

# Obesity

## Title: Obesity Level Predictive Modeling

*Details Dataset: Dataset: Estimation of obesity Level*

Source:

<https://archive.ics.uci.edu/dataset/544/estimation+of+obesity+levels+based+on+eating+habits+and+physical+condition>

This dataset includes data for the estimation of obesity levels in individuals from the countries of Mexico, Peru and Colombia, based on their eating habits and physical condition. It consists of 17 attributes and 2111 records, the records are labeled with the class variable NObesity (Obesity Level), that allows classification of the data using the values of Insufficient Weight, Normal Weight, Overweight Level I, Overweight Level II, Obesity Type I, Obesity Type II and Obesity Type III. 77% of the data was generated synthetically using the Weka tool and the SMOTE filter, 23% of the data was collected directly from users through a web platform.

## Introduction

Obesity is a global health challenge with significant implications for individuals and society. As the prevalence of obesity continues to rise, understanding the factors contributing to obesity and developing effective predictive models are crucial for preventive healthcare interventions. Predictive modeling in the context of obesity aims to anticipate and identify individuals at risk, enabling timely interventions and personalized healthcare strategies.

## Initialization

Import necessary libraries

```
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

library(plotly)
```

```

## Warning: package 'plotly' was built under R version 4.4.2
## Loading required package: ggplot2
## Warning: package 'ggplot2' was built under R version 4.4.3
##
## Attaching package: 'plotly'
## The following object is masked from 'package:ggplot2':
##
##   last_plot
## The following object is masked from 'package:stats':
##
##   filter
## The following object is masked from 'package:graphics':
##
##   layout

library(ggplot2)
library(tidyr)
library(ggcorrplot)

## Warning: package 'ggcorrplot' was built under R version 4.4.3

library(e1071)

## Warning: package 'e1071' was built under R version 4.4.3

library(caTools)

## Warning: package 'caTools' was built under R version 4.4.3

library(tidyverse)

## — Attaching core tidyverse packages ————— tidyverse
2.0.0 —
## ✓ forcats   1.0.0      ✓ readr      2.1.5
## ✓ lubridate 1.9.3      ✓ stringr    1.5.1
## ✓ purrr     1.0.2      ✓ tibble     3.2.1

## — Conflicts —————
tidyverse_conflicts() —
## ✗ plotly::filter() masks dplyr::filter(), stats::filter()
## ✗ dplyr::lag()     masks stats::lag()
## ⓘ Use the conflicted package (<http://conflicted.r-lib.org/>) to force all
conflicts to become errors

library(caret)

## Warning: package 'caret' was built under R version 4.4.3

```

```
## Loading required package: lattice
##
## Attaching package: 'caret'
##
## The following object is masked from 'package:purrr':
##
##     lift
```

## Data Ingestion

Load dataset of Obesity Level as dataframe

```
url <- "https://docs.google.com/spreadsheets/d/e/2PACX-
1vQv7ETe9H0ySeTirE9z67a1X2nGMozFPqYvNxVwXc6-tx3IXX-
Ez0LGppCzoFvTLz6b7NQg_vGA1PLA/pub?output=csv"

# Read the CSV file
obesityDF <- read.csv(url, stringsAsFactors = FALSE)
```

## Data Understanding

Check variables that attribute to the dataset

```
head(obesityDF)

##   Gender Age Height Weight family_history_with_overweight FAVC FCVC NCP
## 1 Female  21   1.62   64.0                      yes     no    2    3
## 2 Female  21   1.52   56.0                      yes     no    3    3
## 3  Male   23   1.80   77.0                      yes     no    2    3
## 4  Male   27   1.80   87.0                      no      no    3    3
## 5  Male   22   1.78   89.8                      no      no    2    1
## 6  Male   29   1.62   53.0                      no     yes    2    3
##      CAEC SMOKE CH20 SCC FAF TUE      CALC      MTRANS
## 1 Sometimes    no    2  no    0    1          no Public_Transportation
## 2 Sometimes   yes    3 yes    3    0 Sometimes Public_Transportation
## 3 Sometimes    no    2  no    2    1 Frequently Public_Transportation
## 4 Sometimes    no    2  no    2    0 Frequently           Walking
## 5 Sometimes    no    2  no    0    0 Sometimes Public_Transportation
## 6 Sometimes    no    2  no    0    0 Sometimes           Automobile
##      NObeyesdad
## 1      Normal_Weight
## 2      Normal_Weight
## 3      Normal_Weight
## 4 Overweight_Level_I
## 5 Overweight_Level_II
## 6      Normal_Weight
```

## Data Preprocessing

1. Check for any NA values

```
colSums(is.na(obesityDF))
```

```
##           Gender           Age
##           0           0
##           Height          Weight
##           0           0
## family_history_with_overweight FAVC
##           0           0
##           FCVC           NCP
##           0           0
##           CAEC           SMOKE
##           0           0
##           CH20           SCC
##           0           0
##           FAF           TUE
##           0           0
##           CALC           MTRANS
##           0           0
##           NObeyesdad
##           0
```

## 2. Do data cleaning process

- Round off the values of age variable from numeric to integer
- Load into new dataframe with age integer

```
obesityDF$Age <- round(obesityDF$Age)
age_obesity <- obesityDF
```

- Check what the dataframe is about

*#not necessary but if want to see what the changes of column name, can use this*

```
str(age_obesity)
```

```
## 'data.frame':    2111 obs. of  17 variables:
## $ Gender          : chr  "Female" "Female" "Male" "Male"
## ...
## $ Age             : num  21 21 23 27 22 29 23 22 24 22 ...
## $ Height          : num  1.62 1.52 1.8 1.8 1.78 1.62 1.5
1.64 1.78 1.72 ...
## $ Weight          : num  64 56 77 87 89.8 53 55 53 64 68
## ...
## $ family_history_with_overweight: chr  "yes" "yes" "yes" "no" ...
## $ FAVC            : chr  "no" "no" "no" "no" ...
## $ FCVC            : num  2 3 2 3 2 2 3 2 3 2 ...
## $ NCP             : num  3 3 3 3 1 3 3 3 3 3 ...
## $ CAEC            : chr  "Sometimes" "Sometimes"
"Sometimes" "Sometimes" ...
## $ SMOKE           : chr  "no" "yes" "no" "no" ...
## $ CH20            : num  2 3 2 2 2 2 2 2 2 2 ...
## $ SCC             : chr  "no" "yes" "no" "no" ...
```

```
## $ FAF : num 0 3 2 2 0 0 1 3 1 1 ...
## $ TUE : num 1 0 1 0 0 0 0 0 1 1 ...
## $ CALC : chr "no" "Sometimes" "Frequently"
"Frequently" ...
## $ MTRANS : chr "Public_Transportation"
"Public_Transportation" "Public_Transportation" "Walking" ...
## $ NObeyesdad : chr "Normal_Weight" "Normal_Weight"
"Normal_Weight" "Overweight_Level_I" ...
```

