat{\theta}^*(b) / B$$

Όταν το πλήθος των B επαναλήψεων τείνει στο άπειρο τότε η εκτίμηση bootstrap του τυπικού σφάλματος του εκτιμητή η οποία στην πραγματικότητα είναι μια εκτίμηση Monte

Carlo της ιδανικής εκτίμησης bootstrap του τυπικού σφάλματος του $\hat{\vartheta}$, δηλαδή του $\sqrt{Var_{\hat{F}}(\hat{\vartheta}^*)}$ ή αλλιώς $se_{\hat{F}}(\hat{\vartheta}^*)$, συγκλίνει στην ιδανική εκτίμηση bootstrap. Συνεπώς, το εμπειρικό τυπικό σφάλμα προσεγγίζει το πληθυσμιακό τυπικό σφάλμα καθώς ο αριθμός των επαναλήψεων αυξάνει. Όμως, ο πληθυσμός σε αυτήν την περίπτωση είναι η εμπειρική κατανομή \hat{F} καθώς βρισκόμαστε στον κόσμο του bootstrap. Στον κόσμο του bootstrap τον ρόλο της παραμέτρου τον έχει το $\hat{\vartheta}$ και τον ρόλο του εκτιμητή το $\hat{\vartheta}^*$.

Σχήμα 1.2



Πηγή: Efron & Tibshirani (1993)

Δηλαδή, όσο το πλήθος των επαναλήψεων τείνει στο άπειρο τότε η εκτίμηση \widehat{se}_{boot} συγκλίνει στην ιδανική εκτίμηση bootstrap δηλαδή στο $se_{\hat{F}}(\hat{\vartheta}^*)$

$$\lim_{B \rightarrow \infty} \widehat{se}_{boot} = se_{\hat{F}}(\hat{\vartheta}^*)$$

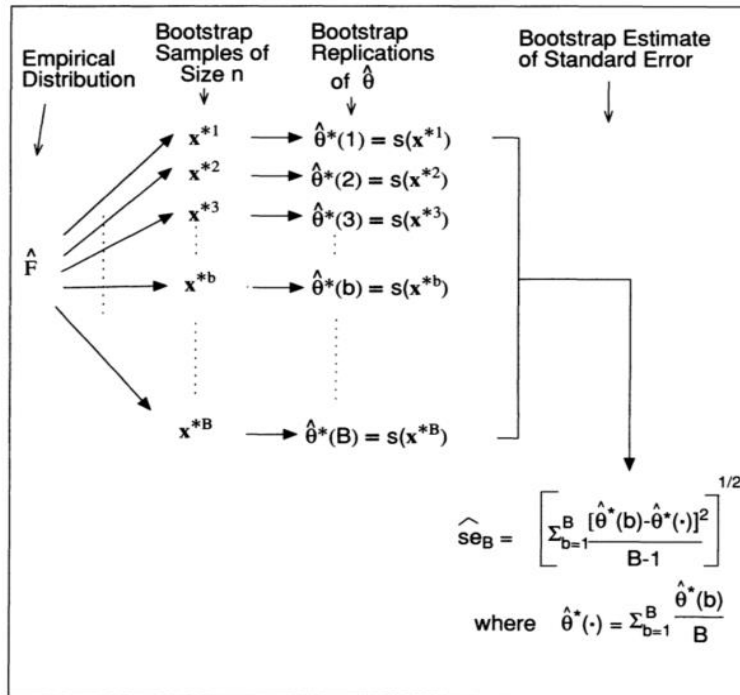
- Για μεγάλο n μέγεθος αρχικού δείγματος $se_{\hat{F}}(\hat{\vartheta}^*) \approx se_F(\hat{\vartheta})$
- Για μεγάλο B προσεγγίζουμε με Monte Carlo το $se_{\hat{F}}(\hat{\vartheta}^*)$ και έχουμε $\widehat{se}_{boot} \approx se_{\hat{F}}(\hat{\vartheta}^*)$.
- Άρα καταλήγουμε ότι για μεγάλα n και B έχουμε $\widehat{se}_{boot} \approx se_{\hat{F}}(\hat{\vartheta}^*) \approx se_F(\hat{\vartheta})$

Τα βήματα του αλγορίθμου για την εκτίμηση bootstrap της μεροληψίας είναι ακριβώς τα ίδια με την εκτίμηση bootstrap του τυπικού σφάλματος με την διαφορά ότι αλλάζει μόνο το βήμα 3. Η εκτίμηση bootstrap της μεροληψίας δίνεται παρακάτω:

$$\widehat{bias}_{boot} = \hat{\vartheta}^*(\cdot) - t(\hat{F})$$

όπου $t(\hat{F})$ είναι ο εκτιμητής plug-in του ϑ .

Σχήμα 1.3



Πηγή: Efron & Tibshirani (1993)

Από το Σχήμα 1.3 μπορούμε να διακρίνουμε την διαδικασία της μεθόδου bootstrap για την εκτίμηση του τυπικού σφάλματος του στατιστικού $s(\mathbf{X})$, μια πιο ολοκληρωμένη απεικόνιση του Σχήματος 1.1 που είδαμε νωρίτερα.

1.7.1 Jackknife after Bootstrap

Οι εκτιμητές bootstrap παρ' όλο που είναι σχεδόν αμερόληπτοι λόγω του τρόπου κατασκευής τους, όπως όλα τα στατιστικά δεν είναι ακριβείς αλλά παρουσιάζουν ένα εγγενές σφάλμα. Αυτό προκύπτει από την μεταβλητότητα της δειγματοληψίας, λόγω του γεγονότος ότι έχουμε μόνο ένα δείγμα μεγέθους n και όχι ολόκληρο τον πληθυσμό, και από την μεταβλητότητα της επαναδειγματοληψίας δηλαδή λόγω του γεγονότος του ότι δημιουργούμε B δείγματα bootstrap αντί να θεωρούμε όλα τα δυνατά υποσύνολα μεγέθους n από τις παρατηρηθείσες τιμές (Efron & Tibshirani, 1993). Παρακάτω, θα παρουσιάσουμε την μέθοδο jackknife after bootstrap η οποία είναι μια απλή μέθοδος για να υπολογίσουμε την μεταβλητότητα των εκτιμητών bootstrap. Ας υποθέσουμε ότι από B δείγματα bootstrap έχουμε υπολογίσει το τυπικό σφάλμα bootstrap του $\hat{\theta}$, το \widehat{se}_{boot} . Θέλουμε να έχουμε ένα μέτρο της αβεβαιότητας του \widehat{se}_{boot} . Το jackknife after bootstrap είναι μια μέθοδος που μας επιτρέπει να

εκτιμήσουμε το $Var(\widehat{se}_{boot})$ χρησιμοποιώντας τα B δείγματα bootstrap που έχουμε στη διάθεση μας. Ας υποθέσουμε τώρα ότι έχουμε έναν δυνατό υπολογιστή για να υπολογίσουμε την jackknife εκτίμηση της διακύμανσης του \widehat{se}_{boot} , τότε θα ακολουθούσαμε τα παρακάτω βήματα.

- Για $i = 1, 2, \dots, n$, αφήνουμε μια παρατήρηση εκτός από το αρχικό δείγμα και επαναυπολογίζουμε το \widehat{se}_{boot} . Το αποτέλεσμα θα είναι το $\widehat{se}_{boot(i)}$ για κάθε i .
- Υπολογίζουμε το $\widehat{Var}_{jack}(\widehat{se}_{boot}) = [(n-1)/n] \sum_{i=1}^n (\widehat{se}_{boot(i)} - \widehat{se}_{boot(\cdot)})^2$
Όπου $\widehat{se}_{boot(\cdot)} = \sum_{i=1}^n \widehat{se}_{boot(i)} / n$

Όμως, εδώ η δυσκολία που προκύπτει είναι στον υπολογισμό του $\widehat{se}_{boot(i)}$ που απαιτεί την χρήση εντελώς διαφορετικών δειγμάτων bootstrap για κάθε i . Για να αντιμετωπίσουμε αυτό το πρόβλημα μπορούμε να χρησιμοποιήσουμε την παρακάτω λογική. Για κάθε παρατήρηση i του αρχικού δείγματος, που θέλουμε να αφήσουμε εκτός, υπάρχουν κάποια δείγματα bootstrap που έχουν προκύψει κατά την διαδικασία του bootstrap για το \widehat{se}_{boot} στα οποία η συγκεκριμένη παρατήρηση δεν εμφανίζεται. Μπορούμε να χρησιμοποιήσουμε αυτά τα δείγματα ώστε να εκτιμήσουμε το $\widehat{se}_{boot(i)}$. Συγκεκριμένα, εκτιμάμε το $\widehat{se}_{boot(i)}$ από την τυπική απόκλιση των τιμών $\hat{\vartheta}^*(b)$ στα δείγματα \mathbf{X}^{*b} που δεν περιέχουν την παρατήρηση i του αρχικού δείγματος. Αν πούμε ότι C_i υποδηλώνει τους δείκτες των δειγμάτων bootstrap που δεν περιέχουν την παρατήρηση i , και ότι υπάρχουν B_i τέτοια δείγματα, τότε

$$\widehat{se}_{boot(i)} = \left[\sum_{b \in C_i} (\hat{\vartheta}^*(b) - \bar{\hat{\vartheta}}_i^*)^2 / B_i \right]^{1/2}$$

όπου $\bar{\hat{\vartheta}}_i^* = \sum_{b \in C_i} \hat{\vartheta}^*(b) / B_i$.

Σημειώνεται ότι το jackknife after bootstrap μπορεί να χρησιμοποιηθεί και για άλλα στατιστικά bootstrap πέραν του τυπικού σφάλματος. Τέλος, στην συγκεκριμένη μέθοδο δημιουργείται πρόβλημα όταν η παρατήρηση i εμφανίζεται σε όλα τα δείγματα bootstrap. Ωστόσο, σύμφωνα με τους Efron and Tibshirani (1993) αυτό είναι εξαιρετικά σπάνιο αν $n \geq 10$ και $B \geq 20$.

1.8 Παραμετρικό Bootstrap

Στην μη παραμετρική μέθοδο bootstrap δεν κάναμε καμία υπόθεση για την κατανομή του πληθυσμού και ο αλγόριθμος bootstrap έκανε δειγματοληψία με επανάθεση B δειγμάτων

bootstrap από τα δεδομένα, δηλαδή από το δείγμα. Στην παραμετρική περίπτωση προσομοιώνουμε B δείγματα μεγέθους n από την \hat{F}_{par} η οποία είναι η παραμετρική εκτίμηση της F .

$$\hat{F}_{par} \rightarrow (x_1^*, x_2^*, \dots, x_n^*)$$

Αφού προσομοιώσουμε B τέτοια δείγματα bootstrap μεγέθους n από την \hat{F}_{par} προχωράμε όπως την μη παραμετρική περίπτωση, δηλαδή υπολογίζουμε τις τιμές $\hat{\theta}_1^*, \hat{\theta}_2^*, \dots, \hat{\theta}_B^*$ του εκτιμητή $\hat{\theta}$ σε κάθε δείγμα bootstrap και για την εκτίμηση bootstrap του τυπικού σφάλματος παίρνουμε την τετραγωνική ρίζα της διασποράς αυτών των τιμών.

1.9 Αποτυχία μεθόδου Bootstrap

Παρακάτω θα αναφέρουμε κάποιες περιπτώσεις που η μέθοδος bootstrap αποτυγχάνει. Αξίζει να σημειωθεί ότι σε πολλά από τα προβλήματα που η μέθοδος bootstrap αποτυγχάνει, έχουν βρεθεί τρόποι αντιμετώπισης. Μια περίπτωση αποτυχίας της μεθόδου bootstrap μπορεί να προκύψει αν προσπαθήσουμε να εκτιμήσουμε ποσότητες οι οποίες δεν έχουν ροπές. Για παράδειγμα, αν τα δεδομένα μας προέρχονται από την κατανομή Cauchy που γνωρίζουμε ότι η αναμενόμενη της τιμή δεν υπάρχει, χρησιμοποιώντας την μέθοδο bootstrap για την εκτίμηση του τυπικού σφάλματος της δειγματικής μέσης τιμής θα οδηγηθούμε σε λανθασμένα αποτελέσματα. Σε αυτήν την περίπτωση η μέθοδος bootstrap θα συμπεριφερθεί απρόβλεπτα ακόμα και αν το μέγεθος του δείγματος είναι τεράστιο. Επιπρόσθετα, η μέθοδος bootstrap αποτυγχάνει όταν ενδιαφερόμαστε για ακραίες τιμές. Αν π.χ. τα δεδομένα μας προέρχονται από ομοιόμορφη κατανομή στο $(0, \vartheta)$, τότε η εκτίμηση μέγιστης πιθανοφάνειας του ϑ είναι η μέγιστη δειγματική τιμή, $t(\mathbf{X}) = \max\{x_1, x_2, \dots, x_n\}$. Σε αυτήν την περίπτωση αν λάβουμε δείγματα bootstrap από τα δεδομένα με σκοπό την εκτίμηση της διασποράς της μέγιστης τιμής, σε κάθε δείγμα bootstrap τα αποτελέσματα δεν θα είναι ικανοποιητικά καθώς θα χρησιμοποιήσουμε στατιστική συνάρτηση που αξιοποιεί μόνο την μεγαλύτερη παρατήρηση ενός δείγματος και η πιθανότητα η μέγιστη παρατήρηση του αρχικού δείγματος μεγέθους n να βρίσκεται εντός του δείγματος bootstrap, καθώς $n \rightarrow \infty$, είναι $1 - (1 - n^{-1})^n \rightarrow 1 - e^{-1} \approx 0.632$ (Davison & Hinkley, 1997). Για την αντιμετώπιση αυτού του προβλήματος μπορούμε να χρησιμοποιήσουμε παραμετρικό bootstrap δηλαδή να προσομοιώσουμε δείγματα bootstrap από την ομοιόμορφη $(0, \hat{\vartheta})$. Επιπλέον, πρόβλημα της μεθόδου αναδύεται αν έχουμε δείγμα από κατανομή με πολύ βαριές ουρές, και σαν αποτέλεσμα $Var(X) = \infty$. Τότε η μέθοδος bootstrap για τον μέσο αποτυγχάνει όπως πρωτοειπώθηκε από τον Babu (1984) και αργότερα από τον Athreya (1987) και Knight (1989). Σύμφωνα με τον Politis (1998) ο λόγος που ίσως η μέθοδος

bootstrap μπορεί να αποτύχει είναι ότι τα δείγματα bootstrap δεν προέρχονται ακριβώς από την κατανομή F αλλά από την εμπειρική \hat{F} και παρ' όλο που είναι ένας συνεπής εκτιμητής της πραγματικής κατανομής F ίσως χάνει κάποια χαρακτηριστικά τα οποία επηρεάζουν σημαντικά την κατανομή του στατιστικού. Να σημειωθεί ότι η ιδιότητα της συνέπειας είναι ασυμπτωτική, δηλαδή είναι για μέγεθος δείγματος $n \rightarrow \infty$.

Επιπρόσθετα, πρόβλημα της μεθόδου bootstrap δημιουργείται όταν το δείγμα μας είναι μικρό. Πιο συγκεκριμένα, η κύρια ανησυχία είναι ότι από ένα δείγμα λίγων παρατηρήσεων τα δείγματα bootstrap θα υποεκτιμούν την πραγματική μεταβλητότητα του $\hat{\theta}$ επειδή οι παρατηρήσεις σε ένα δείγμα bootstrap θα επαναλαμβάνονται και ακόμα και το ίδιο το δείγμα bootstrap μπορεί να επαναληφθεί στις B επαναλήψεις της μεθόδου (Chernick, 2007). Αξίζει να σημειωθεί πως ο αριθμός όλων των δυνατών διαφορετικών δειγμάτων μεγέθους n με επανάθεση από το αρχικό δείγμα σύμφωνα με τον Hall (1992) μπορεί να δοθεί από τον παρακάτω τύπο:

$$\binom{2n-1}{n} = (2n-1)!/[n!(n-1)!]$$

Αυτός ο αριθμός όλων των δυνατών διαφορετικών δειγμάτων bootstrap είναι μεγάλος ακόμα και όταν έχουμε μικρό αρχικό δείγμα, δηλαδή η πιθανότητα να εμφανιστεί ξανά ένα δείγμα bootstrap, αφού έχει εμφανιστεί, είναι πολύ μικρή. Όταν π.χ. έχουμε μέγεθος δείγματος $n = 20$ και ο αριθμός των επαναλήψεων bootstrap είναι $B = 2000$ τότε, σύμφωνα με τον Hall (1992), η πιθανότητα να μην επανεμφανιστεί ένα δείγμα bootstrap δεν είναι μικρότερη από 0.954. Επειδή ο αριθμός των διαφορετικών δειγμάτων bootstrap αυξάνει δραματικά καθώς αυξάνει το μέγεθος του αρχικού δείγματος n , ο ακριβής υπολογισμός του εκτιμητή bootstrap είναι συνήθως αδύνατο να γίνει και η χρήση της μεθόδου Monte Carlo κατά τη διαδικασία του bootstrap κρίνεται απαραίτητη. Σε περίπτωση που το μέγεθος του δείγματος είναι $n \leq 8$, τότε σύμφωνα με τον Fisher and Hall (1990) μπορεί να πραγματοποιηθεί ακριβής υπολογισμός του εκτιμητή bootstrap. Στην περίπτωση ενός πολύ μικρού δείγματος η εμπειρική συνάρτηση κατανομής συνήθως δεν είναι καλή εκτίμηση της κατανομής του πληθυσμού. Δείγματα μικρότερα των 10 παρατηρήσεων ακόμα και σε παραμετρικές περιπτώσεις είναι πολύ μικρά για να παρέχουν εκτιμήσεις για ποσότητες που μας ενδιαφέρουν, πόσο μάλλον όταν τα δείγματα bootstrap είναι μικρότερα των 10 παρατηρήσεων. Σε μη παραμετρικά προβλήματα απαιτούνται μεγαλύτερα μεγέθη δειγμάτων από τις παραμετρικές περιπτώσεις. Σε πολλά πρακτικά προβλήματα το ελάχιστο μέγεθος δείγματος είναι 30, ενώ ο Chernick (2007) προτείνει ότι στις περισσότερες περιπτώσεις θα ήταν καλό να έχουμε $n \geq 50$. Ακόμα, σύμφωνα με τον ίδιο, ο καλύτερος κανόνας όπως πρότεινε αρχικά για την προσέγγιση Monte Carlo κατά τη διαδικασία του bootstrap είναι $B = 100$ (τουλάχιστον) για εκτίμηση bootstrap τυπικού σφάλματος και μεροληψίας εκτιμητή και $B = 1000$ για διαστήματα εμπιστοσύνης μέσω της

μεθόδου bootstrap, όμως στη συνέχεια αναφέρει πως επειδή αυξάνεται η ταχύτητα των υπολογιστών με την πάροδο του χρόνου θα μπορούσαμε να πάρουμε $B = 1000$ για εκτίμηση bootstrap τυπικού σφάλματος και μεροληψίας και $B = 10000$ για διαστήματα εμπιστοσύνης. Τέλος, η μέθοδος bootstrap αποτυγχάνει για εξαρτημένα δεδομένα. Σύμφωνα με τον Singh (1981) υπάρχει ασυνέπεια του εκτιμητή bootstrap όταν τα δεδομένα έχουν κάποιο βαθμό εξάρτησης και δεν είναι iid.

1.10 M-out-of-n Bootstrap

Ο Athreya (1987) επισήμανε ότι οι δυσκολίες της μεθόδου bootstrap όταν τα δεδομένα μας προέρχονται από κάποια κατανομή με πολύ βαριές ουρές μπορούν να ξεπεραστούν χρησιμοποιώντας μικρότερο αριθμό παρατηρήσεων στα δείγματα bootstrap. Η μέθοδος m-out-of-n bootstrap ή αλλιώς moon bootstrap, είναι μια απλή παραλλαγή της αρχικής μη παραμετρικής μεθόδου του Efron (1979) που λύνει όμως πολλά προβλήματα ασυνέπειας της μεθόδου bootstrap που δημιουργούνται σε αρκετές περιπτώσεις. Όταν το μη παραμετρικό bootstrap πρώτοεισήχθη από τον Efron, προτάθηκε το μέγεθος των δειγμάτων bootstrap να είναι ίδιο με το μέγεθος του αρχικού δείγματος n . Ωστόσο, παρ' όλο που αυτό δουλεύει σε αρκετές περιπτώσεις, αποδείχθηκε στην πορεία των χρόνων ότι η μέθοδος έχει καλά αποτελέσματα όταν επιλεγεί μέγεθος δειγμάτων bootstrap $m < n$. Η μέθοδος αυτή έχει προταθεί από πολλούς όπως οι Bickel et al. (1997) και φαίνεται ότι δίνει συνεπείς εκτιμήσεις σε περιπτώσεις που η κλασική μη παραμετρική μέθοδος bootstrap αποτυγχάνει. Συνήθως, η ασυμπτωτική θεωρία απαιτεί το $m \rightarrow \infty$ όσο το $n \rightarrow \infty$, αλλά σε διαφορετικό ρυθμό, πιο αργό τέτοιο ώστε το $m/n \rightarrow 0$. Αξίζει να σημειωθεί ότι αυτή η μέθοδος δουλεύει σε περιπτώσεις και ανεξάρτητων αλλά και εξαρτημένων παρατηρήσεων.

1.11 Διάστημα εμπιστοσύνης Percentile Bootstrap

Ένα διάστημα εμπιστοσύνης είναι ένα τυχαίο διάστημα $[T_1, T_2]$ με άκρα στατιστικές συναρτήσεις, δηλαδή τυχαίες μεταβλητές οι τιμές των οποίων μπορούν να υπολογιστούν από τα δεδομένα.

Το διάστημα $[T_1, T_2]$ είναι ένα $100(1 - \alpha)\%$ ΔΕ για κάποια άγνωστη ποσότητα ϑ αν

$$P(T_1 \leq \vartheta \leq T_2) \geq 1 - \alpha$$

Ιδανικά θα επιθυμούσαμε η παραπάνω πιθανότητα, η οποία καλείται πιθανότητα κάλυψης του διαστήματος εμπιστοσύνης, να είναι ακριβώς ίση με τον ονομαστικό συντελεστή εμπιστοσύνης του διαστήματος $1 - \alpha$. Στην περίπτωση αυτή λέμε ότι το ΔΕ είναι ακριβές.

Ένα ΔΕ για την ποσότητα ϑ είναι μια εναλλακτική μέθοδος εκτίμησής της. Αντί να την εκτιμήσουμε σημειακά με μια μοναδική τιμή, θεωρούμε για την ϑ ένα εύρος τιμών οι οποίες είναι συμβατές με τα παρατηρηθέντα δεδομένα. Τα άκρα του ΔΕ λέγονται όρια εμπιστοσύνης για την ποσότητα ϑ , κάτω και άνω, αντίστοιχα. Όταν υπολογίσουμε τα δύο όρια εμπιστοσύνης από τα δεδομένα μας, παίρνουμε ένα εύρος τιμών και τα άκρα του διαστήματος συνεχίζουμε να τα λέμε κάτω και άνω όρια εμπιστοσύνης. Είναι σύνηθες να κατασκευάζουμε διαστήματα εμπιστοσύνης τα οποία να αφήνουν και στις δύο ουρές την ίδια πιθανότητα $\alpha/2$, δηλαδή $P(T_1 > \vartheta) = P(T_2 < \vartheta) = \alpha/2$. Σε αυτήν την περίπτωση τα διαστήματα εμπιστοσύνης λέγονται ΔΕ ίσων ουρών. Ωστόσο, τα ΔΕ ίσων ουρών δεν είναι τα μόνα που μπορούν να κατασκευαστούν αλλά είναι αυτά με το μικρότερο μήκος αν η κατανομή είναι συμμετρική. Τέλος, ένα σύνηθες κριτήριο αξιολόγησης διαφορετικών $100(1 - \alpha)\%$ ΔΕ για την ίδια ποσότητα, πέραν φυσικά της διατήρησης του ονομαστικού συντελεστή εμπιστοσύνης $1 - \alpha$ είναι το αναμενόμενο μήκος των διαστημάτων. Προτιμάμε ΔΕ με μικρό αναμενόμενο μήκος.

Αν υποθέσουμε ότι ο εκτιμητής $\hat{\vartheta}$ της παραμέτρου ϑ κατανέμεται κανονικά,

$$\hat{\vartheta} \sim N(\vartheta, se(\hat{\vartheta})) \Rightarrow z \equiv \frac{\hat{\vartheta} - \vartheta}{se(\hat{\vartheta})} \sim N(0,1)$$

Για την τυπική κανονική κατανομή μπορούμε εύκολα να ορίσουμε ένα διάστημα $[z_{\alpha/2}, z_{1-\alpha/2}]$ στο οποίο θα ανήκει η z με κάποια δοθείσα πιθανότητα $1 - \alpha$. Ας σημειωθεί ότι το z_α είναι το α -ποσοστιαίο σημείο της τυπικής κανονικής κατανομής.

$$\Phi(z_{\alpha/2}) = P(z < z_{\alpha/2}) = \alpha/2$$

$$\Phi(z_{1-\alpha/2}) = P(z < z_{1-\alpha/2}) = 1 - \alpha/2$$

όπου $\Phi(z)$ είναι η αθροιστική συνάρτηση της τυπικής κανονικής κατανομής. Συνεπώς, η πιθανότητα $z < z_{\alpha/2}$ ή $z > z_{1-\alpha/2}$ είναι ίση με α . Αντίστοιχα,

$$P(z_{\alpha/2} \leq z \leq z_{1-\alpha/2}) = \Phi(z_{1-\alpha/2}) - \Phi(z_{\alpha/2}) = 1 - \alpha$$

Λόγω συμμετρίας της συνάρτησης πυκνότητας της τυπικής κανονικής κατανομής ως προς το μηδέν ισχύει ότι $z_{\alpha/2} = -z_{1-\alpha/2}$ και μπορεί να χρησιμοποιηθεί η μια από τις δύο κρίσιμες τιμές.

Τώρα θα μετασχηματίσουμε το διάστημα $[z_{\alpha/2}, z_{1-\alpha/2}]$ για πιθανότητα $1 - \alpha$ στο αντίστοιχο διάστημα που περιέχει την παράμετρο ϑ . Λύνουμε τις παρακάτω σχέσεις ως προς ϑ

$$z_{\alpha/2} = \frac{\hat{\vartheta} - \vartheta}{se(\hat{\vartheta})}, \quad z_{1-\alpha/2} = \frac{\hat{\vartheta} - \vartheta}{se(\hat{\vartheta})}$$

και προκύπτει ότι το διάστημα εμπιστοσύνης για το ϑ είναι

$$[\hat{\vartheta} - z_{\alpha/2}se(\hat{\vartheta}), \hat{\vartheta} + z_{\alpha/2}se(\hat{\vartheta})]$$

όπου το παραπάνω διάστημα καλείται ΔΕ ίσων ουρών.

Στην πράξη όμως η υπόθεση ότι η κατανομή του εκτιμητή $\hat{\vartheta}$ είναι κανονική μπορεί πολλές φορές να μην ισχύει ή ενδέχεται το τυπικό σφάλμα να είναι άγνωστο, σε αυτές τις περιπτώσεις καταφεύγουμε στην μέθοδο bootstrap. Για σκοπούς της εργασίας θα δούμε μόνο το διάστημα Percentile Bootstrap. Το διάστημα εμπιστοσύνης Percentile Bootstrap είναι το διάστημα που χρησιμοποιεί τα $1 - \alpha/2$, $\alpha/2$ ποσοστιαία σημεία της κατανομής των $\hat{\vartheta}^*$ που έχει προκύψει από το δείγμα bootstrap $\hat{\vartheta}_1^*, \hat{\vartheta}_2^*, \dots, \hat{\vartheta}_B^*$. Δηλαδή ένα $100(1 - \alpha)\%$ percentile διάστημα είναι το

$$[\hat{\vartheta}_{(\alpha/2)}^*, \hat{\vartheta}_{(1-\alpha/2)}^*]$$

Πρόκειται για ένα διάστημα απλό στον υπολογισμό του που προσεγγίζει τον ονομαστικό συντελεστή εμπιστοσύνης με ταχύτητα ανάλογη του \sqrt{n} . Αν χρησιμοποιήσουμε επίπεδο σημαντικότητας $\alpha = 0.10$ και επαναλήψεις bootstrap $B = 1000$, τότε η 50^η και 950^η διατεταγμένη παρατήρηση του δείγματος bootstrap $\hat{\vartheta}_1^*, \hat{\vartheta}_2^*, \dots, \hat{\vartheta}_{1000}^*$ σχηματίζουν το διάστημα Percentile Bootstrap.

Κεφάλαιο 2

2.1 Δεδομένα χρονοσειρών

Η κλασική στατιστική ασχολείται με παρατηρήσεις iid αλλά στα δεδομένα χρονοσειρών αυτή η υπόθεση συνήθως είναι ασταθής. Μια χρονοσειρά ή χρονολογική σειρά, είναι μια ακολουθία διαδοχικών παρατηρήσεων στη διάρκεια του χρόνου που συνήθως παρουσιάζει κάποιο βαθμό εξάρτησης. Το διάστημα μεταξύ των διαδοχικών μετρήσεων στο χρόνο είναι περίπου ίδιο. Δεδομένα που είναι ανεξάρτητα αλλά όχι ισόνομα μπορούν μερικές φορές να αναλυθούν μέσω γραμμικών μοντέλων και από μεθόδους της ανάλυσης παλινδρόμησης. Αν τα δεδομένα είναι εξαρτημένα αλλά ισόνομα τότε πιθανός μπορούν να αναλυθούν χρησιμοποιώντας τεχνικές στάσιμων χρονοσειρών. Σε γενικότερες γραμμές, τα δεδομένα χρονοσειρών μπορεί να μην είναι ούτε ανεξάρτητα και ούτε ισόνομα και συνεπώς η ανάλυση τους να απαιτεί συγκεκριμένες τεχνικές που βασίζονται στην δομή των δεδομένων (Politis & McElroy, 2020). Σε μια χρονοσειρά λόγω της συγκεκριμένης ροής και κατεύθυνσης του χρόνου ο δείκτης i μιας τυχαίας μεταβλητής X_i θα έχει μια φυσική διάταξη, και για να αντικατοπτρίσουμε αυτήν την δομή χρησιμοποιούμε συνήθως τον συμβολισμό X_t όπου t υποδηλώνει χρόνο. Μια χρονοσειρά είναι μια στοχαστική διαδικασία καθώς η τιμή κάθε χρονικής στιγμής συνιστά και μια ξεχωριστή τυχαία μεταβλητή, ενώ οι τιμές της επηρεάζονται από τυχαίους παράγοντες. Μια συλλογή τυχαίων μεταβλητών στη διάρκεια του χρόνου ονομάζεται στοχαστική διαδικασία και συμβολίζεται $\{X_t\}$, η οποία περιέχει άπειρο αριθμό τυχαίων μεταβλητών. Αξίζει να σημειωθεί ότι το σύνολο των παρατηρήσεων μιας χρονοσειράς αποτελεί μια πραγματοποίηση της στοχαστικής διαδικασίας. Σε μια χρονοσειρά έχουμε x_1, x_2, \dots, x_n δηλαδή n διαδοχικές παρατηρήσεις στα διακριτά χρονικά σημεία $t = 1, 2, \dots, n$, ενώ είναι αρκετά σημαντικό να αναφέρουμε ότι η σειρά των παρατηρήσεων μιας χρονοσειράς έχει σημασία. Οι παρατηρήσεις x_1, x_2, \dots, x_n είναι οι τιμές των τυχαίων μεταβλητών X_1, X_2, \dots, X_n δηλαδή η πραγματοποίηση της στοχαστικής διαδικασίας.

2.2 Στασιμότητα

Τα μοντέλα χρονοσειρών συχνά βασίζονται στην υπόθεση της στασιμότητας. Σε μια στάσιμη χρονοσειρά οι στατιστικές της ιδιότητες παραμένουν σταθερές στο χρόνο (Montgomery et al., 2008). Ειδικότερα, σε μια στάσιμη χρονοσειρά δεν παρατηρείται συστηματική αλλαγή του μέσου όρου και της διασποράς στο χρόνο. Η στασιμότητα

υποδηλώνει έναν τύπο στατιστικής ισορροπίας στα δεδομένα. Μια χρονοσειρά λέγεται αυστηρά στάσιμη αν η από κοινού κατανομή των X_{t_1}, \dots, X_{t_k} είναι ίδια με την από κοινού κατανομή των $X_{t_1+\tau}, \dots, X_{t_k+\tau}$ για όλα τα t_1, \dots, t_k, τ . Με άλλα λόγια αν μετακινούσαμε την χρονοσειρά κατά μια χρονική περίοδο τ τότε αυτό δεν θα είχε επίπτωση στις από κοινού κατανομές οι οποίες εξαρτώνται μόνο από την απόσταση μεταξύ t_1, \dots, t_k . Συνεπώς η χρονοσειρά θα φαίνεται περίπου ίδια σε διαφορετικά διαστήματα ιδού μήκους.

Ο μέσος, ή πρώτη ροπή, μιας χρονοσειράς μπορεί να υπολογιστεί για κάθε χρονική στιγμή t δίνοντας μας την ακολουθία των μέσων. Το ίδιο ισχύει και για την συνδιακύμανση μιας χρονοσειράς. Η αυτοσυνδιακύμανση είναι η συνδιακύμανση της χρονοσειράς με την ίδια σε διαφορετική χρονική περίοδο. Ας σημειωθεί ότι υστέρηση τ ορίζεται ως η απόλυτη τιμή της χρονικής διαφοράς $|(t + \tau) - t|$ δύο όρων $X_t, X_{t+\tau}$. Η συνδιακύμανση μεταξύ X_t και $X_{t+\tau}$ ονομάζεται αυτοσυνδιακύμανση υστέρησης τ .

$$\gamma(\tau) = Cov[X_t, X_{t+\tau}] = E\{[X_t - \mu][X_{t+\tau} - \mu]\}$$

Αξίζει να σημειωθεί ότι η αυτοσυνδιακύμανση για $\tau = 0$ είναι απλά η διακύμανση της χρονοσειράς. Δηλαδή, $\gamma(0) = \sigma_x^2$ που είναι σταθερή για στάσιμες χρονοσειρές. Επιπρόσθετα, $\gamma(\tau) = \gamma(-\tau)$, για κάθε τ , κάτι που σημαίνει ότι η αυτοσυνδιακύμανση είναι ταυτόσημη για μια θετική χρονική μετατόπιση και την ίδια αρνητική χρονική μετατόπιση. Η συλλογή των τιμών της $\gamma(\tau)$, $\tau = 0, 1, 2, \dots$ ονομάζεται συνάρτηση αυτοσυνδιακύμανσης.

Ο συντελεστής αυτοσυσχέτισης υστέρησης τ για στάσιμες χρονοσειρές δίνεται από τον παρακάτω τύπο:

$$\rho(\tau) = \frac{E[(x_t - \mu)(x_{t+\tau} - \mu)]}{\sqrt{E[(x_t - \mu)^2]E[(x_{t+\tau} - \mu)^2]}} = \frac{Cov(x_t, x_{t+\tau})}{Var(x_t)} = \frac{\gamma(\tau)}{\gamma(0)}$$

Η συλλογή των τιμών του $\rho(\tau)$, $\tau = 0, 1, 2, \dots$ ονομάζεται συνάρτηση αυτοσυσχέτισης (ACF). Όπως μπορούμε να καταλάβουμε $\rho(0) = 1$ αφού είναι η συσχέτιση μεταξύ X_t και $X_{t+\tau}$ για $\tau = 0$. Επίσης, $\rho(\tau) = \rho(-\tau)$ για κάθε τ και $-1 \leq \rho(\tau) \leq 1$. Σε αντίθεση με τον συντελεστή αυτοσυνδιακύμανσης που εξαρτάται από τις μονάδες μέτρησης η συνάρτηση αυτοσυσχέτισης δεν εξαρτάται από μονάδες μέτρησης. Η αυστηρή στασιμότητα ορίζεται μαθηματικά ως η διατήρηση στο χρόνο t της κοινής κατανομής των $X_{(t_1+\tau)}, \dots, X_{(t_k+\tau)}$. Η συνθήκη στασιμότητας περιορίζεται συνήθως στη διατήρηση της μέσης τιμής και αυτοσυνδιακύμανσης και αναφέρεται ως ασθενής στασιμότητα. Μια χρονοσειρά αν έχει πεπερασμένο μέσο και συνάρτηση αυτοσυνδιακύμανσης τότε έχει ασθενής στασιμότητα ή δεύτερης τάξης στασιμότητα. Ωστόσο, στην πραγματικότητα αντί της αυτοσυνδιακύμανσης εξετάζουμε μόνο την σταθερότητα της διακύμανσης με το χρόνο. Η αυτοσυνδιακύμανση

πρέπει να εξαρτάται μόνο από το χρονικό διάστημα μεταξύ των δυο χρονοσειρών και όχι από την ίδια τη χρονική στιγμή.

$$E[X_t] = \mu$$

και

$$\text{Cov}[X_t, X_{t+\tau}] = \gamma(\tau)$$

Το πρώτο βήμα στην ανάλυση χρονοσειρών είναι να κοιτάξουμε το μοτίβο της συσχέτισης στα διάφορα χρονικά σημεία κάτι το οποίο το καταφέρνουμε από την δειγματική συνάρτηση αυτοσυσχέτισης. Επιπλέον, σημαντικό είναι να αναφέρουμε ότι η αυστηρή στασιμότητα είναι πιο ισχυρή από την ασθενή στασιμότητα αλλά συχνά η τελευταία αρκεί ώστε να θεμελιωθούν ορισμένα βασικά στατιστικά συμπεράσματα. Επίσης, αξίζει να σημειωθεί ότι οι αυστηρά στάσιμες χρονοσειρές με πεπερασμένη διακύμανση είναι ασθενώς στάσιμες χωρίς να ισχύει το αντίστροφο.

Μερικές στοιχειώδεις ιδιότητες μιας αυστηρά στάσιμης χρονοσειράς $\{X_t\}$ σύμφωνα με του Brockwell and Davis (2016) είναι οι παρακάτω

- Οι τυχαίες μεταβλητές X_t είναι ισόνομα κατανεμημένες
- Η από κοινού κατανομή των X_1, \dots, X_n είναι ίδια με την από κοινού κατανομή των $X_{1+\tau}, \dots, X_{n+\tau}$ για όλους τους ακέραιους τ και $n \geq 1$.
- $\{X_t\}$ είναι ασθενώς στάσιμη αν $E(X_t^2) < \infty$ για όλα τα t
- Η ασθενής στασιμότητα δεν συνεπάγεται σε αυστηρή στασιμότητα
- Μια iid ακολουθία είναι αυστηρά στάσιμη

2.3 Λευκός Θόρυβος

Μια διακριτού χρόνου διαδικασία που αποτελείται από μια ακολουθία τυχαίων μεταβλητών που είναι αμοιβαία ανεξάρτητες και ισόνομες, iid τυχαίες μεταβλητές, ονομάζεται λευκός θόρυβος και είναι η πιο απλή στάσιμη χρονοσειρά. Μια χρονοσειρά θεωρείται ότι είναι λευκός θόρυβος αν αποτελείται από τυχαίες μεταβλητές που όλες έχουν μηδενικό μέσο, την ίδια διακύμανση, και είναι ασυσχέτιστες μεταξύ τους (Politis & McElroy, 2020). Γενικότερα, αν μια χρονοσειρά περιέχει μη συσχετισμένες παρατηρήσεις και έχει σταθερή διακύμανση λέμε ότι είναι λευκός θόρυβος ενώ αν αυτές οι παρατηρήσεις είναι κανονικά κατανεμημένες τότε λέμε ότι είναι γκαουσιανός λευκός θόρυβος. Αξίζει να σημειωθεί ότι μια στοχαστική διαδικασία λέγεται γκαουσιανή αν για κάθε πεπερασμένο $m \geq 1$ η m -διάστατη περιθώρια κατανομή ή αλλιώς η από κοινού κατανομή είναι πολυδιάστατη κανονική. Μια υπόθεση που

πραγματοποιούμε συχνά είναι ότι ο λευκός θόρυβος είναι γκαουσιανός. Μερικοί συγγραφείς προτιμούν να κάνουν την ασθενέστερη παραδοχή ότι οι τυχαίες μεταβλητές είναι αμοιβαία ασυσχέτιστες παρά ανεξάρτητες. Αυτό είναι επαρκές για γραμμικές κανονικές διαδικασίες, αλλά η ισχυρότερη υπόθεση της ανεξαρτησίας χρειάζεται όταν μελετώνται μη γραμμικά μοντέλα. Σύμφωνα με τους Chatfield and Xing (2019) όταν ο λευκός θόρυβος ορίζεται ως μια ακολουθία iid τότε μερικές φορές καλείται και ως αυστηρώς λευκός θόρυβος (strict white noise), ενώ όταν οι διαδοχικές τιμές είναι απλώς ασυσχέτιστες και όχι ανεξάρτητες καλείται ασυσχέτιστος λευκός θόρυβος (uncorrelated white noise)

Μια χρονοσειρά $\{Z_t\}$ με μηδενικό μέσο, μηδενική σειριακή συσχέτιση και ίσες διακυμάνσεις ονομάζεται λευκός θόρυβος:

$$Z_t \sim WN(0, \sigma_z^2)$$

Επίσης,

$$\gamma(k) = \text{Cov}(Z_t, Z_{t+k}) = \begin{cases} \sigma_z^2 & k = 0 \\ 0 & k = \pm 1, \pm 2, \dots \end{cases}$$

Αυτό σημαίνει ότι οι διαφορετικές τιμές είναι ασυσχέτιστες δηλαδή:

$$\rho(k) = \begin{cases} 1 & k = 0 \\ 0 & k = \pm 1, \pm 2, \dots \end{cases}$$

2.4 Τυχαίος περίπατος

Υποθέτοντας ότι $\{Z_t\}$ είναι λευκός θόρυβος, τότε μια διαδικασία $\{X_t\}$ ονομάζεται τυχαίος περίπατος αν:

$$X_t = X_{t-1} + Z_t$$

Ξεκινώντας από κάποια τιμή X_0 για $t = 0$ και αντικαθιστώντας επαναληπτικά έως τον χρόνο t τον παραπάνω ορισμό του τυχαίου περιπάτου τότε έχω το άθροισμα όλων των τυχαίων βημάτων ως τη στιγμή t :

$$X_t = \sum_{i=0}^t Z_i$$

Μπορούμε να διαπιστώσουμε ότι $E(X_t) = 0$ και $\text{Var}(X_t) = t\sigma_z^2$, δηλαδή η διασπορά του τυχαίου περιπάτου είναι ανάλογη του χρόνου κάτι το οποίο μας οδηγεί στο ότι η διαδικασία

δεν είναι στάσιμη. Αξίζει να σημειωθεί ότι οι πρώτες διαφορές ενός τυχαίου περιπάτου σχηματίζουν λευκό θόρυβο, δηλαδή μια στάσιμη διαδικασία (Chatfield & Xing, 2019).

$$\nabla X_t = X_t - X_{t-1} = Z_t$$

2.5 Γραμμικές στοχαστικές διαδικασίες

Μια γραμμική στοχαστική διαδικασία ή αλλιώς γραμμική χρονοσειρά ορίζεται για κάθε χρονική στιγμή t ως ένα άθροισμα ασυσχέτιστων τυχαίων μεταβλητών.

$$X_t = \mu + \sum_{i=-\infty}^{\infty} \psi_i Z_{t-i}, \quad Z_t \sim WN(0, \sigma_Z^2)$$

όπου οι συντελεστές ικανοποιούν

$$\sum_{i=-\infty}^{\infty} |\psi_i| < \infty$$

Για ευκολία θέτουμε $\psi_0 = 1$ και $\mu = 0$. Συνεπώς, θεωρώντας τον τελεστή υστέρησης B η γραμμική χρονοσειρά έχει την παρακάτω συμπαγή έκφραση

$$X_t = \psi(B)Z_t = \sum_{i=-\infty}^{\infty} \psi_i B^i Z_t$$

με

$$\psi(B) = \sum_{i=-\infty}^{\infty} \psi_i B^i$$

όπου ο τελεστής του πολωνύμου $\psi(B)$ μπορεί να θεωρηθεί ως γραμμικό φίλτρο το οποίο όταν εφαρμοστεί ως "είσοδο" στην σειρά του λευκού θόρυβου $\{Z_t\}$, ως "έξοδο" παράγει τη γραμμική χρονοσειρά $\{X_t\}$. Αξίζει να σημειωθεί ότι για τα γραμμικά φίλτρα όταν η "είσοδος" είναι οποιαδήποτε στάσιμη χρονοσειρά τότε και η "έξοδος" είναι στάσιμη χρονοσειρά (Brockwell & Davis, 2016).

