

11. Suppose that the data for analysis includes the attribute age. The age values for the data tuples are (in increasing order) 13, 15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70.

Can you find (roughly) the first quartile (Q1) and the third quartile (Q3) of the data?

```

R version 4.4.2 (2024-10-31 ucrt) -- "File of Leaves"
Copyright (C) 2024 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

[Previously saved workspace restored]

> marks<-c(13,15, 16, 16, 19, 20, 20, 21, 22, 22, 25, 25, 25, 25, 30, 33, 33, 35, 35, 35, 35, 36, 40, 45, 46, 52, 70)
> quantile(marks)
 0% 25% 50% 75% 100%
13.0 20.5 25.0 35.0 70.0
  
```

12. Covariance and correlation

Children of three ages are asked to indicate their preference for three photographs of adults. Do the data suggest that there is a significant relationship between age and photograph preference? What is wrong with this study?

		Photograph:		
Age of child		A	B	C
5-6 years:	18		22	20
7-8 years:	2		28	40
9-10 years:	20		10	40

- (i) Use `cov()` to calculate the sample covariance between B and C.
- (ii) Use another call to `cov()` to calculate the sample covariance matrix for the preferences.
- (iii) Use `cor()` to calculate the sample correlation between B and C.
- (iv) Use another call to `cor()` to calculate the sample correlation matrix for the preferences.

```

# Data for photography preferences
age_groups <- c("5", "15", "25", "35", "45")
A <- c(18, 2, 20)
B <- c(22, 28, 10)
C <- c(20, 40, 40)

# Creating a data frame
data <- data.frame(A, B, C)

# (i) Calculate the sample covariance between B and C
cov_B_C <- cov(data$B, data$C)
print(paste("Covariance between B and C:", cov_B_C))

# (ii) Calculate the sample covariance matrix for the preferences
cov_matrix <- cov(data)
print(cov_matrix)

# (iii) Calculate the sample correlation between B and C
cor_B_C <- cor(data$B, data$C)
print(paste("Correlation between B and C:", cor_B_C))

# (iv) Calculate the sample correlation matrix for the preferences
cor_matrix <- cor(data)
print(cor_matrix)

```

```

> # (i) Calculate the sample covariance between B and C
> cov_matrix <- cov(data)
> print("Covariance Matrix:")
[1] "Covariance Matrix:"
> print(cov_matrix)
      A      B      C
A 97.33333 -74 -46.66667
B -74.00000  84 -20.00000
C -46.66667 -20 133.33333

> # (iii) Calculate the sample correlation between B and C
> cor_B_C <- cor(data$B, data$C)
> print(paste("Correlation between B and C:", cor_B_C))
[1] "Correlation between B and C: -0.188982236504614"

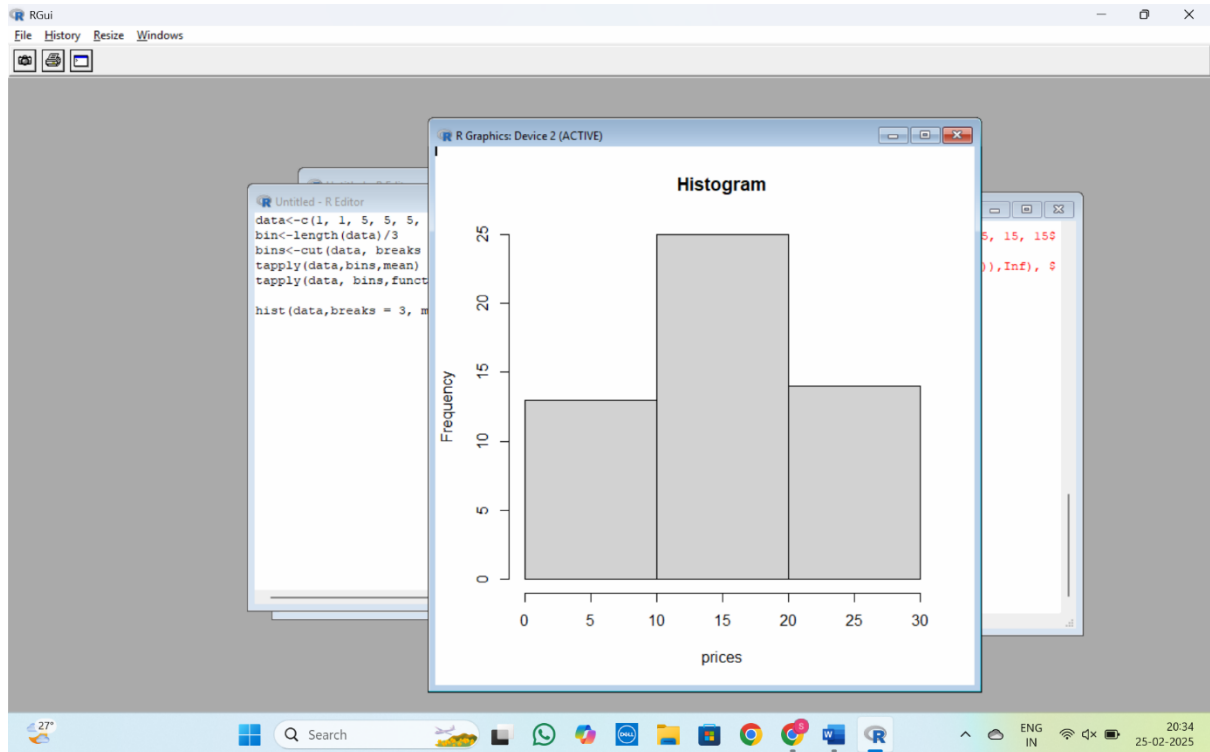
> # (iv) Calculate the sample correlation matrix for the preferences
> cor_matrix <- cor(data)
> print("Correlation Matrix:")
[1] "Correlation Matrix:"
> print(cor_matrix)
      A      B      C
A 1.0000000 -0.8183918 -0.4096440
B -0.8183918 1.0000000 -0.1889822
C -0.4096440 -0.1889822 1.0000000

