

Statistical Learning Final Report

Alberto Calabrese, Eleonora Mesaglio, Greta d'Amore Grelli

2024-05-31

Contents

1	Introduction	1
2	Libraries	1
3	Data	2
3.1	Data Transformation	3
3.2	Data Cleaning	3
3.3	Rename Columns	4
4	Correlation Analysis	4
5	Data Visualization	6
5.1	Histograms	6
5.2	Pairplot	7
5.3	Barplot	9
5.3.1	Beverages Barplot	10
5.4	Boxplot	13
5.5	Scatterplot	14
6	Regression Analysis	17
6.1	Linear Regression	17
6.2	Logistic Regression	17

1 Introduction

Here Eleonora you can write the introduction of the project describing the scope and the data used.

Thank you Albi, I will. What is our project scope though?

I think that we have to analyze the dataset and perform some statistical analysis on it. We can start by calculating the correlation matrix and then we can visualize the data through histograms, pairplots, barplots and boxplots. Finally, we can perform a regression analysis.

2 Libraries

```
library(corrplot)
library(knitr)
```

3 Data

The dataset we will analyze in this project is *Starbucks Beverage Components* from Kaggle, that you can find at the following link: <https://www.kaggle.com/datasets/henryshan/starbucks>.

This data provides a comprehensive guide to the nutritional content of the beverages available on the Starbucks menu. We have a total of 242 samples described by 18 variables. These attributes include the name of the beverage, its categorization and preparation method, the total caloric content and the constituents of the beverage.

In the upcoming code lines, we import the dataset and generate a summary visualization. This initial step allows us to gain a better understanding of the data structure and the variables involved.

```
data <- read.csv("Data/starbucks.csv", header = TRUE, sep = ",")
```

```
# Overview of the data
summary(data)
```

```
## Beverage_category Beverage Beverage_prep Calories
## Length:242 Length:242 Length:242 Min. : 0.0
## Class :character Class :character Class :character 1st Qu.:120.0
## Mode :character Mode :character Mode :character Median :185.0
## Mean :193.9
## 3rd Qu.:260.0
## Max. :510.0
## Total.Fat..g. Trans.Fat..g. Saturated.Fat..g. Sodium..mg.
## Min. : 0.000 Min. :0.000 Min. :0.0000 Min. : 0.000
## 1st Qu.: 0.200 1st Qu.:0.100 1st Qu.:0.0000 1st Qu.: 0.000
## Median : 2.500 Median :0.500 Median :0.0000 Median : 5.000
## Mean : 2.905 Mean :1.307 Mean :0.0376 Mean : 6.364
## 3rd Qu.: 4.500 3rd Qu.:2.000 3rd Qu.:0.1000 3rd Qu.:10.000
## Max. :15.000 Max. :9.000 Max. :0.3000 Max. :40.000
## Total.Carbohydrates..g. Cholesterol..mg. Dietary.Fibre..g. Sugars..g.
## Min. : 0.0 Min. : 0.00 Min. :0.0000 Min. : 0.00
## 1st Qu.: 70.0 1st Qu.:21.00 1st Qu.:0.0000 1st Qu.:18.00
## Median :125.0 Median :34.00 Median :0.0000 Median :32.00
## Mean :128.9 Mean :35.99 Mean :0.8058 Mean :32.96
## 3rd Qu.:170.0 3rd Qu.:50.75 3rd Qu.:1.0000 3rd Qu.:43.75
## Max. :340.0 Max. :90.00 Max. :8.0000 Max. :84.00
## Protein..g. Vitamin.A....DV. Vitamin.C....DV. Calcium....DV.
## Min. : 0.000 Length:242 Length:242 Length:242
## 1st Qu.: 3.000 Class :character Class :character Class :character
## Median : 6.000 Mode :character Mode :character Mode :character
## Mean : 6.979
## 3rd Qu.:10.000
## Max. :20.000
## Iron....DV. Caffeine..mg.
## Length:242 Length:242
## Class :character Class :character
## Mode :character Mode :character
##
##
##
```

```
str(data)
```

```
## 'data.frame': 242 obs. of 18 variables:
```

```
## $ Beverage_category      : chr  "Coffee" "Coffee" "Coffee" "Coffee" ...
## $ Beverage               : chr  "Brewed Coffee" "Brewed Coffee" "Brewed Coffee" "Brewed Coffee" ...
## $ Beverage_prep          : chr  "Short" "Tall" "Grande" "Venti" ...
## $ Calories               : int   3 4 5 5 70 100 70 100 150 110 ...
## $ Total.Fat..g.          : num   0.1 0.1 0.1 0.1 0.1 3.5 2.5 0.2 6 4.5 ...
## $ Trans.Fat..g.          : num   0 0 0 0 0.1 2 0.4 0.2 3 0.5 ...
## $ Saturated.Fat..g.      : num   0 0 0 0 0 0.1 0 0 0.2 0 ...
## $ Sodium..mg.           : int   0 0 0 0 5 15 0 5 25 0 ...
## $ Total.Carbohydrates..g.: int   5 10 10 10 75 85 65 120 135 105 ...
## $ Cholesterol..mg.       : int   0 0 0 0 10 10 6 15 15 10 ...
## $ Dietary.Fibre..g.      : int   0 0 0 0 0 0 1 0 0 1 ...
## $ Sugars..g.             : int   0 0 0 0 9 9 4 14 14 6 ...
## $ Protein..g.            : num   0.3 0.5 1 1 6 6 5 10 10 8 ...
## $ Vitamin.A....DV.       : chr   "0%" "0%" "0%" "0%" ...
## $ Vitamin.C....DV.       : chr   "0%" "0%" "0%" "0%" ...
## $ Calcium....DV.         : chr   "0%" "0%" "0%" "2%" ...
## $ Iron....DV.            : chr   "0%" "0%" "0%" "0%" ...
## $ Caffeine..mg.          : chr   "175" "260" "330" "410" ...
```

3.1 Data Transformation

Note that several variables in our dataset, namely “Vitamin.A...DV.”, “Vitamin.C...DV.”, “Calcium...DV.” and “Iron...DV.”, are represented as percentages. Consequently, the percentage symbol is included in our data. However, when conducting statistical analysis using R, the presence of non-numeric characters such as the percentage symbol can cause complications, interfering with the processing and analysis of the data. Therefore, we proceeded to remove it.

Similarly, as R primarily operates on numeric and categorical data, we also convert all the other numerical variables into numeric format.

