



ΕΘΝΙΚΟ ΜΕΤΣΟΒΙΟ ΠΟΛΥΤΕΧΝΕΙΟ
ΣΧΟΛΗ ΕΦΑΡΜΟΣΜΕΝΩΝ ΜΑΘΗΜΑΤΙΚΩΝ ΚΑΙ
ΦΥΣΙΚΩΝ ΕΠΙΣΤΗΜΩΝ

ΥΠΟΛΟΓΙΣΤΙΚΕΣ ΜΕΘΟΔΟΙ ΣΤΗ ΣΤΑΤΙΣΤΙΚΗ

Μήτσος Αντώνης | ge19079 | anmitsis@hotmail.com

Ιούνιος 2023

Προπτυχιακός Φοιτητής

Πίνακας περιεχομένων

1	Άσκηση 1	2
2	Άσκηση 2	9
3	Άσκηση 3	16
4	Άσκηση 4	33

1 Άσκηση 1

α. Θέλουμε να εκτιμήσουμε, με την απλή μέθοδο Monte Carlo τον όγκο της μοναδιαίας σφαίρας σε d διαστάσεις με $B_d = \{(x_1, \dots, x_d) : x_1^2 + \dots + x_d^2 < 1\}$. Για να κατανοήσουμε το πρόβλημα μπορούμε να το ανάγουμε σε 2 μόνο διαστάσεις, όπου θέλουμε να υπολογίσουμε το εμβαδόν του μοναδιαίου κύκλου. Θεωρούμε επίσης το τετράγωνο με κέντρο το 0, που περιέχει όλα τα σημεία με συντεταγμένες στο $[-1,1]$. Το τετράγωνο αυτό περιέχει μέσα του τον μοναδιαίο κύκλο και έχει την ιδιότητα ότι είναι εύκολο να κάνουμε δειγματοληψία μέσα σε αυτό. Επιλέγοντας λοιπόν ένα τυχαίο σημείο μέσα στο τετράγωνο, το σημείο αυτό θα βρίσκεται μέσα στον κύκλο με πιθανότητα p και θα βρίσκεται εκτός με πιθανότητα $1-p$. Η πιθανότητα p είναι ο λόγος του εμβαδού του κύκλου προς το εμβαδό του τετραγώνου, οπότε δημιουργώντας αρκετά σημεία και μετρώντας πόσα βρίσκονται εντός του κύκλου μπορούμε να εκτιμήσουμε την πιθανότητα και άρα τον λόγο των εμβαδών των δύο σχημάτων. Τελικά η εκτίμηση του του εμβαδού του κύκλου είναι 2 φορές επί την πιθανότητα p αφού το εμβαδόν του τετραγώνου είναι 2. Επιστρέφοντας στο αρχικό πρόβλημα αρκεί να δημιουργήσουμε αρκετά σημεία το καθένα με d -συντεταγμένες στο $[-1,1]$, δηλαδή μέσα στον κύβο d -διαστάσεων και να ελέγξουμε πόσα από αυτά τα σημεία βρίσκονται εντός της μοναδιαίας σφαίρας σε d διαστάσεις.

$$\frac{n_{inside_sphere}}{n_{total}} \approx \frac{V_{sphere}}{V_{cube}} \quad (1.1)$$

Στην αριστερή πλευρά της σχέσης έχουμε το πλήθος των σημείων που βρίσκονται μέσα στη σφαίρα προς το συνολικό πλήθος σημείων και δεξιά έχουμε τους όγκους. Η σχέση τελικά γίνεται:

$$\hat{V}_{sphere} = 2 \cdot \frac{n_{inside_sphere}}{n_{total}} \quad (1.2)$$

```
1 ask1a<- function(max_d,n=10^5) {
2
3   real_vol_calc<-function(d){
4     if (d == 1) {
5       return(2)
6     } else if (d == 2) {
7       return(pi)
8     } else {
9       return((2 * pi / (d-2 + 2)) * real_vol_calc(d - 2))
10    }
11  }
12
13  for(d in 2:max_d){
14
15    real_volume=real_vol_calc(d)
16    count_inside = 0
17
18    for(i in 1:n) {
19      point = runif(d, min = -1, max = 1)
20      if(sum(point^2) <= 1) {
21        count_inside = count_inside + 1
22      }
23    }
24  }
25 }
```

```

23     }
24
25     if(count_inside==0){
26         break
27     }
28 }
29
30 estimated_volume <- 2^d * (count_inside / n)
31 error=abs((real_volume-estimated_volume)/real_volume)
32 print(paste("D: ",d))
33 print(paste("Real_Volume:", real_volume))
34 print(paste("Estimated_Volume:", estimated_volume))
35 print(paste("The_percentage_of_error_is:", round(error,5)*
    ↪ 100,"%",sep = ""))
36
37 }

```

Φτιάχνουμε μια συνάρτηση που δέχεται τη διάσταση που θέλουμε να υπολογίσουμε καθώς και το πλήθος των σημείων που θα προσομοιώσουμε. Αρχικά φτιάχνουμε μια συνάρτηση που θα υπολογίζει τον πραγματικό όγκο της σφαίρας αφού γνωρίζουμε πως $V_1 = 2, V_2 = \pi, V_{d+2} = \frac{2\pi}{d+2} V_d$. Επειδή για μεγάλες διαστάσεις ο όγκος είναι πολύ μικρός αριθμός επιλέγουμε στη συνάρτηση μας να ξεκινάει από διάσταση 2 και να ανεβαίνει διάσταση εως ότου φτάσει στην τελική διάσταση, δηλαδή αυτή που δόθηκε σαν όρισμα, ή μέχρι την στιγμή που για πρώτη φορά κανένα σημείο δεν θα βρεθεί εντός της σφαίρας. Οπότε φτιάχνουμε το for loop που θα τρέξει απο διάσταση 2 έως max_d και μέσα δημιουργούμε n σημεία d διαστάσεων από την ομοιόμορφη κατανομή στο [-1,1]. Έπειτα ελέγχουμε αν το κάθε σημείο βρίσκεται εντός της σφαίρας, δηλαδή το άθροισμα των τετραγώνων των συντεταγμένων του να είναι μικρότερο ή ίσο του 1 και ανάλογα μεγαλώνουμε κατά ένα την τιμή count_inside. Επίσης ελέγχουμε μετα από τα n σημεία αν το count_inside είναι 0 για να τερματίσουμε τη διαδικασία. Τελικά υπολογίζουμε την εκτίμηση του όγκου σύμφωνα με την εξίσωση (1.2) και τυπώνουμε την διάσταση του όγκου, τον πραγματικό όγκο, την εκτίμηση καθώς και το ποσοστό σφάλματος. Καλούμε τη συνάρτηση για διάσταση 6 και πλήθος σημείων 10^5 και παίρνουμε:

```

1 > ask1a(6,10^5)
2 [1] "D: 6"
3 [1] "Real_Volume: 5.16771278004997"
4 [1] "Estimated_Volume: 5.17696"
5 [1] "The_percentage_of_error_is: 0.179%"

```

Αν τώρα επιλέξουμε 17 διαστάσεις παίρνουμε το παρακάτω αποτέλεσμα:

```

1 > ask1a(17,10^5)
2 [1] "D: 16"
3 [1] "Real_Volume: 0.235330630358893"
4 [1] "Estimated_Volume: 0"
5 [1] "The_percentage_of_error_is: 100%"

```

Στις 16 διαστάσεις η διαδικασία σταμάτησε αφού κανένα σημείο δεν βρέθηκε εντός της σφαίρας.

β. Στο προηγούμενο ερώτημα παρατηρήσαμε ότι για μεγάλο αριθμό διαστάσεων, όπου ο όγκος είναι πολύ μικρός η μέθοδος αποτυγχάνει. Μια λύση σε αυτό το πρόβλημα είναι να δημιουργήσουμε μια μαρκοβιανή αλυσίδα η οποία, αντί να επιλέγει τυχαία σημεία, το κάθε σημείο εξαρτάται από το προηγούμενο και η αναλλοίωτη κατανομή της αλυσίδας θα είναι παρόμοια με αυτή που αναζητούμε. Στην περίπτωση μας μπορούμε να θεωρήσουμε για χώρο καταστάσεων της αλυσίδας τον κύλινδρο $C_d = B_{d-1} \times [-1, 1]$ και αναλλοίωτη κατανομή την ομοιόμορφη στον C_d . Συνεπώς η αλυσίδα αυτή θα μας βοηθήσει να επιλέγουμε σημεία στον χώρο του κυλίνδρου προσομοιώνοντας την ομοιόμορφη κατανομή. Για προτεινόμενη κατανομή θα πάρουμε εκείνη που αρχικά από τις πρώτες $d-1$ διαστάσεις επιλέγει μια τυχαία και τη μεταβάλλει κατά μια ομοιόμορφη τυχαία μεταβλητή στο $(-ε, ε)$, ενώ για την τελευταία διάσταση την επιλέγει ομοιόμορφα στο $(-1, 1)$. Τελικά η αλυσίδα μας θα λειτουργεί με τον εξής τρόπο, θα επιλέγει τυχαίες αρχικές τιμές για κάθε διάσταση μέσα στον κύλινδρο, στο επόμενο βήμα μέσω της προτεινόμενης κατανομής θα έχουμε έναν νέο συνδυασμό συντεταγμένων για την αλυσίδα όπου στη συνέχεια θα ελέγχουμε αν πληρεί τις προϋποθέσεις, δηλαδή αν βρίσκεται μέσα στην σφαίρα d - διαστάσεων. Σε περίπτωση που ισχύει αυτό θα μετράμε το πλήθος των επιτυχημένων θέσεων. Μετά από αρκετές επαναλήψεις θα έχουμε μια εκτίμηση του λόγου των όγκων του κυλίνδρου προς τη σφαίρα. Αυτό από το εργοδικό θεώρημα όπου για $h(x)$ θεωρούμε την δείκτρια συνάρτηση της σφαίρας.

$$\frac{1}{T} \sum_{t=1}^T h(X^{(t)}) \rightarrow E_f[h(X)] \quad (1.3)$$

Τελικά έχουμε μια εκτίμηση για τον λόγο

$$\frac{V_d}{2 \times V_{d-1}} \quad (1.4)$$

Μποούμε έτσι επαγωγικά να εκτιμήσουμε τον όγκο V_d για $d = 3, 4, \dots, d_{max}$ όπου d_{max} η διάσταση στην οποία το σφάλμα ξεπερνάει το 10%

```

1 ask1b <- function(d, nreps = 10000, c) {
2   states <- matrix(0, nrow = nreps, ncol = d)
3   states[1,] <- runif(d, -1, 1)
4   accepted = 0
5
6   dims <- 1:(d-1)
7
8   for (i in 2:nreps) {
9     prop_state <- states[i - 1, ]
10
11     dim <- sample(dims, 1)
12     prop_state[dim] <- prop_state[dim] + runif(1, -c, c)
13
14     while (sum(prop_state[1:(d-1)]^2) > 1) {
15       prop_state[dim] <- prop_state[dim] + runif(1, -c, c)
16     }
17
18     prop_state[d] <- runif(1, -1, 1)
19
20

```

```

21   if (sum(prop_state^2) <= 1) {
22     accepted <- accepted + 1
23   }
24   states[i, ] <- prop_state
25 }
26
27 list(states = states, acceptance_rate = accepted / nreps)
28 }

```

Η συνάρτηση που δημιουργήσαμε δέχεται ως ορίσματα τη διάσταση, τον αριθμό βημάτων που θα πραγματοποιήσει καθώς και την σταθερά ϵ , ή c στον κώδικα για την μεταβολή των καταστάσεων. Αρχικοποιούμε το διάνυσμα που θα αποθηκεύει τις καταστάσεις από τις οποίες περνάει η αλυσίδα και επιλέγει τυχαία τις αρχικές τιμές των συντεταγμένων στο $[-1,1]$. Προς το παρόν έχουμε 0 επιτυχημένες καταστάσεις και βάζουμε σε ένα διάνυσμα όλες τις διαστάσεις εκτός της τελευταίας. Ξεκινάμε τη διαδικασία με ένα for loop όσος και ο αριθμός βημάτων και αρχικά η προτεινόμενη κατάσταση είναι η προηγούμενη κατάσταση. Έπειτα μεταβάλλουμε μια τυχαία από τις πρώτες $d-1$ διαστάσεις με τον τρόπο που εξηγήσαμε και πρίν, φροντίζοντας ο χώρος καταστάσεων να είναι ο κύλινδρος. Δηλαδή κάθε κατάσταση που προτείνεται και είναι εκτός της σφαίρας διαστάσεων $d-1$ ξεχνιέται και προτείνεται νέα συντεταγμένη χωρίς να αποθηκευτεί η προηγούμενη. Για την τελευταία διάσταση επιλέγουμε ομοιόμορφα στο $(-1,1)$. Στη συνέχεια ελέγχουμε αν η κατάσταση μετά τις μεταβολές βρίσκεται εντός της σφαίρας με d - διαστάσεις και αυξάνουμε κατά 1 τον μετρητή. Η κατάσταση αποθηκεύεται και συνεχίζουμε με το επόμενο βήμα. Τελικά επιστρέφουμε το σύνολο των καταστάσεων της μαρκοβιανής αλυσίδας καθώς και το σύνολο των σημείων που βρέθηκαν εντός της σφαίρας προς τα συνολικά σημεία που δημιουργήθηκαν.

