



Predicting Heart Attack Risk Through Data Analysis and Machine Learning Models

Nana Firdausi Hassan

Table of Contents

THE DATA	3
Description of Data.....	3
Loading Packages.....	3
Loading Dataset	3
Original Data Structure	4
Data Cleaning	4
Cleaned Data Frame	6
Description of Variables	7
Expectations	8
DATA EXPLORATION AND ANALYSIS	9
Plots for Categorical Variables	9
Correlation Analysis of Heart Attack Risk and Other Variables	15
Distributional Analysis of the Target Variable	18
Summary Statistics of Variables	19
Correlation of Numerical Variables	19
Plots for Correlated Variables	20
Predictive Modeling	23
Logistic Regression.....	24
Hypothesis Testing.....	24
Forward Selection	26
Backward Selection.....	26
Linear Discriminant Analysis	30
Quadratic Discriminant Analysis (QDA)	31
Naive Bayes	31
Regression Analysis	32
Multiple Linear Regression	32
Ridge Regression	34
Lasso Regression.....	34

Principal Component Analysis.....	35
K-Nearest Neighbor	37
K-Means Clustering	37
SUMMARY.....	40

THE DATA

Description of Data

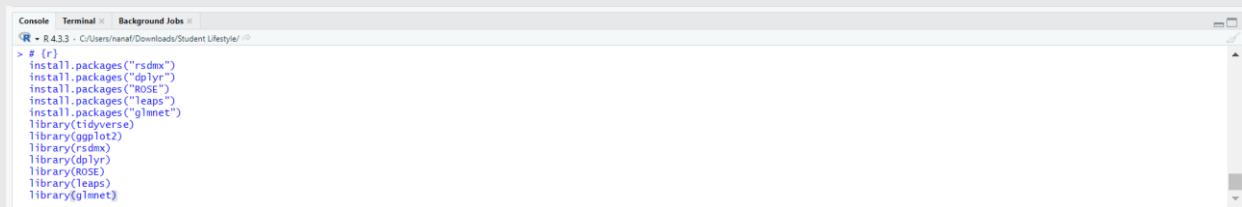
The Heart Attack Prediction Dataset was obtained from Kaggle, it contains detailed information about patients and their risk of heart attack. The dataset includes 8,763 observations and 26 variables, offering a wealth of data on patient demographics, medical history, and diagnostic measurements.

This dataset provides valuable resources for data scientists and healthcare researchers who aim to explore the factors contributing to heart attack risk and develop machine learning models for early detection and prevention.

Link to the dataset: [Heart Attack Prediction Dataset on Kaggle](#)

Loading Packages

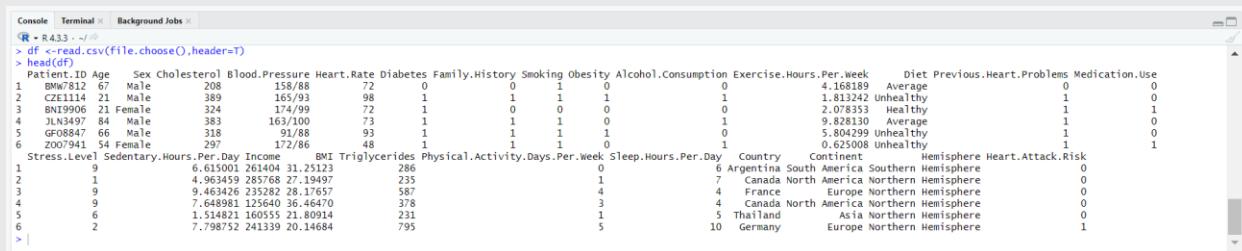
For this analysis, I installed and loaded several R packages to support data processing, visualization, and modeling. dplyr and tidyverse were used for data manipulation and cleaning, while ggplot2 supported data visualization. The ROSE package helped address class imbalance in the dataset. For variable selection and model building, I used leaps for subset selection and glmnet for regularized regression techniques like Lasso and Ridge. Additionally, rsdmx was included to allow access to structured statistical data formats when needed. More packages will be added as the analysis progresses and additional functionality is required.



```
R > # (r)
> install.packages("rsdmx")
> install.packages("dplyr")
> install.packages("ROSE")
> install.packages("leaps")
> install.packages("glmnet")
> library(tidyverse)
> library(ggplot2)
> library(rsdmx)
> library(dplyr)
> library(ROSE)
> library(leaps)
> library(glmnet)
```

Loading Dataset

I imported the dataset using `read.csv()` and used `head()` to view the first few rows. This step helped confirm the data loaded correctly and provided an overview of the variables for further analysis.



```
R > df<-read.csv(file.choose(),header=T)
> head(df)
#> #> #> #> #>
#> #> #> #> #>
#> Patient.ID Age Sex Cholesterol Blood.Pressure Heart.Rate Diabetes Family.History Smoking Obesity Alcohol.Consumption Exercise.Hours.Per.Week Diet Previous.Heart.Problems Medication.Use
#> 1 BMW812 67 Male 208 158/88 72 0 0 1 0 0 4.168189 Average 0 0
#> 2 CZE1114 21 Male 389 165/93 98 1 1 1 1 1 1.813242 Unhealthy 1 0
#> 3 BN9906 21 Female 324 174/99 72 1 0 0 0 0 2.078353 Healthy 1 1
#> 4 JLN3497 84 Male 383 163/109 73 1 1 0 1 1 9.282130 Average 1 0
#> 5 GP0887 50 Male 318 91/88 93 1 1 1 1 0 5.882399 Unhealthy 1 0
#> 6 ZOO7941 54 Female 297 172/86 48 1 1 1 1 0 0.625000 Unhealthy 1 1
#> Stress.Level Sedentary.Hours.Per.Day Income BMI Triglycerides Physical.Activity.Days.Per.Week Sleep.Hours.Per.Day Country Continent Hemisphere Heart.Attack.Risk
#> 1 9 6,615001 2616404 31.25123 286 0 1 Argentina South America Southern Hemisphere 0
#> 2 1 4,963459 285768 27.19497 235 1 7 Canada North America Northern Hemisphere 0
#> 3 9 9,463426 235282 28.17657 587 4 4 France Europe Northern Hemisphere 0
#> 4 9 9,463426 235282 30.16670 378 3 4 Canada North America Northern Hemisphere 0
#> 5 6 1,514821 160555 21.80914 231 2 5 Thailand Asia Northern Hemisphere 0
#> 6 2 7,798752 241339 20.14684 795 5 10 Germany Europe Northern Hemisphere 1
#> |
```

I used `is.data.frame(df)` to confirm the dataset is a data frame and is correctly loaded into R, which returned TRUE. Then, I checked the dimensions with `dim(df)`, revealing that the dataset has 8763 observations and 26 variables.

```
Console Terminal Background Jobs
R - R 4.3.3 - ~/r
> is.data.frame(df)
[1] TRUE
> dim(df)
[1] 8763 26
>
```

Original Data Structure

To obtain information about the variables in the dataset, I used the `str(df)` function, which displays the structure of the dataset, including its mode and data type. By using the `str()` function, I was able to determine that the dataset had variables with different modes, including numeric (num), integer (int), and character (chr) modes. The dataset includes both categorical variables and numerical variables.

```
Console Terminal Background Jobs
R - R 4.3.3 - ~/r
> str(df)
'data.frame': 8763 obs. of 26 variables:
 $ Patient.ID      : chr "8MW7812" "CZE1114" "BN19006" "JLN3497" ...
 $ Age              : int 67 21 21 84 66 54 90 84 20 43 ...
 $ Sex              : chr "Male" "Male" "Female" "Male" ...
 $ Cholesterol     : int 208 389 324 383 318 297 358 220 145 248 ...
 $ Blood.Pressure   : chr "158/88" "165/93" "174/99" "163/100" ...
 $ Heart.Rate        : int 72 90 72 73 93 44 84 107 68 55 ...
 $ Diabetes         : int 0 1 1 1 1 1 0 0 1 0 ...
 $ Family.History   : int 0 1 0 1 1 1 0 0 0 1 ...
 $ Smoking          : int 1 1 0 1 1 1 1 1 1 ...
 $ Obesity           : int 0 1 0 0 1 0 0 1 1 1 ...
 $ Alcohol.Consumption: int 0 1 0 1 0 1 1 1 0 1 ...
 $ Exercise.Hours.Per.Week: num 4.17 1.81 2.08 9.83 5.8 ...
 $ Diet              : chr "Average" "Unhealthy" "Healthy" "Average" ...
 $ Previous.Heart.Problems: int 0 1 1 1 1 1 0 0 0 0 ...
 $ Medication.Use    : int 0 0 0 0 1 0 1 0 ...
 $ Stress.Level      : int 9 1 9 9 6 2 7 4 5 4 ...
 $ Sedentary.Hours.Per.Day: num 6.62 4.96 9.46 7.65 1.51 ...
 $ Income            : num 261404 285768 235282 125640 160555 241339 190450 122093 25086 209703 ...
 $ BMI               : num 31.3 27.2 28.2 36.5 21.8 ...
 $ Triglycerides    : num 286 231 587 378 231 795 284 370 790 232 ...
 $ Physical.Activity.Days.Per.Week: int 0 1 4 3 3 1 5 4 6 7 ...
 $ Sleep.Hours.Per.Day: int 6 7 4 4 5 10 10 7 4 7 ...
 $ Country           : chr "Argentina" "Canada" "France" "Canada" ...
 $ Continent         : chr "South America" "North America" "Europe" "North America" ...
 $ Hemisphere         : chr "Southern Hemisphere" "Northern Hemisphere" "Northern Hemisphere" "Northern Hemisphere" ...
 $ Heart.Attack.Risk  : int 0 0 0 0 1 1 1 0 0 ...
>
```

I used the `colnames(df)` function to retrieve the column names of the dataset.

```
Console Terminal Background Jobs
R - R 4.3.3 - ~/r
> colnames(df)
 [1] "Patient.ID"          "Age"          "Sex"          "Family.History"
 [6] "Cholesterol"          "Smoking"       "Previous.Heart.Problems"
 [11] "Alcohol.Consumption" "Exercise.Hours.Per.Week" "Diet"
 [16] "Stress.Level"         "Sedentary.Hours.Per.Day" "Income"
 [21] "Physical.Activity.Days.Per.Week" "Sleep.Hours.Per.Day" "Country"
 [26] "Heart.Attack.Risk"    "Continent"      "Hemisphere"
>
```

Data Cleaning

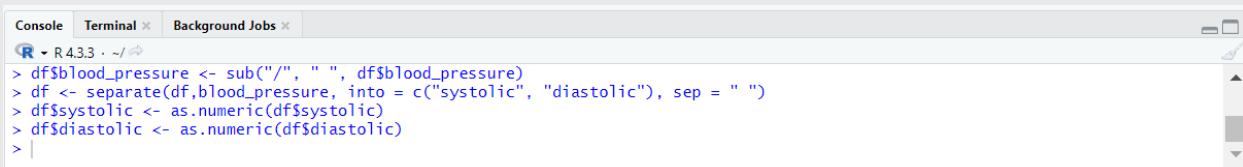
As the first step in my data cleaning process, I used the `clean_names()` function from the `janitor` package to standardize the column names in the dataset. This function reformats all column names to lowercase, replaces spaces and special characters with underscores, and ensures that names are unique and syntactically valid. For example, a

column like "Patient.ID" is converted to "patient_id". This step improves consistency, prevents potential coding errors, and makes data manipulation easier in the following stages of analysis.



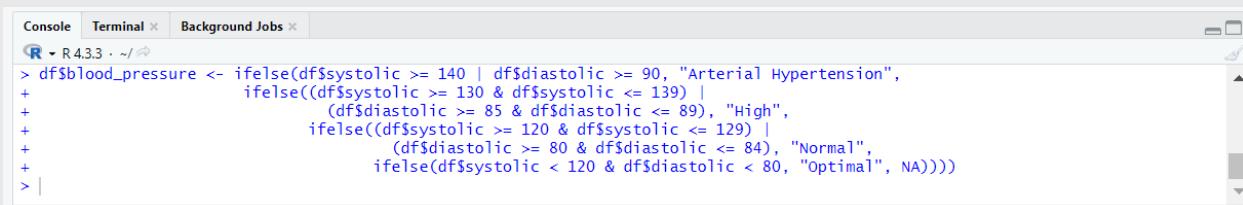
```
Console Terminal Background Jobs
R - R 4.3.3 - ~/Documents
> # Data Cleaning
> df <- clean_names(df)
>
```

I cleaned the blood_pressure column by first replacing the "/" with a space using the sub() function. Then, I split the column into two separate columns, systolic and diastolic, using the separate() function. Finally, I converted both columns to numeric format with as.numeric() to allow for proper analysis. This step made it possible to work with systolic and diastolic values independently in numerical operations.



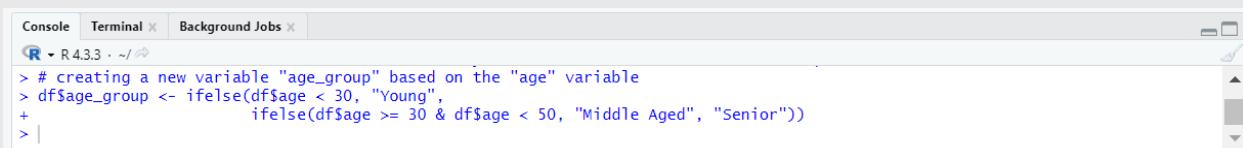
```
Console Terminal Background Jobs
R - R 4.3.3 - ~/Documents
> df$blood_pressure <- sub("/", " ", df$blood_pressure)
> df <- separate(df,blood_pressure, into = c("systolic", "diastolic"), sep = " ")
> df$systolic <- as.numeric(df$systolic)
> df$diastolic <- as.numeric(df$diastolic)
>
```

I created a new classification column for 'blood pressure' using nested ifelse() statements. Readings with systolic ≥ 140 or diastolic ≥ 90 were classified as "Arterial Hypertension." Values with systolic 130–139 or diastolic 85–89 were labeled "High," those with systolic 120–129 or diastolic 80–84 as "Normal," and readings with systolic < 120 and diastolic < 80 as "Optimal." This allowed for standardized categorization of blood pressure levels for analysis.



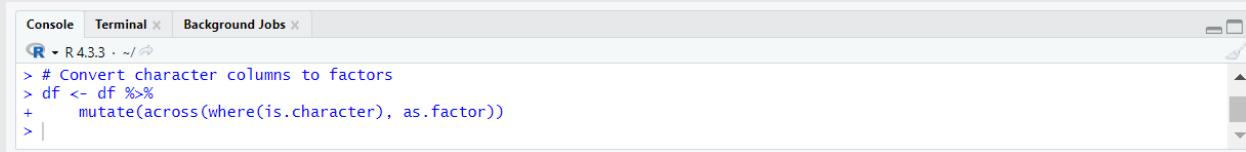
```
Console Terminal Background Jobs
R - R 4.3.3 - ~/Documents
> df$blood_pressure <- ifelse(df$systolic >= 140 | df$diastolic >= 90, "Arterial Hypertension",
+                               ifelse((df$systolic >= 130 & df$systolic <= 139) |
+                                     (df$diastolic >= 85 & df$diastolic <= 89), "High",
+                                     ifelse((df$systolic >= 120 & df$systolic <= 129) |
+                                           (df$diastolic >= 80 & df$diastolic <= 84), "Normal",
+                                           ifelse(df$systolic < 120 & df$diastolic < 80, "Optimal", NA)))
+ )
>
```

I created a new variable, age_group, to categorize individuals based on age: "Young" for ages under 30, "Middle Aged" for ages 30–49, and "Senior" for ages 50 and above.



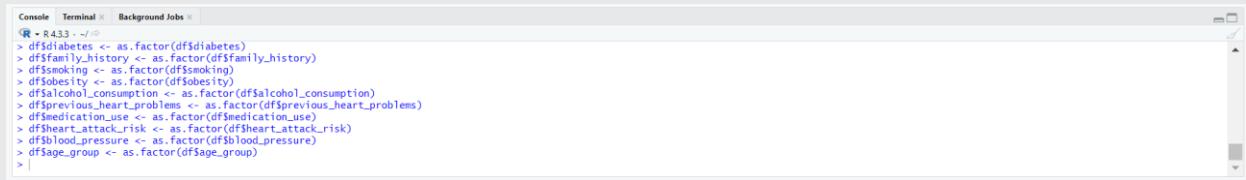
```
Console Terminal Background Jobs
R - R 4.3.3 - ~/Documents
> # creating a new variable "age_group" based on the "age" variable
> df$age_group <- ifelse(df$age < 30, "Young",
+                           ifelse(df$age >= 30 & df$age < 50, "Middle Aged", "Senior"))
>
```

I used 'mutate' with across(where(is.character), as.factor) in R to convert all character columns in my dataset into factors. This transformation ensures that the dataset is optimized for statistical analysis and modeling.