Η $MA(\infty)$ είναι μια γραμμική διαδικασία με $\psi_i = 0$ για $i < 0$. Η έκφραση της γραμμικής χρονοσειράς που μπορούμε να δούμε παρακάτω είναι η διαδικασία του κινητού μέσου άπειρης τάξης $MA(\infty)$

$$X_t = Z_t + \psi_1 Z_{t-1} + \psi_2 Z_{t-2} + \dots = Z_t + \sum_{i=1}^{\infty} \psi_i Z_{t-i}$$

Εκφράζοντας για κάθε χρονική στιγμή t τον αντίστοιχο όρο της χρονοσειράς ως γραμμικό συνδυασμό των προηγούμενων όρων, παίρνουμε την αυτοπαλίνδρομη διαδικασία άπειρης τάξης $AR(\infty)$.

$$X_t = \pi_1 X_{t-1} + \pi_2 X_{t-2} + \dots + Z_t = \sum_{i=1}^{\infty} \pi_i X_{t-i} + Z_t$$

Η συνθήκη,

$$\sum_{i=0}^{\infty} |\pi_i| < \infty$$

με $\pi_0 = 1$, μας οδηγεί στην ιδιότητα της αντιστρεψιμότητας δηλαδή επιτρέπει η Z_t να μπορεί να εκφραστεί ως άπειρο άθροισμα της παρούσας τυχαίας μεταβλητής και προηγούμενων τυχαίων μεταβλητών της χρονοσειράς. Με τη χρήση πάλι του τελεστή υστέρησης

$$\pi(B)X_t = Z_t, \quad \pi(B) = \sum_{i=0}^{\infty} \pi_i B^i$$

και θεωρώντας ότι ισχύει η αντιστρεψιμότητα

$$X_t = \frac{1}{\pi(B)} Z_t, \quad \psi(B) = \frac{1}{\pi(B)}$$

το οποίο δηλώνει την ισοδυναμία των εκφράσεων $MA(\infty)$ και $AR(\infty)$

2.6 Αυτοπαλίνδρομη χρονοσειρά (AR)

Όπως γνωρίζουμε τα μοντέλα παλινδρόμησης ορίζουν μια εξαρτημένη μεταβλητή ως συνάρτηση κάποιων άλλων ανεξάρτητων ή αλλιώς ερμηνευτικών μεταβλητών. Ειδικότερα, στα γραμμικά μοντέλα παλινδρόμησης η εξαρτημένη μεταβλητή είναι γραμμικός συνδυασμός των ανεξάρτητων μεταβλητών. Αντίστοιχα, η διαδικασία $AR(p)$ ορίζει την τυχαία μεταβλητή X_t ως ένα γραμμικό συνδυασμό των προηγούμενων p τυχαίων μεταβλητών, διαταραγμένη από λευκό θόρυβο. Το p υποδηλώνει την τάξη του μοντέλου. Δηλαδή, στα μοντέλα αυτοπαλινδρόμησης θεωρούμε εξαρτημένη μεταβλητή την τυχαία μεταβλητή της χρονοσειράς σε μια χρονική

στιγμή t και ανεξάρτητες μεταβλητές την τυχαία μεταβλητή της χρονοσειράς σε προηγούμενους χρόνους $t - 1, t - 2, \dots, t - p$.

$$X_t = \varphi_1 X_{t-1} + \varphi_2 X_{t-2} + \dots + \varphi_p X_{t-p} + Z_t, \quad Z_t \sim WN(0, \sigma_z^2)$$

Κάνοντας χρήση του τελεστή υστέρησης B προκύπτει η παρακάτω συμπαγής έκφραση για την $AR(p)$

$$\varphi(B)X_t = Z_t$$

όπου, $\varphi(B) = 1 - \sum_{i=1}^p \varphi_i B^i = 1 - \varphi_1 B - \varphi_2 B^2 - \dots - \varphi_p B^p$ το οποίο είναι το χαρακτηριστικό πολυώνυμο της διαδικασίας $AR(p)$. Επιπρόσθετα, η διαδικασία είναι στάσιμη αν οι ρίζες του χαρακτηριστικού πολυωνύμου είναι εκτός του μοναδιαίου κύκλου. Θεωρώντας το B ως μεταβλητή με τιμές στο σύνολο των μιγαδικών αριθμών, το προηγούμενο σημαίνει ότι οι ρίζες του πολυωνύμου έχουν μέτρο μεγαλύτερο της μονάδας.

Αν η μέση τιμή της X_t δεν είναι μηδέν τότε αντικαθιστούμε το X_t με $X_t - \mu$ και έχουμε:

$$X_t - \mu = \varphi_1 (X_{t-1} - \mu) + \varphi_2 (X_{t-2} - \mu) + \dots + \varphi_p (X_{t-p} - \mu) + Z_t$$

ή διαφορετικά

$$X_t = \alpha + \varphi_1 X_{t-1} + \varphi_2 X_{t-2} + \dots + \varphi_p X_{t-p} + Z_t$$

όπου η σταθερά $\alpha = \mu(1 - \varphi_1 - \dots - \varphi_p)$

2.6.1 Αυτοπαλίνδρομη χρονοσειρά τάξης 1

Η πιο απλή αυτοπαλίνδρομη στοχαστική διαδικασία είναι η $AR(1)$ όπου $p = 1$.

$$X_t = \varphi X_{t-1} + Z_t, \quad Z_t \sim WN(0, \sigma_z^2)$$

με συνθήκη στασιμότητας $|\varphi| < 1$.

Αξίζει να σημειωθεί ότι η $AR(1)$ είναι διαδικασία Markov. (Σημειώνεται ότι και οι χρονοσειρές μεγαλύτερης τάξης μπορούν να μετασχηματιστούν σε μαρκοβιανές διαδικασίες θεωρώντας την ακολουθία των διαδοχικών blocks μήκους ίσου με την τάξη της αντίστοιχης χρονοσειράς). Για $|\varphi| = 1$ η διαδικασία είναι αυτή του τυχαίου περιπάτου.

Με διαδοχική αντικατάσταση στη παραπάνω εξίσωση της $AR(1)$ μπορούμε να δούμε

$$X_t = \varphi(\varphi X_{t-2} + Z_{t-1}) + Z_t$$

$$= \varphi^2(\varphi X_{t-3} + Z_{t-2}) + \varphi Z_{t-1} + Z_t$$

και συνεπώς καταλήγουμε στο ότι το X_t μπορεί να εκφραστεί ως άπειρης τάξης διαδικασία MA

$$X_t = Z_t + \varphi Z_{t-1} + \varphi^2 Z_{t-2} + \dots$$

με την προϋπόθεση ότι $-1 < \varphi < +1$ ώστε το άθροισμα να συγκλίνει.

Η δυνατότητα η διαδικασία AR να γραφτεί στη μορφή της MA και το αντίστροφο δηλώνει την δυική σχέση μεταξύ αυτών των δύο διαδικασιών (duality) κάτι το οποίο είναι χρήσιμο για πολλούς λόγους (Chatfield & Xing, 2019). Είναι καλύτερο αντί της διαδοχικής αντικατάστασης να χρησιμοποιήσουμε τον τελεστή υστέρησης B για να το εξερευνήσουμε περαιτέρω. Συνεπώς, η αρχική εξίσωση της $AR(1)$ μπορεί να γραφτεί ως:

$$(1 - \varphi B)X_t = Z_t$$

δηλαδή,

$$X_t = Z_t / (1 - \varphi B) = (1 + \varphi B + \varphi^2 B^2 + \dots) Z_t = Z_t + \varphi Z_{t-1} + \varphi^2 Z_{t-2} + \dots$$

Από την παραπάνω έκφραση προκύπτει

$$E(X_t) = 0$$

$$Var(X_t) = \sigma_z^2 (1 + \varphi^2 + \varphi^4 + \dots)$$

υπό την ανεξαρτησία των Z . Η διακύμανση είναι πεπερασμένη υπό τον όρο $|\varphi|^2 < 1$, οπότε αν $|\varphi| < 1$, σε αυτήν την περίπτωση

$$Var(X_t) = \sigma_x^2 = \sigma_z^2 / (1 - \varphi^2)$$

Προκύπτει μετά από πράξεις ότι η συνάρτηση αυτοσυνδιακύμανσης για υστέρηση k

$$\gamma(k) = Cov(X_t, X_{t+k}) = \frac{\sigma_z^2 \varphi^k}{1 - \varphi^2} = \varphi^k \sigma_x^2$$

Για $k < 0$ ισχύει ότι $\gamma(k) = \gamma(-k)$. Επειδή $\gamma(k)$ δεν εξαρτάται από t τότε η διαδικασία AR τάξης 1 είναι ασθενώς στάσιμη υπό τον όρο $|\varphi| < 1$. Η συνάρτηση αυτοσυσχέτισης της $AR(1)$ δίνεται παρακάτω

$$\rho(k) = \frac{\gamma(k)}{\gamma(0)} = \varphi^k, \quad k \geq 0$$

2.7 Χρονοσειρά κινητού μέσου (MA)

Ας υποθέσουμε ότι $\{Z_t\}$ είναι τυχαίος θόρυβος με μέση τιμή μηδέν και διακύμανση σ_z^2 . Τότε η διαδικασία $\{X_t\}$ λέγεται διαδικασία κινητού μέσου τάξης q , δηλαδή $MA(q)$, αν

$$X_t = Z_t + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \dots + \theta_q Z_{t-q}$$

Μπορούμε επίσης χρησιμοποιώντας τον τελεστή υστέρησης B να γράψουμε την παραπάνω διαδικασία ως

$$X_t = \theta(B)Z_t$$

όπου $\theta(B) = 1 + \theta_1 B + \theta_2 B^2 + \dots + \theta_q B^q$

Σε αντίθεση με την αυτοπαλίνδρομη διαδικασία, η διαδικασία κινητού μέσου είναι στάσιμη για κάθε τιμή των παραμέτρων $\theta_1, \theta_2, \dots, \theta_q$. Πιο συγκεκριμένα η $MA(q)$ είναι πάντα στάσιμη καθώς δίνεται ως πεπερασμένο άθροισμα όρων λευκού θορύβου. Όπως έχουμε δει, μια διαδικασία $AR(p)$ μπορεί να εκφραστεί ως διαδικασία MA άπειρης τάξης, το ίδιο ισχύει και για μια διαδικασία $MA(q)$ η οποία μπορεί να εκφραστεί ως διαδικασία AR άπειρης τάξης αν επιβάλουμε κάποιους περιορισμούς στις παραμέτρους. Τότε μια διαδικασία MA είναι αντιστρέψιμη και έχει κάποιες επιθυμητές μαθηματικές ιδιότητες (Hyndman & Athanasopoulos, 2013). Μια διαδικασία MA τάξης q είναι αντιστρέψιμη, όταν οι ρίζες του χαρακτηριστικού πολυωνύμου $\theta(B)$ βρίσκονται εκτός του μοναδιαίου κύκλου. Θεωρώντας το B ως μεταβλητή με τιμές στο σύνολο των μιγαδικών αριθμών, το προηγούμενο σημαίνει ότι οι ρίζες του πολυωνύμου έχουν μέτρο μεγαλύτερο της μονάδας.

Αν εκφράσουμε μια πρώτης τάξεως διαδικασία ως

$$Z_t = \frac{1}{\theta(B)} X_t$$

Στην περίπτωση πρώτης τάξεως διαδικασίας, έχουμε $\theta(B) = 1 + \theta B$, το οποίο έχει ρίζα $B = -1/\theta$. Υπό την προϋπόθεση ότι $|\theta| < 1$, η ρίζα $B = -1/\theta$ είναι πραγματική και βρίσκεται εκτός του μοναδιαίου κύκλου. Συνεπώς, βλέπουμε ότι το μοντέλο $MA(1)$ είναι αντιστρέψιμο αν $|\theta| < 1$. Αν θεωρήσουμε το B ως μεταβλητή με τιμές στο σύνολο των μιγαδικών αριθμών, ο τελεστής $1/\theta(B)$ μπορεί να εκφραστεί ως

$$\frac{1}{1 + \theta B} = 1 + \sum_{i=1}^{\infty} (-\theta)^i B^i$$

Όταν $|\theta| < 1$, αυτή η άπειρη σειρά συγκλίνει

$$\left| 1 + \sum_{i=1}^{\infty} (-\theta)^i \right| \leq 1 + \sum_{i=1}^{\infty} |\theta|^i = \frac{1}{1 - |\theta|}$$

και επομένως η $Z_t = \frac{1}{\theta(B)} X_t$ γίνεται

$$Z_t = \left(1 + \sum_{i=1}^{\infty} (-\theta)^i B^i \right) X_t = X_t + \sum_{i=1}^{\infty} (-\theta)^i X_{t-1}$$

το οποίο υποδηλώνει αντιστρεψιμότητα.

2.7.1 Χρονοσειρά κινητού μέσου τάξης 1

Παρακάτω μπορούμε να δούμε τη διαδικασία κινητού μέσου τάξης ένα, $MA(1)$:

$$X_t = Z_t + \theta_1 Z_{t-1}, \quad Z_t \sim WN(0, \sigma_z^2)$$

μπορούμε να αναπαραστήσουμε την εξάρτηση που παρουσιάζεται στην εξίσωση παρακάτω

Το μοντέλο $MA(1)$ έχει συνάρτηση αυτοδιακύμανσης

$$\gamma(\kappa) = \begin{cases} (1 + \theta^2)\sigma_z^2, & \kappa = 0 \\ \theta\sigma_z^2, & \kappa = \pm 1 \\ 0, & |\kappa| > 1 \end{cases}$$

και συνάρτηση αυτοσυσχέτισης

$$\rho(\kappa) = \begin{cases} 1, & \kappa = 0 \\ \frac{\theta}{1+\theta^2}, & \kappa = \pm 1 \\ 0, & |\kappa| > 1 \end{cases}$$

Αξίζει να σημειωθεί ότι $|\rho(1)| \leq 1/2$ για όλες τις τιμές του θ . Επίσης, η X_t είναι συσχετισμένη με την X_{t-1} αλλά όχι με τις X_{t-2}, X_{t-3}, \dots . Αυτό έρχεται σε αντίθεση με την διαδικασία $AR(1)$ που η συσχέτιση μεταξύ X_t και $X_{t+\kappa}$ δεν είναι ποτέ μηδενική.

2.8 Αυτοπαλίνδρομα μοντέλα κινητού μέσου ARMA

Μια χρήσιμη οικογένεια μοντέλων προκύπτει όταν συμπεριληφθούν ταυτόχρονα όροι των μοντέλων AR και MA . Οι διαδικασίες AR και MA έχουν διαφορετικές και συμπληρωματικές ιδιότητες κάτι το οποίο μας οδηγεί φυσικά να θεωρήσουμε την σύνθεση των δυο διαδικασιών. Η διαδικασία $ARMA$ έχει τεράστια επιτυχία εδώ και δεκαετίες ως ένα απλό μοντέλο με αξιοσημείωτη απόδοση (Politis & McElroy, 2020). Έγινε δημοφιλής από τους Box and Jenkins

(1970) ενώ αξίζει να σημειωθεί ότι τα αυτοπαλίνδρομα μοντέλα τα οποία αποτελούν ένα σημαντικό υποσύνολο των μοντέλων των Box και Jenkins έγιναν γνωστά από τον Yule (1927). Τα μοντέλα των Box και Jenkins καλύπτουν αδιαμφισβήτητα ένα μεγάλο κομμάτι της ανάλυσης χρονοσειρών και η εκτίμηση των παραμέτρων των μοντέλων συνήθως βασίζεται στην εκτίμηση μέγιστης πιθανοφάνειας ή ελαχίστων τετραγώνων υπό την υπόθεση της κανονικότητας των σφαλμάτων. Αν τα σφάλματα αποκλίνουν σημαντικά από την κανονική κατανομή τότε οι εκτιμήσεις δεν θα είναι ικανοποιητικές.

Η διαδικασία $\{X_t\}$ είναι μια $ARMA(p, q)$ διαδικασία αν $\{X_t\}$ είναι ασθενώς στάσιμη και ικανοποιεί

$$X_t - \varphi_1 X_{t-1} - \varphi_2 X_{t-2} - \dots - \varphi_p X_{t-p} = Z_t + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \dots + \theta_q Z_{t-q}$$

με $Z_t \sim WN(0, \sigma_z^2)$. Επιπλέον, η $\{Y_t\}$ λέγεται ότι είναι $ARMA(p, q)$ με μέσο μ αν $X_t = Y_t - \mu$ ικανοποιεί την παραπάνω. Για την συνέχεια θα υποθέσουμε ότι η διαδικασία $ARMA$ έχει μέση τιμή μηδέν. Συνεπώς μια χρονοσειρά $\{X_t\}$ ακολουθεί μια διαδικασία $ARMA$ τάξης (p, q) όταν

$$X_t = \varphi_1 X_{t-1} + \varphi_2 X_{t-2} + \dots + \varphi_p X_{t-p} + Z_t + \theta_1 Z_{t-1} + \theta_2 Z_{t-2} + \dots + \theta_q Z_{t-q}$$

Με τη χρήση πολυωνύμων υστέρησης παίρνουμε

$$\varphi(B)X_t = \theta(B)Z_t$$

όπου, $\varphi(B) = 1 - \varphi_1 B - \dots - \varphi_p B^p$ και $\theta(B) = 1 + \theta_1 B + \dots + \theta_q B^q$

Οι συνθήκες στις παραμέτρους του μοντέλου για να γίνει η διεργασία στάσιμη και αντιστρέψιμη είναι οι ίδιες με αυτές για μια διεργασία καθαρού AR και καθαρού MA , δηλαδή οι τιμές του $\{\varphi_i\}$ που καθιστούν τη διαδικασία στάσιμη να είναι τέτοιες ώστε οι ρίζες του χαρακτηριστικού πολυωνύμου

$$\varphi(B) = 0$$

να βρίσκονται έκτος του μοναδιαίου κύκλου, ενώ οι τιμές $\{\theta_i\}$ που κάνουν την διαδικασία αντιστρέψιμη, να είναι τέτοιες ώστε οι ρίζες του

$$\theta(B) = 0$$

να είναι εκτός του μοναδιαίου κύκλου.

Συνεπώς, η στασιμότητα της $ARMA$ ορίζεται από το αυτοπαλίνδρομο μέρος της και τότε μπορεί να εκφραστεί ως

$$X_t = \frac{\theta(B)}{\varphi(B)} Z_t$$

ενώ η αντιστρεψιμότητα της διαδικασίας $ARMA$ μπορεί να οριστεί από το μέρος κινητού μέσου της και μπορεί να εκφραστεί ως

$$\frac{\varphi(B)}{\theta(B)}X_t = Z_t$$

Η σημασία της διαδικασίας $ARMA$ εντοπίζεται στο γεγονός ότι μια στάσιμη χρονοσειρά μπορεί συχνά να μοντελοποιηθεί επαρκώς από ένα μοντέλο $ARMA$ που περιλαμβάνει λιγότερες παραμέτρους από όσες μια καθαρή διαδικασία AR ή MA από μόνη της (Chatfield & Xing, 2019). Αυτό είναι ένα αρχικό παράδειγμα της αρχής της οικονομίας (Principle of Parsimony), η οποία μας λέει ότι θέλουμε να βρούμε ένα μοντέλο με όσο το δυνατόν λιγότερες παραμέτρους το οποίο όμως να δίνει επαρκή αναπαράσταση των δεδομένων που έχουμε στη διάθεση μας. Τέλος, αξίζει να σημειωθεί ότι ένα $AR(p)$ μοντέλο είναι μια ειδική περίπτωση ενός $ARMA(p, 0)$, ενώ ένα $MA(q)$ μοντέλο είναι μια ειδική περίπτωση ενός $ARMA(0, q)$.

Ειδικότερα το μοντέλο $ARMA(1,1)$ είναι

$$X_t = \varphi X_{t-1} + Z_t + \theta Z_{t-1}, \quad Z_t \sim WN(0, \sigma_z^2)$$

ή, ισοδύναμα

$$X_t = (1 - \varphi B)^{-1}(1 + \theta B)Z_t$$

με συνθήκη στασιμότητας $|\varphi| < 1$ και συνθήκη αντιστρεψιμότητας $|\theta| < 1$. Αναπτύσσοντας τον όρο $(1 - \varphi B)^{-1}$ στο δεξί μέλος

$$\begin{aligned} X_t &= (1 + \varphi B + \varphi^2 B^2 + \dots)(1 + \theta B)Z_t \\ &= \left(\sum_{i=0}^{\infty} \varphi^i B^i \right) (1 + \theta B)Z_t \\ &= \left(1 + \sum_{i=0}^{\infty} \varphi^{i+1} B^{i+1} + \sum_{i=0}^{\infty} \varphi^i \theta B^{i+1} \right) Z_t \\ &= Z_t + (\varphi + \theta) \sum_{i=1}^{\infty} \varphi^{i-1} Z_{t-1} \end{aligned}$$

Από την παραπάνω εξίσωση προκύπτει ότι $E(X_t) = 0$ και

$$Var(X_t) = Var\left(Z_t + (\varphi + \theta) \sum_{i=1}^{\infty} \varphi^{i-1} Z_{t-1} \right) = \sigma_z^2 + \sigma_z^2 (\varphi + \theta)^2 (1 - \varphi^2)^{-1}$$

Η αυτοσυνδιακύμανση $\gamma(k)$ για $k > 0$, δίνεται παρακάτω

$$\begin{aligned}
Cov(X_t, X_{t+k}) &= (\varphi + \theta)\varphi^{\kappa-1}\sigma_z^2 + (\varphi + \theta)^2\sigma_z^2\varphi^\kappa \sum_{i=1}^{\infty} \varphi^{2i-2} \\
&= (\varphi + \theta)\varphi^{\kappa-1}\sigma_z^2 + (\varphi + \theta)^2\sigma_z^2\varphi^\kappa(1 - \varphi^2)^{-1}
\end{aligned}$$

Η αυτοσυσχέτιση $\rho(\kappa)$ προκύπτει ως

$$\rho(\kappa) = \frac{\gamma(\kappa)}{\gamma(0)} = \frac{Cov(X_t, X_{t+k})}{Var(X_t)} = \frac{\varphi^{\kappa-1}(\varphi + \theta)(1 + \varphi\theta)}{1 + \varphi\theta + \theta^2}$$

Κεφάλαιο 3

3.1 Βάσει μοντέλου Bootstrap στα κατάλοιπα

Ένας στατιστικός όταν έχει διαθέσιμα ιστορικά δεδομένα μπορεί να προσαρμόσει κάποια μοντέλα ανάλογα με τα δεδομένα του, με σκοπό να του δώσουν έγκυρες προβλέψεις. Μερικά κλασικά μοντέλα χρονοσειρών είδαμε στο κεφάλαιο 2. Το κλασικό bootstrap του Efron (1979) φαίνεται να μην δίνει τα επιθυμητά αποτελέσματα όταν τα δεδομένα μας είναι εξαρτημένα σύμφωνα με τον Singh (1981). Σε δεδομένα χρονοσειρών υπάρχουν δύο μέθοδοι εφαρμογής bootstrap. (α) Η προσέγγιση βάσει μοντέλου, στην οποία υποθέτουμε μια παραμετρική προσαρμογή στην υπάρχουσα χρονοσειρά και στη συνέχεια πραγματοποιούμε bootstrap στα σχεδόν iid κατάλοιπα του προσαρμοσμένου μοντέλου. (β) Η μη παραμετρική, ελεύθερου μοντέλου, προσέγγιση όπου η δειγματοληψία με επανάθεση πραγματοποιείται σε blocks παρατηρήσεων της αρχικής χρονοσειράς χωρίς να έχουμε κάνει κάποια παραδοχή. Στη συνέχεια θα αναλύσουμε την προσέγγιση bootstrap βάσει μοντέλου στα κατάλοιπα (residuals) ενός προσαρμοσμένου μοντέλου και συγκεκριμένα θα δούμε την περίπτωση μιας στάσιμης διαδικασίας AR . Πιο συγκεκριμένα, θα δούμε πώς μπορεί να πραγματοποιηθεί bootstrap στα κατάλοιπα και στη συνέχεια να δημιουργήσουμε χρονοσειρές bootstrap από τις οποίες θα υπολογίσουμε αντίστοιχους εκτιμητές bootstrap για τους συντελεστές του μοντέλου. Για να δούμε πώς μπορεί να εφαρμοστεί το bootstrap σε ένα μοντέλο AR θα παρουσιάσουμε την περίπτωση της εφαρμογής του σε ένα πρώτης τάξης $AR(1)$ μοντέλο ενώ ας σημειωθεί ότι οι Efron and Tibshirani (1986) συνέχισαν την εφαρμογή της μεθόδου και σε ένα $AR(2)$ μοντέλο. Παρακάτω, μπορούμε να δούμε ένα $AR(1)$ μοντέλο

$$Y_t = \varphi_1 Y_{t-1} + e_t, \quad e_t \sim WN(0, \sigma_e^2)$$

όπου y_t είναι η παρατήρηση στο χρόνο t πιθανώς κεντραρισμένη στο μηδέν. Αν η μέση τιμή της παρατηρούμενης χρονοσειράς δεν είναι μηδέν, τότε η δειγματική μέση τιμή θα αφαιρεθεί από κάθε όρο της χρονοσειράς ώστε να κεντραριστεί γύρω από το μηδέν. Το φ_1 είναι η άγνωστη παράμετρος του μοντέλου που η τιμή της κυμαίνεται ανάμεσα στο -1 και 1. Συνήθως, σε τέτοια μοντέλα υποθέτουμε κανονικότητα για τα σφάλματα και προχωράμε σε εκτίμηση και συμπερασματολογία βασισμένοι στην υπόθεση της κανονικότητας. Όμως, χρησιμοποιώντας bootstrap μπορούμε να αποφύγουμε την υπόθεση της κανονικότητας. Ας αναφέρουμε ότι στο παραπάνω μοντέλο ότι αν y_n είναι η τελευταία παρατήρηση τότε η πρόβλεψη ένα βήμα μετά, για τη χρονική στιγμή $n + 1$ θα γίνει χρησιμοποιώντας τον εκτιμητή $\hat{\varphi}_1$, δηλαδή η πρόβλεψη θα είναι $\hat{\varphi}_1 y_n$. Είναι σημαντικό να σημειωθεί ότι ο Stine (1987) χρησιμοποίησε μια προσέγγιση

του bootstrap μαζί με την κλασική γκαουσιανή προσέγγιση για να παράσχει προβλέψεις και διαστήματα εμπιστοσύνης και απέδειξε ότι αν και το bootstrap δεν ήταν τόσο αποδοτικό όσο η κλασική εκτίμηση όταν ισχύει η κανονικότητα των σφαλμάτων, σε περιπτώσεις που αυτή δεν ισχύει, τα διαστήματα εμπιστοσύνης bootstrap είναι πολύ καλύτερα και έχουν πιθανότητα κάλυψης που τείνει προς την ονομαστική.

Για να εφαρμόσουμε bootstrap σε ένα μοντέλο $AR(1)$ πρέπει πρώτα να εκτιμήσουμε τον συντελεστή του μοντέλου φ_1 , η εκτίμηση του οποίου θα μπορούσε να πραγματοποιηθεί με την μέθοδο ελαχίστων τετραγώνων ή την μέθοδο μέγιστης πιθανοφάνειας. Οι διαφορές των εκτιμητών που προκύπτουν από τις δύο μεθόδους είναι συνήθως μικρές όπως σημειώνουν οι Efron and Tibshirani (1993). Σύμφωνα με τους Thombs and Schucany (1990) αν χρησιμοποιήσουμε κάποια άλλη μέθοδο εκτίμησης εκτός της μεθόδου ελαχίστων τετραγώνων ενδέχεται τα κατάλοιπα να μην είναι κεντραρισμένα στο μηδέν οπότε θα χρειαστεί να τα κεντράρουμε, δηλαδή να θέσουμε $\tilde{e}_t = \hat{e}_t - \bar{e}$. Σύμφωνα με τον Lahiri (2003), αν δεν χρησιμοποιήσουμε τα κεντραρισμένα κατάλοιπα, τα αποτελέσματα του bootstrap ενδέχεται να έχουν κάποια μεροληψία. Συνεπώς, μετά την προσαρμογή του μοντέλου προκύπτουν τα κατάλοιπα

$$\hat{e}_t = y_t - \hat{\varphi}_1 y_{t-1}, \quad \text{για } t = 2, 3, \dots, n$$

Σημειώνεται ότι δεν μπορούμε να υπολογίσουμε το κατάλοιπο \hat{e}_1 καθώς το y_0 δεν είναι διαθέσιμο. Ένα δείγμα bootstrap $y_1^*, y_2^*, \dots, y_n^*$ δημιουργείται πραγματοποιώντας bootstrap στα κατάλοιπα. Πιο συγκεκριμένα, κάνοντας δειγματοληψία με επανάθεση από τα κατάλοιπα $\hat{e}_2, \hat{e}_3, \dots, \hat{e}_n$ έχουμε ένα δείγμα καταλοίπων bootstrap $e_2^*, e_3^*, \dots, e_n^*$ και χρησιμοποιώντας τη δομή του μοντέλου, με αναδρομή προκύπτει ένα δείγμα bootstrap ή αλλιώς χρονοσειρά bootstrap:

$$\begin{aligned} y_2^* &= \hat{\varphi}_1 y_1^* + e_2^* \\ y_3^* &= \hat{\varphi}_1 y_2^* + e_3^* \\ &\dots \\ y_n^* &= \hat{\varphi}_1 y_{n-1}^* + e_n^* \end{aligned}$$

Σύμφωνα με τον Efron and Tibshirani (1986) παίρνουμε $y_1^* = y_1$ σε κάθε δείγμα bootstrap, δηλαδή σε ένα $AR(1)$ χρησιμοποιούμε την αρχική τιμή της χρονοσειράς ενώ σε ένα μοντέλο $AR(p)$ τις p αρχικές τιμές. Στη συνέχεια από κάθε δείγμα bootstrap, με τιμές $y_1^*, y_2^*, \dots, y_n^*$, προκύπτει προσαρμόζοντας μοντέλο $AR(1)$ ο εκτιμητής $\hat{\varphi}_1^*$. Ο εκτιμητής ελαχίστων τετραγώνων (OLS) της παραμέτρου φ_1 είναι ο $\hat{\varphi}_1$ (ο οποίος συμπίπτει με τον εκτιμητή μέγιστης πιθανοφάνειας MLE σε περίπτωση κανονικότητας των σφαλμάτων), ενώ ο εκτιμητής bootstrap της παραμέτρου σε ένα δείγμα bootstrap είναι ο $\hat{\varphi}_1^*$.

$$\hat{\phi}_1 = \frac{\sum_{t=2}^n Y_t Y_{t-1}}{\sum_{t=2}^n Y_{t-1}^2}, \quad \hat{\phi}_1^* = \frac{\sum_{t=2}^n Y_t^* Y_{t-1}^*}{\sum_{t=2}^n (Y_{t-1}^*)^2}$$

Συνεπώς, για B επαναλήψεις bootstrap στο τέλος της μεθόδου έχουμε B χρονοσειρές bootstrap $y_1^*, y_2^*, \dots, y_n^*$ και B τιμές $\hat{\phi}_1^*$ τις οποίες μπορούμε να χρησιμοποιήσουμε για να εκτιμήσουμε το τυπικό σφάλμα του $\hat{\phi}_1$ και όχι μόνο.

Σύμφωνα με τον Davison and Hinkley (1997) οι χρονοσειρές bootstrap που προέκυψαν σε ένα $AR(1)$ μοντέλο δεν ήταν στάσιμες για αυτό πρότειναν είτε να επιλεγεί μια αρχική τιμή έτσι ώστε να δημιουργήσει ισορροπία στις χρονοσειρές είτε να κάνουμε δειγματοληψία bootstrap για μεγαλύτερο αριθμό καταλοίπων και να χρησιμοποιήσουμε μια περίοδο burn-in, ώστε να πλησιάσει η χρονοσειρά τη στασιμότητα. Γι' αυτό επιλέγουμε κάποιο κατάλληλο k και βάσει αυτού αποκλείουμε τις τιμές y_{-k}^*, \dots, y_0^* και κρατάμε τις $y_1^*, y_2^*, \dots, y_n^*$.

Ο Stine (1987) χρησιμοποίησε ένα ανάπτυγμα Taylor για να εκτιμήσει το MSE της πρόβλεψης που δουλεύει καλά όταν η κατανομή των σφαλμάτων είναι κανονική. Αυτή η προσέγγιση φάνηκε να λειτουργεί καλά για δεδομένα που ακολουθούν κανονική κατανομή. Για μη κανονικές περιπτώσεις ο Stine (1987) παρείχε μια εκτίμηση bootstrap του MSE και παρόλο που ήταν μεροληπτική, έδινε ακριβέστερες προβλέψεις και πιο ικανοποιητικά διαστήματα εμπιστοσύνης. Ας δούμε τώρα την προσέγγιση του Stine, η οποία έχει προταθεί για τέτοια μοντέλα όταν τα σφάλματα δεν είναι κανονικά. Ο Stine έκανε τις παρακάτω υποθέσεις.

- Τα σφάλματα έχουν μηδενικό μέσο
- Τα σφάλματα έχουν πεπερασμένες δεύτερες ροπές
- Τα σφάλματα είναι συμμετρικά γύρω από το μηδέν
- Η αθροιστική συνάρτηση κατανομής είναι γνησίως αύξουσα

Η υπόθεση των πεπερασμένων δεύτερων ροπών εξαιρεί κατανομές σφαλμάτων με βαριές ουρές. Οι παραπάνω υποθέσεις ικανοποιούνται όταν τα σφάλματα είναι κανονικά. Η κύρια διαφορά μεταξύ της προσέγγισης των Efron και Tibshirani και του Stine είναι στην υπόθεση της συμμετρικότητας των σφαλμάτων γύρω από το μηδέν. Ο Stine εισήγαγε τη συνάρτηση κατανομής

$$F_n(x) = 1/2 + (L(x)/\{2(n-p)\}), \quad x \geq 0, t = p+1, \dots, n$$

$$= 1 - F_n(-x), \quad x < 0$$

όπου, $L(x) = \text{πλήθος των } t \text{ για τα οποία } k|\hat{\epsilon}_t| \leq x$ και

$$k = [(n-p)/(n-2p)]^{1/2}$$

και έκανε δειγματοληψία βάσει αυτής αντί βάσει της εμπειρικής κατανομής.

Η επιλογή της F_n παράγει κατάλοιπα bootstrap τα οποία είναι συμμετρικά γύρω από το μηδέν και έχουν ίδια διακύμανση με των αρχικών καταλοίπων. Αξίζει να σημειωθεί πως ο Stine (1987) χρησιμοποίησε μια εναλλακτική μέθοδο για αρχικές τιμές. Όπως είδαμε οι Efron και Tibshirani για ένα $AR(p)$ μοντέλο χρησιμοποιούν τις p πρώτες τιμές της αρχικής χρονοσειράς για να δημιουργήσουν ένα δείγμα bootstrap με τις τιμές $y_1^*, y_2^*, \dots, y_n^*$, ενώ ο Stine στην προσέγγιση του πρότείνει να επιλεγεί ένα τυχαίο μπλοκ p διαδοχικών παρατηρήσεων. Για περισσότερες λεπτομέρειες σχετικά με την προσέγγιση του Stine, ανατρέξτε στο Stine (1987). Αν και η υπόθεση της κανονικότητας στη προσέγγιση bootstrap βάσει μοντέλου δεν χρειάζεται, το γενικότερο πρόβλημα σε αυτήν την μέθοδο προκύπτει στο ότι δεν γνωρίζουμε πραγματικά εάν η βασική στοχαστική δομή εξηγείται από ένα μοντέλο $AR(p)$ και ότι η τάξη του μοντέλου είναι p .

Ένα σημαντικό πρόβλημα της στατιστικής ανάλυσης των χρονοσειρών είναι το εξής: Δεδομένης μιας πραγματοποίησης της διαδικασίας έως το χρόνο n τι θα μπορούσαμε να πούμε για την παρατήρηση στον χρόνο $(n + k)$; Σύμφωνα με τους Thombs and Schucany (1990), η παραπάνω μέθοδος που αναφέραμε λειτουργεί καλά σε ένα μοντέλο $AR(p)$ για την εκτίμηση των τυπικών σφαλμάτων των εκτιμητών $\hat{\phi}_j, j = 1, \dots, p$. Ωστόσο, όταν χρησιμοποιείται bootstrap για την εκτίμηση της δεσμευμένης κατανομής του Y_{n+k} δοθείσης της παρατηρηθείσας χρονοσειράς, τότε για να την προσομοιώσουμε αποτελεσματικά θα πρέπει (α) τα δείγματα bootstrap $Y_1^*, Y_2^*, \dots, Y_B^*$ να παρομοιάζουν τη δομή της συσχέτισης της χρονοσειράς που θέλουμε να προβλέψουμε και (β) αφού ξέρουμε ότι η πρόβλεψη του μοντέλου για $1, 2, \dots, k$ χρονικά βήματα μετά εξαρτάται από τους εκτιμημένους συντελεστές και τις p τελευταίες παρατηρήσεις κάθε φορά, να επιλέγονται οι p τελευταίες τιμές σε κάθε δείγμα bootstrap ως $y_t^* = y_t, y_{t-1}^* = y_{t-1}, \dots, y_{t-p+1}^* = y_{t-p+1}$. Αυτό το πέτυχαν οι Thombs και Schucany εφαρμόζοντας bootstrap σε μια εναλλακτική αναπαράσταση της $AR(p)$, την οπισθοδρομική αναπαράσταση πληροφορίες της οποίας μπορείτε να βρείτε στο άρθρο τους Thombs and Schucany (1990).

Τέλος, αξίζει να σημειωθεί ότι ο Lahiri (2003) παρουσίασε εφαρμογή της μεθόδου bootstrap σε ένα μοντέλο $ARMA$ την οποία και αποκάλεσε ως $ARMA$ Bootstrap ή αλλιώς $ARMAB$. Η συνέπεια της μεθόδου για την εκτίμηση των παραμέτρων ενός $ARMA(p, q)$ μοντέλου μπορεί να δειχθεί αξιοποιώντας το γεγονός ότι η διαδικασία πεπερασμένων παραμέτρων $ARMA(p, q)$ έχει μια άπειρη κινητού μέσου αλλά και αυτοπαλίνδρομη αναπαράσταση. Η συγκεκριμένη μέθοδος προέρχεται από το αρχικό άρθρο του Kreiss and Franke (1992) οι οποίοι εδραίωνουν την συνέπεια της.

3.1.2 Βάσει μοντέλου Bootstrap στα κατάλοιπα ενός AR(p)

Η βάσει μοντέλου μέθοδος bootstrap στα κατάλοιπα μπορεί να γενικευτεί σε μια διαδικασία $AR(p)$. Ας υποθέσουμε ότι $\{Y_t\}_{t \in \mathbb{Z}}$ είναι μια στάσιμη αυτοπαλίνδρομη διαδικασία τάξεως p , $AR(p)$ που ικανοποιεί την παρακάτω εξίσωση

$$Y_t = \varphi_1 Y_{t-1} + \dots + \varphi_p Y_{t-p} + e_t, \quad e_t \sim WN(0, \sigma_e^2)$$

όπου $p \in \mathbb{N}$, $\varphi_1, \dots, \varphi_p$ είναι οι παράμετροι του μοντέλου και $\{e_t\}_{t \in \mathbb{Z}}$ μια ακολουθία μηδενικού μέσου iid τυχαίων μεταβλητών με κοινή κατανομή F . Στη συνέχεια υποθέτουμε ότι οι παράμετροι $\varphi_1, \dots, \varphi_p$ είναι τέτοιες ώστε

$$\varphi(z) = 1 - \sum_{j=1}^p \varphi_j z^j = 1 - \varphi_1 z - \dots - \varphi_p z^p \neq 0 \quad \text{για όλα τα } z \in \mathbb{C} \text{ με } |z| \leq 1$$

Παρόλο που οι τυχαίες μεταβλητές Y_t, Y_{t-1}, \dots κάτω από το $AR(p)$ μοντέλο είναι εξαρτημένες, μπορούμε να χρησιμοποιήσουμε κατάλληλα τη δομή του μοντέλου όπως είδαμε και στην περίπτωση του $AR(1)$ και να προχωρήσουμε σε bootstrap εκτιμήσεις. Η ιδέα όπως έχουμε ήδη πει είναι δειγματοληψία με επανάθεση στα κατάλοιπα ή αλλιώς εκτιμημένα σφάλματα του $AR(p)$ μοντέλου τα οποία είναι περίπου ανεξάρτητα. Ας υποθέσουμε ότι παρατηρείται ένα πεπερασμένο τμήμα Y_1, \dots, Y_n της διαδικασίας $\{Y_t\}_{t \in \mathbb{Z}}$ δηλαδή έχουμε παρατηρήσεις y_1, \dots, y_n . Οι εκτιμητές ελαχίστων τετραγώνων των $\varphi_{1n}, \dots, \varphi_{pn}$ θα είναι οι $\hat{\varphi}_{1n}, \dots, \hat{\varphi}_{pn}$ και δίνονται από την σχέση

$$(\hat{\varphi}_{1n}, \dots, \hat{\varphi}_{pn})' = (\mathbf{V}'_n \mathbf{V}_n)^{-1} \mathbf{V}'_n (y_{p+1}, \dots, y_n)'$$

όπου \mathbf{V}_n είναι πίνακας διάστασης $(n-p) \times p$ με t γραμμή (y_{t+p-1}, \dots, y_t) , $t = 1, \dots, n-p$. Τα $\hat{e}_t = y_t - \hat{\varphi}_{1n} y_{t-1} - \dots - \hat{\varphi}_{pn} y_{t-p}$, $t = p+1, \dots, n$ είναι τα κατάλοιπα του μοντέλου.

Χρησιμοποιώντας bootstrap μπορούμε να κάνουμε δειγματοληψία με επανάθεση στα κατάλοιπα και να ορίσουμε την εκδοχή bootstrap μιας τυχαίας μεταβλητής $T_n = t_n(Y_1, \dots, Y_n; \varphi_1, \dots, \varphi_p, F)$. Σύμφωνα με τον Lahiri (2003) όπως έχουμε ήδη αναφέρει πρέπει πρώτα να κεντράρουμε τα κατάλοιπα και μετά να χρησιμοποιήσουμε bootstrap. Τα κεντραρισμένα κατάλοιπα είναι

$$\tilde{e}_t = \hat{e}_t - \bar{e}, \quad t = p+1, \dots, n$$

όπου $\bar{e} = (n-p)^{-1} \sum_{t=p+1}^n \hat{e}_t$. Στη συνέχεια πραγματοποιούμε δειγματοληψία με επανάθεση στα κεντραρισμένα κατάλοιπα $\{\tilde{e}_{p+1}, \dots, \tilde{e}_n\}$ και σχηματίζουμε δείγματα bootstrap μεταβλητών e_t^* , $t \in \mathbb{Z}$. Οι τυχαίες μεταβλητές e_t^* είναι υπό συνθήκη iid, δεδομένων των (y_1, \dots, y_n) με κοινή κατανομή

$$P_*(e_1^* = \tilde{e}_t) = \frac{1}{n-p}, \quad p+1 \leq t \leq n$$

Από τα δύο παραπάνω προκύπτει $E_*(e_1^*) = (n-p)^{-1} \sum_{t=p+1}^n \tilde{e}_t = 0$. Συνεπώς, οι μεταβλητές bootstrap e_t^* ικανοποιούν μια αναλογία της συνθήκης του μοντέλου $E(e_1) = 0$ αλλά στο επίπεδο του bootstrap. Παρακάτω μπορούμε να δούμε τώρα την εκδοχή bootstrap ενός $AR(p)$

$$Y_t^* = \hat{\varphi}_{1n} Y_{t-1}^* + \dots + \hat{\varphi}_{pn} Y_{t-p}^* + e_t^*, \quad t \in \mathbb{Z}$$

και έστω $\{Y_t^*\}_{t \in \mathbb{Z}}$ μια στάσιμη λύση της παραπάνω εξίσωσης. Αν $\hat{\varphi}_{jn} \rightarrow_p \varphi_j$ καθώς $n \rightarrow \infty$ για $j = 1, \dots, p$, τότε υπάρχει μια τέτοια λύση σε ένα σύνολο των Y_t με πιθανότητα κοντά στο ένα για μεγάλο n .

Συνεπώς, έχοντας πραγματοποιήσει δειγματοληψία με επανάθεση στα κεντραρισμένα κατάλοιπα, ένα δείγμα bootstrap προκύπτει χρησιμοποιώντας τη δομή του μοντέλου με αναδρομή:

$$\begin{aligned} y_{p+1}^* &= \hat{\varphi}_{1n} y_p^* + \dots + \hat{\varphi}_{pn} y_1^* + e_{p+1}^* \\ &\dots \\ &\dots \\ y_n^* &= \hat{\varphi}_{1n} y_{n-1}^* + \dots + \hat{\varphi}_{pn} y_{n-p}^* + e_n^* \end{aligned}$$

Τις p τιμές y_1^*, \dots, y_p^* μπορούμε να τις ορίσουμε ίσες με y_1, \dots, y_p ή ίσες με μηδέν σύμφωνα με τον Lahiri (2003).

Όταν το πολυώνυμο $\hat{\varphi}(z) = 1 - \sum_{j=1}^p \hat{\varphi}_{jn} z^j$ δεν μηδενίζεται στη περιοχή $\{|z| \leq 1\}$ τότε η χρονοσειρά πλησιάζει στη στασιμότητα με εκθετική ταχύτητα και συνεπώς οι αρχικές τιμές έχουν αμελητέα επίδραση μακροπρόθεσμα. Παίρνοντας μεγάλο μέγεθος δείγματος bootstrap από τα κεντραρισμένα κατάλοιπα, μπορούμε να δημιουργήσουμε μια χρονοσειρά μεγάλου μήκους, χρησιμοποιώντας την αναδρομική σχέση, μέχρι να επιτευχθεί η στασιμότητα, με την επίδραση των αρχικών τιμών να είναι αμελητέα, και να επιλέξουμε τις m πιο πρόσφατες τιμές ως χρονοσειρά bootstrap μεγέθους m .