Έχουμε λοιπόν την εκτίμηση για τον λόγο (1.4) και θα υπολογίσουμε επαγωγικά τον όγκο V_d με χρήση της παρακάτω συνάρτησης:

```

1 compute_vd <- function(d, nreps=10000, c=0.1){
2   # Initial values
3   v <- c(2, pi)
4
5   # Compute v(d) for d >= 3
6   for(i in 3:d){
7     x <- ask1b_trace(i, nreps, c)$acceptance_rate
8     v[i] <- 2 * x * v[i-1]
9
10    # Calculate the percentage error
11    real_vd = real_vol_calc(i)
12    percentage_error = abs((v[i] - real_vd) / real_vd) * 100
13
14    # If the percentage error is larger than 10%, stop the
15    ↪ calculation
16    if (percentage_error > 10){
17      print(paste("Stopping at dimension", i, "due to high
18      ↪ error percentage:", percentage_error))
19      break
20    }
21  }
22 }

```

```

20 result <- list(v = tail(v, n=1), real_vd=real_vd, percentage
    ↪ _error = percentage_error)
21 class(result) <- "vd"
22 return(result)
23 }
24
25 print.vd <- function(x) {
26   print(paste("V_=", x$v, ", Real_volume_=", x$real_vd, ", _Error
    ↪ _=", x$percentage_error, "%"))
27 }

```

Η παραπάνω επαναληπτική συνάρτηση ξεκινάει από 3 διαστάσεις και υπολογίζει με τον επαγωγικό τρόπο την εκτίμηση του όγκου καθώς και τον πραγματικό όγκο με τη συνάρτηση *real_vol_calc* του προηγούμενου ερωτήματος. Σε κάθε βήμα της η συνάρτηση ελέγχει το σφάλμα και την πρώτη φορά που υπερβαίνει το 10% τερματίζει. Στη συνέχεια θα πειραματιστούμε με διαφορετικές τιμές των n και ε , για να δούμε τι προτιμάμε και τι συμπεράσματα αντλούμε από τη μέθοδο.

Αρχικά θα μελετήσουμε την επίδραση του πλήθους βημάτων. Επιλέγουμε $d = 4$ και $\varepsilon = 0.1$ και τρέχουμε για πλήθος βημάτων $10^2, 10^3, 10^4, 10^5$.

```

1 print(compute_vd(4, 10^2, 0.1))
2 [1] "Stopping at dimension 3 due to high error percentage: 49
    ↪ %"
3 [1] "V_ = 2.13628300444106, Real_volume_ = 4.18879020478639,
    ↪ Error_ = 49%"
4 print(compute_vd(4, 10^3, 0.1))
5 [1] "Stopping at dimension 4 due to high error percentage:
    ↪ 12.3002787502778"
6 [1] "V_ = 4.32780777410284, Real_volume_ = 4.93480220054468,
    ↪ Error_ = 12.3002787502778%"
7 print(compute_vd(4, 10^4, 0.1))
8 [1] "V_ = 4.65900741500601, Real_volume_ = 4.93480220054468,
    ↪ Error_ = 5.58877082263251%"
9 print(compute_vd(4, 10^5, 0.1))
10 [1] "V_ = 4.88229785958542, Real_volume_ = 4.93480220054468,
    ↪ Error_ = 1.06396039447058%"

```

Παρατηρούμε πως αυξάνοντας τον αριθμό των βημάτων η ακρίβεια αυξάνεται καθώς και η σταθερότητα, αφού για μικρά n οι τιμές της εκτίμησης έχουν τεράστιες διακυμάνσεις. Αυτό οφείλεται στο burn-in, συγκεκριμένα η μέθοδος αυτή χρειάζεται ένα συγκεκριμένο αριθμό βημάτων ώστε να φτάσει στην αναλλοίωτη κατανομή και να έχει την εκτιμητική ικανότητα που χρειαζόμαστε. Συνεπώς χρειαζόμαστε μεγάλο αριθμό βημάτων και τα 10^5 είναι μια καλή ισορροπία σε απόδοση και υπολογιστικό κόστος. Στη συνέχεια θα αλλάξουμε τις τιμές του ε .

```

1 > print(compute_vd(4, 10^5, 0.001))
2 [1] "Stopping at dimension 3 due to high error percentage:
    ↪ 17.8"

```

```

3 [1] "V=3.44318554833441,Real volume=4.18879020478639,
   ↳ Error=17.8%"
4 > print(compute_vd(4,10^5,0.1))
5 [1] "V=4.88229785958542,Real volume=4.93480220054468,
   ↳ Error=1.06396039447058%"
6 > print(compute_vd(4,10^4,0.3))
7 [1] "V=4.92568929855576,Real volume=4.93480220054468,
   ↳ Error=0.184666003186762%"
8 > print(compute_vd(4,10^4,0.7))
9 [1] "V=4.9180380124798,Real volume=4.93480220054468,
   ↳ Error=0.339713475507358%"

```

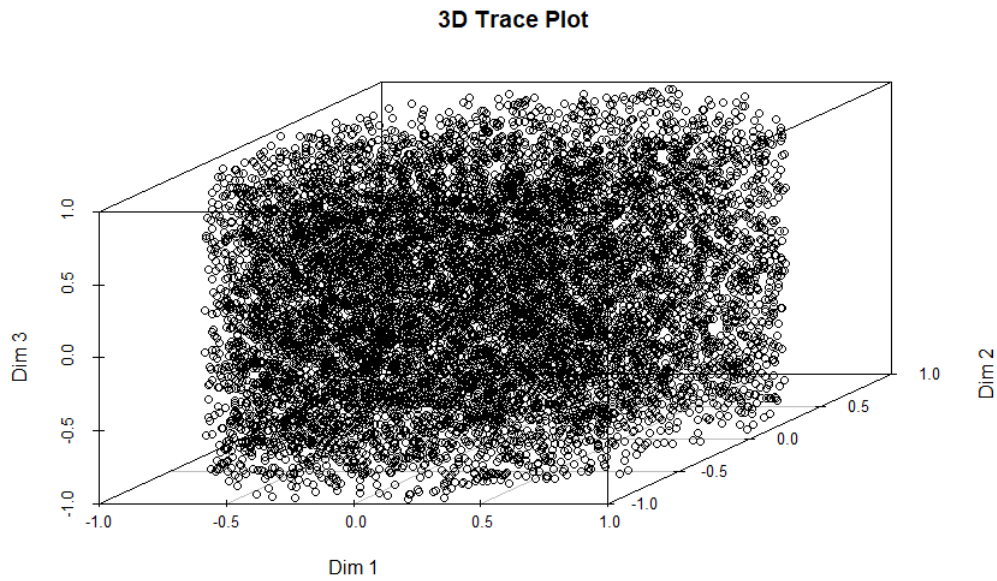
Για τιμές του ϵ μεγαλύτερες από 0.7 η συνάρτηση δεν ολοκληρώνεται καθώς κάθε σημείο, ή τα περισσότερα που προτείνονται βρίσκονται εκτός του κυλίνδρου και άρα απορρίπτονται, με αποτέλεσμα να μην ολοκληρώνεται η διαδικασία. Για τις τιμές που δοκιμάσαμε ξεκινήσαμε από πολύ μικρό ϵ , το οποίο έδωσε μεγάλο σφάλμα και ο λόγος είναι διότι δεν κινείται αρκετά γρήγορα η μαρκοβιανή αλυσίδα μέσα στον χώρο, με άλλα λόγια παραμένει πολύ κοντά στο σημείο εκκίνησης και δεν εξερευνά όλο το χώρο. Οι τιμές 0.1 και 0.3 είναι πολύ καλές, με την 0.3 να ξεχωρίζει ακόμα και σε 10^4 βήματα. Ιδανικά για να βρούμε τις βέλτιστες τιμές για τις παραμέτρους θα έπρεπε να έχουμε ένα πολύ πιο πυκνά διαμερισμένο διάστημα και να δοκιμάσουμε πολλές τιμές, όμως ο υπολογιστικός χρόνος είναι μεγάλος και μια σύντομη σύγκριση των τιμών αρκεί.

Τέλος, θα αναπαραστήσουμε γραφικά, σε 3 διαστάσεις για να είναι κατανοητό, κάθε κατάσταση της μαρκοβιανής αλυσίδας.

```

1 result <- ask1b(3, 10^5, 0.3)
2 scatterplot3d(result$states[, 1], result$states[, 2], result$
   ↳ states[, 3],
3               main="3D Trace Plot", xlab="Dim1", ylab="Dim2
   ↳ ", zlab="Dim3")

```

Με χρήση του παραπάνω κώδικα κάνουμε τη γραφική αναπαράσταση των συντεταγμένων των σημείων της μαρκοβιανής αλυσίδας και όπως βλέπουμε τα σημεία βρίσκονται όλα εντός του κυλίνδρου, ο οποίος είναι ο χώρος καταστάσεων. Επαληθεύουμε έτσι την διαδικασία που πραγματοποιήσαμε καθώς είναι φανερό πως το σύνολο των σημείων μέσα στη σφαίρα προς τα συνολικά σημεία, δηλαδή μέσα στον κύλινδρο, είναι μια εκτίμηση του λόγου των όγκων των δύο αυτών σχημάτων.

2 Άσκηση 2

Αρχικά η εκτιμημένη συνάρτηση πυκνότητας πιθανότητας για ένα τυχαίο δείγμα x_1, x_2, \dots, x_n με την μέθοδο του πυρήνα υπολογίζεται από τον τύπο:

$$\hat{f}(x) = \frac{1}{nh} K\left(\frac{x - x_i}{h}\right) \quad (2.1)$$

Ο όρος $\frac{x-x_i}{h}$ συμβολίζει την απόσταση του σημείου x από την παρατήρηση x_i τυποποιημένο ως προς τη σταθερά h την οποία θα υπολογίσουμε στη συνέχεια. Το μέγεθος του δείγματος συμβολίζεται με n ενώ η συνάρτηση του πυρήνα με K . Θα χρησιμοποιήσουμε τον κανονικό πυρήνα για τον οποίο ισχύει:

$$K(x) = \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{x^2}{2}\right)} \quad (2.2)$$

Για να υπολογίσουμε το h θα μεγιστοποιήσουμε την cross-validated πιθανοφάνεια, η οποία δίνεται από τον τύπο:

$$L(h) = \prod_{i=1}^n \hat{f}_{-i}(x_i) \quad (2.3)$$

όπου με $\hat{f}_{-i}(x)$ συμβολίζουμε την εκτίμηση της πυκνότητας χωρίς την i -οστή παρατήρηση. Τελικά λογαριθμίζοντας την πιθανοφάνεια έχουμε να μεγιστοποιήσουμε ως προς h της εξής σχέση:

$$l(h) = \sum_{i=1}^n \ln \left(\sum_{j=1, j \neq i}^n K\left(\frac{x - x_i}{h}\right) \right) - n \ln((n-1)h) \quad (2.4)$$

Τα δεδομένα της άσκησης λαμβάνονται από την βιβλιοθήκη MASS την οποία θα κατεβάσουμε και θα φορτώσουμε στην R. Τα δεδομένα περιέχουν 133 προσομοιώσεις ατυχήματος μοτοσυκλέτας, όπου κάθε προσομοίωση παρέχει δύο τιμές, την επιτάχυνση κεφαλής σε g που θα είναι η μεταβλητή απόκρισης στο μοντέλο παλινδρόμησης και τον χρόνο σε χιλιοστά δευτερολέπτου από την στιγμή της κρούσης που θα είναι η επεξηγηματική μεταβλητή.

```
1 library(MASS) #load MASS package
2 data(mcycle) #access the mcycle data
3 y<-mcycle$accel #response variable y
4 x<-mcycle$times #predictor variable x
```

α. Αρχικά θα ορίσουμε την συνάρτηση kernel στην R και στη συνέχεια την συνάρτηση likelihood η οποία θα υπολογίζει την λογαριθμική πιθανοφάνεια ενός σημείου x για δοσμένο h .

```
1 kernel_funct<-function(xi,xj,h){
2   x<-(xi-xj)/h
3   return(exp(-0.5*x^2)/sqrt(2*pi)) #Gaussian kernel
4 }
5
6
7 likelihood<-function(x,h){
```

```

8  n<-length(x) #data size
9  kernel_sum<-rep(0,n) #initiate vector
10
11 #Iterate over each element of x
12 for( i in 1:n){
13     temp<-0
14     # Compute the kernel sum for the current data point
15     for(j in x[-i]){
16         temp<-temp + kernel_funct(x[i],j,h)
17     }
18     kernel_sum[i]<-temp #store the value
19 }
20
21 return(sum(log(kernel_sum))-n*log((n-1)*h))
22 }

```

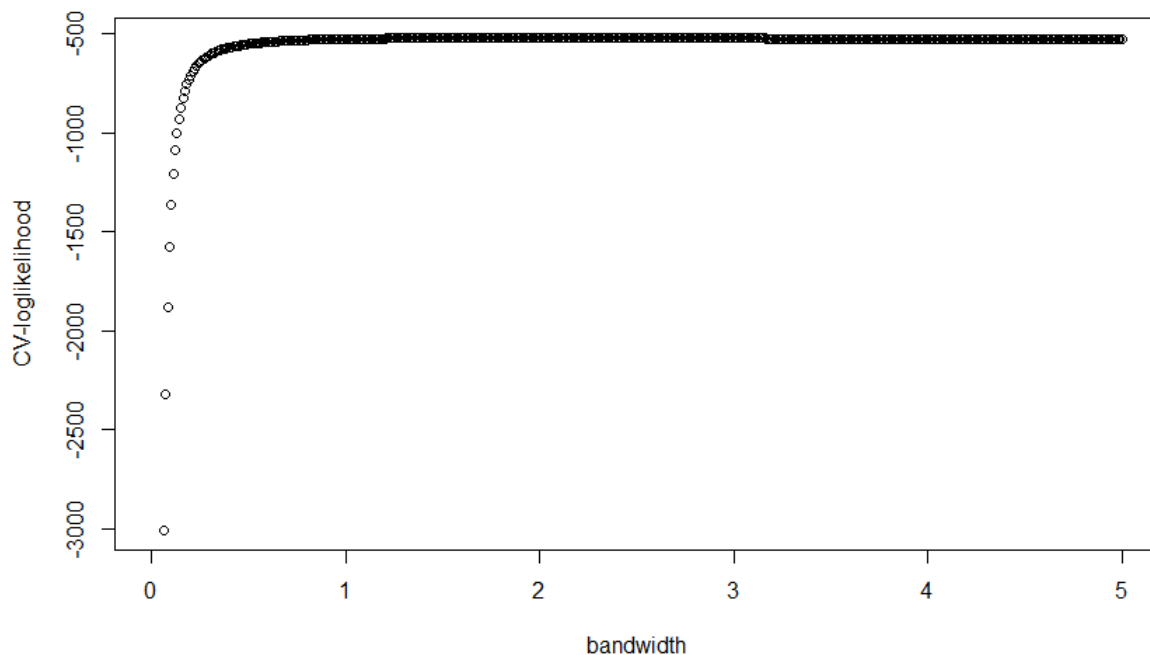
Η συνάρτηση kernel δέχεται ως ορίσματα τα x_i, x_j, h υπολογίζει τον όρο $\frac{x_i - x_j}{h}$ και επιστρέφει την τιμή του κανονικού πυρήνα $K(x)$. Η συνάρτηση likelihood δέχεται ως ορίσματα τα x, h , υπολογίζει το Kernel sum για κάθε i αφού θα χρειαστεί για την επιστροφή της τιμής της λογαριθμικής πιθανοφάνειας για τα x, h .