```

13. Imagine that you have selected data from the All Electronics data warehouse for analysis. The data set will be huge! The following data are a list of All Electronics prices for commonly sold items (rounded to the nearest dollar). The numbers have been sorted: 1, 1, 5, 5, 5, 5, 5, 8, 8, 10, 10, 10, 10, 12, 14, 14, 14, 15, 15, 15, 15, 15, 15, 18, 18, 18, 18, 18,

18, 18, 18, 20, 20, 20, 20, 20, 20, 20, 20, 21, 21, 21, 21, 25, 25, 25, 25, 25, 28, 28, 30, 30, 30.

- (i) Partition the dataset using an equal-frequency partitioning method with bin equal to 3
- (ii) apply data smoothing using bin means and bin boundary.
- (iii) Plot Histogram for the above frequency division



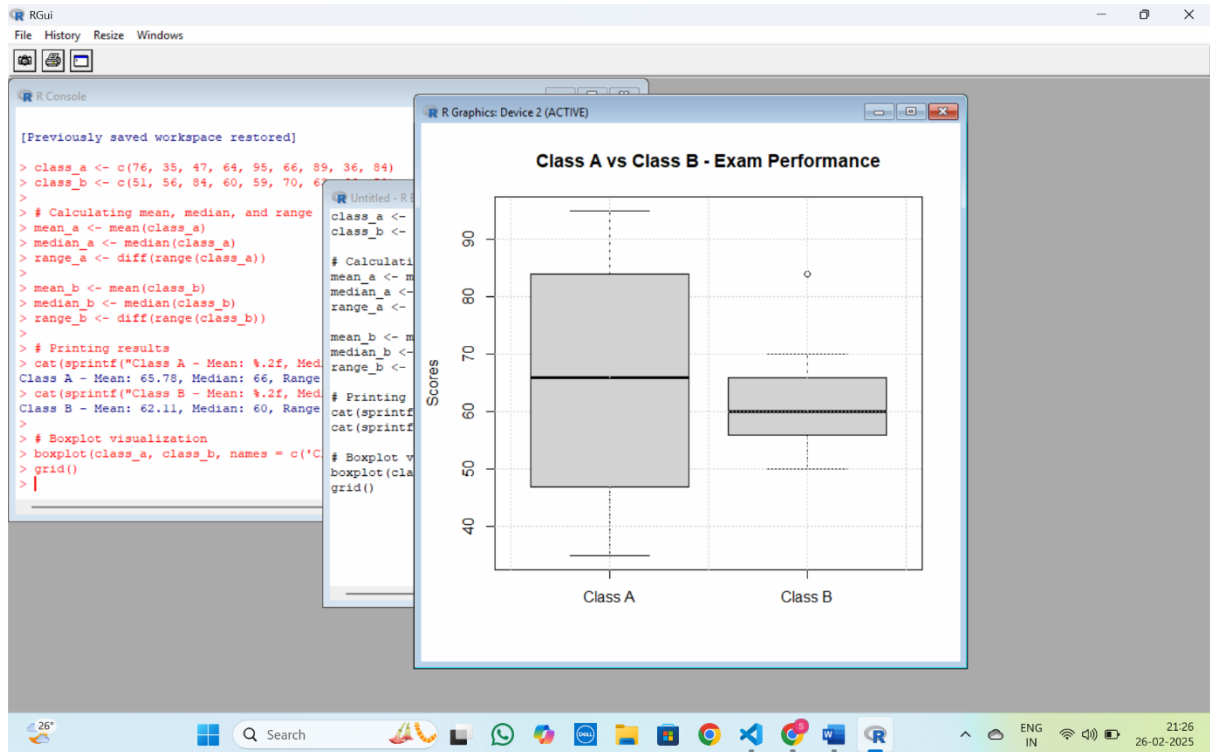
14. Two Maths teachers are comparing how their Year 9 classes performed in the end of year exams. Their results are as follows:

Class A: 76, 35, 47, 64, 95, 66, 89, 36, 84

Class B: 51, 56, 84, 60, 59, 70, 63, 66, 50

- Find which class had scored higher mean, median and range.
- Plot above in boxplot and give the inferences