Η αυτοπαλίνδρομη εκδοχή bootstrap της τυχαίας μεταβλητής $T_n = t_n(Y_1, \dots, Y_n; \varphi_1, \dots, \varphi_p, F)$ βασίζεται τώρα στις χρονοσειρές bootstrap μεγέθους $m > p$ και δίνεται από

$$T_{m,n}^* = t_m(Y_1^*, \dots, Y_m^*; \hat{\varphi}_{1n}, \dots, \hat{\varphi}_{pn}, \hat{F}_n)$$

όπου \hat{F}_n είναι η εμπειρική κατανομή των κεντραρισμένων καταλοίπων. Το μέγεθος των δειγμάτων bootstrap m συνήθως είναι ίσο με το μέγεθος του αρχικού δείγματος n ωστόσο σε

μερικές περιπτώσεις όπως θα δούμε παρακάτω είναι προτιμότερη η επιλογή μικρότερου μεγέθους, $m < n$, δηλαδή εφαρμογή της μεθόδου m -out-of- n bootstrap που περιγράψαμε στο κεφάλαιο 1. Τέλος, αξίζει να σημειωθεί ότι την παραπάνω διαδικασία bootstrap σε ένα αυτοπαλίνδρομο μοντέλο AR ο Lahiri (2003) την ονόμασε ARB .

Ο Bose (1988) διερεύνησε τις ιδιότητες του ARB όταν χρησιμοποιούνται κανονικοποιημένοι εκτιμητές ελαχίστων τετραγώνων των αυτοπαλίνδρομων παραμέτρων και έδειξε ότι το ARB είναι ισχυρά συνεπές (strongly consistent). Πιο συγκεκριμένα, έδειξε ότι κάτω από κατάλληλες συνθήκες αν η διαδικασία $AR(p)$ είναι στάσιμη και τα σφάλματα έχουν πεπερασμένη όγδοη ροπή και ικανοποιούν μια τεχνική συνθήκη η οποία καλείται συνθήκη του Cramer τότε οι εκτιμητές της ARB έχουν ρυθμό σύγκλισης $o(n^{-1/2})$ κάτι που αποδεικνύει ότι πέραν της συνέπειας τους, βελτιώνει και την κανονική προσέγγιση που έχει σφάλμα τάξεως $O(n^{-1/2})$. Αυτό σημαίνει ότι οι εκτιμητές bootstrap των παραμέτρων συγκλίνουν πιο γρήγορα στις σωστές τιμές σε σχέση με την κανονική προσέγγιση και επομένως είναι ακριβέστεροι.

Τέλος, αξίζει να σημειωθεί ότι η ιδέα της αξιοποίησης της δομής της εξίσωσης του αυτοπαλίνδρομου μοντέλου και κατάλληλης εφαρμογής του κλασικού σχήματος bootstrap του Efron (1979) σε μια διαδικασία AR έχει πρωτοαναφερθεί και εφαρμοστεί σε διαφορετικά προβλήματα από τους Freedman and Peters (1984), Efron and Tibshirani (1986) και Swanepoel and Van Wyk (1986).

3.1.3 Ασταθής αυτοπαλίνδρομη διαδικασία

Μια ασταθής αυτοπαλίνδρομη διαδικασία τάξεως p δίνεται παρακάτω

$$Y_t = \sum_{j=1}^p \varphi_j Y_{t-j} + e_t, \quad t \in \mathbb{Z}$$

Σε μια ασταθή διαδικασία τουλάχιστον μια από τις ρίζες του χαρακτηριστικού πολωνύμου $\Psi_p(z) = z^p - \varphi_1 z^{p-1} - \dots - \varphi_p$ βρίσκεται επάνω στον μοναδιαίο κύκλο, ενώ οι υπόλοιπες βρίσκονται μέσα στον μοναδιαίο κύκλο $\{z \in \mathbb{C} : |z| = 1\}$. Παρόλο που το ARB είναι συνεπές όταν έχουμε μια διαδικασία με τις ρίζες του χαρακτηριστικού πολωνύμου εντός του μοναδιαίου κύκλου, δεν είναι δίνει καλά αποτελέσματα στην συγκεκριμένη περίπτωση. Για απλότητα ας δούμε την περίπτωση του $AR(1)$.

$$Y_t = \varphi_1 Y_{t-1} + e_t, \quad t \geq 1$$

με $Y_0 = 0$ και $\varphi_1 \in \{-1, 1\}$ και $\{e_t\}$ για $t = 1, 2, \dots, n, \dots, \infty$ είναι η ακολουθία των σφαλμάτων. Αφού είμαστε σε διαδικασία πρώτης τάξεως, το χαρακτηριστικό πολυώνυμο είναι γραμμικό και έχει μόνο μια ρίζα, και αυτή η ρίζα πρέπει να βρίσκεται επάνω στον μοναδιαίο κύκλο. Το χαρακτηριστικό πολυώνυμο σε αυτήν την περίπτωση είναι $\Psi_1(z) = z - \varphi_1$. Η ρίζα αυτού του πολυωνύμου είναι $z = \varphi_1$. Οι μοναδικές τιμές που μπορεί να πάρει το φ_1 ώστε να βρίσκονται πάνω στον μοναδιαίο κύκλο είναι -1 και 1 . Αν η ρίζα πέφτει πάνω στον μοναδιαίο κύκλο τότε η διαδικασία θεωρείται ότι είναι στα όρια της σταθερότητας. Ο Lahiri (2003) επισημαίνει ότι η εκτίμηση της μεθόδου ελαχίστων τετραγώνων για το φ_1 παραμένει μια συνεπής εκτίμηση, αλλά έχει διαφορετικό ρυθμό σύγκλισης και διαφορετική οριακή κατανομή σε σχέση με στάσιμες περιπτώσεις. Συνεπώς, καταλαβαίνουμε γιατί το βάσει μοντέλου bootstrap στα κατάλοιπα αποτυγχάνει όταν έχουμε μια ασταθή αυτοπαλίνδρομη διαδικασία. Την λύση στο συγκεκριμένο πρόβλημα φαίνεται να την δίνει η μέθοδος m-out-of-n bootstrap την οποία είδαμε στο πρώτο κεφάλαιο. Η Datta (1996) και οι Heimann and Kreiss (1996) έδειξαν ανεξάρτητα ότι το ARB είναι συνεπές όταν το μέγεθος δειγμάτων bootstrap m τείνει στο άπειρο αλλά με πιο αργό ρυθμό από το n . Όμως τα θεωρήματα που έδειξαν απαιτούν συνθήκες για τις ροπές των σφαλμάτων. Το θεώρημα της Datta απαιτεί την ύπαρξη ροπής τάξης $2 + \delta$ για τα σφάλματα, για κάποιο $\delta > 0$, ενώ το θεώρημα των Heimann και Kreiss υποθέτει απλώς την ύπαρξη δεύτερης ροπής. Ωστόσο, η Datta αποδεικνύει έναν πιο αυστηρό ρυθμό σύγκλισης (σχεδόν βεβαίως), ενώ οι Heimann και Kreiss αποδεικνύουν μόνο σύγκλιση κατά πιθανότητα. Εφόσον το πρόβλημα με το bootstrap οφείλεται στην εκτίμηση ελαχίστων τετραγώνων του φ_1 και ο εκτιμημένος συντελεστής $\hat{\varphi}_1$ χρησιμοποιείται για να υπολογισθούν τα κατάλοιπα στα οποία θα πραγματοποιηθεί bootstrap, μια ιδέα είναι απλά να τροποποιηθεί αυτή η εκτίμηση. Οι Datta and Sriram (1997) χρησιμοποιούν μια εκτίμηση συρρίκνωσης για τον συντελεστή αντί για την εκτίμηση ελαχίστων τετραγώνων και επειδή η εκτίμησή τους έχει έναν πιο γρήγορο ρυθμό σύγκλισης, είναι σε θέση να δείξουν ότι αυτό το τροποποιημένο ARB είναι συνεπές. Για το m-out-of-n είναι πάντα ζήτημα το μέγεθος m . Αξίζει να σημειωθεί ότι οι Datta and McCormick (1995) χρησιμοποίησαν μια εκδοχή του Jackknife after Bootstrap για την επιλογή του m .

3.2 Bootstrap σε εξαρτημένα δεδομένα

Ας υποθέσουμε ότι έχουμε X_1, \dots, X_n χρονοσειρά με από κοινού κατανομή P_n . Για να εκτιμήσουμε μια πληθυσμιακή παράμετρο ϑ χρησιμοποιούμε έναν εκτιμητή $\hat{\vartheta}_n$. Ένα σύνηθες πρόβλημα που αντιμετωπίζει ένας στατιστικός είναι ο υπολογισμός ή η εκτίμηση της ακρίβειας αυτού του εκτιμητή, η οποία μπορεί για παράδειγμα να δοθεί με μια εκτίμηση του μέσου τετραγωνικού σφάλματος (MSE). Κάθε μέτρηση της ακρίβειας εξαρτάται από την δειγματική

κατανομή του $\hat{\vartheta}_n - \vartheta$ που είναι τυπικά άγνωστη στην πράξη και μερικές φορές πολύπλοκη. Οι μέθοδοι bootstrap σε χρονοσειρές μας δίνουν την δυνατότητα να προσεγγίσουμε την κατανομή του $\hat{\vartheta}_n$ και συναρτησιακών αυτής χωρίς να κάνουμε παραμετρικές υποθέσεις.

Ας υποθέσουμε τώρα ότι $\{X_1, \dots, X_n\} \equiv \mathbf{X}_n$ με από κοινού κατανομή P_n . Στην συνέχεια παίρνουμε παρατηρήσεις $\{X_1^*, \dots, X_n^*\} \equiv \mathbf{X}_n^*$ από την \hat{P}_n . Αν η \hat{P}_n είναι καλός εκτιμητής της P_n , τότε η σχέση μεταξύ $\{X_1, \dots, X_n\}$ και P_n αναπαράγεται σε μεγάλο βαθμό στον κόσμο του bootstrap από την σχέση των $\{X_1^*, \dots, X_n^*\}$ και \hat{P}_n . Συνεπώς, ορίζουμε την εκδοχή bootstrap του εκτιμητή $\hat{\vartheta}_n$ από το $\hat{\vartheta}_n^*$ αντικαθιστώντας τις X_1, \dots, X_n με X_1^*, \dots, X_n^* και ομοίως ορίζουμε το ϑ^* αντικαθιστώντας στο $\vartheta = \vartheta(P_n)$ το P_n με \hat{P}_n . Τότε, η δεσμευμένη συνάρτηση κατανομής \hat{G}_n ή G_n^* του $\hat{\vartheta}_n^* - \vartheta^*$ δεδομένου του \mathbf{X}_n , είναι η εκτίμηση bootstrap της κατανομής G_n του $\hat{\vartheta}_n - \vartheta$. Για να ορίσουμε τους εκτιμητές bootstrap συναρτησιακών της κατανομής του $\hat{\vartheta}_n - \vartheta$ όπως είναι η διακύμανσή του για παράδειγμα, μπορούμε απλά να χρησιμοποιήσουμε την αρχή της αντικατάστασης (plug-in principle) και να χρησιμοποιήσουμε την αντίστοιχη συνάρτηση στην δεσμευμένη κατανομή του $\hat{\vartheta}_n^* - \vartheta^*$. Συνεπώς, ο εκτιμητής bootstrap του σ_n^2 , δηλαδή της διακύμανσης του $\hat{\vartheta}_n - \vartheta$ δίνεται από την δεσμευμένη διακύμανση του $\hat{\vartheta}_n^* - \vartheta^*$.

$$\hat{\sigma}_n^2 = Var(\hat{\vartheta}_n^* - \vartheta^* | \mathbf{X}_n) = \int x^2 d\hat{G}_n(x) - \left[\int x d\hat{G}_n(x) \right]^2$$

Γενικά, έχοντας επιλέξει μια συγκεκριμένη μέθοδο bootstrap, είναι πολύ δύσκολο και συχνά αδύνατο να οδηγηθούμε σε κλειστής μορφής αναλυτικές εκφράσεις για τους εκτιμητές bootstrap διαφόρων πληθυσμιακών ποσοτήτων. Συνεπώς, εδώ έρχεται και μας βοηθάει ο υπολογιστής. Οι εκτιμητές bootstrap της κατανομής του $\hat{\vartheta}_n - \vartheta$ μπορούν να υπολογιστούν αριθμητικά χρησιμοποιώντας προσομοίωση Monte Carlo. Αρχικά, ένας μεγάλος αριθμός από ανεξάρτητα αντίγραφα bootstrap $\{\hat{\vartheta}_n^{*k} : k = 1, \dots, K\}$ του $\hat{\vartheta}_n^*$ προκύπτουν από επαναλαμβανόμενη δειγματοληψία. Η εμπειρική κατανομή αυτών των αντιγράφων bootstrap δίνει την επιθυμητή προσέγγιση Monte Carlo στην πραγματική κατανομή bootstrap του $\hat{\vartheta}_n^* - \vartheta^*$. Συγκεκριμένα, για την διακύμανση $\sigma_n^2 = Var(\hat{\vartheta}_n - \vartheta)$, η Monte Carlo προσέγγιση του bootstrap εκτιμητή της διακύμανσης $\hat{\sigma}_n^2$ δίνεται από

$$[\hat{\sigma}_n^{MC}]^2 = (K - 1)^{-1} \sum_{k=1}^K [\hat{\vartheta}_n^{*k} - K^{-1} \sum_{j=1}^K \hat{\vartheta}_n^{*j}]^2$$

3.3 Block Bootstrap

Για να χρησιμοποιήσουμε το παραπάνω βάσει μοντέλου bootstrap στα κατάλοιπα που περιγράψαμε, πρέπει να γνωρίζουμε την δομή του μοντέλου και τις παραμέτρους του. Αν

επιλέξουμε λανθασμένο μοντέλο τότε οι χρονοσειρές bootstrap που θα προκύψουν θα έχουν διαφορετική δομή από αυτή των πραγματικών δεδομένων κάτι το οποίο θα οδηγήσει σε διαφορετικά στατιστικά συμπεράσματα από αυτά που θα περιμέναμε. Συνεπώς, αξίζει να σημειωθεί ότι το bootstrap βάσει μοντέλου είναι ανάλογο του παραμετρικού bootstrap που είδαμε στο πρώτο κεφάλαιο. Η κύρια υπόθεση του bootstrap βάσει μοντέλου στα κατάλοιπα ή αλλιώς Residual Bootstrap είναι ότι χρησιμοποιήσαμε την προσεγγιστική ανεξαρτησία των καταλοίπων ώστε να εφαρμόσουμε την κλασική μέθοδο του bootstrap που είδαμε στα iid δεδομένα. Όμως, σε προβλήματα που ο στατιστικός δεν έχει κάποια πρότερη γνώση ώστε να προσδιορίσει κάποιο μοντέλο, αυτή η μέθοδος δεν είναι χρήσιμη. Ο Bose (1988) έδειξε ότι αν μια αυτοπαλίνδρομη διαδικασία είναι το σωστό μοντέλο για τα δεδομένα μας (ή τουλάχιστον σχεδόν σωστό), τότε υπάρχει πλεονέκτημα χρήσης του bootstrap βάσει μοντέλου λόγω των καλών ασυμπτωτικών ιδιοτήτων υψηλής τάξης για μεγάλο αριθμό στατιστικών συναρτήσεων που μπορούν να προκύψουν από το μοντέλο. Αν γνωρίζουμε ότι έχουμε στάσιμη χρονοσειρά αλλά δεν ξέρουμε την δομή της, τότε είμαστε σε μια κατάσταση ανάλογη με την μη παραμετρική προσέγγιση που είδαμε στο πρώτο κεφάλαιο. Το block bootstrap το οποίο πρωτοπαρουσιάστηκε από τον Carlstein (1986) και αναπτύχθηκε από τον Künsch (1989) είναι μια μέθοδος που πραγματοποιεί επαναδειγματοληψία σε blocks παρατηρήσεων και η οποία μπορεί να εφαρμοστεί ακόμα και αν δεν γνωρίζουμε το παραμετρικό μοντέλο. Σε αυτήν τη μέθοδο, προκειμένου να διατηρήσουμε στα δεδομένα μας κάποια από την εξάρτηση που έχουν, αντί να κάνουμε δειγματοληψία από τις παρατηρήσεις κάνουμε δειγματοληψία από blocks παρατηρήσεων τα οποία αν είναι καλά επιλεγμένα διατηρούν αρκετή από την πληροφορία που μας ενδιαφέρει. Για μια ασυσχέτιστη ανταλλάξιμη ακολουθία τυχαίων μεταβλητών η κλασική μη παραμετρική μέθοδος bootstrap του Efron που κάνει δειγματοληψία με επανάθεση μεμονωμένων παρατηρήσεων είναι η κατάλληλη. Για στάσιμες χρονοσειρές οι διαδοχικές παρατηρήσεις είναι συσχετισμένες αλλά οι παρατηρήσεις που απέχουν μεταξύ τους μεγάλο χρονικό διάστημα είναι σχεδόν ασυσχέτιστες. Έτσι, το κλειδί της μεθόδου δειγματοληψίας σε blocks είναι ότι για στάσιμες χρονοσειρές τα μεμονωμένα blocks παρατηρήσεων τα οποία είναι αρκετά μακριά στο χρόνο θα είναι περίπου ασυσχέτιστα και μπορούν να θεωρηθούν ως ανταλλάξιμα. Η ανταλλαξιμότητα που αναφέρθηκε είναι η ιδιότητα ενός τυχαίου δείγματος η οποία είναι ελαφρώς πιο ασθενής από την iid, δηλαδή πιο ασθενής από την ιδιότητα ότι οι παρατηρήσεις είναι ανεξάρτητες και ισόνομα κατανομημένες. Μαθηματικά, μια ακολουθία n παρατηρήσεων είναι ανταλλάξιμη αν η από κοινού κατανομή οποιωνδήποτε p συνεχόμενων παρατηρήσεων δεν αλλάζει όταν αλλάζουμε την σειρά των παρατηρήσεων. Με άλλα λόγια αν τα δεδομένα μπορούν να ανταλλάσσονται μεταξύ τους, δηλαδή να μεταθέσουμε τις παρατηρήσεις σε οποιαδήποτε σειρά, χωρίς να αλλάζει η πιθανοτική δομή του δείγματος, τότε το δείγμα είναι ανταλλάξιμο.

Η βασική ιδέα της δειγματοληψίας σε blocks είναι ότι από το σύνολο των n παρατηρήσεων κατασκευάζουμε k blocks μεγέθους ℓ , δηλαδή ℓ διαδοχικών παρατηρήσεων σε κάθε block. Αυτή η προσέγγιση έχει καλύτερα αποτελέσματα αν η εξάρτηση στην ακολουθία παρατηρήσεων είναι ασθενής και τα blocks είναι μεγάλα σε μέγεθος διατηρώντας πιστά την εξάρτηση που έχουν τα αρχικά δεδομένα (Davison & Hinkley, 1997). Αν υποθέσουμε ότι έχουμε χρονοσειρά μεγέθους $n = k\ell$, τότε μπορούμε να δημιουργήσουμε k μη επικαλυπτόμενα blocks μεγέθους ℓ το κάθε ένα. Αν δεν ισχύει $n = k\ell$ τότε κάποιο block μπορεί να έχει λιγότερες παρατηρήσεις. Δηλαδή, σε αυτήν την περίπτωση αντί να κάνουμε δειγματοληψία με επανάθεση από τις παρατηρήσεις, κάνουμε δειγματοληψία με επανάθεση από τα blocks παρατηρήσεων δημιουργώντας χρονοσειρές από τις οποίες μπορούμε να υπολογίσουμε το στατιστικό που μας ενδιαφέρει. Λόγω του ότι σε κάθε block υπάρχουν διαδοχικές παρατηρήσεις, είναι σαν να έχουμε πληροφορία για εξάρτηση μέχρι τάξεως ℓ .

Υπάρχουν διάφορες παραλλαγές της δειγματοληψίας σε blocks τις οποίες θα αναλύσουμε στη συνέχεια. Παρακάτω μπορούμε να δούμε την διαδικασία των μη επικαλυπτόμενων blocks (Non-overlapping Block Bootstrap ή αλλιώς NBB).

Αν υποθέσουμε ότι έχουμε στη διάθεση μας 15 παρατηρήσεις $(x_1, x_2, \dots, x_{15})$. Τότε, δημιουργώντας $k = 5$ blocks με $\ell = 3$ παρατηρήσεις το κάθε ένα, θα έχουμε $\mathfrak{B}_1 = (x_1, x_2, x_3)$, $\mathfrak{B}_2 = (x_4, x_5, x_6)$, ..., $\mathfrak{B}_5 = (x_{13}, x_{14}, x_{15})$. Συνεχίζοντας για να πάρουμε μια χρονοσειρά bootstrap πραγματοποιούμε τυχαία δειγματοληψία με επανάθεση k φορές από τα blocks $\mathfrak{B}_i, i = 1, 2, 3, 4, 5$. Σύμφωνα με τους Davison and Hinkley (1997) οι χρονοσειρές bootstrap που δημιουργούνται μοιάζουν πιο πολύ με λευκό θόρυβο σε σχέση με την αρχική σειρά λόγω των συνδέσεων μεταξύ των blocks.

Μπορούμε να δούμε την γενική διαδικασία του NBB. Για απλότητα θα υποθέσουμε ότι το ℓ είναι ακέραιος, $1 \leq \ell < n$, που διαιρείται ακριβώς με το n και $k = n/\ell$. Στο NBB για να πάρουμε ένα δείγμα bootstrap κάνουμε τυχαία δειγματοληψία με επανάθεση k blocks από την συλλογή $\{\mathfrak{B}_1, \mathfrak{B}_2, \dots, \mathfrak{B}_k\}$, όπου

$$\begin{aligned} \mathfrak{B}_1 &= (x_1, \dots, x_\ell), \\ \mathfrak{B}_2 &= (x_{\ell+1}, \dots, x_{2\ell}), \\ &\dots \\ \mathfrak{B}_k &= (x_{(k-1)\ell+1}, \dots, x_n) \end{aligned}$$

Μετά την πραγματοποίηση της δειγματοληψίας με επανάθεση έχουμε στη διάθεσή μας τη χρονοσειρά bootstrap $\mathfrak{B}_1^*, \dots, \mathfrak{B}_k^*$, όπου κάθε block περιέχει ℓ στοιχεία όταν $n = k\ell$. Τώρα μπορούμε να ορίσουμε την NBB εκδοχή των εκτιμητών της μορφής $\hat{\theta}_n = T(\hat{F}_n)$ όπου \hat{F}_n η

εμπειρική συνάρτηση κατανομής και $T(\cdot)$ είναι ένα συναρτησιακό της \hat{F}_n . Έχοντας ένα δείγμα NBB $\mathfrak{B}_1^*, \dots, \mathfrak{B}_k^*$ η NBB εκδοχή $\hat{\vartheta}_n^*$ του $\hat{\vartheta}_n$ ορίζεται ως

$$\hat{\vartheta}_n^* = T(\hat{F}_n^*)$$

όπου \hat{F}_n^* είναι η εμπειρική κατανομή των (x_1^*, \dots, x_n^*) .

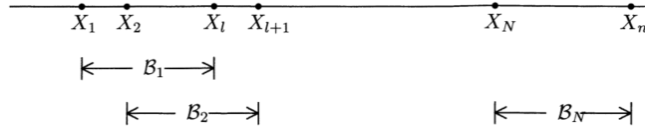
Η παραπάνω διαδικασία επαναλαμβάνεται B φορές αποκτώντας B τιμές $\hat{\vartheta}_n^*$ από τα NBB δείγματα.

3.4 Moving Block Bootstrap (MBB)

Η μέθοδος επαναδειγματοληψίας Moving Block Bootstrap (MBB) είναι μια σημαντική μέθοδος που διατύπωσαν ανεξάρτητα ο Künsch (1989) και οι Liu and Sigh (1992) και η οποία είναι εφαρμόσιμη σε εξαρτημένα δεδομένα χωρίς την απαίτηση κάποιας παραμετρικής υπόθεσης. Όπως είδαμε από την μέθοδο NBB, η οποία να σημειωθεί ότι δημιουργεί λιγότερα blocks από την MBB, η δομή της εξάρτησης των αρχικών δεδομένων περνάει σε κάθε block, και επιπλέον το κοινό μέγεθος των blocks αυξάνει καθώς μεγαλώνει το μέγεθος του δείγματος. Αυτό έχει ως αποτέλεσμα όταν τα αρχικά δεδομένα είναι ασθενώς εξαρτημένα το MBB να αναπαράγει την δομή εξάρτησής τους ασυμπτωτικά. Η γενική ιδέα στο MBB είναι να επιλέξουμε μέγεθος block ℓ αρκετά μεγάλο ώστε οι παρατηρήσεις που απέχουν περισσότερες από ℓ χρονικές μονάδες να είναι σχεδόν ανεξάρτητες και να διατηρείται η συσχέτιση που έχουν παρατηρήσεις που απέχουν λιγότερες από ℓ χρονικές μονάδες (Efron & Tibshirani, 1993). Ας υποθέσουμε ότι X_1, X_2, \dots είναι μια ακολουθία στάσιμων τυχαίων μεταβλητών και $\mathcal{X}_n = \{x_1, x_2, \dots, x_n\}$ είναι οι παρατηρήσεις. Θα προσπαθήσουμε να ορίσουμε την MBB εκδοχή των εκτιμητών της μορφής $\hat{\vartheta}_n = T(\hat{F}_n)$ όπου \hat{F}_n η εμπειρική συνάρτηση κατανομής και $T(\cdot)$ είναι ένα συναρτησιακό της \hat{F}_n . Υποθέτουμε ότι $\ell \equiv \ell_n \in [1, n]$ είναι ένας ακέραιος. Για εξαρτημένα δεδομένα, συνήθως απαιτούμε

$$\ell \rightarrow \infty \text{ και } n^{-1}\ell \rightarrow 0 \text{ καθώς } n \rightarrow \infty$$

δηλαδή το ℓ να αυξάνεται καθώς αυξάνεται το n , αλλά με πιο αργό ρυθμό. Ωστόσο μια περιγραφή του MBB μπορεί να δοθεί και χωρίς αυτούς τους περιορισμούς. Ας συμβολίσουμε το $\mathfrak{B}_i = (x_i, \dots, x_{i+\ell-1})$ το block μεγέθους ℓ που ξεκινάει από το x_i , $1 \leq i \leq N$ όπου $N = n - \ell + 1$



Για να σχηματίσουμε τα MBB δείγματα, επιλέγουμε τυχαία έναν κατάλληλο αριθμό blocks από την συλλογή $\{\mathfrak{B}_1, \mathfrak{B}_2, \dots, \mathfrak{B}_N\}$, όπου

$$\mathfrak{B}_1 = (x_1, x_2, \dots, x_\ell),$$

$$\mathfrak{B}_2 = (x_2, \dots, x_\ell, x_{\ell+1}),$$

....

$$\mathfrak{B}_N = (x_{n-\ell+1}, \dots, x_n)$$

Συνεπώς, πραγματοποιώντας δειγματοληψία με επανάθεση στην συλλογή των blocks παίρνουμε ένα τυχαίο δείγμα bootstrap $\mathfrak{B}_1^*, \dots, \mathfrak{B}_k^*$ όπου κάθε block σημειώνεται ότι περιέχει ℓ στοιχεία. Αν συμβολίσουμε τα στοιχεία ενός block \mathfrak{B}_i^* ως $(x_{(i-1)\ell+1}^*, \dots, x_{i\ell}^*), i = 1, \dots, k$. Τότε οι παρατηρήσεις x_1^*, \dots, x_m^* αποτελούν ένα MBB δείγμα μεγέθους $m \equiv k\ell$. Συνεπώς, η MBB εκδοχή $\hat{\vartheta}_{m,n}^*$ του $\hat{\vartheta}_n$ ορίζεται ως

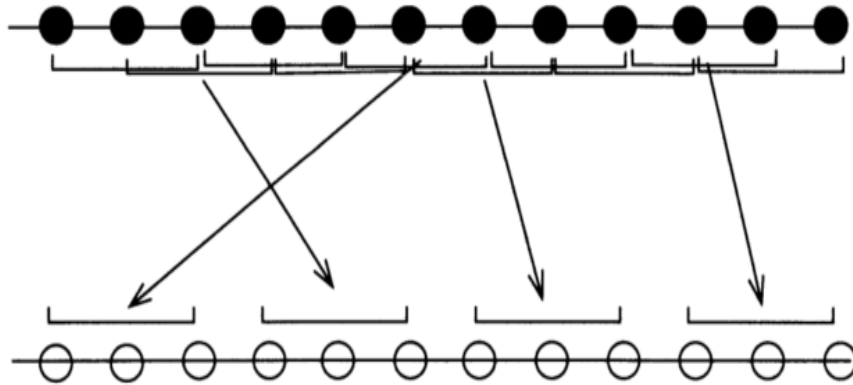
$$\hat{\vartheta}_{m,n}^* = T(\hat{F}_{m,n}^*)$$

όπου $\hat{F}_{m,n}^*$ είναι η εμπειρική κατανομή των (x_1^*, \dots, x_m^*) .

Μερικές τυπικές επιλογές για το ℓ σύμφωνα με τους Kreiss and Lahiri (2012) είναι $\ell = Cn^{1/\delta}$, $\delta = 3, 4$, όπου $C \in \mathbb{R}$ είναι μια σταθερά. Οι Efron and Tibshirani (1993) χρησιμοποίησαν το MBB σε ένα μοντέλο $AR(1)$ με στατιστικό ενδιαφέροντος τον εκτιμημένο συντελεστή $\hat{\varphi}$ του μοντέλου που προέκυψε από τη μέθοδο ελαχίστων τετραγώνων. Πιο συγκεκριμένα, χρησιμοποίησαν MBB ώστε να φτιάξουν B χρονοσειρές bootstrap και σε κάθε μια προσάρμοσαν ένα μοντέλο $AR(1)$ και εκτίμησαν τον συντελεστή του $\hat{\varphi}^*$. Στην αρχή είχαν επιλέξει μέγεθος block $\ell = 3$ και στη συνέχεια $\ell = 5$ και παρατήρησαν ότι με την αύξηση του μεγέθους του block ℓ το τυπικό σφάλμα bootstrap του εκτιμημένου συντελεστή μειωνόταν. Επιπρόσθετα, οι ίδιοι σημείωσαν ότι αν το ℓ δεν διαιρείται ακριβώς με το μέγεθος του δείγματος n τότε πρέπει να πολλαπλασιάσουμε τα τυπικά σφάλματα bootstrap με $(k\ell/n)^{1/2}$ ώστε να προσαρμοστούν στο διαφορετικό μέγεθος των χρονοσειρών. Αν ο n/ℓ δεν είναι ακέραιος αριθμός τότε επιλέγουμε k ως τον μικρότερο ακέραιο που είναι μεγαλύτερος από n/ℓ (Politis & McElroy, 2020). Η χρονοσειρά bootstrap που προκύπτει θα είναι η $x_1^*, \dots, x_{k\ell}^*$.

Παρακάτω στο Σχήμα 3.1 μπορούμε να δούμε την διαδικασία δημιουργίας ενός δείγματος bootstrap με το MBB για $\ell = 3$, $k = 4$ και $n = 12$.

Σχήμα 3.1



Πηγή: Efron & Tibshirani (1993)

Μπορούμε να παρατηρήσουμε ότι τα blocks \mathfrak{B}_i^* 's είναι τυχαία επιλεγμένα από το σύνολο $\{\mathfrak{B}_1, \dots, \mathfrak{B}_N\}$ κάτι το οποίο είναι ισοδύναμο με το να επιλέξουμε k δείκτες τυχαία από το σύνολο $\{1, \dots, N\}$. Συνεπώς, έστω ότι I_1, \dots, I_k iid τυχαίες μεταβλητές με διακριτή ομοιόμορφη κατανομή στο $\{1, \dots, N\}$. Αν θέσουμε $\mathfrak{B}_i^* = \mathfrak{B}_{I_i}^*$ για $i = 1, \dots, k$ τότε το $\mathfrak{B}_1^*, \dots, \mathfrak{B}_k^*$ είναι ένα τυχαίο δείγμα με επανάθεση από το $\{\mathfrak{B}_1, \dots, \mathfrak{B}_N\}$. Το δείγμα bootstrap X_1^*, \dots, X_m^* μπορεί να οριστεί χρησιμοποιώντας τα $\mathfrak{B}_1^*, \dots, \mathfrak{B}_k^*$ όπως πριν. Σημειώνεται ότι δεδομένων των παρατηρήσεων \mathcal{X}_n , τα επαναληπτικά blocks των παρατηρήσεων $(x_1^*, \dots, x_\ell^*)', (x_{(\ell+1)}^*, \dots, x_{2\ell}^*)', \dots, (x_{(k-1)\ell+1}^*, \dots, x_{k\ell}^*)'$ είναι iid ℓ -διάστατα τυχαία διανύσματα που επιλέγονται με πιθανότητα

$$\begin{aligned} P_*((X_1^*, \dots, X_\ell^*)' = (X_j, \dots, X_{j+\ell-1})') \\ &= P_*(I_1 = j) \\ &= N^{-1}, \quad \text{για } 1 \leq j \leq N \end{aligned}$$

όπου P_* συμβολίζει την υπό συνθήκη πιθανότητα δοθέντος \mathcal{X}_n .

Στην ειδική περίπτωση όπου κάθε block αποτελείται από μόνο ένα στοιχείο, δηλαδή $\ell = 1$, τότε σύμφωνα με το παραπάνω προκύπτει ότι τα x_1^*, \dots, x_m^* είναι iid με κοινή κατανομή \hat{F}_n και συνεπώς το MBB γίνεται η κλασική μέθοδος bootstrap του Efron (1979) για iid δεδομένα. Αξίζει να σημειωθεί ότι για $\ell > 1$ η ℓ -διάστατη από κοινού κατανομή της διαδικασίας $\{X_t\}_{t \geq 1}$ διατηρείται μέσα στα επαναδειγματοληπτικά blocks.

Παρ' όλο που οι ορισμοί για τους εκτιμητές bootstrap στο NBB και MBB μοιάζουν πολύ, οι εκτιμητές bootstrap που προκύπτουν $\hat{\vartheta}_n^* = T(\hat{F}_n^*)$ και $\hat{\vartheta}_{m,n}^* = T(\hat{F}_{m,n}^*)$ έχουν πολύ διαφορετικές ιδιότητες ως προς τις κατανομές τους. Παρακάτω θα το αναλύσουμε αυτό στην απλούστερη

περίπτωση που $\hat{\vartheta}_n = n^{-1} \sum_{j=1}^n X_j$ είναι ο δειγματικός μέσος. Η εκδοχή bootstrap του $\hat{\vartheta}_n$ κάτω από αυτές τις δύο μεθόδους είναι $\hat{\vartheta}_n^*$ για το NBB και $\hat{\vartheta}_{m,n}^*$ για το MBB.

$$\hat{\vartheta}_{m,n}^* = m^{-1} \sum_{j=1}^m X_j^* \quad \text{και} \quad \hat{\vartheta}_n^* = n^{-1} \sum_{j=1}^n X_j^*$$

Η πιθανότητα να επιλεγεί ένα block στην MBB είναι N^{-1} , συνεπώς

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

Για να πάρουμε παρόμοια έκφραση για το $E_*(\hat{\vartheta}_n^*)$, σημειώνουμε ότι για το NBB, οι μεταβλητές bootstrap $(X_1^*, \dots, X_{\ell}^*), \dots, (X_{(n-\ell+1)}^*, \dots, X_n^*)$ είναι iid, με

$$P_*((X_1^*, \dots, X_{\ell}^*) = (X_{(j-1)\ell+1}, \dots, X_{j\ell})) = 1/k \quad \text{για} \quad j = 1, \dots, k$$

δηλαδή, η πιθανότητα να επιλεγεί ένα block στην NBB είναι k^{-1} . Επομένως,

$$\begin{aligned} E_*(\hat{\vartheta}_n^*) &= E_*(\ell^{-1} \sum_{i=1}^{\ell} X_i^*) = k^{-1} \sum_{i=1}^k (\ell^{-1} \sum_{i=1}^{\ell} X_{(j-1)\ell+i}) \\ &= (k\ell)^{-1} \{n\bar{X}_n - \sum_{i=k\ell+1}^n X_i\} \end{aligned}$$

Άρα προκύπτει ότι οι δύο εκτιμητές bootstrap έχουν διαφορετικούς μέσους. Ωστόσο, σύμφωνα με τον Lahiri (2003) αν η διαδικασία $\{X_t\}_{t \geq 1}$ ικανοποιεί κάποιες συνθήκες τότε $E\{E_*(\hat{\vartheta}_{m,n}^*) - E_*(\hat{\vartheta}_n^*)\}^2 = O(\ell/n^2)$. Επομένως, η διαφορά μεταξύ των δύο είναι αμελητέα για μεγάλα δείγματα.

3.5 Circular Block Bootstrap (CBB)

Το Circular Block Bootstrap (CBB) των Politis and Romano (1992) έχει ως κύριο σκοπό να βελτιώσει το πρόβλημα της μεθόδου MBB η οποία αποδίδει μικρότερα βάρη στις παρατηρήσεις που βρίσκονται προς την αρχή και το τέλος του συνόλου δεδομένων και μεγάλα βάρη στις κεντρικές παρατηρήσεις. Στο MBB για $\ell \leq j \leq n - \ell + 1$ η παρατήρηση x_j εμφανίζεται σε ακριβώς ℓ blocks από τα $\{\mathfrak{B}_1, \mathfrak{B}_2, \dots, \mathfrak{B}_N\}$, όπου $N = n - \ell + 1$ είναι ο αριθμός

των επικαλυπτόμενων blocks ενώ για $1 \leq j \leq \ell - 1$, η x_j και η x_{n-j+1} εμφανίζεται μόνο σε j blocks. Η ιδέα των Politis and Romano (1992) για το CBB είναι να τοποθετήσουν τα δεδομένα σε κυκλική διάταξη και να σχηματίσουν πρόσθετα blocks. Ως αποτέλεσμα αντί για $n - \ell + 1$ blocks, μπορούν να οριστούν n διαφορετικά blocks $\{\mathfrak{B}_1, \mathfrak{B}_2, \dots, \mathfrak{B}_n\}$ όπου για $t > n - \ell + 1$ το block \mathfrak{B}_t θα είναι της μορφής $(x_t, \dots, x_n, x_1, \dots, x_{\ell-n+t-1})$. Αυτό έχει σαν αποτέλεσμα η εκτίμηση bootstrap του δειγματικού μέσου να είναι σωστά κεντραρισμένη (Politis & McElroy, 2020). Αξίζει να σημειωθεί ότι στην περίπτωση στάσιμης χρονοσειράς, το CBB παρέχει συχνά αμερόληπτους εκτιμητές bootstrap (Dudek, 2017).

3.6 Stationary Bootstrap (SB)

Το Stationary Bootstrap (SB) των Politis and Romano (1994) διαφέρει από τις μεθόδους που αναφέρθηκαν νωρίτερα με την έννοια ότι το μέγεθος των blocks δεν είναι σταθερό αλλά είναι μια τυχαία μεταβλητή που ακολουθεί γεωμετρική κατανομή με αναμενόμενη τιμή ℓ . Λόγω του τυχαίου μεγέθους των blocks, το πλήθος των blocks στο SB είναι επίσης τυχαίο. Όπως και το CBB, το SB έχει σκοπό να εξαλείψει την ανισορροπία του αριθμού εμφάνισης των παρατηρήσεων στα blocks.

Η γενικότερη δυσκολία με τα σχήματα των blocks είναι ότι οι χρονοσειρές bootstrap που δημιουργούνται δεν είναι στάσιμες καθώς η από κοινού κατανομή των παρατηρήσεων που έχουν προκύψει με επανάθεση κοντά σε μια ένωση μεταξύ των blocks διαφέρει από την από κοινού κατανομή των παρατηρήσεων στο κέντρο ενός block. Αυτό το πρόβλημα μπορεί να ξεπεραστεί παίρνοντας blocks τυχαίου μεγέθους. Σε αντίθεση με τις προηγούμενες μεθόδους, στο stationary bootstrap η επαναδειγματοληπτική διαδικασία επαναλαμβάνεται δημιουργώντας στάσιμες χρονοσειρές bootstrap από τις οποίες μπορούμε να πάρουμε μια προσέγγιση της κατανομής του στατιστικού ενδιαφέροντος. Συνεπώς το SB στοχεύει στο να διατηρήσει την στασιμότητα της αρχικής χρονοσειράς στις χρονοσειρές bootstrap που δημιουργεί.

Το stationary bootstrap παίρνει blocks που το μεγέθός τους ακολουθεί γεωμετρική κατανομή με συνάρτηση πιθανότητας

$$P(L = j) = p(1 - p)^{j-1}, \quad j = 1, 2, \dots$$

όπου $p \in (0, 1)$. Αυτό οδηγεί σε χρονοσειρές bootstrap που είναι στάσιμες με μέσο μήκος block $\ell = p^{-1}$. Όσο μικραίνει η πιθανότητα p το μήκος των blocks τείνει να αυξάνεται ενώ όσο αυξάνει η πιθανότητα το μήκος των blocks τείνει να μειώνεται. Η εγκυρότητα του SB εξαρτάται από την επιλογή του p έτσι ώστε $p \rightarrow 0$ αλλά $pn \rightarrow \infty$ καθώς $n \rightarrow \infty$ (Politis & McElroy, 2020).

Ας υποθέσουμε ότι μας δίνεται μια ασθενώς στάσιμη χρονοσειρά και ας συμβολίσουμε το $\mathfrak{B}_{i,\ell} = (x_i, \dots, x_{i+\ell-1})$ το block μεγέθους ℓ που ξεκινάει από το x_i , όπου $x_i = x_{1+(i-1 \bmod n)}$ και $x_0 = x_n$. Επίσης, έστω I_1, I_2, \dots είναι μια ακολουθία iid τυχαίων μεταβλητών από την διακριτή ομοιόμορφη κατανομή στο $\{1, \dots, n\}$, που I_i υποδηλώνει τον δείκτη της τιμής εκκίνησης ενός block i , και L_1, \dots μια ακολουθία iid τυχαίων μεταβλητών από την γεωμετρική κατανομή με μέση τιμή $\ell = p^{-1}$, με L_i να υποδηλώνει το μέγεθος του block i . Το SB δημιουργεί μια χρονοσειρά bootstrap κάνοντας επαναλαμβανόμενη δειγματοληψία blocks τυχαίου μεγέθους $\mathfrak{B}_{I_1, L_1}, \mathfrak{B}_{I_2, L_2}, \dots$ και τοποθετώντας τα ως μια ακολουθία μέχρι ο αριθμός των παρατηρήσεων της χρονοσειράς bootstrap να είναι τουλάχιστον n . Συνεπώς, εδώ οι πρώτες L_1 σε αριθμό παρατηρήσεις της χρονοσειράς bootstrap καθορίζονται από το block $\mathfrak{B}_{I_1, L_1} = (x_{I_1}, \dots, x_{I_1+L_1-1})$. Αν το μήκος της ακολουθίας που δημιουργείται από την μέθοδο ξεπεράσει το μέγεθος δείγματος n τότε σταματάμε την διαδικασία.

Οι τέσσερις μέθοδοι blocks που έχουμε δει έως τώρα, δηλαδή οι NBB, MBB, CBB και SB συγκρίθηκαν θεωρητικά από τον Lahiri (1999), ο οποίος έδειξε ότι και οι τέσσερις μέθοδοι έχουν ίδια μεροληψία ασυμπτωτικά, αλλά οι διακυμάνσεις των NBB και SB είναι μεγαλύτερες, με τις διακυμάνσεις των εκτιμητών SB να είναι οι μεγαλύτερες από όλες. Πιο συγκεκριμένα, έδειξε ότι οι εκτιμητές που προκύπτουν από τα MBB και CBB έχουν ίδιες διακυμάνσεις κάτι που τους καθιστά ασυμπτωτικά ισοδύναμους ως προς το MSE και ότι υπερτερούν των NBB και SB που έχουν μεγαλύτερη διακύμανση στους εκτιμητές τους και επομένως είναι λιγότερο αποδοτικές μέθοδοι. Επιπρόσθετα, σύμφωνα με τον ίδιο, αν επιλεγθεί μέγεθος block ℓ τέτοιο ώστε οι μέθοδοι NBB, MBB και CBB να έχουν μεροληψία και διακύμανση ίδιας τάξης τότε το μέσο τετραγωνικό σφάλμα MSE του εκτιμητή NBB εξακολουθεί να είναι μεγαλύτερο από το MSE των άλλων δύο μεθόδων καθιστώντας την NBB λιγότερο αποτελεσματική. Οι διακυμάνσεις των SB εκτιμητών είναι πάντα τουλάχιστον δύο φορές μεγαλύτερες από τις διακυμάνσεις των εκτιμητών NBB και τουλάχιστον τρεις φορές μεγαλύτερες από τις διακυμάνσεις των MBB και CBB εκτιμητών. Επίσης πραγματοποίησε σύγκριση των μεθόδων ως προς τις καλύτερες δυνατές αποδόσεις τους και οι μέθοδοι MBB και CBB υπερέχουν των μεθόδων NBB και SB, με την μέθοδο NBB να υπερέχει της SB.