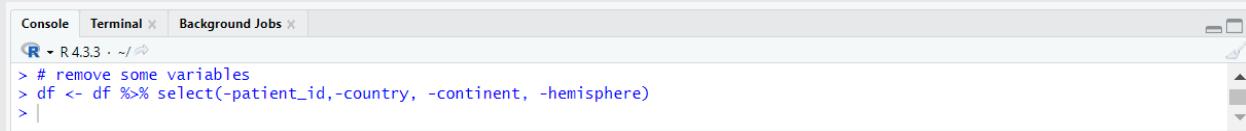
```
Console Terminal Background Jobs
R - R 4.3.3 - ~/ 
> # Convert character columns to factors
> df <- df %>%
+   mutate(across(where(is.character), as.factor))
> |
```

I converted relevant columns to factor type using as.factor() to ensure they are properly recognized as categorical variables during analysis and modeling. This step improves the accuracy of statistical summaries and machine learning algorithms that handle categorical data.



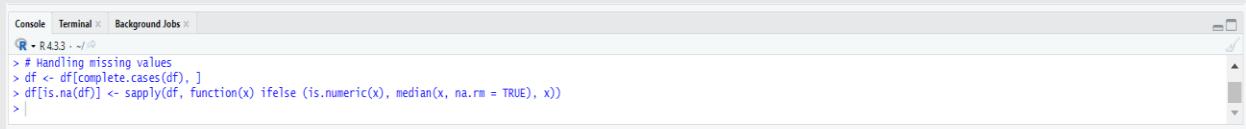
```
Console Terminal Background Jobs
R - R 4.3.3 - ~/ 
> df$diabetes <- as.factor(df$diabetes)
> df$family_history <- as.factor(df$family_history)
> df$smoking <- as.factor(df$smoking)
> df$obesity <- as.factor(df$obesity)
> df$alcohol_consumption <- as.factor(df$alcohol_consumption)
> df$previous_heart_problems <- as.factor(df$previous_heart_problems)
> df$medication_use <- as.factor(df$medication_use)
> df$heart_attack_risk <- as.factor(df$heart_attack_risk)
> df$blood_pressure <- as.factor(df$blood_pressure)
> df$age_group <- as.factor(df$age_group)
> |
```

I removed non-essential columns such as identifiers and geographic details using the select() function with the minus operator. This helped streamline the dataset by keeping only relevant variables for analysis.



```
Console Terminal Background Jobs
R - R 4.3.3 - ~/ 
> # remove some variables
> df <- df %>% select(-patient_id, -country, -continent, -hemisphere)
> |
```

I addressed missing values in the dataset by first removing rows with incomplete cases using complete.cases(). Then, for any remaining NA values, I replaced them with the median for numeric variables, ensuring no data distortion, while non-numeric values were left unchanged.



```
Console Terminal Background Jobs
R - R 4.3.3 - ~/ 
> # Handling missing values
> df <- df[complete.cases(df), ]
> df[is.na(df)] <- sapply(df, function(x) ifelse (is.numeric(x), median(x, na.rm = TRUE), x))
> |
```

Cleaned Data Frame

I checked for any remaining missing values in the dataset using colSums(sapply(df, is.na)). This confirmed that all NA values had been handled, ensuring the dataset is clean and ready for analysis.

```

Console Terminal Background Jobs
R - R 4.3.3 - ~/...
> # Cleaned Data set
> colSums(sapply(df, is.na))
age 0
diabetes 0
family_history 0
previous_heart_problems 0
medication_use 0
stress_level 0
sleep_hours_per_day 0
heart_attack_risk 0
blood_pressure 0
family.history 0
>

```

After cleaning the data, I used `str(df)` to check the structure of the dataset and ensure proper data types. I also reviewed the first few rows with `head(df)` to confirm that the data was correctly processed and formatted, making it ready for further analysis.

```

Console Terminal Background Jobs
R - R 4.3.3 - ~/...
> str(df)
'data.frame': 8763 obs. of 25 variables:
$ age : int 67 21 21 84 66 54 90 84 20 43 ...
$ sex : Factor w/ 2 levels "Female", "Male": 2 2 1 2 2 1 2 2 2 1 ...
$ cholesterol : int 208 389 324 383 318 297 358 220 145 248 ...
$ systolic : num 158 167 174 168 91 172 103 131 144 160 ...
$ diastolic : num 96 99 98 98 98 93 109 109 109 70 ...
$ heart_rate : int 33 59 33 34 54 9 45 56 29 16 ...
$ diabetes : Factor w/ 2 levels "0", "1": 1 2 2 2 2 2 1 1 2 1 ...
$ family.history : Factor w/ 2 levels "1", "2": 1 2 1 2 2 2 2 1 1 2 ...
$ smoking : Factor w/ 2 levels "0", "1": 2 2 1 2 2 2 2 2 2 2 ...
$ obesity : Factor w/ 2 levels "0", "1": 1 2 1 1 2 1 1 2 2 2 ...
$ alcohol_consumption : Factor w/ 2 levels "0", "1": 1 2 1 2 1 2 2 2 1 2 ...
$ exercise_hours_per_week : num 4 4.83 5.88 5.83 5.8 ...
$ diet : Factor w/ 3 levels "Average", "Healthy", "Unhealthy": 1 3 2 1 3 3 2 1 1 3 ...
$ previous_heart_problems : Factor w/ 2 levels "0", "1": 1 2 2 2 2 2 1 1 1 1 ...
$ medication_use : Factor w/ 2 levels "0", "1": 1 1 2 1 1 2 1 2 1 1 ...
$ stress_level : int 9 1 9 9 6 2 7 4 5 ...
$ sedentary_hours_per_day : num 6.62 4.96 9.46 7.65 1.51 ...
$ income : int 261494 285768 352560 125640 160555 241339 190450 122093 25086 209703 ...
$ bmi : num 31.72 27.2 30.5 28.8 ...
$ triglycerides : int 286 235 587 378 231 795 284 370 790 232 ...
$ physical.activity_days_per_week : int 0 1 4 3 1 5 4 6 7 ...
$ sleep_hours_per_day : int 6 7 4 4 10 7 4 7 ...
$ heart_attack_risk : Factor w/ 2 levels "0", "1": 1 1 1 1 2 2 2 1 1 ...
$ blood_pressure : Factor w/ 4 levels "Arterial Hypertension", "Normal", "High", "Senior": 1 1 1 1 2 1 4 2 1 ...
$ age_group : Factor w/ 3 levels "Middle Aged", "Young", "Senior": 2 3 3 2 2 2 2 3 1 ...
>

```

```

Console Terminal Background Jobs
R - R 4.3.3 - ~/...
> head(df)
#> #> #> #> #> #>
#> #> #> #> #> #>
age sex cholesterol systolic diastolic heart_rate diabetes family_history smoking obesity alcohol_consumption exercise_hours_per_week diet previous_heart_problems medication_use
1 67 Male 388 158 88 33 0 1 1 1 0 0 4.963459 Average 0
2 21 Male 385 165 93 59 1 2 2 1 1 1 1 1.813242 Unhealthy 1 0
3 21 Female 324 174 99 33 1 1 0 0 0 0 2.078353 Healthy 1 1
4 84 Male 383 163 100 34 1 2 1 0 1 1 9.828130 Average 1 0
5 66 Male 318 91 88 54 1 2 1 1 1 0 5.804299 Unhealthy 1 0
6 54 Female 297 172 86 9 1 2 1 0 1 1 0.625008 Unhealthy 1 1
#> stress_level sedentary_hours_per_day income bmi triglycerides physical.activity_days_per_week sleep_hours_per_day heart_attack_risk blood_pressure age_group
1 6.615001 261494 31.25123 286 0 6 Arterial Hypertension Senior
2 1 4.963459 285768 27.19497 235 1 7 0 Arterial Hypertension Young
3 9 9.463426 235282 28.17657 587 4 4 0 Arterial Hypertension Young
4 9 7.648981 125640 36.46470 378 3 4 0 Arterial Hypertension Senior
5 6 1.514821 160555 21.80914 231 1 5 0 High Senior
6 2 7.798752 241339 20.14684 795 5 10 1 Arterial Hypertension Senior
>

```

Description of Variables

This table outlines the variables in the "Heart Attack Risk Prediction Dataset". It provides descriptions and data types for each variable.

Variable	Description	Class
Patient ID	Unique identifier for each patient	Character
Age	Age of the patient	Integer
Sex	Gender of the patient (Male/Female)	Character
Cholesterol	Cholesterol levels of the patient	Integer
Blood Pressure	Blood pressure of the patient (systolic/diastolic)	Character
Heart Rate	Heart rate of the patient	Integer
Diabetes	Whether the patient has diabetes (Yes/No)	Integer

Family History	Family history of heart-related problems (1: Yes, 0: No)	Integer
Smoking	Smoking status of the patient (1: Smoker, 0: Non-smoker)	Integer
Obesity	Obesity status of the patient (1: Obese, 0: Not obese)	Integer
Alcohol Consumption	Level of alcohol consumption by the patient (None/Light/Moderate/Heavy)	Integer
Exercise Hours Per Week	Number of exercise hours per week	Numeric
Diet	Dietary habits of the patient (Healthy/Average/Unhealthy)	Character
Previous Heart Problems	Previous heart problems of the patient (1: Yes, 0: No)	Integer
Medication Use	Medication usage by the patient (1: Yes, 0: No)	Integer
Stress Level	Stress level reported by the patient (1-10)	Integer
Sedentary Hours Per Day	Hours of sedentary activity per day	Numeric
Income	Income level of the patient	Integer
BMI	Body Mass Index (BMI) of the patient	Numeric
Triglycerides	Triglyceride levels of the patient	Integer
Physical Activity Days Per Week	Days of physical activity per week	Integer
Sleep Hours Per Day	Hours of sleep per day	Integer
Country	Country of the patient	Character
Continent	Continent where the patient resides	Character
Hemisphere	Hemisphere where the patient resides	Character
Heart Attack Risk	Presence of heart attack risk (1: Yes, 0: No)	Integer

Expectations

This analysis aims to explore the heart attack risk dataset in detail, focusing on identifying the key variables that contribute to predicting heart attack risk. I will examine factors such as age, cholesterol levels, BMI, family history, blood pressure, alcohol consumption, diet, physical activity, etc. These variables will be evaluated using statistical methods, including logistic regression and classification analysis, to determine the most significant predictors of heart attack risk.

To enhance understanding, I will utilize data visualization techniques, such as scatter plots, box plots, and other charts, to identify patterns and trends within the data. These visualizations will aid in discovering the relationships between various health factors and the risk of heart attack.

The goal of this analysis is to generate actionable insights and recommendations for heart attack prevention. By accurately pinpointing the critical risk factors and creating reliable predictive models, the analysis will support the development of strategies to improve heart health and reduce the likelihood of heart attacks.

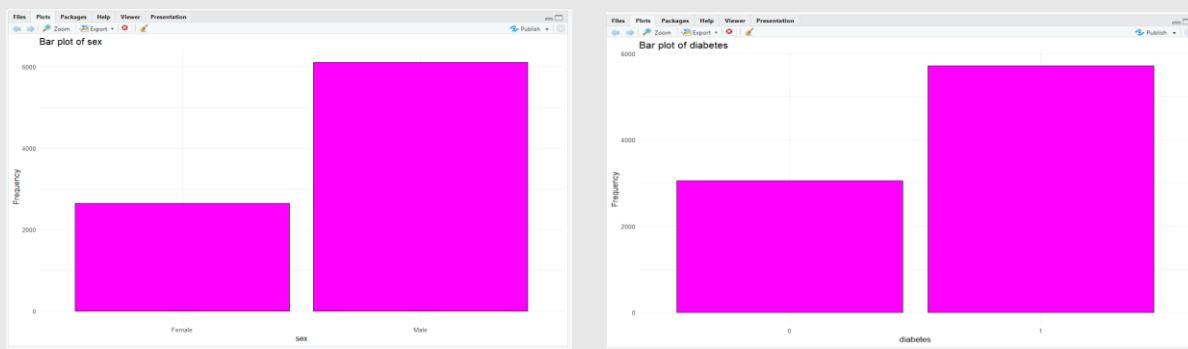
DATA EXPLORATION AND ANALYSIS

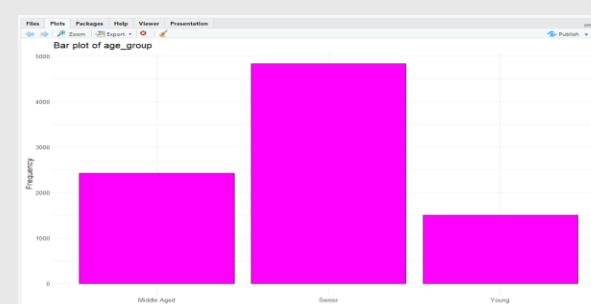
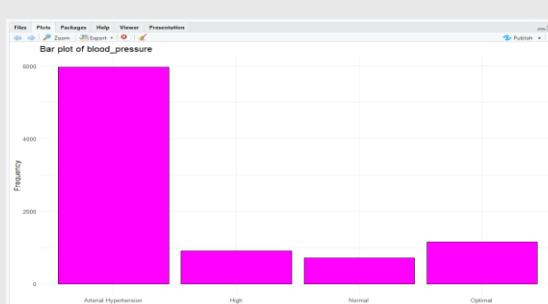
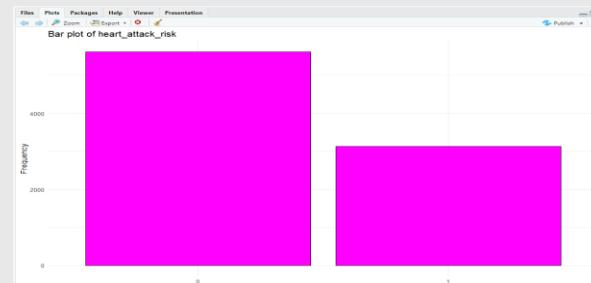
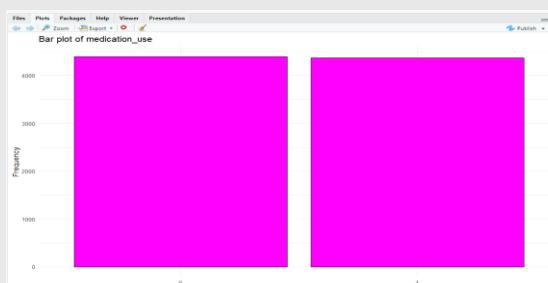
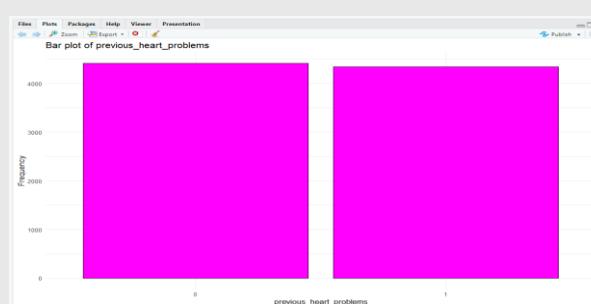
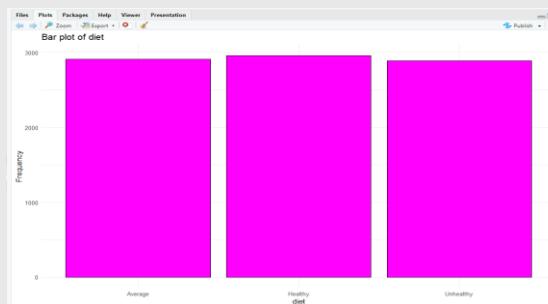
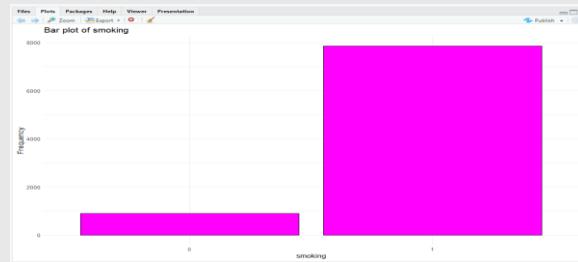
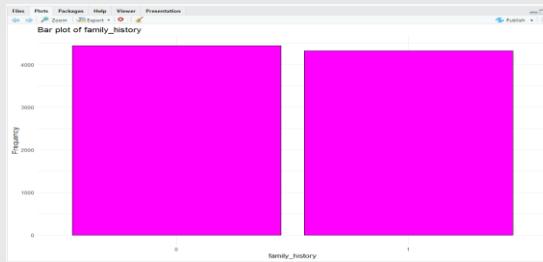
To begin my data exploration of the Heart Attack Risk Prediction Dataset, I used the `summary()` function to examine key variables such as age, cholesterol, heart rate, and diet. This provided basic statistical insights, including the 1st quartile, 3rd quartile minimum, maximum, median, and mean values for the numerical variables. For the categorical variable diet, I observed the frequency distribution across its categories: Average, Healthy, and Unhealthy. This initial analysis gave me a general understanding of the dataset's structure and helped identify patterns worth exploring further.

```
Console Terminal × Background Jobs ×
R - R4.3.3 - ~/r
> #Data Analysis
> summary(select(df, age, cholesterol, heart_rate, diet))
      age      cholesterol      heart_rate      diet
  Min. :18.00  Min. :120.0  Min. :1.00  Average :2912
  1st Qu.:35.00  1st Qu.:192.0  1st Qu.:18.00  Healthy  :2960
  Median :54.00  Median :259.0  Median :36.00  Unhealthy:2891
  Mean   :53.71  Mean   :259.9  Mean   :36.02
  3rd Qu.:72.00  3rd Qu.:330.0  3rd Qu.:54.00
  Max.   :90.00  Max.   :400.0  Max.   :71.00
>
```

