

Στατιστική στην Πληροφορική - Εργασία 1

Αυγεράκης Βασίλης - p3210013
Παντελίδης Ιπποκράτης - p3210150

Άσκηση 1)

Πριν προχωρήσουμε στην επίλυση των επιμέρους ερωτημάτων θα υπολογίσουμε την σύνοψη των 5 αριθμών, την μέση τιμή και την τυπική απόκλιση για κάθε ομάδα δεδομένων καθώς θα τις χρησιμοποιήσουμε στα ερωτήματα α) (boxplots) και β) (σύγκριση καταλληλότητας).

Δεδομένα I

5-num sum -> (min = 30.3, Q1 = 31.1, m = 32.65, Q3 = 33.6, max = 34.5)
mean -> 32.55
standard deviation -> 1.419898

Δεδομένα II

5-num sum -> (min = 0.0, Q1 = 0.2, m = 1.3, Q3 = 4.2, max = 9.0)
mean -> 2.64
standard deviation -> 3.059121

Δεδομένα III

5-num sum -> (min = 0.0, Q1 = 17.5, m = 39.5, Q3 = 59.0, max = 96.0)
mean -> 41.15
standard deviation -> 28.26754

Stemplots

Δεδομένα I

30	3
31	0 1
32	1 6 7
33	4 6
34	2 5

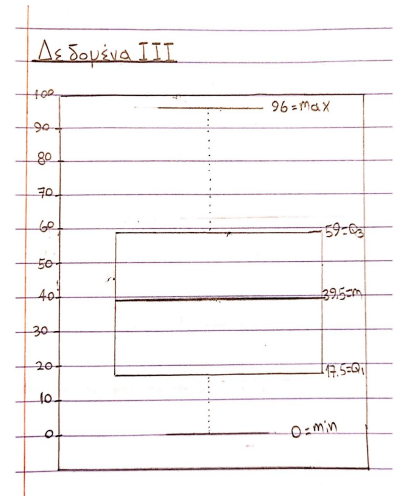
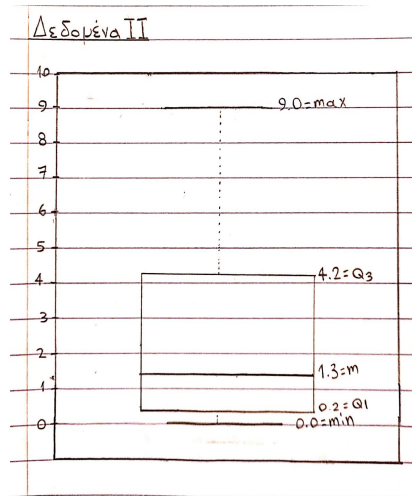
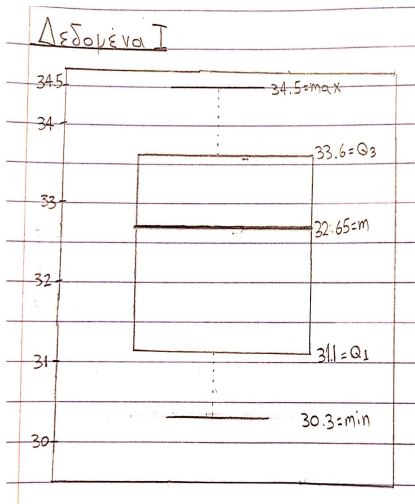
Δεδομένα II

0	0 0 2 8
1	2 4
2	
3	2
4	2
5	
6	4
7	
8	
9	0

Δεδομένα III

0	0 1 6 8
1	0 3 5 6 7 7 8 8
2	0 0 1 5 6
3	0 5 9
4	0 1 3 4 6 8
5	2 4 8 9 9
6	0 6
7	
8	1 6 7 8 9
9	4 6

Boxplots



b) Η μέση τιμή και η τυπική απόκλιση συνοψίζουν καλύτερα την κατανομή για την ομάδα δεδομένων I, όπου τα δεδομένα είναι συμμετρικά κατανεμημένα γύρω από το \bar{x} και δεν υπάρχουν ατυπικές τιμές. Αντίθετα, οι ομάδες δεδομένων II και III έχουν μη συμμετρικά δεδομένα με outliers, καθιστώντας την μέση τιμή και την τυπική απόκλιση λιγότερο αποτελεσματική για την περιγραφή τους. Επομένως, η σύνοψη των 5 αριθμών συνοψίζει καλύτερα την κατανομή των δύο τελευταίων ομάδων, καθώς αυτή δεν επηρεάζεται από τις ατυπικές τιμές.

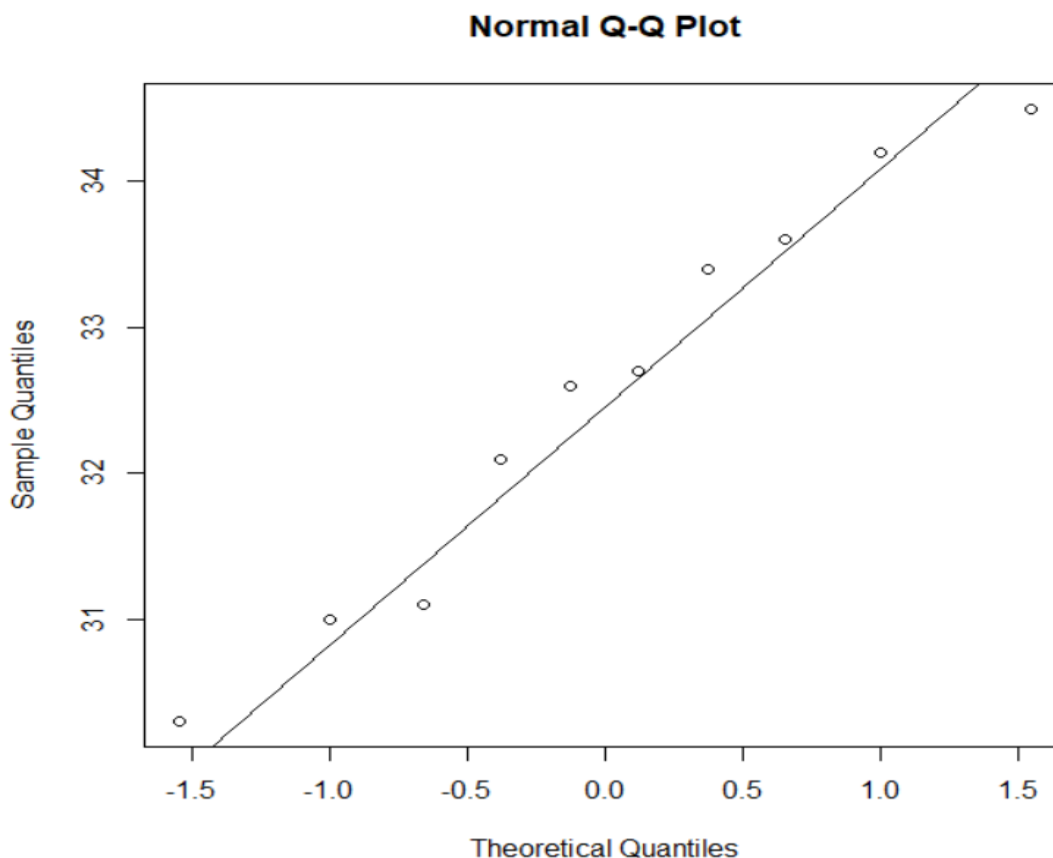
c) Για το ερώτημα αυτό θα χρησιμοποιήσουμε τον κανόνα 68-95-99.7 σύμφωνα με τον οποίο για κάθε κανονική κατανομή $N(\mu, \sigma)$, όπου μ = μέση τιμή = διάμεση τιμή και σ = τυπική απόκλιση ισχύει ότι :

