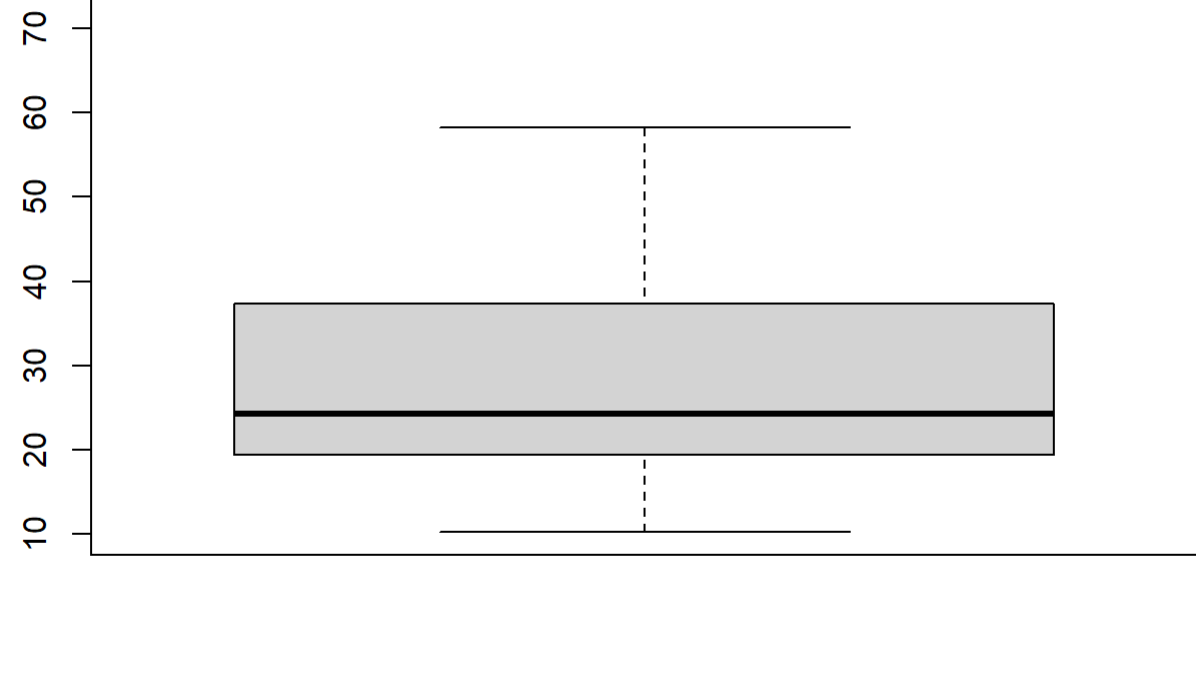


Q1) a.

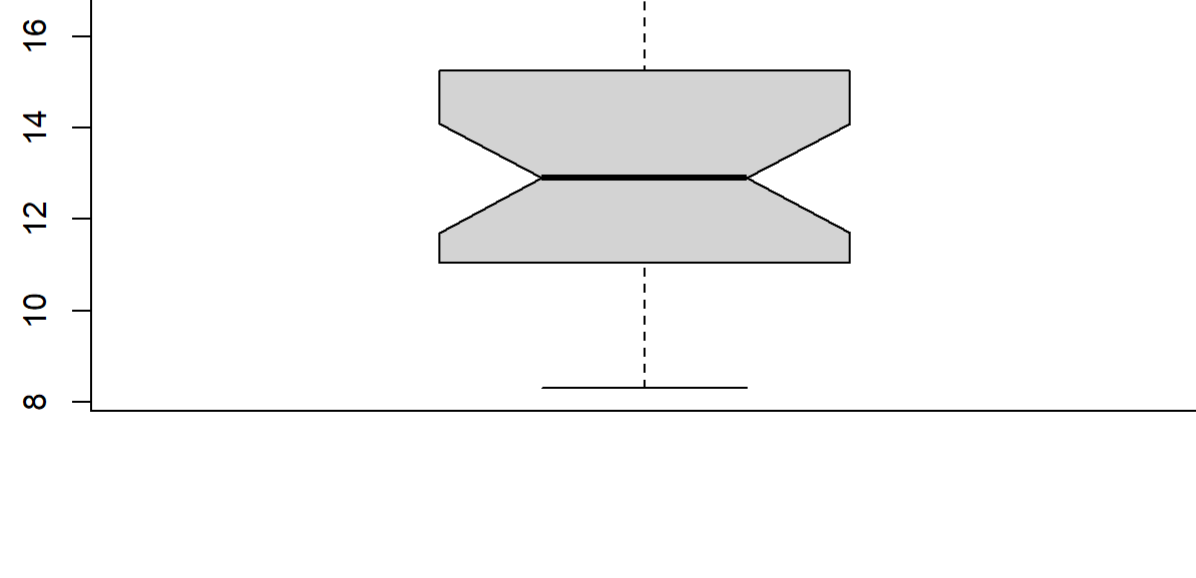
```
1 data(trees)
2 boxplot(trees$Volume, varwidth=TRUE)
```