Αξίζει να σημειωθεί ότι ο Lahiri (2003) διατύπωσε το Generalized Block Bootstrap (GBB), μια γενίκευση της ιδέας των Politis and Romano (1992, 1994) η οποία περιέχει ένα ενοποιημένο πλαίσιο για να περιγράψει διάφορες μεθόδους block bootstrap συμπεριλαμβανομένων και των CBB και SB.

3.7 Post-blackening Bootstrap

Ένα αρκετά σημαντικό πρόβλημα που προκύπτει από τις μεθόδους των blocks και μερικές φορές οδηγούν σε λανθασμένες προσεγγίσεις επαναδειγματοληψίας είναι ότι δημιουργούνται χρονοσειρές bootstrap λιγότερο εξαρτημένες από τα αρχικά δεδομένα. Ο αριθμός και το μέγεθος των blocks λόγω του μήκους της αρχικής χρονοσειράς μπορεί να μην είναι αρκετός ώστε να διατηρήσει την δομή εξάρτησης των αρχικών παρατηρήσεων. Αυτό έχει ως συνέπεια να δημιουργούνται χρονοσειρές bootstrap που οι παρατηρήσεις είναι λιγότερο εξαρτημένες από τις παρατηρήσεις της αρχικής χρονοσειράς. Για παράδειγμα, αν χρησιμοποιήσουμε μέγεθος block $l = 1$, δηλαδή είμαστε στο κλασικό bootstrap του Efron (1979), η χρονοσειρά bootstrap που θα προκύψει θα είναι λευκός θόρυβος καθώς οι παρατηρήσεις θα είναι τυχαία επιλεγμένες με επανάθεση χωρίς καμία δομή εξάρτησης μεταξύ τους. Επιπλέον, η χρονοσειρά bootstrap συχνά εμφανίζει σφάλματα που προκύπτουν από την ένωση τυχαία επιλεγμένων blocks. Ο Chernick (2007) αναφέρει ότι κάποια από τα μειονεκτήματα των block μεθόδων είναι, ότι τα επαναδειγματοληπτικά blocks δεν μιμούνται ακριβώς τη συμπεριφορά της χρονοσειράς και ότι έχουν την τάση να αποδυναμώνουν την εξάρτηση στη σειρά. Για την αντιμετώπιση αυτών των προβλημάτων προέκυψε μια στρατηγική ενδιάμεση του bootstrap βάσει μοντέλου και του block bootstrap από τους Davison and Hinkley (1997). Η ιδέα προέρχεται από μια διαδικασία “προ λεύκανσης” (pre-whitening) της αρχικής εξαρτημένης χρονοσειράς με την προσαρμογή ενός μοντέλου που έχει σκοπό να αφαιρέσει μεγάλο μέρος της εξάρτησης μεταξύ των αρχικών παρατηρήσεων. Μια νέα χρονοσειρά bootstrap δημιουργείται αφού πρώτα έχει πραγματοποιηθεί δειγματοληψία με επανάθεση από blocks των καταλοίπων του προσαρμοσμένου μοντέλου. Δηλαδή, αφού έχουμε μια σειρά καταλοίπων που προέκυψε από δειγματοληψία με επανάθεση στα blocks των καταλοίπων, εφαρμόζουμε αναδρομικά το εκτιμημένο μοντέλο χρησιμοποιώντας τα κατάλοιπα και έτσι προκύπτει η χρονοσειρά bootstrap. Για παράδειγμα αν ένα $AR(1)$ μοντέλο είχε χρησιμοποιηθεί για “προ λεύκανση” των δεδομένων, μια νέα χρονοσειρά προκύπτει, αφού βέβαια πρώτα έχει πραγματοποιηθεί η επαναδειγματοληψία των καταλοίπων από τα blocks των καταλοίπων του προσαρμοσμένου μοντέλου, από $Y_t^* = \hat{\varphi}_1 Y_{t-1}^* + e_t^*$, $t = 2, \dots, n$ με $Y_1^* = y_1$. Αυτή η post-blackened προσέγγιση φαίνεται να λειτουργεί πιο σταθερά στην πράξη σύμφωνα με τους Davison and Hinkley (1997).

3.8 Matched Block Bootstrap (MaBB) και Tapered Block Bootstrap (TBB)

Όπως είδαμε, στις μεθόδους block bootstrap εφαρμόζουμε δειγματοληψία με επανάθεση σε blocks δεδομένων της αρχικής χρονοσειράς με αποτέλεσμα να προκύψουν προσομοιωμένες εκδοχές της αρχικής χρονοσειράς όταν ενώσουμε τα blocks. Τα blocks δεδομένων που στην αρχική χρονοσειρά είναι εξαρτημένα, στην χρονοσειρά bootstrap είναι ανεξάρτητα κάτι το οποίο προκαλεί μεροληψία στη διακύμανση ενός εκτιμητή bootstrap, ειδικότερα αν η εξάρτηση

στα αρχικά δεδομένα είναι ισχυρή. Το σκεπτικό στο Matched Block Bootstrap ώστε να βελτιωθεί η απόδοση των εκτιμητών bootstrap είναι να αντιστοιχίσουμε τα blocks μεταξύ τους με κάποιον τρόπο, δηλαδή να χρησιμοποιήσουμε έναν κανόνα μετάβασης των blocks ο οποίος θα ευνοούσε τα blocks που θα ήταν εκ των προτέρων πιο πιθανό να είναι κοντά με κάποια άλλα, ώστε να βελτιωθεί η απόδοση των εκτιμητών bootstrap. Σύμφωνα με τους Carlstein et al. (1998), στο Matched Block Bootstrap (MaBB) η επαναδειγματοληψία των blocks πραγματοποιείται χρησιμοποιώντας μια αλυσίδα Markov και αυτό έχει σαν αποτέλεσμα τα επαναδειγματοληπτικά blocks να είναι εξαρτημένα. Οι Carlstein et al. (1998) δείχνουν ότι ο εκτιμητής MaBB της διακύμανσης του δειγματικού μέσου έχει τιμή η οποία είναι συγκρίσιμης τάξης με την διακύμανση του NBB εκτιμητή και έχει μεροληψία μικρότερης τάξης από αυτή του NBB.

Όπως είπαμε, το MaBB δημιουργεί μια ακολουθία από blocks επιλέγοντας blocks που είναι εκ των προτέρων πιο πιθανό να είναι κοντά το ένα με το άλλο. Τα blocks αρχικά κατασκευάζονται χρησιμοποιώντας μια από τις μεθόδους block bootstrap που περιγράψαμε στις προηγούμενες ενότητες και αργότερα παίρνουμε δείγμα χρησιμοποιώντας μια μαρκοβιανή αλυσίδα η οποία ευνοεί τα blocks των οποίων τα άκρα τους ταιριάζουν με το τέλος του προηγούμενου block (rank matching). Για την MaBB ας υποθέσουμε ότι το block i μεγέθους ℓ είναι $\mathfrak{B}_i = (x_{i,1}, \dots, x_{i,\ell})$, και $x_{i,j}$ υποδηλώνει την j -παρατήρηση του i -οστού block. Αν χρησιμοποιήσουμε την μέθοδο επικαλυπτόμενων blocks MBB, δημιουργούνται $N = n - \ell + 1$ διαφορετικά blocks. Συμβολίζουμε με R_i το rank του τέλους του block \mathfrak{B}_i , δηλαδή της τελευταίας παρατήρησης $x_{i,\ell}$ μεταξύ των τελευταίων παρατηρήσεων όλων των blocks, δηλαδή των $x_{1,\ell}, \dots, x_{N,\ell}$. Αρχικά θα πραγματοποιηθεί δειγματοληψία ενός block \mathfrak{B}_1 από την συλλογή των blocks. Στην συνέχεια ένα ακόμα block j θα επιλεγεί τυχαία από τα $2\kappa + 1$ blocks που κατατάσσονται μεταξύ $R_1 - \kappa$ και $R_1 + \kappa$, όπου R_1 είναι το rank του block \mathfrak{B}_1 , δηλαδή της τελευταίας του παρατήρησης, και κ είναι ένας μικρός θετικός ακέραιος αριθμός. Το επόμενο block που προστίθεται στην χρονοσειρά bootstrap είναι τώρα το block που ακολουθεί το block j στην αρχική πραγματική χρονοσειρά. Η παραπάνω διαδικασία επαναλαμβάνεται έως ότου η χρονοσειρά bootstrap φτάσει το επιθυμητό μήκος. Με αυτόν τον τρόπο έχουμε στη διάθεση μας μία MaBB χρονοσειρά από την οποία μπορούμε να υπολογίσουμε το στατιστικό που επιθυμούμε. Η διαδικασία επαναλαμβάνεται B φορές.

Το Tapered Block Bootstrap (TBB) των Paparoditis and Politis (2001) είναι μια παραλλαγή της μεθόδου block bootstrap που περιλαμβάνει την εφαρμογή της τεχνικής του tapering στα blocks. Η τεχνική εφαρμόζει ένα βάρος σε κάθε τιμή του block. Πιο συγκεκριμένα, συρρικνώνει τις οριακές τιμές σε ένα block προς μια κοινή τιμή όπως ο δειγματικός μέσος. Γενικά σε αυτή τη μέθοδο οι παρατηρήσεις των blocks που βρίσκονται στις άκρες, δηλαδή τα τελικά σημεία των blocks, συρρικνώνονται προς μια τιμή-στόχο πριν συνδεθούν τα blocks και

σχηματίσουν μια χρονοσειρά bootstrap (Politis, 2003). Συνεπώς, έχοντας για παράδειγμα επικαλυπτόμενα blocks MBB, εφαρμόζουμε την αλλαγή κλίμακας με τη μέθοδο tapering και κάνουμε δειγματοληψία με επανάθεση από τα tapered blocks. Αν και αυτές οι μέθοδοι είναι κάπως πιο περίπλοκες από αυτές που έχουμε δει μέχρι τώρα, το MaBB και το TBB οδηγούν σε πιο ακριβείς εκτιμητές bootstrap της διακύμανσης ενός εκτιμητή (Kreiss & Lahiri, 2012).

3.9 Επιλογή μεγέθους block

Η επιτυχία της μεθόδου block bootstrap εξαρτάται σε μεγάλο βαθμό από την κατάλληλη επιλογή του μεγέθους block ℓ το οποίο λειτουργεί ως παράμετρος ρύθμισης (tuning parameter). Το ιδανικό μέγεθος block εξαρτάται από το μέγεθος του αρχικού δείγματος και την δομή της συσχέτισης και διαφέρει για διαφορετικές μεθόδους block bootstrap. Η χρήση πολύ μικρής τιμής ℓ θα καταστρέψει την δομή εξάρτησης με αποτέλεσμα να αυξηθεί η μεροληψία της μεθόδου bootstrap. Αν το ℓ είναι πολύ μεγάλο, π.χ. κοντά στο μέγεθος του δείγματος, όλα τα στατιστικά bootstrap θα είναι σχεδόν ίσα με τον εκτιμητή και ως αποτέλεσμα η κατανομή των εκτιμητών bootstrap θα είναι πολύ συγκεντρωμένη γύρω από την τιμή του.

Σύμφωνα με τον Bühlmann (2002), ένα βέλτιστο μήκος block ℓ εξαρτάται από τουλάχιστον τρία πράγματα: την διαδικασία που παράγει τα δεδομένα, το στατιστικό για το οποίο θα εφαρμόσουμε bootstrap, και τον σκοπό για τον οποίο χρησιμοποιούμε bootstrap (για παράδειγμα bootstrap για εκτίμηση μεροληψίας, διακύμανσης ή κατανομής του στατιστικού ενδιαφέροντος). Η ασυμπτωτική θεωρία μάς λέει μόνο ότι το μέγεθος ℓ πρέπει να αυξάνεται καθώς αυξάνεται το μέγεθος του αρχικού δείγματος και ότι το βέλτιστο μέγεθος block για εκτίμηση τυπικών σφαλμάτων bootstrap είναι της τάξεως $Cn^{1/3}$, όπου η σταθερά C εξαρτάται από το στατιστικό ενδιαφέροντος και την εξάρτηση μεταξύ των παρατηρήσεων. Σύμφωνα με την ασυμπτωτική θεωρία οι Hall et al. (1995) αναφέρουν γενικά ότι $\ell \sim Cn^{1/k}$ με $k = 3, 4, 5$, ανάλογα με το πλαίσιο, δηλαδή εκτίμηση μεροληψίας ή διασποράς εκτιμητή ή της κατανομής της στατιστικής συνάρτησης για μονόπλευρο και δίπλευρο έλεγχο αντίστοιχα. Σύμφωνα με τον Politis (2003), υπάρχουν δύο κύριες προσεγγίσεις για την κατάλληλη επιλογή του μεγέθους των blocks. Η πρώτη είναι μέσω μιας προσέγγισης υποδειγματοληψίας/διασταυρούμενης επικύρωσης και η δεύτερη χρησιμοποιώντας μεθόδους αντικατάστασης. Πιο συγκεκριμένα, μια προσέγγιση υποδειγματοληψίας/διασταυρούμενης επικύρωσης για το μέγεθος του block προτάθηκε από τους Hall et al. (1995). Οι Hall et al. (1995) πρότειναν ότι το βέλτιστο μέγεθος block να είναι $n^{1/3}$, $n^{1/4}$, $n^{1/5}$ ανάλογα από το πλαίσιο, με το $n^{1/3}$ να είναι το βέλτιστο μέγεθος block για διακύμανση και μεροληψία.

Η προσέγγιση της μεθόδου της αντικατάστασης περιλαμβάνει τον υπολογισμό μίας έκφρασης για την βέλτιστη βάσει κάποιου κριτηρίου τιμή του ℓ και στη συνέχεια την εκτίμηση και αντικατάσταση όλων των άγνωστων παραμέτρων σε αυτήν την έκφραση. Μια εναλλακτική μέθοδος βασιζόμενη στο Jackknife after Bootstrap (JAB method) προτάθηκε από τους Lahiri et al. (2007). Την ονόμασαν μη παραμετρική μέθοδο αντικατάστασης (nonparametric plug-in, NPPI), καθώς δουλεύει όπως η μέθοδος της αντικατάστασης (plug-in method) αλλά ταυτόχρονα δεν απαιτείται από τον χρήστη να βρει αναλυτικά μια ακριβή έκφραση για το βέλτιστο μέγεθος block. Η μέθοδος NPPI είναι εφαρμόσιμη σε προβλήματα εκτίμησης block bootstrap συμπεριλαμβανομένων της διακύμανσης, της συνάρτησης κατανομής και των ποσοστημόριων. Ωστόσο, είναι σημαντικό να επισημάνουμε ότι πρόκειται για μια υπολογιστικά επίπονη μέθοδο (computer-intensive method) καθώς χρησιμοποιεί έναν συνδυασμό bootstrap και jackknife. Για περισσότερες πληροφορίες σχετικά με την μέθοδο, ανατρέξτε στους Lahiri et al. (2007), Lahiri (2003).

Σύμφωνα με τον Carlstein (1986), όταν αυξάνεται το μέγεθος του block ℓ , η μεροληψία των εκτιμητών μειώνεται και επίσης όταν η εξάρτηση μεταξύ των τυχαίων μεταβλητών γίνεται πιο ισχυρή τότε χρειαζόμαστε μεγαλύτερο ℓ . Βασισμένος σε αυτούς τους ισχυρισμούς κατέληξε ότι το βέλτιστο μέγεθος block για το μοντέλο $AR(1)$, $Y_t = \varphi Y_{t-1} + e_t$, είναι $\ell^* = (2\varphi/(1 - \varphi^2))^{2/3} n^{1/3}$. Να σημειωθεί ότι το ℓ^* είναι αύξουσα συνάρτηση του φ . Για παράδειγμα, για μέγεθος δείγματος $n = 200$ και $\varphi = .5, .8$ και $.9$ παίρνουμε αντίστοιχα $\ell^* = 7.08, 15.81$ και 26.18 . Στην πράξη αυτοί οι αριθμοί στρογγυλοποιούνται.

3.10 AR-Sieve Bootstrap

Μια ακόμα μέθοδος με την οποία μπορούμε να εφαρμόσουμε bootstrap σε μια στάσιμη χρονοσειρά είναι η AR-sieve bootstrap. Αυτή η μέθοδος χρησιμοποιήθηκε από τους Swanepoel and Van Wyk (1986) και διερευνήθηκε περαιτέρω από τους Kreiss (1988, 1992), Paparoditis (1992) και Bühlmann (1997). Ας υποθέσουμε ότι P είναι η άγνωστη από κοινού κατανομή πιθανότητας για μια απείρως μεγάλη ακολουθία $\{X_1, X_2, \dots, X_n, \dots\}$. Στην ειδική περίπτωση των iid δεδομένων, χρησιμοποιούμε την εμπειρική συνάρτηση κατανομής \hat{F}_n για τους πρώτους n όρους της που θα παρατηρήσουμε. Δεδομένου ότι \hat{F}_n είναι η μονοδιάστατη εκτίμηση της πραγματικής κατανομής F που είναι κοινή για κάθε X_i , η από κοινού κατανομή των πρώτων n παρατηρήσεων είναι λόγω ανεξαρτησίας το γινόμενο F^n . Μια φυσική εκτίμηση της από κοινού κατανομής θα ήταν η $(\hat{F}_n)^n$. Ωστόσο, στην περίπτωση των εξαρτημένων δεδομένων, η από κοινού κατανομή δεν είναι το γινόμενο των μονοδιάστατων περιθωρίων κατανομών. Η ιδέα του sieve είναι να προσεγγίσει αυτήν την από κοινού κατανομή από κατανομές που πλησιάζουν

την P καθώς το μέγεθος n αυξάνει. Σε αυτήν την περίπτωση εξετάζουμε μια μεγάλη κατηγορία στάσιμων μοντέλων χρονοσειρών που μπορεί να αναπαρασταθεί ως αυτοπαλίνδρομη διαδικασία άπειρης τάξης. Σημαντικό να αναφερθεί ότι ο Bühlmann (1997) χρησιμοποίησε την μέθοδο sieve προσεγγίζοντας την P από την από κοινού κατανομή ενός πεπερασμένου μοντέλου $AR(p)$. Για κάθε τέτοιο στάσιμο μοντέλο, η προσέγγιση πλησιάζει την P καθώς $n \rightarrow \infty$.

Γενικότερα, στο πλαίσιο των στάσιμων χρονοσειρών, έχουν προταθεί δύο διαφορετικές προσεγγίσεις bootstrap, τις οποίες έχουμε ήδη δει. Η μια είναι η προσέγγιση βάσει μοντέλου στην οποία η επαναδειγματοληψία πραγματοποιείται στα σχεδόν iid κατάλοιπα και η άλλη είναι η ελεύθερη μοντέλου προσέγγιση όπου η δειγματοληψία πραγματοποιείται από blocks παρατηρήσεων. Το sieve bootstrap χρησιμοποιεί την ιδέα του βάσει μοντέλου bootstrap, δηλαδή της προσαρμογής ενός παραμετρικού μοντέλου και επαναδειγματοληψίας στα κατάλοιπα, αλλά αντί να θεωρεί ένα σταθερό μοντέλο πεπερασμένης διάστασης, προσεγγίζει ένα μη παραμετρικό μοντέλο άπειρης διάστασης από μια ακολουθία παραμετρικών μοντέλων πεπερασμένης διάστασης. Αυτή η προσέγγιση στην πράξη πραγματοποιείται όταν επιλέγουμε ένα μοντέλο χρησιμοποιώντας κάποιο κριτήριο όπως είναι το Akaike Information Criterion το οποίο μας βοηθάει να επιλέξουμε την τάξη p ενός αυτοπαλίνδρομου μοντέλου, δηλαδή προσαρμόζοντας το καλύτερο σύμφωνα με το κριτήριο αυτοπαλίνδρομο μοντέλο p τάξης. Επιπρόσθετα, το σημαντικό πλεονέκτημα αυτής της μεθόδου είναι ότι δεν απαιτεί την πρότερη γνώση της τάξης του AR μοντέλου αλλά αυτή επιλέγεται σύμφωνα με το κριτήριο. Αξίζει να σημειωθεί ότι οι μέθοδοι sieve bootstrap βασίζονται στην ιδέα της sieve προσέγγισης του Grenander (1981).

Μία χρονοσειρά λέμε ότι είναι γραμμική και αντιστρέψιμη αν μπορεί να αναπαρασταθεί ως αυτοπαλίνδρομη διαδικασία άπειρης τάξης $AR(\infty)$. Η προσέγγιση του AR -sieve bootstrap υποθέτει ότι η χρονοσειρά έχει την μορφή $AR(\infty)$

$$X_t - \mu_X = \sum_{j=1}^{\infty} \varphi_j (X_{t-j} - \mu_X) + e_t, \quad t \in \mathbb{Z}$$

όπου $\mu_X = E(X_t)$, $\{e_t\}_{t \in \mathbb{Z}}$ η ακολουθία των σφαλμάτων iid τυχαίων μεταβλητών μηδενικού μέσου $E(e_t) = 0$. Ωστόσο, η προσαρμογή του μοντέλου πρέπει να πραγματοποιηθεί με χρήση μιας πεπερασμένης τάξης p . Με αυτήν την αναπαράσταση, ο μηχανισμός δημιουργίας δεδομένων X_t μπορεί να προσεγγιστεί από τα σφάλματα e_t . Αυτό επιτυγχάνεται προσαρμόζοντας ένα κατάλληλο μεγάλης τάξης $AR(p)$, με αυξανόμενης τάξης $p(n)$ όσο το μέγεθος του δείγματος n αυξάνει, $p = p(n) \rightarrow \infty$ με $p(n) = o(n)$ ($n \rightarrow \infty$), δηλαδή το p αυξάνει συναρτήσει του μεγέθους του δείγματος αλλά πιο αργά από αυτό. Η προσέγγιση AR -sieve πραγματοποιείται με $AR(p)$ μοντέλα όπως έχει αναφερθεί.

$$X_t - \mu_X = \sum_{j=1}^p \varphi_j (X_{t-j} - \mu_X) + e_t, \quad t \in \mathbb{Z}$$

όπου $\mu_X = E(X_t)$, $\{e_t\}_{t \in \mathbb{Z}}$ η ακολουθία των σφαλμάτων iid τυχαίων μεταβλητών μηδενικού μέσου $E(e_t) = 0$. Έχοντας στη διάθεση μας τα δεδομένα, επιλέγουμε πρώτα την τάξη \hat{p} του αυτοπαλίνδρομου μοντέλου με κάποιο κριτήριο όπως είναι το Akaike Information Criterion (AIC). Εδώ οι παράμετροι ενδιαφέροντος είναι $\eta_{\hat{p}} = (\mu_X, (\varphi_1, \dots, \varphi_{\hat{p}}), F_e)$. Με το F_e συμβολίζουμε την κατανομή των iid σφαλμάτων e_t . Οι εκτιμήσεις αυτών των παραμέτρων δίνονται ακολούθως

$$\hat{\mu}_X = \bar{X}_n = n^{-1} \sum_{t=1}^n X_t$$

Οι εκτιμημένοι συντελεστές $\hat{\varphi}_{\hat{p}} = (\hat{\varphi}_{1,n}, \dots, \hat{\varphi}_{\hat{p},n})$, συνήθως αλλά όχι απαραίτητα, υπολογίζονται από τη μέθοδο Yule-Walker, όπου αν $\hat{\varphi}_{\hat{p}}$ οι εκτιμημένοι συντελεστές με την μέθοδο Yule-Walker τότε

$$\hat{\varphi}_{\hat{p}} = \hat{\Gamma}_{\hat{p}}^{-1} \hat{\gamma}_{\hat{p}}$$

$$\hat{\gamma}_h = \frac{1}{n} \sum_{t=1}^{n-|h|} (X_t - \bar{X}_n)(X_{t-|h|} - \bar{X}_n), \quad 0 \leq h \leq \hat{p}$$

όπου $\bar{X}_n = n^{-1} \sum_{t=1}^n X_t$, $\hat{\Gamma}_{\hat{p}} = \hat{\gamma}(|r-s|)$ για $r, s = 1, \dots, \hat{p}$ και $\hat{\gamma}_{\hat{p}} = (\hat{\gamma}_1, \dots, \hat{\gamma}_{\hat{p}})'$

και

$$\hat{F}_e = \hat{P}(e_t \leq x) = (n - \hat{p})^{-1} \sum_{t=\hat{p}+1}^n \mathbb{1}_{[R_t - \bar{R} \leq x]}$$

με $R_t = \hat{e}_t = X_t - \sum_{j=1}^{\hat{p}} \hat{\varphi}_j X_{t-j}$ να είναι τα κατάλοιπα και $\bar{R} = \bar{\hat{e}}_t = (n - \hat{p})^{-1} \sum_{t=\hat{p}+1}^n \hat{e}_t$ ο μέσος των διαθέσιμων καταλοίπων R_t , ($t = \hat{p} + 1, \dots, n$) και συνεπώς από $R_t - \bar{R}$, ($t = \hat{p} + 1, \dots, n$) να προκύπτουν τα κεντραρισμένα κατάλοιπα \tilde{e}_t , ($t = \hat{p} + 1, \dots, n$).

Στη συνέχεια για το AR-sieve bootstrap, δημιουργούμε χρονοσειρές bootstrap πραγματοποιώντας δειγματοληψία με επανάθεση στα κεντραρισμένα κατάλοιπα $\{\tilde{e}_t\}_{t=\hat{p}+1}^n$, δηλαδή από την εμπειρική συνάρτηση κατανομής \hat{F}_e , και έπειτα χρησιμοποιείται το μοντέλο AR για να παραχθεί η νέα χρονοσειρά μέσω αναδρομής. Δηλαδή, οι χρονοσειρές bootstrap είναι βασισμένες στο μοντέλο $AR(\hat{p})$

$$X_t^{*AR-S} - \hat{\mu}_X = \sum_{j=1}^{\hat{p}} \hat{\varphi}_j(X_{t-j}^{*AR-S} - \hat{\mu}_X) + e_t^*, \quad t \in \mathbb{Z}$$

με $\{e_t^*\}_{t \in \mathbb{Z}}$ iid ακολουθία των καταλοίπων με περιθώρια κατανομή $e_t^* \sim \hat{F}_e$. Το AR-sieve bootstrap δείγμα είναι ένα πεπερασμένο δείγμα X_1^*, \dots, X_n^* από την παραπάνω διαδικασία η οποία έχει από κατανομή $\hat{P}_{n,AR}$. Επειδή οι πρώτες \hat{p} τιμές δεν μπορούν να εκτιμηθούν, μπορούμε να επιλέξουμε τιμές εκκίνησης ίσες με $\hat{\mu}_X$ και να χρησιμοποιήσουμε την παραπάνω αναδρομή έως ότου επιτευχθεί στασιμότητα στην σειρά και στη συνέχεια να διαγράψουμε τις αρχικές παρατηρήσεις.

Οι παραγόμενες χρονοσειρές bootstrap μπορούν να χρησιμοποιηθούν για διάφορους σκοπούς. Ας υποθέσουμε τώρα οποιοδήποτε στατιστικό $T_n = T_n(X_1, \dots, X_n)$ όπου T_n είναι μια μετρήσιμη συνάρτηση των παρατηρήσεων. Ορίζουμε την εκτίμηση bootstrap του στατιστικού ως

$$T_n^* = T_n(X_1^*, \dots, X_n^*)$$

Με λίγα λόγια στο sieve bootstrap υποθέτουμε ότι η χρονοσειρά μας έχει την μορφή του $AR(\infty)$ και επιλέγεται χρησιμοποιώντας κάποιο κριτήριο το καλύτερο προσαρμοσμένο αυτοπαλίνδρομο μοντέλο p τάξης, όπου $p < n$, και υπολογίζονται τα $n - p$ κατάλοιπα. Το bootstrap πραγματοποιείται στα κεντραρισμένα $n - p$ κατάλοιπα και χρησιμοποιώντας την δομή του μοντέλου AR προκύπτουν οι χρονοσειρές sieve bootstrap.

Κεφάλαιο 4

4.1 Περιγραφή μελέτης προσομοίωσης

Η εφαρμογή των μεθόδων bootstrap στο κεφάλαιο 4 και 5 έχει πραγματοποιηθεί στην γλώσσα προγραμματισμού R.

Θα μελετήσουμε την απόδοση 5 διαφορετικών μεθόδων bootstrap σε προσομοιωμένα δεδομένα χρονοσειρών, και συγκεκριμένα σε δεδομένα που προκύπτουν από τη διαδικασία

- $AR(1)$

$$X_t = \varphi X_{t-1} + e_t$$

όπου τα σφάλματα $e_t \sim N(0,1)$.

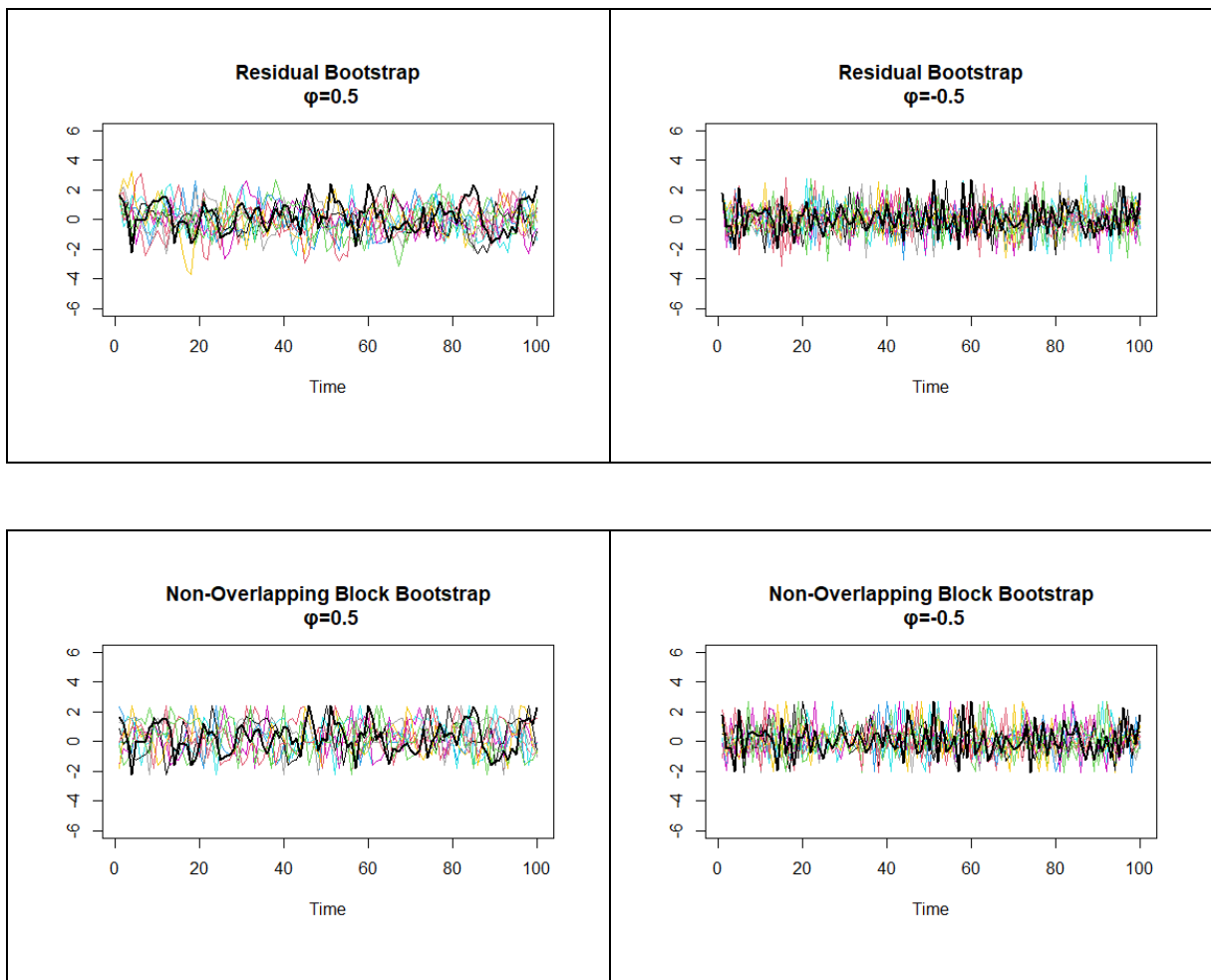
Συγκεκριμένα, θα χρησιμοποιηθεί η βάσει μοντέλου μέθοδος bootstrap στα σχεδόν iid κατάλοιπα, την οποία καλούμε και Residual Bootstrap, η Non-Overlapping Block Bootstrap (NBB), η Moving Block Bootstrap (MBB), η Circular Block Bootstrap (CBB) και η Stationary Bootstrap (SB). Θα πραγματοποιήσουμε προσομοίωση Monte Carlo διαφορετικών σεναρίων για διαφορετικά μεγέθη αρχικού δείγματος $n = 100$ και 200 . Επίσης, να επισημάνουμε ότι για τις μεθόδους block bootstrap θα χρησιμοποιηθούν διάφορα μεγέθη για τα blocks. Σε κάθε επανάληψη Monte Carlo θα προσομοιώνονται δεδομένα από την παραπάνω διαδικασία στα οποία θα εφαρμόζεται η μέθοδος bootstrap της επιλογής μας. Αφού εφαρμοστεί μια μέθοδος bootstrap στα προσομοιωμένα δεδομένα και προκύψουν οι χρονοσειρές bootstrap, σε κάθε χρονοσειρά bootstrap θα προσαρμόζεται μοντέλο $AR(1)$. Να σημειωθεί ότι θα εφαρμοστούν οι μέθοδοι για τιμές $\varphi = -0.75, -0.50, -0.25, -0.10, 0.10, 0.25, 0.50, 0.75$. Τα αποτελέσματα που προκύπτουν για κάθε μέθοδο bootstrap από την μελέτη προσομοίωσης θα σχολιαστούν καταλλήλως.

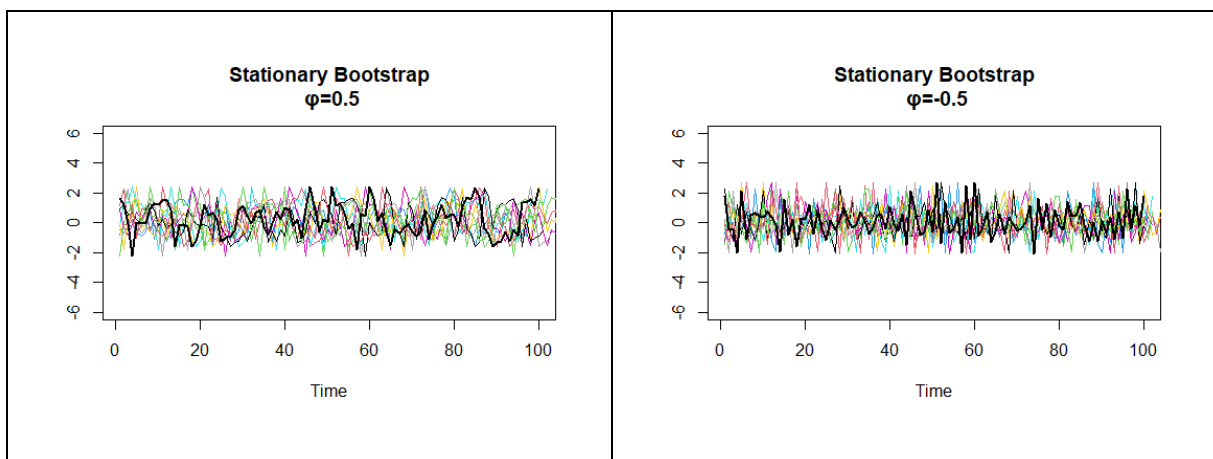
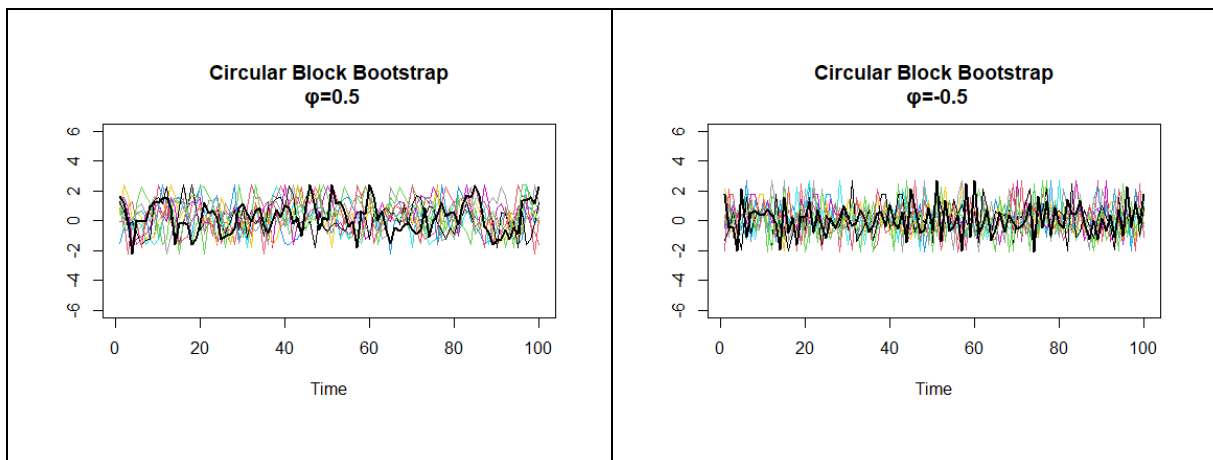
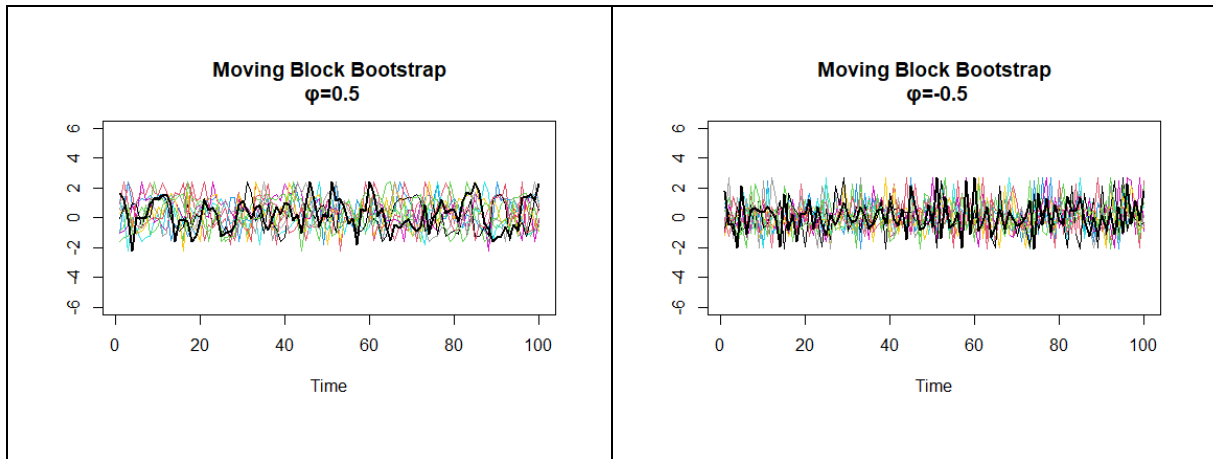
4.2 Οπτικοποίηση εφαρμογής μεθόδων Bootstrap

Αρχικά, πριν την προσομοίωση Monte Carlo, θα πραγματοποιήσουμε μεμονωμένες εφαρμογές των διαφορετικών μεθόδων bootstrap για $B = 10$ επαναλήψεις, σε προσομοιωμένα δεδομένα μεγέθους $n = 100$ από $AR(1)$ με $\varphi = 0.5$ και $\varphi = -0.5$. Η εφαρμογή πραγματοποιήθηκε ώστε να δοθεί μια πρώτη εικόνα για τη συμπεριφορά των μεθόδων. Η έντονη μαύρη γραμμή υποδηλώνει την αρχική προσομοιωμένη χρονοσειρά AR , ενώ οι

πολύχρωμες γραμμές είναι οι 10 χρονοσειρές bootstrap που προκύπτουν από κάθε μέθοδο. Ας σημειωθεί ότι, το μέγεθος των blocks στις μεθόδους block bootstrap έχει επιλεγεί να είναι $\ell = 4$ για την συγκεκριμένη οπτικοποίηση. Για να έχουν οι στάσιμες χρονοσειρές που προκύπτουν από την μέθοδο Stationary Bootstrap μέσω μήκος block ίσο με 4, η πιθανότητα που χρησιμοποιείται για την γεωμετρική κατανομή είναι 0.25. Παρατηρώντας τα παρακάτω γραφήματα μπορούμε να διαπιστώσουμε ότι όλες οι μέθοδοι δημιουργούν χρονοσειρές bootstrap με χαρακτηριστικά που ομοιάζουν με αυτά της αρχικής χρονοσειράς.