- το 68% των περιπτώσεων βρίσκεται στο διάστημα $(\mu - \sigma, \mu + \sigma)$
- το 95% των περιπτώσεων βρίσκεται στο διάστημα $(\mu - 2\sigma, \mu + 2\sigma)$
- Το 99.7% των περιπτώσεων βρίσκεται στο διάστημα $(\mu - 3\sigma, \mu + 3\sigma)$

Μπορούμε επίσης να χρησιμοποιήσουμε το Normal-quantile plot, ένα οπτικό μέσο για να συγκρίνουμε την κατανομή των δεδομένων μας με την κανονική κατανομή. Το διάγραμμα αυτό μας λέει ότι όσο περισσότερο συγγραμμικά είναι τα σημεία, δηλαδή κοντά πάνω στην ευθεία, τόσο πιο πολύ ταυτίζεται η κατανομή μας (τα δεδομένα στο άξονα y) με την κανονική.

Δεδομένα I

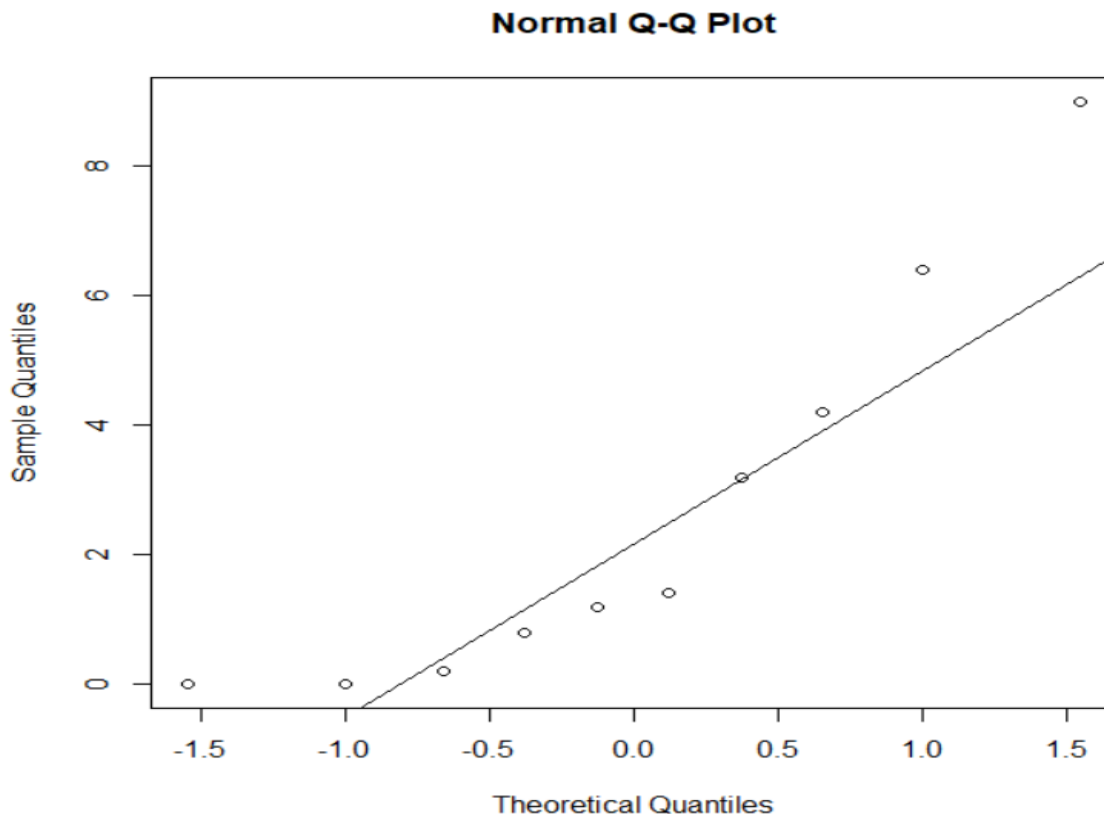
- Στο διάστημα $(\mu - \sigma, \mu + \sigma)$ δηλαδή το (31.1301, 33.9699) ανήκουν 5 από τις 10 τιμές της ομάδας, δηλαδή το 50%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 68%, οπότε έχουμε απόκλιση -18%.
- Στο διάστημα $(\mu - 2\sigma, \mu + 2\sigma)$ δηλαδή το (29.7102, 35.3898) ανήκουν και οι 10 από τις 10 τιμές της ομάδας, δηλαδή το 100%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 95% οπότε έχουμε απόκλιση +5%.
- Στο διάστημα $(\mu - 3\sigma, \mu + 3\sigma)$ δηλαδή το (28.29031, 36.80969) ανήκουν πάλι οι 10 από τις 10 τιμές της ομάδας, δηλαδή το 100%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 99,7% οπότε έχουμε απόκλιση +0,03%.



Παρατηρούμε ότι τα σημεία είναι αρκετά συγγραμμικά μεταξύ και σε συνδυασμό με τα σχετικά χαμηλά ποσοστά αποκλίσεων μπορούμε να συμπεράνουμε ότι η κατανομή των δεδομένων της ομάδας I προσεγγίζει αρκετά καλά την κανονική κατανομή.