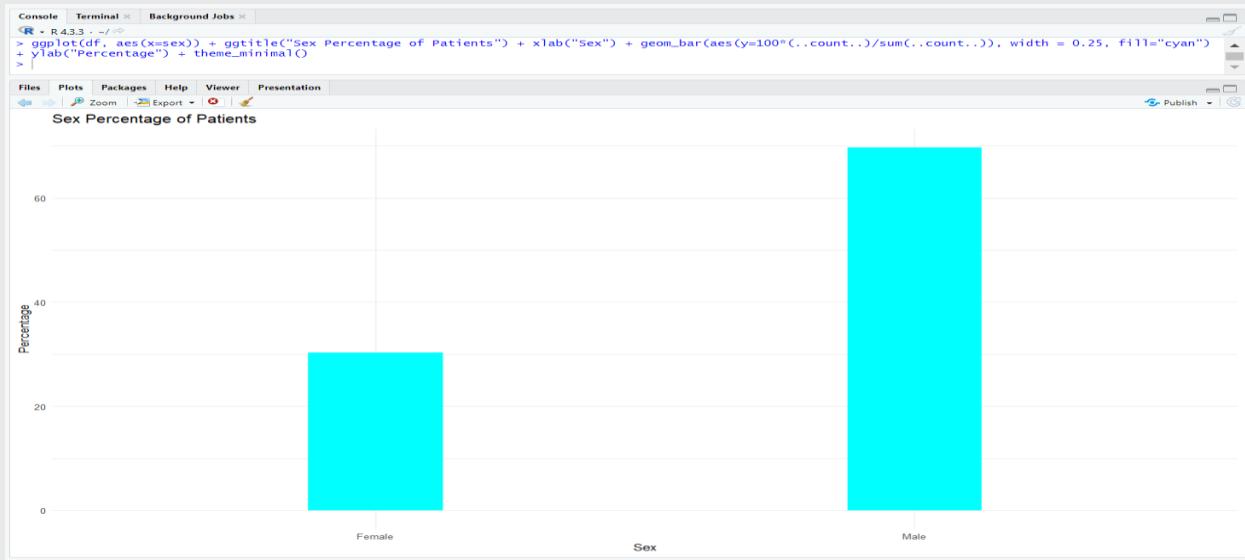
Plots for Categorical Variables

```
Console Terminal × Background Jobs ×
R - R 4.3.3 - ~/r
> # plot for categorical variables
> for (var_name in names(df)) {
+   if(is.factor(df[[var_name]])) {
+     #create a bar plot
+     p <- ggplot(df, aes_string(x = var_name)) +
+       geom_bar(fill = "magenta", color = "black") +
+       theme_minimal() +
+       labs(title = paste("Bar plot of", var_name), x = var_name, y ="Frequency")
+     print(p)
+   }
+ }
```

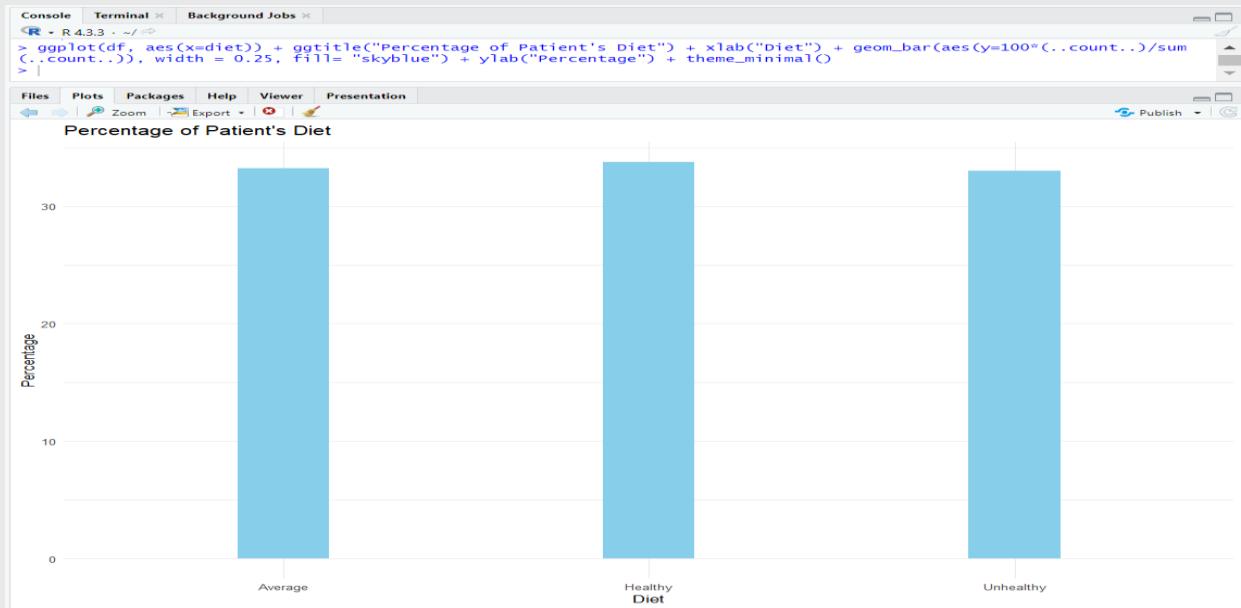




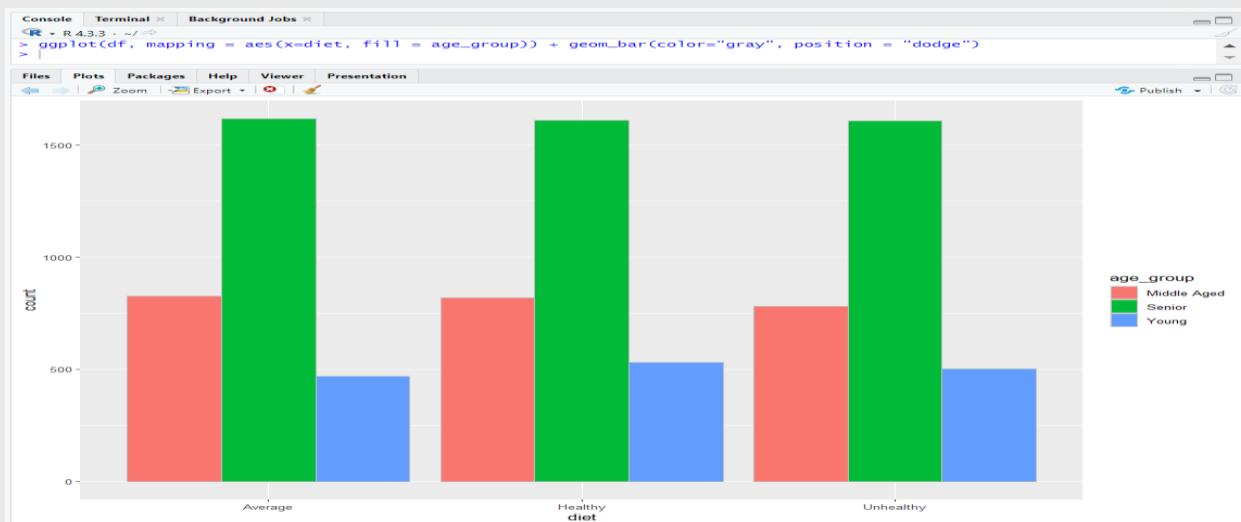
To gain an insight into the sex distribution of patients, I generated a graph that depicts the percentage distribution in the ‘Heart Attack Risk Prediction Dataset’. The results reveal that male patients represent a larger percentage compared to female patients, highlighting a noticeable disparity in gender representation within the dataset. This observation underscores the predominance of male patients in the analyzed data



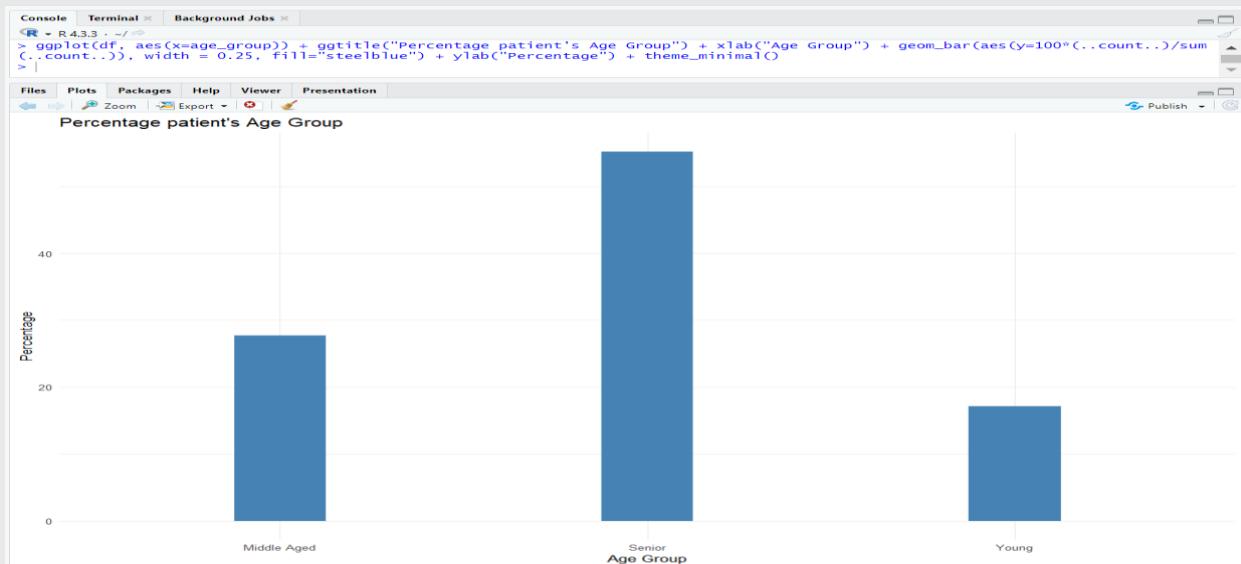
The bar chart below illustrates the percentage breakdown of a patient's dietary intake across three categories. Approximately 35% of the patient's diet consists of food classified as average. Similarly, healthy food choices also constitute about 35% of the patient's diet. The remaining portion, roughly 34%, is categorized as unhealthy. Therefore, the patient's diet is almost evenly split between average, healthy, and unhealthy food options.



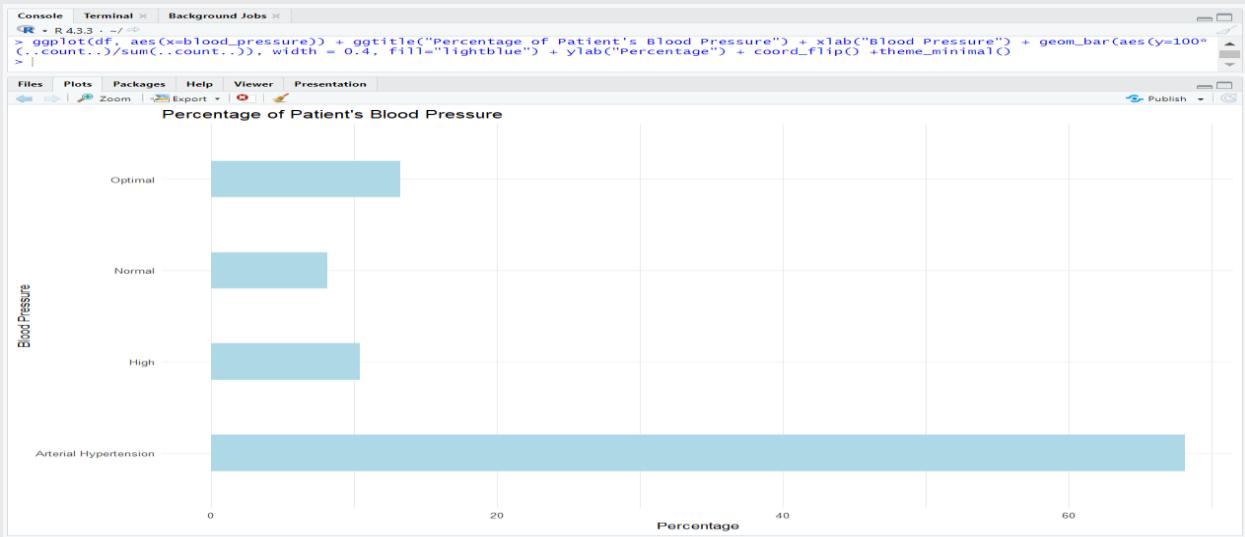
The bar chart shows the distribution of individuals across diet categories (Average, Healthy, Unhealthy), segmented by age group (Young, Middle Aged, Senior). In all diet categories, Seniors have the highest count, followed by Middle Aged individuals, while Young individuals have the lowest. This pattern is consistent across all diet groups.



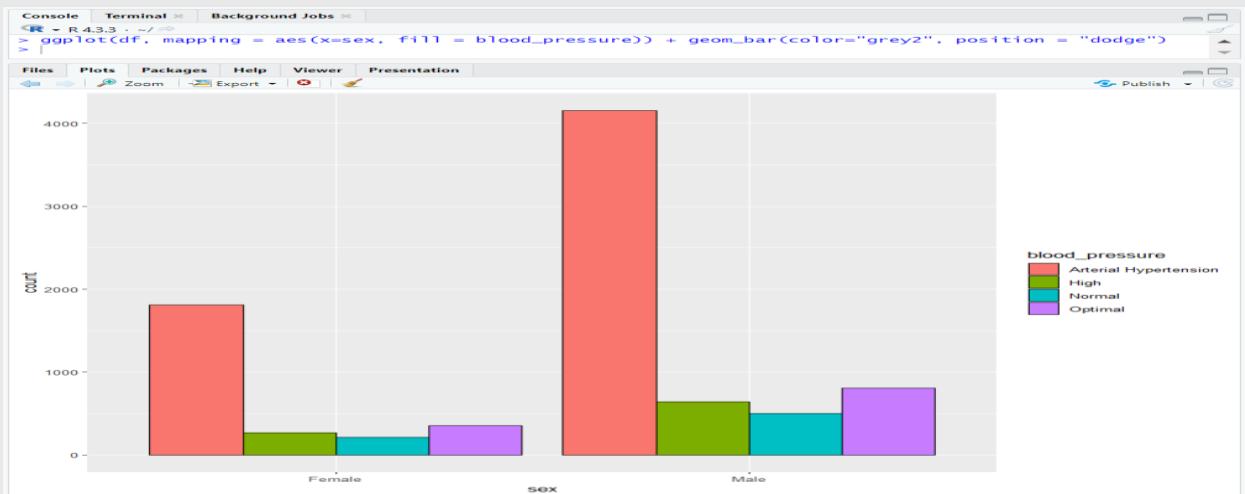
The bar chart displays the percentage distribution of patients across three age groups. The "Middle Aged" group represents approximately 28% of the patients. A significantly larger proportion, around 55%, falls into the "Senior" age group. The "Young" age group constitutes the smallest percentage, at approximately 17%. Thus, most of the patients in this dataset are seniors, followed by middle-aged individuals, with the smallest representation being young patients.



The results indicate that most patients fall into the "Arterial Hypertension" category, accounting for the highest percentage among all blood pressure groups. In comparison, the proportions of patients with "Optimal," "Normal," and "High" blood pressure are significantly lower. This suggests that arterial hypertension is the most prevalent blood pressure condition in the sample, highlighting a potential area of concern for patient health management.



As the graph below illustrates, arterial hypertension is the most prevalent blood pressure category among the patients, with a noticeably higher proportion compared to the "Optimal," "Normal," and "High" blood pressure groups. This pattern is observed in both males and females, emphasizing the widespread nature of hypertension within the studied population.



The graph illustrates patient counts by age group and sex. The Senior age group has the highest patient count, with males significantly outnumbering females. The Middle Aged group also has more males than females, while the Young group has the fewest patients, again with a male predominance.



Based on the two charts, arterial hypertension is a prevalent condition across the patient population, and it appears to disproportionately affect senior males, who represent the largest group of patients overall. Males are the dominant group in both charts: the blood pressure chart shows a higher prevalence of arterial hypertension among males, and the age group chart indicates that males outnumber females in the Senior, Middle Aged, and Young categories. This consistent male predominance suggests a potential correlation between sex and both the overall patient population and the prevalence of arterial hypertension.