Οπτικοποίηση:

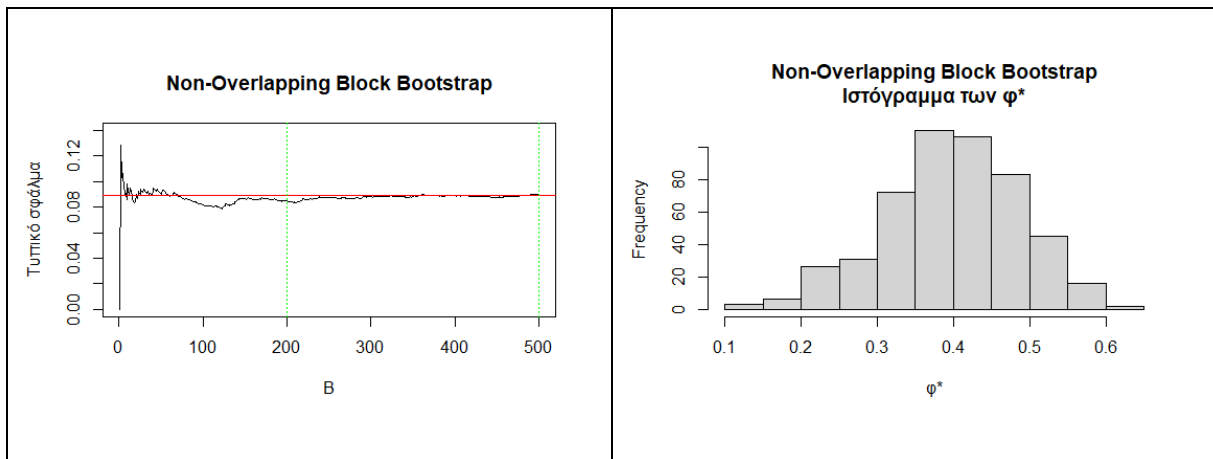
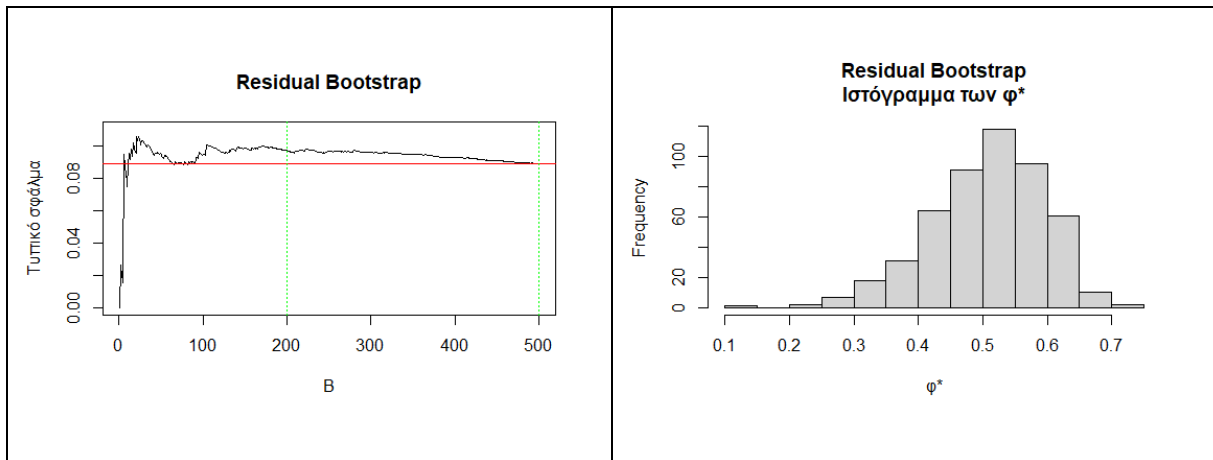


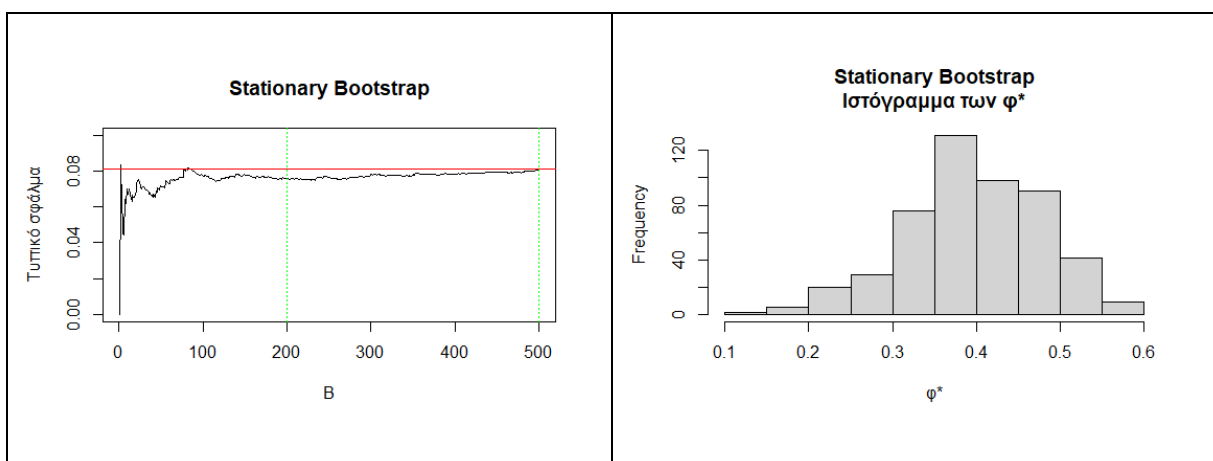
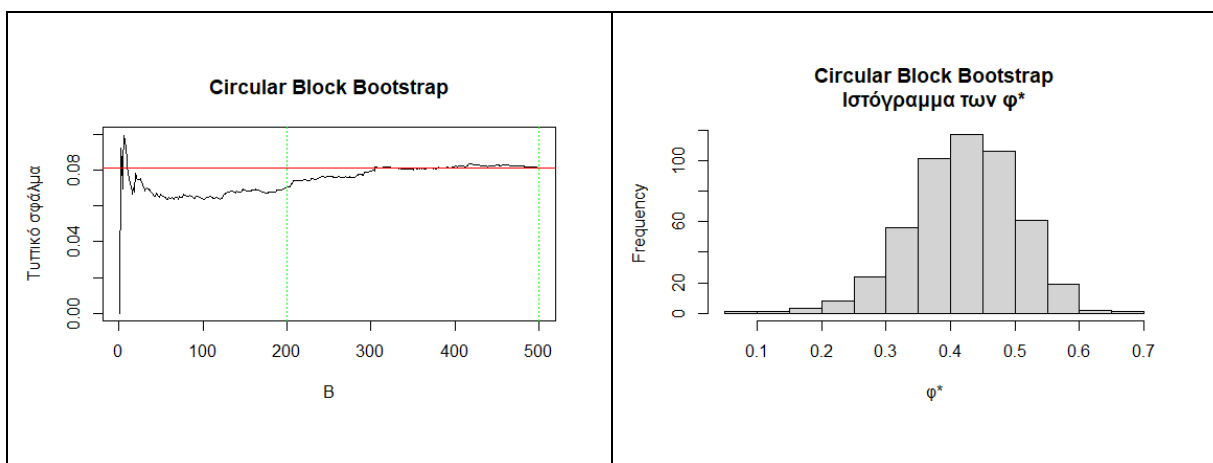
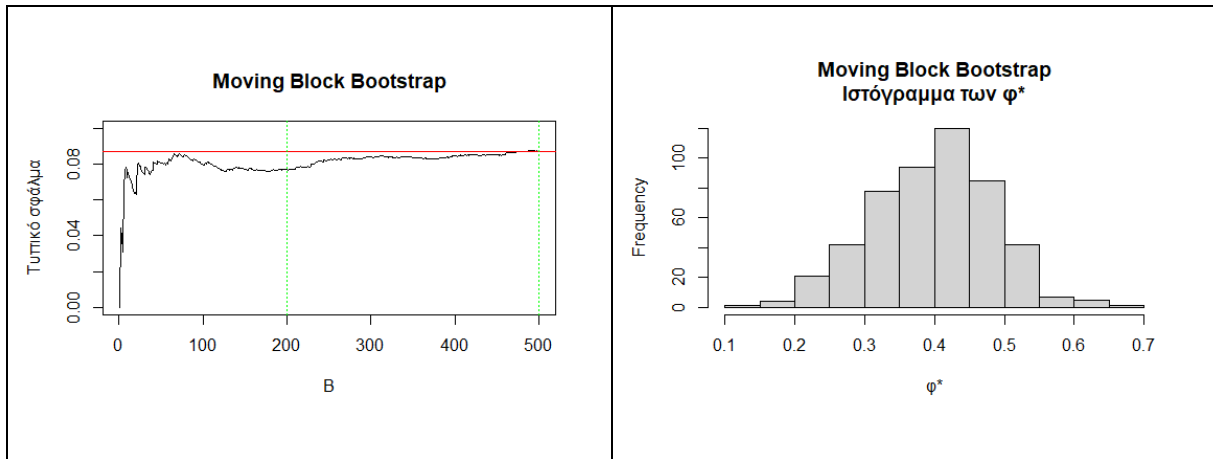


Επίσης, για τα ίδια δεδομένα $AR(1)$ που χρησιμοποιήθηκαν παραπάνω θα εφαρμόσουμε τις 5 μεθόδους bootstrap για $B = 500$ και σε κάθε χρονοσειρά bootstrap που δημιουργείται από την εκάστοτε μέθοδο θα προσαρμόσουμε ένα μοντέλο $AR(1)$. Ως αποτέλεσμα στο τέλος κάθε μεθόδου θα έχουμε στη διάθεση μας B χρονοσειρές bootstrap και B τιμές $\hat{\varphi}^*$. Θα

χρησιμοποιήσουμε τους εκτιμητές bootstrap της παραμέτρου που προκύπτουν από κάθε μέθοδο ώστε να απεικονίσουμε την ακολουθία των τυπικών σφαλμάτων του εκτιμητή της παραμέτρου, καθώς και ένα ιστόγραμμα των τιμών αυτών. Παρακάτω θα δούμε τις μεθόδους bootstrap χρησιμοποιώντας ως αρχικά δεδομένα, προσομοιωμένα δεδομένα μεγέθους $n = 100$ από $AR(1)$ με $\varphi = 0.5$

Η κόκκινη γραμμή υποδηλώνει το τυπικό σφάλμα του εκτιμημένου συντελεστή που προκύπτει από το bootstrap. Οι κάθετες πράσινες διακεκομμένες γραμμές εστιάζουν στις επαναλήψεις bootstrap, 200 και 500 αντίστοιχα. Από την ακολουθία των τυπικών σφαλμάτων μπορούμε να αντλήσουμε πληροφορία για το αν και πότε σταθεροποιείται το τυπικό σφάλμα.

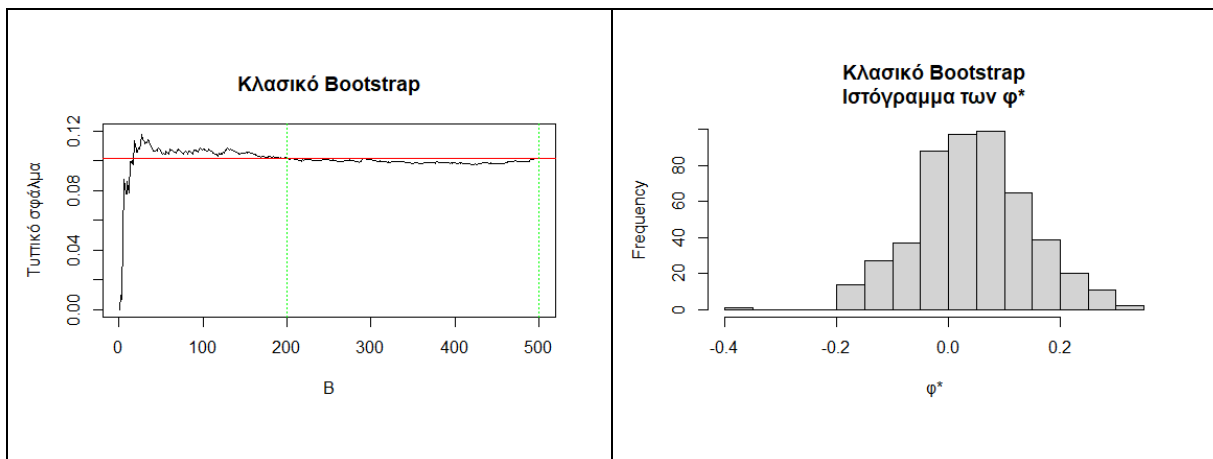




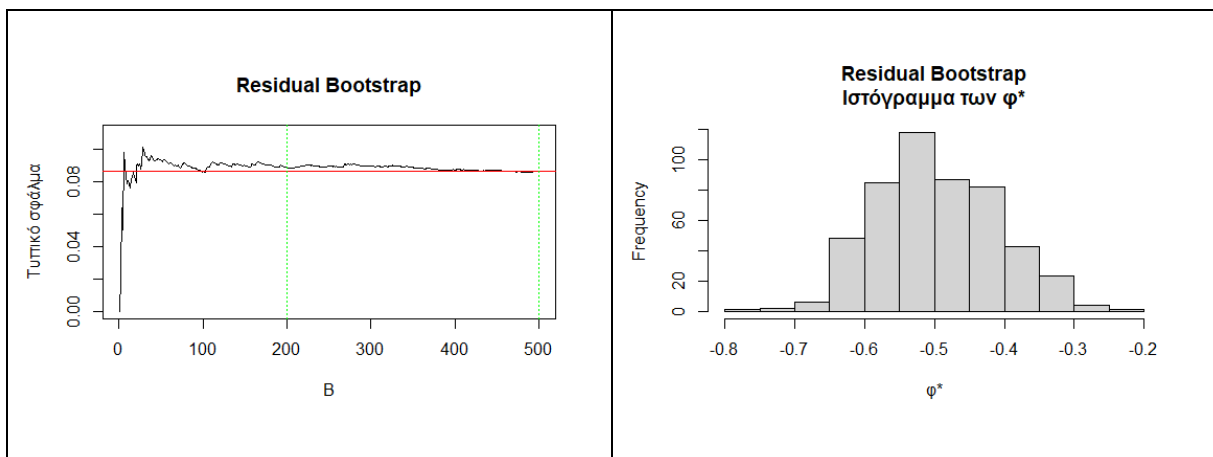
Παρατηρούμε από τα παραπάνω αποτελέσματα ότι η ακολουθία των τυπικών σφαλμάτων φαίνεται σε όλες τις μεθόδους bootstrap να σταθεροποιείται σχετικά νωρίς. Επιπρόσθετα, οι τιμές των $\hat{\varphi}^*$ φαίνεται να προσεγγίζουν την πραγματική τιμή της παραμέτρου, δηλαδή 0.5 που εδώ την γνωρίζουμε. Ας σημειωθεί ότι τα αποτελέσματα των μεθόδων block bootstrap

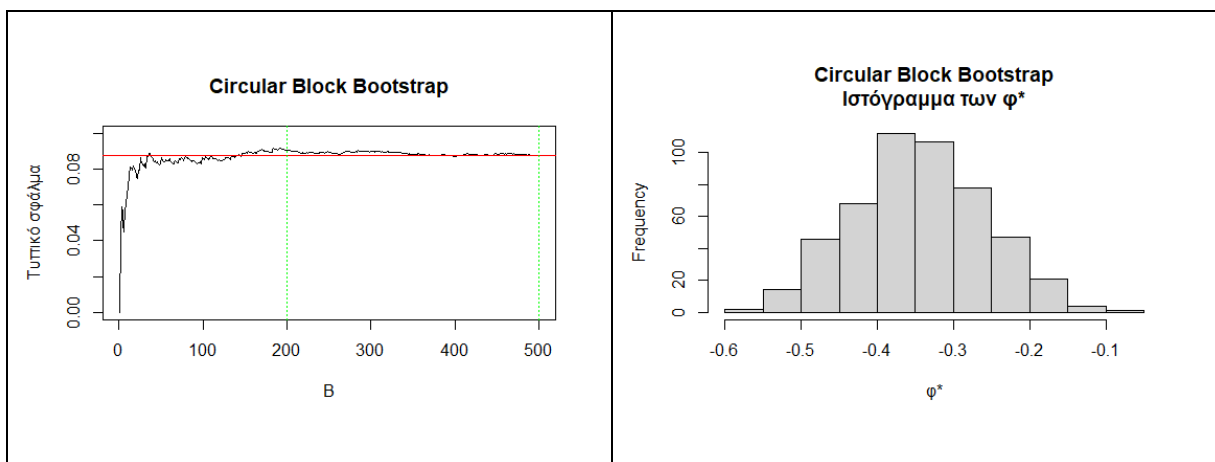
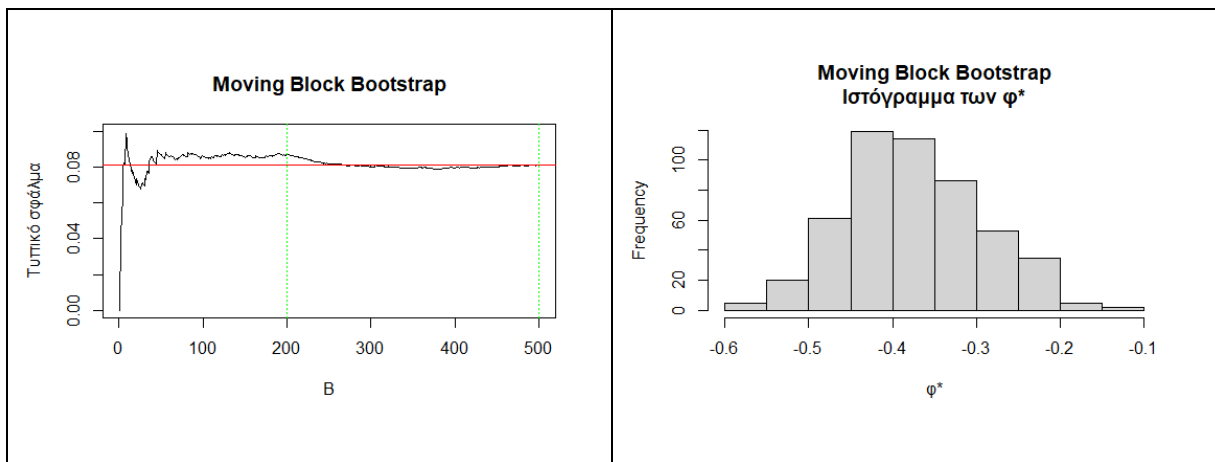
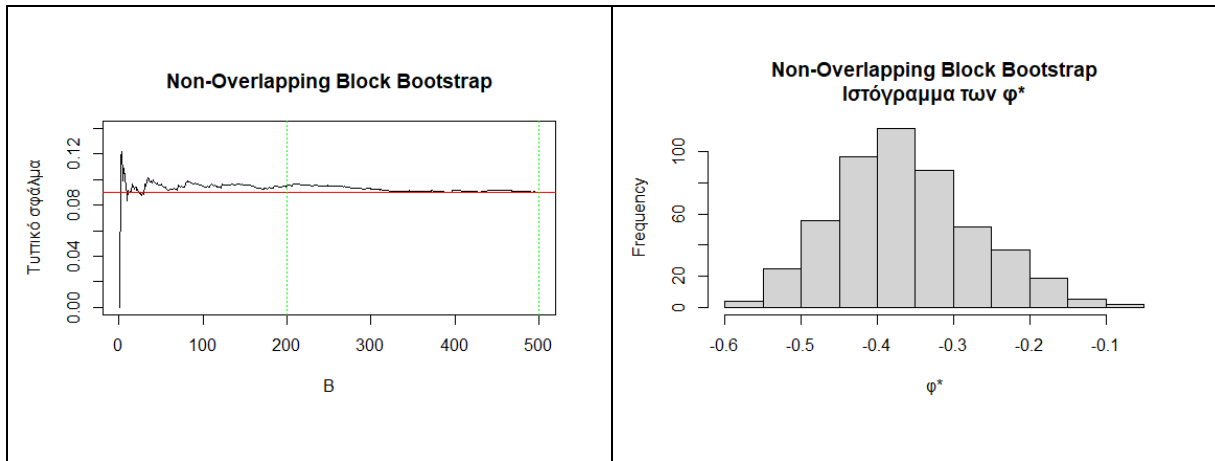
εξαρτώνται πολύ από την τιμή του μεγέθους του block και όπως θα δούμε παρακάτω για μεγαλύτερο μέγεθος block οι τιμές των $\hat{\varphi}^*$ προσεγγίζουν ακόμη περισσότερο την πραγματική τιμή της παραμέτρου 0.5.

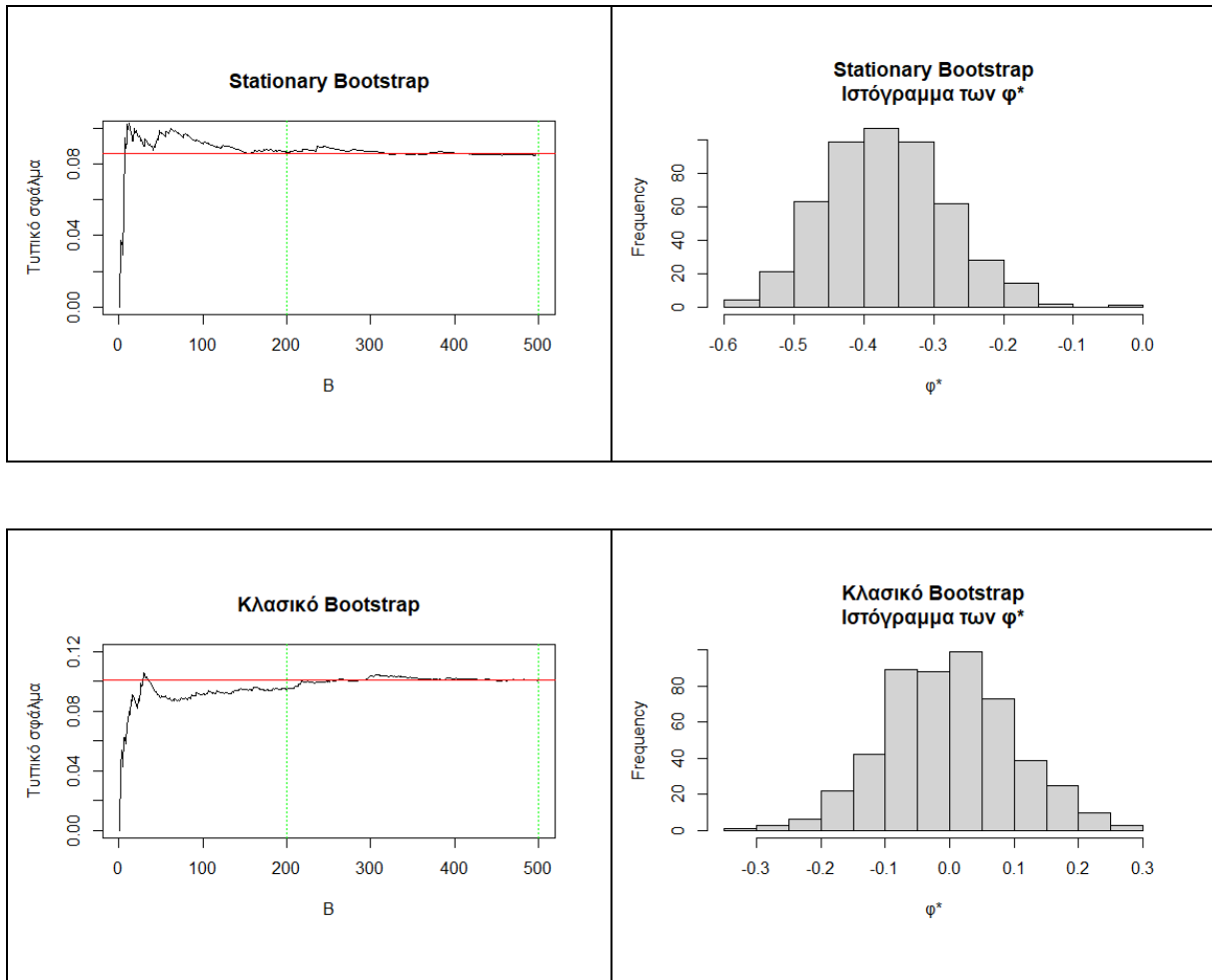
Τέλος, για $\varphi = 0.5$ θα παρουσιάσουμε την εφαρμογή του κλασικού bootstrap του Efron (1979) το οποίο αποτυγχάνει σε δεδομένα χρονοσειρών. Όπως μπορούμε να παρατηρήσουμε το ιστόγραμμα των $\hat{\varphi}^*$ δεν κρίνεται ικανοποιητικό βάσει της τιμής της παραμέτρου. Το αποτέλεσμα αυτό ήταν αναμενόμενο καθώς στο iid bootstrap όλη η πληροφορία της εξάρτησης των δεδομένων χάνεται λόγω του ολικού ανακατέματος των δεδομένων.



Παρακάτω θα δούμε τις μεθόδους bootstrap χρησιμοποιώντας ως αρχικά δεδομένα, προσομοιωμένα δεδομένα μεγέθους $n = 100$ από $AR(1)$ με $\varphi = -0.5$







Παρατηρούμε ότι και στην προκειμένη περίπτωση με $\varphi = -0.5$ το κλασικό bootstrap φαίνεται να αποτυγχάνει, ενώ οι μέθοδοι bootstrap για χρονοσειρές να δίνουν μια ικανοποιητική εικόνα. Στα πλαίσια της εξέτασης των διαφορετικών μεθόδων bootstrap στη συνέχεια πραγματοποιήθηκε προσομοίωση Monte Carlo 500 επαναλήψεων, με μεγέθη αρχικού δείγματος $n = 100$ και 200 , και 200 επαναλήψεων bootstrap. Οι επαναλήψεις bootstrap επιλέχθηκαν να είναι 200 καθώς είδαμε και διαγραμματικά ότι 200 επαναλήψεις αρκούν ώστε να σταθεροποιηθεί η εκτίμηση του τυπικού σφάλματος και επίσης σύμφωνα με τους Efron and Tibshirani (1993) για εκτίμηση τυπικού σφάλματος κρίνονται αρκετές.

4.3 Προσομοίωση για $n = 100$

Παρακάτω πραγματοποιήθηκε προσομοίωση Monte Carlo 500 επαναλήψεων, όπου σε κάθε επανάληψη προσομοιώθηκαν n παρατηρήσεις από ένα μοντέλο $AR(1)$ με το ανάλογο φ και στις οποίες εφαρμόστηκε η εκάστοτε μέθοδος bootstrap. Για τις μεθόδους block bootstrap

πραγματοποιήθηκαν δοκιμές για διάφορα ℓ . Σημαντικό να σημειωθεί, ότι για την μέθοδο Residual Bootstrap έχει χρησιμοποιηθεί η μέθοδος ελαχίστων τετραγώνων ενώ για τις μεθόδους των block bootstrap στην προσαρμογή του μοντέλου η μέθοδος μέγιστης πιθανοφάνειας.

Πίνακας 4.1

φ	$n = 100$						
	MLE		OLS		Residual Bootstrap		
	SE	Bias	SE	Bias	SE	Bias	Perc. CI
-0.75	0.0667	0.0178 (-0.7321)	0.0678	0.0141 (-0.7358)	0.0697	0.0120 (-0.7238)	0.916
-0.50	0.0867	0.0091 (-0.4908)	0.0871	0.0081 (-0.4918)	0.0866	0.0048 (-0.4869)	0.924
-0.25	0.0969	0.0033 (-0.2466)	0.0969	0.0033 (-0.2466)	0.0955	-0.0025 (-0.2492)	0.92
-0.10	0.0995	0.00007 (-0.0999)	0.0995	0.0001 (-0.0998)	0.0980	-0.0071 (-0.1069)	0.91
0.10	0.0995	-0.0044 (0.0955)	0.0995	-0.0044 (0.0955)	0.0983	-0.0131 (0.0823)	0.908
0.25	0.0969	-0.0078 (0.2421)	0.0970	-0.0078 (0.2421)	0.0964	-0.0177 (0.2244)	0.9
0.50	0.0868	-0.0133 (0.4866)	0.0873	-0.0124 (0.4875)	0.0885	-0.0254 (0.4621)	0.88
0.75	0.0666	-0.0186 (0.7313)	0.0678	-0.0148 (0.7351)	0.0731	-0.0337 (0.7014)	0.824

Μετά την εκτέλεση 500 επαναλήψεων Monte Carlo για κάθε φ , αφού δηλαδή για ένα φ προσομοιώσουμε 500 διαφορετικές χρονοσειρές μεγέθους 100, από τον παραπάνω πίνακα μπορούμε να δούμε το τυπικό σφάλμα και την μεροληψία των εκτιμητών της μεθόδου μέγιστης πιθανοφάνειας (MLE) και της μεθόδου ελαχίστων τετραγώνων (OLS), καθώς επίσης και τις τιμές των εκτιμήσεων των μεθόδων αυτών οι οποίες βρίσκονται στις παρενθέσεις. Να σημειωθεί ότι οι εκτιμήσεις των τυπικών σφαλμάτων του εκτιμητή, που προκύπτουν από την MLE και OLS βασίζονται στην πληροφορία του Fisher υπό την υπόθεση της κανονικότητας των σφαλμάτων. Από την εκτίμηση με τις μεθόδους MLE και OLS παρατηρούμε ένα μοτίβο στις τιμές των τυπικών σφαλμάτων και της μεροληψίας και συγκεκριμένα ότι, όταν η τιμή του φ μειώνεται κατά απόλυτη τιμή το τυπικό σφάλμα αυξάνει και η μεροληψία μειώνει και το αντίστροφο όταν η τιμή του φ αυξάνει κατά απόλυτη τιμή. Αυτό είναι αποτέλεσμα της εξάρτησης η οποία γίνεται πιο ισχυρή όσο μεγαλώνει η τιμή του φ και πιο ασθενής όσο μικραίνει. Στο δεξί μέρος του πίνακα μπορούμε να δούμε τα αποτελέσματα της μεθόδου Residual Bootstrap τα οποία εκ πρώτης όψεως φαίνονται πολύ ικανοποιητικά όπως και θα αναμέναμε, καθώς γνωρίζουμε το παραμετρικό μοντέλο που ερμηνεύει τα δεδομένα μας. Οι

εκτιμήσεις bootstrap φαίνεται να είναι πολύ κοντά στις εκτιμήσεις της μεθόδου OLS και να ακολουθούν το μοτίβο που προαναφέραμε. Τα τυπικά σφάλματα της Residual Bootstrap φαίνονται ελαφρώς μειωμένα για μικρές κατά απόλυτες τιμές του συντελεστή φ σε σύγκριση με της μεθόδου OLS. Επιπρόσθετα, δίνονται οι πιθανότητες κάλυψης του 90% διαστήματος Percentile Bootstrap για το φ . Φαίνεται η πιθανότητα κάλυψης να αυξάνεται όταν το φ κατά απόλυτη τιμή μειώνεται. Αυτό συμβαίνει λόγω του ότι μειώνεται η μεροληψία και αυξάνεται η τυπική απόκλιση των τιμών $\hat{\varphi}^*$, το οποίο έχει σαν αποτέλεσμα οι εκτιμήσεις bootstrap να είναι πιο στοχευμένες στην πραγματική τιμή και το εύρος των τιμών τους μεγαλύτερο δηλαδή πιο διεσπαρμένες.

Πίνακας 4.2

$n = 100$ και $\ell = 4$												
φ	NBB			MBB			CBB			SB		
	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0737	0.1693 (-0.5627)	0.070	0.0740	0.1829 (-0.5491)	0.004	0.0746	0.1873 (-0.5447)	0.002	0.0890	0.1900 (-0.5420)	0.020
-0.50	0.0867	0.1188 (-0.3719)	0.568	0.0885	0.1256 (-0.3652)	0.530	0.0884	0.1284 (-0.3624)	0.518	0.0919	0.1298 (-0.3609)	0.538
-0.25	0.0927	0.0599 (-0.1867)	0.840	0.0946	0.0647 (-0.1819)	0.880	0.0944	0.0658 (-0.1807)	0.878	0.0933	0.0670 (-0.1795)	0.870
-0.10	0.0941	0.0236 (-0.0763)	0.922	0.0960	0.0277 (-0.0722)	0.954	0.0958	0.0279 (-0.0720)	0.956	0.0935	0.0289 (-0.0709)	0.944
0.10	0.0942	-0.0248 (0.0707)	0.908	0.0960	-0.0214 (0.0741)	0.948	0.0957	-0.0225 (0.0730)	0.942	0.0933	-0.0216 (0.0739)	0.942
0.25	0.0930	-0.0608 (0.1812)	0.826	0.0946	-0.0579 (0.1842)	0.866	0.0944	-0.0599 (0.1821)	0.868	0.0928	-0.0591 (0.1829)	0.852
0.50	0.0877	-0.1193 (0.3673)	0.542	0.0888	-0.1166 (0.3700)	0.576	0.0889	-0.1202 (0.3663)	0.542	0.0913	-0.1198 (0.3668)	0.546
0.75	0.0743	-0.1706 (0.5606)	0.082	0.0749	-0.1675 (0.5638)	0.076	0.0756	-0.1727 (0.5586)	0.068	0.0874	-0.1735 (0.5578)	0.136

Στον Πίνακα 4.2 μπορούμε να δούμε τα αποτελέσματα των μεθόδων των blocks για μέγεθος block $\ell = 4$. Η μέθοδος NBB δημιουργεί $k = n/\ell = 25$ μη επικαλυπτόμενα blocks μεγέθους 4 στα οποία πραγματοποιούμε δειγματοληψία με επανάθεση k blocks για τον σχηματισμό μιας χρονοσειράς bootstrap. Η μέθοδος MBB δημιουργεί $N = n - \ell + 1 = 97$ επικαλυπτόμενα blocks μεγέθους 4 στα οποία πραγματοποιούμε δειγματοληψία με επανάθεση k blocks, ενώ η μέθοδος CBB τυλίγοντας τα δεδομένα γύρω από έναν νοητό κύκλο δημιουργεί $n = 100$ blocks μεγέθους 4 στα οποία πραγματοποιούμε και πάλι δειγματοληψία με επανάθεση k blocks ώστε να σχηματίσουμε μια χρονοσειρά bootstrap. Στην περίπτωση της μεθόδου SB έχει

χρησιμοποιηθεί $p = 1/\ell = 0.25$ ώστε το αναμενόμενο μήκος block να είναι 4. Παρατηρούμε ότι οι 4 μέθοδοι δίνουν παρόμοια αποτελέσματα τα οποία ακολουθούν το μοτίβο των εκτιμήσεων της MLE. Πιο συγκεκριμένα φαίνεται το τυπικό σφάλμα των μεθόδων για τα διάφορα φ να είναι περίπου ίδιο με της μεθόδου MLE ενώ για μικρά κατά απόλυτη τιμή φ φαίνεται να είναι και ελαφρώς μικρότερο. Ιδανικά θα επιθυμούσαμε η πιθανότητα κάλυψης του διαστήματος εμπιστοσύνης να είναι ακριβώς ίση με τον ονομαστικό συντελεστή εμπιστοσύνης $1 - \alpha$, όμως μπορούμε να διακρίνουμε ότι η πιθανότητα κάλυψης του διαστήματος Percentile Bootstrap είναι πολύ μικρή για μεγάλες κατά απόλυτες τιμές του φ , δηλαδή για ισχυρή εξάρτηση στην χρονοσειρά. Αυτό συμβαίνει λόγω της αυξημένης μεροληψίας που έχουν οι εκτιμήσεις για μεγάλα φ και του μειωμένου τυπικού σφάλματος, δηλαδή μειωμένου εύρους τιμών $\hat{\varphi}^*$. Πιο συγκεκριμένα, όσο το φ πλησιάζει το 1 το διάστημα Percentile Bootstrap φαίνεται να μετακινείται αριστερότερα του πραγματικού φ ενώ όταν το φ πλησιάζει το -1 το διάστημα φαίνεται να μετακινείται δεξιότερα του πραγματικού φ . Για καλύτερα αποτελέσματα για μεγάλες κατά απόλυτες τιμές του φ πρέπει να αυξήσουμε το μέγεθος ℓ των blocks ώστε να διατηρείται πιστά η εξάρτηση που έχουν τα αρχικά δεδομένα. Όπως γνωρίζουμε η προσέγγιση των blocks έχει καλύτερα αποτελέσματα αν η εξάρτηση στην ακολουθία παρατηρήσεων είναι ασθενής και τα blocks είναι μεγάλα σε μέγεθος. Για μικρά κατά απόλυτη τιμή φ , δηλαδή μικρή εξάρτηση στην χρονοσειρά, μπορούμε να πούμε ότι η επιλογή μεγέθους block ίσο με 4 είναι πολύ ικανοποιητική. Για μεγάλα κατά απόλυτη τιμή φ , ισχυρότερη εξάρτηση, περιμένουμε ότι αυξημένες τιμές μεγέθους blocks θα βελτιώσουν την εικόνα. Γενικότερα, επιθυμούμε μια μέθοδος να δίνει μικρή μεροληψία και μικρό τυπικό σφάλμα τέτοια ώστε η πιθανότητα κάλυψης των Percentile διαστημάτων εμπιστοσύνης να είναι όσο το δυνατόν μεγαλύτερη. Δηλαδή, στις 500 επαναλήψεις Monte Carlo η πραγματική τιμή του φ να βρίσκεται όσο το δυνατόν πιο πολλές φορές εντός του 90% διαστήματος. Παρακάτω μπορούμε να παρατηρήσουμε την συμπεριφορά των μεθόδων block bootstrap για μεγαλύτερο μέγεθος blocks.

Πίνακας 4.3

$n = 100$ και $\ell = 8$												
φ	NBB			MBB			CBB			SB		
	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0723	0.0892 (-0.6428)	0.524	0.0715	0.0922 (-0.6398)	0.512	0.0725	0.0980 (-0.6340)	0.47	0.0766	0.0997 (-0.6323)	0.49
-0.50	0.0847	0.0616 (-0.4291)	0.764	0.0852	0.0645 (-0.4262)	0.788	0.0853	0.0679 (-0.4229)	0.78	0.0841	0.0697 (-0.4211)	0.762
-0.25	0.0911	0.0306 (-0.2159)	0.854	0.0911	0.0336 (-0.2129)	0.898	0.0912	0.0349 (-0.2116)	0.906	0.0880	0.0368 (-0.2098)	0.862

-0.10	0.0928	0.0116 (-0.0882)	0.89	0.0923	0.0146 (-0.0853)	0.896	0.0925	0.0147 (-0.0851)	0.914	0.0889	0.0165 (-0.0833)	0.884
0.10	0.0928	-0.0135 (0.0820)	0.876	0.0917	-0.0107 (0.0847)	0.89	0.0922	-0.0120 (0.0834)	0.896	0.0886	-0.0105 (0.0850)	0.878
0.25	0.0912	-0.0323 (0.2098)	0.836	0.0899	-0.0295 (0.2126)	0.862	0.0905	-0.0319 (0.2101)	0.868	0.0875	-0.0307 (0.2114)	0.848
0.50	0.0850	-0.0630 (0.4235)	0.714	0.0839	-0.0592 (0.4273)	0.748	0.0847	-0.0638 (0.4228)	0.722	0.0835	-0.0628 (0.4237)	0.724
0.75	0.0721	-0.0900 (0.6412)	0.502	0.0711	-0.0844 (0.6469)	0.532	0.0724	-0.0912 (0.6400)	0.508	0.0753	-0.0904 (0.6408)	0.504

Στον Πίνακα 4.3 μπορούμε να δούμε τα αποτελέσματα των μεθόδων των blocks για μέγεθος block $\ell = 8$. Λόγω του ότι n/ℓ δεν είναι ακέραιος αριθμός επιλέγουμε ως k τον μικρότερο ακέραιο που είναι μεγαλύτερος από n/ℓ . Συνεπώς, η μέθοδος NBB δημιουργεί $k = \lceil n/\ell \rceil = 13$ μη επικαλυπτόμενα blocks μεγέθους 8 στα οποία πραγματοποιούμε δειγματοληψία με επανάθεση k blocks για τον σχηματισμό μιας χρονοσειράς bootstrap. Η μέθοδος MBB δημιουργεί $N = n - \ell + 1 = 93$ επικαλυπτόμενα blocks μεγέθους 8 στα οποία πραγματοποιούμε δειγματοληψία με επανάθεση k blocks, ενώ η μέθοδος CBB τυλίγοντας τα δεδομένα γύρω από έναν νοητό κύκλο δημιουργεί $n = 100$ blocks μεγέθους 8 στα οποία πραγματοποιούμε και πάλι δειγματοληψία με επανάθεση k blocks ώστε να σχηματίσουμε μια χρονοσειρά bootstrap. Στην περίπτωση της μεθόδου SB έχει χρησιμοποιηθεί $p = 1/\ell = 0.125$. Να σημειωθεί ότι στην συγκεκριμένη περίπτωση και γενικότερα στις περιπτώσεις που το ℓ δεν διαιρείται ακριβώς με το μέγεθος του δείγματος τα τυπικά σφάλματα bootstrap έχουν πολλαπλασιαστεί με $(k\ell/n)^{1/2}$. Από τα παραπάνω αποτελέσματα μπορούμε να διακρίνουμε ότι το τυπικό σφάλμα των εκτιμήσεων bootstrap είναι ελαφρώς μειωμένο για όλες τις μεθόδους σε σχέση με τον Πίνακα 4.2 ενώ η μεροληψία φαίνεται να υποδιπλασιάζεται, δηλαδή να έχουμε πιο στοχευμένες εκτιμήσεις. Απόρροια αυτού είναι οι εκτιμήσεις που προκύπτουν από τις μεθόδους για μεγάλα κατά απόλυτη τιμή φ , $|\varphi| = 0.5$ και $|\varphi| = 0.75$, να βελτιώνονται κατά πολύ περισσότερο συγκριτικά με τις εκτιμήσεις για μικρά κατά απόλυτη τιμή φ στις οποίες δείχνει η πιθανότητα κάλυψης ελαφρώς να μειώνεται. Αυτό συμβαίνει επειδή ο υποδιπλασιασμός σε μεγάλο μέγεθος μεροληψία είναι πιο έντονος από τον υποδιπλασιασμό μικρών μεροληψιών. Σαν αποτέλεσμα οι πιθανότητες κάλυψης των Percentile διαστημάτων εμπιστοσύνης όταν $\ell = 8$ μεγαλώνουν έντονα για μεγάλα κατά απόλυτη τιμή φ .

Πίνακας 4.4

$n = 100$ και $\ell = 16$												
	NBB			MBB			CBB			SB		
φ	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI

-0.75	0.0694	0.0466 (-0.6854)	0.708	0.0654	0.0472 (-0.6848)	0.684	0.0668	0.0524 (-0.6796)	0.672	0.0665	0.0546 (-0.6774)	0.656
-0.50	0.0837	0.0324 (-0.4583)	0.838	0.0788	0.0345 (-0.4563)	0.794	0.0799	0.0374 (-0.4534)	0.824	0.0763	0.0396 (-0.4512)	0.756
-0.25	0.0905	0.0163 (-0.2302)	0.862	0.0852	0.0189 (-0.2277)	0.822	0.0863	0.0198 (-0.2268)	0.858	0.0812	0.0216 (-0.2249)	0.806
-0.10	0.0920	0.0064 (-0.0934)	0.856	0.0864	0.0088 (-0.0910)	0.81	0.0877	0.0088 (-0.0911)	0.852	0.0822	0.0102 (-0.0896)	0.814
0.10	0.0916	-0.0068 (0.0887)	0.862	0.0854	-0.0047 (0.0908)	0.804	0.0872	-0.0059 (0.0896)	0.832	0.0819	-0.0052 (0.0903)	0.804
0.25	0.0897	-0.0167 (0.2253)	0.816	0.0832	-0.0148 (0.2273)	0.768	0.0853	-0.0169 (0.2252)	0.808	0.0804	-0.0166 (0.2254)	0.778
0.50	0.0830	-0.0330 (0.4535)	0.752	0.0765	-0.0304 (0.4561)	0.736	0.0791	-0.0342 (0.4523)	0.766	0.0754	-0.0346 (0.4520)	0.718
0.75	0.0690	-0.0475 (0.6838)	0.672	0.0636	-0.0428 (0.6885)	0.662	0.0664	-0.0487 (0.6826)	0.638	0.0652	-0.0492 (0.6821)	0.624

Στον Πίνακα 4.4 μπορούμε να δούμε τα αποτελέσματα των μεθόδων των blocks για μέγεθος block $\ell = 16$. Είναι εμφανές ότι οι 4 μέθοδοι δίνουν παρόμοια αποτελέσματα τα οποία ακολουθούν το μοτίβο των εκτιμήσεων της MLE. Ωστόσο, μπορούμε να διακρίνουμε ότι τα τυπικά σφάλματα των μεθόδων MBB, CBB και SB είναι μικρότερα από της μεθόδου NBB η οποία χρησιμοποιεί 7 μη επικαλυπτόμενα blocks. Μπορούμε να πούμε ότι τα αποτελέσματα των μεθόδων βελτιώνονται για μεγάλα κατά απόλυτη τιμή φ σε σχέση με τα μικρότερα μεγέθη blocks που είδαμε παραπάνω στους Πίνακες 4.2 και 4.3 αλλά δεν μπορούμε να πούμε το ίδιο για μικρά κατά απόλυτη τιμή φ αφού οι πιθανότητες κάλυψης των διαστημάτων εμπιστοσύνης μειώνονται. Το τυπικό σφάλμα των εκτιμήσεων bootstrap φαίνεται να είναι ελαφρώς μειωμένο για όλες τις μεθόδους σε σχέση με τους Πίνακες 4.2 και 4.3, χαμηλότερο και από το τυπικό σφάλμα της MLE του Πίνακα 4.1, ενώ η μεροληψία φαίνεται να υποδιπλασιάζεται σε σχέση με τον Πίνακα 4.3 και να υποτετραπλασιάζεται σε σχέση με τον Πίνακα 4.2. Αυτό οδηγεί τις εκτιμήσεις που προκύπτουν από τις μεθόδους για μεγάλα κατά απόλυτη τιμή φ να βελτιώνονται κατά πολύ περισσότερο. Στον παρακάτω πίνακα μπορούμε να δούμε τα αποτελέσματα των μεθόδων για μέγεθος block $\ell = 25$.

Πίνακας 4.5

$n = 100$ και $\ell = 25$												
NBB				MBB			CBB			SB		
φ	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0576	0.0245 (-0.7075)	0.714	0.0585	0.0313 (-0.7007)	0.726	0.0616	0.0354 (-0.6966)	0.73	0.0560	0.0397 (-0.6923)	0.668

-0.50	0.0701	0.0177 (-0.4731)	0.768	0.0714	0.0234 (-0.4673)	0.77	0.0745	0.0258 (-0.4649)	0.822	0.0653	0.0297 (-0.4611)	0.736
-0.25	0.0766	0.0092 (-0.2374)	0.766	0.0774	0.0132 (-0.2334)	0.78	0.0805	0.0136 (-0.2329)	0.832	0.0700	0.0169 (-0.2297)	0.76
-0.10	0.0782	0.0036 (-0.0962)	0.74	0.0785	0.0064 (-0.0934)	0.77	0.0818	0.0058 (-0.0940)	0.836	0.0710	0.0085 (-0.0913)	0.762
0.10	0.0781	-0.0040 (0.0915)	0.748	0.0779	-0.0029 (0.0926)	0.762	0.0814	-0.0046 (0.0909)	0.82	0.0707	-0.0029 (0.0926)	0.744
0.25	0.0764	-0.0098 (0.2322)	0.72	0.0760	-0.0098 (0.2323)	0.734	0.0796	-0.0123 (0.2298)	0.788	0.0693	-0.0113 (0.2307)	0.714
0.50	0.0703	-0.0190 (0.4676)	0.708	0.0700	-0.0204 (0.4661)	0.726	0.0735	-0.0243 (0.4622)	0.748	0.0645	-0.0243 (0.4622)	0.674
0.75	0.0583	-0.0265 (0.7048)	0.69	0.0580	-0.0281 (0.7032)	0.684	0.0614	-0.0339 (0.6974)	0.694	0.0548	-0.0345 (0.6967)	0.62

Αξίζει να σημειωθεί ότι για $\ell = 25$ η μέθοδος NBB δημιουργεί $k = n/\ell = 4$ μη επικαλυπτόμενα blocks μεγέθους 25 στα οποία πραγματοποιούμε δειγματοληψία με επανάθεση k blocks ενώ ότι η μέθοδος MBB δημιουργεί $N = n - \ell + 1 = 76$ επικαλυπτόμενα blocks μεγέθους 25. Το τυπικό σφάλμα των εκτιμήσεων bootstrap φαίνεται να είναι μειωμένο για όλες τις μεθόδους σε σχέση με τον Πίνακα 4.4 και αισθητά μειωμένο σε σχέση με τους Πίνακες 4.2 και 4.3 ενώ η μεροληψία φαίνεται και πάλι να υποδιπλασιάζεται σε σχέση με τον Πίνακα 4.4. Αξίζει να σημειωθεί, ότι το τυπικό σφάλμα της CBB είναι ελαφρώς αυξημένο συγκριτικά με των NBB, MBB και SB.