Δεδομένα II

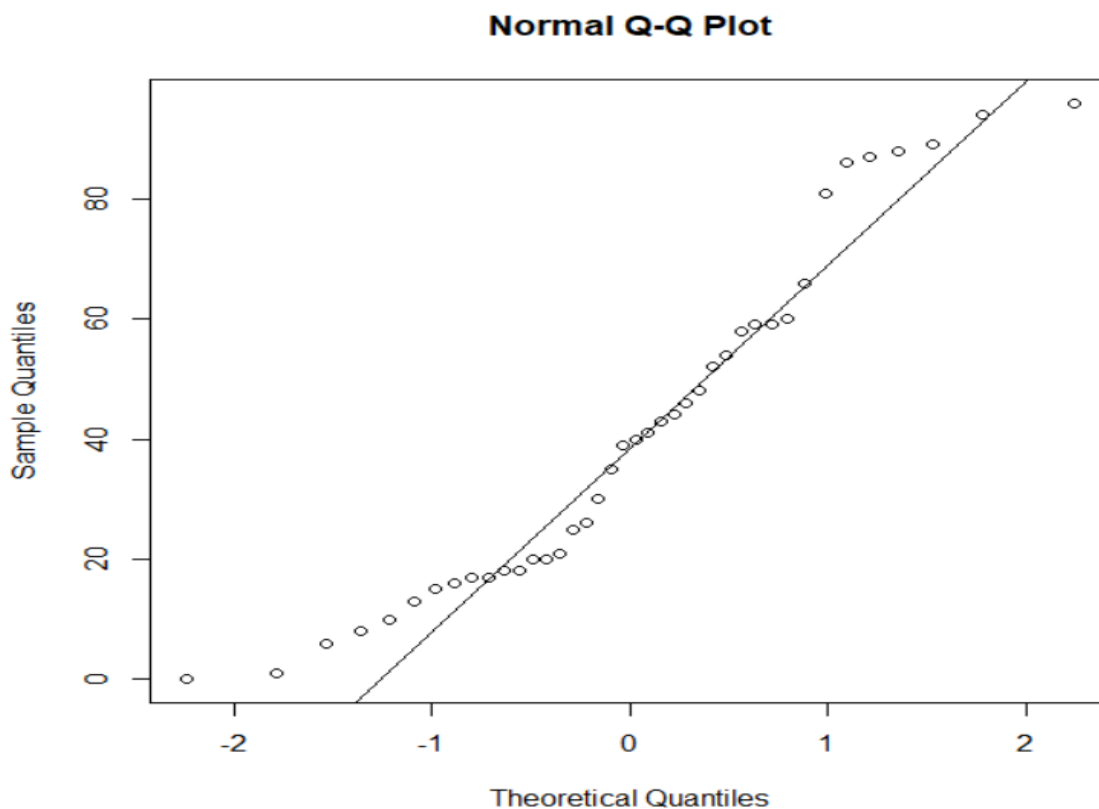
- Στο διάστημα $(\mu - \sigma, \mu + \sigma)$ δηλαδή το $(-0.4191212, 5.699121)$ ανήκουν οι 8 από τις 10 τιμές της ομάδας, δηλαδή το 80%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 68%, οπότε έχουμε απόκλιση +12%.
- Στο διάστημα $(\mu - 2\sigma, \mu + 2\sigma)$ δηλαδή το $(-3.478242, 8.758242)$ ανήκουν οι 9 από τις 10 τιμές της ομάδας, δηλαδή το 90%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 95%, οπότε έχουμε απόκλιση -5%.
- Στο διάστημα $(\mu - 3\sigma, \mu + 3\sigma)$ δηλαδή το $(-6.537363, 11.81736)$ ανήκουν και οι 10 από τις 10 τιμές της ομάδας, δηλαδή το 100%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 99,7%, οπότε έχουμε απόκλιση +0,03%.



Παρατηρούμε ότι τα σημεία δεν είναι και τόσο συγγραμικά, αφού δεν βρίσκονται τόσο κοντά στην ευθεία καθώς επίσης μοιάζουν να έχουν την μορφή μιας κυρτής συνάρτησης. Από το διάγραμμα μπορούμε να συμπεράνουμε ότι τα δεδομένα της ομάδας II δεν προσεγγίζουν τόσο καλά την καμπύλη πυκνότητας της κανονικής κατανομής.

Δεδομένα III

- Στο διάστημα $(\mu - \sigma, \mu + \sigma)$ δηλαδή το $(12.88246, 69.41754)$ ανήκουν οι 28 από τις 40 τιμές της ομάδας, δηλαδή το 70%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 68%, οπότε υπάρχει απόκλιση +2%.
- Στο διάστημα $(\mu - 2\sigma, \mu + 2\sigma)$ δηλαδή το $(-15.38508, 97.68508)$ ανήκουν και οι 40 από τις 40 τιμές της ομάδας, δηλαδή το 100%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 95%, οπότε έχουμε απόκλιση +5%.
- Στο διάστημα $(\mu - 3\sigma, \mu + 3\sigma)$ δηλαδή το $(-43.65262, 125.9526)$ ανήκουν και οι 40 από τις 40 τιμές της ομάδας, δηλαδή το 100%. Σύμφωνα με την κανονική κατανομή θα έπρεπε να είχαμε το 99.7%, οπότε έχουμε απόκλιση +0,03%.



Παρατηρούμε ότι το $\frac{1}{3}$ των σημείων και συγκεκριμένα τα κεντρικά σημεία είναι συγγραμμικά όποτε συνολικά μπορούμε να καταλήξουμε στο συμπέρασμα ότι τα συνολικά δεδομένα της ομάδας III δεν προσεγγίζουν αρκετά καλά την καμπύλη πυκνότητας της κανονικής κατανομής.

Άσκηση 2)

a) Τα δεδομένα μας προέρχονται από το “[Sports Reference | Sports Stats, fast, easy, and up-to-date | Sports-Reference.com](https://www.sports-reference.com/)” και αφορούν στατιστικά ομάδων του NBA μόνο για την κανονική περίοδο από την σεζόν 1946 μέχρι και την προηγούμενη ολοκληρωμένη. Περιέχονται 1633 περιπτώσεις και τα δεδομένα μας τυχαίνει να είναι ταξινομημένα σε φθίνουσα σειρά ως προς τον αριθμό νικών. Σε μερικές περιπτώσεις λείπουν ορισμένες τιμές και έχουν τιμή NA.

b) Τα δεδομένα μας έχουν 2 κατηγορηματικές μεταβλητές :

- Season (κατηγορηματική) : είναι η σεζόν και έχει την ακόλουθη μορφή ‘y-y+1’ όπου y είναι κάποιο έτος.
- State (κατηγορηματική) : είναι ο γεωγραφικός διαχωρισμός στον οποίο ανήκει η ομάδα και λαμβάνει τιμές East ή West.
- Team (κατηγορηματική) : είναι το υποκοριστικό της ομάδας με 3 γράμματα.