These preprocessing steps ensure a smooth and efficient analysis, making it easier to explore, visualize, and understand our data.

```
# Remove percentage sign from the data
data$Vitamin.C....DV. <- as.numeric(gsub("%", "", data$Vitamin.C....DV.))
data$Calcium....DV. <- as.numeric(gsub("%", "", data$Calcium....DV.))
data$Iron....DV. <- as.numeric(gsub("%", "", data$Iron....DV.))
data$Vitamin.A....DV. <- as.numeric(gsub("%", "", data$Vitamin.A....DV.))

# Set the other variables as numeric
data$Calories <- as.numeric(data$Calories)
data$Trans.Fat..g. <- as.numeric(data$Trans.Fat..g.)
data$Total.Fat..g. <- as.numeric(data$Total.Fat..g.)
data$Cholesterol..mg. <- as.numeric(data$Cholesterol..mg.)
data$Sodium..mg. <- as.numeric(data$Sodium..mg.)
data$Total.Carbohydrates..g. <- as.numeric(data$Total.Carbohydrates..g.)
data$Dietary.Fibre..g. <- as.numeric(data$Dietary.Fibre..g.)
data$Sugars..g. <- as.numeric(data$Sugars..g.)
data$Caffeine..mg. <- as.numeric(data$Caffeine..mg.)
```

3.2 Data Cleaning

Another challenge we have to face is the presence of missing data. Indeed, in “Caffeine..mg.” column there are some NA values. This is a common issue in data analysis and needs to be addressed appropriately to ensure the validity of our statistical results.

There are different approaches to deal with this issue.

One way to handle this issue is removing the observations with NA values.

A different approach is to fill in the NA values with the average or the median of the observed values for that specific attribute. This method helps to preserve the overall data distribution while addressing the missing data points.

In our case, we decided to impute the NA values in the “Caffeine..mg.” column with the median of the observed values for that attribute. This choice is suitable because we have observed that the data is skewed, and the median is a more robust measure of central tendency in the presence of outliers.

```
data_cleaned <- data
data_cleaned$Caffeine..mg.[is.na(data_cleaned$Caffeine..mg.)] <- median(
  data_cleaned$Caffeine..mg., na.rm = TRUE)

summary(data_cleaned$Caffeine..mg.)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      0.00   70.00   75.00   88.14  130.00  410.00
```

3.3 Rename Columns

Lastly, taking in consideration our cleaned data, we rename the columns by removing dots and units of measure, in order to obtain a more readable dataset.

```
colnames(data_cleaned) <- c("Beverage_category", "Beverage",
  "Beverage_prep", "Calories",
  "Total_Fat", "Trans_Fat",
  "Saturated_Fat", "Sodium",
  "Total_Carbohydrates", "Cholesterol",
  "Dietary_Fibre", "Sugars",
  "Protein", "Vitamin_A",
  "Vitamin_C", "Calcium",
  "Iron", "Caffeine")
```

4 Correlation Analysis

After completing these preliminary preprocessing steps, we calculate the correlation matrix for our dataset. This computation helps us in comprehending the interrelationships among the dataset’s variables. In the correlation matrix, a value near to 1 at the ij position indicates a strong positive correlation between the i -th and j -th variables. Conversely, a value close to -1 signifies a strong negative correlation. A value near 0 suggests that the two variables do not significantly influence each other.

Observe that the first three columns of our data are categorical features, therefore for these we cannot compute Pearson’s correlation coefficient. In the following code lines we remove them to compute and plot such matrix.

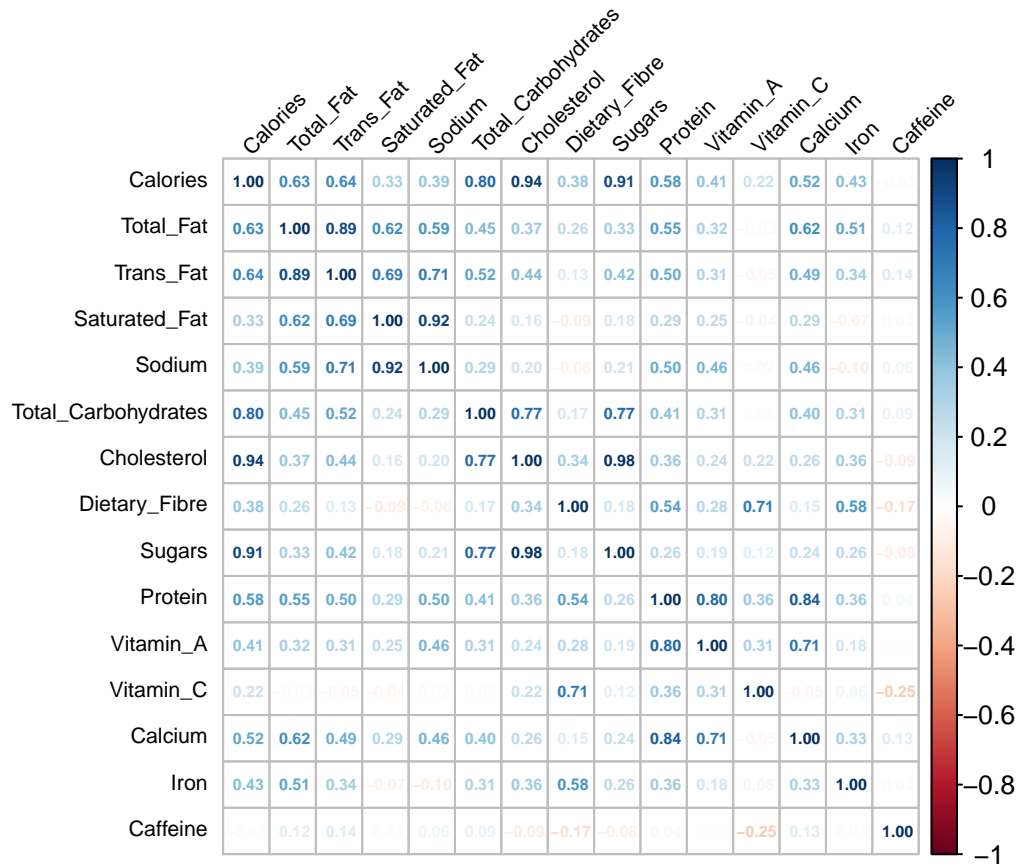
```
# Remove first 3 columns for the correlation matrix since they are categorical
data_num <- data_cleaned[, -c(1:3)]

# Calculate the correlation matrix

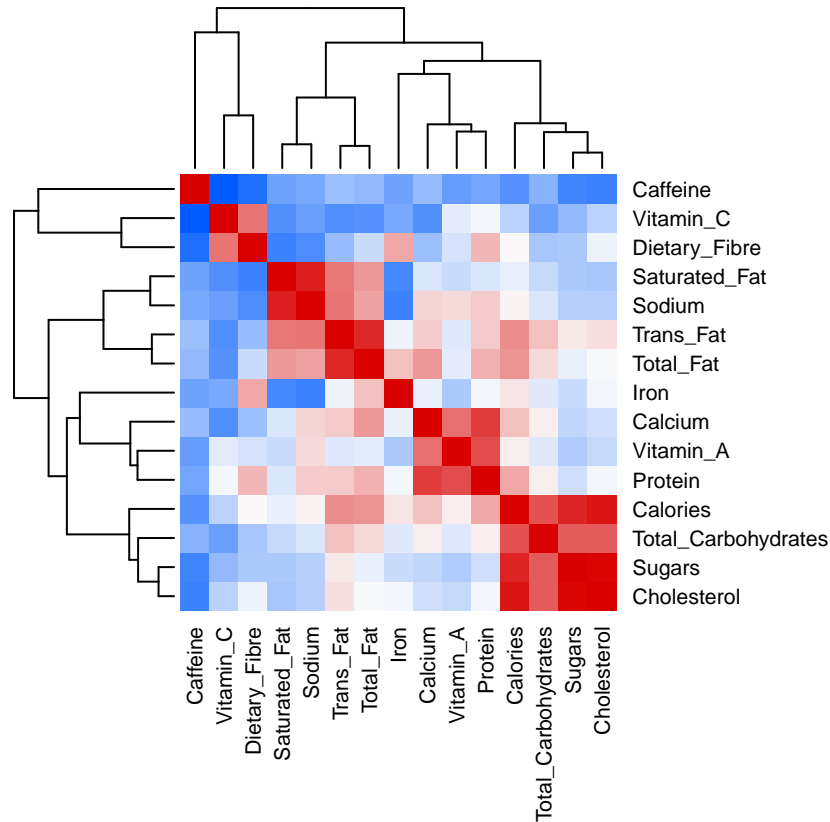
correlation_matrix <- cor(data_num)

# Plot the correlation matrix using corrplot
```

```
corrplot(correlation_matrix, method = "number", tl.col = "black",
         tl.srt = 45, addCoef.col = "black", number.cex = 0.5, tl.cex = 0.7)
```



```
# Heatmap of the correlation matrix
heatmap(cor(data_num),
       col = colorRampPalette(c("#005cfc", "#fbfbfb", "#d90000"))(100),
       symm = TRUE,
       margins = c(8, 8),
       cexRow = 0.8,
       cexCol = 0.8)
```