Όπως μπορούμε να παρατηρήσουμε από τον παραπάνω πίνακα το αυξημένο μέγεθος των blocks φαίνεται να δίνει μια βελτιωμένη εικόνα για μεγάλα κατά απόλυτη τιμή φ σε σχέση με μικρότερα μεγέθη blocks καθώς μειώνεται η μεροληψία και το τυπικό σφάλμα και αυξάνεται η πιθανότητα κάλυψης. Δεν φαίνεται όμως το ίδιο για μικρά κατά απόλυτη τιμή φ καθώς οι πιθανότητες κάλυψης των 90% διαστημάτων μειώνονται, δηλαδή οι φορές που το διάστημα “πιάνει” την πραγματική τιμή της παραμέτρου στο εύρος τιμών των $\hat{\varphi}^*$ μειώνονται. Εμείς επιθυμούμε έναν συνδυασμό μικρού τυπικού σφάλματος και μεροληψίας με όσο το δυνατόν μεγαλύτερη πιθανότητα κάλυψης, δηλαδή σε κάθε εφαρμογή του bootstrap η μάζα των τιμών $\hat{\varphi}^*$ να βρίσκεται γύρω από την πραγματική τιμή της παραμέτρου. Από τα παραπάνω μπορούμε να συμπεράνουμε ότι αυξημένες τιμές των blocks καταφέρνουν και εγκλωβίζουν την εξάρτηση της αρχικής χρονοσειράς καλύτερα και φαίνεται ότι όταν η εξάρτηση γίνεται ισχυρότερη, μεγαλύτερες τιμές ℓ οδηγούν σε καλύτερα αποτελέσματα. Αυτό το συμπέρασμα φαίνεται να συμφωνεί με τον Carlstein (1986) που επισήμανε ότι όταν η εξάρτηση μεταξύ των τυχαίων μεταβλητών γίνεται πιο ισχυρή τότε χρειαζόμαστε μεγαλύτερο ℓ . Για μικρότερα κατά απόλυτη τιμή μεγέθη φ καταλήγουμε ότι μικρά ℓ είναι προτιμότερα.

4.4 Προσομοίωση για $n = 200$

Πίνακας 4.6

φ	$n = 200$						
	MLE		OLS		Residual Bootstrap		
	SE	Bias	SE	Bias	SE	Bias	Perc. CI
-0.75	0.0469	0.0082 (-0.7417)	0.0472	0.0063 (-0.7436)	0.0481	0.0061 (-0.7375)	0.9
-0.50	0.0612	0.0034 (-0.4965)	0.0613	0.0029 (-0.4970)	0.0614	0.0026 (-0.4943)	0.922
-0.25	0.0684	0.0006 (-0.2493)	0.0684	0.0006 (-0.2493)	0.0683	-0.0011 (-0.2505)	0.912
-0.10	0.0703	-0.0007 (-0.1007)	0.0703	-0.0007 (-0.1007)	0.0702	-0.0034 (-0.1041)	0.914
0.10	0.0703	-0.0022 (0.0977)	0.0703	-0.0023 (0.0976)	0.0703	-0.0065 (0.0911)	0.916
0.25	0.0684	-0.0033 (0.2466)	0.0685	-0.0033 (0.2466)	0.0687	-0.0088 (0.2377)	0.906
0.50	0.0612	-0.0052 (0.4947)	0.0614	-0.0047 (0.4952)	0.0621	-0.0128 (0.4824)	0.896
0.75	0.0467	-0.0079 (0.7420)	0.0472	-0.0059 (0.7440)	0.0494	-0.0167 (0.7272)	0.896

Παραπάνω μπορούμε να δούμε την περίπτωση που η αρχικές χρονοσειρές είναι μεγέθους 200. Από τον Πίνακα 4.6 παρατηρούμε ότι οι εκτιμήσεις των μεθόδων MLE και OLS έχουν μικρότερη μεροληψία και τυπικό σφάλμα σε σχέση με τον Πίνακα 4.1 για $n = 100$. Στο δεξί μέρος του πίνακα μπορούμε να διακρίνουμε τα αποτελέσματα της μεθόδου Residual Bootstrap. Οι εκτιμήσεις bootstrap φαίνεται να είναι πολύ κοντά στις εκτιμήσεις της μεθόδου OLS και να ακολουθούν το μοτίβο που έχουμε ήδη αναφέρει. Συγκριτικά με τα αποτελέσματα της ίδιας μεθόδου που βρίσκονται στον Πίνακα 4.1 μπορούμε να παρατηρήσουμε σαφή μείωση στις τιμές του τυπικού σφάλματος και μεροληψίας κάτι το οποίο είναι αναμενόμενο καθώς οι εκτιμητές ελαχίστων τετραγώνων $\hat{\varphi}^*$ των χρονοσειρών bootstrap είναι βελτιωμένοι λόγω της αύξησης του μεγέθους n (και επειδή ο εκτιμητής $\hat{\varphi}$ OLS που χρησιμοποιούμε για να φτιάξουμε μια χρονοσειρά bootstrap είναι βελτιωμένος λόγω του $n = 200$). Επιπρόσθετα, δίνονται οι πιθανότητες κάλυψης του 90% διαστήματος Percentile Bootstrap για το φ . Παρακάτω, μπορούμε να δούμε τους πίνακες με τα αποτελέσματα των μεθόδων blocks για $\ell = 8, 16, 25, 50$. Η αύξηση του μεγέθους της αρχικής χρονοσειράς μας οδήγησε στο να επιλεγούν μεγαλύτερα μεγέθη για τα blocks και να γίνει δοκιμή των μεθόδων για $\ell = 50$.

Πίνακας 4.7

$n = 200$ και $\ell = 8$												
φ	NBB			MBB			CBB			SB		
	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0501	0.0889 (-0.6528)	0.298	0.0506	0.0930 (-0.6486)	0.27	0.0511	0.0963 (-0.6454)	0.234	0.0548	0.0968 (-0.6448)	0.3
-0.50	0.0601	0.0604 (-0.4360)	0.708	0.0612	0.0637 (-0.4327)	0.726	0.0614	0.0657 (-0.4307)	0.706	0.0610	0.0661 (-0.4304)	0.702
-0.25	0.0654	0.0301 (-0.2192)	0.868	0.0662	0.0328 (-0.2165)	0.89	0.0663	0.0336 (-0.2157)	0.9	0.0644	0.0340 (-0.2153)	0.868
-0.10	0.0668	0.0116 (-0.0890)	0.894	0.0675	0.0139 (-0.0867)	0.924	0.0675	0.0140 (-0.0866)	0.916	0.0653	0.0145 (-0.0861)	0.914
0.10	0.0670	-0.0129 (0.0847)	0.876	0.0676	-0.0114 (0.0863)	0.926	0.0674	-0.0120 (0.0856)	0.92	0.0653	-0.0115 (0.0861)	0.896
0.25	0.0659	-0.0313 (0.2153)	0.848	0.0664	-0.0303 (0.2162)	0.866	0.0662	-0.0315 (0.2150)	0.882	0.0645	-0.0310 (0.2155)	0.868
0.50	0.0610	-0.0615 (0.4332)	0.698	0.0616	-0.0612 (0.4335)	0.718	0.0616	-0.0632 (0.4315)	0.702	0.0613	-0.0628 (0.4319)	0.704
0.75	0.0505	-0.0896 (0.6524)	0.326	0.0512	-0.0893 (0.6527)	0.308	0.0515	-0.0923 (0.6497)	0.294	0.0551	-0.0923 (0.6497)	0.35

Πίνακας 4.8

$n = 200$ και $\ell = 16$												
φ	NBB			MBB			CBB			SB		
	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0479	0.0459 (-0.6957)	0.67	0.0476	0.0469 (-0.6948)	0.676	0.0480	0.0500 (-0.6916)	0.646	0.0481	0.0506 (-0.6910)	0.634
-0.50	0.0591	0.0313 (-0.4651)	0.818	0.0586	0.0328 (-0.4637)	0.822	0.0588	0.0345 (-0.4620)	0.826	0.0566	0.0351 (-0.4613)	0.786
-0.25	0.0650	0.0158 (-0.2335)	0.878	0.0641	0.0173 (-0.2320)	0.886	0.0643	0.0176 (-0.2317)	0.886	0.0611	0.0183 (-0.2310)	0.856
-0.10	0.0665	0.0063 (-0.0943)	0.884	0.0655	0.0076 (-0.0930)	0.89	0.0657	0.0073 (-0.0934)	0.89	0.0622	0.0079 (-0.0927)	0.87
0.10	0.0664	-0.0062 (0.0914)	0.856	0.0653	-0.0054 (0.0922)	0.868	0.0657	-0.0064 (0.0912)	0.89	0.0622	-0.0059 (0.0917)	0.862
0.25	0.0650	-0.0156 (0.2309)	0.854	0.0638	-0.0151 (0.2314)	0.86	0.0643	-0.0167 (0.2299)	0.872	0.0610	-0.0162 (0.2303)	0.846
0.50	0.0593	-0.0311 (0.4636)	0.828	0.0585	-0.0307 (0.4640)	0.812	0.0592	-0.0331 (0.4615)	0.814	0.0567	-0.0328 (0.4619)	0.782
0.75	0.0478	-0.0456 (0.6964)	0.672	0.0478	-0.0445 (0.6975)	0.694	0.0487	-0.0481 (0.6939)	0.66	0.0484	-0.0477 (0.6943)	0.646

Πίνακας 4.9

$n = 200$ και $\ell = 25$												
	NBB			MBB			CBB			SB		
φ	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0436	0.0257 (-0.7159)	0.764	0.0449	0.0303 (-0.7113)	0.76	0.0454	0.0335 (-0.7081)	0.756	0.0425	0.0339 (-0.7077)	0.696
-0.50	0.0548	0.0181 (-0.4784)	0.8	0.0563	0.0216 (-0.4749)	0.83	0.0565	0.0236 (-0.4729)	0.83	0.0513	0.0240 (-0.4724)	0.798
-0.25	0.0607	0.0092 (-0.2401)	0.840	0.0619	0.0117 (-0.2376)	0.864	0.0620	0.0123 (-0.2370)	0.86	0.0561	0.0129 (-0.2364)	0.826
-0.10	0.0622	0.0036 (-0.0970)	0.858	0.0632	0.0054 (-0.0952)	0.862	0.0634	0.0052 (-0.0954)	0.872	0.0573	0.0059 (-0.0947)	0.818
0.10	0.0622	-0.0040 (0.0936)	0.846	0.0631	-0.0030 (0.0946)	0.848	0.0633	-0.0041 (0.0935)	0.878	0.0573	-0.0033 (0.0943)	0.816
0.25	0.0609	-0.0098 (0.2368)	0.84	0.0616	-0.0094 (0.2371)	0.84	0.0619	-0.0111 (0.2354)	0.872	0.0561	-0.0103 (0.2363)	0.81
0.50	0.0558	0.0192 (0.4755)	0.822	0.0561	-0.0196 (0.4751)	0.82	0.0567	-0.0223 (0.4724)	0.836	0.0516	-0.0214 (0.4732)	0.784
0.75	0.0453	-0.0282 (0.7138)	0.736	0.0451	-0.0283 (0.7137)	0.772	0.0462	-0.0321 (0.7099)	0.738	0.0429	-0.0316 (0.7104)	0.708

Πίνακας 4.10

$n = 200$ και $\ell = 50$												
	NBB			MBB			CBB			SB		
φ	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI	SE	Bias	Perc. CI
-0.75	0.0376	0.0122 (-0.7294)	0.732	0.0384	0.0154 (-0.7262)	0.732	0.0401	0.0178 (-0.7239)	0.76	0.0360	0.0196 (-0.7220)	0.714
-0.50	0.0480	0.0088 (-0.4877)	0.726	0.0485	0.0116 (-0.4848)	0.752	0.0506	0.0128 (-0.4837)	0.812	0.0445	0.0144 (-0.4821)	0.724
-0.25	0.0534	0.0045 (-0.2448)	0.76	0.0537	0.0068 (-0.2425)	0.73	0.0560	0.0066 (-0.2427)	0.816	0.0491	0.0079 (-0.2414)	0.736
-0.10	0.0548	0.0018 (-0.0989)	0.754	0.0549	0.0035 (-0.0971)	0.73	0.0574	0.0027 (-0.0979)	0.824	0.0503	0.0037 (-0.0969)	0.74
0.10	0.0548	-0.0019 (0.0957)	0.754	0.0548	-0.0009 (0.0967)	0.716	0.0573	-0.0024 (0.0952)	0.804	0.0503	-0.0017 (0.0959)	0.736
0.25	0.0537	-0.0047 (0.2418)	0.744	0.0534	-0.0043 (0.2423)	0.722	0.0559	-0.0062 (0.2403)	0.816	0.0491	-0.0057 (0.2408)	0.744
0.50	0.0492	-0.0091 (0.4855)	0.752	0.0485	-0.0094 (0.4853)	0.734	0.0507	-0.0122 (0.4825)	0.8	0.0448	-0.0121 (0.4826)	0.736
0.75	0.0394	-0.0131 (0.7289)	0.734	0.0386	-0.0135 (0.7285)	0.736	0.0408	-0.0172 (0.7248)	0.764	0.0364	-0.0179 (0.7241)	0.718

Από τους παραπάνω πίνακες μπορούμε να συμπεράνουμε ότι αυξημένες τιμές των blocks βελτιώνουν τα αποτελέσματα των μεθόδων όταν το φ μεγαλώνει κατά απόλυτη τιμή, όπως συνέβαινε και για $n = 100$, καθώς καταφέρνουν και εγκλωβίζουν την εξάρτηση της αρχικής χρονοσειράς καλύτερα, και φαίνεται ότι όταν η εξάρτηση γίνεται ισχυρότερη μεγαλύτερες τιμές ℓ οδηγούν σε καλύτερα αποτελέσματα. Η χρήση πολύ μικρής τιμής του ℓ καταστρέφει την δομή της εξάρτησης με αποτέλεσμα η μεροληψία των μεθόδων των blocks να αυξάνεται έντονα για μεγάλα φ . Για μικρά φ , αυξημένες τιμές του μεγέθους των blocks επηρεάζουν αρνητικά τις πιθανότητες κάλυψης. Οι πιθανότητες κάλυψης δείχνουν να μειώνονται όταν αυξάνονται οι τιμές των ℓ , λόγω της μείωσης του τυπικού σφάλματος και της μικρής σε ποσότητα μείωσης της μεροληψίας. Αξίζει να σημειωθεί ότι στον Πίνακα 4.10 η μέθοδος NBB δημιουργεί $k = n/\ell = 4$ μη επικαλυπτόμενα blocks ενώ η μέθοδος MBB δημιουργεί $N = n - \ell + 1 = 151$ επικαλυπτόμενα. Συγκρίνοντας τα αποτελέσματα των μεθόδων των blocks για $n = 100$ και $n = 200$ στα κοινά $\ell = 8, 16, 24$ μπορούμε να δούμε ότι η μεροληψία βρίσκεται στα ίδια περίπου επίπεδα ενώ το τυπικό σφάλμα για $n = 200$ είναι μειωμένο κάτι το οποίο αναμέναμε λόγω της αύξησης του n . Οι πιθανότητες κάλυψης των μεθόδων για μεγαλύτερα φ , $|\varphi| = 0.5$ και $|\varphi| = 0.75$, όταν $n = 200$ και $\ell = 8$ φαίνεται να μειώνονται συγκριτικά με την περίπτωση που $n = 100$ και $\ell = 8$. Δηλαδή, συμπεραίνουμε ότι μεγαλύτερα κατά απόλυτη τιμή φ όταν αυξάνεται το μέγεθος n χρειάζονται μεγαλύτερα μεγέθη blocks. Επιπρόσθετα, η εκτίμηση των μεθόδων που βρίσκεται εντός των παρενθέσεων φαίνεται να είναι πιο κοντά στην πραγματική τιμή της παραμέτρου για $n = 200$ σε σχέση με $n = 100$. Κρίνεται σημαντικό να αναφερθεί ότι μεγαλύτερου μήκους χρονοσειρές μας δίνουν την δυνατότητα να χρησιμοποιήσουμε μεγαλύτερο ℓ στις μεθόδους των blocks, όπως και κάναμε στον Πίνακα 4.10. Τέλος, μπορούμε να πούμε ότι υπάρχει δυσκολία σύγκρισής μεταξύ των μεθόδων block bootstrap καθώς κάθε μέθοδος αποδίδει καλύτερα σε διαφορετικό μήκος block και διαφορετικό μέγεθος αρχικής χρονοσειράς.

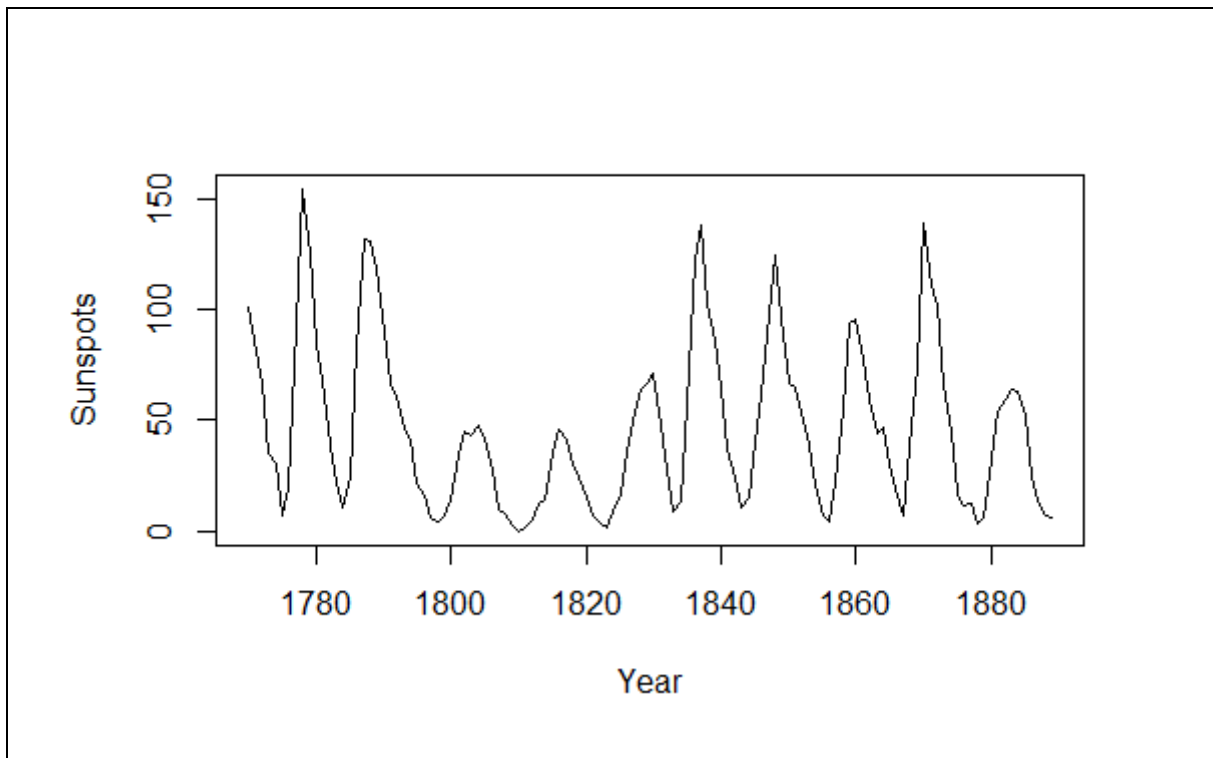
Κεφάλαιο 5

5.1 Εφαρμογή μεθόδων Bootstrap σε πραγματικά δεδομένα

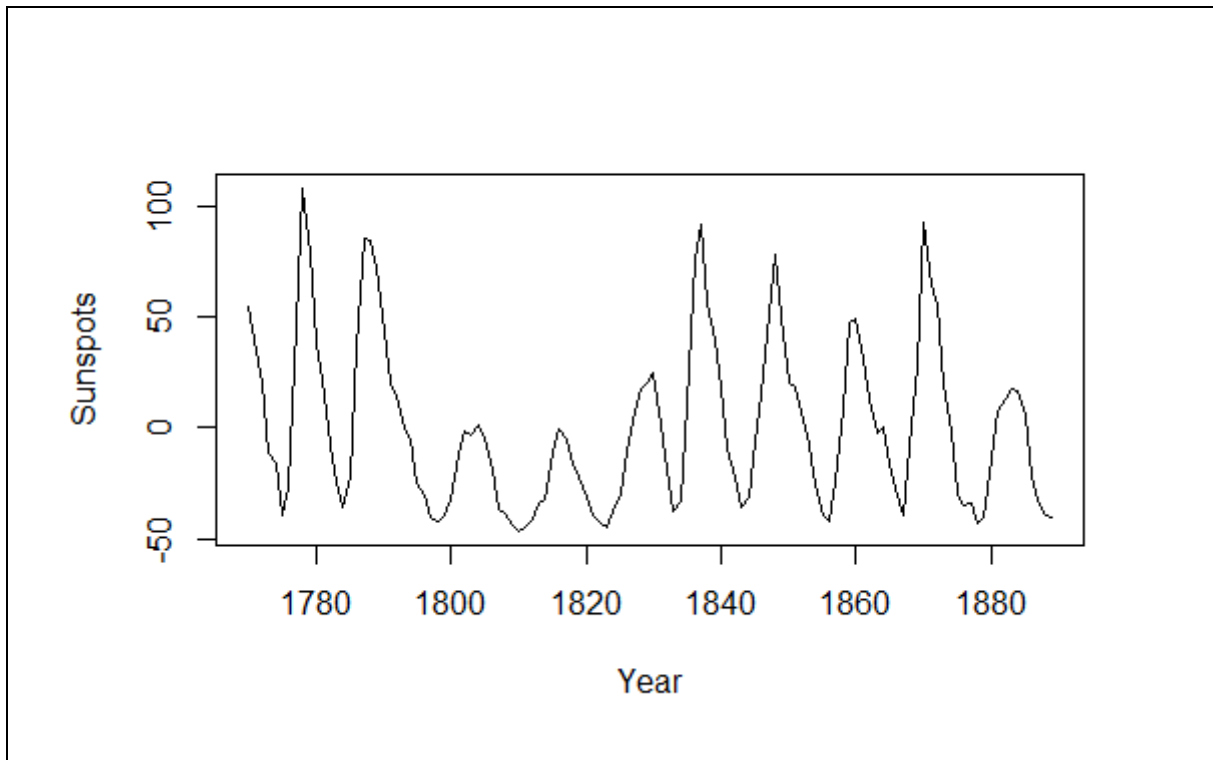
Στην παρακάτω ανάλυση χρησιμοποιήθηκαν δεδομένα που αφορούν τον μέσο αριθμό ηλιακών κηλίδων για τις χρονιές 1770 έως 1889. Τα δεδομένα βρίσκονται στο dataset `wolfer sunspot numbers` του πακέτου `BSDA`.

Οι αστρονόμοι παρατηρούν τον ήλιο και καταγράφουν πληροφορίες σχετικά με τις ηλιακές κηλίδες από την εμφάνιση του τηλεσκοπίου το 1609. Η πρώτη παρατήρηση των ηλιακών κηλίδων έγινε το 1610 από τον Γαλιλαίο και από τότε συνεχίζει η παρακολούθηση αυτού του φαινομένου. Αξίζει να σημειωθεί πως οι ηλιακές κηλίδες είναι παροδικά φαινόμενα τα οποία εμφανίζονται στην επιφάνεια του ηλίου, τη λεγόμενη φωτόσφαιρα. Πρόκειται για μικρές μαύρες περιοχές που θεωρούνται οι περισσότερο εντυπωσιακοί και ενδιαφέροντες σχηματισμοί της φωτόσφαιρας. Ο λόγος που οι ηλιακές κηλίδες φαίνονται μαύρες είναι λόγω της χαμηλής θερμοκρασίας της φωτόσφαιρας που τις περιβάλλει.

Μπορούμε να δούμε παρακάτω την χρονοσειρά του μέσου αριθμού ηλιακών κηλίδων.



Λόγω του ότι η μέση τιμή της χρονοσειράς δεν είναι μηδέν, η μέση τιμή θα αφαιρεθεί από κάθε παρατήρηση της χρονοσειράς ώστε η χρονοσειρά να κεντραριστεί γύρω από το μηδέν.

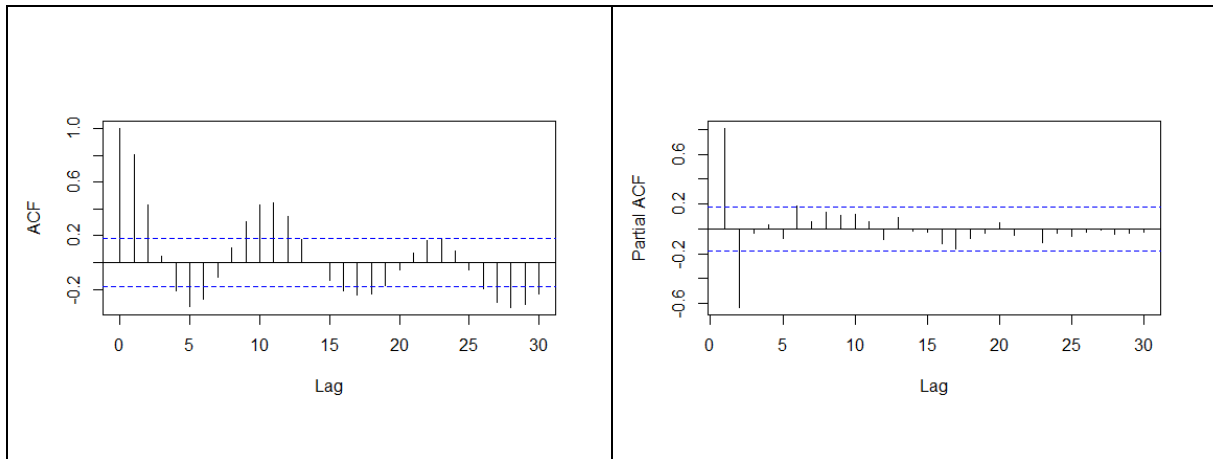


Οι στατιστικές ιδιότητες της χρονοσειράς φαίνεται να παραμένουν σχετικά αναλλοίωτες στον χρόνο, με επιφύλαξη τη διακύμανση, ωστόσο πραγματοποιήσαμε τον έλεγχο Dickey Fuller (ADF) ώστε να ελεγχθεί αν πρόκειται για στάσιμη χρονοσειρά. Ο έλεγχος πραγματοποιήθηκε σε επίπεδο σημαντικότητας $\alpha = 0.05$

H_0 : η χρονοσειρά έχει μοναδιαία ρίζα (unit root) H_1 : η χρονοσειρά είναι στάσιμη

```
Augmented Dickey-Fuller Test
data: sun
Dickey-Fuller = -4.8061, Lag order = 4, p-value = 0.01
alternative hypothesis: stationary
```

Το p-value του ελέγχου είναι 0.01, μικρότερο δηλαδή από το επίπεδο σημαντικότητας 5%, επομένως οδηγούμαστε σε απόρριψη της μηδενικής υπόθεσης της μη στασιμότητας. Συνεπώς, μπορούμε να θεωρήσουμε ότι η χρονοσειράς μας είναι στάσιμη. Πριν από την εφαρμογή των μεθόδων bootstrap στα δεδομένα μας θα απεικονίσουμε το διάγραμμα αυτοσυσχέτισης και μερικής αυτοσυσχέτισης ώστε να αναγνωρίσουμε το καταλληλότερο μοντέλο που ερμηνεύει τα δεδομένα μας.



Μπορούμε να παρατηρήσουμε ότι από το διάγραμμα αυτοσυσχετίσεων ότι οι τιμές των συντελεστών αυτοσυσχετίσης φθίνουν προς το μηδέν ακολουθώντας ημιτονοειδή πορεία. Από το διάγραμμα των μερικών αυτοσυσχετίσεων οι τιμές των συντελεστών μερικής αυτοσυσχετίσης μηδενίζονται απότομα μετά από 2 περιόδους υστέρησης. Η παραπάνω εικόνα μας οδηγεί στο συμπέρασμα ότι η χρονοσειρά μας περιγράφεται ικανοποιητικά από ένα μοντέλο $AR(2)$. Ωστόσο, για την εφαρμογή των μεθόδων bootstrap θα προσαρμοστεί αρχικά μοντέλο $AR(1)$ και εν συνεχεία $AR(2)$. Να σημειώσουμε ότι για την εφαρμογή της μεθόδου Residual Bootstrap έχει χρησιμοποιηθεί η μέθοδος OLS ενώ για τις μεθόδους των blocks η μέθοδος MLE.

Προσαρμόζοντας αυτοπαλίνδρομο μοντέλο πρώτης τάξης

$$Y_t = \varphi_1 Y_{t-1} + e_t,$$

υπό την υπόθεση $e_t \sim iid N(0, \sigma_e^2)$, οι εκτιμήσεις των μεθόδων MLE και OLS δίνονται ακολούθως. Οι εκτιμήσεις των τυπικών σφαλμάτων που βλέπουμε στον Πίνακα 5.1 βασίζονται στην πληροφορία του Fisher.

Πίνακας 5.1

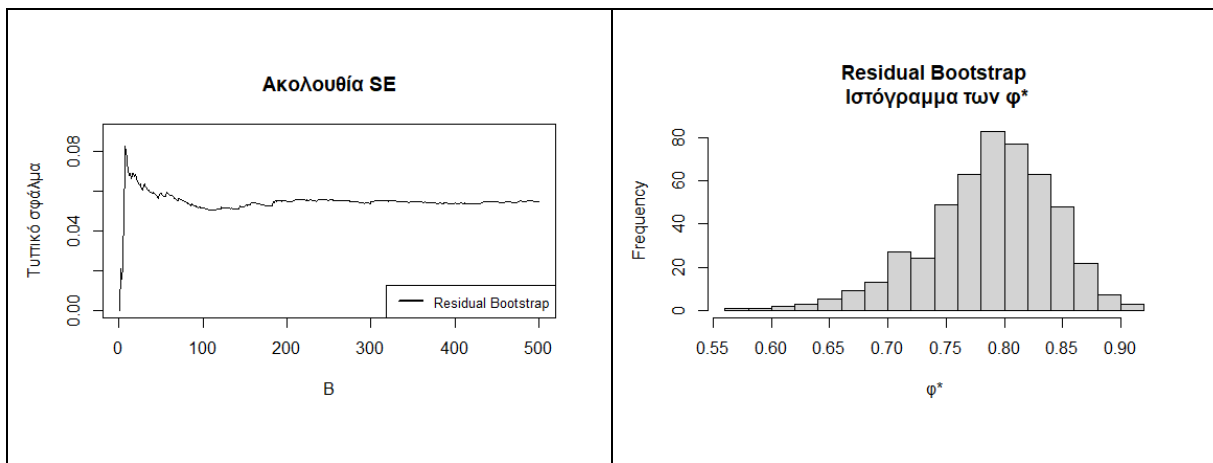
Μοντέλο $AR(1)$			
MLE		OLS	
Εκτιμητής $\hat{\varphi}$	SE	Εκτιμητής $\hat{\varphi}$	SE
0.8224	0.0520	0.8149	0.0524

Από τα παραπάνω αποτελέσματα του Πίνακα 5.1 και συγκεκριμένα από την τιμή του εκτιμητή της μεθόδου MLE, οδηγηθήκαμε στην επιλογή μεγάλου μεγέθους block για τις μεθόδους των blocks, λόγω της αυξημένης εξάρτησης των τυχαίων μεταβλητών στη

χρονοσειρά. Στον Πίνακα 5.2 μπορούμε να δούμε τα αποτελέσματα της εφαρμογής του βάσει μοντέλου bootstrap ενώ στον Πίνακα 5.3 των μεθόδων block bootstrap.

Πίνακας 5.2

Residual Bootstrap για AR(1)	
SE	Bias
0.0546	-0.0279 (0.7869)

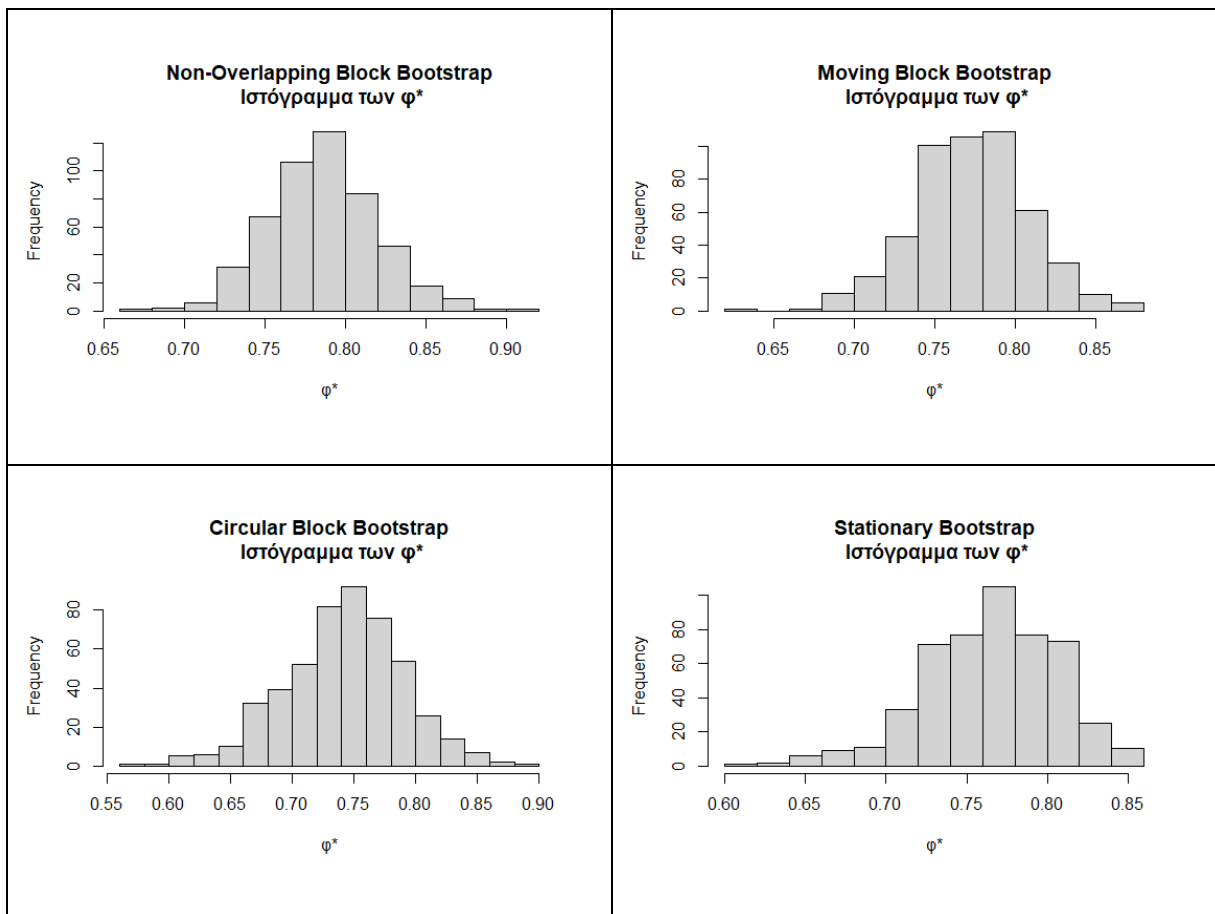
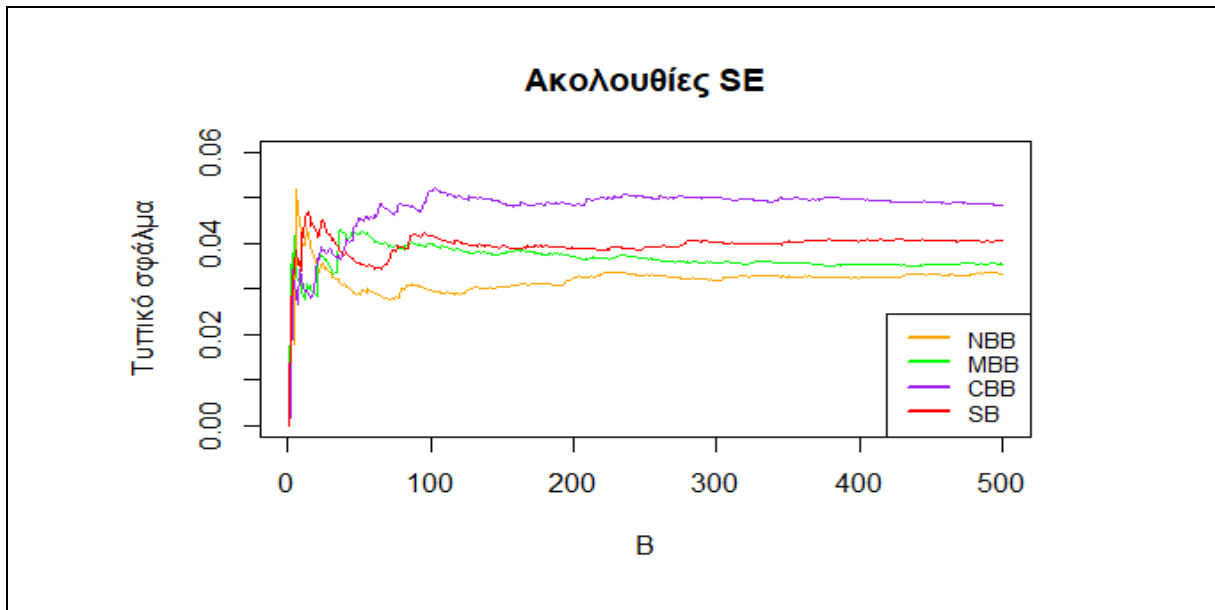


Για τις μεθόδους block bootstrap θα χρησιμοποιηθεί, βάσει των προσομοιώσεων που πραγματοποιήθηκαν στο κεφάλαιο 4, μέγεθος block 25. Επίσης, βασισμένοι στο βέλτιστο μέγεθος block του Carlstein για ένα AR(1) θα εφαρμοστούν οι μέθοδοι και με μέγεθος block 14.

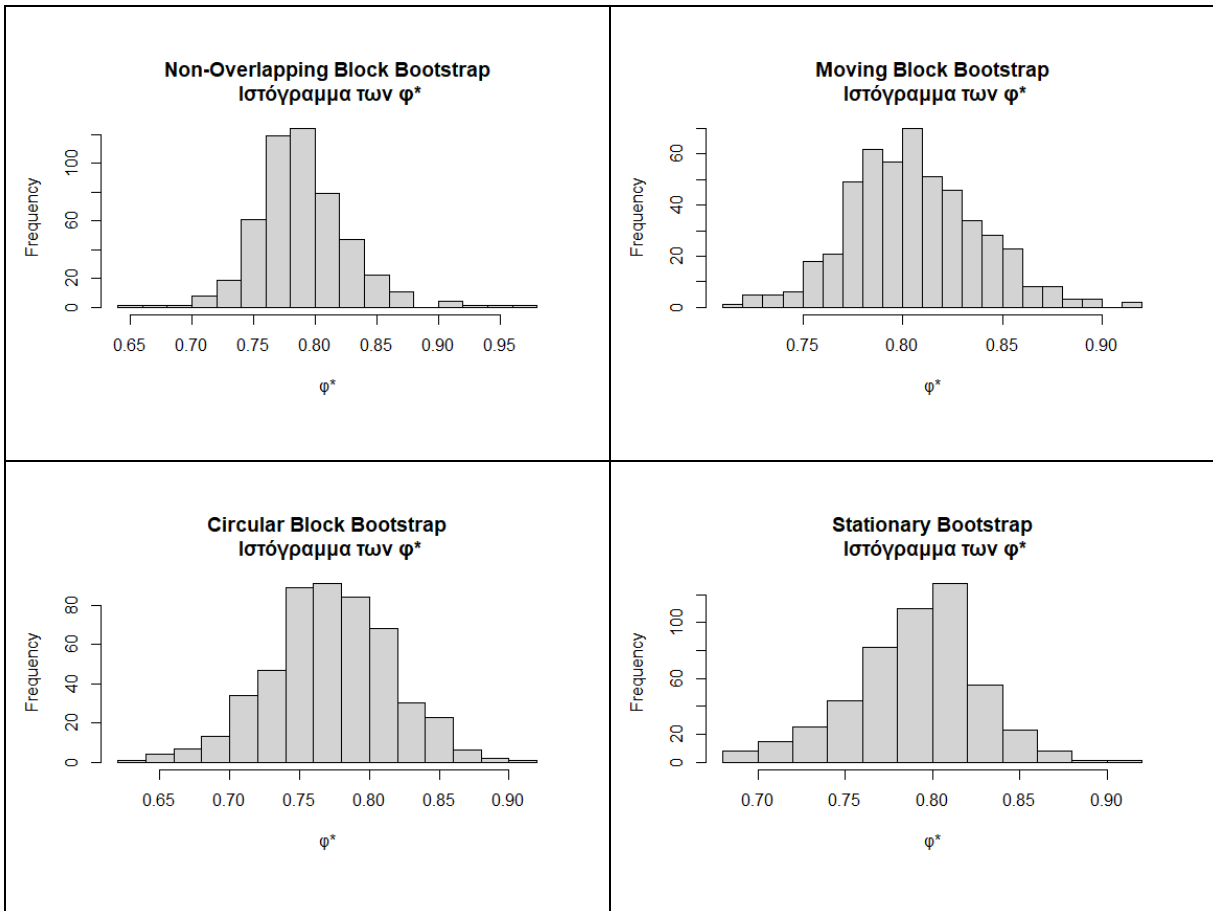
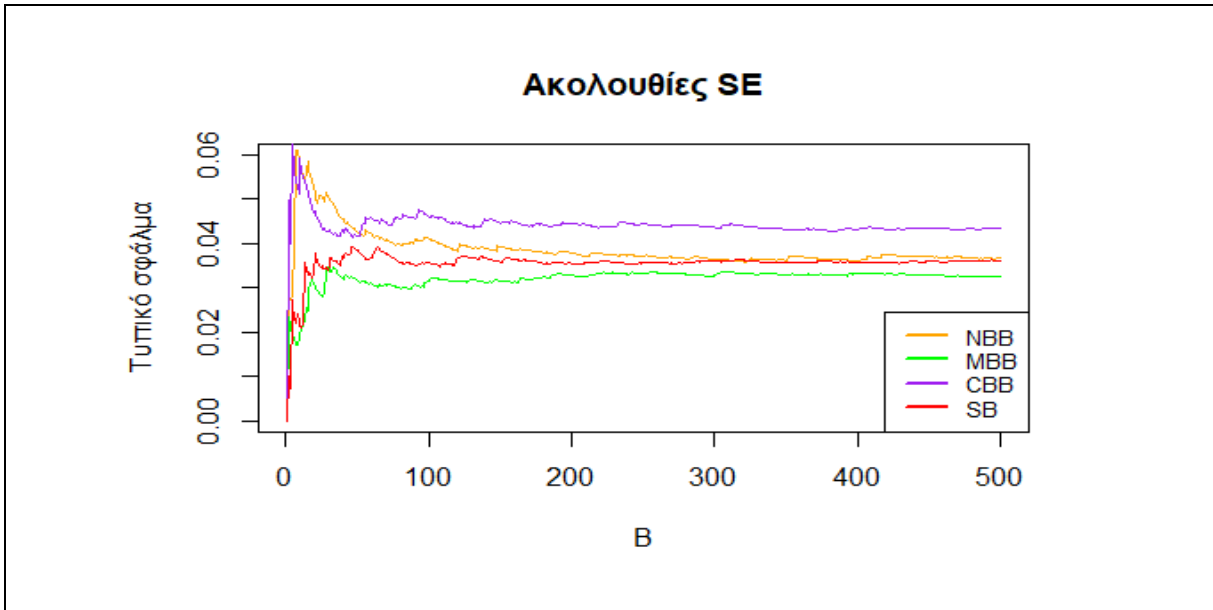
Πίνακας 5.3

ℓ	Μέθοδοι block bootstrap για AR(1)							
	NBB		MBB		CBB		SB	
	SE	Bias	SE	Bias	SE	Bias	SE	Bias
14	0.0334	-0.0364 (0.7860)	0.0355	-0.0501 (0.7722)	0.0484	-0.0798 (0.7426)	0.0407	-0.0575 (0.7649)
25	0.0368	-0.0331 (0.7892)	0.0325	-0.0166 (0.8058)	0.0434	-0.0503 (0.7720)	0.0362	-0.0320 (0.7903)

- Για $\ell = 14$



- Για $\ell = 25$



Αναλύοντας τα παραπάνω αποτελέσματα μπορούμε να διακρίνουμε από τον Πίνακα 5.2 ότι το τυπικό σφάλμα του εκτιμητή της παραμέτρου που προκύπτει από την μέθοδο Residual Bootstrap είναι 0.564, δηλαδή κοντά στην τιμή του τυπικού σφάλματος που προκύπτει από την συνήθη φόρμουλα. Η μεροληψία της μεθόδου Residual Bootstrap είναι -0.0279 . Δηλαδή, η μέση τιμή των 500^{ov} τιμών $\hat{\varphi}^*$, $\hat{\varphi}_1^*$, $\hat{\varphi}_2^*$, ..., $\hat{\varphi}_{500}^*$ (από κάθε προσαρμοσμένη χρονοσειρά bootstrap προκύπτει ένα $\hat{\varphi}^*$), είναι πολύ κοντά στην εκτίμηση $\hat{\varphi}$ της OLS. Για την ακρίβεια είναι 0.7869 ενώ ο OLS εκτιμητής 0.8149. Ακριβώς κάτω από τον Πίνακα 5.2 μπορούμε να δούμε την ακολουθία των τυπικών σφαλμάτων που φαίνεται να σταθεροποιείται λίγο πριν τις 200 επαναλήψεις bootstrap. Από το ιστόγραμμα του Residual Bootstrap μπορούμε να διακρίνουμε αριστερή λοξότητα στην κατανομή των bootstrap εκτιμητών $\hat{\varphi}^*$. Αξίζει να αναφερθεί ότι στο κεντρικό κομμάτι των τιμών του ιστογράμματος φαίνεται να βρίσκεται η εκτίμηση της OLS. Στον Πίνακα 5.3 βλέπουμε τα αποτελέσματα των μεθόδων block bootstrap για $\ell = 14$ και $\ell = 25$ ενώ παρακάτω μπορούμε να δούμε τα γραφήματα που αντιστοιχούν σε κάθε μέγεθος block και κάθε μέθοδο. Να σημειώσουμε ότι η σύγκριση των τυπικών σφαλμάτων για τις μεθόδους των blocks μπορεί να πραγματοποιηθεί παρατηρώντας και τις ακολουθίες των τυπικών σφαλμάτων. Για την μέθοδο NBB αύξηση του μεγέθους των blocks οδηγεί σε μείωση της μεροληψίας αλλά αύξηση του τυπικού σφάλματος. Τα ιστογράμματα φαίνεται να είναι συμμετρικά για την NBB και για τα δύο μεγέθη blocks με ελαφριά δεξιά λοξότητα για $\ell = 25$. Οι εκτιμήσεις της μεθόδου NBB είναι κοντά στην εκτίμηση της MLE ενώ το τυπικό σφάλμα χαμηλότερο. Για την μέθοδο MBB η αύξηση του μεγέθους των blocks σε $\ell = 25$ οδηγεί σε πολύ μεγάλη αναλογικά μείωση της μεροληψίας. Η εκτίμηση της μεθόδου είναι 0.8058 δηλαδή πολύ κοντά στην τιμή 0.8224 του MLE εκτιμητή. Αξίζει να σημειωθεί ότι το κεντρικό κομμάτι των τιμών του ιστογράμματος στο MBB φαίνεται να περιέχει την εκτίμηση της MLE, (ίσως βρίσκεται ελαφρώς αριστερότερα), σε αντίθεση με τις άλλες μεθόδους που το κεντρικό κομμάτι των κατανομών, των εκτιμητών bootstrap $\hat{\varphi}^*$, βρίσκεται αριστερότερα από τον εκτιμητή $\hat{\varphi}$. Η μέθοδος CBB είναι αυτή με την μεγαλύτερη μεροληψία και τυπικό σφάλμα συγκριτικά με τις άλλες μεθόδους. Αυτό μπορούμε να το διακρίνουμε και από τις ακολουθίες των τυπικών σφαλμάτων για $\ell = 14$ και $\ell = 25$ αλλά και από τον Πίνακα 5.3. Αύξηση του ℓ στην μέθοδο SB προκαλεί μεγάλη μείωση της μεροληψίας και συνεπώς μεγάλη βελτίωση στην εκτίμηση bootstrap του $\hat{\varphi}$. Η τιμή της εκτίμησης που δίνει το SB είναι 0.7903 δηλαδή αρκετά κοντά στην MLE εκτίμηση 0.8224. Επίσης, για την SB η αύξηση του ℓ προκαλεί μείωση στο τυπικό σφάλμα και ένα πιο μικρό τυπικό σφάλμα οδηγεί σε πιο μικρή διασπορά των τιμών $\hat{\varphi}^*$. Γενικότερα οι μέθοδοι block bootstrap φαίνεται να υποεκτιμάνε τον εκτιμητή MLE, και η μέθοδος Residual Bootstrap να υποεκτιμάει τον εκτιμητή OLS.