Στη συνέχεια για να υπολογίσουμε το βέλτιστο h θα διαμερίσουμε το διάστημα $[0.01, 10]$ με βήμα 0.01 και έτσι θα έχουμε 500 h για τα οποία θα υπολογίσουμε την τιμή της συνάρτησης likelihood με δεδομένα τις παρατηρήσεις της επεξηγηματικής μας μεταβλητής x . Το διάστημα διαμέρισης επιλέγεται από εμάς και καλό είναι να ξεκινήσουμε από ένα μεγάλο διάστημα με λίγο μεγαλύτερο βήμα και σταδιακά να συγκλινουμε σε μικρότερο διάστημα μειώνοντας το βήμα και άρα αυξάνοντας την ακρίβεια. Παρακάτω δίνεται ο κώδικας που υπολογίζει το βέλτιστο h καθώς και η γραφική αναπαράσταση της λογαριθμικής πιθανοφάνειας συναρτήσει του h .

```

1  h<-seq(from = 0.01, to = 5, by = 0.01)
2  length(h)
3  likelihood_vec<-rep(0,length(h)) #initialise vector
4  #compute likelihood for each h
5  for(k in 1:length(h)){
6      likelihood_vec[k]<-likelihood(x,h[k])
7  }
8  plot(h,likelihood_vec,xlab = "bandwidth",ylab = "CV-
   ↪ loglikelihood")
9  h_opt<-h[which.max(likelihood_vec)]

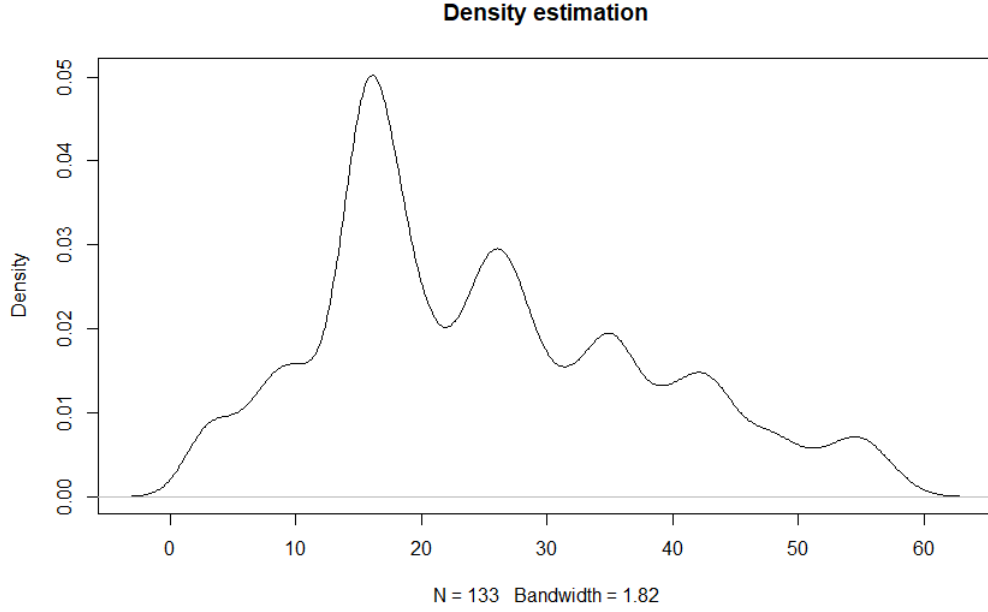
```



Η εντολή `which.max()` μας δίνει το index του h για το οποίο μεγιστοποιείται η πιθανοφάνεια και έτσι παίρνουμε το βέλτιστο h . Στη συνέχεια για να αναπαραστήσουμε γραφικά την εκτιμώμενη $f(x)$ με το βέλτιστο h θα χρησιμοποιήσουμε την παρακάτω εντολή:

```
1 plot(density(x, kernel = "gaussian", bw=h_opt), main = "Density
   ↪      estimation")
```

Η συνάρτηση `density()` της R κάνει αυτό ακριβώς που χρειαζόμαστε. Για ορίσματα βάζουμε τα δεδομένα x , το είδος του πυρήνα και το βέλτιστο h .



Η εικόνα είναι αρκετά ικανοποιητική καθώς δεν παρουσιάζεται πρόβλημα underfitting ή overfitting. Αυτό μπορούμε να το ελέγξουμε αν πειραματιστούμε με την τιμή του h όπου για μικρότερο h το γράφημα γίνεται πολύ απότομο και χάνει την προβλεπτική του ικανότητα για τιμές που δεν έχουν παρατηρηθεί, ενώ για μεγαλύτερο h το γράφημα γίνεται πολύ λείο και δεν παρέχει αρκετές πληροφορίες για την εκτίμηση. Επίσης στον άξονα x αναγράφεται και η τιμή του βέλτιστου h η οποία είναι $bandwidth = 1.82$.

β. Στη συνέχεια θα προσαρμόσουμε το μοντέλο μη παραμετρικής παλινδρόμησης Nadaraya-Watson με τις ίδιες μεταβλητές όπως και πριν. Το h θα υπολογιστεί με την μέθοδο του leave-one-out CV.

Ας υποθέσουμε πως έχουμε δύο τυχαίες μεταβλητές Y, X , τότε το απλό γραμμικό μοντέλο ορίζεται ως εξής:

$$m(x) = E[Y|X = x] = \int y f(y|x) dy = \int y \frac{f(x, y)}{f_X(x)} dy \quad (2.5)$$

Όμως οι $f(x, y)$ και $f_X(x)$ είναι άγνωστες, οπότε θα χρησιμοποιήσουμε τις ακόλουθες προσεγγίσεις.

$$\hat{f}(x, y) = \frac{1}{nh_x h_y} \sum_{i=1}^n K_X\left(\frac{x - x_i}{h_x}\right) K_Y\left(\frac{y - y_i}{h_y}\right) \quad (2.6)$$

$$\hat{f}_X(x) = \frac{1}{nh_x} \sum_{i=1}^n K_X\left(\frac{x - x_i}{h_x}\right) \quad (2.7)$$

Εύκολα αποδυνκνείται πως το μοντέλο Nadaraya-Watson είναι το εξής:

$$\hat{m}(x) = \frac{\sum_{i=1}^n K_X\left(\frac{x - x_i}{h_x}\right) y_i}{\sum_{i=1}^n K_X\left(\frac{x - x_i}{h_x}\right)} = \sum_{i=1}^n w_i y_i = \hat{m}_{NW}(x) \quad (2.8)$$

Για την εκτιμήτρια της διασποράς του εν λόγω μοντέλου ισχύει:

$$\hat{\sigma}^2(x) = \frac{\sum_{i=1}^n K_X\left(\frac{x-x_i}{h_x}\right) e_i^2}{\sum_{i=1}^n K_X\left(\frac{x-x_i}{h_x}\right)}, \quad e_i^2 = y_i - \hat{m}(x) \quad (2.9)$$

Για τον υπολογισμό του βέλτιστου h όπως αναφέρθηκε και πριν θα χρησιμοποιήσουμε τη μέθοδο leave-one-out CV.

```

1 h<-seq(1.01,10,0.01) #sequence of h to choose from
2 cv_mse<-rep(0,length(h)) #Initialise mse vector for each h
3 n<-length(x)
4
5 for( i in 1:length(h)){ #loopoing over every h
6   total_mse<-0
7
8   for(j in 1:n){      #leave-one-out CV
9     x_train<-x[-j]
10    y_train<-y[-j]
11    x_test<-x[j]
12    y_test<-y[j]
13
14    w<-(1/sqrt(2*pi))*exp(-0.5*((x_train-x_test)/h[i])^2) #
      ↳ calculating the kernel
15    estimate<-sum(w*y_train)/sum(w)
16
17    mse<-(y_test-estimate)^2
18    total_mse<-total_mse + mse #sum of mse for each
      ↳ observation missing
19  }
20  cv_mse[i]<-total_mse/n #average of the total sum for
      ↳ current h
21 }
22 min_mse<-which.min(cv_mse)
23 h_opt<-h[min_mse]
```

Αρχικά δημιουργούμε ένα διάνυσμα με h ξεκινώντας από το 1.01 έως το 10 με βήμα 0.01 όπως ζητείται στην άσκηση. Στη συνέχεια θα δημιουργήσουμε ένα κένο διάνυσμα μεγέθους ίσου με το μήκος των h για να αποθηκεύσει τις τιμές των MSE για κάθε h . Ξεκινάμε τη διαδικασία με loop για κάθε h και δημιουργούμε εσωτερικό βρόγχο που αφαιρεί κάθε φορά μια διαφορετική παρατήρηση από το x και υπολογίζει το MSE για αυτά τα δεδομένα. Τέλος παίρνει το άθροισμα και βρίσκει τη μέση του τιμή αποθηκεύοντας στο διάνυσμα που θα περιέχει τα MSE για κάθε h . Οι τελευταίες δύο εντολές μας δίνουν το βέλτιστο h το οποίο είναι μάλιστα και πρώτο στο διάνυσμα και είναι το 1.01. Τώρα μπορώ να προσαρμόσω το τελικό μοντέλο και να κάνω τη γραφική αναπαράσταση.

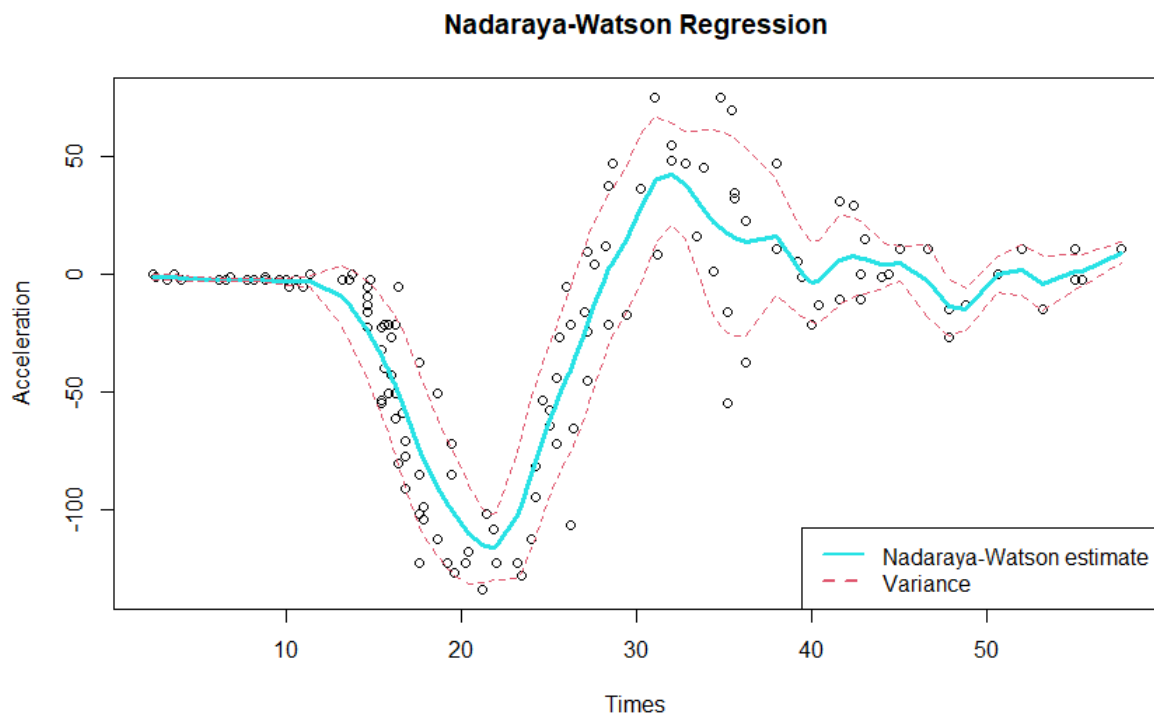
```

1 pred_values<-rep(0,length(x)) #Initialize vectors for
      ↳ predicted values
2 variance<-rep(0,length(x))      #and variance
3
```

```

4 for( i in 1:length(x)){
5   weights<-(1/sqrt(2*pi))*exp(-0.5*((x-x[i])/h_opt)^2)
6
7   pred_values[i]<-sum(weights*y)/sum(weights)
8   variance[i]<-sum(weights*(y-pred_values[i])^2)/sum(weights)
9 }
10 #plotting data points
11 plot(x,y,main="Nadaraya-Watson Regression",xlab = "Times",
12      → ylab = "Acceleration")
13 lines(x,pred_values,col=5,lw=3) #regression line on cyan
14      → colour
15 lines(x,pred_values+sqrt(variance),col=2,lty=2) #variance
16      → dotted lines on red
17 lines(x,pred_values-sqrt(variance),col=2,lty=2)
18 legend("bottomright",legend = c("Nadaraya-Watson estimate",
19      → "Variance"),col=c(5,2),
20      lty = c(1,2),lwd = 2)

```



Στον άξονα x έχουμε την επεξηγηματική μεταβλητή, χρόνος σε χιλιοστά του δευτερολέπτου από την στιγμή της κρούσης και στον άξονα y την μεταβλητή απόκρισης, επιτάχυνση κεφαλής σε g . Για να δημιουργήσουμε αυτήν την εικόνα αρχικά δημιουργήσαμε δύο διανύσματα που θα αποθηκεύσουν τις εκτιμήσεις και τις διασπορές. Στη συνέχεια προσαρμόζουμε το μοντέλο μη παραμετρικής παλινδρόμησης Nadaraya-Watson με το βέλτιστο h που υπολογίσαμε και αποθηκεύουμε για κάθε παρατήρηση την εκτίμηση και την διασπορά της. Για την γραφική αναπαράσταση βάζουμε πρώτα τις παρατηρήσεις μας. Στη συνέχεια τις εκ-

τιμήσεις σε μορφή συνεχούς γραμμής και προσθέτουμε και αφαιρούμε μια διασπορά στις εκτιμήσεις. Το μοντέλο μας φαίνεται να προσαρμόζεται πολύ καλά στα δεδομένα.