I see a 5 number summary of the Volume of the trees present in the dataset with one outlier which has value greater than $1.5 \times \text{IQR}(Q3-Q1)$, where $Q3 \sim 35$ & $Q1 \sim 15$ -20. The median of the data points lies between 20 to 25. The data seems to be right skewed as the boxplot has slight long whisker towards the right.

Q1) b.

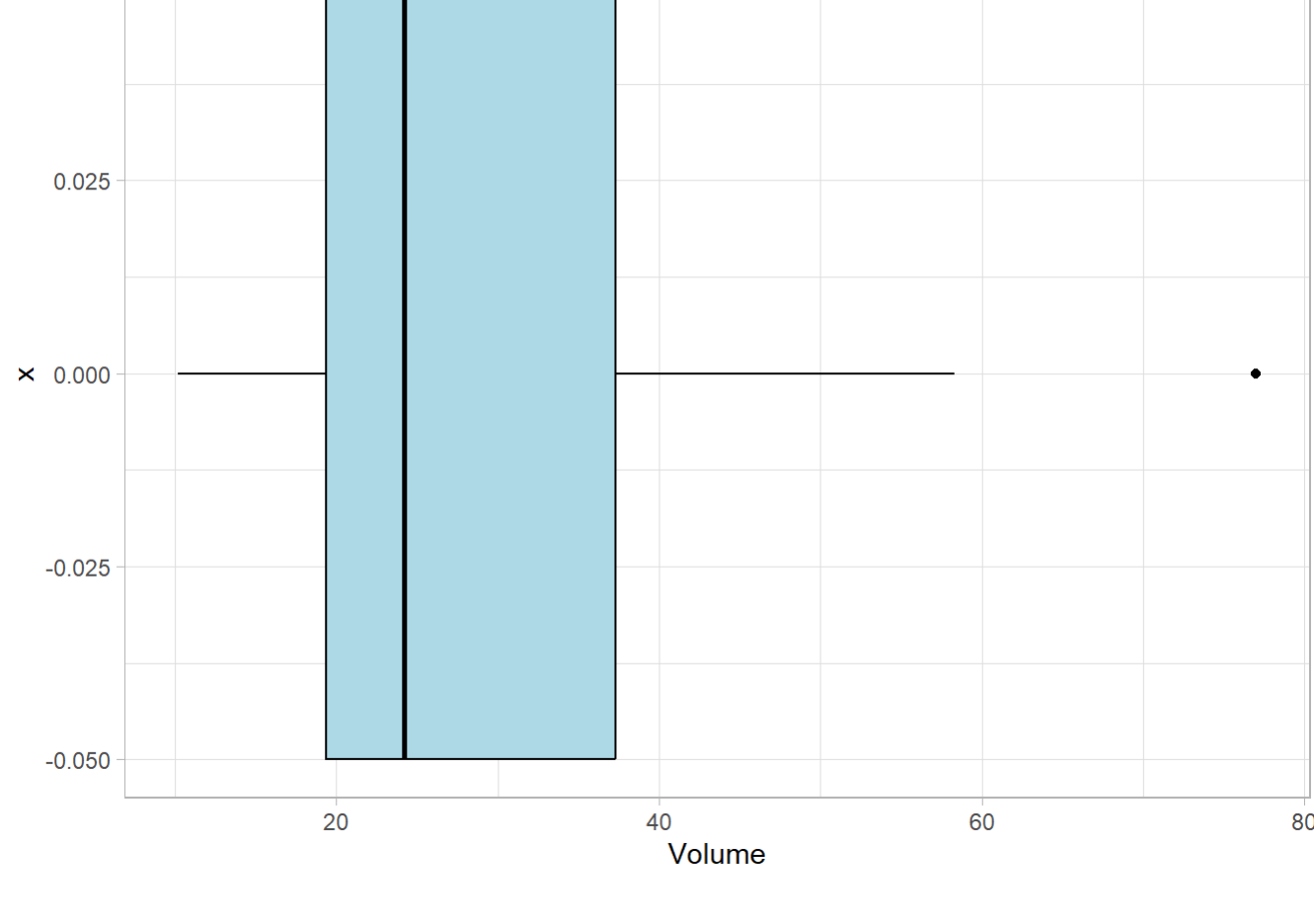
```
1 boxplot(trees$Girth, notch = TRUE)
```