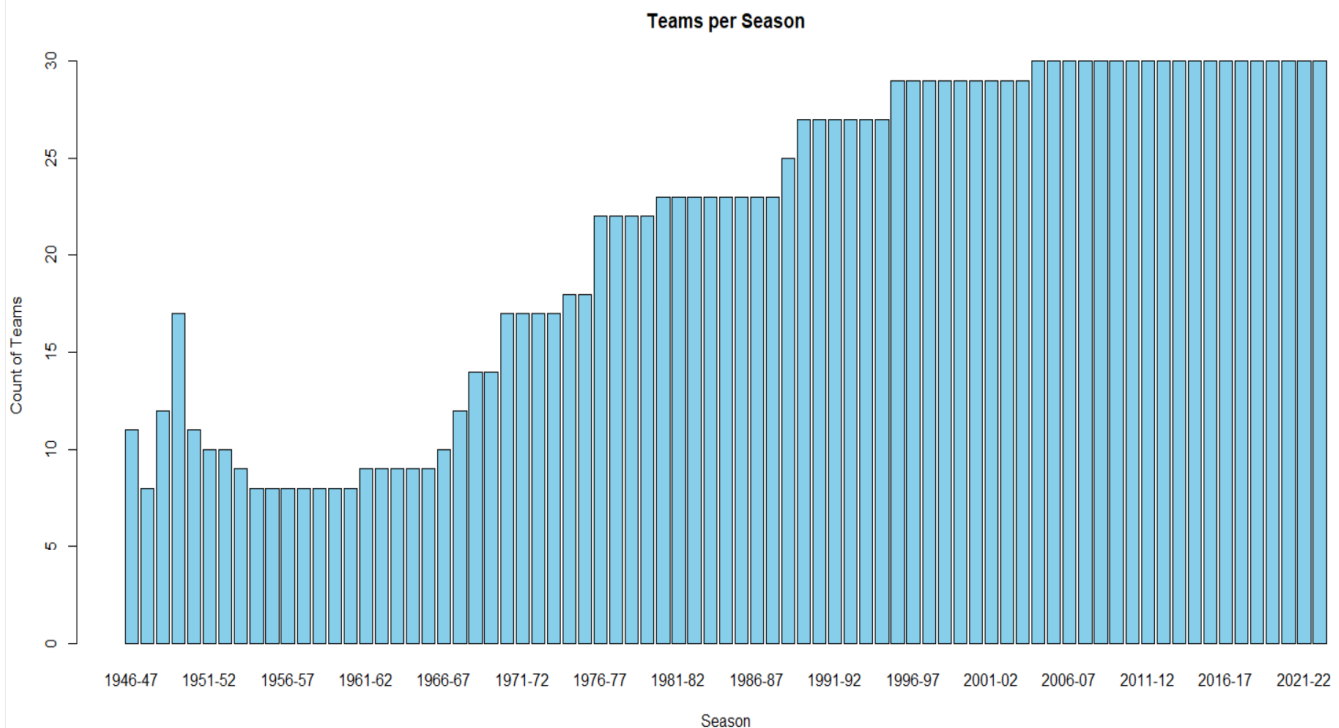
και 20 ποσοτικές εκ των οποίων οι σημαντικότερες είναι οι :

- W/L% : είναι το ποσοστό νικών-ηττών.
- G : το συνολικό πλήθος των αγώνων που έπαιξαν.
- MP : είναι τα συνολικά λεπτά που αγωνίστηκαν όλοι οι παίκτες της ομάδας.
- FG : είναι τα συνολικά εύστοχα καλάθια (υπάρχει επίσης το FGA που είναι οι συνολικές προσπάθειες).
- PTS : οι συνολικοί πόντοι που πέτυχαν.

Υπάρχουν επίσης και άλλες μεταβλητές που αντιστοιχούν σε άλλα στατιστικά του μπάσκετ όπως (τα συνολικά rebound -> TRB οι συνολικές assist-> AST και πολλά άλλα).

c) Σε αυτό το ερώτημα θα αναπαραστήσουμε μερικές από τις μεταβλητές μας και θα εξηγήσουμε την μορφή των κατανομών τους.

1) Season (κατηγορηματική)

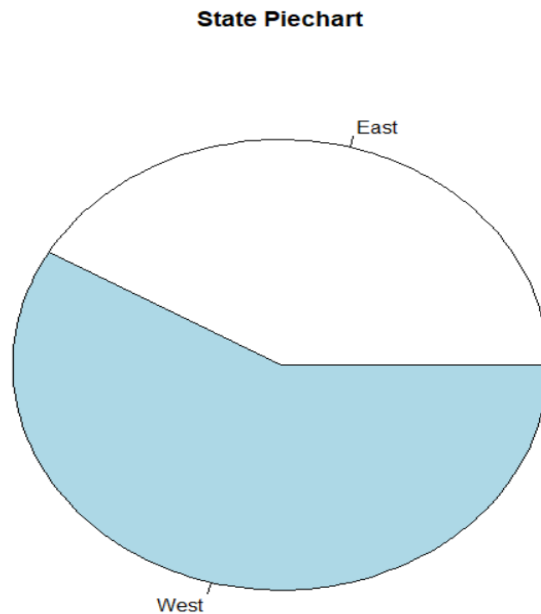


Στο παραπάνω barplot έχουμε απεικονίσει την κατανομή της κατηγορηματικής μεταβλητής “Season” και αυτή οφείλεται στο γεγονός ότι το NBA ξεκίνησε με 11 ομάδες και με την πάροδο των χρόνων εισάγονταν όλο και περισσότερες μέχρι που φτάσαμε στις 30 όπου είναι και ο αριθμός των ομάδων που είναι στο πρωτάθλημα τις τελευταίες σχεδόν 20 σεζόν. Παρατηρούμε επίσης ότι την σεζόν “1949-50” υπήρξαν 17 ομάδες, αρκετές περισσότερες από αυτές των γειτονικών σεζόν, γεγονός που οφείλεται στην ένωση δύο κατηγοριών σε μία, απόφαση που κράτησε μόνο για την συγκεκριμένη χρονιά.

Εντολές που ακολουθήσαμε στην R :

```
> barplot(table(Season),  
+         main = "Teams per Season",  
+         xlab = "Season",  
+         ylab = "Count of Teams",  
+         col = "skyblue")
```

2) State (κατηγορηματική)

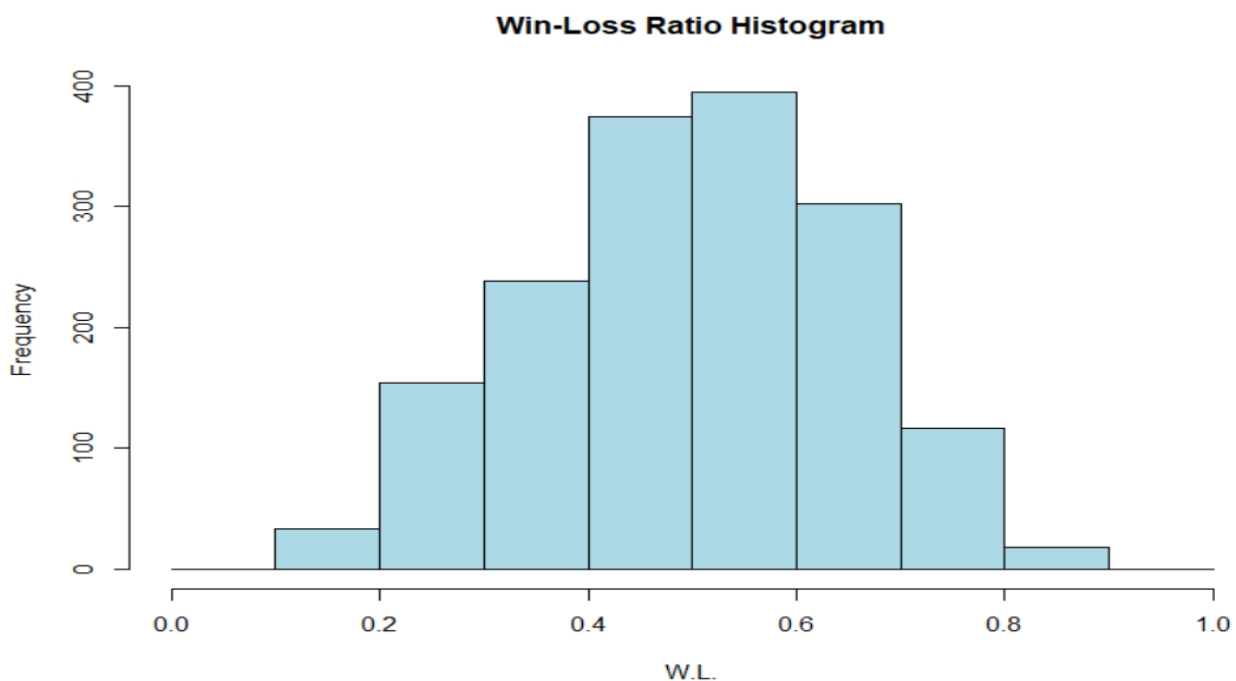


Στο παραπάνω piechart παρατηρούμε ότι σχεδόν το 60% των ομάδων στα δεδομένα μας ανήκει στον όμιλο της Δύσης και το υπόλοιπο 40% σε αυτόν της Ανατολής.