Class B: 51, 56, 84, 60, 59, 70, 63, 66, 50

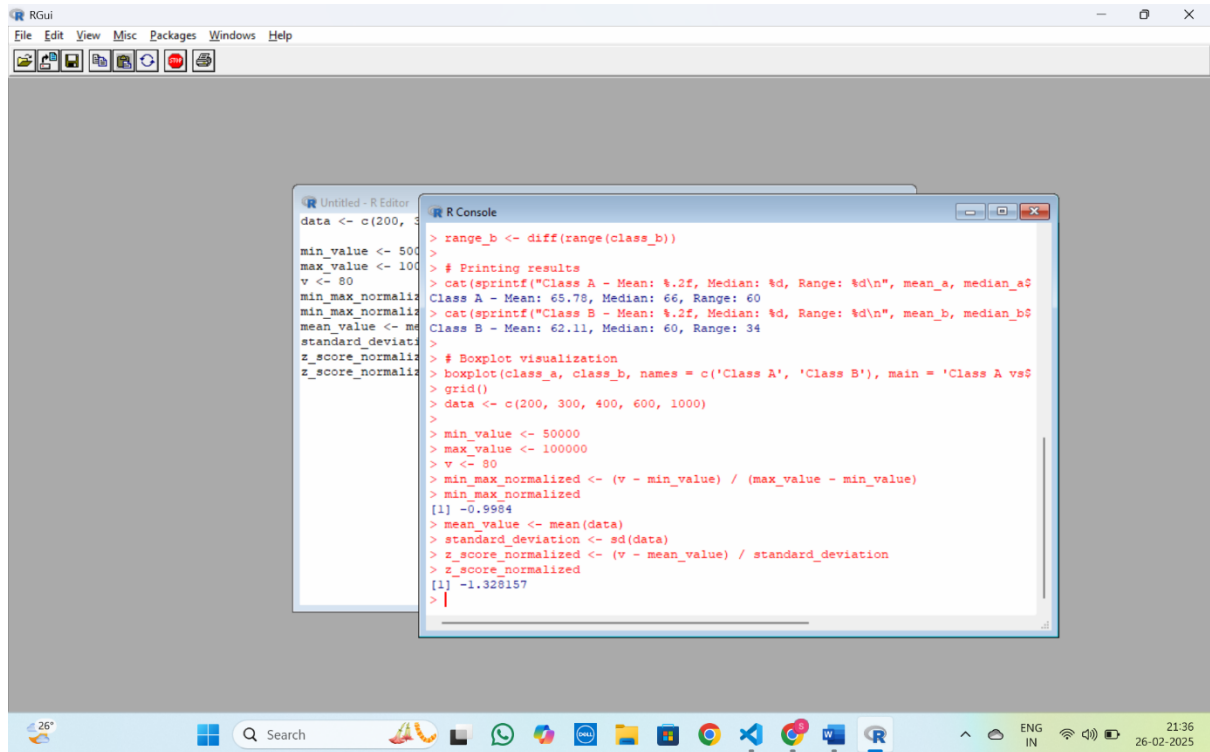


15. Let us consider one example to make the calculation method clear. Assume that the minimum and maximum values for the feature F are \$50,000 and \$100,000 