In the above Boxplot, the notch=TRUE parameter is used to add notches to the box plot which helps in comparing the medians if multiple boxplots were there in the same visual. Notches are used to compare groups; if the notches of two boxes do not overlap, this is a strong evidence that the medians differ.

Q2) a.

```
1 library(ggplot2)
2 ggplot(data = trees, mapping = aes(y=Volume, x=0)) +
3 geom_boxplot(col="black", fill="lightblue", width=0.1) + coord_flip() + theme_light()
```

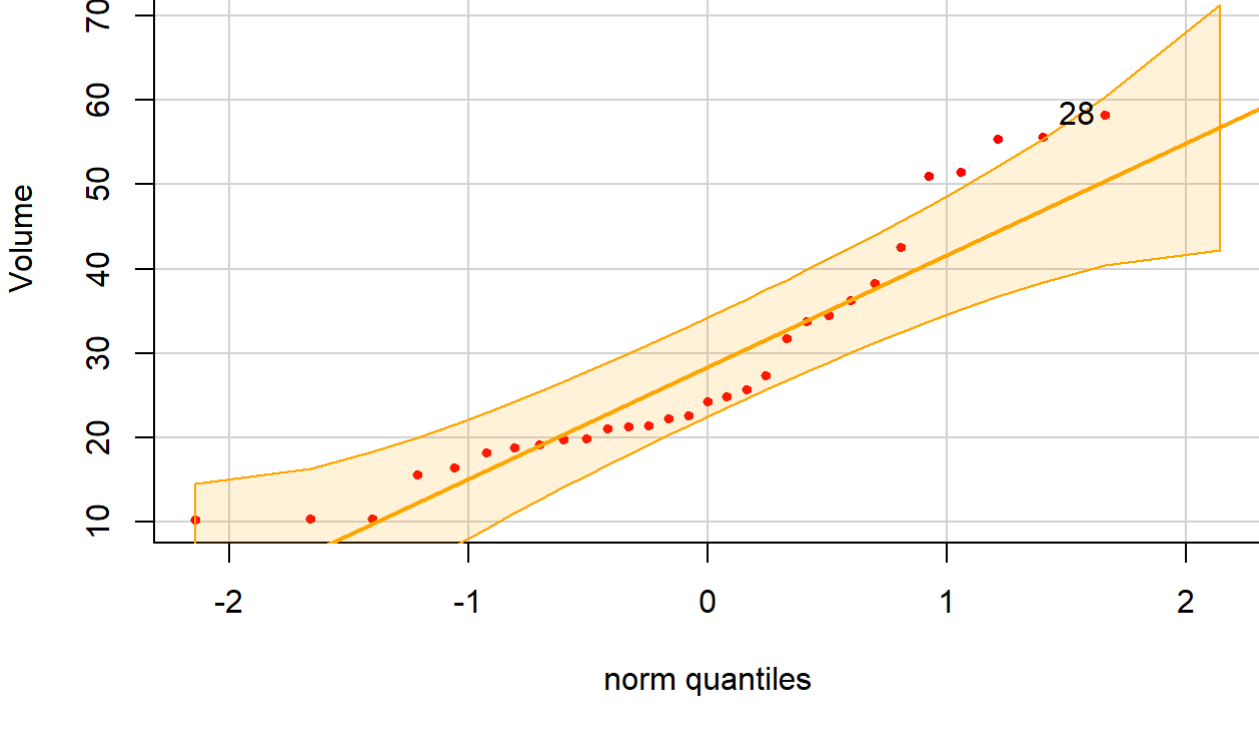


Q2) b.

```
1 library(car)
```

Loading required package: carData

```
1 qqPlot(trees$Volume, main="", ylab="Volume", cex=0.6, pch=19, col="red",
2 col.lines = "orange")
```