3 Άσκηση 3

α. Έστω x_1, x_2, \dots, x_n ανεξάρτητες αυτόνομες τυχαίες μεταβλητές που ακολουθούν την κανονική κατανομή 2 διαστάσεων $N(m, \Sigma)$. Θέλουμε να υπολογίσουμε τις εκτιμήτριες μέγιστης πιθανοφάνειας των παραμέτρων, m και Σ της κανονικής κατανομής σε 2 διαστάσεις. Η συνάρτηση πυκνότητας πιθανότητας της διδιάστατης κανονικής κατανομής δίνεται από τον τύπο:

$$f_X(x_j; m, \Sigma) = (2\pi)^{-1} |\Sigma|^{-1/2} \exp \left(-\frac{1}{2} (x_j - m)^T \Sigma^{-1} (x_j - m) \right) \quad (3.1)$$

Όπου m το 2×1 διάνυσμα των μέσων τιμών, Σ ο 2×2 πίνακας συνδιακύμανσης, $|\Sigma|$ η οριζουσα του και x^T συμβολίζουμε τον ανάστροφο πίνακα. Η συνάρτηση πιθανοφάνειας υπολογίζεται ως εξής:

$$\begin{aligned} L(m, \Sigma; x_1, x_2, \dots, x_n) &= \prod_{i=1}^n f_X(x_i; m, \Sigma) \\ &= \prod_{i=1}^n (2\pi)^{-1} |\Sigma|^{-1/2} \exp \left(-\frac{1}{2} (x_i - m)^T \Sigma^{-1} (x_i - m) \right) \\ &= (2\pi)^{-n} |\Sigma|^{-n/2} \exp \left(-\frac{1}{2} \sum_{i=1}^n (x_i - m)^T \Sigma^{-1} (x_i - m) \right) \end{aligned} \quad (3.2)$$

Λογαριμίζοντας έχουμε την log-likelihood συνάρτηση:

$$l(m, \Sigma; x_1, x_2, \dots, x_n) = -\frac{n}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma|^{-1} - \frac{1}{2} \sum_{i=1}^n (x_i - m)^T \Sigma^{-1} (x_i - m) \quad (3.3)$$

Πρωτού υπολογίσουμε τις εκτιμήτριες υπενθυμίζουμε μερικά πράγματα.

- $\text{tr}(a) = a$, όταν a αριθμός
- $\text{tr}(AB) = \text{tr}(BA)$
- $\text{tr}(aA + bB) = a\text{tr}(A) + b\text{tr}(B)$
- $\nabla_A(\ln |A|) = (A^{-1})^T$
- $\nabla_A(\text{tr}(BA)) = B^T$
- $\nabla_x(x^T A x) = 2Ax$

Για να βρούμε τις εκτιμήτριες μέγιστης πιθανοφάνειας των m και Σ θα μηδενίσω τις παραγώγους ως προς κάθε παράμετρο και θα λύσω ως προς αυτές. Ξεκινάω για την \hat{m} .

$$\begin{aligned}
& \nabla_m (l(m, \Sigma; x_1, x_2, \dots, x_n)) \\
&= \nabla_m \left(-\frac{n}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma|^{-1} - \frac{1}{2} \sum_{i=1}^n (x_i - m)^T \Sigma^{-1} (x_i - m) \right) \\
&= \sum_{i=1}^n \Sigma^{-1} (x_i - m) \\
&= \Sigma^{-1} \sum_{i=1}^n (x_i - m) = 0
\end{aligned} \tag{3.4}$$

Το οποίο ισχύει μόνο όταν

$$\begin{aligned}
& \sum_{i=1}^n x_i - nm = 0 \\
& \hat{m} = \frac{1}{n} \sum_{i=1}^n x_i
\end{aligned} \tag{3.5}$$

Ακολουθώ την ίδια διαδικασία και για το $\hat{\Sigma}$.

$$\begin{aligned}
& \nabla_{\Sigma^{-1}} (l(m, \Sigma; x_1, x_2, \dots, x_n)) \\
&= \nabla_{\Sigma^{-1}} \left(-\frac{n}{2} \ln 2\pi + \frac{n}{2} \ln |\Sigma|^{-1} - \frac{1}{2} \sum_{i=1}^n (x_i - m)^T \Sigma^{-1} (x_i - m) \right) \\
&= \frac{n}{2} \nabla_{\Sigma^{-1}} (\ln |\Sigma|^{-1}) - \frac{1}{2} \nabla_{\Sigma^{-1}} \left(\sum_{i=1}^n \text{tr}[(x_i - m)^T \Sigma^{-1} (x_i - m)] \right) \\
&= \frac{n}{2} \Sigma^T - \frac{1}{2} \nabla_{\Sigma^{-1}} \left(\sum_{i=1}^n \text{tr}[(x_i - m)(x_i - m)^T \Sigma^{-1}] \right) \\
&= \frac{n}{2} \Sigma^T - \frac{1}{2} \nabla_{\Sigma^{-1}} \left(\text{tr} \left[\sum_{i=1}^n (x_i - m)(x_i - m)^T \Sigma^{-1} \right] \right) \\
&= \frac{n}{2} \Sigma^T - \frac{1}{2} \left(\sum_{i=1}^n (x_i - m)(x_i - m)^T \right)^T \\
&= \frac{n}{2} \Sigma - \frac{1}{2} \left(\sum_{i=1}^n (x_i - m)(x_i - m)^T \right) = 0
\end{aligned}$$

Οπότε τελικά έχω:

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{m})(x_i - \hat{m})^T \tag{3.6}$$

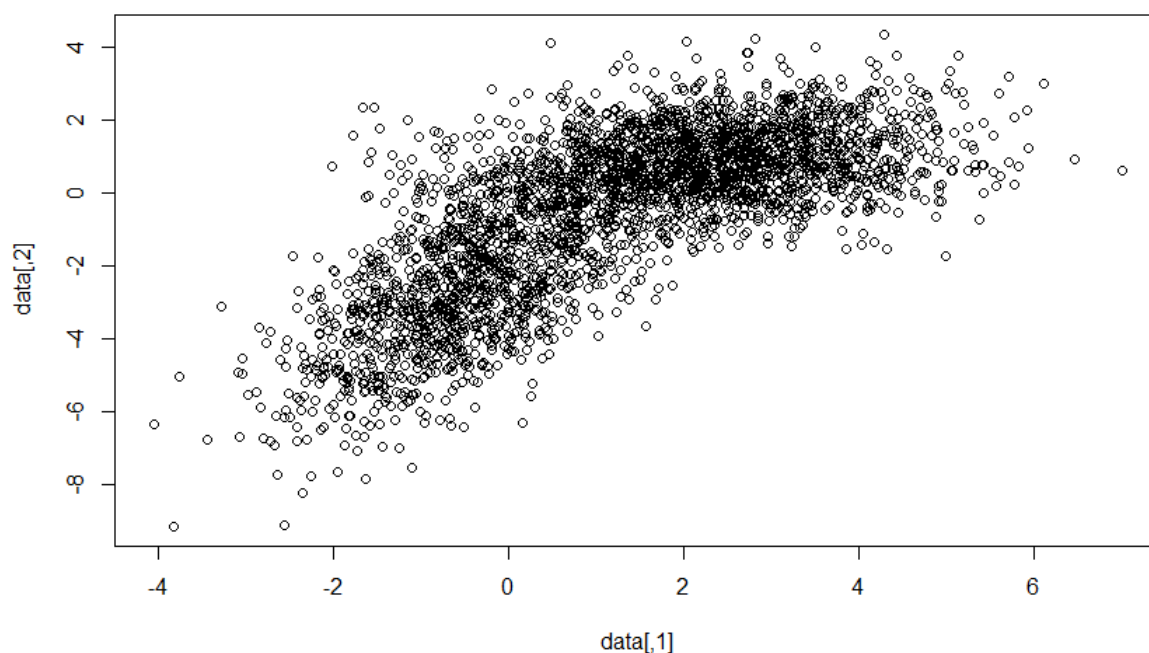
Αφού βρήκαμε αναλυτικά τις εκτιμήτριες μέγιστης πιθανοφάνειας, θα προσαρμόσουμε στα δεδομένα μας ένα μοντέλο κανονικής κατανομής δύο διαστάσεων. Αρχικά φορτώνουμε τα δεδομένα καθώς και την βιβλιοθήκη MASS στην R.

```

1 library("MASS")
2 data<-unname(as.matrix(read.csv2("datafile.csv"))))

```

Τα δεδομένα μας περιέχουν συντεταγμένες 3000 σημείων τα οποία απεικονίζονται στο παρακάτω σχήμα.



```

1 means<-c(mean(data[,1]),mean(data[,2]))
2 s11<-sum((data[,1]-means[1])^2)/n
3 s22<-sum((data[,2]-means[2])^2)/n
4 s12<-sum((data[,1]-means[1])*(data[,2]-means[2]))/n
5 cov_matrix<-matrix(c(s11,s12,s12,s22),nrow = 2)

```

Με τον παραπάνω κώδικα αρχικά δημιουργούμε το διάνυσμα που περιέχει τους μέσους των στηλών των παρατηρήσεων και στη συνέχεια υπολογίζουμε τη διασπορά κάθε στήλης καθώς και τη συνδιασπορά τους. Τελικά τα βάζουμε σε έναν πίνακα και έχουμε τον πίνακα συνδιακύμανσης.

```

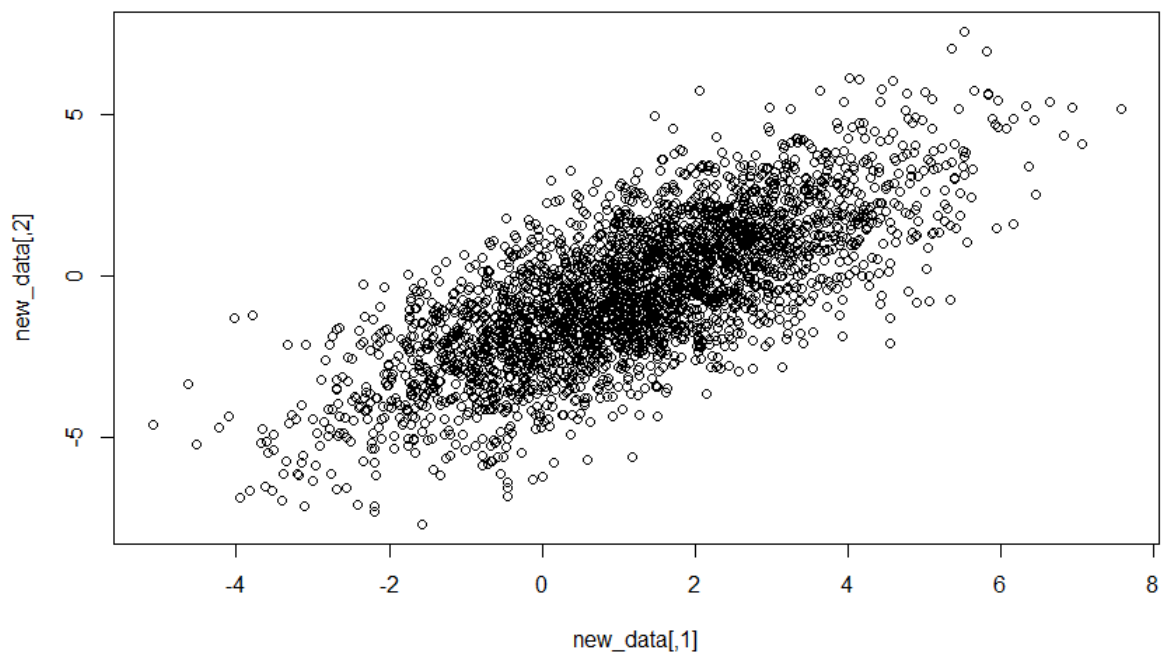
1 > means
2 [1] 1.1797288 -0.5346472
3 > cov_matrix
4      [,1]      [,2]
5 [1,] 3.412038 3.068844
6 [2,] 3.068844 5.054024

```

Οπότε τα δεδομένα μας ακολουθούν διδιάσταση κανονική κατανομή με διάνυσμα μέσων το means και πίνακα συνδιακύμανσης τον cov_matrix.

β. Στη συνέχεια θέλουμε να προσομοιώσουμε 3000 νέες παρατηρήσεις από την κατανομή που εκτιμήσαμε στο προηγούμενο ερώτημα. Θα χρησιμοποιήσουμε την εντολή `mvrnorm` η οποία κάνει την ίδια δουλειά με την `rnorm` αλλά για πολυδιάστατη κανονική κατανομή. Για ορίσματα θα βάλουμε το μέγεθος του δείγματος που θέλουμε, μετά το διάνυσμα των μέσων και τέλος τον πίνακα συνδιακύμανσης.

```
1 n<-nrow(data)
2 new_data <- mvrnorm(n = n, mu = means, Sigma = cov_matrix)
3 plot(new_data)
```



Οι εικόνες είναι παρόμοιες και υπάρχει σίγουρα μια συσχέτιση αλλά όπως βλέπουμε η καμπύλη που υπήρχε στην πρώτη γραφική δεν υπάρχει εδώ κάτι που μας οδηγεί στο να πιστέψουμε πως τα δεδομένα δεν ακολουθούν κανονική κατανομή.