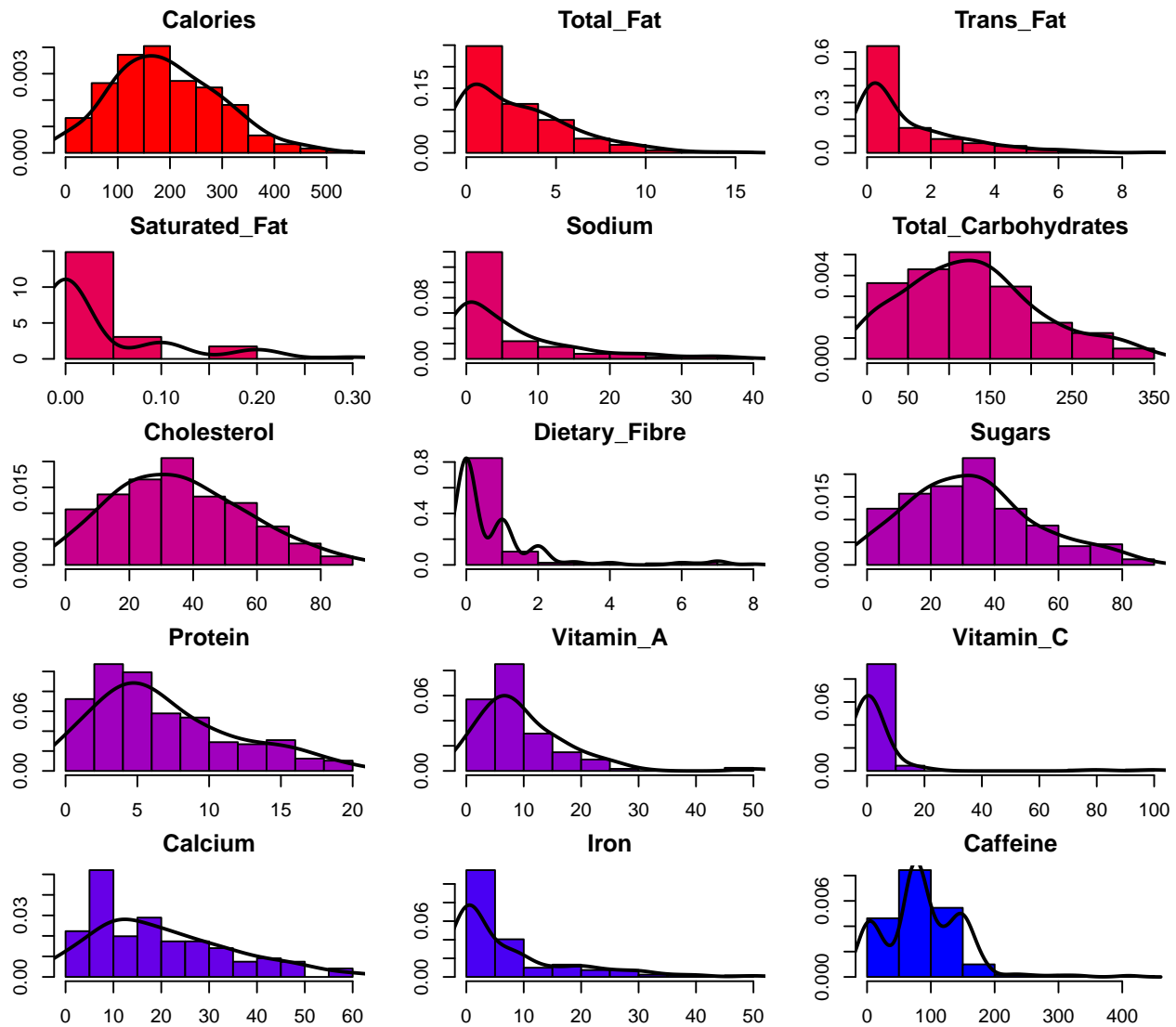
Moreover, we visualized the correlation matrix through a heatmap. The heatmap provides a visual representation of the correlation matrix, making it easier to identify patterns and relationships between the variables. The color gradient helps to distinguish between positive and negative correlations, with darker shades indicating stronger correlations.

5 Data Visualization

5.1 Histograms

We will plot some histograms to visualize the data.

```
# Histogram of the data with density distribution
par(mfrow = c(5, 3), mar = c(2, 2, 2, 2))
col <- c('#ff0000', '#f70028', '#ee0040', '#e50055', '#dc0069',
        '#d2007b', '#c7008d', '#bb009e', '#ae00ae', '#a000be',
        '#8f00cc', '#7d00da', '#6700e7', '#4900f3', '#0000ff')
for (i in 1:ncol(data_num)) {
  hist(data_num[, i], main = colnames(data_num)[i],
        xlab = colnames(data_num)[i], col = col[i], freq = FALSE)
  dens <- density(data_num[, i], na.rm=TRUE, adjust=1.25)
  lines(dens, col = "black", lwd = 2)
}
```



ADD COMMENTS ON THE GRAPH

5.2 Pairplot

We will plot a pairplot to visualize the relationship between the variables. The pairplot is a grid of scatterplots that shows the relationship between each pair of variables in the dataset. This visualization helps us to identify patterns and correlations between the variables.