[1] 31 28

The Volume of trees does not follow normal distribution, as we can see that the points do not lie on the straight line and some points are even going out of the confidence band.

Shapiro-Wilk Test

```
1 shapiro.test(trees$Volume)
```

Shapiro-Wilk normality test

data: trees\$Volume
W = 0.88757, p-value = 0.003579

As we can see the p-value for the Shapiro test comes out to be 0.003579 which is less than 0.05(alpha), therefore we can say that trees\$Volume does not follow normal distribution.

Chi-Square Test

```
1 library(nortest)
2 pearson.test(trees$Volume)
```

Pearson chi-square normality test

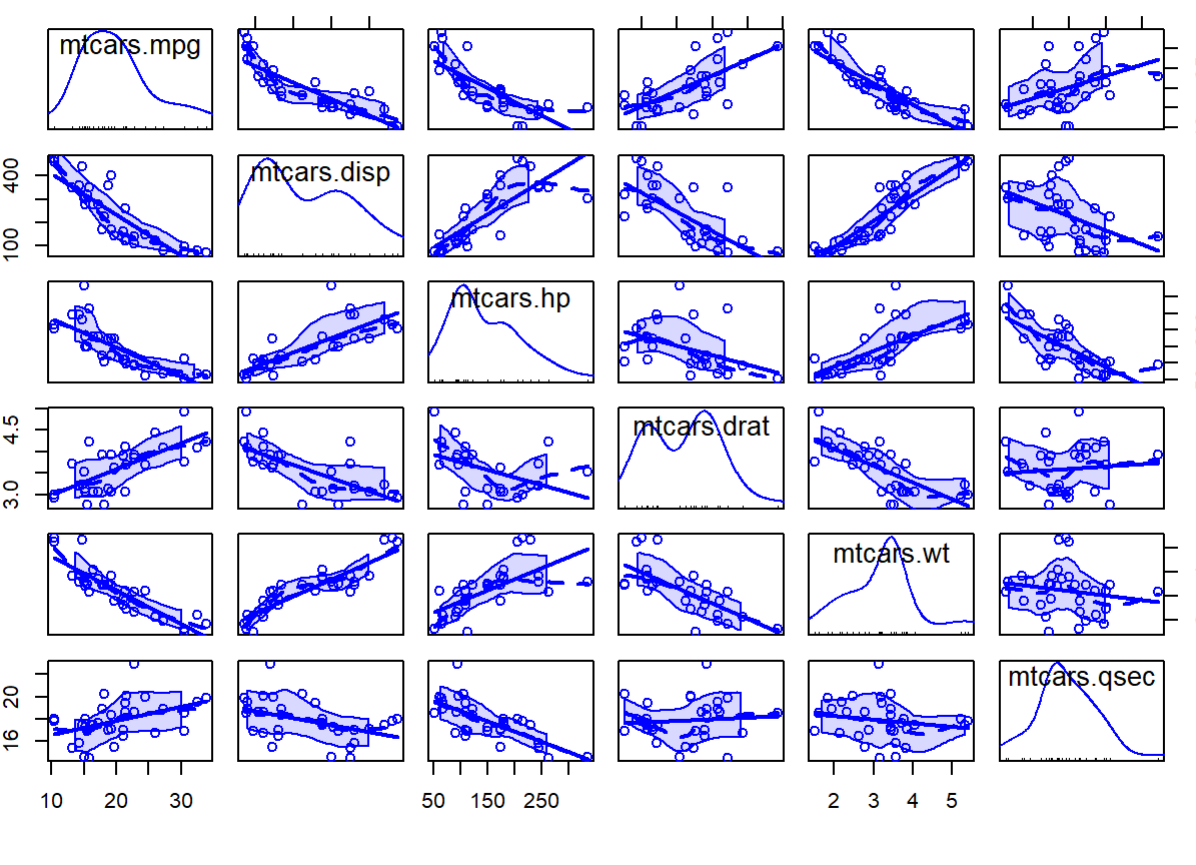
data: trees\$Volume
P = 15.194, p-value = 0.009567