Correlation Analysis of Heart Attack Risk and Other Variables

I used the `xtabs()` function in R to create a three-way contingency table showing the frequency of diet, cholesterol levels, and heart attack risk in the Heart Attack Risk Prediction Dataset.

```

Console Terminal Background Jobs
[R - R4.3.3 - ~/]
> diet.cho.df <- xtabs(~diet+cholesterol+heart_attack_risk, data = df)
> diet.cho.df
  heart_attack_risk = 0

  cholesterol
diet   120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166
  Average    7   9   7   9   7   10   9   6   6   6   10   6   8   8   4   5   4   4   6   9   4   8   5   11   3   12   7   4   11   10   7   8   10   3   7   7   6   3   5   8   4   11   3   5   6   8   5
  Healthy    4   5   4   5   12   11   7   9   8   13   10   2   8   8   7   7   3   4   1   3   7   6   6   5   12   13   6   7   5   10   6   10   5   5   9   7   7   8   6   11   7   8   10   11   4   10
  Unhealthy   6   8   7   7   6   4   9   8   4   8   5   5   2   2   7   8   6   9   9   4   5   5   10   4   5   8   8   6   12   10   5   5   10   9   6   5   7   5   9   2   6   7   5   10   6   8

  cholesterol
diet   167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213
  Average    6   11   8   3   3   12   6   4   4   4   7   4   8   9   9   10   9   4   3   10   3   4   7   4   11   4   8   9   11   5   9   7   8   6   4   15   2   6   3   10   6   5   8   7   8   7   8   9
  Healthy    5   5   5   6   6   4   5   10   5   6   9   9   1   7   6   5   5   7   8   7   7   5   4   6   6   2   9   10   5   9   8   10   3   9   8   7   4   6   10   6   9   8   12   6   4   3   2   7   5   6
  Unhealthy   5   5   5   6   6   4   5   10   5   6   9   9   1   7   6   5   5   7   8   7   7   5   4   6   6   2   9   10   5   9   8   10   3   9   8   7   4   6   10   6   9   8   12   6   4   3   2   7   5   6

  cholesterol
diet   214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260
  Average    11   7   9   7   8   6   6   9   8   1   6   8   7   4   6   4   9   4   5   4   6   10   9   3   5   8   6   13   8   7   11   6   5   5   8   1   10   6   8   12   7   5   4   5   11   8   6
  Healthy    3   3   8   8   17   2   6   11   2   6   9   8   4   9   7   6   7   6   11   11   13   11   10   5   5   10   4   8   10   6   4   9   6   2   10   9   11   7   8   7   5   14   8   5   8   7
  Unhealthy   7   2   8   6   5   7   6   5   3   9   4   7   7   6   4   7   2   9   12   7   4   7   6   8   6   3   5   7   3   7   9   8   9   6   13   5   10   9   6   9   7

  cholesterol
diet   261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307
  Average    6   9   6   5   10   5   10   2   7   6   4   7   10   8   5   5   8   4   8   6   5   11   5   6   7   8   4   7   6   5   10   8   10   8   6   5   1   4   8   5   10   2   5   7
  Healthy    6   8   7   6   4   5   5   7   6   5   5   11   10   9   8   4   7   2   4   4   7   2   5   8   5   3   8   4   9   7   9   3   4   10   11   8   10   4   9   9   8   12   9   6   4   8   2
  Unhealthy   10   4   7   8   5   10   9   7   4   12   6   11   4   7   8   4   5   3   6   3   7   6   5   11   7   3   4   9   7   8   6   3   5   9   4   12   7   8   6   5   7   11   5   9   2   5

  cholesterol
diet   308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 354
  Average    15   7   8   9   7   1   3   8   5   3   8   9   6   10   8   9   5   3   7   7   3   6   3   6   11   10   9   6   10   4   3   5   7   5   12   5   6   8   4   8   3   7   4   8   4
  Healthy    3   2   8   6   7   8   5   5   8   7   3   5   10   9   6   11   4   6   4   4   6   5   8   7   5   3   5   7   3   6   10   11   2   3   11   2   3   7   8   8   2   4   6   4   10   9
  Unhealthy   4   5   6   5   8   7   7   8   6   2   10   4   7   12   2   6   2   10   7   5   6   9   6   4   7   5   12   5   5   6   6   6   6   7   5   5   3   8   6   5   2   9   5   7

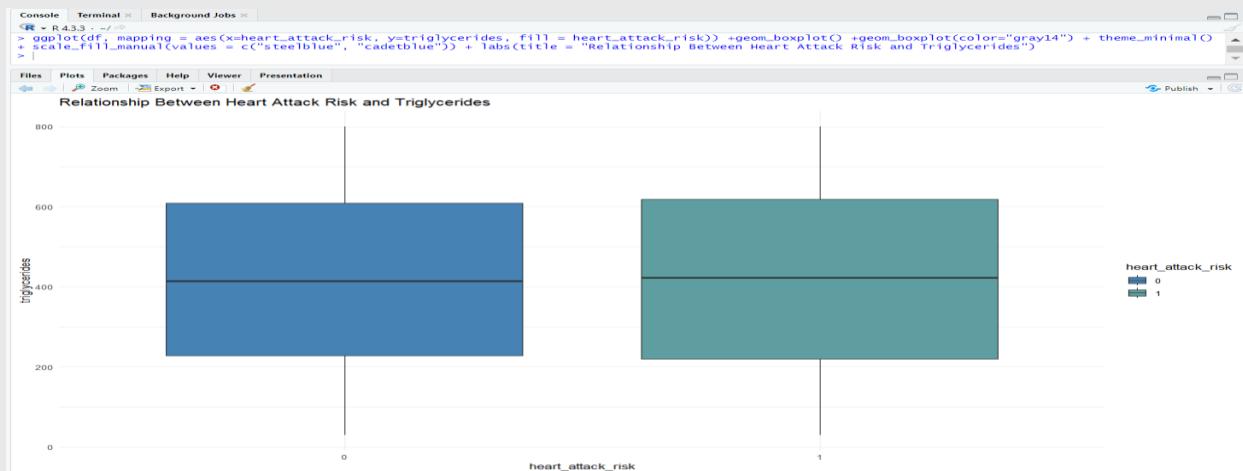
  cholesterol
diet   355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400
  Average    7   4   12   8   7   13   3   12   3   7   10   3   6   5   6   5   10   9   9   5   5   8   5   6   7   6   4   6   8   4   11   4   4   2   4   9   4   11   5   7   4   3   6   7   16   6
  Healthy    4   4   5   7   5   4   8   5   9   12   6   7   8   5   6   7   1   7   6   6   7   5   9   7   6   12   8   7   10   9   6   4   3   8   3   3   4   4   4   11   4   9   1   3   9   8
  Unhealthy   6   8   5   9   8   12   7   9   8   4   6   3   4   4   7   10   8   8   12   3   9   8   6   10   8   9   6   10   5   9   9   6   7   6   6   8   10   4   3   4   4
```

The marginal frequencies of diet types were extracted using `margin.table(diet.cho.df, 1)`. The resulting data frame shows that 2,960 individuals follow a healthy diet, 2,912 follow an average diet, and 2,891 follow an unhealthy diet. The distribution is relatively balanced across the three diet categories.

```

Console Terminal Background Jobs
[R - R4.3.3 - ~/]
> diet.cho.heart.df <- as.data.frame(margin.table(diet.cho.df, 1))
> diet.cho.heart.df
  diet
  1 Average 2912
  2 Healthy 2960
  3 Unhealthy 2891
>
```

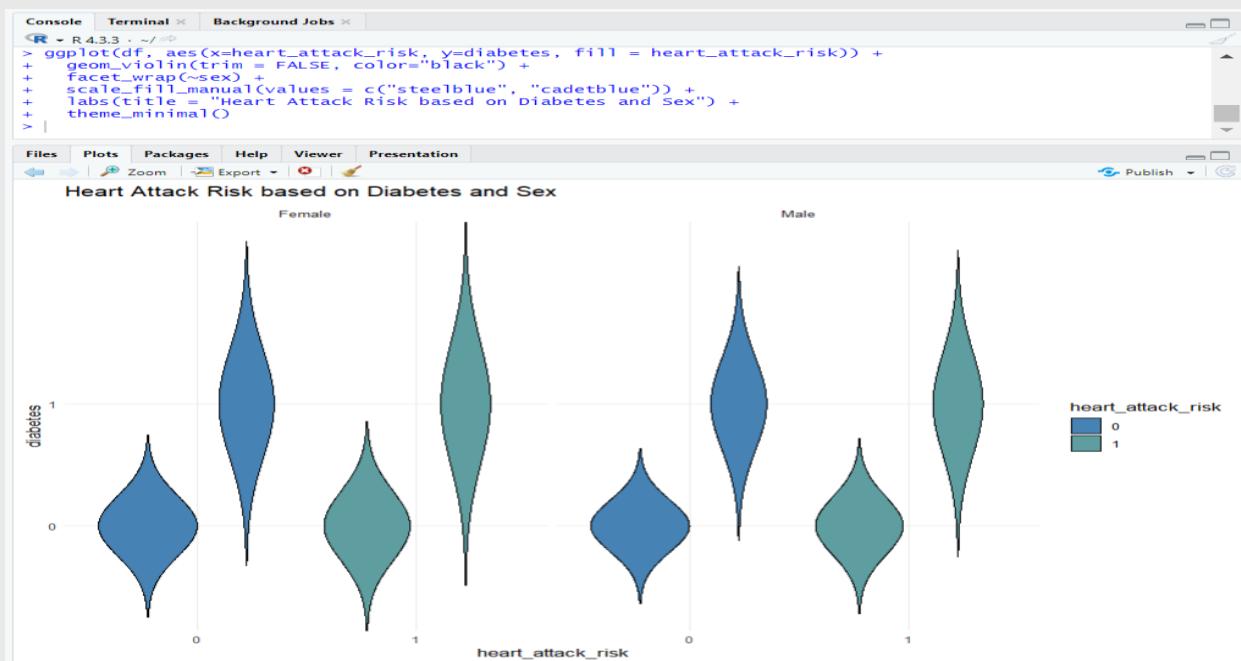
Based on the boxplot, similar median triglyceride levels are observed for both groups, indicating that those with heart attack risk (1) do not have markedly different triglyceride levels compared to those without (0). While the group with heart attack risk displays a slightly wider range of triglyceride values, this difference isn't substantial enough to suggest a strong relationship.



The box plot compares cholesterol levels across heart attack risk categories (0 = no risk, 1 = risk) by sex. It shows similar median cholesterol levels for both males and females, regardless of risk. However, females with risk exhibit slightly more variability in cholesterol levels. No strong correlation is visually evident.

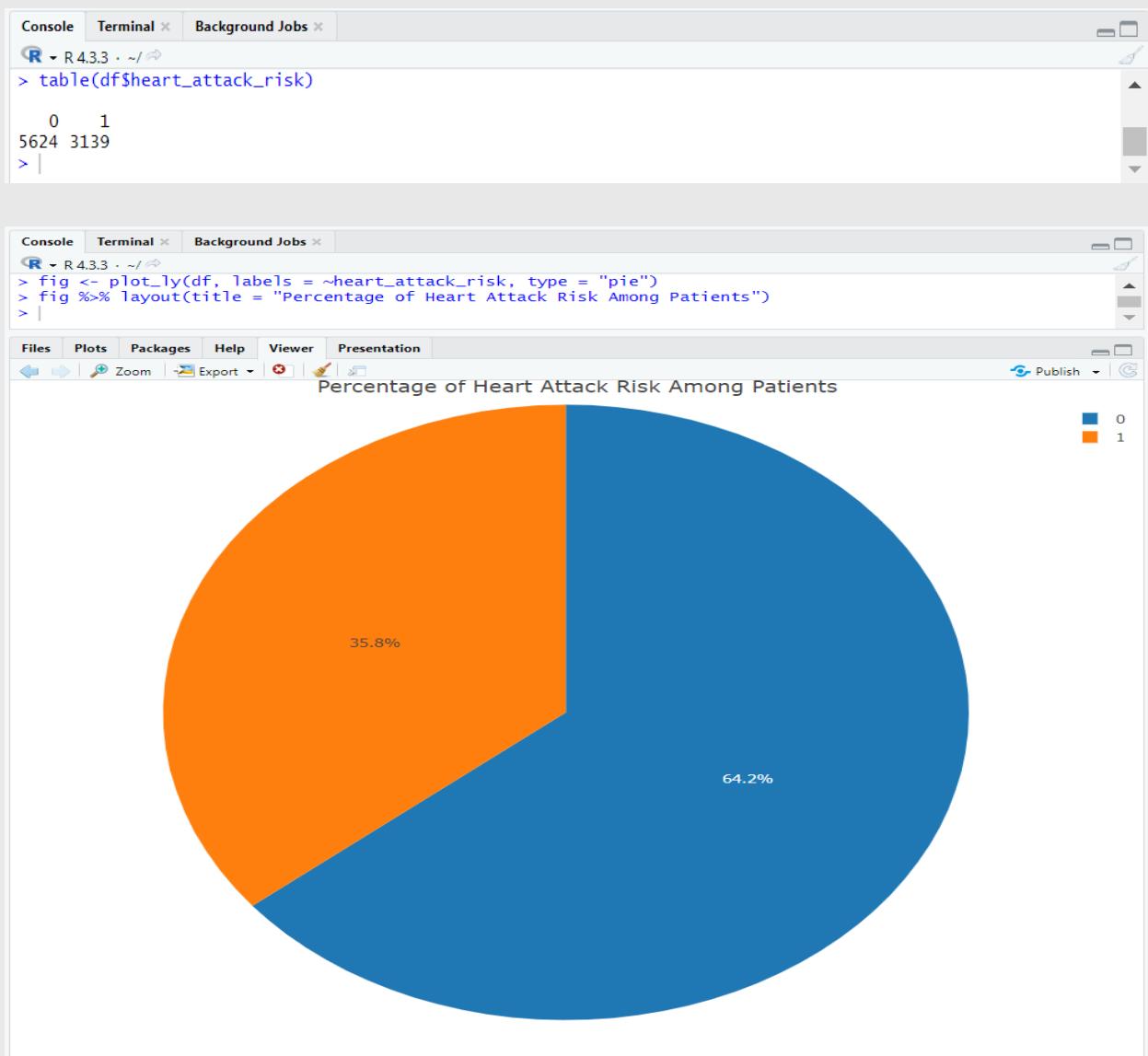


The plot below shows that individuals with high heart attack risk (1) have higher diabetes levels compared to those with no risk (0). In males, diabetes levels are more variable in the high-risk group, while females show a narrower distribution. The trend of higher diabetes prevalence in high-risk individuals is consistent across both sexes.



Distributional Analysis of the Target Variable

To better understand the "heart_attack_risk" variable, I used the `table(df$heart_attack_risk)` function to examine its distribution. The variable is binary, with two categories: 0 for low risk and 1 for high risk. The data reveals a higher proportion of individuals in the low-risk group (0) compared to the high-risk group (1). This insight helps set the stage for further analysis, particularly when exploring factors influencing heart attack risk.