First of all we have to define the function for the pairplot. We will define a function for the histogram, the correlation and the smooth line.

```
# Histogram function
panel.hist <- function(x, colour, ...)
{
  usr <- par("usr")
  on.exit(par(usr))
  par(usr = c(usr[1:2], 0, 1.5))
  h <- hist(x, plot = FALSE)
  breaks <- h$breaks
  nB <- length(breaks)
```

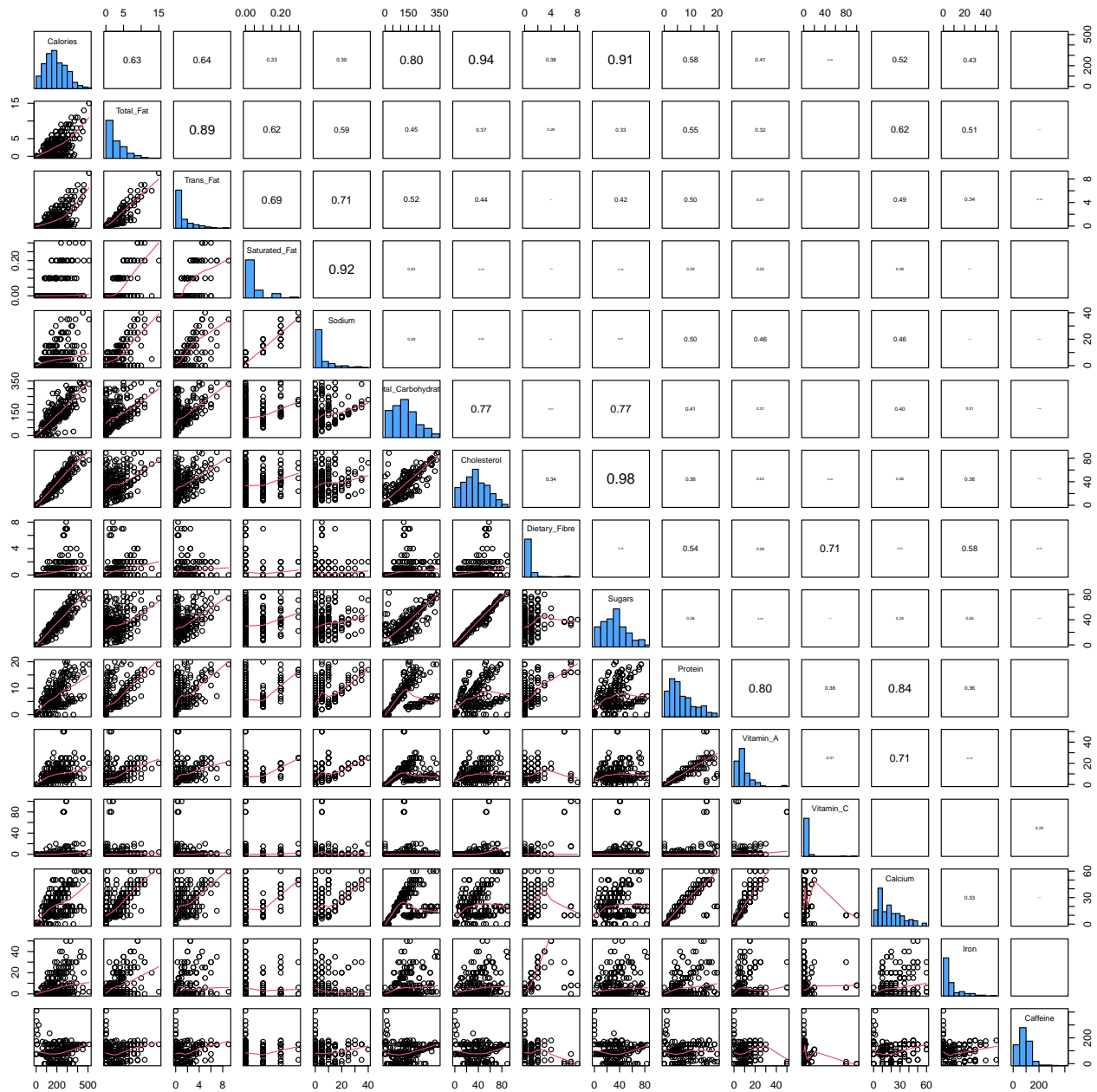
```

y <- h$counts
y <- y / max(y)
rect(breaks[-nB], 0, breaks[-1], y, col = colour, ...)
}

# Correlations function
panel.cor <- function(x, y, digits = 2, prefix = "", cex.cor, ...)
{
  usr <- par("usr")
  on.exit(par(usr))
  par(usr = c(0, 1, 0, 1))
  r <- abs(cor(x, y))
  txt <- format(c(r, 0.123456789), digits = digits)[1]
  txt <- paste0(prefix, txt)
  if(missing(cex.cor)) cex.cor <- 0.5/strwidth(txt)
  text(0.5, 0.5, txt, cex = cex.cor * r)
}

pairs(data_num,
      diag.panel = panel.hist,
      upper.panel = panel.cor,
      lower.panel = panel.smooth,
      colour = "#4ea5ff")