γ. Θέλουμε να προσαρμόσουμε τώρα στα δεδομένα μας ένα μοντέλο μείζης κανονικών κατανομών

$$f_x(x) = pf_1(x) + (1 - p)f_2(x) \quad (3.7)$$

Όπου f_1, f_2 είναι συναρτήσεις πυκνότητας πιθανότητας κανονικών κατανομών δύο διαστάσεων $N(m_1, \Sigma_1)$ και $N(m_2, \Sigma_2)$ και $p \in (0, 1)$. Επειδή δεν μπορούμε να υπολογίσουμε αναλυτικά τις εκτιμήτριες μέγιστης πιθανοφάνειας των παραμέτρων με τις παρατηρήσεις που έχουμε, θα θεωρήσουμε την λανθάνουσα μεταβλητή $Z \sim \text{Bernulli}(p)$, η οποία κωδικοποιεί την κατανομή μείζης από την οποία προέρχεται η κάθε παρατήρηση. Για να εκτιμήσουμε τις παραμέτρους θα υλοποιήσουμε την μέθοδο Expectation-Maximization. Για την Z ισχύει

ότι

$$z = \begin{cases} 1, & \text{αν η παρατήρηση } x \text{ ανήκει στην κατανομή } f_1 \\ 0, & \text{αλλιώς} \end{cases} \quad (3.8)$$

με $P[z = 1] = p$ και $P[z = 2] = 1 - p$. Επίσης θα ισχύει $(X|z = i) \sim N(m_i, \Sigma_i), i = 1, 2$, οπότε η από κοινού συνάρτηση πυκνότητας πιθανότητας των X, Z είναι

$$f(X, Z) = [f(X, z = 1)P(z = 1)]^z [f(X, z = 2)P(z = 2)]^{1-z} \quad (3.9)$$

Θέλουμε να υπολογίσουμε την πιθανοφάνεια αυτής της κατανομής.

$$\begin{aligned} L(p, m_1, m_2, \Sigma_1, \Sigma_2; x_1, x_2, \dots, x_n) \\ &= \prod_{i=1}^n f(x_i, z_i) \\ &= \prod_{i=1}^n [f(x_i, z_i = 1)P(z_i = 1)]^{z_i} [f(x_i, z_i = 2)P(z_i = 2)]^{1-z_i} \\ &= \prod_{i=1}^n [pN(x_i; m_1, \Sigma_1)]^{z_i} [(1-p)N(x_i; m_2, \Sigma_2)]^{1-z_i} \end{aligned} \quad (3.10)$$

Παίρνω τον λογάριθμο της πιθανοφάνειας.

$$\begin{aligned} l(p, m_1, m_2, \Sigma_1, \Sigma_2; x_1, x_2, \dots, x_n) \\ &= \sum_{i=1}^n z_i \ln N(x_i; m_1, \Sigma_1) + \sum_{i=1}^n (1 - z_i) \ln N(x_i; m_2, \Sigma_2) + \ln p \sum_{i=1}^n z_i + \ln(1 - p) \sum_{i=1}^n (1 - z_i) \end{aligned} \quad (3.11)$$

Επειδή δεν γνωρίζω τα z_i θα τα αντικαταστήσω με την τιμή $E[z_i|x_i, \theta^{(0)}]$, όπου $\theta^{(0)}$ είναι οι αρχικές τιμές των παραμέτρων που επιλέγουμε. Από το θεώρημα του Bayes έχουμε

$$\begin{aligned} E[z_i|x_i, \theta^{(0)}] &= P[z_i = 1|x_i, \theta^{(0)}] \\ &= \frac{P[z_i = 1]f(X_i, z_i = 1)}{\sum_{j=1}^2 P[z_i = j]f(X_i, z_i = j)} \\ &= \frac{p^{(0)}N(x_i; m_1^{(0)}, \Sigma_1^{(0)})}{p^{(0)}N(x_i; m_1^{(0)}, \Sigma_1^{(0)}) + (1 - p^{(0)})N(x_i; m_2^{(0)}, \Sigma_2^{(0)})} = \gamma_i \end{aligned} \quad (3.12)$$

Τον όρο γ_i τον ονομάζουμε responsibility και το αντικαθιστούμε στη σχέση (3.11)

$$\begin{aligned} Q(\theta|X, \theta^{(0)}) \\ &= \sum_{i=1}^n \gamma_i \ln N(x_i; m_1, \Sigma_1) + \sum_{i=1}^n (1 - \gamma_i) \ln N(x_i; m_2, \Sigma_2) + \ln p \sum_{i=1}^n \gamma_i + \ln(1 - p) \sum_{i=1}^n (1 - \gamma_i) \end{aligned} \quad (3.13)$$

Θα βρούμε τώρα τις παραμέτρους που μεγιστοποιούν την ποσότητα Q παραγωγίζοντας και μηδενίζοντας ως προς κάθε παράμετρο. Οι πράξεις είναι παρόμοιες με του ερωτήματος α και

τελικά καταλήγουμε:

$$\hat{m}_1 = \frac{\sum_{i=1}^n \gamma_i x_i}{\sum_{i=1}^n \gamma_i} \quad (3.14)$$

$$\hat{m}_2 = \frac{\sum_{i=1}^n (1 - \gamma_i) x_i}{\sum_{i=1}^n \gamma_i} \quad (3.15)$$

$$\hat{\Sigma}_1 = \frac{\sum_{i=1}^n \gamma_i (x_i - \hat{m}_1)(x_i - \hat{m}_1)^T}{\sum_{i=1}^n \gamma_i} \quad (3.16)$$

$$\hat{\Sigma}_2 = \frac{\sum_{i=1}^n (1 - \gamma_i) (x_i - \hat{m}_2)(x_i - \hat{m}_2)^T}{\sum_{i=1}^n \gamma_i} \quad (3.17)$$

$$\hat{p} = \frac{\sum_{i=1}^n \gamma_i}{n} \quad (3.18)$$

Το μόνο που έχει αλλάξει είναι ότι αθροίζουμε κάθε παρατήρηση επί την ευθύνη (responsibility) της προς το συνολικό άθροισμα των γ_i . Για την εκτιμήτρια του p παίρνουμε την μέση τιμή των γ_i . Οπότε ξεκινάμε τον αλγόριθμό μας ορίζοντας τις αρχικές τιμές, στη συνέχεια υπολογίζουμε τα γ_i βάση των αρχικών μας παραμέτρων και μετά ανανεώνουμε τις παραμέτρους βάση των σχέσεων (3.14) έως (3.18) για να μεγιστοποιήσουμε την πιθανοφάνεια. Παρακάτω δίνεται ο κώδικας που υλοποιεί αυτή τη μέθοδο.

```

1 library(mvtnorm)
2 EM <- function(data, tolerance = 1e-6){
3
4   n <- nrow(data)
5   d <- ncol(data)
6   m1<-c(-1,-2) #initial values
7   m2<-c(1.5,1)
8   S1<-matrix(c(0.7,1,1,3),nrow = 2)
9   S2<-matrix(c(1.8,0.5,0.5,1),nrow = 2)
10  lambda<-0.5
11
12  converged <- FALSE
13  old.params <- c(lambda, m1, as.vector(S1), m2, as.vector(S2
14    ↪ ))
15  params.history <- list()
16
17  while (!converged){
18
19    params.history[[length(params.history) + 1]] <- list("
20      ↪ lambda" = lambda, "m1" = m1, "S1" = S1, "m2" = m2,
21      ↪ "S2" = S2) #for error calculation
22    # E-step: calculate the responsibility
23    prob1 <- lambda * dmvtorm(data, mean = m1, sigma = S1,log
24      ↪ = TRUE) #log likelihoods
25    prob2 <- (1 - lambda) * dmvtorm(data, mean = m2, sigma =
26      ↪ S2,log = TRUE)
27    tau1 <- prob1 / (prob1 + prob2) #responsibilities

```

```

23   tau2 <- 1-tau1
24
25   # M-Step: update the parameters to maximize the expected
      ↪ log-likelihood
26   lambda <- sum(tau1)/n
27   m1 <- colSums(tau1 * data) / sum(tau1)
28   m2 <- colSums(tau2 * data) / sum(tau2)
29   S1 <- t(tau1 * (data - m1)) %*% (data - m1) / sum(tau1)
30   S2 <- t(tau2 * (data - m2)) %*% (data - m2) / sum(tau2)
31   new.params <- c(lambda, m1, as.vector(S1), m2, as.vector(
      ↪ S2))
32   param.diff <- sum(abs(new.params - old.params))
33   if (param.diff < tolerance){
34     converged <- TRUE
35   }
36   old.params <- new.params
37 }
38
39
40   return(list("final" = list("lambda" = lambda, "m1" = m1, "
      ↪ S1" = S1, "m2" = m2, "S2" = S2), "history" = params.
      ↪ history))
41 }
42
43 results<-EM(data)
44 results$final
45 results$history
46 length(results$history)

```

Για αρχικές τιμές έχουμε πολλές επιλογές, μπορούμε να ξεκινήσουμε με τους μέσους κάθε στηλης συν μιας τυχαίας μεταβλητής, μπορούμε να εφαρμόσουμε την μέθοδο k-clusters ή ακόμα και να ξεκινήσουμε από το 0. Ένας πιο αντιαισθητικός τρόπος και αναξιόπιστος, αλλά χρήσιμος στην περίπτωση μας είναι να παρατηρήσουμε τη γραφική παράσταση των σημείων. Με λίγο πειραματισμό για τα νούμερα μπορούμε να έχουμε μια αρχική εκτίμηση των παραμέτρων δημιουργώντας νέα σημεία και βλέποντας αν το νέο σχήμα μοιάζει σε αυτό που θέλουμε. Η μέθοδος θα εξηγηθεί καλύτερα στο ερώτημα δ αλλά επιλέξαμε για αρχικές τιμές των παραμέτρων αυτές που φαίνονται στον κώδικα. Αφού επιλέξαμε αρχικές τιμές ξεκινάμε ένα while loop που θα τρέχει εως ότου το άθροισμα των σφαλμάτων κάθε παραμέτρου είναι μικρότερο απο 10^{-6} , δημιουργούμε μια λίστα με τις παραμέτρους του μοντέλου ώστε να τις αποθηκεύουμε σε κάθε βήμα. Στη συνέχεια για το E-step της μεθόδου υπολογίζουμε τα responsibilities με χρήση της εξίσωσης (3.12), για την συνάρτηση πυκνότητας πιθανότητας διδιάστατης κανονικής κατανομής χρησιμοποιούμε τη βιβλιοθήκη mtnorm που περιέχει την εντολή dmtnorm η οποία όταν δέχεται το log = TRUE επιστρέφει τις πιθανοφάνειες για τα δεδομένα. Για το M-step υπολογίζουμε τις νέες τιμές των παραμέτρων. Αποθηκεύουμε τις νέες τιμές των παραμέτρων στη λίστα μας και υπολογίζουμε το άθροισμα των απολύτων διαφορών των παραμέτρων και ελέγχουμε αν είναι μικρότερο του 10^{-6} σε περίπτωση που είναι, τερματίζουμε τη διαδικασία και επιστρέφουμε τη λίστα με τις παραμέτρους σε κάθε βήμα και τις τελευταίες παραμέτρους που υπολόγισε πριν τερματίσει, αλλιώς συνεχίζει η

διαδικασία. Έπειτα καλούμε τη συνάρτηση και παίρνουμε τα αποτελέσματα:

```
1 > results$final
2 $lambda
3 [1] 0.5723955
4
5 $m1
6 [1] 1.1797288 -0.5346472
7
8 $S1
9           [,1]      [,2]
10 [1,] 4.948984 3.153103
11 [2,] 3.153103 6.624679
12
13 $m2
14 [1] 1.1797288 -0.5346472
15
16 $S2
17           [,1]      [,2]
18 [1,] 4.948984 3.153103
19 [2,] 3.153103 6.624679
20
21 #-----
22 > results$history
23 [[1]]
24 [[1]]$lambda
25 [1] 0.5
26
27 [[1]]$m1
28 [1] -1 -2
29
30 [[1]]$S1
31           [,1] [,2]
32 [1,] 0.7      1
33 [2,] 1.0      3
34
35 [[1]]$m2
36 [1] 1.5 1.0
37
38 [[1]]$S2
39           [,1] [,2]
40 [1,] 1.8 0.5
41 [2,] 0.5 1.0
42
43
44 [[2]]
45 [[2]]$lambda
46 [1] 0.5798949
47
```



```

48 [[2]]$m1
49 [1] 1.919493 0.302922
50
51 [[2]]$S1
52           [,1]      [,2]
53 [1,] 4.011192 1.768793
54 [2,] 1.768793 4.323695
55
56 [[2]]$m2
57 [1] 0.1585909 -1.6907915
58
59 [[2]]$S2
60           [,1]      [,2]
61 [1,] 4.419030 2.990689
62 [2,] 2.990689 7.478789
63
64
65 [[3]]
66 [[3]]$lambda
67 [1] 0.5706126
68
69 [[3]]$m1
70 [1] 1.0064879 -0.7471182
71
72 [[3]]$S1
73           [,1]      [,2]
74 [1,] 5.043982 3.331804
75 [2,] 3.331804 7.074240
76
77 [[3]]$m2
78 [1] 1.4099485 -0.2522946
79
80 [[3]]$S2
81           [,1]      [,2]
82 [1,] 4.729293 2.799306
83 [2,] 2.799306 5.888067
84
85
86 [[4]]
87 [[4]]$lambda
88 [1] 0.5730394
89
90 [[4]]$m1
91 [1] 1.2172759 -0.4850478
92
93 [[4]]$S1
94           [,1]      [,2]
95 [1,] 4.904907 3.090795