- Add BMI column
- By calculate using “BMI = weight (kg) ÷ height2 (meters)”

```
age_obesity$BMI = age_obesity$Weight/(age_obesity$Height^2)
str(age_obesity)

## 'data.frame': 2111 obs. of 18 variables:
## $ Gender : chr "Female" "Female" "Male" "Male"
...
## $ Age : num 21 21 23 27 22 29 23 22 24 22 ...
## $ Height : num 1.62 1.52 1.8 1.8 1.78 1.62 1.5
1.64 1.78 1.72 ...
## $ Weight : num 64 56 77 87 89.8 53 55 53 64 68
...
## $ family_history_with_overweight: chr "yes" "yes" "yes" "no" ...
## $ FAVC : chr "no" "no" "no" "no" ...
## $ FCVC : num 2 3 2 3 2 2 3 2 3 2 ...
## $ NCP : num 3 3 3 3 1 3 3 3 3 3 ...
## $ CAEC : chr "Sometimes" "Sometimes"
"Sometimes" "Sometimes" ...
## $ SMOKE : chr "no" "yes" "no" "no" ...
## $ CH20 : num 2 3 2 2 2 2 2 2 2 2 ...
## $ SCC : chr "no" "yes" "no" "no" ...
## $ FAF : num 0 3 2 2 0 0 1 3 1 1 ...
## $ TUE : num 1 0 1 0 0 0 0 0 1 1 ...
## $ CALC : chr "no" "Sometimes" "Frequently"
"Frequently" ...
## $ MTRANS : chr "Public_Transportation"
"Public_Transportation" "Public_Transportation" "Walking" ...
## $ NObeyesdad : chr "Normal_Weight" "Normal_Weight"
"Normal_Weight" "Overweight_Level_I" ...
## $ BMI : num 24.4 24.2 23.8 26.9 28.3 ...
```

- Reorder columns to put BMI immediately after Height and Weight

```
obesity_new <- age_obesity[, c("Gender", "Age", "Height", "Weight", "BMI",
"family_history_with_overweight", "FAVC", "FCVC", "NCP", "CAEC", "SMOKE",
"CH20", "SCC", "FAF", "TUE", "CALC", "MTRANS", "NObeyesdad")]
str(obesity_new)

## 'data.frame': 2111 obs. of 18 variables:
## $ Gender : chr "Female" "Female" "Male" "Male"
...
```

```
## $ Age : num 21 21 23 27 22 29 23 22 24 22 ...
## $ Height : num 1.62 1.52 1.8 1.8 1.78 1.62 1.5
1.64 1.78 1.72 ...
## $ Weight : num 64 56 77 87 89.8 53 55 53 64 68
...
## $ BMI : num 24.4 24.2 23.8 26.9 28.3 ...
## $ family_history_with_overweight: chr "yes" "yes" "yes" "no" ...
## $ FAVC : chr "no" "no" "no" "no" ...
## $ FCVC : num 2 3 2 3 2 2 3 2 3 2 ...
## $ NCP : num 3 3 3 3 1 3 3 3 3 3 ...
## $ CAEC : chr "Sometimes" "Sometimes"
"Sometimes" "Sometimes" ...
## $ SMOKE : chr "no" "yes" "no" "no" ...
## $ CH2O : num 2 3 2 2 2 2 2 2 2 2 ...
## $ SCC : chr "no" "yes" "no" "no" ...
## $ FAF : num 0 3 2 2 0 0 1 3 1 1 ...
## $ TUE : num 1 0 1 0 0 0 0 0 1 1 ...
## $ CALC : chr "no" "Sometimes" "Frequently"
"Frequently" ...
## $ MTRANS : chr "Public_Transportation"
"Public_Transportation" "Public_Transportation" "Walking" ...
## $ NObeyesdad : chr "Normal_Weight" "Normal_Weight"
"Normal_Weight" "Overweight_Level_I" ...
```

Rename the column name for better reading

```
names(obesity_new) <- c("Gender",
                        "Age",
                        "Height",
                        "Weight",
                        "BMI",
                        "Family_History_with_Overweight",
                        "High_Caloric_Food_Consumption",
                        "Frequency_Consumption_of_Vegetables",
                        "Number_of_Main_Meals",
                        "Consumption_of_Food_Between_Meals",
                        "Smoke",
                        "Consumption_of_Water_Daily",
                        "Calories_Consumption_Monitoring",
                        "Physical_Activity_Frequency",
                        "Time_Using_Technology",
                        "Consumption_of_Alcohol",
                        "Transportation_Used",
                        "Obesity")
```

- Remove '\_' underscore in values of Obesity column

```
obesity_new$Obesity <- gsub("_", " ", obesity_new$Obesity)
head(obesity_new)
```

```

##   Gender Age Height Weight      BMI Family_History_with_Overweight
## 1 Female  21   1.62  64.0 24.38653                      yes
## 2 Female  21   1.52  56.0 24.23823                      yes
## 3 Male    23   1.80  77.0 23.76543                      yes
## 4 Male    27   1.80  87.0 26.85185                      no
## 5 Male    22   1.78  89.8 28.34238                      no
## 6 Male    29   1.62  53.0 20.19509                      no
##   High_Caloric_Food_Consumption Frequency_Consumption_of_Vegetables
## 1                                no                                2
## 2                                no                                3
## 3                                no                                2
## 4                                no                                3
## 5                                no                                2
## 6                                yes                                2
##   Number_of_Main_Meals Consumption_of_Food_Between_Meals Smoke
## 1                    3                               Sometimes no
## 2                    3                               Sometimes yes
## 3                    3                               Sometimes no
## 4                    3                               Sometimes no
## 5                    1                               Sometimes no
## 6                    3                               Sometimes no
##   Consumption_of_Water_Daily Calories_Consumption_Monitoring
## 1                            2                                no
## 2                            3                                yes
## 3                            2                                no
## 4                            2                                no
## 5                            2                                no
## 6                            2                                no
##   Physical_Activity_Frequency Time_Using_Technology Consumption_of_Alcohol
## 1                            0                            1                no
## 2                            3                            0            Sometimes
## 3                            2                            1            Frequently
## 4                            2                            0            Frequently
## 5                            0                            0            Sometimes
## 6                            0                            0            Sometimes
##   Transportation_Used      Obesity
## 1 Public_Transportation Normal Weight
## 2 Public_Transportation Normal Weight
## 3 Public_Transportation Normal Weight
## 4           Walking      Overweight Level I
## 5 Public_Transportation Overweight Level II
## 6           Automobile      Normal Weight