```

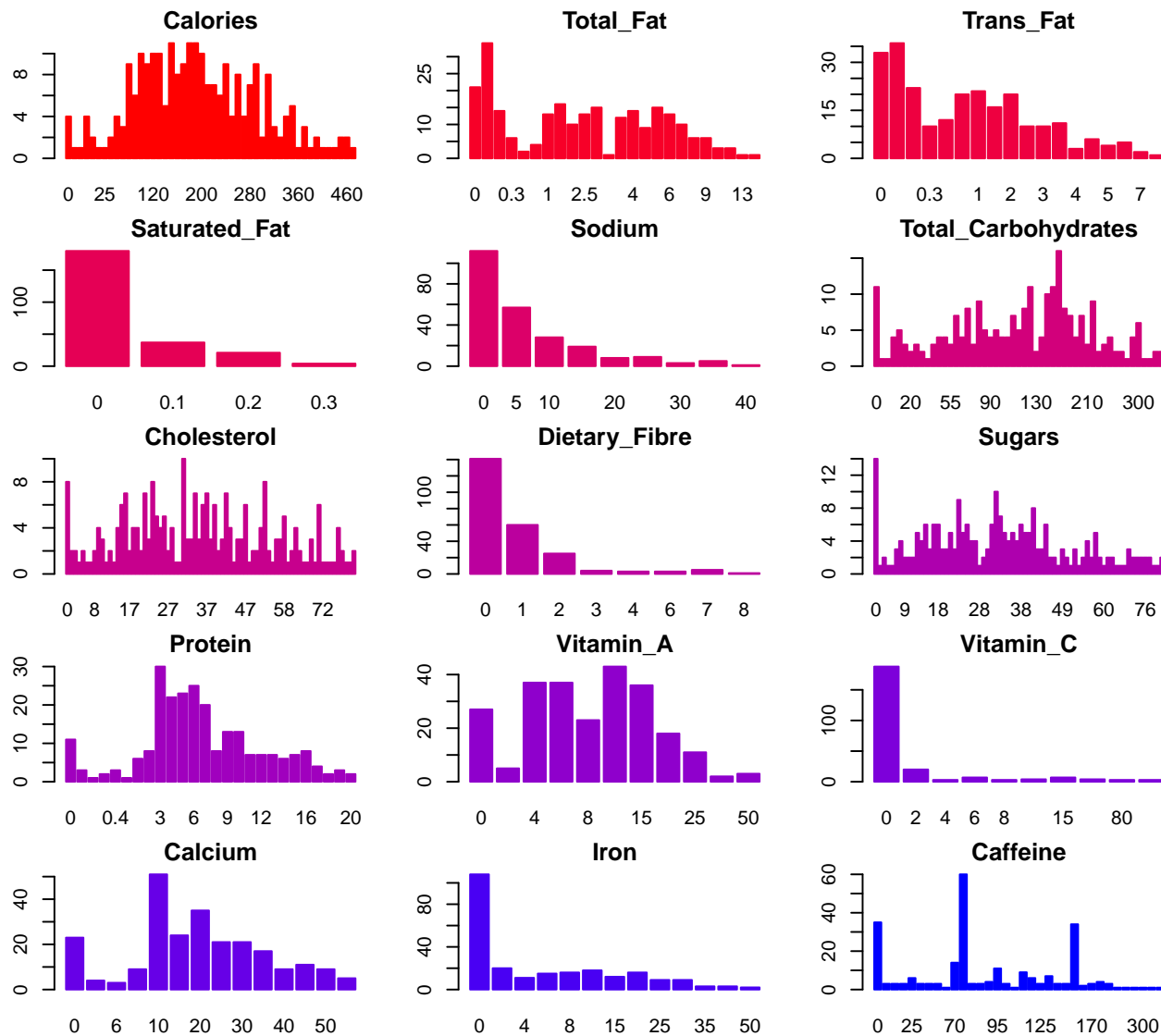



ADD COMMENTS ON THE GRAPH

5.3 Barplot

We will plot a barplot of the data. The barplot is a graphical representation of the data that displays the frequency of each category in a categorical variable. This visualization helps us to understand the distribution of the data and identify the most common categories in the dataset.

```
# Barplot of the data
par(mfrow = c(5, 3), mar = c(2, 2, 2, 2))
for (i in 1:ncol(data_num)) {
  barplot(table(data_num[, i]), main = colnames(data_num)[i],
    xlab = colnames(data_num)[i], col = col[i], border = col[i])
}
```



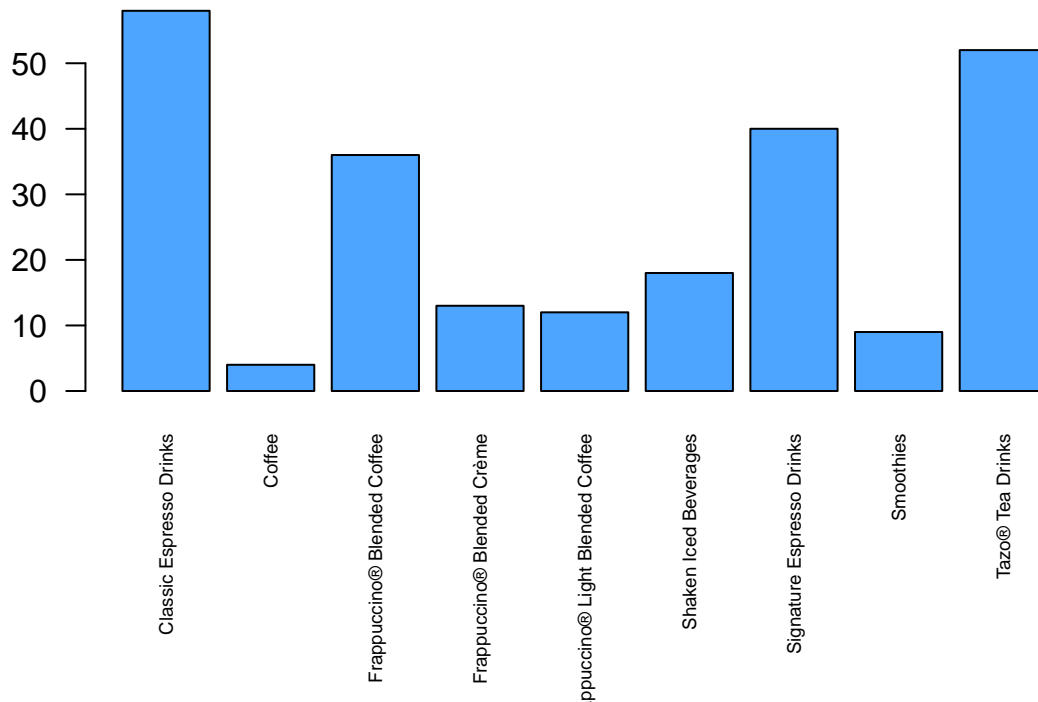
ADD COMMENTS ON THE GRAPH

5.3.1 Beverages Barplot

We create a barplot to visualize the distribution of the 'Beverage_category' variable and the 'Beverage_prep' variable in order to understand the most common beverages and preparation methods.

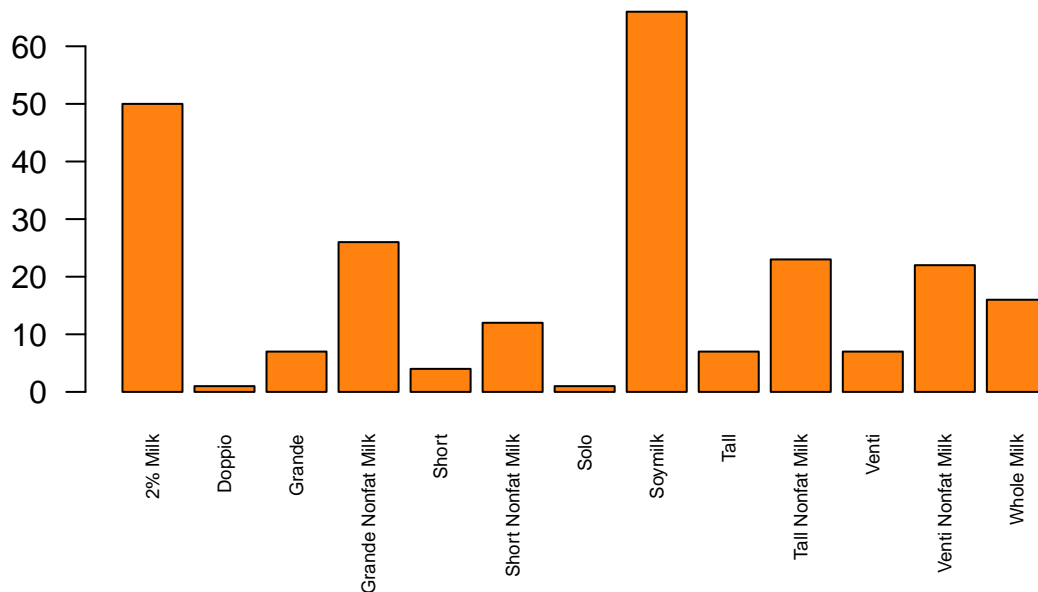
```
# Beverage category
par(mfrow = c(1, 1), mar = c(8, 2, 2, 2))
barplot(table(data$Beverage_category),
        main = "Distribution of Beverage Categories",
        ylab = "Count",
        col = "#4ea5ff",
        las = 2,
        cex.names = 0.6)
```