```

```

96 [2,] 3.090795 6.486403
97
98 [[4]]$m2
99 [1] 1.1293355 -0.6012163
100
101 [[4]]$S2
102      [,1]      [,2]
103 [1,] 5.003782 3.230759
104 [2,] 3.230759 6.802684
105
106
107 [[5]]
108 [[5]]$lambda
109 [1] 0.5722337
110
111 [[5]]$m1
112 [1] 1.171511 -0.545906
113
114 [[5]]$S1
115      [,1]      [,2]
116 [1,] 4.960235 3.167631
117 [2,] 3.167631 6.657321
118
119 [[5]]$m2
120 [1] 1.1907214 -0.5195859
121
122 [[5]]$S2
123      [,1]      [,2]
124 [1,] 4.933730 3.133374
125 [2,] 3.133374 6.580629
126
127
128 [[6]]
129 [[6]]$lambda
130 [1] 0.5724354
131
132 [[6]]$m1
133 [1] 1.1815745 -0.5320576
134
135 [[6]]$S1
136      [,1]      [,2]
137 [1,] 4.946146 3.149607
138 [2,] 3.149607 6.616890
139
140 [[6]]$m2
141 [1] 1.1772577 -0.5381141
142
143 [[6]]$S2

```

```

144         [,1]      [,2]
145 [1,]  4.952774  3.157768
146 [2,]  3.157768  6.635089
147
148
149 [[7]]
150 [[7]]$lambda
151 [1]  0.572386
152
153 [[7]]$m1
154 [1]  1.1793096 -0.5352435
155
156 [[7]]$S1
157         [,1]      [,2]
158 [1,]  4.949668  3.153922
159 [2,]  3.153922  6.626501
160
161 [[7]]$m2
162 [1]  1.180290 -0.533849
163
164 [[7]]$S2
165         [,1]      [,2]
166 [1,]  4.948069  3.152006
167 [2,]  3.152006  6.622240
168
169
170 [[8]]
171 [[8]]$lambda
172 [1]  0.5723977
173
174 [[8]]$m1
175 [1]  1.1798248 -0.5345096
176
177 [[8]]$S1
178         [,1]      [,2]
179 [1,]  4.948823  3.152912
180 [2,]  3.152912  6.624255
181
182 [[8]]$m2
183 [1]  1.1796004 -0.5348313
184
185 [[8]]$S2
186         [,1]      [,2]
187 [1,]  4.949201  3.153359
188 [2,]  3.153359  6.625247
189
190
191 [[9]]

```

```

192 [[9]]$lambda
193 [1] 0.572395
194
195 [[9]]$m1
196 [1] 1.1797068 -0.5346789
197
198 [[9]]$S1
199           [,1]      [,2]
200 [1,] 4.949022 3.153147
201 [2,] 3.153147 6.624778
202
203 [[9]]$m2
204 [1] 1.1797583 -0.5346047
205
206 [[9]]$S2
207           [,1]      [,2]
208 [1,] 4.948934 3.153044
209 [2,] 3.153044 6.624548
210
211
212 [[10]]
213 [[10]]$lambda
214 [1] 0.5723956
215
216 [[10]]$m1
217 [1] 1.1797339 -0.5346398
218
219 [[10]]$S1
220           [,1]      [,2]
221 [1,] 4.948976 3.153093
222 [2,] 3.153093 6.624657
223
224 [[10]]$m2
225 [1] 1.179722 -0.534657
226
227 [[10]]$S2
228           [,1]      [,2]
229 [1,] 4.948996 3.153117
230 [2,] 3.153117 6.624710
231
232
233 [[11]]
234 [[11]]$lambda
235 [1] 0.5723955
236
237 [[11]]$m1
238 [1] 1.1797276 -0.5346489
239

```

```

240 [[11]]$S1
241           [,1]      [,2]
242 [1,]  4.948986  3.153105
243 [2,]  3.153105  6.624685
244
245 [[11]]$m2
246 [1]  1.1797304 -0.5346449
247
248 [[11]]$S2
249           [,1]      [,2]
250 [1,]  4.948982  3.153100
251 [2,]  3.153100  6.624672
252
253
254 [[12]]
255 [[12]]$lambda
256 [1]  0.5723955
257
258 [[12]]$m1
259 [1]  1.1797291 -0.5346468
260
261 [[12]]$S1
262           [,1]      [,2]
263 [1,]  4.948984  3.153103
264 [2,]  3.153103  6.624678
265
266 [[12]]$m2
267 [1]  1.1797285 -0.5346477
268
269 [[12]]$S2
270           [,1]      [,2]
271 [1,]  4.948985  3.153104
272 [2,]  3.153104  6.624681
273
274
275 [[13]]
276 [[13]]$lambda
277 [1]  0.5723955
278
279 [[13]]$m1
280 [1]  1.1797288 -0.5346473
281
282 [[13]]$S1
283           [,1]      [,2]
284 [1,]  4.948984  3.153103
285 [2,]  3.153103  6.624680
286
287 [[13]]$m2

```

```

288 [1] 1.179729 -0.534647
289
290 [[13]]$S2
291      [,1]      [,2]
292 [1,] 4.948984 3.153103
293 [2,] 3.153103 6.624679
294
295
296 [[14]]
297 [[14]]$lambda
298 [1] 0.5723955
299
300 [[14]]$m1
301 [1] 1.1797288 -0.5346471
302
303 [[14]]$S1
304      [,1]      [,2]
305 [1,] 4.948984 3.153103
306 [2,] 3.153103 6.624679
307
308 [[14]]$m2
309 [1] 1.1797288 -0.5346472
310
311 [[14]]$S2
312      [,1]      [,2]
313 [1,] 4.948984 3.153103
314 [2,] 3.153103 6.624680
315
316 # -----
317 > length(results$history)
318 [1] 14

```

Παρατηρούμε πως η μέθοδος κάνει τις παραμέτρους να συγκλίνουν προς τις τιμές του ερωτήματος α σαν να έχουμε ένα μοντέλο διδιάστατης κανονικής κατανομής. Μια άλλη εναλλακτική είναι στη συνάρτηση `dmnorm` να μην επιλέξουμε `log = TRUE` και έτσι λαμβάνουμε τα παρακάτω τελικά αποτελέσματα.

```

1 > results$final
2 $lambda
3 [1] 4.072829e-08
4
5 $m1
6 [1] 0.4746379 -1.6557931
7
8 $S1
9      [,1]      [,2]
10 [1,] 6.244159 4.538616
11 [2,] 4.538616 10.394352
12

```

```

13 $m2
14 [1] 1.1797288 -0.5346471
15
16 $S2
17      [,1]      [,2]
18 [1,] 4.948984 3.153103
19 [2,] 3.153103 6.624679

```

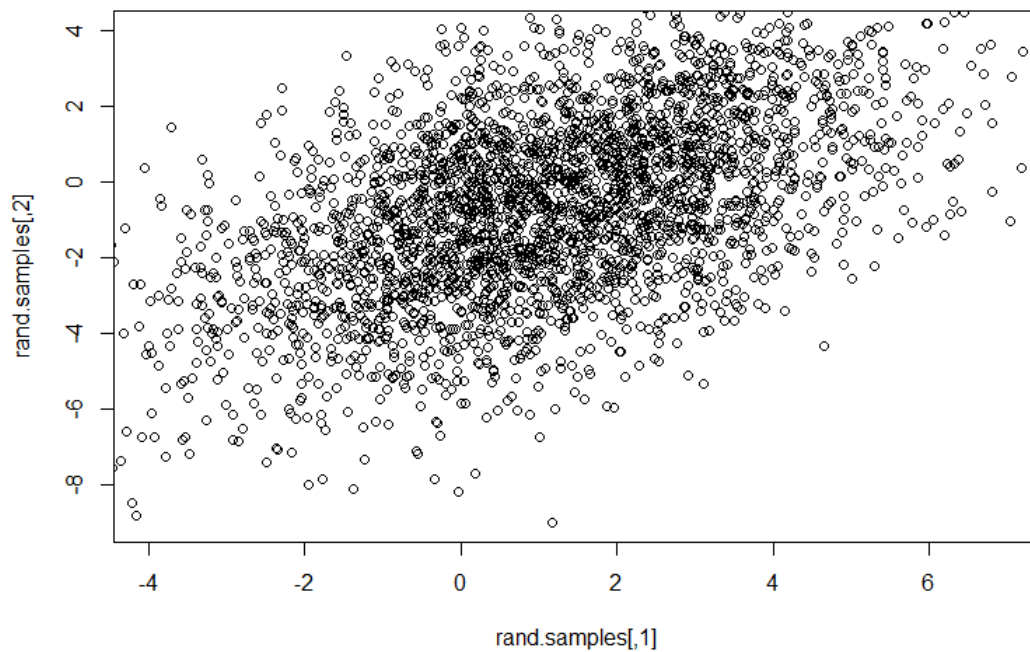
Εδώ παρατηρούμε πως διαφοροποιούνται οι παράμετροι των κανονικών κατανομών αλλά η παράμετρος μείξης μηδενίζεται και πάλι προσομοιώνει μια κατανομή αντί για μείξη.

δ. Για να προσομοιώσουμε σημεία από μείξη κανονικών κατανομών θα δημιουργήσουμε 3000 τυχαίους αριθμούς μεταξύ 0 και 1 και όταν είναι μικρότερος του λ θα προσομοιώνει σημείο από την πρώτη κατανομή και αλλιώς από τη δεύτερη. Τελικά θα έχουμε 3000 νέα σημεία και θα τα αναπαραστήσουμε γραφικά.

```

1 N=3000
2 U =runif(N)
3 rand.samples = matrix(NA,nrow = N,ncol=2)
4
5 for(i in 1:N){
6   if(U[i]<0.5){
7     rand.samples[i,] = mvrnorm(1,results$final[[2]],results$
      ↪ final[[3]])#results$final[[3]]
8   }else{
9     rand.samples[i,] = mvrnorm(1,results$final[[4]],results$
      ↪ final[[5]])#results$final[[5]]
10  }
11 }
12 plot(rand.samples,xlim = c(-4,7),ylim = c(-9,4))

```

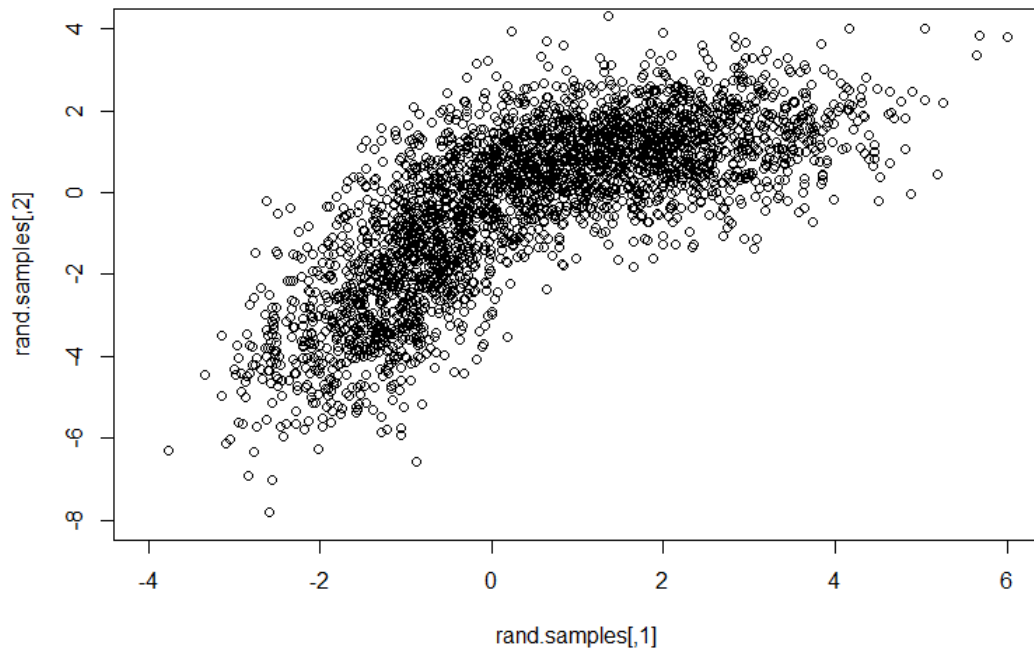


Η εικόνα που δημιουργείται δεν είναι ικανοποιητική για τα δεδομένα και μπορούμε να δούμε πως για τις αρχικές τιμές που εκτιμήσαμε η εικόνα είναι πολύ πιο κοντά στα δεδομένα.

```

1 N=3000
2 U =runif(N)
3 rand.samples = matrix(NA,nrow = N,ncol=2)
4 l1<-c(-1,-2)
5 l2<-c(1.5,1)
6 sigma1<-matrix(c(0.7,1,1,3),nrow = 2)
7 sigma2<-matrix(c(1.8,0.5,0.5,1),nrow = 2)
8 for(i in 1:N){
9   if(U[i]<0.4){
10     rand.samples[i,] = mvrnorm(1,l1,sigma1)
11   }else{
12     rand.samples[i,] = mvrnorm(1,l2,sigma2)
13   }
14 }
15 plot(rand.samples,xlim = c(-4,6),ylim = c(-8,4))

```

4 Άσκηση 4

Στην βιβλιοθήκη `faraway` θα βρούμε τα δεδομένα `fat` τα οποία θα χρησιμοποιήσουμε για αυτήν την άσκηση. Τα δεδομένα περιέχουν πληροφορίες για 252 άνδρες, σχετικά με το ποσοστό σωματικού λίπους χρησιμοποιώντας την εξίσωση Brozek και 17 ακόμη επεξηγηματικές μεταβλητές. Αρχικά θα φορτώσουμε τα δεδομένα μας και θα αφαιρέσουμε τις μεταβλητές `siri`, `density`, `free` μιας και δεν θα χρησιμοποιηθούν για την άσκηση.

```
1 library(faraway) #load the library
2 data<-faraway::fat #store our data
3 data<-data[,c(-2,-3,-8)]#remove unused columns
4
5 for (i in 2:ncol(data)){ #normalising
6   data[,i]<-data[,i]-mean(data[,i]) #sum = 0
7   data[,i]<-data[,i]/sqrt(sum((data[,i]-mean(data[,i]))^2)) #
      ↳ sum of squares = 1
8 }
```

Αφού φορτώσαμε τα δεδομένα στη συνέχεια, μόνο για τις τιμές των επεξηγηματικών μεταβλητών πραγματοποιήθηκε τυποποίηση. Συγκεκριμένα, μετά την αλλαγή κάθε μεταβλητή έχει άθροισμα ίσο με 0 και άθροισμα τετραγώνων ίσο με 1. Αυτό το καταφέραμε αφαιρώντας από κάθε παρατήρηση τη μέση τιμή της στήλης που ανήκει και διαιρώντας με την τετραγωνική ρίζα του αθροίσματος τετραγώνων της στήλης.