Παρακάτω θα προχωρήσουμε στην εφαρμογή των μεθόδων σε ένα $AR(2)$ μοντέλο. Προσαρμόζοντας αυτοπαλίνδρομο μοντέλο δεύτερης τάξης

$$Y_t = \varphi_1 Y_{t-1} + \varphi_2 Y_{t-2} + e_t,$$

υπό την υπόθεση $e_t \sim iid N(0, \sigma_e^2)$, οι εκτιμήσεις των μεθόδων MLE και OLS δίνονται ακολούθως. Οι εκτιμήσεις των τυπικών σφαλμάτων που βλέπουμε στον Πίνακα 5.4 βασίζονται στην πληροφορία του Fisher.

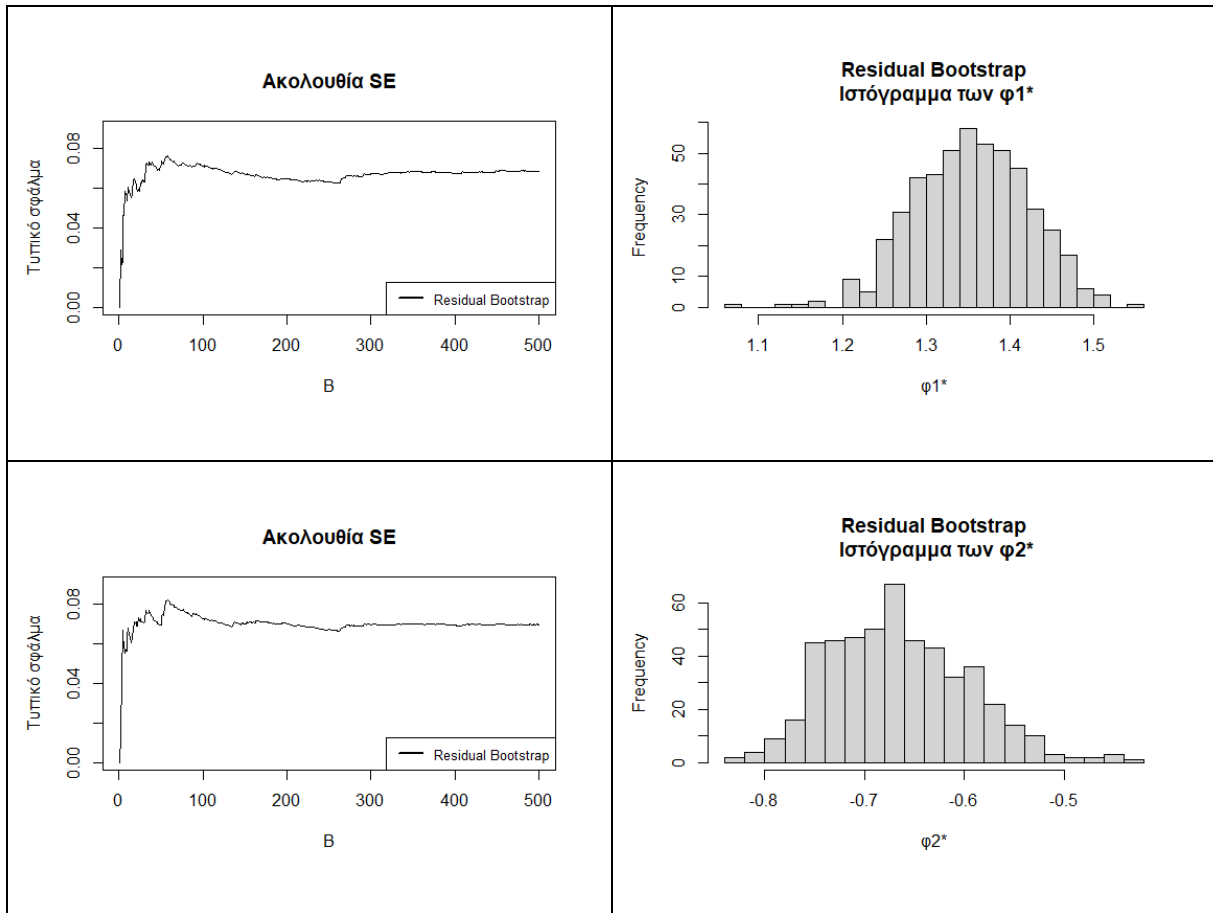
Πίνακας 5.4

Μοντέλο AR(2)							
MLE				OLS			
Εκτιμητής	SE	Εκτιμητής	SE	Εκτιμητής	SE	Εκτιμητής	SE
$\hat{\varphi}_1$	$\hat{\varphi}_1$	$\hat{\varphi}_2$	$\hat{\varphi}_2$	$\hat{\varphi}_1$	$\hat{\varphi}_1$	$\hat{\varphi}_2$	$\hat{\varphi}_2$
1.3755	0.0664	-0.6782	0.0665	1.3726	0.0675	-0.6765	0.0672

Πίνακας 5.5

Residual Bootstrap για AR(2)			
$\hat{\varphi}_1$		$\hat{\varphi}_2$	
SE	Bias	SE	Bias
0.0683	-0.0189 (1.3536)	0.0696	0.0099 (-0.6665)

Στον Πίνακα 5.5 μπορούμε να δούμε τα αποτελέσματα της μεθόδου Residual Bootstrap τα οποία είναι πολύ κοντά με τα αποτελέσματα της μεθόδου ελαχίστων τετραγώνων του Πίνακα 5.4. Από τα παρακάτω γραφήματα που αφορούν την μέθοδο Residual Bootstrap, πάνω αριστερά μπορούμε να διακρίνουμε την ακολουθία των τυπικών σφαλμάτων του $\hat{\varphi}_1$ η οποία δείχνει να σταθεροποιείται σχετικά νωρίς. Το ίδιο φαίνεται να συμβαίνει και στην ακολουθία τυπικών σφαλμάτων του εκτιμητή $\hat{\varphi}_2$. Επιπρόσθετα, μπορούμε να παρατηρήσουμε το παραγόμενο ιστόγραμμα των $\hat{\varphi}_1^{*(1)}, \hat{\varphi}_1^{*(2)}, \dots, \hat{\varphi}_1^{*(500)}$ πάνω δεξιά και το ιστόγραμμα των $\hat{\varphi}_2^{*(1)}, \hat{\varphi}_2^{*(2)}, \dots, \hat{\varphi}_2^{*(500)}$ από κάτω του. Ο επιπλέον όρος που προσθέσαμε στο AR(2) μοντέλο φαίνεται να μειώνει την λοξότητα των $\hat{\varphi}_1^*$ που είδαμε στην περίπτωση του AR(1) μοντέλου. Το ιστόγραμμα των $\hat{\varphi}_2^*$ φαίνεται να παρουσιάζει δεξιά λοξότητα.

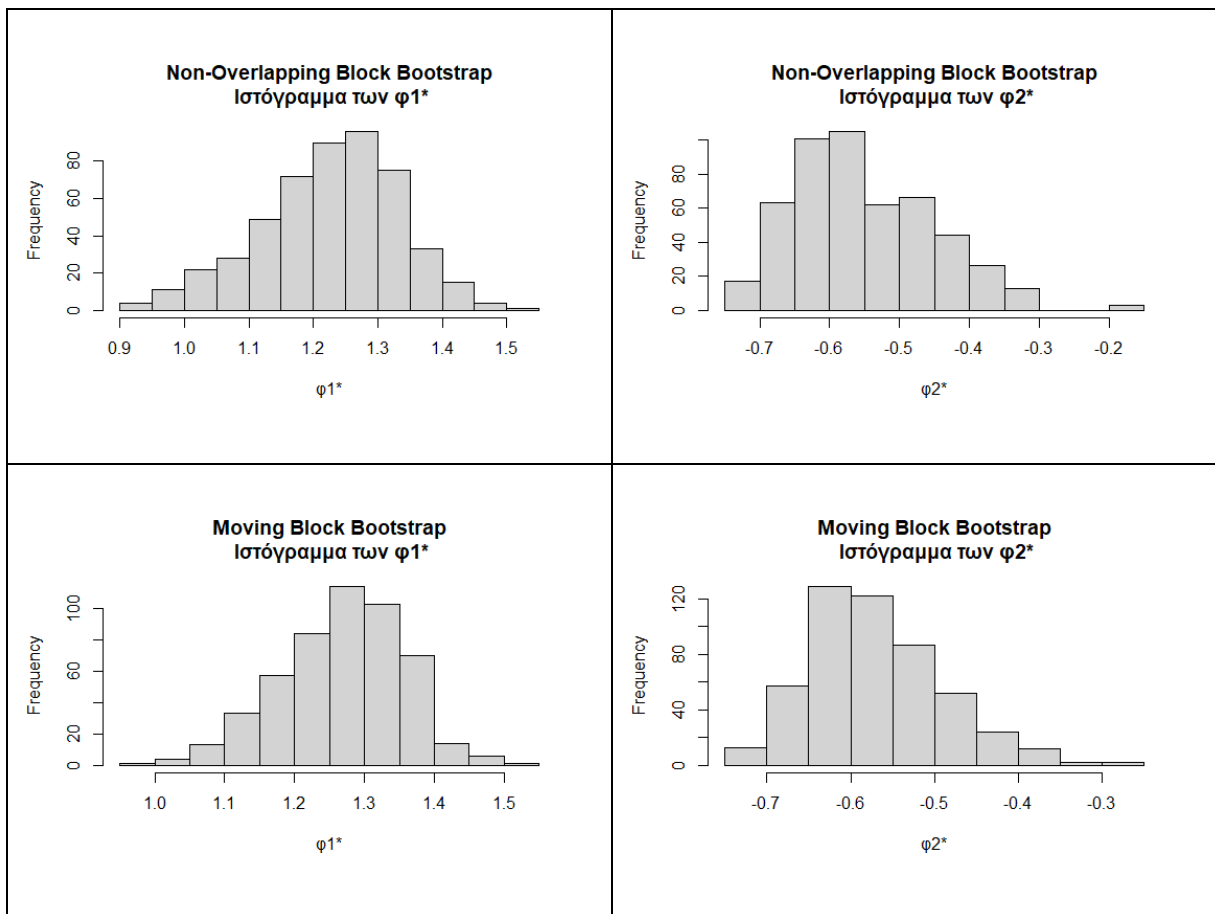
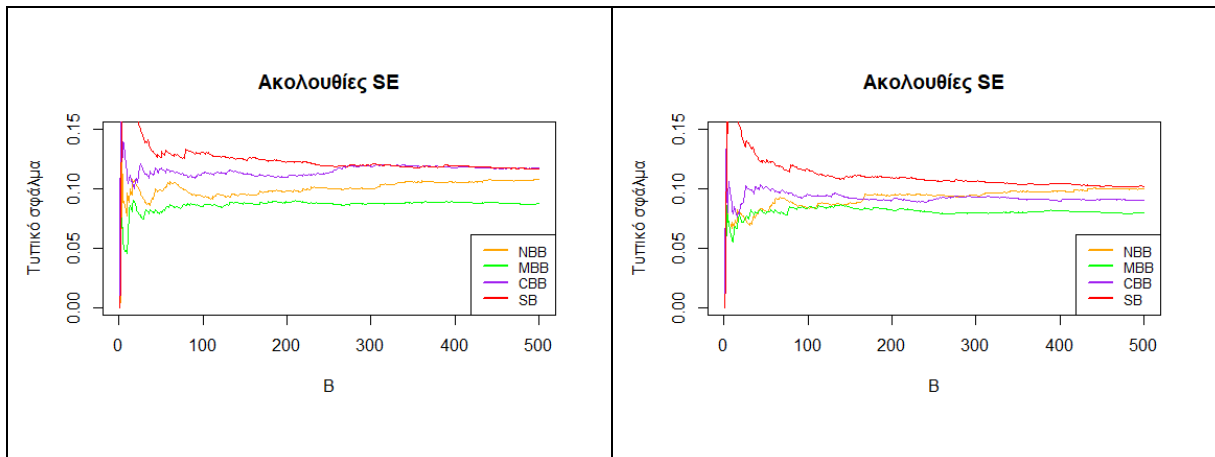


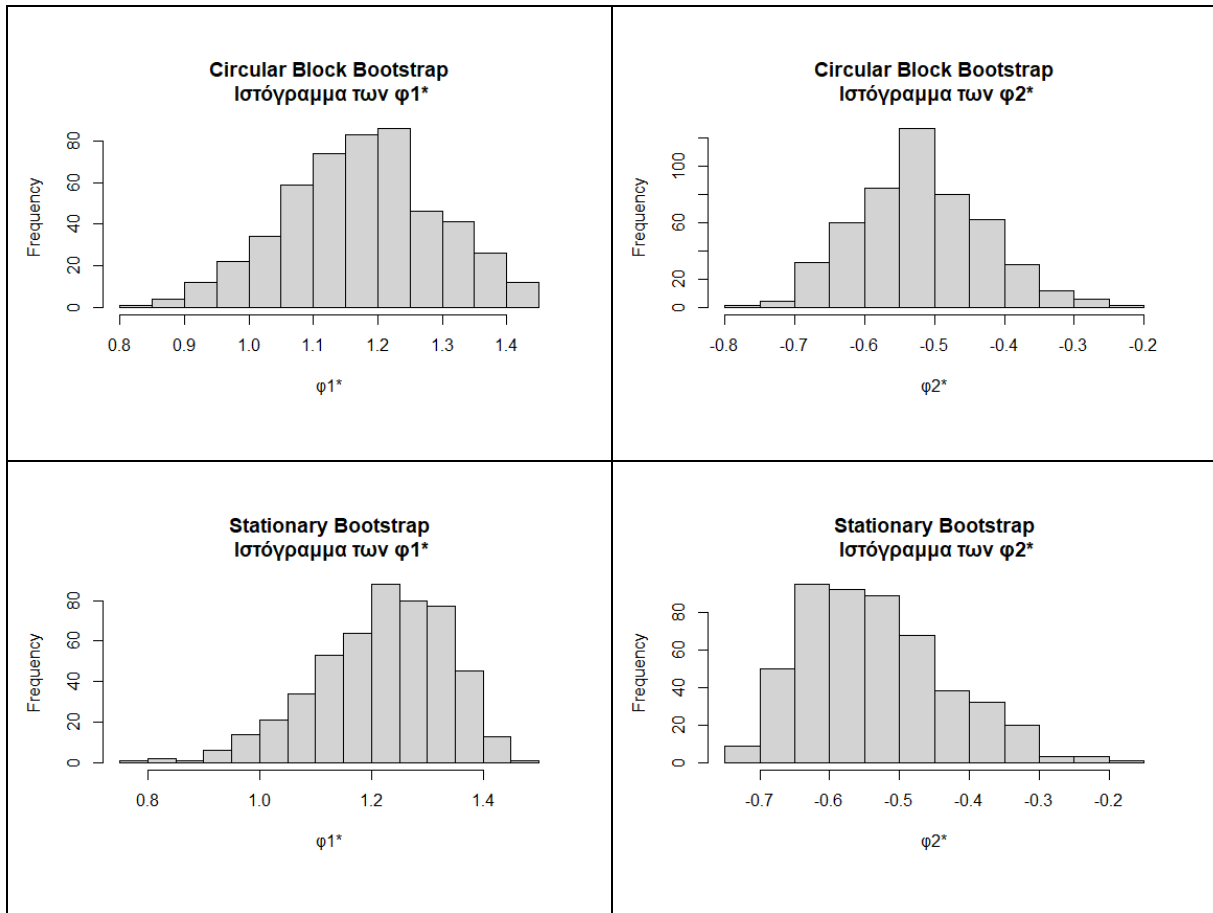
Πίνακας 5.6

	Μέθοδοι block bootstrap για AR(2)							
	NBB		MBB		CBB		SB	
$\ell =$	SE	Bias	SE	Bias	SE	Bias	SE	Bias
25								
$\hat{\phi}_1$	0.1078	-0.1493 (1.2260)	0.0873	-0.1066 (1.2688)	0.1170	-0.2023 (1.173)	0.1164	-0.1581 (1.217)
$\hat{\phi}_2$	0.0997	0.1283 (-0.5498)	0.0795	0.1083 (-0.5698)	0.0902	0.1591 (-0.5190)	0.1019	0.1425 (-0.5356)

Από τον Πίνακα 5.6 βλέπουμε ότι οι μέθοδοι των blocks δεν έχουν τόσο καλή απόδοση όσο η μέθοδος Residual Bootstrap. Παρατηρώντας τον Πίνακα 5.6 μπορούμε να παρατηρήσουμε ότι για $\ell = 25$, η μέθοδος MBB έχει το μικρότερο τυπικό σφάλμα και μεροληψία συγκριτικά με τις άλλες μεθόδους. Η εκτίμηση της μεθόδου MLE για το $\hat{\phi}_1$ είναι 1.3755 ενώ η εκτίμηση

bootstrap του $\hat{\varphi}_1$ από την μέθοδο MBB είναι 1.2688. Αξίζει να σημειωθεί, ότι η μέθοδος SB εμφανίζει την μεγαλύτερη μεροληψία. Παρακάτω παρατηρούμε τις ακολουθίες των τυπικών σφαλμάτων των εκτιμητών. Στα αριστερά βρίσκεται η περίπτωση του $\hat{\varphi}_1$ και δεξιά του $\hat{\varphi}_2$. Όπως μπορούμε να διακρίνουμε και από τον πίνακα η μέθοδος SB εμφανίζει μεγαλύτερο τυπικό σφάλμα για το $\hat{\varphi}_2$ συγκριτικά με τις άλλες μεθόδους ενώ περίπου ίδιο με την CBB για το $\hat{\varphi}_1$.





Από τον Πίνακα 5.6 αλλά και τα παραπάνω συγκρίσιμα ιστογράμματα γίνεται εμφανές ότι οι μέθοδοι block bootstrap υποεκτιμάνε τον εκτιμητή $\hat{\varphi}_1$ και υπερεκτιμάνε τον εκτιμητή $\hat{\varphi}_2$. Επίσης, αξίζει να σημειωθεί ότι τα ιστογράμματα των $\hat{\varphi}_1^*$ εμφανίζουν αριστερή λοξότητα σε όλες τις μεθόδους των blocks. Συμπερασματικά μπορούμε να πούμε ότι καταλληλότερη μέθοδος bootstrap κρίνεται η μέθοδος Residual Bootstrap κάτι το οποίο και περιμέναμε καθώς το παραμετρικό μοντέλο που ερμήνευε τα δεδομένα μας φάνηκε να είναι ένα $AR(2)$. Μπορεί οι μέθοδοι block bootstrap να μην επιτυγχάνουν την ακρίβεια της Residual Bootstrap ωστόσο πρόκειται για μεθόδους που μας παρέχουν πολύ ικανοποιητικά αποτελέσματα. Η ευαισθησία των μεθόδων των blocks στην επιλογή του ℓ είναι μεγάλη και επηρεάζει την απόδοση των μεθόδων γι' αυτό και κρίνεται πολύ σημαντική η επιλογή του, και όπως ήδη έχει αναφερθεί το ιδανικό μέγεθος block εξαρτάται από το μέγεθος του αρχικού δείγματος και την δομή της συσχέτισης του και διαφέρει για διαφορετικές μεθόδους block bootstrap.

Παραρτήματα

Π1 Συναρτήσεις μεθόδων και μελέτη προσομοίωσης

```
library(stats)
library(tseries)
library(BSDA)
```

#nbb function για δείκτες μιας χρονοσειράς nbb

```
blockfun<-function(data,lblock){
  block_index_matrix<-matrix(nrow=ceiling(length(data)/lblock),ncol=lblock) #καθε γραμμη
  αποτελεί ενα block (τους δεικτες του)
  for (i in seq(0,length(data)-1,lblock)){
    block_index<-(i+1):(i+lblock)
    if (i==seq(0,length(data)-1,lblock)[length(seq(0,length(data)-1,lblock))] &&
(length(data)+1) %in% block_index){
      block_index<-block_index[block_index<=length(data)]
    }
    while(length(block_index)!=lblock){block_index<-c(block_index,NA)}
    block_index_matrix[i/lblock + 1, ]<-block_index
  }
  #print(block_index_matrix) #πινακας με nnb blocks
  m<-sample(1:ceiling(length(data)/lblock),ceiling(length(data)/lblock),replace=T)
  #print(m)
  newseries_index<-c(t(block_index_matrix[m,]))
  newseries_index<-newseries_index[!is.na(newseries_index)]
  return(newseries_index)
}
#x<-arima.sim(model=list(ar=0.3),n=100) #αν x δεδομένα (εδώ ar(1) με φ=0.3)
#blockfun(x,4)# δεικτες μιας nbb χρονοσειρας με blocklength 4
#x[blockfun(x,4)]#μια nbb χρονοσειρα
```

#mbb function για δείκτες μιας χρονοσειράς mbb

```
mbb_boot_index<-function(data,lamda){
  N=length(data)-lamda+1;N #ποσα συνοδικα mbb (overlapping blocks) δημιουργουνται
  indexes<-c() #δεικτες στοιχειων αρχικης χρονοσειρας
  for (i in 1:length(data)){
    indexes<-c(indexes,i)
  }
  mbb_index_matrix<-matrix(nrow=N,ncol=lamda)#καθε γραμμη αποτελεί ενα block
  μεγεθους lamda (τους δεικτες του)
  j=0
  for (i in 1:N){ #απο το block=1 εως το block=N
    mbb_index_matrix[i,]<-indexes[1+j]:indexes[lamda+j]
    j=j+1
  }
}
```

```

}
k<-ceiling(length(data)/lamda)#ποσα blocks θα τραβηξουμε (μεγεθους lamda)
#print(mbb_index_matrix) #πινακας με mbb blocks
m<-sample(1:N,k,replace=T)
#print(m)
mbb_newseries_index<-c(t(mbb_index_matrix[m,]))
return(mbb_newseries_index)
}
#mbb_boot_index(x,4) # δεικτες μιας mbb χρονοσειρας με blocklength 4
#x[mbb_boot_index(x,4)] #mbb χρονοσειρα

```

#cbb function για δείκτες μιας χρονοσειράς cbb

```

cbb_boot_index<-function(data,lamda){
  N=length(data)-lamda+1;N #ποσα συνολικα mbb (overlapping blocks) δημιουργουνται
  indexes<-c(1:length(data)) #δεικτες στοιχειων αρχικης χρονοσειρας
  cbb_index_matrix<-matrix(nrow=length(data),ncol=lamda);cbb_index_matrix
  j=0
  for (i in 1:N){ #απο το block=1 εως το block=N
    #print(indexes[1+j]:indexes[lamda+j])
    cbb_index_matrix[i,]<-indexes[1+j]:indexes[lamda+j]
    j=j+1
  }
  #γεμισαμε τον πινακα με τα mbb block, τωρα θα γεμισουμε με cbb
  for (t in (N+1):length(data)){
    cbb_index_matrix[t,]<-c(c(indexes[t]:indexes[length(data)]),c(indexes[1]:indexes[lamda-
length(data)+t-1])) #εφαρμογη θεωριας (xt..xn,x1..x1-n+t-1)
    #print(cbb_index_matrix) #πινακας cbb blocks
  }
  k<-ceiling(length(data)/lamda)#ποσα blocks θα τραβηξουμε (μεγεθους lamda) .
  m<-sample(1:length(data),k,replace=T) #ειναι διαφορετικο απο του mbb καθως τωρα
  τραβαμε απο 1:length(data) πριν απο 1:N
  cbb_newseries_index<-c(t(cbb_index_matrix[m,]))
  return(cbb_newseries_index)
}
#cbb_boot_index(x,4) #δεικτες μιας cbb χρονοσειρας με blocklength 4
#x[cbb_boot_index(x,4)] #cbb χρονοσειρα

```

#sb function

#ένα sbb block

```

sbb_block_index<-function(data,lamda){
  i<-sample(1:length(data),1)
  #xi=x(1+(i-1 mod n))
  start_point<-1+((i)-1%%(length(data)))
  #γεωμετρικη
  p<-1/lamda
  blocklen_geom<-rgeom(1,prob=p)+1
  #αν τυχει και η γεωμετρικη δωσει blocklen τοσο μεγαλο που με ενα θα καλυψουμε ολη τη
  σειρα θα το τρεχουμε μεχρι να παρουμε μικροτερο

```



```

while (blocklen_geom>=length(data)){
  blocklen_geom<-rgeom(1,prob=p)+1
}
#(xi,..., x(i+lamda-1)) ενα block
sbb_block<-start_point:(start_point+blocklen_geom-1)
if ((length(data)+1)%in%sbb_block){ #αν οι δείκτες βγαιουν εκτος του πληθους στοιχειων
της χρονοσειρας, τοτε κραταμε <length(data)
  sbb_block<-sbb_block[sbb_block<=length(data)]
}
return(sbb_block)#ενα block sbb
}
#sbbfun για δείκτες μιας χρονοσειράς sbb
sbbfun<-function(data1, lamda1){
  sbb_boot_index<-c()
  while (length(sbb_boot_index)<length(data1)){
    sbb_boot_index<-c(sbb_boot_index,sbb_block_index(data1, lamda1))
  }
  return(sbb_boot_index)
}
#sbbfun(x,4) #δείκτες μιας sbb χρονοσειράς με μέσο μηκος block 4
#x[sbbfun(x,4)] #sbb χρονοσειρα

```

residual bootstrap function επιστρέφει $B \hat{\varphi}^*$.hat εκτιμήσεις bootstrap

```

# resboot_fun(order_μοντελου,intercept=0 (δεν εχει σταθερο),δεδομενα, επαναληψεις
bootstrap)
resboot_fun<-function(ord,inter,timedata,Brep){
  #set.seed(10)
  data<-timedata #δεδομενα
  if (inter==0){ #αν δεν εχει σταθερο ορο προσαρμοζουμε μοντελο με ols χωρις σταθερο ορο
  arfit1<-ar.ols(data,order.max = ord,demean =F,intercept = F,aic=F)#ols
  }else{
  arfit1<-ar.ols(data,order.max = ord,demean =F,intercept = T,aic=F)
  } #συνεπως στο arfit1 ειναι το προσαρμοσμένο μοντέλο ταξεως που επιθυμούμε (απο arfit1
θα παρουμε residuals)
#bootstrap στα residuals
B<-Brep;
bootstrap_samples<-matrix(0,nrow=B,ncol=length(data))#πινακα Bxlength(data) (καθε
γραμμη θα ειναι μια boot χρονοσειρα)
bootstrap_samples
for (b in 1:B){
  bootresid<-sample(arfit1$resid[!is.na(arfit1$resid)],replace=T) #bootstrap δειγμα
καταλοιπων (εχω αφαιρεσει NA τιμες)
  bootresid<-bootresid-mean(bootresid) #κεντραρισμα
  bootstrap_data<-numeric(length(data))#φτιαχνουμε ενα διανυσμα μηδενικων το οποιο οταν
το γεμισουμε θα ειναι η νεα χρονοσειρα bootstrap
  bootstrap_data[1:(arfit1$order)]<-data[1:(arfit1$order)]# p πρωτες τιμες της bootstrap
χρονοσειρας θα ειναι ιδια με της αρχικης χρονοσειρας

```

```

bootresid<-c(rep(0,arfit1$order),bootresid) #το εβαλα ετσι για να ταιριαζει με την θεωρια
και να ξεκινησει απο e*(t) η αναδρομη (οι προηγουμ τιμες e*(t-1),.. θα ειναι 0)
#δεν θα χρησιμοποιησω τα 0 που εβαλα (οσα και αν ειναι)
for (t in (arfit1$order+1):length(data)){ #ξεκιναι η κατασκευη της (b) χρονοσειρας
bootstrap
  bootstrap_y<-0 #η τιμη αυτη μετα απο αναδρομη θα δωσει το y*(t) της εκαστοτε
χρονικης στιγμης
  for (i in 1:arfit1$order) { #υπολογιζεται τιμη χρονικης στιγμης t με βαση την
προηγουμενη/προηγουμενες αν τρεξει το loop πολλες φορες (εξαρταται την ταξη order)
    bootstrap_y<-bootstrap_y+arfit1$ar[i]*bootstrap_data[t-i] #
  }
  if (inter==0){
    bootstrap_data[t]<-bootstrap_y+bootresid[t] #βαζουμε στη χρονοσειρα bootstrap στη
θεση t την τιμη y*(t)
  }else{
    bootstrap_data[t]<-arfit1$x.intercept+bootstrap_y+bootresid[t] #με σταθερο ορο
  }
}
bootstrap_samples[b,]<-bootstrap_data #βαζουμε στον πινακα στην θεση b (γραμμη) την
bootstrap χρονοσειρα
}
bootstrap_phis<-matrix(0,nrow=B,ncol=arfit1$order) #πινακας οπου καθε γραμμη θα ειναι
το/τα φ που προκυπτει απο καθε χρονοσειρα bootstrap
for (i in 1:nrow(bootstrap_samples)){
  if (inter==0){
    bootstrap_phis[i,]<-unname(ar(bootstrap_samples[i,],method = "ols",order.max =
arfit1$order,demean =F,intercept = F,aic=F)$ar)
  }else{
    bootstrap_phis[i,]<-unname(ar(bootstrap_samples[i,],method = "ols",order.max =
arfit1$order,demean =F,intercept =T,aic=F)$ar)
  }
}
return(bootstrap_phis)
}
#res_phis<-resboot_fun(1,0,x,500);res_phis #πρώτης ταξης, δεν εχει σταθερο, δεδομενα, B
#sd(res_phis)
#mean(res_phis)

```

#Monte Carlo

```

M<-500
n<-200
errors<-matrix(rnorm(M*n),nrow=M);errors
dim(errors)

```

#nbb

```

mc_nbb<-function(M_nbb,n_nbb,phi_nbb,B_nbb,blocklen_nbb){
  set.seed(10)
  monte_carlo_nbb_se<-c() #1

```

```

monte_carlo_nbb_mean<-c() #2
monte_carlo_nbb_bias<-c() #3
sum<-0
for (m in 1:M_nbb){
  #y<-arima.sim(model=list(ar=phi_nbb),n_nbb)
  y<-c(errors[m,1],rep(0,n_nbb-1))
  for (i in 2:n_nbb){
    y[i]<-y[i-1]*phi_nbb+errors[m,i]
  }
  nbb_phis_star<-c()
  for (b in 1:B_nbb){
    nbb_indexes<-blockfun(y,blocklen_nbb) #συναρτηση που επιστρεφει τους δεικτες για να
φτιαξουμε 1 χρονοσειρα nbb
    nbb_phis_star<-c(nbb_phis_star,unname(arima(y[nbb_indexes], order = c(1, 0, 0),
include.mean = F)$coef))
    #nbb_phis_star<-c(nbb_phis_star,unname(ar(y[nbb_indexes],method="mle",order.max =
1,demean =F,intercept = F,aic=F)$ar))
  }
  monte_carlo_nbb_se<-
c(monte_carlo_nbb_se,sd(nbb_phis_star)*(((ceiling(n_nbb/blocklen_nbb)*blocklen_nbb)/n_n
bb))*(1/2))) #προσαρμογη se σε πιθανο διαφ.μεγεθος
  monte_carlo_nbb_mean<-c(monte_carlo_nbb_mean,mean(nbb_phis_star))
  bias<-mean(nbb_phis_star)-unname(arima(y, order = c(1, 0, 0), include.mean = F)$coef)
  #bias<-mean(nbb_phis_star)-unname(ar(y,method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar)
  monte_carlo_nbb_bias<-c(monte_carlo_nbb_bias,bias)
  percentile_ci<-unname(c(quantile(nbb_phis_star,0.05),quantile(nbb_phis_star,0.95)))
  if (phi_nbb>percentile_ci[1]&&phi_nbb<percentile_ci[2]){
    sum=sum+1
  }
}
}
prob<-sum/M_nbb

return(list(mean(monte_carlo_nbb_se),mean(monte_carlo_nbb_mean),mean(monte_carlo_nb
b_bias),prob))
}
#προσαρμόζεται με τα επιθυμητά ορίσματα
#M,n,phi,B,blocklen
#mc_nbb(500,100,0.5,200,8) #mc=500,n=100,φ=0.5,B=200,lamda=8

#προσαρμόζεται με τα επιθυμητά ορίσματα
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_nbb(500,200,f,200,4)
  print(result)
  cat("#####\n")
}

```

#mbb

```
mc_mbb<-function(M_mbb,n_mbb,phi_mbb,B_mbb,blocklen_mbb){
  set.seed(10)
  mbb_matrix<-matrix(nrow=M_mbb,ncol=3)
  colnames(mbb_matrix)<-c("LowerBound","UpperBound","φ in CI?")
  monte_carlo_mbb_se<-c() #1
  monte_carlo_mbb_mean<-c() #2
  monte_carlo_mbb_bias<-c() #3
  sum<-0
  for (m in 1:M_mbb){
    #y<-arima.sim(model=list(ar=phi_mbb),n_mbb)
    y<-c(errors[m,1],rep(0,n_mbb-1))
    for (i in 2:n_mbb){
      y[i]<-y[i-1]*phi_mbb+errors[m,i]
    }
    mbb_phis_star<-c()
    for (b in 1:B_mbb){
      mbb_indexes<-mbb_boot_index(y,blocklen_mbb) #συναρτηση που επιστρεφει τους
      δεικτες για να φτιαξουμε 1 χρονοσειρα mbb
      mbb_phis_star<-c(mbb_phis_star,unname(arima(y[mbb_indexes], order = c(1, 0, 0),
      include.mean = F)$coef))
      #mbb_phis_star<-c(mbb_phis_star,unname(ar(y[mbb_indexes],method="mle",order.max
      = 1,demean = F,intercept = F,aic=F)$ar))
    }
    monte_carlo_mbb_se<-
    c(monte_carlo_mbb_se,sd(mbb_phis_star)*(((ceiling(n_mbb/blocklen_mbb)*blocklen_mbb)/
    n_mbb)**(1/2)))#προσαρμογη σε σε πιθανο διαφ.μεγεθος
    monte_carlo_mbb_mean<-c(monte_carlo_mbb_mean,mean(mbb_phis_star))
    bias<-mean(mbb_phis_star)-unname(arima(y, order = c(1, 0, 0), include.mean = F)$coef)
    #bias<-mean(mbb_phis_star)-unname(ar(y,method="mle",order.max = 1,demean
    =F,intercept = F,aic=F)$ar)
    monte_carlo_mbb_bias<-c(monte_carlo_mbb_bias,bias)
    #percentile_ci<-unname(c(quantile(mbb_phis_star,0.025),quantile(mbb_phis_star,0.975)))
    percentile_ci<-unname(c(quantile(mbb_phis_star,0.05),quantile(mbb_phis_star,0.95)))
    mbb_matrix[m,1]<-round(percentile_ci[1],4)
    mbb_matrix[m,2]<-round(percentile_ci[2],4)
    if (phi_mbb>percentile_ci[1]&&phi_mbb<percentile_ci[2]){
      sum=sum+1
      mbb_matrix[m,3]<-"Yes"
    }else{
      mbb_matrix[m,3]<-"No"
    }
  }
  }
  prob<-sum/M_mbb

  return(list(mean(monte_carlo_mbb_se),mean(monte_carlo_mbb_mean),mean(monte_carlo_m
  bb_bias),prob))#,prob,mbb_matrix))
}
```

```

#M,n,phi,B,blocklen
#προσαρμολογείται με τα επιθυμητα ορισματα
#mc_mbb(500,100,0.5,200,8) #mc=500,n=100,φ=0.5,B=200,lamda=8
#προσαρμολογείται με τα επιθυμητα ορισματα
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_mbb(500,200,f,200,4)
  print(result)
  cat("#####\n")
}

#cbb
mc_cbb<-function(M_cbb,n_cbb,phi_cbb,B_cbb,blocklen_cbb){
  set.seed(10)
  monte_carlo_cbb_se<-c() #1
  monte_carlo_cbb_mean<-c() #2
  monte_carlo_cbb_bias<-c() #3
  sum<-0
  for (m in 1:M_cbb){
    y<-c(errors[m,1],rep(0,n_cbb-1))
    for (i in 2:n_cbb){
      y[i]<-y[i-1]*phi_cbb+errors[m,i]
    }
    cbb_phis_star<-c()
    for (b in 1:B_cbb){
      cbb_indexes<-cbb_boot_index(y,blocklen_cbb) #συναρτηση που επιστρεφει τους δεικτες
      για να φτιαξουμε 1 χρονοσειρα cbb
      cbb_phis_star<-c(cbb_phis_star,unname(arima(y[cbb_indexes], order = c(1, 0, 0),
include.mean = F)$coef))
      #cbb_phis_star<-c(cbb_phis_star,unname(ar(y[cbb_indexes],method="mle",order.max =
1,demean = F,intercept = F,aic=F)$ar))
    }
    monte_carlo_cbb_se<-
c(monte_carlo_cbb_se,sd(cbb_phis_star)*(((ceiling(n_cbb/blocklen_cbb)*blocklen_cbb)/n_c
bb)**(1/2))) #προσαρμογη σε σε πιθανο διαφ.μεγεθος
    monte_carlo_cbb_mean<-c(monte_carlo_cbb_mean,mean(cbb_phis_star))
    bias<-mean(cbb_phis_star)-unname(arima(y, order = c(1, 0, 0), include.mean = F)$coef)
    #bias<-mean(cbb_phis_star)-unname(ar(y,method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar)
    monte_carlo_cbb_bias<-c(monte_carlo_cbb_bias,bias)
    #percentile_ci<-unname(c(quantile(cbb_phis_star,0.025),quantile(cbb_phis_star,0.975)))
    percentile_ci<-unname(c(quantile(cbb_phis_star,0.05),quantile(cbb_phis_star,0.95)))
    if (phi_cbb>percentile_ci[1]&&phi_cbb<percentile_ci[2]){
      sum=sum+1
    }
  }
}
prob<-sum/M_cbb

```

```

return(list(mean(monte_carlo_cbb_se),mean(monte_carlo_cbb_mean),mean(monte_carlo_cbb
_bias),prob))
}

```

```

#M,n,phi,B,blocklen
#mc_cbb(500,100,0.5,200,8) #mc=500,n=100,φ=0.5,B=200,lamda=8
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_cbb(500,200,f,200,4)
  print(result)
  cat("#####\n")
}

```

#sb

```

mc_sb<-function(M_sb,n_sb,phi_sb,B_sb,blocklen_sb){
  set.seed(10)
  monte_carlo_sb_se<-c() #1
  monte_carlo_sb_mean<-c() #2
  monte_carlo_sb_bias<-c() #3
  sum<-0
  for (m in 1:M_sb){
    y<-c(errors[m,1],rep(0,n_sb-1))
    for (i in 2:n_sb){
      y[i]<-y[i-1]*phi_sb+errors[m,i]
    }
    sb_phis_star<-c()
    for (b in 1:B_sb){
      sb_indexes<-sbbfun(y,blocklen_sb) #συναρτηση που επιστρεφει τους δεικτες για να
φτιαξουμε 1 χρονοσειρα sb
      sb_phis_star<-c(sb_phis_star,unname(arima(y[sb_indexes], order = c(1, 0, 0),
include.mean = F)$coef))
      #sb_phis_star<-c(sb_phis_star,unname(ar(y[sb_indexes],method="mle",order.max =
1,demean =F,intercept = F,aic=F)$ar))
    }
    monte_carlo_sb_se<-
c(monte_carlo_sb_se,sd(sb_phis_star)*(((ceiling(n_sb/blocklen_sb)*blocklen_sb)/n_sb)**(1/
2))) #προσαρμογη σε σε πιθανο διαφ.μεγεθος
    monte_carlo_sb_mean<-c(monte_carlo_sb_mean,mean(sb_phis_star))
    bias<-mean(sb_phis_star)-unname(arima(y, order = c(1, 0, 0), include.mean = F)$coef)
    #bias<-mean(sb_phis_star)-unname(ar(y,method="mle",order.max = 1,demean =F,intercept
= F,aic=F)$ar)
    monte_carlo_sb_bias<-c(monte_carlo_sb_bias,bias)
    #percentile_ci<-unname(c(quantile(sb_phis_star,0.025),quantile(sb_phis_star,0.975)))
    percentile_ci<-unname(c(quantile(sb_phis_star,0.05),quantile(sb_phis_star,0.95)))
    if (phi_sb>percentile_ci[1]&&phi_sb<percentile_ci[2]){
      sum=sum+1
    }
  }
}

```

```

}
}
prob<-sum/M_sb

return(list(mean(monte_carlo_sb_se),mean(monte_carlo_sb_mean),mean(monte_carlo_sb_bi
as),prob))
}
#mc_sb(500,100,0.5,200,8) #mc=500,n=100,φ=0.5,B=200,lamda=8
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_sb(500,100,f,200,25)
  print(result)
  cat("#####\n")
}

#residual bootstrap
#M επαναληψεις, n μεγεθος χρονοσειρας, phi, B, μεσος η οχι (0 οχι)
mc_res<-function(M_res,n_res,phi_res,B_res,intercept_res){
  set.seed(10)
  monte_carlo_res_se<-c() #1
  monte_carlo_res_mean<-c() #2
  monte_carlo_res_bias<-c() #3
  sum<-0
  for (m in 1:M_res){
    y<-c(errors[m,1],rep(0,n_res-1))
    for (i in 2:n_res){
      y[i]<-y[i-1]*phi_res+errors[m,i]
    }
    #ταξη, μεσος ή όχι, δεδομενα, B επαναληψεις
    resboot_phis_star<-resboot_fun(1,intercept_res,y,B_res)
    monte_carlo_res_se<-c(monte_carlo_res_se,sd(resboot_phis_star)) #προσαρμογη σε σε
πιθανο διαφ.μεγεθος
    monte_carlo_res_mean<-c(monte_carlo_res_mean,mean(resboot_phis_star))
    bias<-mean(resboot_phis_star)-unnname(ar.ols(y,order.max=1,demean =F,intercept =
F,aic=F)$ar)#ols το αλλαζουμε αν εχει μεσο η οχι..
    monte_carlo_res_bias<-c(monte_carlo_res_bias,bias)
    percentile_ci<-
unnname(c(quantile(resboot_phis_star,0.05),quantile(resboot_phis_star,0.95)))
    if (phi_res>percentile_ci[1]&&phi_res<percentile_ci[2]){
      sum=sum+1
    }
  }
}
prob<-sum/M_res

return(list(mean(monte_carlo_res_se),mean(monte_carlo_res_mean),mean(monte_carlo_res_
bias),prob))
}