Distribution of Beverage Categories



```
# Beverage preparation
barplot(table(data$Beverage_prep),
  main = "Distribution of Beverage Preparation",
  ylab = "Count",
  col = "#ff810f",
  las = 2,
  cex.names = 0.6)
```

Distribution of Beverage Preparation



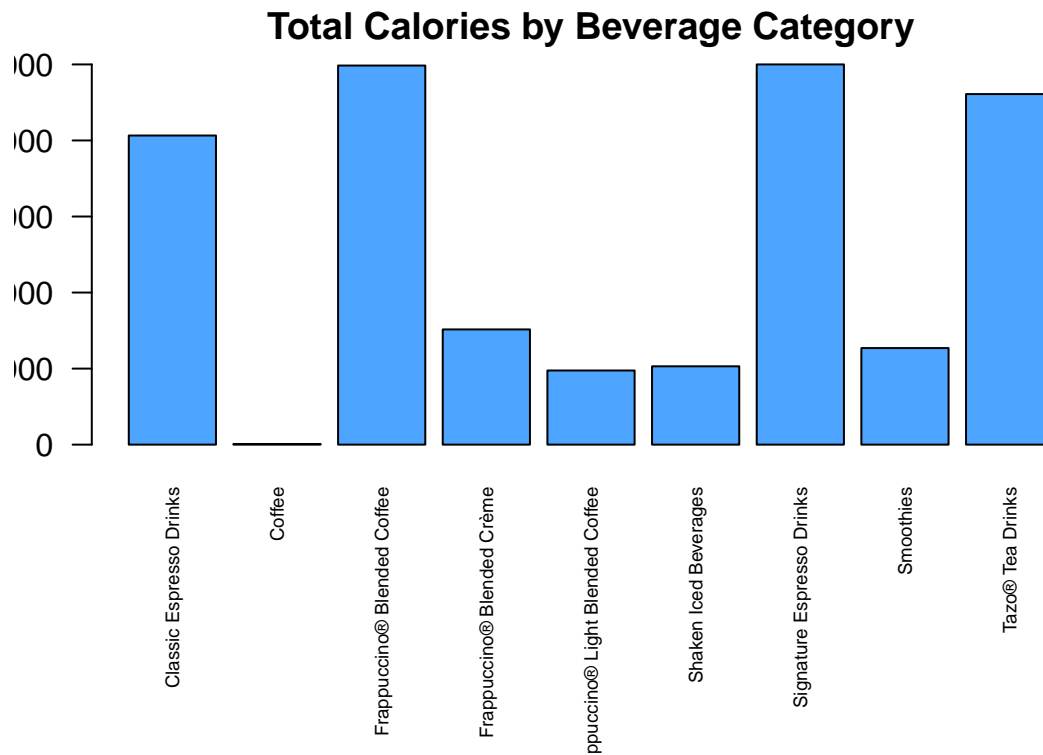
Now we want to compare the total calories for each categories of bevarage. First we agggrate the data to obtain the total calories for each categories of bevarage and then we create a barplot to visualize the results.

```

par(mfrow = c(1, 1), mar = c(8, 2, 2, 2))
total_calories_by_category <- aggregate(Calories ~ Beverage_category,
                                         data = data_cleaned, sum)

barplot(height = total_calories_by_category$Calories,
        names.arg = total_calories_by_category$Beverage_category,
        main = "Total Calories by Beverage Category",
        ylab = "Total Calories",
        col = "#4ea5ff",
        las = 2,
        cex.names = 0.6)

```



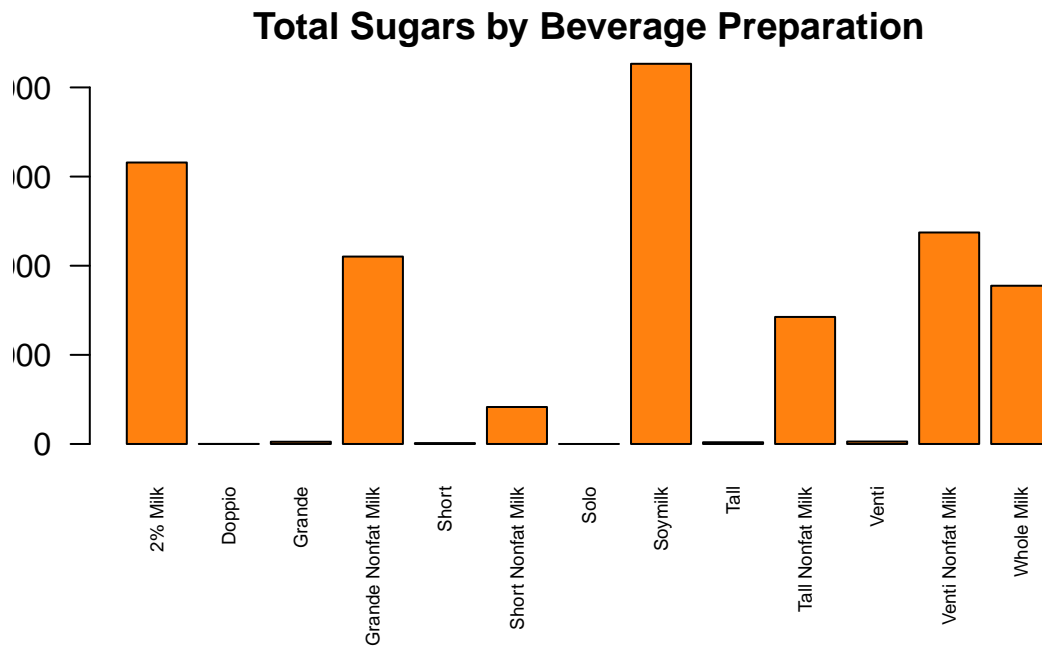
Now we want to compare the total sugars for each preparation of beverage. First we aggregate the data to obtain the total sugars for each preparation of beverage and then we create a barplot to visualize the results.

```

par(mfrow = c(1, 1), mar = c(8, 2, 2, 2))
total_sugar_by_prep <- aggregate(Total_Carbohydrates ~ Beverage_prep,
                                  data = data_cleaned, sum)

barplot(height = total_sugar_by_prep$Total_Carbohydrates,
        names.arg = total_sugar_by_prep$Beverage_prep,
        main = "Total Sugars by Beverage Preparation",
        ylab = "Total Sugars (g)",
        col = "#ff810f",
        las = 2,
        cex.names = 0.6)

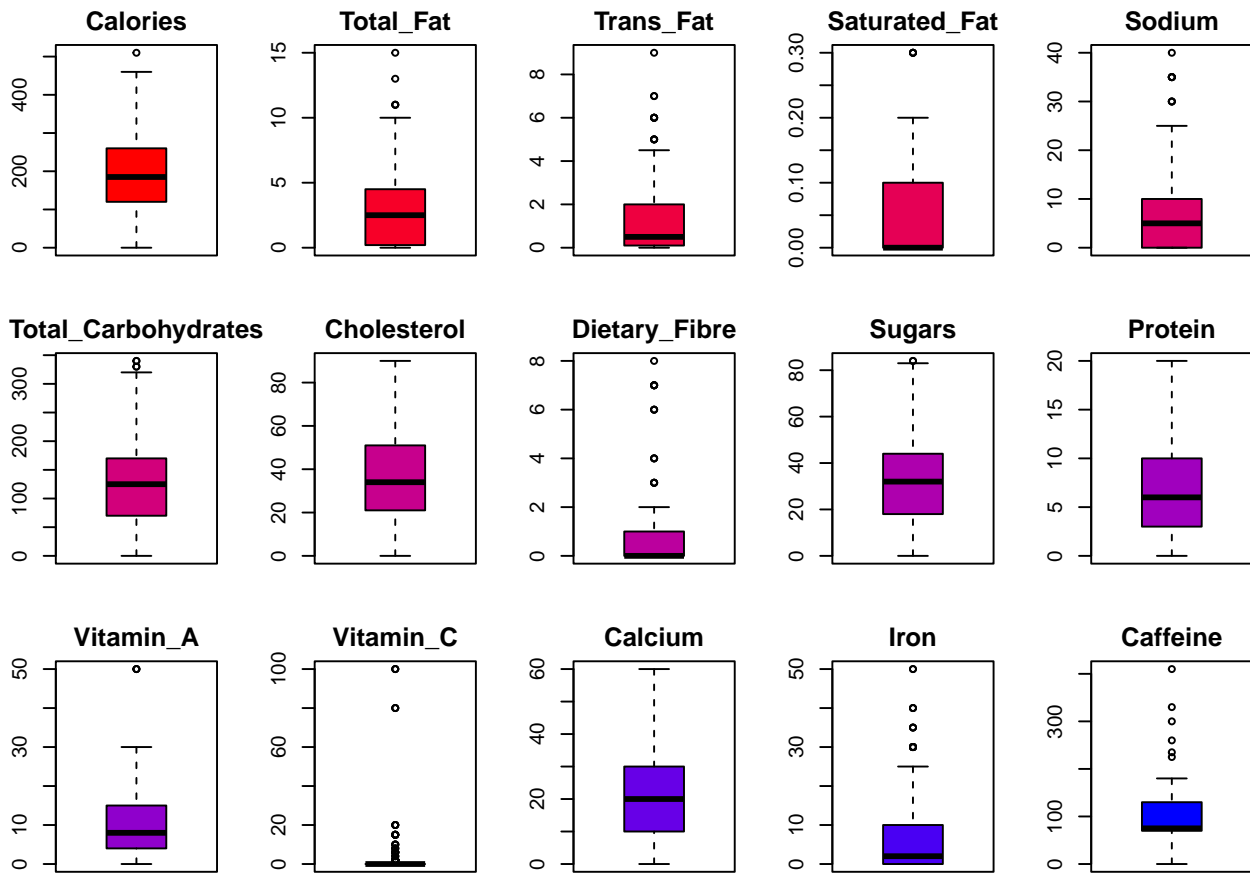
```