correspondingly. It needs to range F from 0 to 1. In accordance with min-max normalization, $v = \$80$,

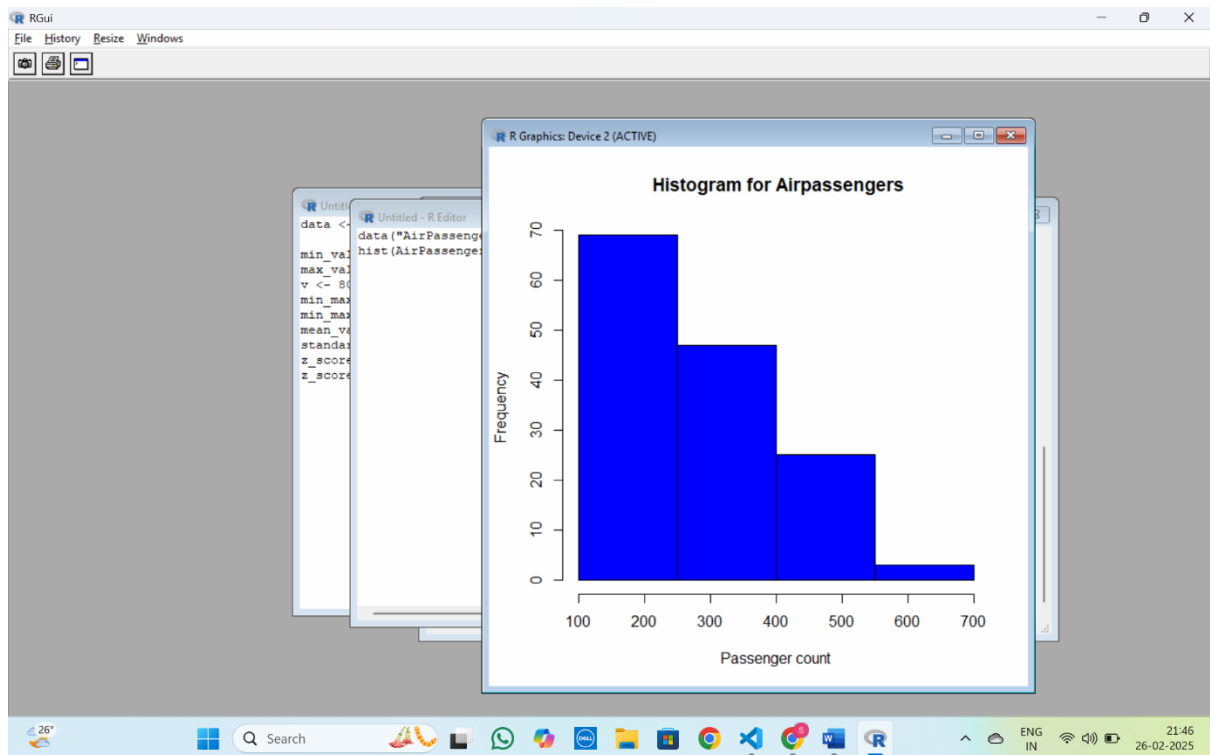
b) Use the two methods below to normalize the following group of data: 200, 300, 400, 600, 1000