Summary Statistics of Variables

I used the `summary(df)` function in R to get an overview of the dataset. This provided key statistics for each variable, including minimum, maximum, mean, and quartiles for numeric variables, as well as frequency counts for categorical ones.

```
Console Terminal X Background Jobs X
[R - R4.3.3 - ~]
> summary(df)
#> #> age sex cholesterol systolic diastolic heart_rate diabetes family_history smoking obesity alcohol_consumption exercise_hours_per_week diet
#> Min. :18.00 Female:2652 Min. :120.0 Min. :90.0 Min. :60.00 Min. :1.00 0:3047 1:4443 0: 904 0:4369 0:3522 Min. : 0.002442 Average :2912
#> 1st Qu.:35.00 Male:6111 1st Qu.:112.0 1st Qu.:111.0 1st Qu.:79.00 1st Qu.:18.00 1st Qu.:5716 2:4320 1:7859 1:4394 1:5241 1st Qu.: 0.002442 Average :2912
#> Median :45.00 Median:135.0 Median:131.0 Median:105.00 Median:25.00 Median:36.00 Median:17.00 Median:69.559 Healthy :2960
#> Mean :53.71 Mean :259.9 Mean :135.1 Mean :85.00 Mean :36.02 Mean :10.04284 Unhealthy:2891
#> 3rd Qu.:72.00 3rd Qu.:330.0 3rd Qu.:158.0 3rd Qu.:98.00 3rd Qu.:54.00 3rd Qu.:15.0008 Max. :19.998709
#> Max. :90.00 Max. :400.0 Max. :180.0 Max. :71.00 Max. :7.00
#> previous_heart_problems medication_use stress_level sedentary_hours_per_day income bmi triglycerides physical_activity_days_per_week sleep_hours_per_day heart_attack_risk
#> 0:418 0:4396 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367
#> 1:4345 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367 1:4367
#>
#> blood_pressure age_group
#> Arterial Hypertension:591 Middle Age:2426
#> High : 912 Senior :4835
#> Normal : 717 Young :1502
#> Optimal :1163 Max. :11.999313
#>
#> |
```

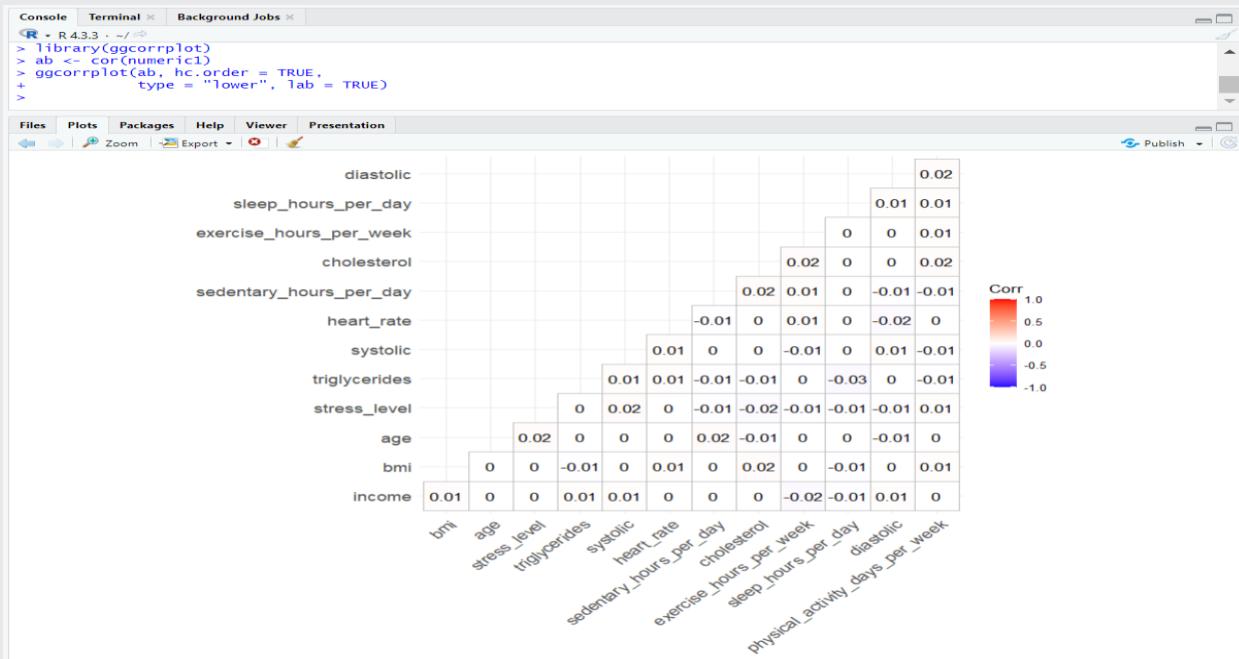
Correlation of Numerical Variables

I analyzed correlations among numeric variables in the dataset using the `cor()` function and visualized them with `ggcorrplot`. The matrix revealed mostly weak relationships, with no coefficients exceeding ± 0.7 .

There was a mild positive correlation between systolic and diastolic blood pressure, and a weak association between exercise hours and physical activity days. BMI and income also showed a slight positive correlation. On the negative side, stress level and sleep hours had a weak inverse relationship, as did sedentary hours and physical activity days.

Overall, the data showed weak linear associations.

```
Console Terminal X Background Jobs X
[R - R4.3.3 - ~]
> numeric1 <- select_if(df, is.numeric)
> cor(numeric1)
#> #> age cholesterol systolic diastolic heart_rate exercise_hours_per_week stress_level sedentary_hours_per_day income bmi
#> age 1.00000000 -9.107011e-03 0.0030702087 -0.009825617 -0.003840129 0.001205639 0.018306646 1.728013e-02 -1.732790e-03 -2.611846e-03
#> cholesterol -0.009107011 1.0000000e+00 0.0001326976 0.002083335 0.0003149083 0.021517136 -0.024487111 1.891449e-02 6.750208e-06 1.729187e-02
#> systolic -0.009107009 1.326976e-04 1.00000000 0.01337009 0.0084818882 -0.009505361 0.017847831 3.392607e-03 1.041436e-02 4.279043e-03
#> diastolic -0.009107011 0.0001326976 1.00000000 0.002083335 0.0003149083 -0.009505361 0.017847831 3.392607e-03 1.041436e-02 4.279043e-03
#> heart_rate -0.003844013 2.114964e-04 0.0084818882 -0.018113057 1.000000000 -0.008276329 -0.004546769 1.000000000 -0.009102419 8.755601e-03 -2.341385e-02 3.776921e-03
#> exercise_hours_per_week 0.001205639 2.151714e-02 -0.0095055610 -0.003468859 0.0082761293 1.000000000 -0.009102419 1.000000000 -0.009102419 8.755601e-03 -2.341385e-02 3.776921e-03
#> stress_level 0.018306646 -2.448711e-02 0.0178478307 -0.008445057 -0.004546768 -0.009102419 1.000000000 5.397241e-03 -2.760451e-03 -3.250447e-03
#> sedentary_hours_per_day 0.017280134 1.891449e-02 0.0033926071 -0.006606069 -0.0102320484 0.008755601 -0.005397241 1.000000000 3.510621e-03 -2.356074e-05
#> income -0.0026118946 1.729187e-02 0.0042790432 0.000805527 0.0052985748 0.003776921 -0.002250447 3.510621e-03 8.835838e-03 8.110375e-04
#> bmi -0.003921303 0.003414957 -0.453721e-03 0.0051206953 0.000544911 0.0122436948 0.001716949 -0.003921303 -2.156074e-05 8.835838e-03 0.000000e+00
#> triglycerides 0.001383668 1.605594e-02 -0.0075739574 0.016294383 0.0008343817 0.007725186 0.007404630 5.784609e-03 1.073856e-02 -5.963607e-03
#> physical_activity_days_per_week 0.001383668 1.605594e-02 -0.0075739574 0.016294383 0.0008343817 0.007725186 0.0074046306 -6.178012e-03 1.302733e-04 8.110375e-03
#> sleep_hours_per_day -0.002184704 4.456229e-03 -0.0046276733 0.010679456 0.0018112469 0.001245336 -0.014205407 4.792013e-03 -6.598343e-03 -1.003041e-02
#>
#> Total correlation matrix physical_activity_days_per_week sleep_hours_per_day
#> age 0.002416957 0.001383668 -0.002184704
#> cholesterol -0.004543721 0.0160559355 -0.004456229
#> systolic 0.005120695 -0.0075739574 -0.004627763
#> diastolic 0.005544911 0.0162943837 -0.010679456
#> heart_rate 0.003921303 0.003414957 -0.003844013
#> exercise_hours_per_week 0.001716949 0.003921303 -0.007251861
#> stress_level -0.003921303 0.003414957 -0.0074046302
#> sedentary_hours_per_day -0.005784609 -0.0061780115 0.004792013
#> income 0.001383668 0.001383668 -0.006539817
#> bmi -0.003921303 0.008755601 -0.010679456
#> triglycerides 1.000000000 -0.0075564192 -0.029215971
#> physical_activity_days_per_week -0.007556419 1.000000000 0.0140334348
#> sleep_hours_per_day -0.029215971 1.000000000 0.0140334348
#>
```



Plots for Correlated Variables

I analyzed the relationship between systolic and diastolic blood pressure using Pearson's correlation. The result ($r = 0.013$) indicates no significant linear association between the two. This is unusual, as systolic and diastolic values are typically related, which may suggest data quality issues or a non-linear relationship. The result is shown in the plot below.

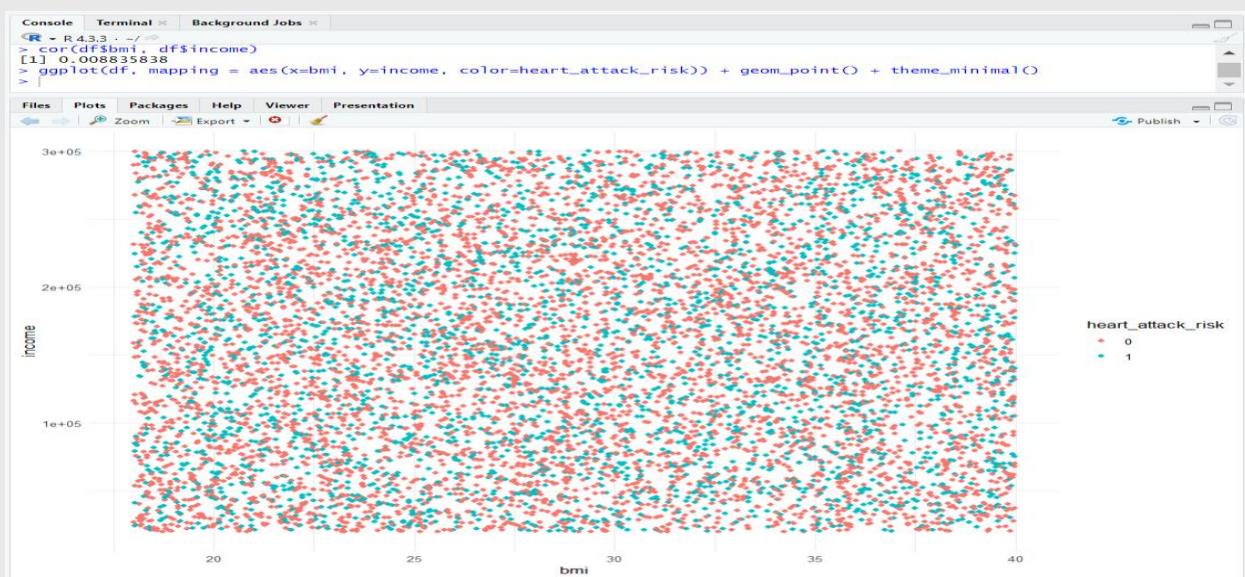


Based on the scatterplot of exercise hours per week versus physical activity days, there is no visible trend. The data points are widely scattered across the 0–6 days and 0–20 hours range. This visual observation aligns with the Pearson correlation coefficient of 0.008, indicating no linear relationship.



The plot below shows a near-zero linear relationship between BMI and income. BMI ranges from 20 to 40, and income values are unclear due to formatting issues. Heart attack risk is indicated by color, but no clear clustering by risk group is observed.

The weak correlation is counterintuitive, suggesting potential non-linear relationships.



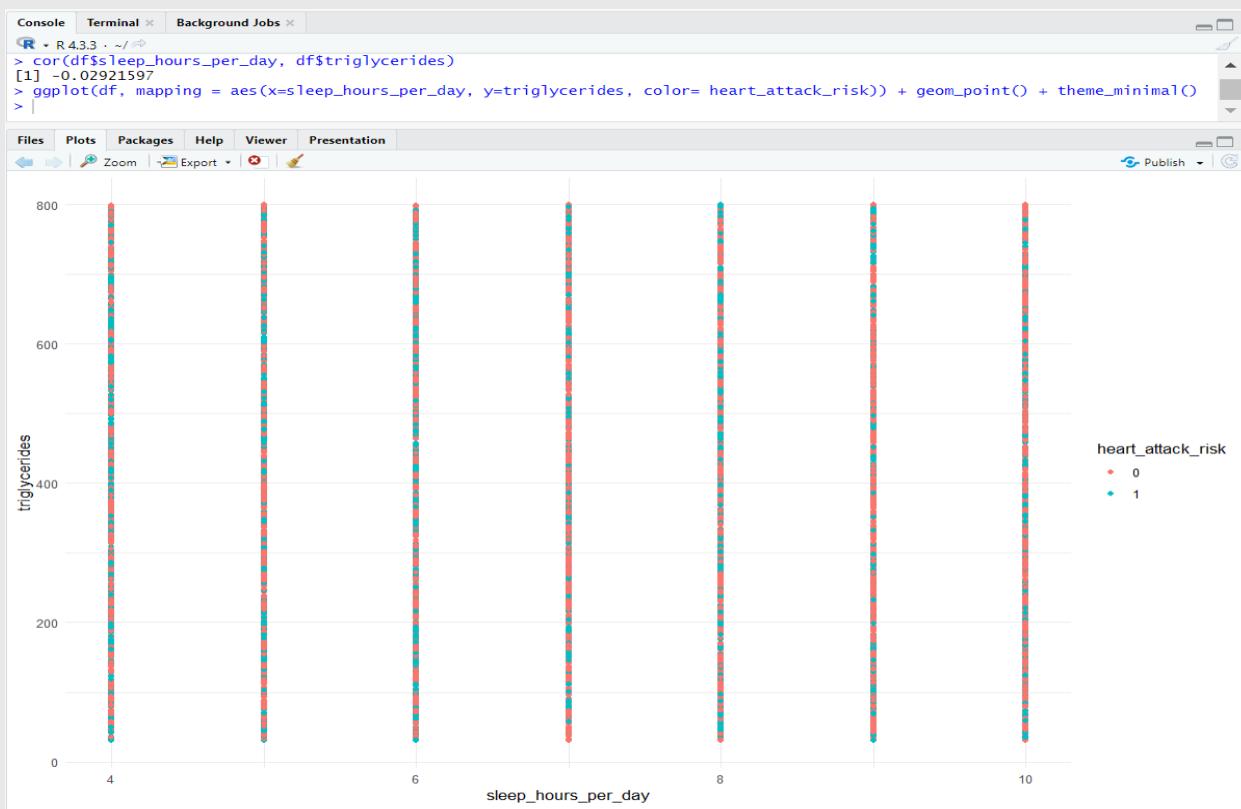
The analysis found a near-zero negative correlation ($r = -0.014$) between stress levels and sleep hours per day, indicating no linear relationship. The scatterplot (stress: 0–10; sleep: 2.5–10) showed no clear trend, even when colored by heart attack risk.



A near-zero correlation ($r = -0.006$) between sedentary hours and activity days per week indicates no linear relationship.

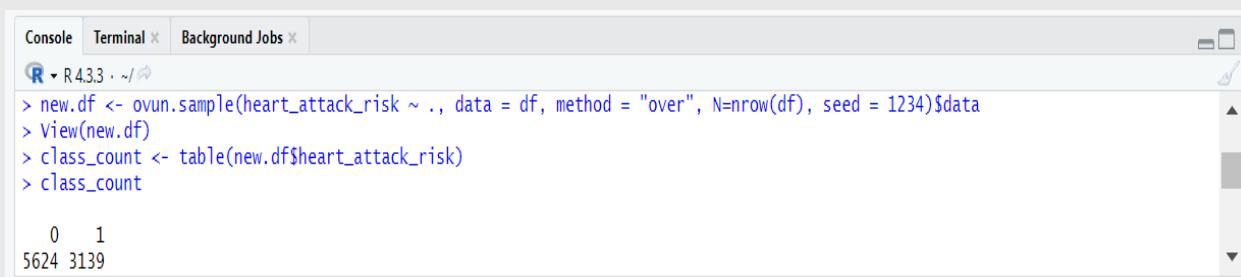


A near-zero negative correlation ($r = -0.029$) was found between sleep hours and triglycerides, suggesting no linear relationship. The scatterplot showed no clear trend.



Predictive Modeling

To support the analysis, I loaded the ROSE package, which had been previously installed at the start of this project. This package was employed to address class imbalance in the Heart Attack Risk variable by oversampling the minority class using the `ovun.sample()` function. After resampling, I used the `table()` function to examine the distribution of instances across the different Heart Attack Risk classes.



The code below was used to divide the dataset into training and testing subsets. The split enabled model development on the training data to predict the target variable, "Heart Attack Risk" and subsequent evaluation of the model's accuracy on the test data.



A screenshot of the RStudio interface showing the Console tab. The console window displays R code for splitting a dataset. The code includes setting a seed, creating a sample index, and splitting the dataset into training and testing subsets. The R logo icon and the text "R 4.3.3 · ~/..." are visible at the top left of the console area.

```
R 4.3.3 · ~/...
> set.seed(1234)
> new.df_split <- sample(x=nrow(new.df), size = .70*nrow(new.df))
> train <- new.df[new.df_split,]
> test<- new.df[-new.df_split,]
>
```

Logistic Regression

To identify the most effective logistic regression model for predicting Heart Attack Risk in the Heart Attack Risk Prediction Dataset, I developed multiple models using various predictor variables. The objective was to compare their performance using Mean Squared Error (MSE). The initial model incorporated selected features from the dataset, with Heart Attack Risk as the target variable. The assumption was that the predictors in the "new.df" dataset would contribute meaningfully to the prediction. Model performance was assessed using the validation set approach, based on MSE values.

Hypothesis Testing

This analysis aims to predict Heart Attack Risk using the available variables in the dataset, where the response variable is categorical.