```

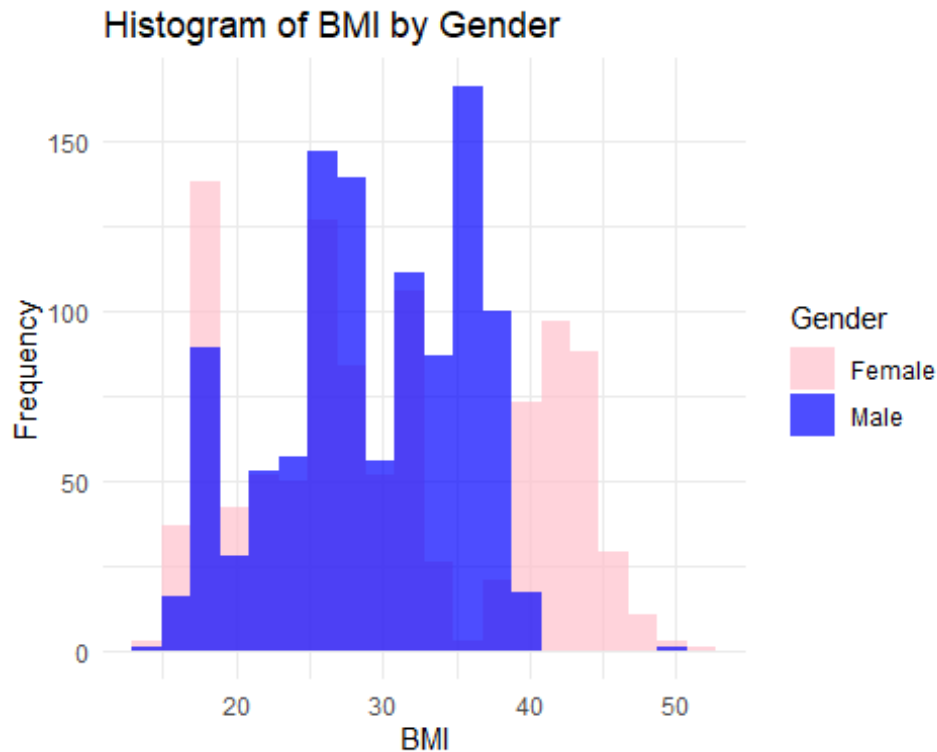
- save as new file

```
write.csv(obesity_new, "obesity_new1.csv", row.names = FALSE)
```

## Exploratory Data Analysis

1. Plot histogram BMI with gender differentiation

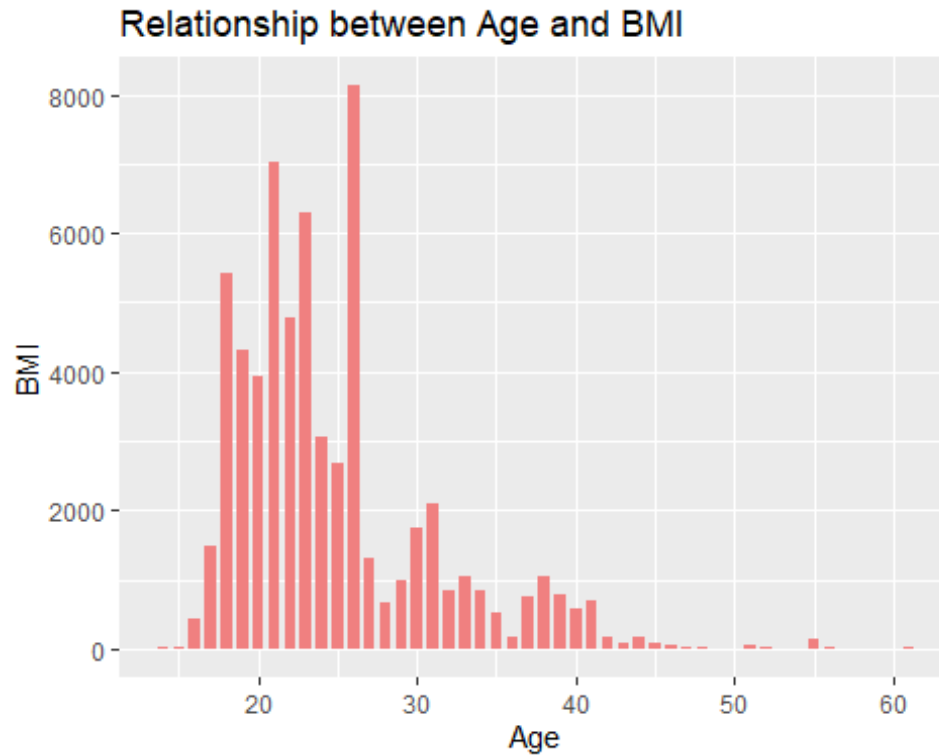
```
ggplot(obesity_new, aes(x = BMI, fill = Gender)) +
  geom_histogram(position = "identity", alpha = 0.7, bins = 20) +
  labs(title = "Histogram of BMI by Gender",
       x = "BMI",
       y = "Frequency") +
  scale_fill_manual(values = c("Male" = "blue", "Female" = "pink")) +
  theme_minimal()
```



## 2. Plot correlation between age and BMI

```
ggplot(obesity_new, aes(x = Age, y = BMI)) +
  geom_bar(stat = "identity", fill = "lightcoral", width = 0.7) +
  labs(title = "Relationship between Age and BMI",
       x = "Age",
       y = "BMI")
```