(a) min-max normalization by setting $\min = 0$ and $\max = 1$

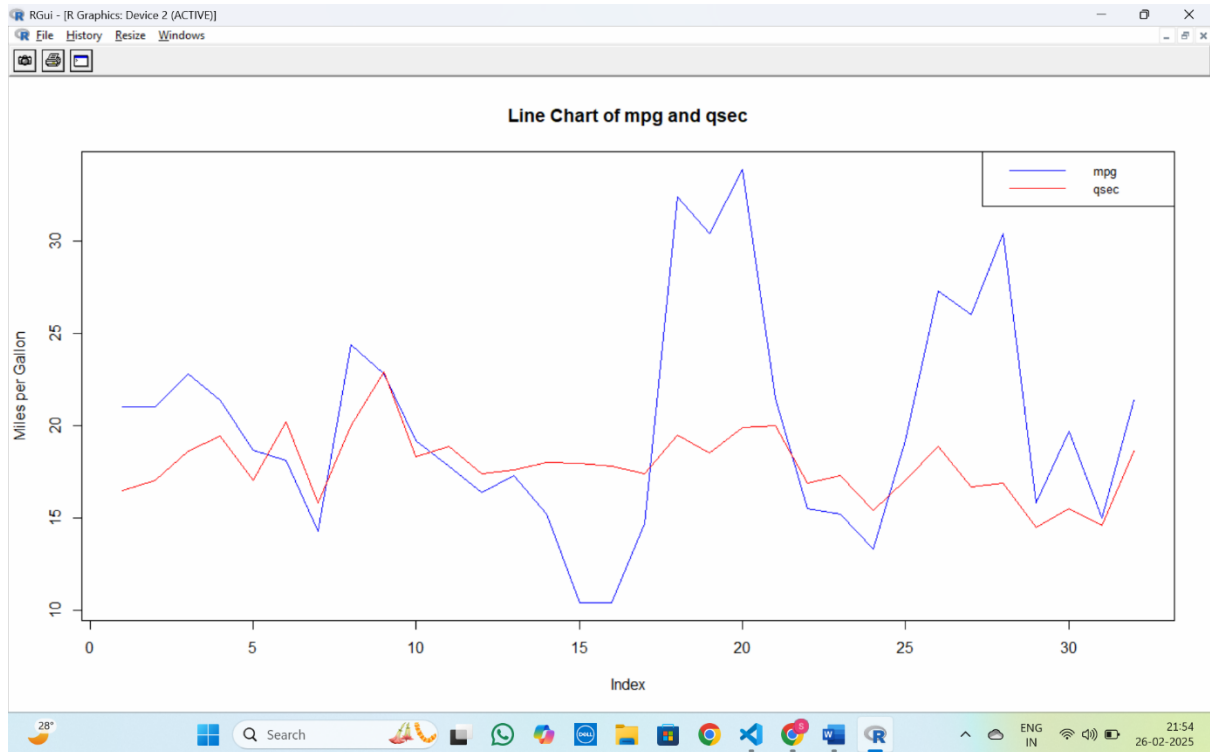
(b) z-score normalization



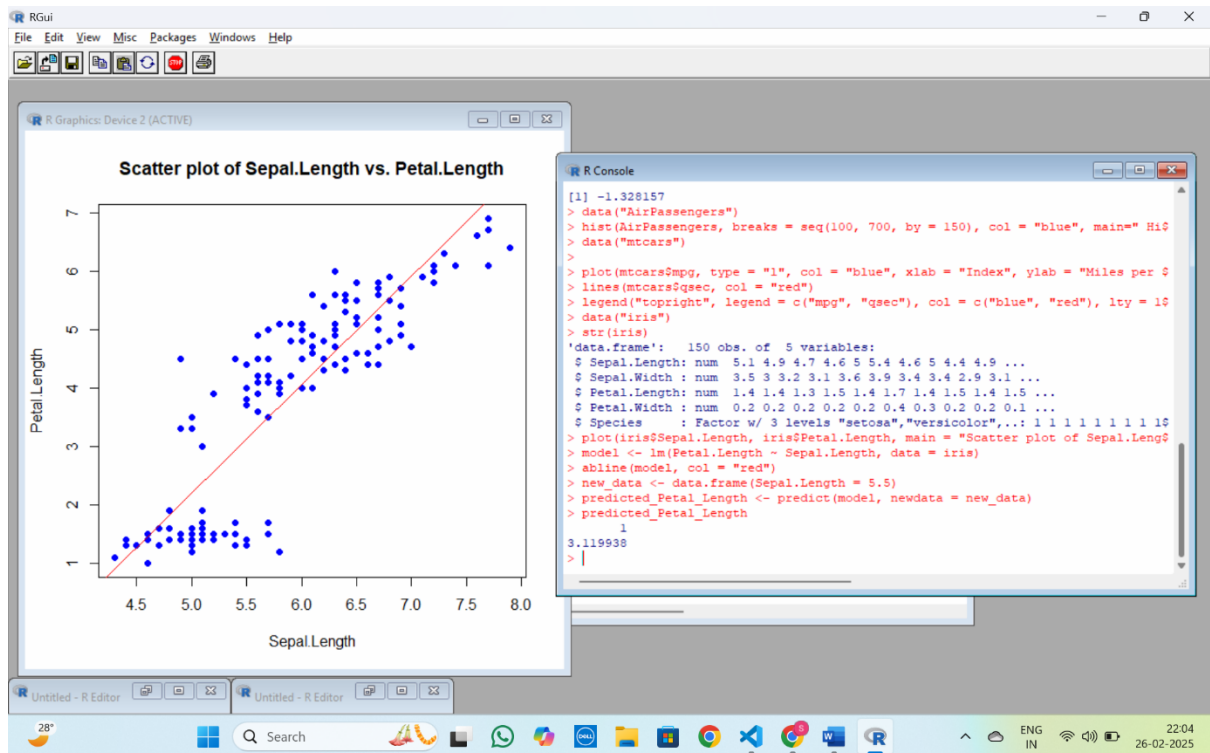
16. Make a histogram for the “AirPassengers” dataset, start at 100 