Model 1

The model did not come with deviance residual, so I created it separately.

```

Console Terminal Background Jobs
R - R 4.3.3 . ~/ ◊
> glm.fit <- glm(heart_attack_risk~., data = train, family = "binomial")
> summary(glm.fit)

Call:
glm(formula = heart_attack_risk ~ ., family = "binomial", data = train)

Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) -8.881e-01 4.452e-01 -1.995 0.046088 *
age 1.144e-02 2.871e-03 3.983 6.8e-05 ***
sexMale 8.914e-02 7.007e-02 1.272 0.203291
cholesterol 1.008e-03 3.305e-04 3.049 0.002299 **
systolic -5.054e-05 1.367e-03 -0.037 0.970511
diastolic -2.899e-03 2.343e-03 -1.237 0.215904
heart_rate 7.117e-04 1.300e-03 0.547 0.584177
diabetes1 1.581e-02 5.611e-02 0.282 0.778119
family_history -1.312e-01 5.370e-02 -2.443 0.014556 *
smoking1 -9.383e-02 1.166e-01 -0.805 0.420844
obesity1 1.325e-02 5.369e-02 0.247 0.805108
alcohol_consumption1 -1.782e-01 5.427e-02 -3.285 0.001021 **
exercise_hours_per_week 1.044e-03 4.590e-03 0.227 0.820085
dietHealthy -3.930e-02 6.438e-02 -0.610 0.541605
dietUnhealthy -2.306e-01 6.643e-02 -3.472 0.000517 ***
previous_heart_problems1 3.300e-02 5.367e-02 0.615 0.538613
medication_use1 9.758e-02 5.372e-02 1.816 0.069320 .
stress_level 5.753e-03 9.352e-03 0.615 0.538463
sedentary_hours_per_day 1.098e-04 7.794e-03 0.014 0.988758
income 1.391e-07 3.328e-07 0.418 0.675935
bmi 8.731e-03 4.225e-03 2.066 0.038802 *
triglycerides -8.047e-05 1.201e-04 -0.670 0.502807
physical_activity_days_per_week -9.139e-03 1.177e-02 -0.776 0.437687
sleep_hours_per_day -1.335e-02 1.355e-02 -0.986 0.324279
blood_pressureHigh -1.319e-01 1.024e-01 -1.289 0.197466
blood_pressureNormal -2.567e-01 1.171e-01 -2.193 0.028335 *
blood_pressureOptimal -8.541e-02 1.188e-01 -0.719 0.472085
age_groupSenior -3.060e-01 1.097e-01 -2.788 0.005305 **
age_groupYoung 2.033e-01 9.563e-02 2.126 0.033544 *
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 8013.6 on 6133 degrees of freedom
Residual deviance: 7936.9 on 6105 degrees of freedom
AIC: 7994.9

Number of Fisher Scoring iterations: 4

```

```

Console Terminal Background Jobs
R - R 4.3.3 . ~/ ◊
> # Capture the summary
> model_summary <- summary(glm.fit)
> # Print deviance residuals manually (mimicking what summary usually shows)
> cat("Deviance Residuals:\n")
Deviance Residuals:
> print(summary(residuals(glm.fit, type = "deviance")))
  Min. 1st Qu. Median Mean 3rd Qu. Max.
-1.2425 -0.9606 -0.8651 -0.0894 1.3656 1.7407

```

```

Console Terminal Background Jobs
R - R 4.3.3 . ~/ ◊
> probs <- predict(glm.fit, test, type = "response")
> predict <- rep("No", length(probs))
> predict[probs > 0.5] <- "Yes"
> table(predict, test$heart_attack_risk)

predict   0   1
  No 1688 924
  Yes 8   9
> mean(predict != test$heart_attack_risk)
[1] 1
>

```

A logistic regression model was developed using the training dataset to predict Heart Attack Risk based on all available variables. The model was fit using the binomial family. Several predictors were found to be statistically significant at the 5% level. These include age, cholesterol, family history, alcohol consumption, unhealthy diet, BMI, blood pressure (normal), age group (senior), and age group (young). Specifically, age and cholesterol were positively associated with heart attack risk. Interestingly, alcohol

consumption and being in the senior age group were associated with reduced risk, which may reflect confounding lifestyle or health behavior factors. An unhealthy diet and higher BMI were also linked to an increased risk of heart attack.

Forward Selection

```
Console Terminal × Background Jobs ×
R 4.3.3 · ~/r
> fwd.set = regsubsets(heart_attack_risk~, data = new.df, nvmax = 10, method = "forward")
>
> summary(fwd.set)
Subset selection object
Call: regsubsets.formula(heart_attack_risk ~ ., data = new.df, nvmax = 10,
  method = "forward")
28 variables (and intercept)
```

Backward Selection

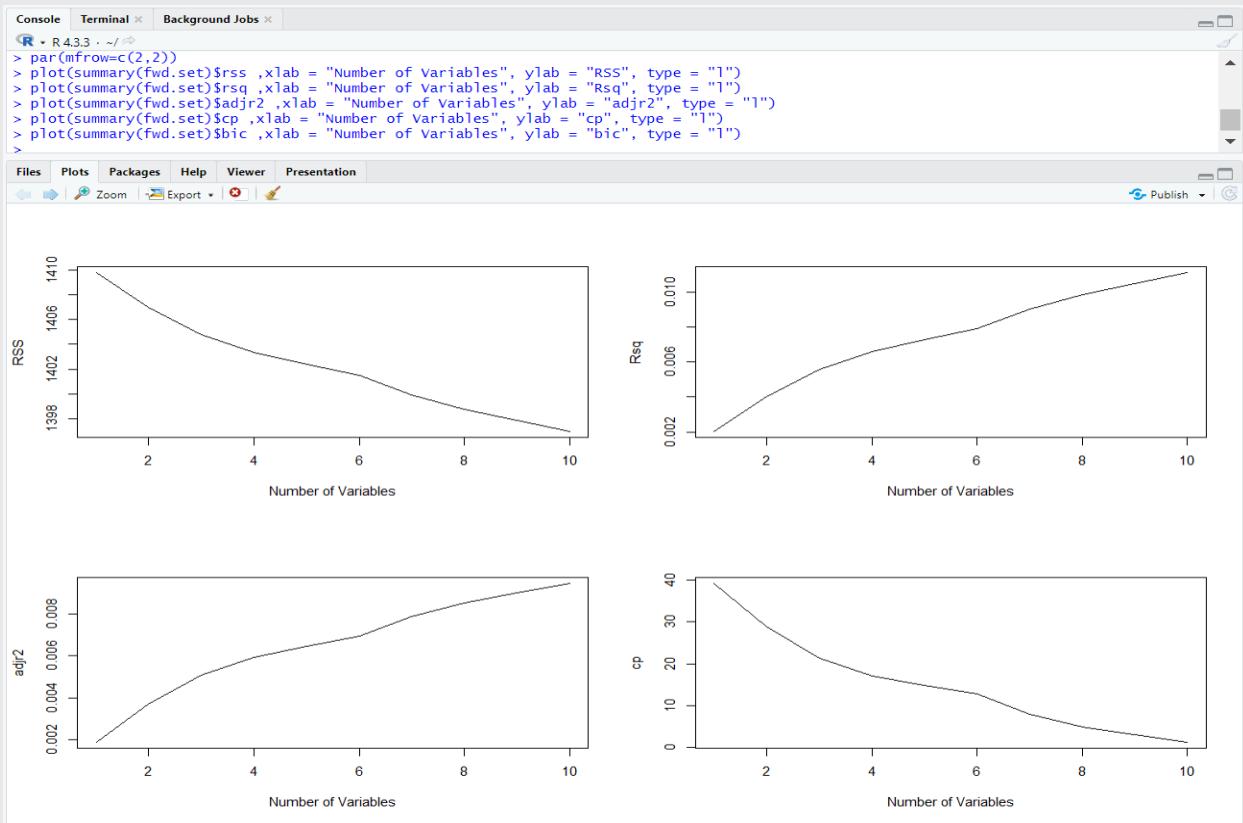
```
Console Terminal × Background Jobs ×
R 4.3.3 · ~/r
> bkd.set <- regsubsets(heart_attack_risk~, data = new.df, nvmax = 10, method = "backward")
> summary(bkd.set)
Subset selection object
Call: regsubsets.formula(heart_attack_risk ~ ., data = new.df, nvmax = 10,
  method = "backward")
28 variables (and intercept)
```

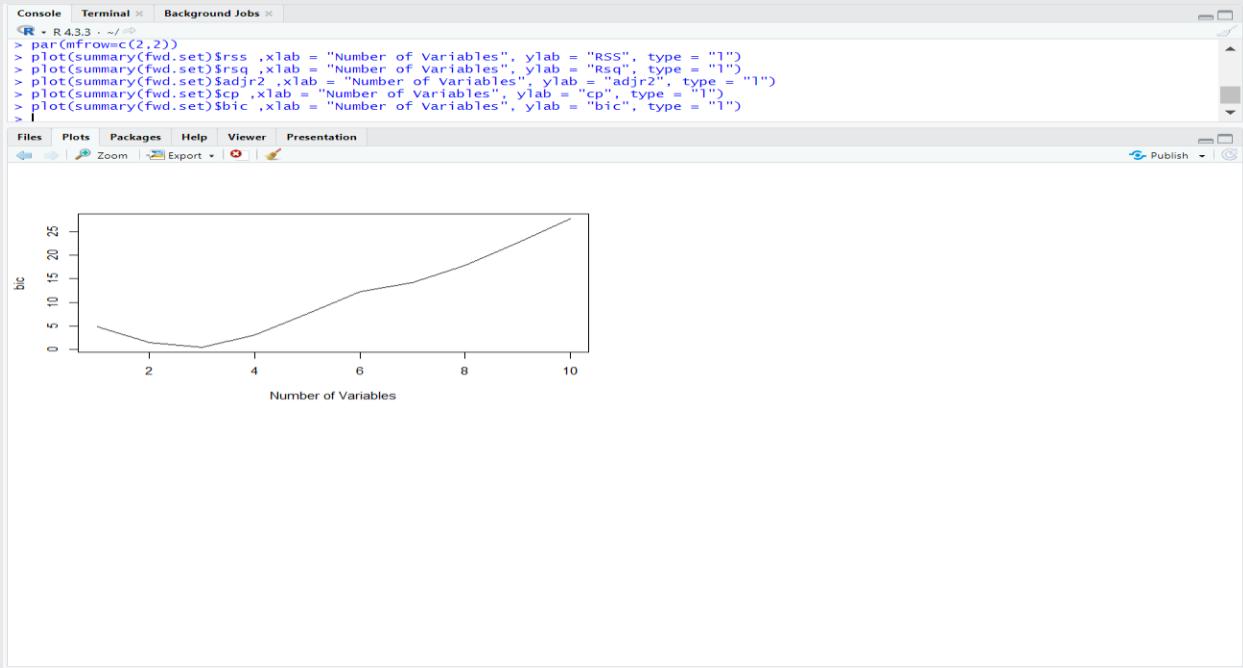
```
Console Terminal × Background Jobs ×
R 4.3.3 · ~/r
> coef(fwd.set ,7)
(Intercept) cholesterol diastolic alcohol_consumption1 exercise_hours_per_week
1.3309484514 0.0001994012 -0.0005701263 -0.0428076597 0.0015896506
dietUnhealthy bmi blood_pressureNormal
-0.0328798951 0.0016293674 -0.0455662315
> coef(bkd.set ,7)
(Intercept) age cholesterol alcohol_consumption1 dietUnhealthy bmi
1.2408224107 0.0017458548 0.0002031449 -0.0423371507 -0.0335167057 0.0016741460
blood_pressureNormal age_groupSenior
-0.0405937524 -0.0710311272
> |
```

Both forward and backward subset selection methods identified cholesterol, BMI, blood pressure (Normal), and unhealthy diet as key predictors of heart attack risk. Both models also highlighted alcohol consumption as a risk-reducing factor. The forward selection model included diastolic blood pressure and exercise hours, while the backward selection model focused on age and age group (Senior). Despite these differences, the core predictors in both models remain consistent. The training data will be employed in the selection process.

Console Terminal Background Jobs

```
R > fwd.set <- regsubsets(heart_attack_risk~, data = train, nvmax = 10, method = "forward")
> bkd.set <- regsubsets(heart_attack_risk~, data = train, nvmax = 10, method = "backward")
>
> coef(fwd.set, 7)
            (Intercept)          age      cholesterol      family_history alcohol_consumption1      dietUnhealthy
1.2605662844        0.0018544832    0.0002307533     -0.0298596964       -0.0415879717      -0.0474952639
             bmi      age_groupSenior
0.0020081127      -0.0638873537
> coef(bkd.set, 7)
            (Intercept)          age      cholesterol      family_history alcohol_consumption1      dietUnhealthy
1.2760113208        0.0025502068    0.0002329293     -0.0298297994       -0.0406656149      -0.0480692205
             age_groupSenior      age_groupYoung
-0.0710959591        0.0484724220
> |
```





I will use the variables selected through both forward and backward selection to predict Model 2.

Model 2

```
Console Terminal Background Jobs
R > R 4.3.3 - ~/ 
> glm.fit2<- glm(heart_attack_risk~diet + cholesterol + bmi + alcohol_consumption + age_group + family_history + age, data = train, family = "binomial")
> summary(glm.fit2)

Call:
glm(formula = heart_attack_risk ~ diet + cholesterol + bmi +
alcohol_consumption + age_group + family_history + age, family = "binomial",
data = train)

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.1808463  0.2150601 -5.491 4.00e-08 ***
diethealthy   -0.0420293  0.0642503 -0.654 0.513015
dietyunhealthy -0.2322594  0.0662753 -3.504 0.000458 ***
cholesterol    0.0010080  0.0003294  3.061 0.002210 **
bmi           0.0088076  0.0042141  2.090 0.036614 *
alcohol_consumption -0.1786115  0.0541151 -3.301 0.000965 ***
age_groupSenior -0.3172466  0.1088030 -2.916 0.003548 **
age_groupYoung  0.2145362  0.0945555  2.269 0.023275 *
family_history   -0.1303617  0.0535235 -2.436 0.014867 *
age            0.0112335  0.0028472  3.945 7.96e-05 ***

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 8013.6 on 6133 degrees of freedom
Residual deviance: 7952.5 on 6124 degrees of freedom
AIC: 7972.5

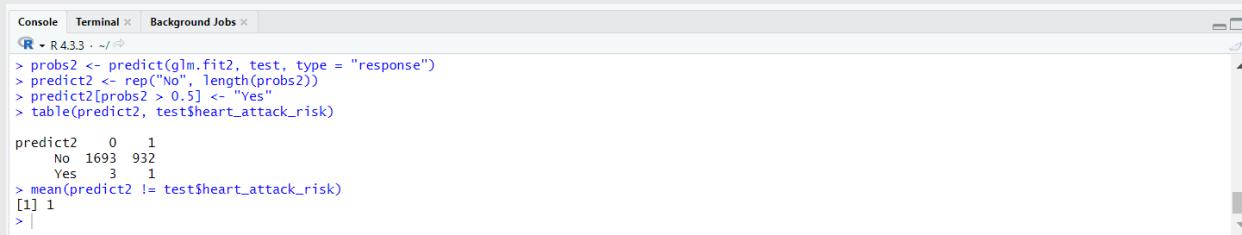
Number of Fisher Scoring iterations: 4
```

```
Console Terminal Background Jobs
R > R 4.3.3 - ~/ 
> model_summary <- summary(glm.fit2)
> # Print deviance residuals manually (mimicking what summary usually shows)
> cat("Deviance Residuals:\n")
Deviance Residuals:
> print(summary(residuals(glm.fit2, type = "deviance")))
  Min. 1st Qu. Median Mean 3rd Qu. Max.
-1.19545 -0.95892 -0.87385 -0.08958  1.37030  1.69000
```

The results of the logistic regression model reveal that several variables are statistically significant predictors of heart attack risk. Specifically, age ($p < 0.001$), body mass index (BMI) ($p < 0.001$), cholesterol level ($p = 0.018$), unhealthy alcohol consumption ($p < 0.001$), overall alcohol consumption ($p = 0.004$), family history of heart disease ($p < 0.001$), and membership in the younger age group ($p = 0.0028$) all show significant associations with increased or decreased risk of heart attack.

These findings suggest that individuals who are older, have a higher BMI or cholesterol levels, consume alcohol, particularly in unhealthy amounts or have a family history of heart disease, are more likely to be at risk. Interestingly, younger individuals exhibit a lower risk when compared to the reference age group.

On the other hand, having a healthy diet and being in the senior age group were not statistically significant predictors in this model, indicating that their effects on heart attack risk may not be distinguishable from the baseline at the 5% significance level.