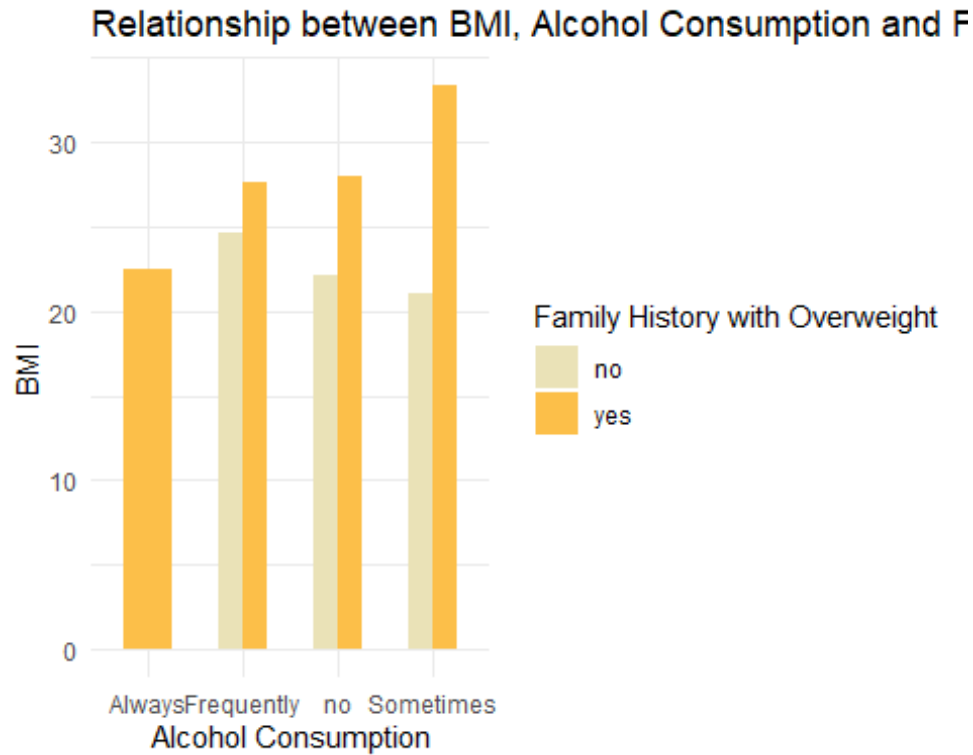




### 3. Alcohol Consumption, Family History with Overweight vs BMI

```
library(ggplot2)

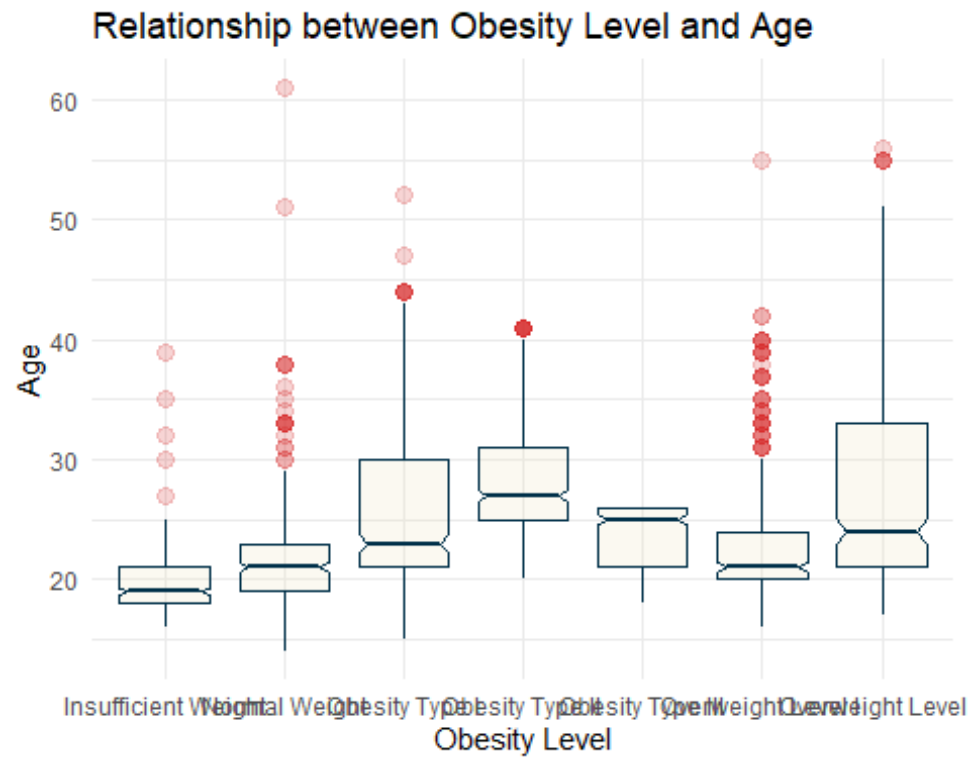
ggplot(obesity_new, aes(x = as.factor(Consumption_of_Alcohol),
                        y = BMI,
                        fill = Family_History_with_Overweight)) +
  geom_bar(stat = "summary", fun = "mean", position = "dodge", width = 0.5) +
  scale_fill_manual(values = c("#eae2b7", "#fcbf49")) +
  theme_minimal() +
  labs(
    title = "Relationship between BMI, Alcohol Consumption and Family History
with Overweight",
    x = "Alcohol Consumption",
    y = "BMI",
    fill = "Family History with Overweight"
  )
```



#### 4. Obesity vs Age

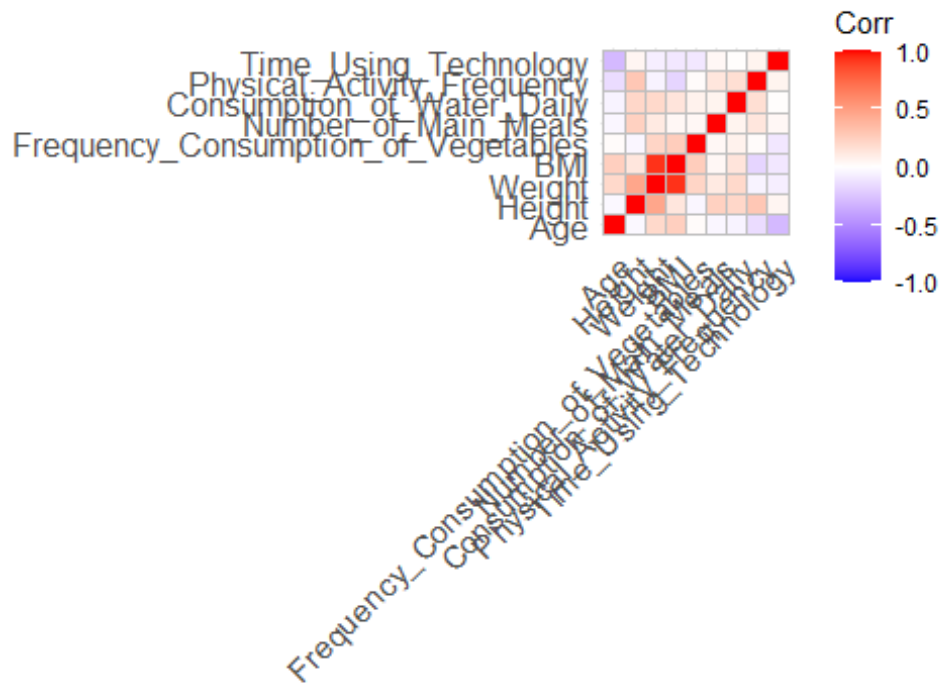
```
library(ggplot2)

ggplot(obesity_new, aes(x = as.factor(Obesity), y = Age)) +
  geom_boxplot(
    color = "#003049",
    fill = "#eae2b7",
    alpha = 0.2,
    notch = TRUE,
    notchwidth = 0.8,
    outlier.colour = "#d62828",
    outlier.fill = "#d62828",
    outlier.size = 3
  ) +
  theme_minimal() +
  labs(title = "Relationship between Obesity Level and Age",
       x = "Obesity Level",
       y = "Age")
```