Εντολές που ακολουθήσαμε στην R :

```
> pie(table(State), main = "State Piechart")
```

3) Win-Loss Ratio or W-L% (ποσοτική)

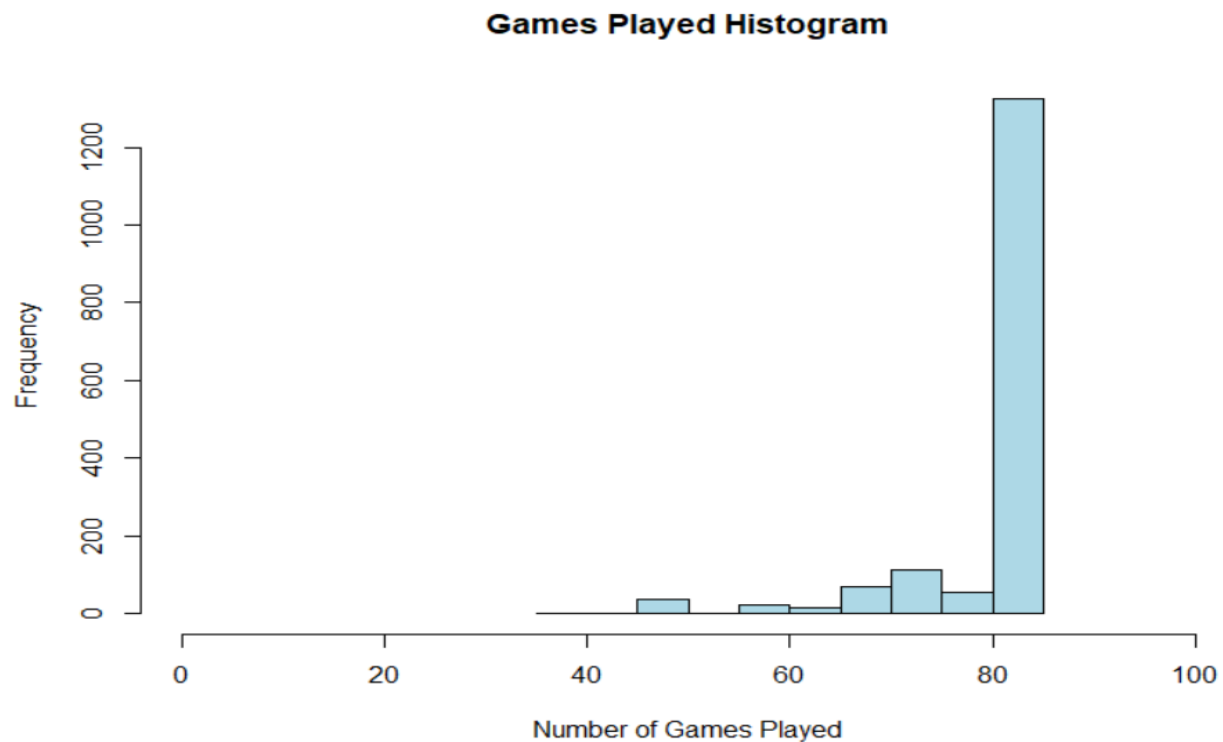


Από το παραπάνω ιστόγραμμα παρατηρούμε ότι η κατανομή των δεδομένων είναι πολύ κοντά στην κανονική, και το διάγραμμα δικαιολογείται από το γεγονός ότι υπάρχουν λίγες ομάδες που πετυχαίνουν κάθε σεζόν πολύ καλά ή πολύ κακά ($>80\%$ ή $<20\%$) ποσοστά νίκης-ήττας όμως οι περισσότερες ομάδες κυμαίνονται σε ένα ποσοστό γύρω στο 40% με 60% και για αυτόν τον λόγο υπάρχει τόσο μεγάλη συχνότητα στις κεντρικές τιμές.

Εντολές που ακολουθήσαμε στην R :

```
> hist(W.L.,  
+      main = "Win-Loss Ratio Histogram",  
+      xlim = c(0,1),  
+      breaks = seq(0, 1, by = 0.1),  
+      col = "lightblue"  
+ )
```

4) Games Played or G (ποσοτική)



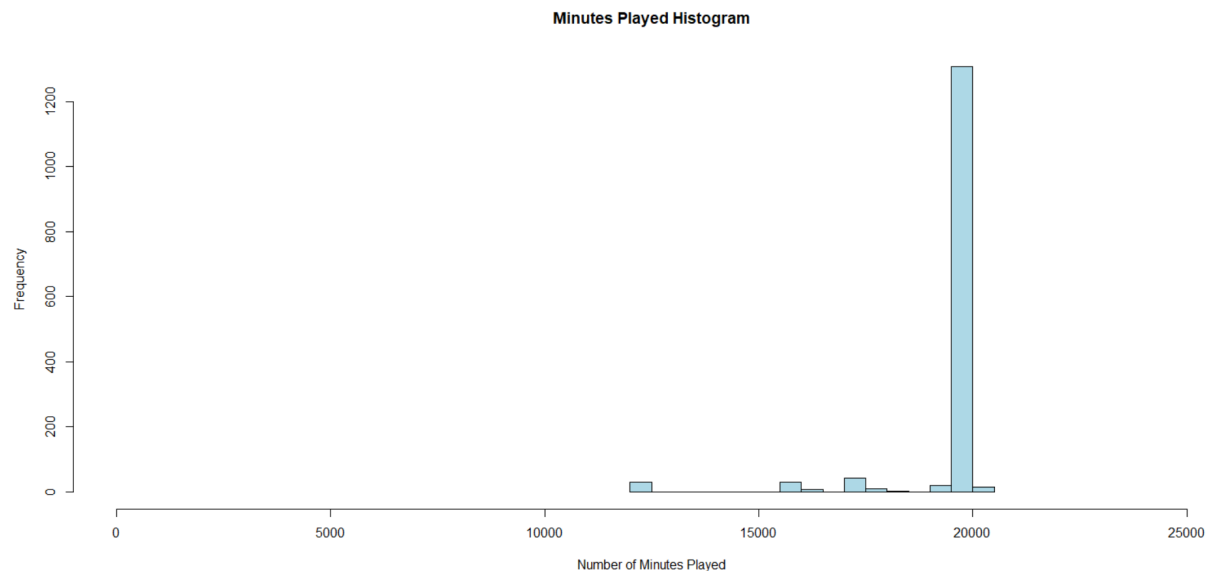
Από το παραπάνω histogram παρατηρούμε μια τεράστια συγκριτικά με τις άλλες συχνότητα στην τιμή 82, γεγονός που είναι λογικό αφού η κανονική διάρκεια μιας σεζόν στο NBA έχει 82 αγώνες για κάθε μία από τις ομάδες. Ωστόσο μπορούμε να δούμε ότι υπάρχουν και ορισμένα outliers και αυτό συμβαίνει επειδή μία ομάδα μπορεί να παραιτήθηκε από το πρωτάθλημα

στα μέσα της σεζόν έχοντας παίξει λιγότερους από της κανονικής διάρκειας αγώνες (> 40 και < 80).