5.4 Boxplot

We will plot a boxplot of the data. The boxplot is a graphical representation of the data that displays the distribution of the data, including the median, quartiles, and outliers. This visualization helps us to identify the spread and variability of the data.

```
# Boxplot of the data
par(mfrow = c(3, 5), mar = c(2, 2, 2, 2))
for (i in 1:ncol(data_num)) {
  boxplot(data_num[, i], main = colnames(data_num)[i],
          xlab = colnames(data_num)[i], col = col[i])
}
```



5.5 Scatterplot

We will plot a scatterplot of the data. The scatterplot is a graphical representation of the data that displays the relationship between two variables. This visualization helps us to identify patterns and correlations between the variables.

We create a scatterplot to compare the amounts of calories and fat for each categories of bevarage. We assign distinct colors to each beverage category and create a legend to identify each category.

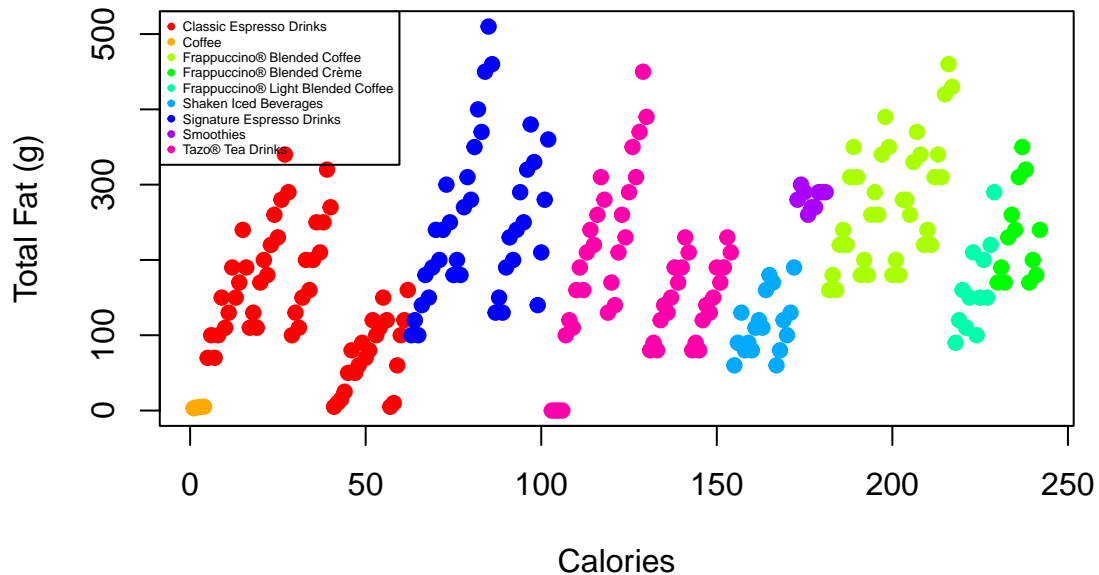
```
# Set the variable as factor
data_cleaned$Beverage_category <- as.factor(data_cleaned$Beverage_category)

# Assign distinct colors to each beverage category
colors <- rainbow(length(unique(data_cleaned$Beverage_category)))
color_map <- setNames(colors, levels(data_cleaned$Beverage_category))

# Create a scatterplot to compare amounts of calories and fat
# for each categories of bevarage
par(mfrow = c(1, 1))
plot(data_cleaned$Calories,
     data_cleaned$Total_Fat_g,
     col = color_map[data_cleaned$Beverage_category],
     pch = 19,
     xlab = "Calories",
     ylab = "Total Fat (g)",
     main = "Calories vs Total Fat")
```

```
# Legend
legend("topleft", legend = levels(data_cleaned$Beverage_category),
      col = colors, cex = 0.4, pch = 19)
```