## 5. Find the correlation between the numerical fields

```
numericFields <- dplyr::select_if(obesity_new, is.numeric)
r <- cor(numericFields, use="complete.obs")
ggcorrplot(r)
```



```
# Identify numeric columns
numeric_cols <- sapply(obesity_new, is.numeric)
numeric_data <- obesity_new[, numeric_cols]

# Function to detect outliers using IQR
detect_outliers <- function(x) {
  Q1 <- quantile(x, 0.25, na.rm = TRUE)
  Q3 <- quantile(x, 0.75, na.rm = TRUE)
  IQR_val <- Q3 - Q1
  outliers <- which(x < (Q1 - 1.5 * IQR_val) | x > (Q3 + 1.5 * IQR_val))
  return(outliers)
}

# Apply to each numeric column
outlier_summary <- sapply(numeric_data, function(col)
length(detect_outliers(col)))

# View summary table
outlier_summary_df <- data.frame(
  Variable = names(outlier_summary),
  Outlier_Count = as.integer(outlier_summary)
)
print(outlier_summary_df)

##                               Variable Outlier_Count
## 1                               Age                160
```

```

## 2          Height          1
## 3          Weight          1
## 4          BMI            0
## 5 Frequency_Consumption_of_Vegetables 0
## 6          Number_of_Main_Meals      579
## 7          Consumption_of_Water_Daily 0
## 8          Physical_Activity_Frequency 0
## 9          Time_Using_Technology      0

remove_outliers_iqr <- function(df) {
  numeric_cols <- sapply(df, is.numeric)
  for (col in names(df)[numeric_cols]) {
    Q1 <- quantile(df[[col]], 0.25, na.rm = TRUE)
    Q3 <- quantile(df[[col]], 0.75, na.rm = TRUE)
    IQR_val <- Q3 - Q1
    lower <- Q1 - 1.5 * IQR_val
    upper <- Q3 + 1.5 * IQR_val
    df <- df[df[[col]] >= lower & df[[col]] <= upper, ]
  }
  return(df)
}

obesity_clean <- remove_outliers_iqr(obesity_new)
cat("Original rows:", nrow(obesity_new), "\n")

## Original rows: 2111

cat("After outlier removal:", nrow(obesity_clean), "\n")

## After outlier removal: 1407

# Copy the dataset to avoid modifying the original
obesity_label_encoded <- obesity_clean

# Identify categorical columns
categorical_vars <- sapply(obesity_label_encoded, function(x) is.factor(x) ||
is.character(x))

# Apply label encoding to categorical columns
obesity_label_encoded[categorical_vars] <-
lapply(obesity_label_encoded[categorical_vars], function(x)
as.numeric(factor(x)))

# View the encoded dataset
head(obesity_label_encoded)

##   Gender Age Height Weight      BMI Family_History_with_Overweight
## 1      1  21   1.62    64 24.38653                2
## 2      1  21   1.52    56 24.23823                2
## 3      2  23   1.80    77 23.76543                2
## 4      2  27   1.80    87 26.85185                1

```

```

## 6      2  29  1.62    53 20.19509      1
## 7      1  23  1.50    55 24.44444      2
##   High_Caloric_Food_Consumption Frequency_Consumption_of_Vegetables
## 1              1              2
## 2              1              3
## 3              1              2
## 4              1              3
## 6              2              2
## 7              2              3
##   Number_of_Main_Meals Consumption_of_Food_Between_Meals Smoke
## 1              3              4      1
## 2              3              4      2
## 3              3              4      1
## 4              3              4      1
## 6              3              4      1
## 7              3              4      1
##   Consumption_of_Water_Daily Calories_Consumption_Monitoring
## 1              2              1
## 2              3              2
## 3              2              1
## 4              2              1
## 6              2              1
## 7              2              1
##   Physical_Activity_Frequency Time_Using_Technology Consumption_of_Alcohol
## 1              0              1              2
## 2              3              0              3
## 3              2              1              1
## 4              2              0              1
## 6              0              0              3
## 7              1              0              3
##   Transportation_Used Obesity
## 1              4          2
## 2              4          2
## 3              4          2
## 4              5          6
## 6              1          2
## 7              3          2

```

```
obesity_standardized <- obesity_label_encoded
```

```
obesity_standardized[numeric_cols] <-
```

```
scale(obesity_label_encoded[numeric_cols])
```

```
str(obesity_standardized)
```

```
## 'data.frame': 1407 obs. of 18 variables:
```

```
## $ Gender : num 1 1 2 2 2 1 2 2 2 2 ...
```

```
## $ Age : num -0.5212 -0.5212 -0.0488
```

```
0.8958 1.3681 ...
```

```
## $ Height : num -1.084 -2.23 0.978 0.978 -
```

```
1.084 ...
```

```
## $ Weight : num -1.015 -1.306 -0.541 -0.177 -
1.415 ...
## $ BMI : num -0.788 -0.806 -0.862 -0.497 -
1.284 ...
## $ Family_History_with_Overweight : num 2 2 2 1 1 2 1 2 2 2 ...
## $ High_Caloric_Food_Consumption : num 1 1 1 1 2 2 1 2 2 2 ...
## $ Frequency_Consumption_of_Vegetables: num -0.826 0.994 -0.826 0.994 -
0.826 ...
## $ Number_of_Main_Meals : num 0.209 0.209 0.209 0.209 0.209
...
## $ Consumption_of_Food_Between_Meals : num 4 4 4 4 4 4 4 4 4 2 ...
## $ Smoke : num 1 2 1 1 1 1 1 1 1 1 ...
## $ Consumption_of_Water_Daily : num -0.0723 1.5791 -0.0723 -
0.0723 -0.0723 ...
## $ Calories_Consumption_Monitoring : num 1 2 1 1 1 1 1 1 1 1 ...
## $ Physical_Activity_Frequency : num -1.22 2.33 1.15 1.15 -1.22
...
## $ Time_Using_Technology : num 0.51 -1.2 0.51 -1.2 -1.2 ...
## $ Consumption_of_Alcohol : num 2 3 1 1 3 3 3 1 2 3 ...
## $ Transportation_Used : num 4 4 4 5 1 3 4 4 4 4 ...
## $ Obesity : num 2 2 2 6 2 2 2 2 2 3 ...
```