α. Αρχικά θέλουμε να χωρίσουμε το δείγμα μας σε training data και test data με μεγέθη 4/5 του δείγματος και 1/5 αντίστοιχα.

```
1 set.seed(0) #for reproducibility
2 deigma<-sample(1:nrow(data),round(4/5*nrow(data)),replace =
      ↳ FALSE)
3 training_sample<-data[deigma,]
4 test_sample<-data[-deigma,]
5
6 X<-training_sample[,-1]
7 Y<-training_sample$brozek
8 x_test<-test_sample[,-1]
9 y_test<-test_sample[,1]
```

Χρησιμοποιούμε τη συνάρτηση `sample()` της R η οποία θα μας δώσει τα indices των παρατηρήσεων που θα μπουν στο training sample. Η συνάρτηση δέχεται ως ορίσματα το εύρος των αριθμών από τους οποίους θα επιλέξει, που σε μας είναι το πλήθος των γραμμών, έπειτα πόσους αριθμούς θα επιστρέψει και ζητάμε τα 4/5 των παρατηρήσεων στρογγυλοποιημένα καθώς μιλάμε για διακριτές τιμές. Επίσης θα επιλέξουμε `replace = FALSE` για να μην επιλαγεί μια παρατήρηση πολλαπλές φορές. Στη συνέχεια ορίζουμε τα δύο δείγματα παίρνοντας μόνο τις παρατηρήσεις που μας έδωσε η τυχαία συνάρτηση `sample()`.

β. Η μέθοδος Lasso είναι μια μέθοδος συρρίκνωσης που χρησιμοποιείται για να μειώσει την επίδραση συγκεκριμένων εμταβλητών ενός μοντέλου παλινδρόμησης, ενώ ταυτόχρονα πραγματοποιεί και έμμεση επιλογή μεταβλητών μοντέλου. Ο τρόπος που λειτουργεί είναι πως

ποινικοποιείτους συντελεστές β_j , μέσω της l_1 νόρμας του διανύσματος των συντελεστών. Οι συντελεστές του μοντέλου Lasso βρίσκονται ελαχιστοποιώντας την παρακάτω σχέση:

$$\min(y - Z\beta)^T(y - Z\beta), \quad \sum_{i=1}^p |\beta_j| \leq t \quad (4.1)$$

Όπου y το διάνυσμα μεταβλητής απόκρισης, Z ο πίνακας των τυποποιημένων επεξηγηματικών μεταβλητών, β το διάνυσμα των συντελεστών του μοντέλου και t είναι η παράμετρος ρύθμισης (regularization parameter). Η σχέση ισοδύναμα γράφεται και ως εξής:

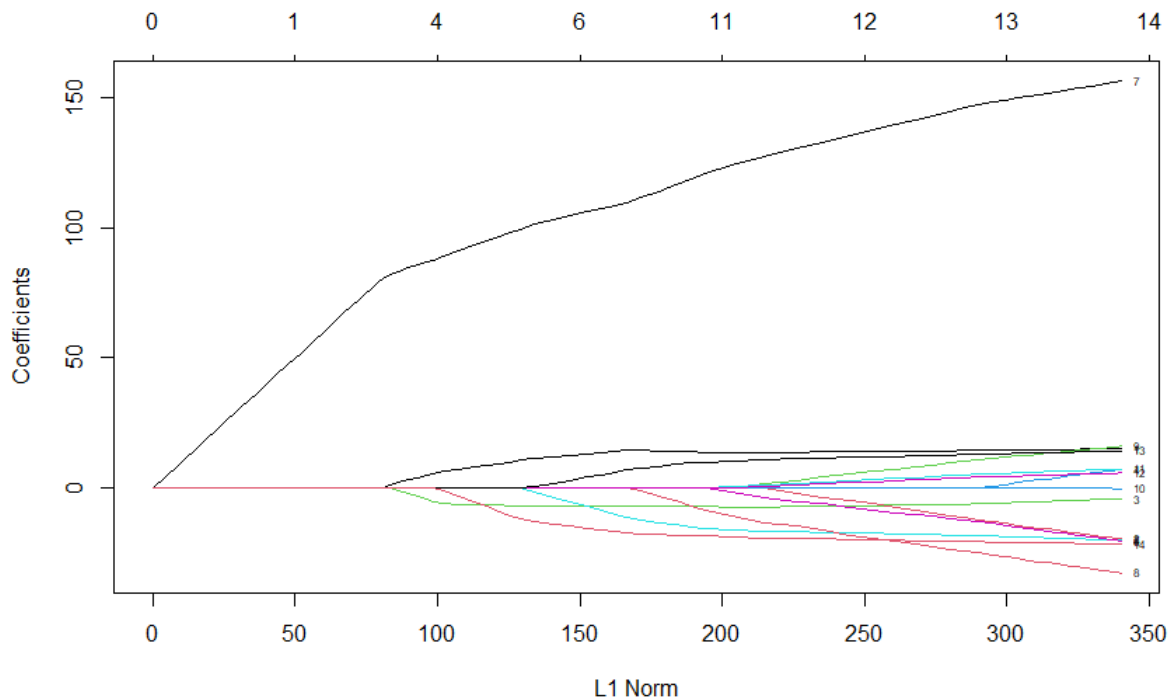
$$\min \left\{ (y - Z\beta)^T(y - Z\beta) + \lambda \sum_{i=1}^p |\beta_j| \right\} \quad (4.2)$$

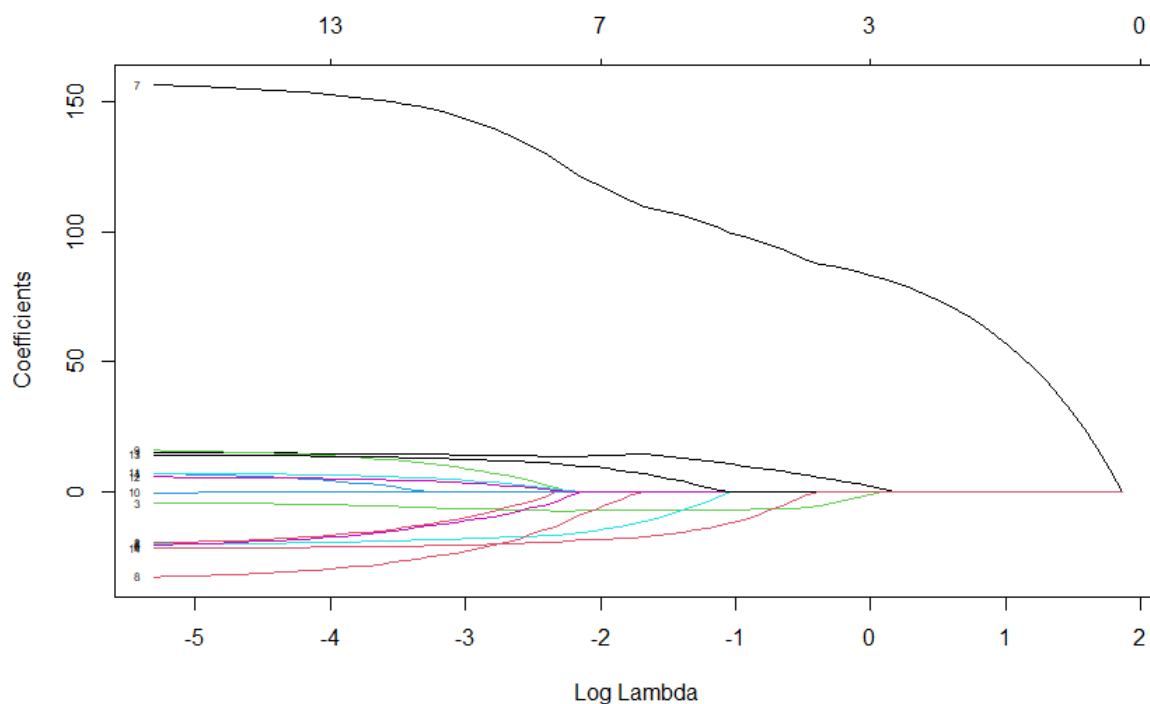
με λ πάλι η παράμετρος σύγκρισης σε άλλη μορφή.

Για να εφαρμόσουμε την μέθοδο θα χρησιμοποιήσουμε τη βιβλιοθήκη `glmnet` της R στο training sample.

```
library(glmnet) #load the library
lasso1<-glmnet(X,Y)
plot(lasso1,label = T)
plot(lasso1,xvar = 'lambda',label = T) #plotting with log scaling
```

Αρχικά αφού καλέσουμε τη βιβλιοθήκη κάνουμε χρήση της ομόνυμης εντολής `glmnet()` με ορίσματα τις επεξηγηματικές μεταβλητές και την μεταβλητή απόκρισης. Στη συνέχεια αναπαριστούμε γραφικά το μοντέλο συναρτήσει της l_1 νόρμας και τον λογάριθμο του λ όπως φαίνεται στη συνέχεια.





Τα συμπεράσματα για τα δύο γραφήματα είναι ίδια, αφού βλέπουμε το σταδιακό μηδενισμό των συντελεστών των μεταβλητών καθώς το t μειώνεται στην πρώτη γραφική και καθώς ο λογάριθμος του λ αυξάνεται στη δεύτερη. Οι πιο σημαντικές μεταβλητές στο μοντέλο είναι αυτές που μηδενίζονται τελευταίες, όπως οι 7,13 και 3.

γ. Στη συνέχεια θέλουμε να επιλέξουμε, χρησιμοποιώντας cross-validation, την τιμή της παραμέτρου λ που ελαχιστοποιεί το CV-MSE. Θα χρησιμοποιήσουμε την εντολή `cv.glmnet()` και θα καταλήξουμε σε ένα μοντέλο με λιγότερες επεξηγηματικές μεταβλητές, που θα το καλέσουμε M1. Έπειτα με παρόμοια διαδικασία θα βρούμε την τιμή της λ που ελαχιστοποιεί το CV-MSE με σφάλμα εντός μιας τυπικής απόκλισης από την ελάχιστη τιμή, αυτό το μοντέλο θα το καλέσουμε M2. Σκοπός είναι να συγκρίνουμε τα δύο αυτά μοντέλα στα test data και να επιλέξουμε αυτό με το μικρότερο MSE.

```

1 X_matrix<-as.matrix(X)
2 lasso2<-cv.glmnet(X_matrix,Y)
3 M1<-coef(lasso2,s="lambda.min") #find the coefficients of the
  ↳ model
4 M2<-coef(lasso2,s="lambda.1se")
5 M1.ind<-c(1,3,5,7,8,13,14) #store the indices to fit the
  ↳ model
6 M2.ind<-c(1,3,7,14)

```

Αρχικά μετατρέπω τις επεξηγηματικές μου μεταβλητές σε μορφή πίνακα και καλώ την εντολή `cv.glmnet()`. Έπειτα βρίσκω τους συντελεστές του μοντέλου για κάθε περίπτωση με την εντολή `coef()` και αποθηκεύω σε μορφή διανύσματος τους δείκτες των μη μηδενικών μεταβλητών. Αυτό το κάνω για να μπορέσω να προσαρμόσω στη συνέχεια τα μοντέλα M1 και M2. Παρακάτω βλέπουμε τους συντελεστές των δύο μοντέλων.

```

1 > M1
2 15 x 1 sparse Matrix of class "dgCMatrix"
3           s1
4 (Intercept) 18.809049
5 age         14.079613
6 weight      .
7 height      -7.050080
8 adipos      .
9 neck        -14.019852
10 chest      .
11 abdom      116.538211
12 hip        -4.974351
13 thigh      .
14 knee       .
15 ankle      .
16 biceps     .
17 forearm    9.003045
18 wrist     -18.123980
19
20 > M2
21 15 x 1 sparse Matrix of class "dgCMatrix"
22           s1
23 (Intercept) 18.840156
24 age         9.450551
25 weight      .
26 height      -6.630711
27 adipos      .
28 neck        .
29 chest      .
30 abdom      96.501029
31 hip         .
32 thigh      .
33 knee       .
34 ankle      .
35 biceps     .
36 forearm    .
37 wrist     -8.915406

```

Έπειτα προσαρμόζουμε τα δύο μοντέλα παλινδρόμησης βάσει των μεταβλητών που δεν μηδενίστηκαν και υπολογίζουμε τα MSE.

```

1 model1<-lm(Y~.,data = X[,M1.ind]) #fit the models
2 model2<-lm(Y~.,data = X[,M2.ind])
3 MSE.M1<-sum((predict(model1,x_test[,M1.ind])-y_test)^2)/
  ↪ length(y_test)
4 MSE.M2<-sum((predict(model2,x_test[,M2.ind])-y_test)^2)/
  ↪ length(y_test)
5 M_Lasso<-model2

```

Το MSE υπολογίζεται ως η διαφορά των προβλεπόμενων τιμών με τις πραγματικές στο τετράγωνο, προς το πλήθος των παρατηρήσεων, για το test sample. Για τις προβλεπόμενες τιμές χρησιμοποιούμε την εντολή predict().

```
1 > MSE.M1
2 [1] 15.76858
3 > MSE.M2
4 [1] 15.43327
```

Τελικά επιλέγουμε το M2 ως M_Lasso αφού έχει μικρότερο MSE.