Calories vs Total Fat

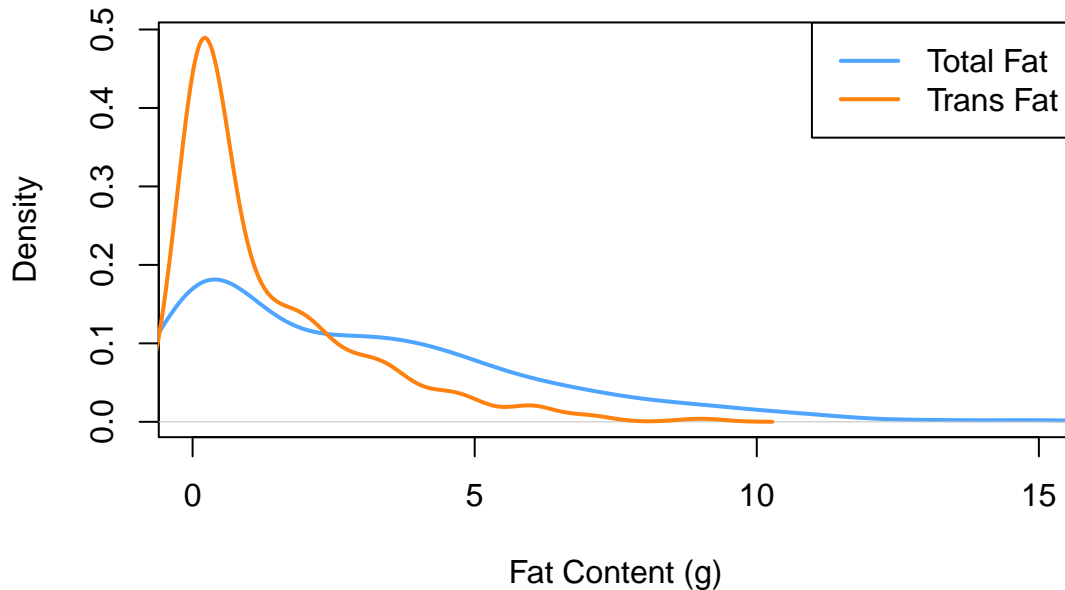


```
# Comparision between total fat and trans fat ( che cazzo sono?)
```

```
# Numeric variable -> calculate density
total_fat_density <- density(data_cleaned$Total_Fat)
trans_fat_density <- density(data_cleaned$Trans_Fat)

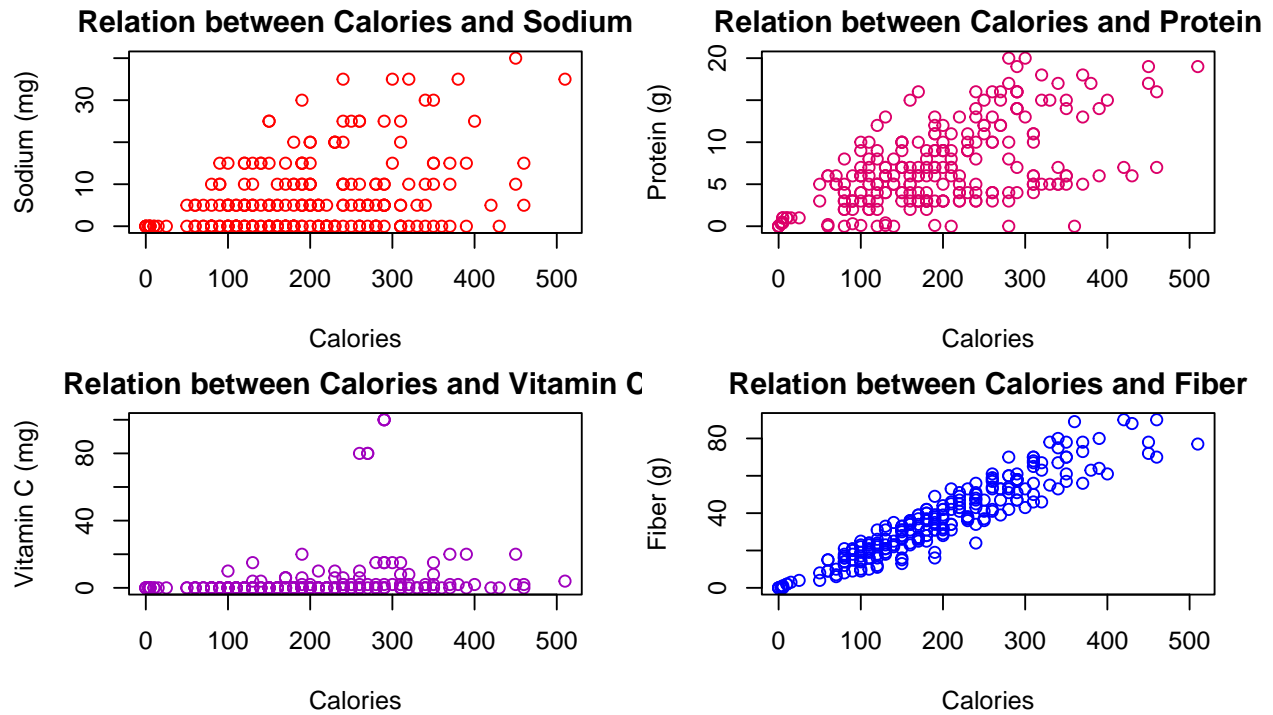
plot(total_fat_density, col = "#4ea5ff",
     main = "Comparison of Total Fat and Trans Fat Distributions",
     xlab = "Fat Content (g)", ylab = "Density",
     ylim = c(0, max(total_fat_density$y, trans_fat_density$y)),
     xlim = range(data_cleaned$Total_Fat, data_cleaned$Trans_Fat),
     lwd = 2, lty = 1)
lines(trans_fat_density, col = "#ff810f", lwd = 2, lty = 1)
legend("topright", legend = c("Total Fat", "Trans Fat"),
      col = c("#4ea5ff", "#ff810f"), lwd = 2, lty = 1)
```

Comparison of Total Fat and Trans Fat Distributions



Create scatterplot to look into relationship between calories and other variables. We will plot the relationship between calories and sodium, protein, vitamin C and fiber.

```
par(mfrow = c(2, 2), mar = c(4, 4, 2, 2))
with(data_cleaned, {
  plot(Calories, Sodium , main = "Relation between Calories and Sodium",
       xlab = "Calories", ylab = "Sodium (mg)", col = col[1])
  plot(Calories, Protein , main = "Relation between Calories and Protein",
       xlab = "Calories", ylab = "Protein (g)", col = col[5])
  plot(Calories, Vitamin_C , main = "Relation between Calories and Vitamin C",
       xlab = "Calories", ylab = "Vitamin C (mg)", col = col[10])
  plot(Calories, Cholesterol , main = "Relation between Calories and Fiber",
       xlab = "Calories", ylab = "Fiber (g)", col = col[15])
})
```

There's increase in every feature with increase in calories. Features like proteins and fiber rapidly increase, instead vitamin and cholesterol more flat growing. Confirmed by correlation coefficients

ADD COMMENTS ON THE GRAPH

6 Regression Analysis

6.1 Linear Regression

6.2 Logistic Regression