## Modeling

### 1. BMI Prediction (Regression Model)

- Define a Rsquared function and remove rows with Null values

```
df_reg <- obesity_new
df_reg <- df_reg %>%
  mutate_if(is.character, as.factor) %>%
  na.omit()
set.seed(123)
```

- Train, Test, Split

```
splitIndex <- createDataPartition(obesity_standardized$BMI, p = 0.8, list =
FALSE)
training_data <- obesity_standardized[splitIndex, ]
testing_data <- obesity_standardized[-splitIndex, ]
```

- Create Linear Regression model and train based on Training Data

```
model <- lm(BMI ~ ., data = training_data)
```

- Make predictions with created model using Testing Data

```
pred_reg <- round(as.numeric(predict(model, newdata = testing_data)), digits =
2)
pred_reg

## [1] -0.95 -1.33 -0.95 -0.57 -1.47 -0.95 -0.72 -0.91 -1.18 0.20 -0.42 -
0.96
```

```
## [13] -0.97 -0.35 -0.83 -0.42 -1.22 -0.92 -1.74 -0.78 -1.10 -1.44 -0.58
0.11
## [25] -0.74 -1.19 -1.10 -0.29 -0.68 -1.25 -0.41 0.22 0.13 -1.91 -0.98 -
0.83
## [37] -1.29 -0.95 -0.53 -1.36 -1.12 -1.24 0.46 -0.81 -1.08 -1.08 -1.14 -
0.75
## [49] -0.33 -1.62 -0.57 -0.51 0.24 -0.85 -0.97 -0.96 -0.89 -1.09 -1.34 -
0.72
## [61] -1.08 -1.19 -1.55 -1.53 -1.83 -1.24 -1.66 -1.63 -1.41 -2.22 -1.82 -
1.69
## [73] -1.44 -1.38 -1.44 -1.55 -1.49 -1.77 -1.62 -1.81 -1.89 -1.47 -1.61 -
1.56
## [85] -1.63 -1.68 -1.92 -0.54 -0.64 -0.70 -0.63 -0.57 -0.31 -0.55 -0.57 -
0.58
## [97] -0.66 -0.48 -0.51 -0.59 -0.55 -0.62 -0.62 -0.65 -0.55 -0.59 -0.39 -
0.25
## [109] -0.38 -0.41 -0.29 -0.27 -0.16 -0.12 -0.26 -0.50 -0.50 -0.27 -0.21 -
0.47
## [121] -0.21 -0.16 -0.41 -0.30 -0.26 -0.27 -0.27 -0.41 -0.25 -0.33 -0.21 -
0.52
## [133] -0.38 -0.18 -0.45 -0.53 -0.29 -0.36 -0.04 -0.11 0.50 -0.01 -0.06
0.20
## [145] 0.13 0.17 0.18 0.06 -0.05 0.14 -0.01 -0.01 0.03 0.05 0.14
0.30
## [157] -0.02 0.03 0.14 -0.07 -0.08 0.34 0.31 0.12 0.07 0.09 0.07
0.24
## [169] 0.42 -0.05 0.31 0.35 0.77 0.64 0.57 0.51 0.69 0.63 0.87
0.64
## [181] 0.61 0.69 0.53 0.73 0.64 0.55 0.66 0.67 0.63 0.70 0.67
0.62
## [193] 0.91 0.91 0.76 0.58 0.57 0.59 0.77 0.74 0.87 0.93 0.80
0.64
## [205] 0.61 0.65 0.66 0.70 0.55 0.72 0.95 0.87 0.78 0.54 0.65
0.54
## [217] 0.63 1.49 1.92 1.07 1.43 1.13 1.51 1.53 0.79 0.92 1.50
1.59
## [229] 1.16 1.48 1.51 2.08 0.90 1.01 1.13 1.19 1.13 1.50 2.11
1.53
## [241] 1.17 1.13 0.99 1.13 1.52 0.91 0.88 0.94 0.78 0.74 1.36
1.39
## [253] 1.66 1.15 1.52 1.99 1.05 1.52 0.89 1.53 1.96 1.90 1.45
1.26
## [265] 1.50 1.53 0.88 0.89 2.07 1.54 1.55 1.50 1.14 1.19 1.06
1.13
## [277] 1.12 1.03 1.58 1.55
```

-Check for prediction accuracy using Mean Square Error

```
mse_reg <- mean((testing_data$BMI - pred_reg)^2)
summary(model)
```



```
##
## Call:
## lm(formula = BMI ~ ., data = training_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.37936 -0.04572 -0.00695  0.05251  0.35345
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -0.0870413   0.0452446   -1.924  0.054636
## .
## Gender         0.0232701   0.0078587    2.961  0.003131
## **
## Age           0.0026383   0.0034608    0.762  0.446016
## Height       -0.3512360   0.0042865  -81.939 < 2e-16
## ***
## Weight        1.0682883   0.0043099  247.866 < 2e-16
## ***
## Family_History_with_Overweight 0.0497732   0.0091584    5.435 6.75e-08
## ***
## High_Caloric_Food_Consumption 0.0307038   0.0092401    3.323 0.000920
## ***
## Frequency_Consumption_of_Vegetables 0.0252844   0.0031453    8.039 2.31e-15
## ***
## Number_of_Main_Meals 0.0103980   0.0027789    3.742 0.000192
## ***
## Consumption_of_Food_Between_Meals 0.0065591   0.0041870    1.567 0.117513
## Smoke        -0.0547496   0.0188798   -2.900 0.003806
## **
## Consumption_of_Water_Daily 0.0004806   0.0028952    0.166 0.868195
## Calories_Consumption_Monitoring -0.0579424   0.0138969   -4.169 3.29e-05
## ***
## Physical_Activity_Frequency -0.0147006   0.0029734   -4.944 8.83e-07
## ***
## Time_Using_Technology -0.0030970   0.0027469   -1.127 0.259801
## Consumption_of_Alcohol -0.0169783   0.0056455   -3.007 0.002694
## **
## Transportation_Used 0.0021635   0.0028578    0.757 0.449188
## Obesity       0.0079122   0.0017310    4.571 5.40e-06
## ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08818 on 1109 degrees of freedom
## Multiple R-squared:  0.9923, Adjusted R-squared:  0.9922
## F-statistic: 8432 on 17 and 1109 DF, p-value: < 2.2e-16
```

- Generate random elements without replacement
- Convert the variables into factors