Εντολές που ακολουθήσαμε στην R :

```
> hist(G,  
+      main = "Games Played Histogram",  
+      xlab = "Number of Games Played",  
+      xlim = c(0, 100),  
+      col = "lightblue"  
+ )
```

5) Minutes Played or MP (ποσοτική)

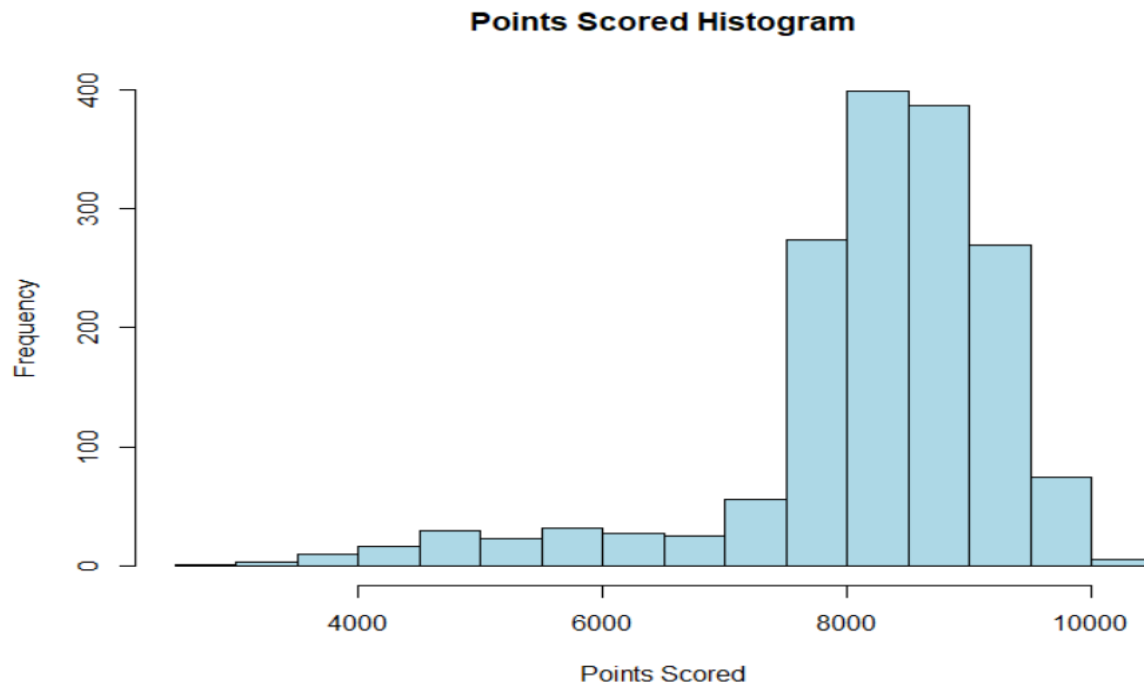


Όμοια με το παραπάνω ιστόγραμμα μπορούμε να παρατηρήσουμε ότι οι συντριπτική πλειονότητα των ομάδων τελειώνουν την σεζόν με τους παίκτες τους να έχουν παίξει από 19.750 μέχρι 19.900 λεπτά. Ο λόγος που έχουμε αποκλίσεις είναι ότι ορισμένοι αγώνες της κανονικής διάρκειας μπορεί να έχουν πάει στην παράταση οπότε οι παίκτες των ομάδων πρέπει να αγωνιστούν παραπάνω χρόνο. Όπως και πριν υπάρχουν outliers αφού οι παίκτες των ομάδων που εγκατέλειψαν την σεζόν είναι λογικό να αγωνίστηκαν λιγότερο χρόνο.

Εντολές που ακολουθήσαμε στην R :

```
> hist(  
+      MP,  
+      main = "Minutes Played Histogram",  
+      xlab = "Number of Minutes Played",  
+      ylab = "Frequency",  
+      col = "lightblue",  
+      xlim = c(0, 25000)  
+ )
```

6) Points scored or PTS (ποσοτική)



Παρατηρώντας το ιστόγραμμα βλέπουμε ότι το κύριο σώμα των τιμών είναι συμ- μετρικά κατανεμημένες στο διάστημα $\Lambda = (7000-1000]$. Μία παρατήρηση είναι ότι τα outliers πιθανότατα προκλήθηκαν από ομάδες που αποσύρθηκαν κάτι που θα ήταν λογικό αφού λιγότεροι αγώνες σημαίνει και λιγότεροι πόντοι συνολικά. Η μεγάλη συγκέντρωση στο διάστημα Λ είναι λογική αφού από έναν αγώνα ηba περιμένουμε γύρω στους 100 πόντους σε συνδυασμό με τα 82 παιχνίδια δηλαδή περίπου 8000 πόντους.

Εντολές που ακολουθήσαμε στην R :

```
> hist(PTS,  
+      main = "Points Scored Histogram",  
+      xlab = "Points Scored",  
+      col = "lightblue"  
+ )
```

d) Στην συνέχεια θα εξετάσουμε τις ποσοτικές μεταβλητές που έχουμε δώσει την μορφή των κατανομών τους στο προηγούμενο ερώτημα

- W-L(Win-Loss Ratio) :

1) mean -> 0.499823, standard deviation -> 0.149548

2) 5-num -> (min = 0.106, Q1 = 0.390, m = 0.512, Q3 = 0.610, max = 0.890)

Όπως φαίνεται από το ιστόγραμμα και το qq-plot η μεταβλητή ακολουθεί αρκετά πιστά την κανονική κατανομή που σημαίνει ότι οι τιμές W to L ratio είναι συμμετρικά κατανομημένες ως προς την μέση τιμή. Άρα η μέση τιμή και η τυπική απόκλιση είναι πιο κατάλληλες για τον χαρακτηρισμό αυτού του μεγέθους.

- G(Games-played):

1) mean -> 79.40233, standard deviation -> 6.652375

2) 5-num -> (min = 35, Q1 = 82, m = 82, Q3 = 82, max = 82)

Οι σεζόν στο NBA απαρτίζεται συνολικά από 82 αγώνες. Επομένως η πλειοψηφία των ομάδων θα παίξει 82 αγώνες ή ορισμένες λιγότερες σε περίπτωση που αποσυρθούν. Αυτή η παρατήρηση μαζί με το ιστόγραμμα του παραπάνω ερωτήματος μας οδηγεί στο συμπέρασμα ότι η κατανομή των τιμών δεν είναι συμμετρική γύρω από την μέση τιμή καθιστώντας την επιλογή της σύνοψης των 5 αριθμών καλύτερο μετρικό για την κατανομή της G. Επίσης μπορούμε να συμπεράνουμε ότι πάνω από το 75% των ομάδων δεν αποσύρθηκαν πρόωρα και η συντομότερη απόσυρση ήταν μετά τον 35ο αγώνα.