The Chi-square goodness of fit test is a statistical hypothesis test used to determine whether a variable is likely to come from a specified distribution (in this case normal distribution) or not.

As we can see, the p-value of the Pearson Chi-Square test is less than 0.05, we can say that trees\$Volume does not follow normal distribution.

Q3) a.

```
1 scatterplotMatrix(~mtcars$mpg+mtcars$dis+mtcars$hp+mtcars$drat+mtcars$wt+mtcars$qsec,
```



By Looking at the first row of the scatter plot matrix, we can see that **mpg** is most **Positively associated** with **drat** and most **negatively associated** with **wt**.

Q3) b

```
1 data(mtcars)
2
3 selected_columns <- mtcars[, c("mpg", "disp", "hp", "drat", "wt", "qsec")]
4
5 correlation_matrix <- cor(selected_columns)
6
7 diag(correlation_matrix) <- NA
8
9 max_correlation <- max(correlation_matrix, na.rm = TRUE)
10
11 indices <- which(correlation_matrix == max_correlation, arr.ind = TRUE)
12
13 column1 <- rownames(correlation_matrix)[indices[1, 1]]
14 column2 <- colnames(correlation_matrix)[indices[1, 2]]
15
16 cat("Pair with maximum correlation:", column1, "and", column2, "with correlation =",
```

Pair with maximum correlation: wt and disp with correlation = 0.8879799

The Correlation between wt and disp is the highest with correlation = 0.8879799

Q4)

```
1 library(dplyr)
```

Attaching package: 'dplyr'

The following object is masked from 'package:car':

recode

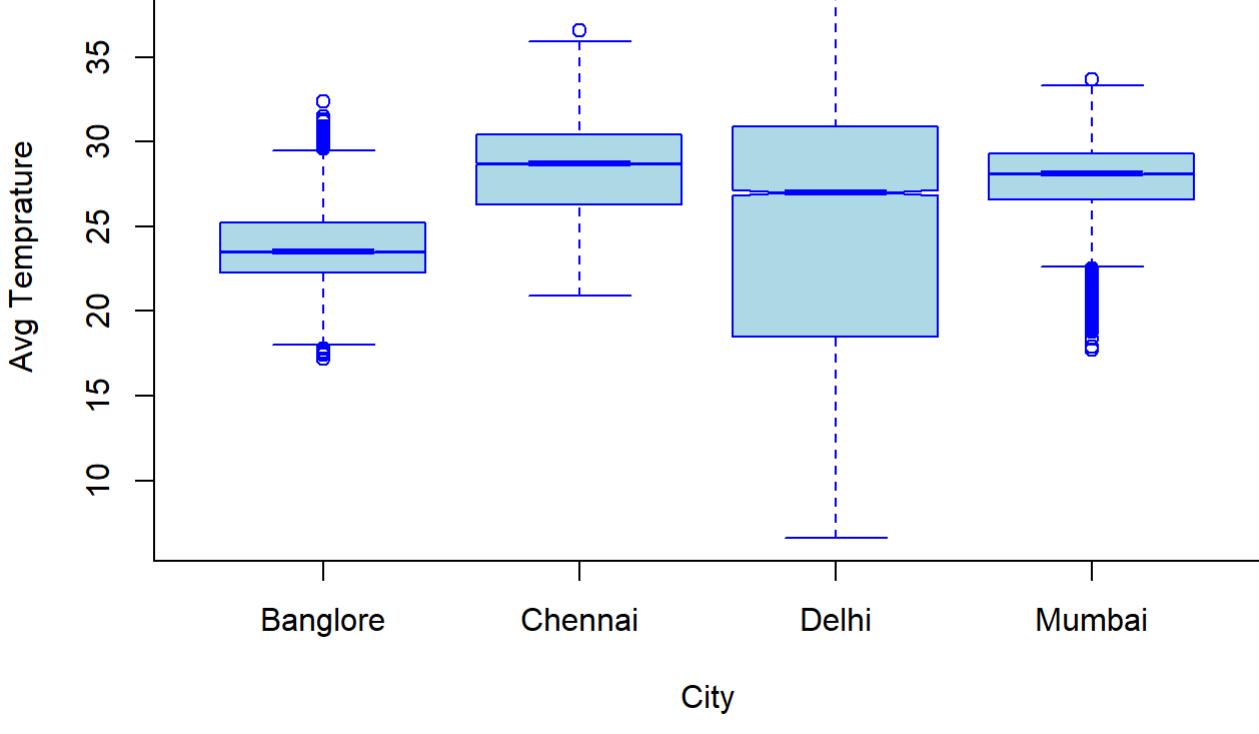
The following objects are masked from 'package:stats':