on the x-axis, and from values 200 to 700, make the bins 150 wide



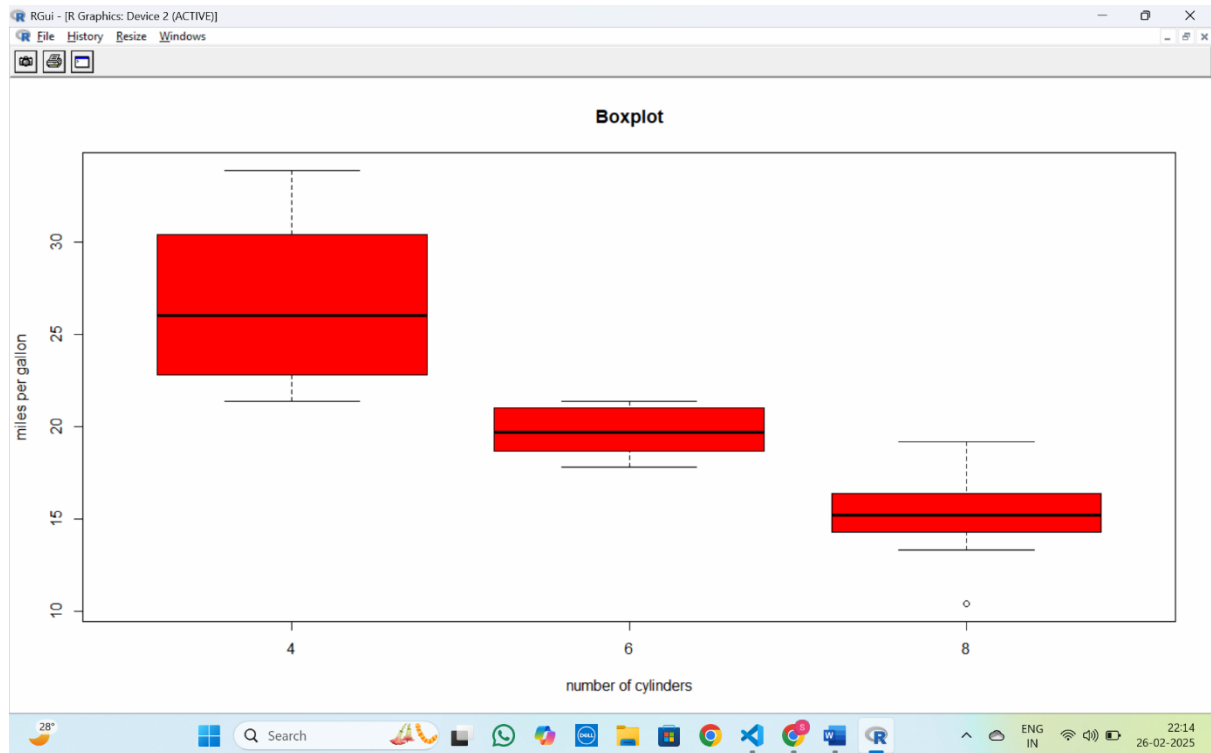
17. Obtain Multiple Lines in Line Chart using a single Plot Function in R. Use attributes “mpg” and “qsec” of the dataset “mtcars”