- Train, Test, Split

```
splitIndex <- createDataPartition(obesity_standardized$obesity, p = 0.8, list
= FALSE)
training_data <- obesity_standardized[splitIndex, ]
testing_data <- obesity_standardized[-splitIndex, ]
```

- Make predictions with created model using Testing Data

[illegible]

```

7
## 1009 1010 1036 1050 1054 1057 1058 1066 1070 1072 1103 1109 1111 1116 1129
1137
##    7    7    7    7    7    7    7    7    7    7    7    7    7    7    7
7
## 1143 1144 1146 1154 1157 1162 1164 1169 1177 1179 1198 1203 1205 1229 1237
1240
##    7    7    7    7    7    7    7    7    7    7    7    7    7    4    3
3
## 1259 1261 1273 1278 1280 1293 1295 1299 1305 1314 1346 1361 1366 1367 1375
1400
##    3    3    3    3    3    3    3    3    4    3    3    4    3    3    3
3
## 1407 1429 1443 1447 1449 1452 1458 1465 1466 1467 1499 1501 1529 1532 1540
1547
##    3    3    3    3    3    3    3    3    3    3    3    3    4    4    4
4
## 1554 1555 1558 1571 1575 1577 1587 1588 1595 1612 1636 1639 1647 1648 1649
1654
##    4    4    4    4    4    4    4    4    4    4    4    4    4    4    4
4
## 1668 1675 1677 1684 1685 1689 1695 1696 1697 1702 1715 1717 1720 1734 1739
1747
##    4    4    4    4    4    4    4    4    4    4    4    4    4    4    4
4
## 1751 1752 1759 1761 1770 1774 1803 1805 1819 1821 1829 1831 1836 1838 1840
1841
##    4    4    4    4    4    4    5    5    5    5    5    5    5    5    5
5
## 1843 1853 1872 1876 1878 1884 1887 1899 1907 1909 1914 1921 1924 1930 1934
1945
##    5    5    5    5    5    5    5    5    5    5    5    5    5    5    5
5
## 1947 1948 1961 1966 1967 1970 1974 1976 1977 1978 1993 1994 2004 2011 2013
2017
##    5    5    5    5    5    5    5    5    5    5    5    5    5    5    5
5
## 2021 2025 2032 2033 2036 2037 2043 2052 2054 2066 2067 2069 2074 2076 2078
2088
##    5    5    5    5    5    5    5    5    5    5    5    5    5    5    5
5
## 2093 2098 2106 2109 2110 2111
##    5    5    5    5    5    5
## Levels: 1 2 3 4 5 6 7

```

- Check for prediction accuracy using Confusion Matrix

```

mse_svm <- confusionMatrix(pred_svm, testing_data$Obesity)
mse_svm

```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  1  2  3  4  5  6  7
##           1 28  0  0  0  0  0  0
##           2  1 38  0  0  0  3  0
##           3  0  0 35  0  0  0  0
##           4  0  0  4 43  0  0  0
##           5  0  0  0  0 64  0  0
##           6  0  1  0  0  0 26  1
##           7  0  0  0  0  0  1 33
##
## Overall Statistics
##
##           Accuracy : 0.9604
##           95% CI : (0.9303, 0.9801)
##           No Information Rate : 0.2302
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9532
##
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.9655  0.9744  0.8974  1.0000  1.0000  0.86667
## Specificity      1.0000  0.9833  1.0000  0.9830  1.0000  0.99194
## Pos Pred Value   1.0000  0.9048  1.0000  0.9149  1.0000  0.92857
## Neg Pred Value   0.9960  0.9958  0.9835  1.0000  1.0000  0.98400
## Prevalence       0.1043  0.1403  0.1403  0.1547  0.2302  0.10791
## Detection Rate   0.1007  0.1367  0.1259  0.1547  0.2302  0.09353
## Detection Prevalence 0.1007  0.1511  0.1259  0.1691  0.2302  0.10072
## Balanced Accuracy 0.9828  0.9788  0.9487  0.9915  1.0000  0.92930
##           Class: 7
## Sensitivity      0.9706
## Specificity      0.9959
## Pos Pred Value   0.9706
## Neg Pred Value   0.9959
## Prevalence       0.1223
## Detection Rate   0.1187
## Detection Prevalence 0.1223
## Balanced Accuracy 0.9832
```

### 3. Obesity Level Prediction (Logistic Regression Model)

```
library(nnet)
library(caret)
```

```
# Ensure target is a factor
```

```

training_data$Obesity <- as.factor(training_data$Obesity)
testing_data$Obesity <- as.factor(testing_data$Obesity)

model <- train(
  Obesity ~ .,
  data = training_data,
  method = "multinom",
  trControl = trainControl(method = "cv", number = 5),
  trace = FALSE
)
logit_pred <- predict(model, testing_data)
logit_conf <- confusionMatrix(logit_pred, testing_data$Obesity)
print(logit_conf)

```

```
## Confusion Matrix and Statistics
```

```
##
##              Reference
## Prediction  1  2  3  4  5  6  7
##           1 28  0  0  0  0  0  0
##           2  1 37  0  0  0  4  0
##           3  0  0 34  0  0  0  0
##           4  0  0  3 43  0  0  2
##           5  0  1  1  0 64  0  0
##           6  0  1  0  0  0 24  1
##           7  0  0  1  0  0  2 31
##
```

```
## Overall Statistics
```

```
##
##              Accuracy : 0.9388
##              95% CI : (0.9039, 0.964)
##      No Information Rate : 0.2302
##      P-Value [Acc > NIR] : < 2.2e-16
##
```

```
##              Kappa : 0.9276
##
```

```
## McNemar's Test P-Value : NA
##
```

```
## Statistics by Class:
```

```
##
##              Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.9655   0.9487   0.8718   1.0000   1.0000   0.80000
## Specificity      1.0000   0.9791   1.0000   0.9787   0.9907   0.99194
## Pos Pred Value    1.0000   0.8810   1.0000   0.8958   0.9697   0.92308
## Neg Pred Value    0.9960   0.9915   0.9795   1.0000   1.0000   0.97619
## Prevalence        0.1043   0.1403   0.1403   0.1547   0.2302   0.10791
## Detection Rate    0.1007   0.1331   0.1223   0.1547   0.2302   0.08633
## Detection Prevalence 0.1007   0.1511   0.1223   0.1727   0.2374   0.09353
## Balanced Accuracy  0.9828   0.9639   0.9359   0.9894   0.9953   0.89597
##
##              Class: 7
## Sensitivity      0.9118
```

```
## Specificity          0.9877
## Pos Pred Value      0.9118
## Neg Pred Value      0.9877
## Prevalence          0.1223
## Detection Rate      0.1115
## Detection Prevalence 0.1223
## Balanced Accuracy    0.9497
```

## 4. Obesity Level Prediction (Random Forest Model)

- Prepare dataset (reuse cleaned data)