filter, lag

The following objects are masked from 'package:base':

intersect, setdiff, setequal, union

```
1 mumbai_data <- read.csv('Temperature_And_Precipitation_Cities_IN/Mumbai_1990_2022_S
2
3
4 chennai_data <- read.csv('Temperature_And_Precipitation_Cities_IN/Chennai_1990_2022
5
6 delhi_data <- read.csv('Temperature_And_Precipitation_Cities_IN/Delhi_NCR_1990_2022
7
8
9 banglore_data <- read.csv('Temperature_And_Precipitation_Cities_IN/Bangalore_1990_2
10
11
12 # Combine the datasets with a grouping variable
13
14 combined_data <- bind_rows(
15   mutate(mumbai_data, Group = "Mumbai"),
16   mutate(chennai_data, Group = "Chennai"),
17   mutate(delhi_data, Group = "Delhi"),
18   mutate(banglore_data, Group = "Banglore")
19 )
20
21
22 boxplot(combined_data$avg ~ combined_data$Group, xlab="City",
23   ylab="Avg Temperature", border = "blue", col="lightblue", notch=TRUE)
```



As we can see from th above visual, the median Average daily temperature for Bangalore, Chennai, Delhi and Mumbai over the years 1990 to 2022 is in between 23~30 degrees, also Bangalore and Mumbai are having a lot of outliers, i.e people in those cities had to go through extreme weather conditions compared to normal days.

Q5)

Some of the examples of unethical data visualization might be:

1) **Showing lots of variables in a pie chart**: The ideal number is in between 2 and 7, anything more than that will make the visual untidy and will over-populate the visual which make it difficult to interpret.

2) **Truncating the Y-axis** in the graphs: It will mislead or manipulate the users perception of data. Truncating the Y-axis refers to displaying a chart in a way that omits a portion of the Y-axis, making differences between data points appear more significant or less significant than they actually are.

3) **Unusual Coloring**: Good outcomes are associated with green color and bad outcomes with red, if we deviate from these then coloring is misleading and can create misunderstandings.

4) **Lack of Labels**: If charts are lacking labels, percentages, or any specific information about the categories then viewers cannot discern the exact values or proportions represented by each category.

5) **No Data Source**: Transparency is crucial in data visualization, and viewers need to know where the data comes from and how it was collected.

6) **Misleading Labels**: Using misleading labels or titles that don't accurately represent the data, its source, or its context.

7) **Data Manipulation**: Manipulating the underlying data before visualization, such as altering data values, outliers, or summary statistics to support a particular agenda.

8) **Overlaying Charts**: Overlaying multiple datasets with different units of measurement on a single chart without proper scaling or labeling.