18. Download the Dataset "water" From R dataset Link. Find out whether there is a linear relation between attributes "mortality" and "hardness" by plot function. Fit the Data into the Linear Regression model. Predict the mortality for the hardness=88.



19. Create a Boxplot graph for the relation between "mpg"(miles per galloon) and "cyl"(number of Cylinders) for the dataset "mtcars" available in R Environment.



20. Assume the Tennis coach wants to determine if any of his team players are scoring outliers. To visualize the distribution of points scored by his players, then how can he decide to develop the box plot? Give suitable example using Boxplot visualization technique.

```
RGU - [R Console]
File Edit View Misc Packages Windows Help

> min_value <- 50000
> max_value <- 100000
> v <- 80
> min_max_normalized <- (v - min_value) / (max_value - min_value)
> min_max_normalized
[1] -0.9984
> mean_value <- mean(data)
> standard_deviation <- sd(data)
> z_score_normalized <- (v - mean_value) / standard_deviation
> z_score_normalized
[1] -1.325157
> data("AirPassengers")
> hist(AirPassengers, breaks = seq(100, 700, by = 150), col = "blue", main="Histogram for Airpassengers", xlab = "Passenger count", ylab = "Frequency")
> data("mtcars")
>
> plot(mtcars$mpg, type = "l", col = "blue", xlab = "Index", ylab = "Miles per Gallon", main = "Line Chart of mpg and qsec")
> lines(mtcars$qsec, col = "red")
> legend("topright", legend = c("mpg", "qsec"), col = c("blue", "red"), lty = 1, cex = 0.8)
> data("iris")
> str(iris)
'data.frame': 150 obs. of 5 variables:
 $ Sepal.Length: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
 $ Sepal.Width : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
 $ Petal.Length: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
 $ Petal.Width : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
 $ Species      : Factor w/ 3 levels "setosa","versicolour",...: 1 1 1 1 1 1 1 1 1 1 ...
> plot(iris$Sepal.Length, iris$Petal.Length, main = "Scatter plot of Sepal.Length vs. Petal.Length", xlab = "Sepal.Length", ylab = "Petal.Length", col = "blue")
> model <- lm(Petal.Length ~ Sepal.Length, data = iris)
> abline(model, col = "red")
> new_data <- data.frame(Sepal.Length = 5.5)
> predicted_Petal_Length <- predict(model, newdata = new_data)
> predicted_Petal_Length
1
3.119938
> data("mtcars")
> boxplot(mpg ~ cyl, data = mtcars, main = "Boxplot", xlab = "number of cylinders", ylab = "miles per gallon", col = "red")
> score <- c(20, 25, 30, 32, 35, 38, 40, 45, 50, 52, 55, 56, 58, 59, 60, 62, 65, 70, 75, 80, 85)
>
> boxplot(score, col = "lightblue", main = "Box Plot of Points Scored by Tennis Players", ylab = "Points Scored")
>
>
```