```
set.seed(789)
df_rf <- obesity_standardized
df_rf$Gender <- as.factor(df_rf$Gender)
df_rf$Obesity <- as.factor(df_rf$Obesity)
```

- Train/test split

```
splitIndex <- createDataPartition(df_rf$Obesity, p = 0.7, list = FALSE)
train_data_rf <- df_rf[splitIndex, ]
test_data_rf <- df_rf[-splitIndex, ]
```

- Train Random Forest model

```
library(randomForest)

## Warning: package 'randomForest' was built under R version 4.4.2
## randomForest 4.7-1.2
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
##
## The following object is masked from 'package:ggplot2':
##
##     margin
##
## The following object is masked from 'package:dplyr':
##
##     combine

rf_model <- randomForest(Obesity ~ ., data = train_data_rf, ntree = 100,
importance = TRUE)
```

- Predict

```
pred_rf <- predict(rf_model, newdata = test_data_rf)
```

- Confusion Matrix

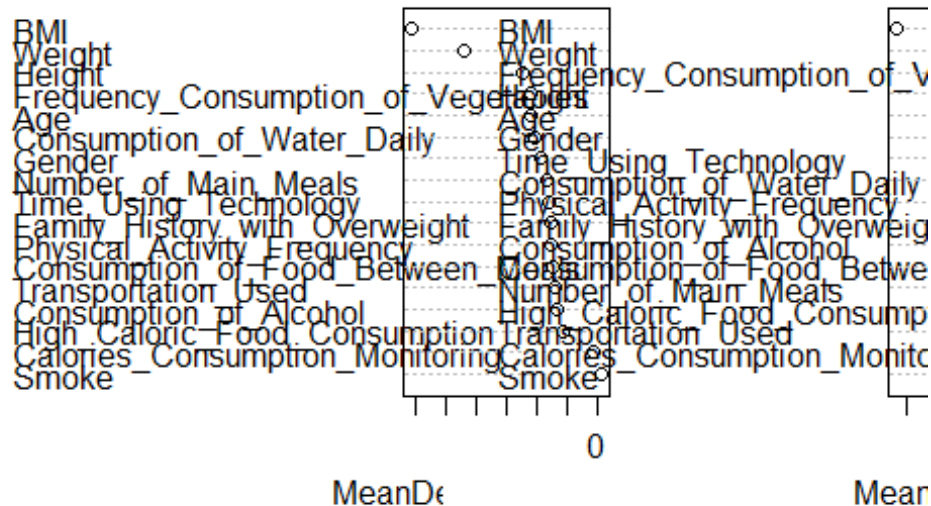
```
rf_conf <- confusionMatrix(pred_rf, test_data_rf$Obesity)
print(rf_conf)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  1  2  3  4  5  6  7
##           1 42  0  0  0  0  0
##           2  1 59  0  0  0  1
##           3  0  0 58  0  0  0
##           4  0  0  0 65  0  0
##           5  0  0  0  0 97  0
##           6  0  0  0  0  0 44
##           7  0  0  0  0  0  0 50
##
## Overall Statistics
##
##           Accuracy : 0.9928
##           95% CI : (0.9792, 0.9985)
##           No Information Rate : 0.2321
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9915
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: 1 Class: 2 Class: 3 Class: 4 Class: 5 Class: 6
## Sensitivity      0.9767   1.0000   1.0000   1.0000   1.0000   0.9778
## Specificity      1.0000   0.9916   1.0000   1.0000   1.0000   1.0000
## Pos Pred Value   1.0000   0.9516   1.0000   1.0000   1.0000   1.0000
## Neg Pred Value   0.9973   1.0000   1.0000   1.0000   1.0000   0.9973
## Prevalence       0.1029   0.1411   0.1388   0.1555   0.2321   0.1077
## Detection Rate   0.1005   0.1411   0.1388   0.1555   0.2321   0.1053
## Detection Prevalence 0.1005   0.1483   0.1388   0.1555   0.2321   0.1053
## Balanced Accuracy 0.9884   0.9958   1.0000   1.0000   1.0000   0.9889
##
##           Class: 7
## Sensitivity      0.9804
## Specificity      1.0000
## Pos Pred Value   1.0000
## Neg Pred Value   0.9973
## Prevalence       0.1220
## Detection Rate   0.1196
## Detection Prevalence 0.1196
## Balanced Accuracy 0.9902
```

- Plot variable importance

```
varImpPlot(rf_model, main = "Variable Importance (Random Forest)")
```

## Variable Importance (Random Forest)



```
svm_acc <- confusionMatrix(pred_svm,
testing_data$obesity)$overall['Accuracy']
rf_acc <- confusionMatrix(pred_rf, test_data_rf$obesity)$overall['Accuracy']
logit_acc<- confusionMatrix(logit_pred,
testing_data$obesity)$overall['Accuracy']
accuracy_results <- data.frame(
  Model = c("Multinomial Logistic Regression", "Random Forest", "SVM"),
  Accuracy = c(logit_acc, rf_acc, svm_acc)
)
print(accuracy_results)

##               Model  Accuracy
## 1 Multinomial Logistic Regression 0.9388489
## 2               Random Forest 0.9928230
## 3                   SVM 0.9604317
```

## Conclusion

- **BMI Prediction** (Regression Model) The linear regression model for predicting BMI showed excellent performance, achieving a high R-squared value of 0.99. This indicates that the model explains nearly all the variance in BMI using the input features. Key predictors included height, weight, family history of overweight, and dietary habits, all of which had statistically significant effects.
- **Obesity Level Prediction**



### Support Vector Machine (SVM) Model

The improved SVM model, after kernel adjustment and hyperparameter tuning, achieved a strong accuracy of approximately **96.04%**. This is a significant improvement from the earlier 15% accuracy when using a basic linear kernel. The enhanced model is now capable of effectively handling multi-class classification, though further optimization and advanced feature engineering could push its performance even higher.

### Logistic Regression Model

The Multinomial Logistic Regression model also performed well, achieving an accuracy of **93.8%**. It correctly classified most obesity levels and showed reliable consistency across categories. While slightly less accurate than the Random Forest and SVM models, it remains a robust and interpretable choice for multi-class classification problems.

### Random Forest Model

In contrast, the Random Forest model demonstrated the highest performance, with an accuracy of **99.2%** and a Kappa statistic of **0.9915**, indicating near-perfect agreement between predictions and actual labels. It successfully classified all obesity categories with high sensitivity and specificity and offered valuable insights into feature importance, making it both powerful and interpretable for classification tasks.