δ. Χρησιμοποιώντας όλα τα δεδομένα θα προσαρμόσουμε το γραμμικό μοντέλο με επεξηγηματικές μεταβλητές ίδιες με αυτές του M_Lasso.

```
1 x<-data[, -1]
2 y<-data[, 1]
3 M3<-lm(y~., data=x[, M2.ind])
4
5 > M3
6 Call:
7 lm(formula = y ~ ., data = x[, M2.ind])
8
9 Coefficients:
10 (Intercept)          age          height          abdom
11      ↪ wrist
18.938          13.391          -9.552          115.433
      ↪ -28.651
```

ε. Ψάχνουμε στον χώρο όλων των πιθανών μοντέλων, εκείνο το μοντέλο το οποίο ελαχιστοποιεί την τιμή του κριτηρίου BIC. Θέλουμε να προσαρμόσουμε δηλαδή μοντέλα με κάθε δυνατό συνδυασμό μεταβλητών και να υπολογίσουμε το BIC για κάθε ένα από αυτά, επιλέγοντας τελικά το μοντέλο που δίνει την μικρότερη τιμή. Αρχικά θα αποθηκεύσω όλες τις μεταβλητές ξεχωριστά καθώς θα μου χρειαστεί στην αναζήτηση μοντέλων.

```
1 brozek<-data$brozek
2 age<-x$age
3 weight<-x$weight
4 height<-x$height
5 adipos<-x$adipos
6 neck<-x$neck
7 chest<-x$chest
8 abdom<-x$abdom
9 hip<-x$hip
10 thigh<-x$thigh
11 knee<-x$knee
12 ankle<-x$ankle
13 biceps<-x$biceps
14 forearm<-x$forearm
15 wrist<-x$wrist
16
```

```

17 predictors<-colnames(x)
18 response<-"brozek"
19 best_bic <- Inf
20 best_model <- NULL
21
22 # Loop over all subset sizes
23 for (k in 1:length(predictors)) {
24   # Loop over all subsets of size k
25   subsets <- combn(predictors, k, simplify = FALSE)
26   for (subset in subsets) {
27     # Construct the formula for this subset
28     predictors_str <- paste(subset, collapse = "+")
29     formula_str <- paste(response, "~", predictors_str)
30     formula <- as.formula(formula_str)
31     # Fit the model and compute its BIC
32     model <- lm(formula, data = x)
33     bic <- BIC(model)
34
35
36     # If this is the best model so far, remember it
37     if (bic < best_bic) {
38       best_bic <- bic
39       M4 <- model
40     }
41   }
42 }
43
44 summary(M4)
45 print(best_bic)

```

Έπειτα δημιουργώ ένα διάνυσμα που περιέχει τα ονόματα των επεξηγηματικών μεταβλητών, το ίδιο για την μεταβλητή απόκρισης και ορίζω αρχική τιμή για τι BIC το άπειρο. Το βήμα αυτό όπως και το `best_model<- NULL` είναι προαιρετικό και βοηθάει στην αρχικοποίηση των δεδομένων. Έπειτα κάνουμε `iterate` σε κάθε πιθανό μέγεθος μοντέλου, πέρα από το μηδενικό, και σε κάθε στάδιο του βρόγχου με χρήση της εντολής `combn` παίρνουμε όλους τους δυνατούς συνδυασμούς επεξηγηματικών μεταβλητών εκείνου του μεγέθους. Στη συνέχεια φτιάχνουμε τη φόρμουλα του μοντέλου ενώνοντας τα string πρώτα της μεταβλητής απόκρισης, έπειτα την `~` και τέλος τον συγκεκριμένο συνδυασμό επεξηγηματικών μεταβλητών που τις χωρίζουμε μεταξύ τους με `+`. Έτσι έχουμε καταφέρει να έχουμε για συγκεκριμένο μέγεθος όλα τα πιθανά μοντέλα προσαρμοσμένα και αποθηκεύουμε σε ένα διάνυσμα το κάθε BIC. Αυτό πραγματοποιείται για μεγέθη μοντέλων από το 1 μέχρι το πλήρες μοντέλο και άρα έχουμε κάνει αναζήτηση στον χώρο όλων των πιθανών μοντέλων. Τέλος με ένα απλό `if statement` έχουμε την μικρότερη τιμή του κριτηρίου BIC καθώς και το μοντέλο που το αποκαλούμε `M4`.

```

1 > summary(M4)
2 Call:
3 lm(formula = formula, data = x)
4
5 Residuals:
6      Min       1Q   Median       3Q      Max
7 -9.8002 -2.8728 -0.1545  2.8980  8.3845
8
9 Coefficients:
10             Estimate Std. Error t value Pr(>|t|)
11 (Intercept)   18.9385     0.2533   74.761 < 2e-16 ***
12 weight       -58.4648    10.6704   -5.479 1.05e-07 ***
13 abdom        157.4037     8.8692   17.747 < 2e-16 ***
14 forearm      14.2904     5.3853    2.654 0.008480 **
15 wrist        -20.5853     6.0628   -3.395 0.000799 ***
16 ---
17 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 1
18
19 Residual standard error: 4.021 on 247 degrees of freedom
20 Multiple R-squared:  0.7351,    Adjusted R-squared:  0.7308
21 F-statistic: 171.4 on 4 and 247 DF,  p-value: < 2.2e-16
22
23 > print(best_bic)
24 [1] 1444.647

```

Το μοντέλο M4 περιέχει για επεξηγηματικές μεταβλητές τις weight, abdom, forearm, wrist και είναι διαφορετικό από το M3.

στ. Χρησιμοποιώντας την μέθοδο 5-fold cross-validation και την (within fold) ARMSE (Average Root Mean Square Error) συνάρτηση, θα εξεταστεί ποιο από τα μοντέλα M3 και M4 έχει την καλύτερη προβλεπτική ικανότητα. Αρχικά θα χωρίσω το δείγμα σε 5 folds περίπου ίδιου μεγέθους και θα χρησιμοποιήσω τα τέσσερα από αυτά για την προσαρμογή του μοντέλου ενώ το τελευταίο για πρόβλεψη και για τον υπολογισμό του MSE. Θα το κάνω αυτό για κάθε fold και θα υπολογίσω την μέση τιμή των MSE για να προβώ σε σύγκριση των μοντέλων. Έστω T_k ένα από τα folds, n_k το μέγεθος του, y_i είναι οι παρατηρήσεις της μεταβλητής απόκρισης του fold αυτού και $y_{i,-k}$ είναι η πρόβλεψη του y_i με όλα τα δεδομένα εκτός του T_k . Τότε έχουμε:

$$MSE(T_k) = \frac{1}{n_k} \sum_{i \in T_k} (y_i - y_{i,-k})^2 \quad (4.3)$$

$$ARMSE = \frac{1}{5} \sqrt{\sum_{i=1}^5 MSE(T_k)} \quad (4.4)$$

Δίνεται ο κώδικας στην R.

```
1 n<-length(y)
2 set.seed(155)
3 randomise<-sample(1:n,n)
4 F1<-randomise[1:50]
5 F2<-randomise[51:100]
6 F3<-randomise[101:150]
7 F4<-randomise[151:200]
8 F5<-randomise[201:252]
9 folds<-list(F1,F2,F3,F4,F5)
10 M3.ind<-c(1,3,7,14)
11 M4.ind<-c(2,7,13,14)
12
13 MSE.M3<-rep(0,5)
14 MSE.M4<-rep(0,5)
15
16 for(i in 1:5){
17   fold<-unlist(folds[i])
18   xnew_train<-x[-fold,]
19   ynew_train<-y[-fold]
20   xnew_test<-x[fold,]
21   ynew_test<-y[fold]
22
23   M3.mod<-lm(ynew_train~.,data = xnew_train[,M3.ind])
24   M4.mod<-lm(ynew_train~.,data = xnew_train[,M4.ind])
25
26   MSE.M3[i]<-sum((predict(M3.mod,xnew_test[,M3.ind])-ynew_
27     ↪ test)^2)/length(ynew_test)
28   MSE.M4[i]<-sum((predict(M4.mod,xnew_test[,M4.ind])-ynew_
29     ↪ test)^2)/length(ynew_test)
30 }
31
32 ARMSE.M3<-sqrt(sum(MSE.M3))/5
33 ARMSE.M4<-sqrt(sum(MSE.M4))/5
```

Αρχικά για να δημιουργήσουμε τα folds θα χρειαστούμε τη συνάρτηση `sample()` όπου θα επιλέξει από τις παρατηρήσεις μας ένα τυχαίο δείγμα ίδιου μεγέθους και άρα απλά έχει ανακατέψει τις παρατηρήσεις. Στη συνέχεια χωρίζουμε το τυχαίοποιημένο δείγμα σε 5 μέρη και τα τοποθετούμε σε μια λίστα. Θα χρειαστούμε και τους δείκτες των επεξηγηματικών μεταβλητών των δύο μοντέλων και αρχικοποιούμε τα διανύσματα που θα αποθηκεύσουν τις τιμές των MSE. Στη συνέχεια για κάθε i επιλέγουμε το i -οστό fold και το αφαιρούμε από τις παρατηρήσεις για την προσαρμογή του μοντέλου, ενώ θα το χρειαστούμε για προβλέψεις. Προσαρμόζουμε τα δύο μοντέλα στα δεδομένα και στη συνέχεια υπολογίζουμε το MSE για το συγκεκριμένο fold. Τέλος αφού έχουμε υπολογίσει 5 τιμές για κάθε MSE θα υπολογίσουμε τις μέσες τιμές τους. Το M3 έχει $ARMSE = 1.877849$ ενώ το M4 έχει $ARMSE = 1.814728$. Παρατηρούμε πως το μοντέλο M4 έχει μικρότερη τιμή για το ARMSE και άρα έχει την καλύτερη προβλεπτική ικανότητα.

ζ. Χρησιμοποιώντας το M4 ως τελικό μοντέλο σε όλα τα δεδομένα και 1000 Bootstrap δείγματα από τα υπόλοιπα θέλουμε να δώσουμε εκτιμήτριες και τυπικά σφάλματα για τους συντελεστές του μοντέλου. Θα χρησιμοποιήσουμε μη παραμετρικό Bootstrap καθώς δεν γνωρίζουμε τη μέση τιμή και τη διασπορά των υπολοίπων. Συγκεκριμένα από το δείγμα των υπολοίπων θα δημιουργήσουμε 1000 νέα δείγματα ίδιου μεγέθους κάνοντας δειγματοληψία με επανατοποθέτηση. Για το καινούριο σύνολο δεδομένων θα ισχύει:

$$\bar{\hat{\theta}}^* \equiv \frac{1}{B} \sum_{i=1}^B \hat{\theta}^* \rightarrow E[\hat{\theta}] \quad (4.5)$$

$$\sqrt{\frac{1}{B-1} \sum_{i=1}^B (\hat{\theta}_i^* - \bar{\hat{\theta}}^*)^2} \rightarrow se(\hat{\theta}) \quad (4.6)$$

Όπου B είναι το πλήθος των δειγμάτων. Οι πράξεις αυτές θα γίνουν με την βοήθεια της βιβλιοθήκης boot της R.

```

1 M_final<-M4
2 library(boot)
3 res<-M_final$residuals
4
5 bootstrap_func<-function(data,indices){
6
7     bootstrap_res<-res[indices]
8
9     bootstrap_response <- brozek + bootstrap_res
10    bootstrap_model <- lm(bootstrap_response ~., data = x[M4.
      ↪ ind])
11    return(coef(bootstrap_model))
12 }
13
14 results1<-boot(data=res,statistic = bootstrap_func , R=1000)
15
16 bootstrap_coefficients <- colMeans(results1$t) #estimates of
      ↪ the coefficients
17 print(bootstrap_coefficients)
18 print(M4$coefficients)
19
20 bootstrap_se <- apply(results1$t, 2, sd) #estimates of
      ↪ standard error
21 print(bootstrap_se)
22 print(summary(M4)$coefficients[, "Std. Error"])
```

Αρχικά αποθηκεύουμε σε ένα διάνυσμα τα υπόλοιπα του τελικού μοντέλου, στη συνέχεια θα φτιάξουμε τη συνάρτηση που θα χρησιμοποιήσει η εντολή boot. Η συνάρτηση δέχεται τα αρχικά δεδομένα και τους δείκτες που παράγει η εντολή και αφού επιλέξει το νέο δείγμα υπολοίπων το προσθέτει στην μεταβλητή απόκρισης και προσαρμόζει καινούριο μοντέλο παλινδρόμησης με τις ίδιες επεξηγηματικές μεταβλητές και αποθηκεύει τους συντελεστές του μοντέλου. Έπειτα καλούμε τη συνάρτηση boot με δεδομένα τα υπόλοιπα, συνάρτηση

αυτή που δημιουργήσαμε και πλήθος δειγμάτων 1000. Για να συγκρίνουμε τα αποτελέσματα θα τυπώσουμε τους εκτιμητές καθώς και τους πραγματικούς συντελεστές του μοντέλου και μετά τις εκτιμήσεις των τυπικών σφαλμάτων μαζί με τις πραγματικές τιμές του αρχικού μοντέλου.

```
1 > print(bootstrap_coefficients)
2 [1] 18.94283 -57.80087 157.04111 14.28471 -20.82949
3 > print(M4$coefficients)
4 (Intercept)      weight      abdom      forearm      wrist
5 18.93849 -58.46475 157.40365 14.29038 -20.58529
6
7 > print(bootstrap_se)
8 [1] 0.2464339 10.3167843 8.3938198 5.1239055 6.0259312
9 > print(summary(M4)$coefficients[, "Std. Error"])
10 (Intercept)      weight      abdom      forearm      wrist
11 0.2533216 10.6704481 8.8692358 5.3852752 6.0628484
```

Όπως βλέπουμε οι εκτιμήσεις είναι πολύ κοντά στις πραγματικές τιμές και για τους συντελεστές και για τα τυπικά σφάλματα. Σε περίπτωση που θέλουμε να έχουμε μικρότερο σφάλμα στις εκτιμήσεις μας μπορούμε να μεγαλώσουμε το R ή να εφαρμόσουμε κάποια μέθοδο ελάττωσης διασποράς στις εκτιμήσεις.