- MP(Minutes-Played):

1) mean -> 19467.32, standard deviation -> 1289.573

2) 5-num -> (min = 12025, Q1 = 19755, m = 19805, Q3 = 19855, max = 20080)

Όμοια με το παραπάνω μέγεθος λόγω της μη συμμετρικότητας των τιμών γύρω από την μέση τιμή και το γεγονός ότι υπάρχουν ατυπικές τιμές καλύτερα χαρακτηρίζουν τα minutes played η σύνοψη των 5 αριθμών. Μάλιστα λόγω ατυπικών τιμών ακόμα και η μέση τιμή δεν μας δίνει σωστή πληροφορία για την μέση διάρκεια των αγώνων κάθε σεζόν, αφού το mean είναι εκτός του

διαστήματος της τυπικής διάρκειας 19.750 μέχρι 19.900 λεπτών που αναφέρθηκε σε προηγούμενο ερώτημα.

- PTS(Points scored):

1) mean -> 8201.246, standard deviation -> 1147.32

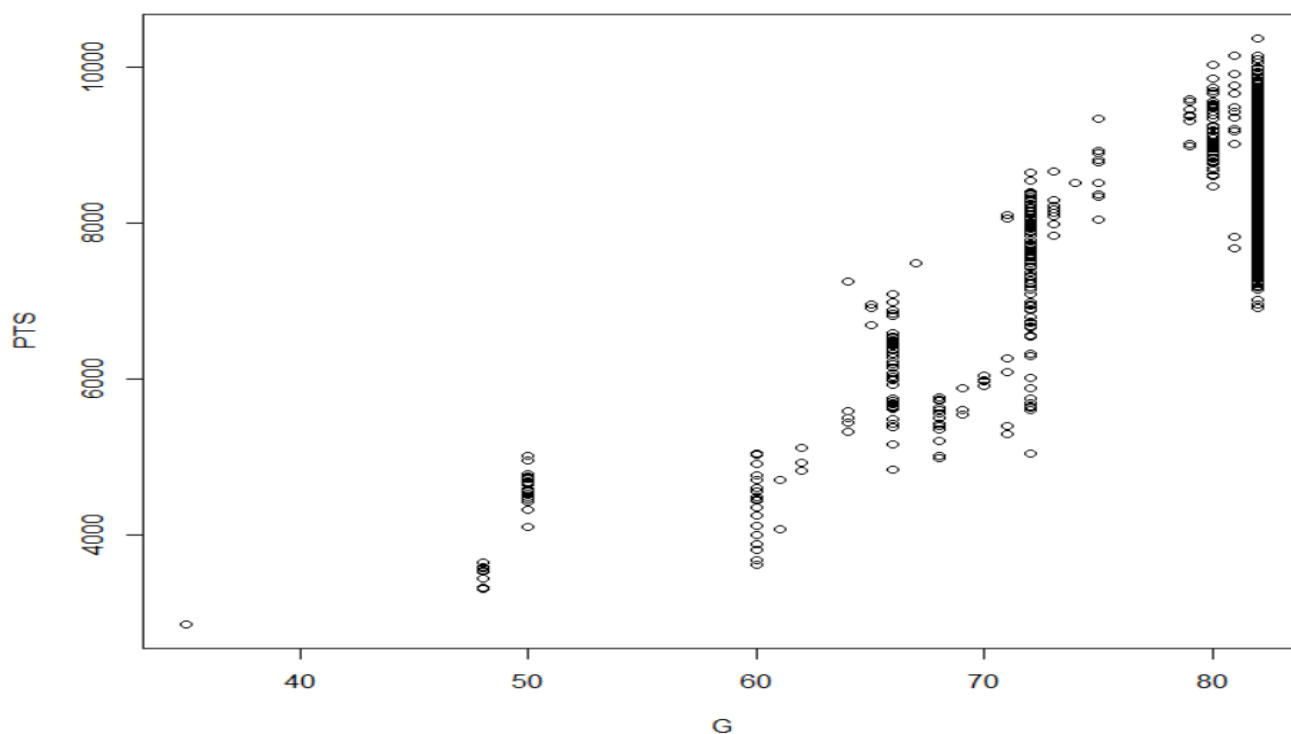
2) 5-num -> (min = 2844, Q1 = 7892, m = 8398, Q3 = 8906, max = 10371)

Λόγω της μη συμμετρικότητας των τιμών γύρω από την μέση τιμή και το γεγονός ότι υπάρχουν ατυπικές τιμές καλύτερος είναι ο χαρακτηρισμός του five number summary. Επίσης μας οδηγεί στο συμπέρασμα που καταλήξαμε ότι το 75% των αγώνων έχει από 95 έως περίπου 126 πόντους. Παρόλα αυτά τυχαίνει να μπορούμε επίσης εύστοχα να χαρακτηρίσουμε και τους πόντους και με την μέση τιμή και την απόκλιση αφού μας δίνει την διαίσθηση ότι οι περισσότερες ομάδες έχουν ανά σεζόν περίπου 8000 πόντους με τυπική απόκλιση 1100 κάτι που συμβαδίζει με το ιστόγραμμα. Εν τέλη ίσως πιο κατατοπιστική είναι η σύνοψη των 5 αριθμών αλλά η διαφορά είναι θέμα επιλογής και τι συμπεράσματα προσπαθούμε να αντλήσουμε από αυτά τα μετρικά.

e) Σχέση Games Played - Points Scored :

Θα επιλέξουμε την μεταβλητή G και την μεταβλητή PTS και θα διερευνήσουμε αν υπάρχει σχέση μεταξύ τους, δηλαδή αν τα συνολικά παιχνίδια που παίζει μία ομάδα σε μία σεζόν σχετίζονται με τους πόντους που πετυχαίνει σε αυτήν. Ας θεωρήσουμε λοιπόν, ότι η επεξηγηματική μεταβλητή είναι η G ενώ μεταβλητή απόκρισης η PTS. Αφού μιλάμε για 2 ποσοτικές μεταβλητές είναι χρήσιμο να δούμε το scatter plot.

Απο το παρακάτω scatterplot μπορούμε να διακρίνουμε μια όχι τόσο γραμμική αλλά αύξουσα σχέση μεταξύ των δύο μεταβλητών οπότε μένει να βρούμε τον συντελεστή συσχέτισης τους για να δούμε κατά πόσο υπάρχει με σιγουριά και αν ναι πόσο ισχυρή είναι.