```

```

#mc_res(500,100,0.5,200,0) #φ=0.5
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_res(500,200,f,200,0)
  print(result)
  cat("#####\n")
}

```

#mle

```

#M επαναληψεις, n μεγεθος χρονοσειρας, phi πραγματικο
mc_mle<-function(M_mle,n_mle,phi_mle){
  monte_carlo_mle_se<-c() #1
  monte_carlo_mle_coef<-c() #2
  monte_carlo_mle_bias<-c() #3
  for (m in 1:M_mle){
    y<-c(errors[m,1],rep(0,n_mle-1))
    for (i in 2:n_mle){
      y[i]<-y[i-1]*phi_mle+errors[m,i]
    }
    monte_carlo_mle_se<-c(monte_carlo_mle_se,sqrt(arima(y, order = c(1, 0, 0), include.mean
= F)$var.coef)) #se mle εκτιμητη
    mle_coef<-arima(y, order = c(1, 0, 0), include.mean = F)$coef #coef mle εκτιμητη
    monte_carlo_mle_coef<-c(monte_carlo_mle_coef,mle_coef) #coef mle εκτιμητη
    mle_bias<-mle_coef-phi_mle
    monte_carlo_mle_bias<-c(monte_carlo_mle_bias,mle_bias)
  }
}

```

```

return(list(mean(monte_carlo_mle_se),mean(monte_carlo_mle_coef),mean(monte_carlo_mle
_bias)))
}
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_mle(500,100,f)
  print(result)
  cat("#####\n")
}

```

#ols

```

#M επαναληψεις, n μεγεθος χρονοσειρας, phi πραγματικο
mc_ols<-function(M_ols,n_ols,phi_ols){
  monte_carlo_ols_se<-c() #1
  monte_carlo_ols_coef<-c() #2
  monte_carlo_ols_bias<-c() #3
  for (m in 1:M_ols){
    y<-c(errors[m,1],rep(0,n_ols-1))
  }
}

```



```

for (i in 2:n_ols){
  y[i]<-y[i-1]*phi_ols+errors[m,i]
}
monte_carlo_ols_se<-c(monte_carlo_ols_se,ar.ols(y,order.max=1,demean =F,intercept =
F,aic=F)$asy.se.coef[2]$ar) #se ols εκτιμητη
ols_coef<-ar.ols(y,order.max=1,demean =F,intercept = F,aic=F)$ar[1] #coef ols εκτιμητη
monte_carlo_ols_coef<-c(monte_carlo_ols_coef,ols_coef) #coef ols εκτιμητη
ols_bias<-ols_coef-phi_ols
monte_carlo_ols_bias<-c(monte_carlo_ols_bias,ols_bias)
}

return(list(mean(monte_carlo_ols_se),mean(monte_carlo_ols_coef),mean(monte_carlo_ols_bi
as)))
}
fs<-c(-0.75,-0.5,-0.25,-0.1,0.1,0.25,0.5,0.75)
for (f in fs){
  cat("Για φ =",f,"\n")
  result<-mc_ols(500,100,f)
  print(result)
  cat("#####\n")
}

```

Π2 Οπτικοποίηση

#visualization 10 χρονοσειρών bootstrap

```

set.seed(1)
x<- arima.sim(model=list(ar=0.5),n=100) #simulation διαδικασιας/χρονοσειρας ar(1) με 0.5 φ
set.seed(1)
x1<- arima.sim(model=list(ar=-0.5),n=100) #simulation διαδικασιας/χρονοσειρας ar(1) με -
0.5 φ
par(mfcol=c(1,1))

```

#nbb_visualization

```

vis1 fun<-function(data,lamda,Breps){
  B<-Breps
  for (i in 1:Breps){
    nbb_indexes<-blockfun(data,lamda)
    lines(data[nbb_indexes],col=i+1,lwd=1) #μια χρονοσειρα nbb_bootstrap προστίθεται σε
καθε επανάληψη
  }
  lines(data,lwd=2)
}
#φ=0.5
plot(x,ylim=c(-6,6),ylab="",lwd=2, main="Non-Overlapping Block Bootstrap\n φ=0.5")
vis1 fun(x,4,10)
#φ=-0.5
plot(x1,ylim=c(-6,6),ylab="",lwd=2, main="Non-Overlapping Block Bootstrap\n φ=-0.5")
vis1 fun(x1,4,10)

```

#mbb_visualization

```
vis2fun<-function(data,lamda,Breps){
  B<-Breps
  for (i in 1:Breps){
    mbb_indexes<-mbb_boot_index(data,lamda)
    lines(data[mbb_indexes],col=i+1,lwd=1) #μια χρονοσειρα mbb_bootstrap προστίθεται σε
καθε επανάληψη
  }
  lines(data,lwd=2)
}
plot(x,ylim=c(-6,6),ylab="",lwd=2, main="Moving Block Bootstrap\n φ=0.5")
vis2fun(x,4,10)
#
plot(x1,ylim=c(-6,6),ylab="",lwd=2, main="Moving Block Bootstrap\n φ=-0.5")
vis2fun(x1,4,10)
```

#cbb_visualization

```
vis3fun<-function(data,lamda,Breps){
  B<-Breps
  for (i in 1:Breps){
    cbb_indexes<-cbb_boot_index(data,lamda)
    lines(data[cbb_indexes],col=i+1,lwd=1) #μια χρονοσειρα cbb_bootstrap προστίθεται σε
καθε επανάληψη
  }
  lines(data,lwd=2)
}
plot(x,ylim=c(-6,6),ylab="",lwd=2, main="Circular Block Bootstrap\n φ=0.5")
vis3fun(x,4,10)
#
plot(x1,ylim=c(-6,6),ylab="",lwd=2, main="Circular Block Bootstrap\n φ=-0.5")
vis3fun(x1,4,10)
```

#sb_visualization

```
vis4fun<-function(data,lamda,Breps){
  B<-Breps
  for (i in 1:Breps){
    sbb_indexes<-sbbfun(data,lamda)
    lines(data[sbb_indexes],col=i+1,lwd=1) #μια χρονοσειρα sbb_bootstrap προστίθεται σε
καθε επανάληψη
  }
  lines(data,lwd=2)
}
plot(x,ylim=c(-6,6),ylab="",lwd=2, main="Stationary Bootstrap\n φ=0.5")
vis4fun(x,4,10)
#
plot(x1,ylim=c(-6,6),ylab="",lwd=2, main="Stationary Bootstrap\n φ=-0.5")
vis4fun(x1,4,10)
```

#residual fun bootstrap για visualization

```
#συναρτηση resboot_visual (order_μοντελου,intercept=0 (δεν εχει σταθερο),δεδομενα,
επαναληψεις bootstrap)
resboot_visual<-function(ord,inter,timedata,Brep){
  data<-timedata #
  par(mfcol=c(1,1))
  plot(data,ylim=c(-6,6),ylab="",lwd=2) #για να δωσει ευρος στο plot
  if (inter==0){ #αν δεν εχουμε σταθερο ορο προσαρμοζουμε μοντελο με ols χωρις σταθερο
ορο
  arfit1<-ar.ols(data,order.max = ord,demean =F,intercept = F,aic=F)#ols
  }else{
  arfit1<-ar.ols(data,order.max = ord,demean =F,intercept = T,aic=F)
  } #συνεπως στο arfit1 ειναι το προσαρμοσμένο μοντέλο ταξεως που επιθυμώ
#bootstrap στα residuals
B<-Brep;
bootstrap_samples<-matrix(0,nrow=B,ncol=length(data))#πινακα Bxlength(data) (καθε
γραμμη θα ειναι μια boot χρονοσειρα)
bootstrap_samples
for (b in 1:B){
  bootresid<-sample(arfit1$resid[!is.na(arfit1$resid)],replace=T) #bootstrap δειγμα
καταλοιπων (εχω αφαιρεσει NA τιμες)
  bootresid<-bootresid-mean(bootresid) #κεντραρισμα καταλοιπων
  bootstrap_data<-numeric(length(data))#φτιαχνουμε ενα διανυσμα μηδενικων το οποιο οταν
το γεμισουμε θα ειναι η νεα χρονοσειρα bootstrap
  bootstrap_data[1:(arfit1$order)]<-data[1:(arfit1$order)]# p πρωτες τιμες της bootstrap
χρονοσειρας θα ειναι ιδια με της αρχικης χρονοσειρας
  bootresid<-c(rep(0,arfit1$order),bootresid) #το εβαλα ετσι για να ταιριαζει με την θεωρια
και να ξεκινησει απο  $e^*(t)$  η αναδρομη (οι προηγουμε τιμες  $e^*(t-1)$ ,.. θα ειναι 0)
  #δεν θα χρησιμοποιησω τα 0 που εβαλα (οσα και αν ειναι)
  for (t in (arfit1$order+1):length(data)){ #ξεκιναει η κατασκευη της (b) χρονοσειρας
bootstrap
  bootstrap_y<-0 #η τιμη αυτη μετα απο αναδρομη θα δωσει το  $y^*(t)$  της εκαστοτε
χρονικης στιγμης
  for (i in 1:arfit1$order){ #υπολογιζουμε τιμη χρονικης στιγμης t με βαση την
προηγουμενη/προηγουμενες αν τρεξει το loop πολλες φορες (εξαρταται την ταξη order)
  bootstrap_y<-bootstrap_y+arfit1$ar[i]*bootstrap_data[t-i] #
  }
  if (inter==0){
  bootstrap_data[t]<-bootstrap_y+bootresid[t] #βαζουμε στη χρονοσειρα bootstrap στη
θεση t την τιμη  $y^*(t)$ 
  }else{
  bootstrap_data[t]<-arfit1$x.intercept+bootstrap_y+bootresid[t] #με σταθερο ορο
  }
  }
  bootstrap_samples[b,]<-bootstrap_data #βαζουμε στον πινακα στην θεση b (γραμμη) την
bootstrap χρονοσειρα
```

```

lines(ts(bootstrap_data,start=start(data)[1],end=end(data)[1],frequency =
1),col=b+1,lwd=1) #βάζω σε ένα σχήμα τις χρονοσειρές bootstrap (μια μια σε κάθε
επανάληψη)
  #το κάνουμε ts παραπάνω για να μπορεί να ταιριάζει στο σχήμα σε περίπτωση που έχω μια
αλλη χρονοσειρά που έχει αρχή πχ το 1980 και τέλος το 2023
}
}
#φ=0.5
resboot_visual(1,0,x,10)
lines(x,lwd=2)
title("Residual Bootstrap\n φ=0.5")
#φ=-0.5
resboot_visual(1,0,x1,10)
lines(x1,lwd=2)
title("Residual Bootstrap\n φ=-0.5")

```

#ακολουθίες τυπικών σφαλμάτων και ιστογράμματα

```
#nbb
```

```
B<-500
```

```
nbbphis_star<-c()
```

```
for (b in 1:B){
```

```
  newindex<-blockfun(x,4)
```

```
  nbbphis_star<-c(nbbphis_star,unname(ar(x[newindex],method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar))
```

```
  #nbbseries<-x1[blockfun(x1,4)]
```

```
  #nbbphis_star<-c(nbbphis_star,unname(ar(nbbseries,method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar))
```

```
}
```

```
hist(nbbphis_star,xlab="φ*",main="Non-Overlapping Block Bootstrap\n Ιστόγραμμα των φ*
")
```

```
#ακολουθία τυπικών σφαλμάτων για να δούμε αν σταθεροποιείται το τυπικό σφάλμα με B
επανάληψεις bootstrap
```

```
nbb_seqphis<-c(0);nbb_seqphis
```

```
for (i in 2:length(nbbphis_star)){
```

```
  nbb_seqphis<-c(nbb_seqphis,sd(nbbphis_star[1:i]))
```

```
}
```

```
plot(nbb_seqphis,type="l",ylim=c(0,0.14),ylab="Τυπικό σφάλμα",xlab="B",main="Non-
Overlapping Block Bootstrap")
```

```
abline(h=sd(nbbphis_star),lwd=1,col="red")
```

```
abline(v=200,col="green",lty=3)
```

```
abline(v=B,col="green",lty=3)
```

```
#mbb
```

```
mbb_phis<-c();B<-500
```

```
for (b in 1:B){
```

```
  mbbseries<-x[mbb_boot_index(x,4)] #mbb_χρονοσειρά bootstrap (μεγεθος block 4)
```

```
  mbb_phis<-c(mbb_phis,unname(ar(mbbseries,order.max = 1,method="mle",demean
=F,intercept = F,aic=F)$ar))
```

```

#mbbseries<-x1[mbb_boot_index(x1,4)] #mbb_χρονοσειρα bootstrap (μεγεθος block 4)
#mbb_phis<-c(mbb_phis,unname(ar(mbbseries,order.max = 1,method="mle",demean
=F,intercept = F,aic=F)$ar))
}
mbb_phis #τα φ* (εχουμε B γιατι εχουμε B χρονοσειρες bootstrap)
hist(mbb_phis,xlab="φ*",main="Moving Block Bootstrap\n Ιστόγραμμα των φ* ")
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαγμα με B
επαναληψεις bootstrap)
mbb_seqphis<-c(0);mbb_seqphis
for (i in 2:length(mbb_phis)){
  mbb_seqphis<-c(mbb_seqphis,sd(mbb_phis[1:i]))
}
plot(mbb_seqphis,type="l",ylim=c(0,0.1),ylab="Τυπικό σφάλμα",xlab="B",main="Moving
Block Bootstrap")
abline(h=sd(mbb_phis),lwd=1,col="red")
abline(v=200,col="green",lty=3)
abline(v=B,col="green",lty=3)

```

#cbb

```

cbb_phis<-c();B<-500
for (i in 1:B){
  cbbseries<-x[cbb_boot_index(x,4)] #cbb_χρονοσειρα bootstrap (blocklength=4)
  cbb_phis<-c(cbb_phis,unname(ar(cbbseries,method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar))
  #cbbseries<-x1[cbb_boot_index(x1,4)] #cbb_χρονοσειρα bootstrap (blocklength=4)
  #cbb_phis<-c(cbb_phis,unname(ar(cbbseries,method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar))
}
cbb_phis
hist(cbb_phis,xlab="φ*",main="Circular Block Bootstrap\n Ιστόγραμμα των φ* ")
sd(cbb_phis) #se
mean(cbb_phis) #mean of φ*(b), b=1,2...B

```

```

#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαγμα με B
επαναληψεις bootstrap)
cbb_seqphis<-c(0);cbb_seqphis
for (i in 2:length(cbb_phis)){
  cbb_seqphis<-c(cbb_seqphis,sd(cbb_phis[1:i]))
}
plot(cbb_seqphis,type="l",ylim=c(0,0.1),ylab="Τυπικό σφάλμα",xlab="B",main="Circular
Block Bootstrap")
abline(h=sd(cbb_phis),lwd=1,col="red")
abline(v=200,col="green",lty=3)
abline(v=B,col="green",lty=3)

```

#sb

```

sbb_phis<-c();B<-500
for (b in 1:B){

```

```

sbbseries<-x[sbbfun(x,4)]
#sbbseries<-x1[sbbfun(x1,4)]
sbb_phis<-c(sbb_phis,unname(ar(sbbseries,method="mle",order.max = 1,demean
=F,intercept = F,aic=F)$ar))
}
sbb_phis
sd(sbb_phis)
mean(sbb_phis)
hist(sbb_phis,xlab="φ*",main="Stationary Bootstrap\n Ιστόγραμμα των φ* ")
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαλαμα με B
επαναληψεις bootstrap)
sbb_seqphis<-c(0);sbb_seqphis
for (i in 2:length(sbb_phis)){
  sbb_seqphis<-c(sbb_seqphis,sd(sbb_phis[1:i]))
}
plot(sbb_seqphis,ylim=c(0,0.1),type="l",ylab="Τυπικό σφάλμα",xlab="B",main="Stationary
Bootstrap")
abline(h=sd(sbb_phis),lwd=1,col="red")
abline(v=200,col="green",lty=3)
abline(v=B,col="green",lty=3)

```

#residual bootstrap

```

#συναρτηση resboot(order_μοντελου,intercept=0 (δεν εχει σταθερο),δεδομενα, επαναληψεις
bootstrap)
resboot<-function(ord,inter,timedata,Brep){
  set.seed(1)
  data<-timedata #
  par(mfcol=c(1,ord+1))
  if (min(data)<0){
    plot(data,ylim=c(min(data)+0.5*min(data),max(data)+0.5*max(data)),ylab="",lwd=2,lty=6)
    #για να δωσει ευρος στο plot
  }else{
    plot(data,ylim=c(min(data)-0.5*min(data),max(data)+0.5*max(data)),ylab="",lwd=2,lty=6)
    #για να δωσει ευρος στο plot
  }
  if (inter==0){ #αν δεν εχω σταθερο ορο προσαρμοζουμε μοντελο με ols χωρις σταθερο ορο
    arfit1<-ar.ols(data,order.max = ord,demean =F,intercept = F,aic=F)#ols
  }else{
    arfit1<-ar.ols(data,order.max = ord,demean =F,intercept = T,aic=F)
  } #συνεπως στο arfit1 ειναι το προσαρμοσμένο μοντέλο ταξεως που επιθυμούμε (απο arfit1
θα παρουμε residuals)
#bootstrap στα residuals
B<-Brep;
bootstrap_samples<-matrix(0,nrow=B,ncol=length(data))#πινακα Bxlength(data) (καθε
γραμμη θα ειναι μια boot χρονοσειρα)
bootstrap_samples
for (b in 1:B){

```

```

bootresid<-sample(arfit1$resid[!is.na(arfit1$resid)],replace=T) #bootstrap δειγμα
καταλοιπων (εχουμε αφαιρεσει NA τιμες)
bootresid<-bootresid-mean(bootresid) #κεντραρισμα
bootstrap_data<-numeric(length(data))#φτιαχνουμε ενα διανυσμα μηδενικων το οποιο οταν
το γεμισουμε θα ειναι η νεα χρονοσειρα bootstrap
bootstrap_data[1:(arfit1$order)]<-data[1:(arfit1$order)]# p πρωτες τιμες της bootstrap
χρονοσειρας θα ειναι ιδια με της αρχικης χρονοσειρας
bootresid<-c(rep(0,arfit1$order),bootresid) #το εβαλα ετσι για να ταιριαζει με την θεωρια
και να ξεκινήσει απο  $e^*(t)$  η αναδρομη (οι προηγουμε τιμες  $e^*(t-1)$ ,.. θα ειναι 0)
#δεν θα χρησιμοποιησω τα 0 που εβαλα (οσα και αν ειναι)
for (t in (arfit1$order+1):length(data)){ #ξεκινει η κατασκευη της (b) χρονοσειρας
bootstrap
bootstrap_y<-0 #η τιμη αυτη μετα απο αναδρομη θα δωσει το  $y^*(t)$  της εκαστοτε
χρονικης στιγμης
for (i in 1:arfit1$order){ #υπολογιζω τιμη χρονικης στιγμης t με βαση την
προηγουμενη/προηγουμενες αν τρεξει το loop πολλες φορες (εξαρταται την ταξη order)
bootstrap_y<-bootstrap_y+arfit1$ar[i]*bootstrap_data[t-i] #
}
if (inter==0){
bootstrap_data[t]<-bootstrap_y+bootresid[t] #βαζουμε στη χρονοσειρα bootstrap στη
θεση t την τιμη  $y^*(t)$ 
}else{
bootstrap_data[t]<-arfit1$x.intercept+bootstrap_y+bootresid[t] #με σταθερο ορο
}
}
bootstrap_samples[b,]<-bootstrap_data #βαζουμε στον πινακα, στην θεση b (γραμμη) την
bootstrap χρονοσειρα
lines(ts(bootstrap_data,start=start(data)[1],end=end(data)[1],frequency = 1),col=b,lwd=1)
#βαζω σε ενα σχημα τις χρονοσειρες bootstrap (μια μια σε καθε επαναληψη)
#το κανω ts παραπανω για να μπορει να ταιριαξει στο σχημα σε περιπτωση που εχω μια
αλλη χρονοσειρα που εχει αρχη πχ το 1980 και τελος το 2023
}
bootstrap_phis<-matrix(0,nrow=B,ncol=arfit1$order) #πινακας οπου καθε γραμμη θα ειναι
το/τα φ που προκυπτει απο καθε χρονοσειρα bootstrap
for (i in 1:nrow(bootstrap_samples)){
if (inter==0){
bootstrap_phis[i,]<-unname(ar(bootstrap_samples[i,],method = "ols",order.max =
arfit1$order,demean =F,intercept = F,aic=F)$ar)
}else{
bootstrap_phis[i,]<-unname(ar(bootstrap_samples[i,],method = "ols",order.max =
arfit1$order,demean =F,intercept =T,aic=F)$ar)
}
}
}
for (j in 1:ord){
cat("----το μοντέλο ήταν τάξης: ",ord,"ης----\n")
cat("sd_bootstrap_phi_",j,"=",sd(bootstrap_phis[,j])," mean_bootstrap_phi_",j,"=
",mean(bootstrap_phis[,j]))
cat("\n#####")
}

```

```

    hist(bootstrap_phis[,j],xlab=paste("Histogram of boot_phi", j),main="")
  }
  return(bootstrap_phis)
}
res_phis<-resboot(1,0,x,500)
#res_phis<-resboot(1,0,x1,500) #αν θελω να βαλω x1 δηλαδη χρονοσειρα με φ=-0.5
mean(res_phis)
sd(res_phis)
par(mfcol=c(1,1))
hist(res_phis,xlab="φ*",main="Residual Bootstrap\n Ιστόγραμμα των φ* ")
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαλαμα με B
επαναληψεις bootstrap)
res_seqphis<-c(0);res_seqphis
for (i in 2:length(res_phis)){
  res_seqphis<-c(res_seqphis,sd(res_phis[1:i]))
}
plot(res_seqphis,ylim=c(0,0.11),type="l",ylab="Τυπικό σφάλμα",xlab="B",main="Residual
Bootstrap")
abline(h=sd(res_phis),lwd=1,col="red")
abline(v=200,col="green",lty=3)
abline(v=500,col="green",lty=3)

```

#κλασικό iid bootstrap

```

efron_boot_phis<-c()
for(b in 1:B){
  m<-sample(1:length(x),length(x),replace=T)
  efron_boot_phis<-c(efron_boot_phis,unname(ar(x[m],method = "mle",order.max =
1,demean =F,intercept = F,aic=F)$ar))
  #m<-sample(1:length(x1),length(x1),replace=T)
  #efron_boot_phis<-c(efron_boot_phis,unname(ar(x1[m],method = "mle",order.max =
1,demean =F,intercept = F,aic=F)$ar))
}
mean(efron_boot_phis)
sd(efron_boot_phis) #se
hist(efron_boot_phis,xlab="φ*",main="Κλασικό Bootstrap\n Ιστόγραμμα των φ* ")
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαλαμα με B
επαναληψεις bootstrap)
efronboot_seqphis<-c(0);efronboot_seqphis
for (i in 2:length(efron_boot_phis)){
  efronboot_seqphis<-c(efronboot_seqphis,sd(efron_boot_phis[1:i]))
}
plot(efronboot_seqphis,type="l",ylim=c(0,.12),ylab="Τυπικό
σφάλμα",xlab="B",main="Κλασικό Bootstrap")
abline(h=sd(efron_boot_phis),lwd=1,col="red")
abline(v=200,col="green",lty=3)
abline(v=B,col="green",lty=3)

```


Π3 Εφαρμογή σε πραγματικά δεδομένα

```
attach(Sunspot)
Sunspot
#1770-1889
sun<-sunspots[71:190]
mean(sun)
length(sun)
par(mfcol=c(1,1))
plot(year[71:190],sun,type = "l",xlab="Year",ylab="Sunspots")
sun<-sun-mean(sun);sun
mean(sun)
adf_test_result <- adf.test(sun)
print(adf_test_result)
#acf
acf(sun, lag.max = 30, main="")
#pacf
pacf(sun, lag.max = 30, main="")
#
#ols
ar.ols(sun,order.max=1,demean =F,intercept = F,aic=F)$ar[1]#ols
unname(ar.ols(sun,order.max=1,demean =F,intercept = F,aic=F)$asy.se.coef[2])
#mle
arima(sun, order = c(1, 0, 0), include.mean = F)

#nbb
B<-500
nbbphis_star1<-c()
for (b in 1:B){
  nbbindexes<-blockfun(sun,25) #nbb indexes χρονοσειρας bootstrap
  nbbphis_star1<-c(nbbphis_star1,unname(arima(sun[nbbindexes], order = c(1, 0, 0),
include.mean = F)$coef))
}
nbbphis_star1
sd(nbbphis_star1)
mean(nbbphis_star1)
mean(nbbphis_star1)-unname(arima(sun, order = c(1, 0, 0), include.mean = F)$coef)
#ιστόγραμμα των φ*
hist(nbbphis_star1,main="Non-Overlapping Block Bootstrap\n Ιστόγραμμα των
φ*",breaks=15,xlab="φ*")
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαλμα με B
επαναληψεις bootstrap)
nbb_seqphis<-c(0);nbb_seqphis
for (i in 2:length(nbbphis_star1)){
  nbb_seqphis<-c(nbb_seqphis,sd(nbbphis_star1[1:i]))
}
par(mfcol=c(1,1))
```

```
plot(nbb_seqphis,type="l",col="orange",ylim=c(0,0.06),ylab="Τυπικό
σφάλμα",xlab="B",main="Ακολουθίες SE")
legend("bottomright", legend = c("NBB", "MBB", "CBB", "SB"), col =
c("orange", "green", "purple", "red"), lty = 1, lwd = 2, xpd = TRUE, inset = c(0, 0), cex=0.8)
```

#mbb

```
B<-500
mbbphis_star1<-c()
for (b in 1:B){
  mbbindexes<-mbb_boot_index(sun,25) #δεικτες mbb_χρονοσειρας bootstrap
  mbbphis_star1<-c(mbbphis_star1,unname(arima(sun[mbbindexes], order = c(1, 0, 0),
include.mean = F)$coef))
}
mbbphis_star1 #τα φ* (εχω B γιατι εχω B χρονοσειρες bootstrap)
sd(mbbphis_star1)
mean(mbbphis_star1)
mean(mbbphis_star1)-unname(arima(sun, order = c(1, 0, 0), include.mean = F)$coef)
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαλμα με B
επαναληψεις bootstrap)
mbb_seqphis<-c(0);mbb_seqphis
for (i in 2:length(mbbphis_star1)){
  mbb_seqphis<-c(mbb_seqphis,sd(mbbphis_star1[1:i]))
}
lines(mbb_seqphis,type="l",col="green")
#ιστόγραμμα των φ*
hist(mbbphis_star1)
hist(mbbphis_star1,main="Moving Block Bootstrap\n Ιστόγραμμα των
φ*",breaks=15,xlab="φ*")
```

#cbb

```
B<-500
cbbphis_star1<-c()
for (b in 1:B){
  cbbindexes<-cbb_boot_index(sun,25) #indexes cbb_χρονοσειρας bootstrap
  cbbphis_star1<-c(cbbphis_star1,unname(arima(sun[cbbindexes], order = c(1, 0, 0),
include.mean = F)$coef))
}
cbbphis_star1
sd(cbbphis_star1)
mean(cbbphis_star1)
mean(cbbphis_star1)-unname(arima(sun, order = c(1, 0, 0), include.mean = F)$coef)
#ακολουθία τυπικών σφαλμάτων (για να δουμε αν σταθεροποιείται το τυπικο σφαλμα με B
επαναληψεις bootstrap)
cbb_seqphis<-c(0);cbb_seqphis
for (i in 2:length(cbbphis_star1)){
  cbb_seqphis<-c(cbb_seqphis,sd(cbbphis_star1[1:i]))
}
lines(cbb_seqphis,type="l",col="purple")
```

```

#ιστογραμμα των φ*
hist(cbbphis_star1)
hist(cbbphis_star1,main="Circular Block Bootstrap\n Ιστόγραμμα των
φ*",breaks=15,xlab="φ*")

#sb
sbbphis_star1<-c()
for (b in 1:B){
  sbbindexes<-sbbfun(sun,25) #sbb indexes sbb χρονοσειρας bootstrap
  sbbphis_star1<-c(sbbphis_star1,unname(arima(sun[sbbindexes], order = c(1, 0, 0),
include.mean = F)$coef))
}
sbbphis_star1
sd(sbbphis_star1)
mean(sbbphis_star1)
mean(sbbphis_star1)-unname(arima(sun, order = c(1, 0, 0), include.mean = F)$coef)
#ιστογραμμα των φ*
hist(sbbphis_star1)
hist(sbbphis_star1,main="Stationary Bootstrap\n Ιστόγραμμα των φ*",breaks=15,xlab="φ*")
#ακολουθία τυπικών σφαλμάτων (για να δω αν σταθεροποιείται το τυπικο σφαλμα με B
επαναληψεις bootstrap)
sbb_seqphis<-c(0);sbb_seqphis
for (i in 2:length(sbbphis_star1)){
  sbb_seqphis<-c(sbb_seqphis,sd(sbbphis_star1[1:i]))
}
lines(sbb_seqphis,type="l",col="red")

#residual bootstrap
#συναρτηση resboot(order_μοντελου,intercept=0 (δεν εχει σταθερο),δεδομενα, επαναληψεις
bootstrap)
res_phis<-resboot(1,0,sun,500)
res_phis
sd(res_phis)
mean(res_phis)-ar.ols(sun,order.max=1,demean =F,intercept = F,aic=F)$ar[1]
mean(res_phis)
#ακολουθία τυπικών σφαλμάτων (για να δούμε αν σταθεροποιείται το τυπικο σφαλμα με B
επαναληψεις bootstrap)
res_seqphis<-c(0);res_seqphis
for (i in 2:length(res_phis)){
  res_seqphis<-c(res_seqphis,sd(res_phis[1:i]))
}
par(mfcol=c(1,1))
plot(res_seqphis,type="l",ylim=c(0,0.09),ylab="Τυπικό σφάλμα",xlab="B",main="Ακολουθία
SE")
legend("bottomright", legend ="Residual Bootstrap", col = "black", lty = 1, lwd = 2, xpd =
TRUE,inset = c(0, 0),cex=0.8)
hist(res_phis,main="Residual Bootstrap Ιστόγραμμα των φ*",breaks=15,xlab="φ*")

```

```

hist(res_phis,main="Residual Bootstrap\n Ιστόγραμμα των φ*",breaks=15,xlab="φ*")

#για δεύτερης τάξης παρατίθεται τμήμα κώδικα
B<-500
phis_star11<-c();phis_star12<-c()
for (b in 1:B){
  indexes<-blockfun(sun,25) #nbb indexes χρονοσειρας bootstrap
  #indexes<-mbb_boot_index(sun,25) #mbb
  #indexes<-cbb_boot_index(sun,25) #cbb
  #indexes<-sbbfun(sun,25) #sb
  phis_star11<-c(phis_star11,unname(arima(sun[indexes], order = c(2, 0, 0), include.mean =
F)$coef[1]))
  phis_star12<-c(phis_star12,unname(arima(sun[indexes], order = c(2, 0, 0), include.mean =
F)$coef[2]))
}
#se φ1
sd(phis_star11)
#mean φ1
mean(phis_star11)
#bias φ1
mean(phis_star11)-unname(arima(sun, order = c(2, 0, 0), include.mean = F)$coef)[1]
#hist φ1
hist(phis_star11,main="Non-Overlapping Block Bootstrap\n Ιστόγραμμα των
φ1*",breaks=15,xlab="φ1*")
#se φ2
sd(phis_star12)
#mean φ2
mean(phis_star12)
#bias φ2
mean(phis_star12)-(arima(sun, order = c(2, 0, 0), include.mean = F)$coef)[2]
#
hist(phis_star12,main="Non-Overlapping Block Bootstrap\n Ιστόγραμμα των
φ2*",breaks=15,xlab="φ2*")

res_phis<-resboot(2,0,sun,500)
#coeff1
mean(res_phis[,1])
mean(res_phis[,1])-ar.ols(sun,order.max=2,demean =F,intercept = F,aic=F)$ar[1]
sd(res_phis[,1])
#coeff2
mean(res_phis[,2])
sd(res_phis[,2])
mean(res_phis[,2])-ar.ols(sun,order.max=2,demean =F,intercept = F,aic=F)$ar[2]
par(mfcol=c(1,1))
hist(res_phis[,1],main="Residual Bootstrap\n Ιστόγραμμα των φ1*",breaks=18,xlab="φ1*")
hist(res_phis[,2],main="Residual Bootstrap\n Ιστόγραμμα των φ2*",breaks=15,xlab="φ2*")

```

Βιβλιογραφία

Ελληνική

Ηλιόπουλος, Γ. (2021). *Υπολογιστικές στατιστικές τεχνικές*, Πανεπιστημιακές Σημειώσεις για το ΠΜΣ Εφαρμοσμένης Στατιστικής του Πανεπιστημίου Πειραιά.

Μπούτσικας, Μ. (2021). *Μέθοδοι προσομοίωσης*, Πανεπιστημιακές Σημειώσεις για το ΠΜΣ Εφαρμοσμένης Στατιστικής του Πανεπιστημίου Πειραιά.

Ξένη

Allen, D. M. (1974). The relationship between variable selection and data augmentation and a method for prediction, *Technometrics*, 16(1), 125–127. doi:10.1080/00401706.1974.10489157

Athreya, K. B. (1987). Bootstrap of the mean in the infinite variance case, *Annals of Statistics*, 15(2), 724-731. doi:10.1214/aos/1176350371

Babu, G. J. (1984). Bootstrapping Statistics with Linear Combinations of Chi-Squares as Weak Limit, *Sankhyā: The Indian Journal of Statistics, Series A (1961-2002)*, 46(1), 85–93.

Bickel, P. J., Götze, F., & van Zwet, W. R. (1997). Resampling fewer than n observations, gains, losses, and remedies for losses, *Statistica Sinica*, 7(1), 1–31.

Bose, A. (1988). Edgeworth correction by Bootstrap in autoregressions, *Annals of Statistics*, 16(4). doi:10.1214/aos/1176351063

Box, G. E. P., & Jenkins, G. M. (1970). *Time Series Analysis Forecasting and Control*, Holden-Day, San Francisco. - References - Scientific Research Publishing.

Brockwell, P. J., & Davis, R. A. (2016). *Introduction to Time Series and Forecasting*, 3rd ed., Springer Texts in Statistics.

Bühlmann, P. (1997). Sieve Bootstrap for time series, *Bernoulli*, 3(2), 123-148. doi:10.2307/3318584

Bühlmann, P. (2002). Bootstraps for time series, *Statistical Science*, 17(1), 52-72. doi:10.1214/ss/1023798998

- Carlstein, E. (1986). The Use of Subseries Values for Estimating the Variance of a General Statistic from a Stationary Sequence, *Annals of Statistics*, 14(3), 1171-1179. doi:10.1214/aos/1176350057
- Carlstein, E., Do, K., Hall, P., Hesterberg, T., Künsch, H. R., & Künsch, H. R. (1998). Matched-Block Bootstrap for dependent data, *Bernoulli*, 4(3), 305-328. doi:10.2307/3318719
- Chatfield, C., & Xing, H. (2019). *The Analysis of Time Series: An Introduction with R*, Chapman and Hall/CRC.
- Chernick, M. R. (2007). *Bootstrap Methods: A guide for practitioners and researchers*, John Wiley & Sons. doi:10.1002/9780470192573
- Cowpertwait, P. S., & Metcalfe, A. (2009). *Introductory Time Series with R*, In Springer eBooks. doi:10.1007/978-0-387-88698-5
- Datta, S. (1996). On asymptotic properties of bootstrap for AR(1) processes, *Journal of Statistical Planning and Inference*, 53(3), 361–374. doi:10.1016/0378-3758(95)00147-6
- Datta, S., & McCormick, W. P. (1995). Bootstrap Inference for a First-Order Autoregression with Positive Innovations, *Journal of the American Statistical Association*, 90(432), 1289–1300. doi:10.1080/01621459.1995.10476633
- Datta, S., & Sriram, T. N. (1997). A modified bootstrap for autoregression without stationarity, *Journal of Statistical Planning and Inference*, 59(1), 19–30. doi:10.1016/s0378-3758(96)00092-4
- Davison, A. C., & Hinkley, D. V. (1997). *Bootstrap Methods and their Application*, Cambridge University Press. doi:10.1017/cbo9780511802843
- Dudek, A. E. (2017). Block bootstrap for periodic characteristics of periodically correlated time series, *Journal of Nonparametric Statistics*, 30(1), 87–124. doi:10.1080/10485252.2017.1404060
- Efron, B. (1979). Bootstrap methods: Another look at the jackknife, *Annals of Statistics*, 7(1), 1-26. doi:10.1214/aos/1176344552
- Efron, B. (1982). *The Jackknife, the Bootstrap and other resampling plans*, Society for Industrial and Applied Mathematics, Philadelphia, Pennsylvania. doi:10.1137/1.9781611970319
- Efron, B. (1992). Jackknife-After-Bootstrap standard errors and influence functions, *Journal of the Royal Statistical Society: Series B-Methodological*, 54(1), 83–111. doi:10.1111/j.2517-6161.1992.tb01866.x
- Efron, B., & Hastie, T. (2016). *Computer age statistical inference*, Cambridge University Press. doi:10.1017/cbo9781316576533

- Efron, B., & Tibshirani, R. (1986). Bootstrap methods for standard errors, confidence intervals, and other measures of statistical accuracy, *Statistical Science*, 1(1), 54-75. doi:10.1214/ss/1177013815
- Efron, B., & Tibshirani, R. (1993). *An introduction to the Bootstrap*, In Chapman and Hall/CRC.
- Fisher, N. I., & Hall, P. (1990). On bootstrap hypothesis testing, *Australian Journal of Statistics*, 32(2), 177-190.
- Freedman, D. A., & Peters, S. C. (1984). Bootstrapping an econometric model: Some empirical results, *Journal of Business & Economic Statistics*, 2(2), 150-158.
- Hall, P. (1992). *The Bootstrap and Edgeworth expansion*, In Springer series in statistics. doi:10.1007/978-1-4612-4384-7
- Hall, P., Horowitz, J. L., & Jing, B. (1995). On blocking rules for the bootstrap with dependent data, *Biometrika*, 82(3), 561–574. doi:10.1093/biomet/82.3.561
- Heimann, G., & Kreiss, J. (1996). Bootstrapping general first order autoregression, *Statistics & Probability Letters*, 30(1), 87–98. doi:10.1016/0167-7152(95)00205-7
- Hyndman, R. J., & Athanasopoulos, G. (2013). *Forecasting: principles and practice*, 2nd ed., OTexts, Australia.
- Jentsch, C., & Politis, D. N. (2015). Covariance matrix estimation and linear process bootstrap for multivariate time series of possibly increasing dimension, *Annals of Statistics*, 43(3). doi:10.1214/14-aos1301
- Knight, K. (1989). On the Bootstrap of the Sample Mean in the Infinite Variance Case, *Annals of Statistics*, 17(3). doi:10.1214/aos/1176347262
- Kreiss, J.P. (1992). Bootstrap procedures for AR (∞) — processes, In: Jöckel, KH., Rothe, G., Sendler, W. (eds) *Bootstrapping and Related Techniques*, Lecture Notes in Economics and Mathematical Systems, 107–113.
- Kreiss, J. P., & Franke, J. (1992). Bootstrapping stationary autoregressive moving-average models, *Journal of Time Series Analysis*, 13(4), 297-317.
- Kreiss, J.P., & Lahiri, S. N. (2012). Bootstrap methods for time series, *In Handbook of Statistics*, 3–26. doi:10.1016/b978-0-444-53858-1.00001-6
- Kreiss, J.P. (1988). Asymptotic statistical inference for a class of stochastic processes, Habilitationsschrift, Faculty of Mathematics, Univ. Hamburg, Germany
- Künsch, H. R. (1989). The jackknife and the bootstrap for general stationary observations, *Annals of Statistics*, 17(3), 1217-1241. doi:10.1214/aos/1176347265

- Lahiri, S. N. (1999). Theoretical comparisons of block bootstrap methods, *Annals of Statistics*, 27(1), 386-404. doi:10.1214/aos/1018031117
- Lahiri, S. N. (2003). *Resampling methods for dependent data*, In Springer series in statistics. doi:10.1007/978-1-4757-3803-2
- Lahiri, S. N., Furukawa, K., & Lee, Y. (2007). A nonparametric plug-in rule for selecting optimal block lengths for block bootstrap methods, *Statistical Methodology*, 4(3), 292–321. doi:10.1016/j.stamet.2006.08.002
- Liu, R. Y., & Singh, K. (1992). Moving blocks jackknife and bootstrap capture weak dependence, *In Exploring the limits of bootstrap*, John Wiley, New York, 225-248.
- McMurry, T. L., & Politis, D. N. (2010). Banded and tapered estimates for autocovariance matrices and the linear process bootstrap, *Journal of Time Series Analysis*, 31(6), 471–482. doi:10.1111/j.1467-9892.2010.00679.x
- Montgomery, D. C., Jennings, C., & Kulahci, M. (2008). *Introduction to time series analysis and forecasting*, Wiley.
- Paparoditis, E. (1992). Bootstrapping Some Statistics Useful in Identifying ARMA Models, In: Jöckel, KH., Rothe, G., Sendler, W. (eds) *Bootstrapping and Related Techniques*, Lecture Notes in Economics and Mathematical Systems, 115-119. doi.org/10.1007/978-3-642-48850-4_15
- Paparoditis, E., & Politis, D. N. (2001). Tapered block bootstrap, *Biometrika*, 88(4), 1105–1119. doi:10.1093/biomet/88.4.1105
- Politis, D. N. (1998). Computer-intensive methods in statistical analysis, *IEEE Signal Processing Magazine*, 15(1), 39–55. doi:10.1109/79.647042
- Politis, D. N. (2003). The impact of Bootstrap methods on time series analysis, *Statistical Science*, 18(2), 219-230. doi:10.1214/ss/1063994977
- Politis, D. N. , and Romano, J. P. (1992). A circular block-resampling procedure for stationary data, *In Exploring the Limits of Bootstrap*, 263 – 270, Wiley , New York.
- Politis, D. N., & Romano, J. P. (1994). The stationary Bootstrap, *Journal of the American Statistical Association*, 89(428), 1303–1313. doi:10.1080/01621459.1994.10476870
- Politis, D. N., Romano, J. P., & Wolf, M. (1999). *Subsampling*, In Springer series in statistics. doi:10.1007/978-1-4612-1554-7
- Politis, D.N., & McElroy, T.S. (2020). *Time Series: A First Course with Bootstrap Starter*, 1st ed., Chapman and Hall/CRC. doi:10.1201/9780429109553
- Quenouille, M. H. (1949). Approximate Tests of Correlation in Time-Series, *Journal of the Royal Statistical Society: Series B-Methodological*, 11(1), 68–84.

Rao, C. R. (1989). *Statistics and Truth: Putting Chance to Work*, International Co-Operative Publishing House.

Shao, J., & Tu, D. (1995). *The Jackknife and Bootstrap*, In Springer series in statistics. doi:10.1007/978-1-4612-0795-5

Shao, J., & Wu, C. (1989). A general theory for Jackknife variance estimation, *Annals of Statistics*, 17(3). doi:10.1214/aos/1176347263

Singh, K. (1981). On the Asymptotic Accuracy of Efron's Bootstrap, *Annals of Statistics*, 9(6), 1187-1195. doi:10.1214/aos/1176345636

Stine, R. A. (1987). Estimating properties of autoregressive forecasts, *Journal of the American Statistical Association*, 82(400), 1072–1078. doi:10.1080/01621459.1987.10478542

Stone, M. (1974). Cross-Validatory choice and assessment of statistical predictions, *Journal of the Royal Statistical Society: Series B-Methodological*, 36(2), 111–133. doi:10.1111/j.2517-6161.1974.tb00994.x

Swanepoel, J., & Van Wyk, J. (1986). The bootstrap applied to power spectral density function estimation, *Biometrika*, 73(1), 135–141. doi:10.1093/biomet/73.1.135

Thombs, L. A., & Schucany, W. R. (1990). Bootstrap prediction intervals for autoregression, *Journal of the American Statistical Association*, 85(410), 486–492. doi:10.1080/01621459.1990.10476225

Trosset, M. W. (2009). *An Introduction to Statistical Inference and Its Applications with R*, In Chapman and Hall/CRC. doi:10.1201/9781584889489

Tukey, J. (1958). Bias and confidence in not quite large samples, *Annals of Mathematical Statistics*, 29(2), 614.

Wu, C. F. J. (1986). Jackknife, Bootstrap and other resampling methods in regression analysis, *Annals of Statistics*, 14(4), 1261-1295. doi:10.1214/aos/1176350142

Yule, G. (1927). On a method of investigating periodicities disturbed series, with special reference to Wolfer's sunspot numbers, *Philosophical Transactions of the Royal Society of London*, 226(636–646), 267–298. doi:10.1098/rsta.1927.0007