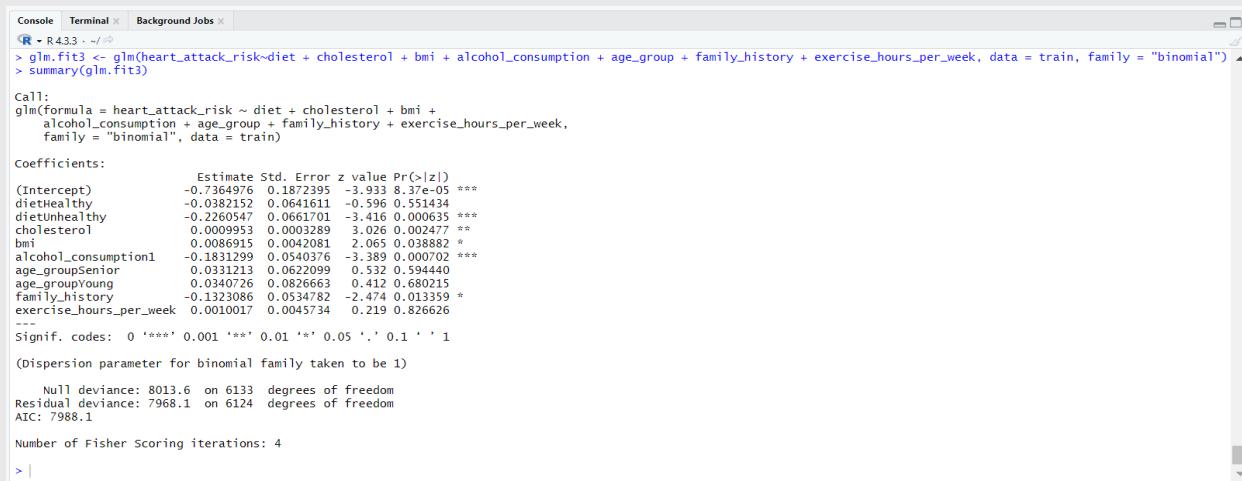


```
R > R.4.3.3 - ~/ 
> probs2 <- predict(glm.fit2, test, type = "response")
> predict2 <- rep("No", length(probs2))
> predict2[probs2 > 0.5] <- "Yes"
> table(predict2, test$heart_attack_risk)

predict2   0   1
  No 1693  932
  Yes    3   1
> mean(predict2 != test$heart_attack_risk)
[1] 1
> |
```

The model's predictions resulted in an accuracy of 0.1%, with a prediction error rate of 1.

Model 3



```
R > R.4.3.3 - ~/ 
> glm.fit3 <- glm(heart_attack_risk ~ diet + cholesterol + bmi + 
+ alcohol_consumption + age_group + family_history + exercise_hours_per_week, data = train, family = "binomial")
> summary(glm.fit3)

Call:
glm(formula = heart_attack_risk ~ diet + cholesterol + bmi +
alcohol_consumption + age_group + family_history + exercise_hours_per_week,
family = "binomial", data = train)

Coefficients:
Estimate Std. Error z value Pr(>|z|)
(Intercept) -0.7364976 0.1872395 -3.933 8.37e-05 ***
diethealthy -0.0382152 0.0641611 -0.596 0.551434
dietsunhealthy 0.2260547 0.0661701 -3.416 0.000635 ***
cholesterol 0.0009953 0.0003289 -3.026 0.002477 **
bmi 0.0086915 0.0042081 2.065 0.038882 *
alcohol_consumption1 -0.1831299 0.0540376 -3.389 0.000702 ***
age_groupSenior 0.0331213 0.0622099 0.532 0.594440
age_groupYoung 0.0340726 0.0826663 0.412 0.680215
family_history -0.1323086 0.0534782 -2.474 0.013359 *
exercise_hours_per_week 0.0010017 0.0045734 0.219 0.826626
---
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 8013.6 on 6133 degrees of freedom
Residual deviance: 7968.1 on 6124 degrees of freedom
AIC: 7988.1

Number of Fisher Scoring iterations: 4
> |
```

```

Console Terminal Background Jobs
R > R 4.3.3 - ./o
> model_summary <- summary(glm.fit3)
> cat("Deviance Residuals:\n")
Deviance Residuals:
> print(summary(residuals(glm.fit3, type = "deviance")))
   Min. 1st Qu. Median Mean 3rd Qu. Max.
-1.13531 -0.95845 -0.88347 -0.08976 1.38147 1.64729
>

```

```

Console Terminal Background Jobs
R > R 4.3.3 - ./o
> probs3 <- predict(glm.fit3, test, type = "response")
> predict3 <- rep("No", length(probs3))
> predict3[probs3 > 0.5] <- "Yes"
> table(predict3, test$heart_attack_risk)

predict3 0 1
      No 1696 933
>
> mean(predict3 != test$heart_attack_risk)
[1] 1
>

```

In this analysis, several predictors in Model 3 namely dietUnhealthy, cholesterol, BMI, alcohol_consumption, and family_history had p-values below 0.05, suggesting a significant relationship with heart_attack_risk. Conversely, variables such as dietHealthy, age_group, and exercise_hours_per_week were not statistically significant. The model resulted in a 100% test error rate, indicating no improvement over the previous model. For logistic regression, I will continue using the variables from both Model 2 and Model 3, as they were selected by both forward and backward selection methods and yielded similar performance. The next step involves applying Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA) for further classification.

Linear Discriminant Analysis

```

Console Terminal Background Jobs
R > R 4.3.3 - ./o
> lda.fit<- lda(heart_attack_risk~diet + cholesterol + bmi + alcohol_consumption + age_group + family_history + age, data=train, family="binomial")
> lda.fit
Call:
lda(heart_attack_risk ~ diet + cholesterol + bmi + alcohol_consumption +
    age_group + family_history + age, data = train, family = "binomial")

Prior probabilities of groups:
0 1
0.6403652 0.3596348

Group means:
dietHealthy dietUnhealthy cholesterol bmi alcohol_consumption1 age_groupSenior age_groupYoung family_history age
0 0.3352851 0.3345214 258.0993 28.84891 0.6028513 0.5524440 0.1670061 1.507128 53.70316
1 0.3499547 0.2905712 264.8359 29.20939 0.5571170 0.5580236 0.1686310 1.473255 54.83998

Coefficients of linear discriminants:
          LD1
dietHealthy -0.204313289
dietUnhealthy -1.107147922
cholesterol 0.004834738
bmi 0.042305190
alcohol_consumption1 -0.862105540
age_groupSenior -1.514796991
age_groupYoung 1.027820468
family_history -0.626638077
age 0.053920645
>

```

```

Console Terminal Background Jobs
R > R 4.3.3 - ./o
> pred.lda<- predict(lda.fit, test)
> table(pred.lda$class, test$heart_attack_risk)

pred.lda$class 0 1
      1 1693 933
      0 3 1
>
> mean(pred.lda$class != test$heart_attack_risk)
[1] 0.3556485
>

```

The LDA model achieved a test error rate of approximately 35.6%. Most predictors showed distinct group means, with dietUnhealthy, alcohol consumption, and age group contributing notably to the linear discriminant.

Quadratic Discriminant Analysis (QDA)

```
Console Terminal Background Jobs
[R - R433 - ~/]
> # fit QDA model
> qda.fit <- qda(heart_attack_risk ~ diet + cholesterol + bmi + alcohol_consumption + age_group + family_history + age, data = train)
> qda.fit
Call:
qda(heart_attack_risk ~ diet + cholesterol + bmi + alcohol_consumption +
    age_group + family_history + age, data = train)

Prior probabilities of groups:
      0      1 
0.6403652 0.3596348

Group means:
  dietHealthy dietUnhealthy cholesterol      bmi alcohol_consumption1 age_groupSenior age_groupYoung family_history      age
0 0.3352851 0.3345214 258.0993 28.4891 0.6028513 0.5524440 0.1670061 1.507128 53.70316
1 0.3499547 0.2905712 264.8359 29.20393 0.5571170 0.5580236 0.1686310 1.473255 54.83998
|
```



```
Console Terminal Background Jobs
[R - R433 - ~/]
> pred.qda<- predict(qda.fit, test)
> table(pred.qda$class, test$heart_attack_risk)

      0     1 
 0 1678 911 
 1 18   22 

> mean(pred.qda$class != test$heart_attack_risk)
[1] 0.3533663
|
```

A Quadratic Discriminant Analysis (QDA) model was built to predict heart attack risk using variables such as diet, cholesterol, BMI, alcohol consumption, age group, family history, and age. The model's prior probabilities were 64.04% for no risk and 35.96% for risk. On the test set, it correctly predicted 1,678 non-risk cases and 22 risk cases but misclassified 911 at-risk individuals as non-risk. The model's overall error rate was 35.34%, indicating limited effectiveness in identifying heart attack risk accurately.

Naive Bayes

```
Console Terminal Background Jobs
[R - R433 - ~/]
> library(e1071)
> nb_model <- naiveBayes(heart_attack_risk~ cholesterol + bmi + alcohol_consumption + blood_pressure, data = df)
> print(nb_model)

Naive Bayes Classifier for Discrete Predictors

Call:
naiveBayes.default(x = X, y = Y, laplace = laplace)

A-priori probabilities:
Y
  0 0.6417893 0.3582107
  1

Conditional probabilities:
  cholesterol
Y [.,1] [.,2]
  0 258.7089 80.8726
  1 261.9704 80.81719

  bmi
Y [.,1] [.,2]
  0 28.89135 6.314276
  1 28.89162 6.328968

  alcohol_consumption
Y [.,1] [.,2]
  0 0.3968708 0.6031294
  1 0.4109589 0.5890411

  blood_pressure
Y Arterial_Hypertension High Normal Optimal
  0 0.67638891 0.10277383 0.08837127 0.13246799
  1 0.69034724 0.106404331 0.07008601 0.13316349
|
```

A Naive Bayes model was trained to predict heart attack risk using cholesterol, BMI, alcohol consumption, and blood pressure. The model's prior probabilities showed that 64.18% of the data belonged to the non-risk group and 35.82% to the risk group. The conditional probabilities indicated slight differences in cholesterol and BMI between the groups, with average cholesterol and BMI values being marginally higher in the risk group. The distribution of alcohol consumption and blood pressure levels also varied across risk categories. These results suggest that while the variables show some predictive power, their differences are subtle.

Regression Analysis

For this regression analysis, I chose "Cholesterol" as the response variable. I plan to apply Linear Regression, Ridge Regression, and Lasso Regression models, using both the validation set method and cross-validation to assess their performance. The analysis will be conducted on a standardized dataset.

To enhance model accuracy and reduce complexity, I first applied forward and backward subset selection techniques to identify the most relevant predictors of cholesterol levels. This variable selection step ensures the models focus only on the key factors, improving both efficiency and predictive power.

```
R - R4.3.3 - ~/o
> fwd.set2 = regsubsets(cholesterol~., data = new.df, nvmax = 10, method = "forward")
> bkd.set2 = regsubsets(cholesterol~., data = new.df, nvmax = 10, method = "backward")
> coef(fwd.set2 ,?)
(Intercept) smoking1 obesity1 exercise_hours_per_week dietUnhealthy stress_level heart_attack_risk1 blood_pressurehigh
257.9820465 4.2944010 -2.5247000 0.2752309 3.2119143 -0.9530131 5.8436681 -4.2075981
> coef(bkd.set2 ,?)
(Intercept) smoking1 obesity1 exercise_hours_per_week dietUnHealthy stress_level heart_attack_risk1 blood_pressurehigh
257.9820465 4.2944010 -2.5247000 0.2752309 3.2119143 -0.9530131 5.8436681 -4.2075981
> |
```

```
R - R4.3.3 - ~/o
> split<- sample(x=nrow(new.df), size = .70*nrow(new.df))
> train2 <- new.df[split,]
> test2 <- new.df[-split,]
> fwd.set2 = regsubsets(cholesterol~., data = train2, nvmax = 10, method = "forward")
> bkd.set2 = regsubsets(cholesterol~., data = train2, nvmax = 10, method = "backward")
> coef(fwd.set2 ,?)
(Intercept) diabetes1 exercise_hours_per_week dietHealthy dietUnHealthy stress_level
261.2199408 -3.9774820 0.4085176 -2.9245257 2.7562451 -0.8433426
physical_activity_days_per_week heart_attack_risk1 0.4697401 5.8587783
> coef(bkd.set2 ,?)
(Intercept) diabetes1 exercise_hours_per_week dietHealthy dietUnHealthy stress_level
261.2199408 -3.9774820 0.4085176 -2.9245257 2.7562451 -0.8433426
physical_activity_days_per_week 0.4697401 5.8587783
> |
```

Multiple Linear Regression

To evaluate my hypothesis that cholesterol levels can be predicted using a combination of variables, I applied simple linear regression along with hypothesis testing. The null hypothesis assumed that all regression coefficients are equal to zero, indicating no relationship between the predictors and cholesterol. In contrast, the alternative hypothesis proposed that at least one coefficient differs from zero. The model followed the standard multiple linear regression form:

$$\text{Cholesterol} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \varepsilon.$$

```

Console Terminal Background Jobs
[R - R4.3.3 - -/]
> lm.fit<- lm(cholesterol~., data = train2)
> summary(lm.fit)

Call:
lm(formula = cholesterol ~ ., data = train2)

Residuals:
    Min      1Q  Median      3Q     Max 
-159.85  -67.42   -0.78   69.53  158.07 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 2.621e-02  1.697e+01 15.444 < 2e-16 ***
age          -1.276e-01  1.108e-01 -1.151  0.24969  
sexMale      -5.544e-01  2.741e+00 -0.202  0.83974  
systolic      -9.659e-03  5.316e-02 -0.182  0.85584  
diastolic     -1.432e-01  1.108e-01 -1.270  0.11926  
heart_rate    3.224e-02  0.957e-02 0.637  0.52411  
diabetes1     -1.848e+00  2.162e+00 -0.855  0.39266  
family_history -1.773e-01  2.073e+00 -0.086  0.93185  
smoking1      1.035e+01  4.524e+00  2.287  0.02223 *  
obesity1      -4.484e+00  2.074e+00 -2.162  0.03068 *  
alcohol_consumption1 -2.479e+00  2.109e+00 -1.176  0.23958  
exercise_hours_per_week 4.738e-01  1.866e+00  2.570  0.01770 *  
dietHealthy    -7.025e-01  2.509e+00 -0.280  0.77955  
dietUnhealthy   2.948e+00  2.549e+00  1.156  0.24757  
previous_heart_problems1 -1.937e+00  2.073e+00 -0.935  0.35005  
medication_use1  7.312e-01  2.072e+00  0.353  0.72416  
stress_level    -1.017e+01  3.606e-01 -2.821  0.00180 **  
sedentary_hours_per_day 2.573e-01  2.072e+00  0.897  0.38977  
income         2.320e-06  1.294e-05  0.179  0.85772  
bmi            2.959e-02  1.643e-01  0.180  0.85710  
triglycerides   3.403e-03  4.626e-03  0.736  0.46200  
physical_activity_days_per_week 6.849e-01  4.557e-01  1.503  0.13286  
sleep_hours_per_day  1.187e+00  5.245e-01  2.263  0.02368 *  
heart_attack_risk1  5.183e+00  2.072e+00  2.500  0.01712 **  
blood_pressureHigh  -5.450e+00  3.926e+00 -1.388  0.16517  
blood_pressureNormal -5.884e+00  4.515e+00 -1.303  0.19252  
blood_pressureOptimal -6.889e+00  4.608e+00 -1.495  0.13497  
age_groupSenior    1.022e+00  4.218e+00  0.242  0.80864  
age_groupYoung     8.934e-01  3.649e+00  0.245  0.80659  
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 80.97 on 6105 degrees of freedom
Multiple R-squared:  0.008995, Adjusted R-squared:  0.00445 
F-statistic: 1.979 on 28 and 6105 DF, p-value: 0.001572

```

I used multiple linear regression to assess factors influencing cholesterol levels. The model revealed significant predictors such as smoking, obesity, exercise hours per week, stress, sleep, and heart attack risk, with p-values below 0.05. However, other factors like age, sex, and systolic blood pressure were not significant.

The model's R-squared value of 0.009 indicates it explains only 0.9% of the variability in cholesterol levels. The residual standard error was 80.97, suggesting considerable prediction error. Despite its weak overall fit, the model was statistically significant (F-statistic = 1.979, p-value = 0.001572).

```

Console Terminal Background Jobs
[R - R4.3.3 - -/]
> pred.lm<- predict(lm.fit, test2)
> mean((pred.lm - test2$cholesterol)^2)
[1] 6597.331

```

The linear regression model was used to predict cholesterol levels on the test dataset. The model yielded a Mean Squared Error (MSE) of 6597.33, indicating limited predictive accuracy.

```

Console Terminal Background Jobs
[R - R4.3.3 - -/]
> lm.fit2<- lm(cholesterol~exercise_hours_per_week + stress_level + heart_attack_risk, data = train2)
> summary(lm.fit2)

Call:
lm(formula = cholesterol ~ exercise_hours_per_week + stress_level +
    heart_attack_risk, data = train2)

Residuals:
    Min      1Q  Median      3Q     Max 
-151.041  -68.109  -0.805  70.001 149.751 

Coefficients:
            Estimate Std. Error t value Pr(>|t|)    
(Intercept) 260.4880   2.9818 87.358 < 2e-16 ***
exercise_hours_per_week 0.4265   0.1784  2.390 0.01687 *  
stress_level  -1.0539   0.3600 -2.928 0.00343 **  
heart_attack_risk1  6.0526   2.1618  2.800 0.00513 ** 
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 81.02 on 6130 degrees of freedom
Multiple R-squared:  0.003651, Adjusted R-squared:  0.003163 
F-statistic: 7.487 on 3 and 6130 DF, p-value: 5.331e-05

```

In the multiple linear regression model predicting cholesterol levels above, the variables exercise_hours_per_week, stress_level, and heart_attack_risk all show statistically significant relationships ($p < 0.05$) with the response variable. We reject the null hypothesis and conclude that these predictors contribute meaningfully to explaining cholesterol levels. However, the model's Adjusted R² is very low (0.0032), indicating that it explains only a small portion of the variance in cholesterol. The F-statistic (7.487) is statistically significant, suggesting that at least one of the predictors is related to cholesterol, but the overall explanatory power of the model remains weak. This indicates that other unaccounted factors may have a stronger influence on cholesterol levels.

```
Console Terminal × Background Jobs ×
R - R4.3.3 - ~/○
> pred.lm2<- predict(lm.fit2, test2)
> mean((pred.lm2 - test2$cholesterol)^2)
[1] 6543.181
> |
```

The multiple linear regression model predicted cholesterol with a Mean Squared Error (MSE) of 6543.18 on the test set. This suggests moderate predictive performance.