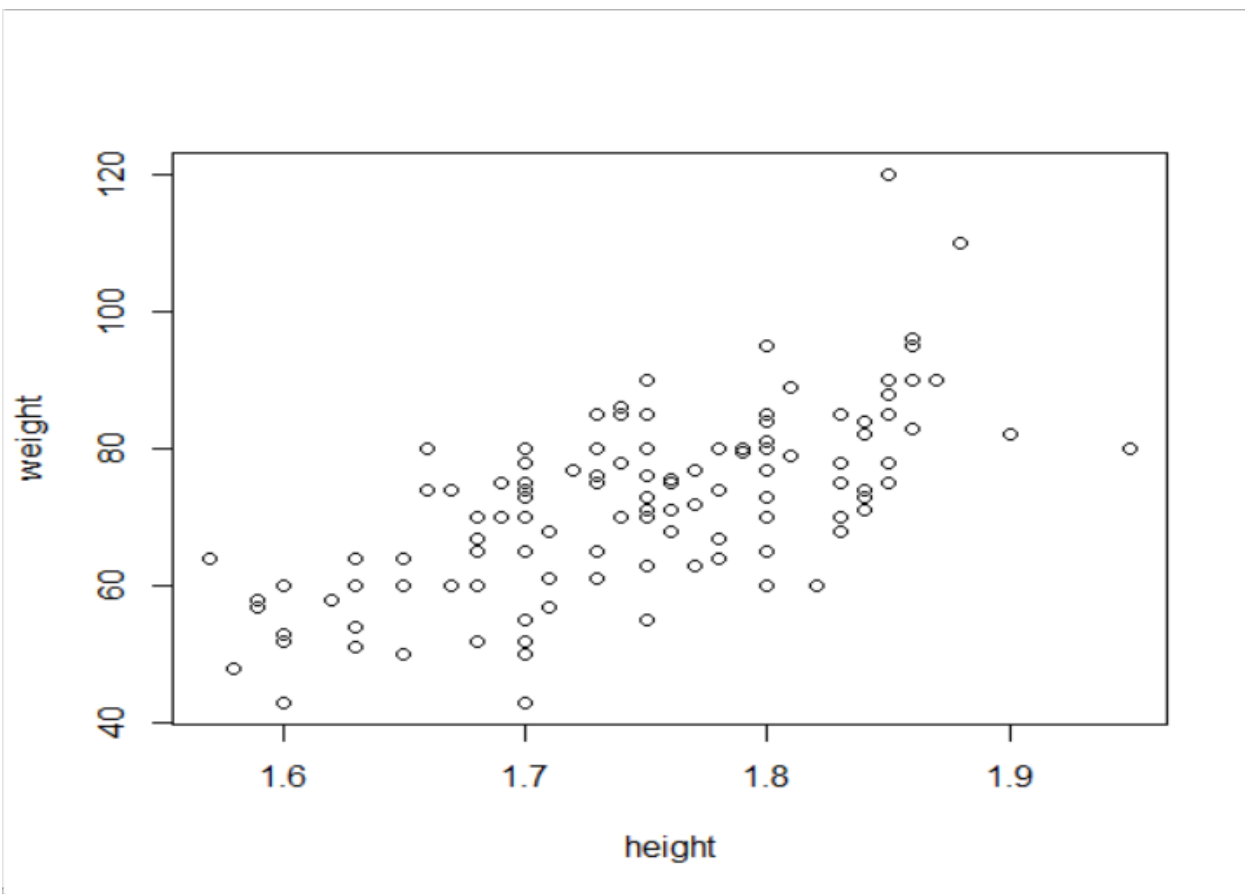
Με την χρήση της R `> cor(G, PTS)` βρίσκουμε ότι ο συντελεστής συσχέτισης είναι $r = 0.8178434$, γεγονός που μας αποδεικνύει ότι η σχέση είναι αρκετά ισχυρή. Η σχέση των δύο μεταβλητών είναι αιτιατή, αφού αν μεταβάλλουμε το G είναι σίγουρο ότι θα μεταβληθεί και το PTS και αυτό συμβαίνει επειδή οι πόντοι που πετυχαίνει κατά μέσο όρο μια ομάδα ανα αγώνα στο NBA είναι από 90 μέχρι 110 πόντους. Επομένως το γεγονός να μην αγωνιστεί σε ένα ή περισσότερους αγώνες έχει μεγάλη επίπτωση στους συνολικούς πόντους που θα πετύχει στην σεζόν.

Άσκηση 3)

- a) Από τα δεδομένα του ερωτηματολογίου του 2023 επιλέξαμε τις ποσοτικές μεταβλητές height και weight για το συγκεκριμένο ερώτημα. Για την κατασκευή του scatterplot ορίσαμε ως επεξηγηματική μεταβλητή το ύψος (x axis) και ως απόκριση το βάρος (y axis) καθώς με βάση το ύψος προκύπτει το βάρος και καλέσαμε στην R την εντολή :

```
> plot(height, weight)
```

Κάθε περίπτωση i μας δίνει ένα σημείο (x_i, y_i) όπου x_i = ύψος και y_i = βάρος. Επομένως έχουμε το ακόλουθο scatterplot :



Όπως φαίνεται και από το διάγραμμα :

- η μορφή της σχέσης είναι γραμμική
- η κατεύθυνση της σχέσης είναι αύξουσα
- η δύναμη της σχέσης είναι θετική και ισχυρή

b) Για να υπολογίσουμε τον συντελεστή συσχέτισης των δύο μεταβλητών θα χρησιμοποιήσουμε τον τύπο :

$$r = \frac{1}{n-1} \sum_i \left(\frac{x_i - \bar{x}}{s_x} \right) \left(\frac{y_i - \bar{y}}{s_y} \right)$$

όπου έχουμε :

- n = πλήθος των περιπτώσεων
- x_i και y_i οι τιμές του x και y στην i -οστή περίπτωση
- \bar{x} και \bar{y} η μέση τιμή του x και του y αντίστοιχα
- s_x και s_y η τυπική απόκλιση του x και του y αντίστοιχα

Κάνοντας τους υπολογισμούς για $x = \text{height}$ και $y = \text{weight}$ στην περίπτωση μας καταλήγουμε στις παρακάτω τιμές.

$n = 129$, $\bar{x} = 1.746744$, $\bar{y} = 71.85938$, $s_x = 0.08070164$ και $s_y = 12.59702$

Οπότε συνδυάζοντας τις παραπάνω τιμές και καλώντας στην R την :

```
> cor(height, weight, use="complete.obs")
```

Προκύπτει λοιπόν, ο συντελεστής συσχέτισης $r = 0.6718937$, γεγονός που αποδεικνύει τα συμπεράσματα μας για την σχέση των δύο μεταβλητών.

Στη συνέχεια, για να εκτελέσουμε την γραμμική παλινδρόμηση ελαχίστων τετραγώνων ακολουθήσαμε στην R τις εντολές :

```
> lm(weight ~ height) -> m  
> abline(m)
```

Η εκτέλεση της πρώτης εντολής μας επιστρέφει και αναθέτει στο “m” τους συντελεστές της γραμμικής συνάρτησης που ελαχιστοποιεί τα τετράγωνα των αποστάσεων μεταξύ του σημείου και της ευθείας. Και η δεύτερη εντολή εμφανίζει αυτή την ευθεία στο ήδη υπάρχον scatterplot. Επομένως έχουμε το ακόλουθο διάγραμμα :