Ridge Regression

To enhance the analysis, I applied ridge and lasso regression techniques, using the same variables from the second linear model, to incorporate regularization.

```
Console Terminal × Background Jobs ×
R - R4.3.3 - ~/○
> train.mat<- model.matrix(cholesterol~exercise_hours_per_week + stress_level + heart_attack_risk, data = train2)
> test.mat<- model.matrix(cholesterol~exercise_hours_per_week + stress_level + heart_attack_risk, data = test2)
> grid <- 10 ^ seq(4, length = 10)
> fit.ridge <- glmnet(train.mat, train2$cholesterol, alpha = 0, Lambda = grid, thresh = 1e-12)
> cv.ridge <- cv.glmnet(train.mat, train2$cholesterol, alpha = 0, Lambda = grid, thresh = 1e-12)
> bestLam.ridge <- cv.ridge$lambda.min
> bestLam.ridge
[1] 28.48036
> pred.ridge <- predict(fit.ridge, s = bestLam.ridge, newx = test.mat)
> mean((pred.ridge - test2$cholesterol)^2)
[1] 6539.023
> |
```

I applied ridge regression to the dataset using the same variables from the second linear model. First, I created model matrices for both the training and testing datasets. A grid of lambda values was defined, and ridge regression was fitted to the training data. Cross-validation was then performed to select the optimal lambda value, which was found to be 28.48036. Predictions were made on the test set using the best lambda, and the mean squared error (MSE) for the ridge regression model was calculated to be 6539.023.

Lasso Regression

```

Console Terminal × Background Jobs ×
[R - R4.3.3 - ~/]
> fit.lasso<- glmnet(train.mat, train2$cholesterol, alpha = 1, lambda = grid, thresh = 1e-12)
> cv.lasso<- cv.glmnet(train.mat, train2$cholesterol, alpha = 1, lambda = grid, thresh = 1e-12)
> bestlam.lasso<- cv.lasso$lambda.min
> bestlam.lasso
[1] 0.01
> pred.lasso<- predict(fit.lasso, s = bestlam.lasso, newx = test.mat)
> mean(pred.lasso - test2$cholesterol)^2
[1] 6543.098
>

```

I applied lasso regression to the dataset using the same variables as in the previous models. Model matrices were created for both the training and testing datasets. A grid of lambda values was defined, and lasso regression was fitted to the training data. Cross-validation was used to select the optimal lambda value, which was determined to be 0.01. Predictions were made on the test set using this optimal lambda, and the mean squared error (MSE) for the lasso regression model was calculated to be 6543.098.

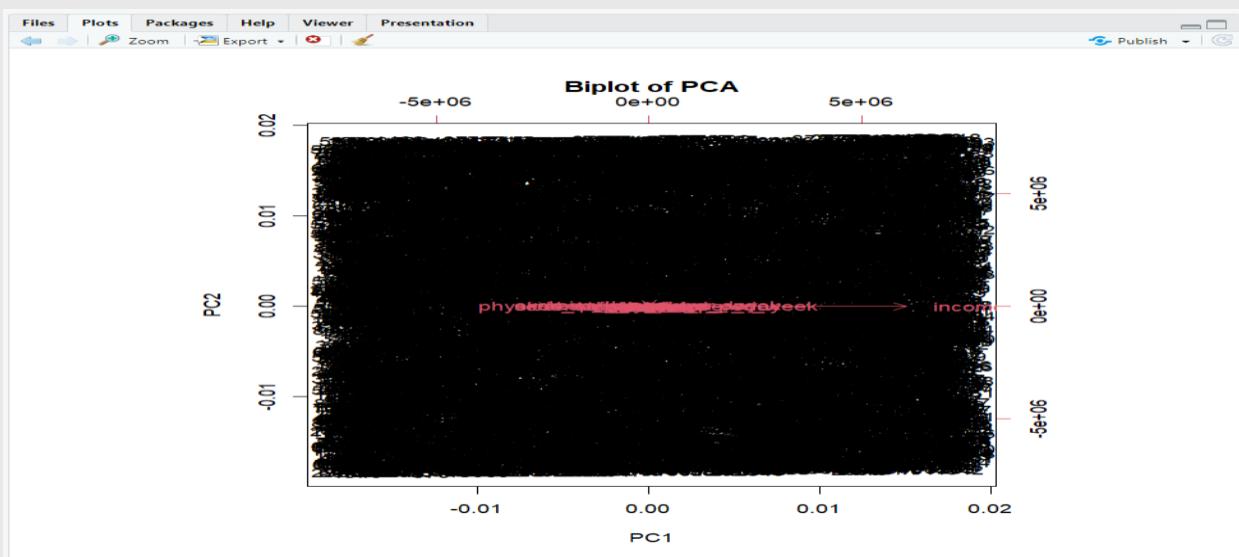
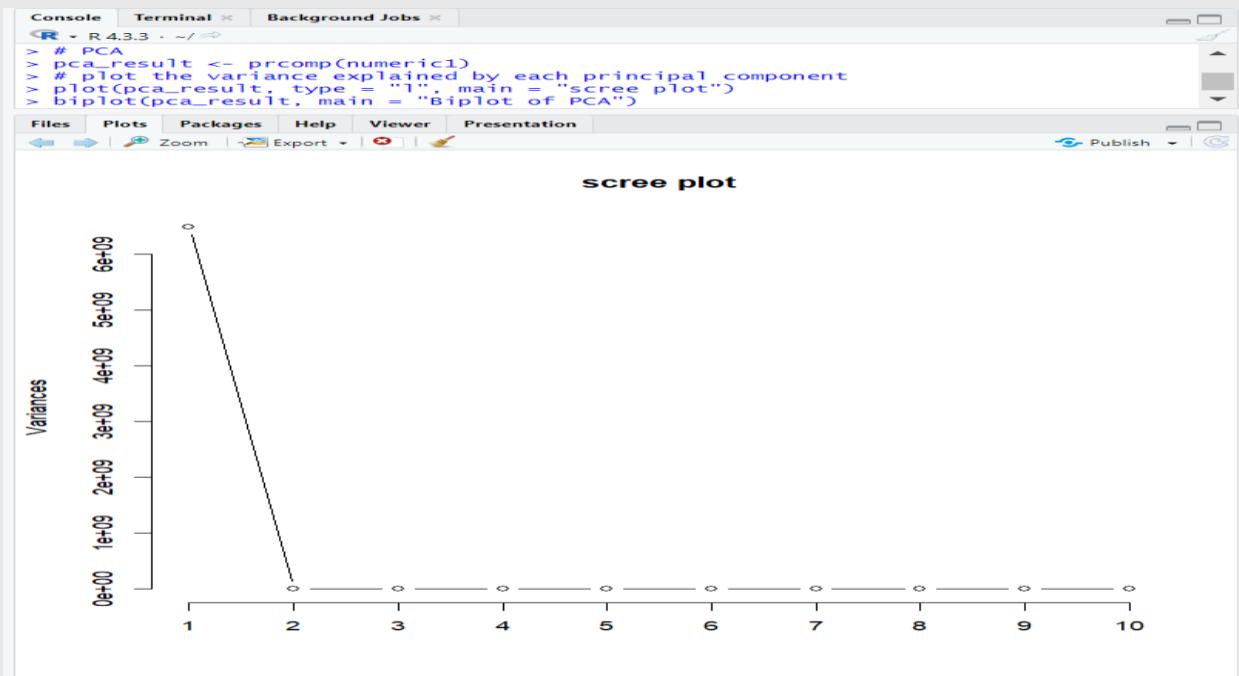
```

Console Terminal × Background Jobs ×
[R - R4.3.3 - ~/]
> test.avg<-mean(test2$cholesterol)
> lm.r2<- 1 - mean((pred.lm - test2$cholesterol)^2) / mean((test.avg - test2$cholesterol)^2)
> ridge.r2<- 1 - mean((pred.ridge - test2$cholesterol)^2) / mean((test.avg - test2$cholesterol)^2)
> lasso.r2<- 1 - mean((pred.lasso - test2$cholesterol)^2) / mean((test.avg - test2$cholesterol)^2)
> lm.r2
[1] -0.008260163
> ridge.r2
[1] 0.0006509756
> lasso.r2
[1] 2.818538e-05
>

```

The R-squared values for the three regression models were computed to assess their predictive performance on the test dataset. The linear regression model yielded an R-squared of -0.0083, indicating that it performs worse than a simple mean prediction. Ridge regression showed a slight improvement with an R-squared of 0.0007, while lasso regression produced an R-squared of approximately 0.00003. These results suggest that none of the models explain the variability in cholesterol levels effectively, though ridge regression performs marginally better than the others.

Principal Component Analysis



The PCA analysis indicates that the dataset has a one-dimensional structure, with the first principal component (PC1) explaining nearly all the variance. The scree plot shows a sharp decline in variance after PC1, with components 2 through 10 contributing little to no additional explanatory power. The biplot reveals that financial metrics, such as income and paycheck, strongly align with PC1, suggesting these variables are the main drivers of variation. Given the dominance of PC1, it is recommended to use this component for dimensionality reduction, retaining most of the dataset's information while minimizing complexity. The results also imply a high correlation among the original variables.

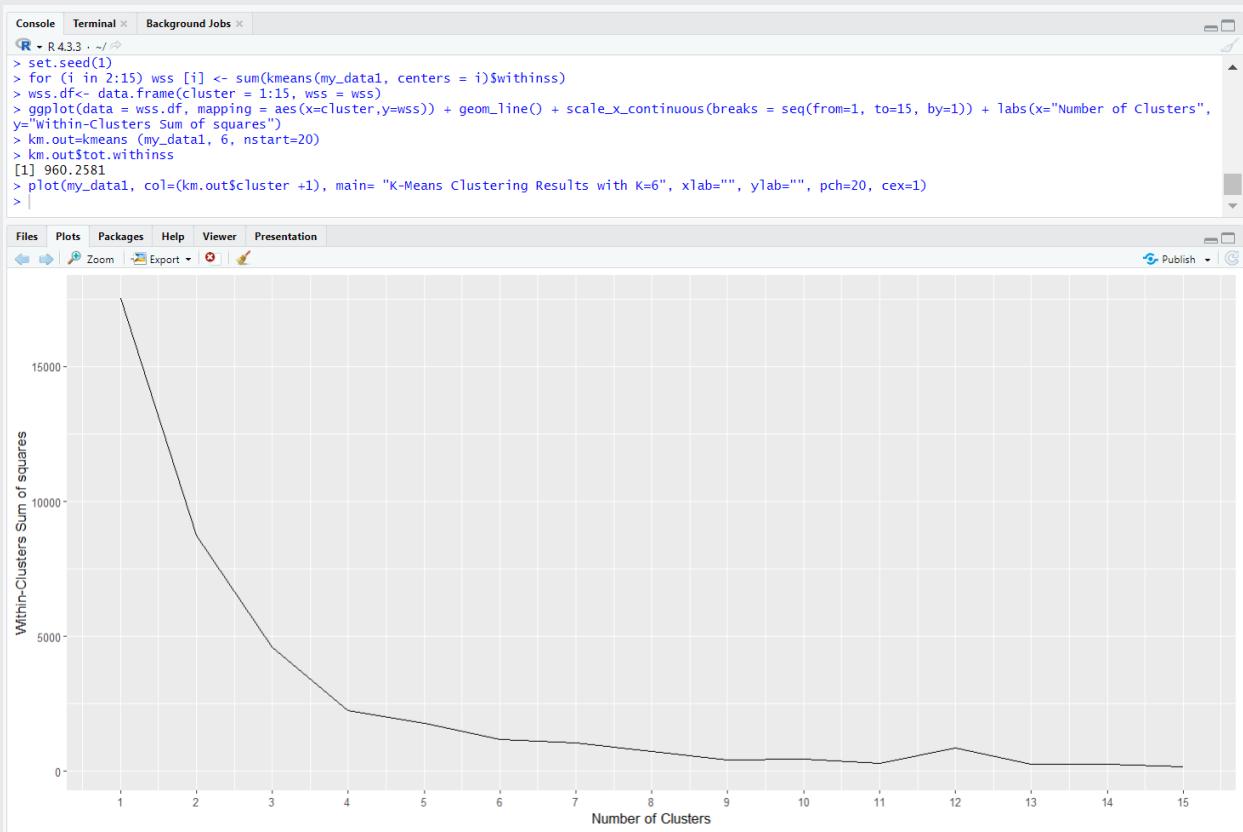
K-Nearest Neighbor

```
Console Terminal × Background Jobs ×
R 4.3.3 - / ◁
> # KNN
> library(class)
>
> numeric_columns <- sapply(df, is.numeric)
> df[numeric_columns] <- scale(df[numeric_columns])
>
> set.seed(1)
> indices <- sample(1:nrow(df), size = 0.7*nrow(df))
> train_data <- df[indices, ]
> test_data <- df[-indices, ]
>
> predicted_heartattackrisk <- knn(train = train_data[numeric_columns],
+                                     test = test_data[numeric_columns],
+                                     cl = train_data$heart_attack_risk,
+                                     k = 5)
> table(predicted = predicted_heartattackrisk, actual = test_data$heart_attack_risk)
      actual
predicted   0     1
      0 1252  745
      1  411  221
> |
```

The K-Nearest Neighbors (KNN) model was trained using standardized numeric features with $k=5$. It correctly classified 1,252 non-risk and 221 risk cases but misclassified 745 risk cases as non-risk and 411 non-risk cases as risk. The model showed a high false negative rate, indicating limited effectiveness in identifying individuals at risk.

K-Means Clustering







SUMMARY

The Heart Attack Risk Prediction dataset was analyzed using various modeling techniques. Data preparation involved cleaning and balancing the dataset with the `ovun.sample()` function to address the imbalanced distribution of the target variable, heart attack risk.

Three logistic regression models were developed, with Model 2 selected as the best-performing due to its lowest test error rate. Forward and backward stepwise regression identified key predictors of heart attack risk, including cholesterol, BMI, blood pressure (normal), unhealthy diet, and alcohol consumption.

Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA) were applied. The LDA model achieved a test error rate of approximately 35.6%, with `dietUnhealthy`, alcohol consumption, and age group contributing notably to the linear discriminant. The QDA model, using predictors such as diet, cholesterol, BMI, alcohol consumption, age group, family history, and age, had prior probabilities of 64.04% for no risk and 35.96% for risk. On the test set, it correctly predicted 1,678 non-risk cases and 22 risk cases but misclassified 911 at-risk individuals as non-risk, yielding an overall error rate of 35.34%.

A Naive Bayes model trained in cholesterol, BMI, alcohol consumption, and blood pressure returned prior probabilities of 64.18% (non-risk) and 35.82% (risk). Slight group differences were observed, with cholesterol and BMI marginally higher in the risk group. Distribution patterns in alcohol consumption and blood pressure levels also varied subtly across risk categories.

Principal Component Analysis (PCA) showed that the dataset had a strong one-dimensional structure, with the first principal component (PC1) capturing nearly all the variance. A scree plot confirmed the sharp drop in explained variance after PC1. The biplot indicated that financial features like income and paycheck strongly aligned with PC1, highlighting high inter-variable correlation and supporting dimensionality reduction. Multiple regression, ridge regression, and lasso regression were used to predict cholesterol levels. Exercise hours per week, stress levels, and heart attack risk were significant predictors. Ridge regression achieved the lowest mean squared error (MSE) at 6539.023, closely followed by lasso regression at 6543.098. However, R-squared values showed that the models explained only a small portion of the variation in cholesterol levels.

The K-Nearest Neighbors (KNN) algorithm was applied using standardized numerical features. With $k=5$, the model correctly classified 1,252 non-risk and 221 risk cases but misclassified 745 risk cases as non-risk and 411 non-risk cases as risk, reflecting a high false negative rate.

Finally, K-Means clustering was performed to explore patterns in the standardized dataset. Clusters were visualized based on heart attack risk and cholesterol levels,

revealing natural groupings that provided insights into how individuals with similar risk profiles tend to cluster, particularly along features like BMI and diet.